Provable Low-Frequency Bias of In-Context Learning of Representations

Editors: List of editors' names

Abstract

In-context learning (ICL) enables large language models (LLMs) to acquire new behaviors from the input sequence alone without any parameter updates. Recent studies have shown that ICL can surpass the original meaning learned in pretraining stage through internalizing the structure of the data-generating process (DGP) of the prompt into the hidden representations. However, the mechanisms by which LLMs achieve this ability is left open. In this paper, we present the first rigorous explanation of such phenomena by introducing a unified framework of double convergence, where hidden representations converge both over context and across layers. This double convergence process leads to an implicit bias towards smooth (low-frequency) representations, which we prove analytically and verify empirically. Our theory explains several open empirical observations, including why learned representations exhibit globally structured but locally distorted geometry, and why their total energy decays without vanishing. Moreover, our theory predicts that ICL has an intrinsic robustness towards high-frequency noise, which we empirically confirm. These results provide new insights into the underlying mechanisms of ICL, and a theoretical foundation to study it that hopefully extends to more general data distributions and settings.

Keywords: In-Context Learning, Inference Dynamics, Graph Spectral Methods

1. Introduction

It has become a major challenge for today's machine learning community to understand how large language models (LLMs) perform in-context learning (ICL), i.e. the ability to learn patterns or tasks solely from input sequences without any gradient updates (Brown et al., 2020; Min et al., 2022; Garg et al., 2022; Akyürek et al., 2022a). Empirical studies have demonstrated that LLMs can carry out a variety of tasks, including logical reasoning (Wei et al., 2022), programming (Gao et al., 2023), and solving mathematical problems (Hendrycks et al., 2021); recent work also suggests that LLMs are able to stay robust against noisy prompts (Cheng et al., 2025; Alazraki et al., 2025). However, the mechanisms underlying these capabilities remain largely elusive. A particularly striking phenomenon, recently highlighted by Park et al. (2024), shows that when a pre-trained LLM is fed a sequence generated by a random walk on a planar graph, with each node corresponds to a word (see "Data Generating Process" part of Figure 1), the model's hidden representations converge to a state that reflects the original graph structure, even though the graph itself was never explicitly provided. We refer to this emergent behavior as In-Context Learning of Representations (ICLR).

The ICLR phenomenon suggests an important mechanism of ICL: the model (in-contextly) learns to embed the information of the data-generating process (DGP) into the hidden representations. Therefore, understanding the ICLR phenomenon is a crucial step towards a deeper understanding of the mechanisms of ICL. Moreover, Park et al. (2024) also shows that an energy function decays over the course of this process, suggesting there might be an

underlying principle that drives the emergence of ICLR. However, it is left open what is the nature of this principle and how does it applies to the representations.

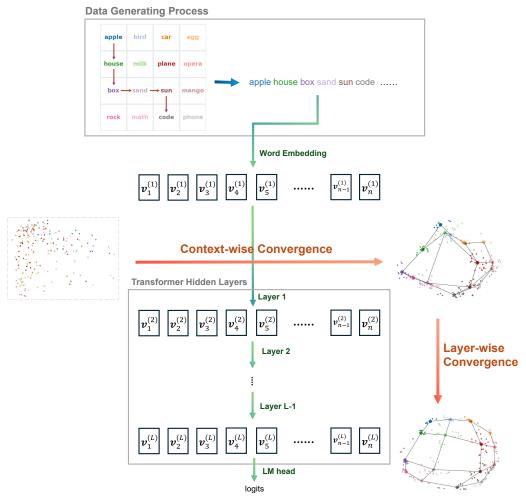


Figure 1: Overview of the DGP and the Double Convergence Process. The input sequence is generated by a random walk over a graph defined on the vocabulary (top), and then passed into a pre-trained Transformer model. Context-wise convergence occurs within each layer: token representations associated with the same word converge into tight clusters. Layer-wise convergence occurs across layers: the cluster centers gradually evolve towards a structure aligned with the underlying graph.

In this paper, we present the first theoretical explanation of the ICLR phenomenon, showing that it arises as a consequence of an intrinsic bias in LLMs towards low-frequency hidden representations. The central idea of our theoretical framework is a process we call **Double Convergence**, wherein hidden representations converge both along the context length and across layers. Specifically, the low-frequency bias emerges from the interaction of the **Context-wise Process** and the **Layer-wise Process**.

1. **Context-wise Process**: If the attention map "reflects" a function that is depends only on token identities (formally defined in Theorem 2 and Theorem 3), and the representations converge to a set of *latent representations* as the context length increases

Extended Abstract Track

(formally defined in Theorem 1), then the attention effectively applies the reflected function to the latent representations. This result is formally shown in Theorem 4;

2. Layer-wise Process: The *latent representations* across layers, and eventually converge to a state that captures the distributional properties of the input sequence. Under the specific DGP considered in Park et al. (2024), we prove that this convergence exactly yields the representations characteristic of the ICLR phenomenon. This result is formally stated in Theorem 7.

The concept of Double Convergence is illustrated in Figure 1. While our main analysis focuses on the specific DGP used in Park et al. (2024), the techniques and theoretical insights we develop are able to be extended beyond this particular setting. In Appendix I, we provide a generalized framework that is decoupled from any specific DGP, highlighting the broader applicability of our results.

Our theoretical framework provides a comprehensive explanation for several phenomena and open questions raised in Park et al. (2024). Furthermore, as both a validation and application of our theoretical results, we also demonstrate that ICL exhibits implicit robustness against high-frequency noise in the input data, which is consistent with recent empirical findings (Cheng et al., 2025; Alazraki et al., 2025).

In summary, our main contributions in this paper are as follows:

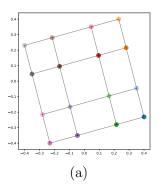
- 1. We identify the double convergence process, which serves as a general framework for studying the evolution of the representation in ICL (Appendices D and E);
- 2. We provide theoretical explanations for several previously unexplained phenomena in Park et al. (2024), including why ICL can suppress the original meaning of each word learned in pre-training (Appendix B.1), why the learned representations form an apparently regular yet slightly distorted structure (Appendices B.2 and B.3), and why the energy decreases but does not converge to zero (Appendix B.4);
- 3. Our theory highlights an implicit low-frequency bias in ICL, and predicts that LLMs are naturally robust to errors in the input prompts, which we have verified empirically (Appendix B.5).

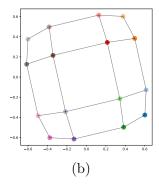
Due to space limit, we defer all the detailed theoretical results to appendix. Please see Appendices C to E for the formal definitions and assumptions used in our theoretical proof and the theoretical results. In the next section, we briefly discuss the implications of our theory and defer the full discussion to Appendix B.

2. A Brief Discussion About the Theoretical Results

Due to space limitation, in this section, we only briefly discuss the implications of our theory, and defer the full discussion to Appendix B.

Our main theoretical result (Theorem 7) predicts that hidden representations converge to the top eigenvectors of the graph matrix M, which correspond to smooth, low-frequency structures. This aligns with the classic literature on graph learning and spectral methods (Kipf and Welling, 2016; Li et al., 2018; Wu et al., 2019; Yang et al., 2021; Spielman, 2019; Trevisan, 2013), and explains why the ICL phenomenon captures global graph structure. Such low-frequency embeddings are also known to match human intuition in graph visualization (Tutte, 1963), naturally connecting theory with practice.





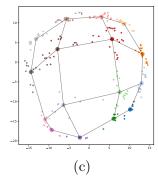
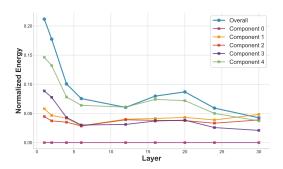


Figure 2: Comparison between empirical observations and theoretical predictions. (a): Idealized grid structure; (b): Theoretical prediction by Theorem 7; (c): Principal components of hidden representations from Llama-3.1-8B. All panels show strong alignment.

We empirically validate this prediction by comparing the principal components of hidden representations with the top eigenvectors of W, following the setup of Park et al. (2024). As shown in Figure 5, the theoretical predictions closely align with the empirical results. Interestingly, while the overall grid structure is preserved, distortions such as compression at the periphery appear. This effect can be explained by the reweighting induced by the stationary distribution of the attention process, rather than just uneven visitation frequency.

A surprising finding of Park et al. (2024) is that ICL can suppress the original semantic meaning of words: embeddings align with the data-generating process rather than natural semantics. Our perspective is that higher-frequency components, which often carry semantic detail, are gradually attenuated as depth and context length grow, leaving only the low-frequency structural signals.



Supplied to the state of the st

Figure 3: Normalized energy evolution across layers. Overall energy decreases, but low-frequency components persist.

Figure 4: **Predicted robustness** against noise. Despite noisy inputs, ICL restores the underlying structure.

Finally, our framework predicts LLM's robustness against noise: since high-frequency components decay naturally, ICL models should tolerate and correct moderate high-frequency noise in the input. Experiments with noisy sequences confirm this, as shown in Figure 7, where the model gradually restores the clean underlying graph structure. This observation is consistent with earlier empirical findings that ICL can implicitly denoise input (Cheng et al., 2025; Alazraki et al., 2025).

References

- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing Systems*, 36:45614–45650, 2023.
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. arXiv preprint arXiv:2211.15661, 2022a.
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. arXiv preprint arXiv:2211.15661, 2022b.
- Lisa Alazraki, Maximilian Mozes, Jon Ander Campos, Yi Chern Tan, Marek Rei, and Max Bartolo. Llms can implicitly learn from mistakes in-context. arXiv preprint arXiv:2502.08550, 2025.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Chen Cheng, Xinzhi Yu, Haodong Wen, Jingsong Sun, Guanzhang Yue, Yihao Zhang, and Zeming Wei. Exploring the robustness of in-context learning with noisy labels. In *ICASSP* 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 1–5. IEEE, 2025.
- Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. Why can gpt learn in-context? language models implicitly perform gradient descent as meta-optimizers. arXiv preprint arXiv:2212.10559, 2022.
- Jianqing Fan, Bai Jiang, and Qiang Sun. Hoeffding's inequality for general markov chains and its applications to statistical learning. *Journal of Machine Learning Research*, 22(139): 1–35, 2021.
- David J Field. Relations between the statistics of natural images and the response properties of cortical cells. *Journal of the Optical Society of America A*, 4(12):2379–2394, 1987.
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models. In *International Conference on Machine Learning*, pages 10764–10799. PMLR, 2023.
- Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- Borjan Geshkovski, Cyril Letrouit, Yury Polyanskiy, and Philippe Rigollet. A mathematical perspective on transformers. arXiv preprint arXiv:2312.10794, 2023.

- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. arXiv preprint arXiv:2103.03874, 2021.
- Yunzhe Hu, Difan Zou, and Dong Xu. Hyperspherical energy transformer with recurrent depth. arXiv preprint arXiv:2502.11646, 2025.
- Yu Huang, Yuan Cheng, and Yingbin Liang. In-context convergence of transformers. arXiv preprint arXiv:2310.05249, 2023.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. arXiv preprint arXiv:1609.02907, 2016.
- Gen Li, Yuchen Jiao, Yu Huang, Yuting Wei, and Yuxin Chen. Transformers meet in-context learning: A universal approximation theory. arXiv preprint arXiv:2506.05200, 2025.
- Qimai Li, Zhichao Han, and Xiao-Ming Wu. Deeper insights into graph convolutional networks for semi-supervised learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- Yue M Lu, Mary I Letey, Jacob A Zavatone-Veth, Anindita Maiti, and Cengiz Pehlevan. Asymptotic theory of in-context learning by linear attention. arXiv preprint arXiv:2405.11751, 2024.
- Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? arXiv preprint arXiv:2202.12837, 2022.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. In-context learning and induction heads. *Transformer Circuits Thread*, 2022. https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html.
- Core Francisco Park, Andrew Lee, Ekdeep Singh Lubana, Yongyi Yang, Maya Okawa, Kento Nishi, Martin Wattenberg, and Hidenori Tanaka. Iclr: In-context learning of representations. arXiv preprint arXiv:2501.00070, 2024.
- 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.
- D Spielman. Spectral and algebraic graph theory, incomplete draft, dated december 4, 2019, 2019.
- Akiyoshi Tomihari and Ryo Karakida. Recurrent self-attention dynamics: An energy-agnostic perspective from jacobians. arXiv preprint arXiv:2505.19458, 2025.

Extended Abstract Track

- Luca Trevisan. Lecture notes on expansion, sparsest cut, and spectral graph theory. *URl:* https://people.eecs.berkeley.edu/~luca/books/expanders.pdf, 2013.
- William Thomas Tutte. How to draw a graph. Proceedings of the London Mathematical Society, 3(1):743–767, 1963.
- Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pages 35151–35174. PMLR, 2023.
- Mingze Wang, Ruoxi Yu, Lei Wu, et al. How transformers implement induction heads: Approximation and optimization analysis. arXiv preprint arXiv:2410.11474, 2024.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplifying graph convolutional networks. In *International conference on machine learning*, pages 6861–6871. Pmlr, 2019.
- Yongyi Yang, Tang Liu, Yangkun Wang, Jinjing Zhou, Quan Gan, Zhewei Wei, Zheng Zhang, Zengfeng Huang, and David Wipf. Graph neural networks inspired by classical iterative algorithms. In *International Conference on Machine Learning*, pages 11773–11783. PMLR, 2021.
- Yongyi Yang, David P Wipf, et al. Transformers from an optimization perspective. Advances in Neural Information Processing Systems, 35:36958–36971, 2022.
- Yaodong Yu, Sam Buchanan, Druv Pai, Tianzhe Chu, Ziyang Wu, Shengbang Tong, Benjamin Haeffele, and Yi Ma. White-box transformers via sparse rate reduction. *Advances in Neural Information Processing Systems*, 36:9422–9457, 2023.

Appendix A. Background

In this section, we briefly review the current theoretical understanding of the mechanisms underlying ICL and the structure of hidden representations in Transformers. We highlight the limitations and challenges faced by existing studies and explain how the perspective adopted in this paper offers a potential path forward.

Mechanisms of In-Context Learning Understanding the mechanisms behind ICL has become a central topic in deep learning research. However, progress has been limited by the highly complex architecture and learning dynamics of neural networks. Existing theoretical approaches can be broadly categorized into two lines of work: 1) Existential results: This line of work constructs specific Transformer implementations that implement certain in-context algorithms, hence proving the existence of Transformers capable of performing in-context learning (Akyürek et al., 2022b; Dai et al., 2022; Von Oswald et al., 2023; Li et al., 2025). 2) Learning dynamics or loss landscape of simplified models: Another line of work studies learning dynamics or the structure of the loss landscape in simplified Transformer settings, typically with a small number of layers or restricted architectures. These studies show that the models can converge to configurations that exhibit in-context learning behaviors. However, due to the complexity of Transformer training dynamics, these analyses are usually restricted to one-layer Transformers (Lu et al., 2024; Huang et al., 2023), linearized models (Ahn et al., 2023), or two-layer models with controlled training setups (Wang et al., 2024).

Structure of Hidden Representations Another related line of theoretical researches investigates how hidden representations evolve across layers during inference (this topic is sometimes referred to as "inference dynamics") (Ramsauer et al., 2020; Yang et al., 2022; Yu et al., 2023; Geshkovski et al., 2023; Tomihari and Karakida, 2025; Hu et al., 2025). Despite their insights, these studies also typically require simplifying or modifying the Transformer architecture due to the non-linear and heterogeneous nature of Transformers. Yang et al. (2022) outlined four major challenges in analyzing inference dynamics, many of which remain unsolved.

Structure of Attention Maps Olsson et al. (2022), identified a specific attention mechanism known as induction heads. Generally speaking, they are attention heads that implement a form of token copying: they identify a previous occurrence of the current token and attend to the token that follows it. Formally, let the input tokens be $\{x_k\}_{k=1}^n$, generated by a Markov process, and the attention layer be defined as $\mathbf{u}_k = \sum_{j=1}^k a_{k,j} \mathbf{v}_k$, where $\{\mathbf{v}_k\}_{k=1}^n$ and $\{\mathbf{u}_k\}_{k=1}^n$ are input and output representations respectively, and $\{a_{k,j}\}_{k,j\in[n]}$ denotes the attention weights, then the induction heads can be defined as attention maps satisfying the following condition:

$$a_{k,j} > 0 \implies x_j \in \mathcal{N}(x_k),$$
 (1)

where $\mathcal{N}(x)$ is the set of all possible next tokens of x.

A major challenge in existing methods when analyzing inference dynamics arises from the interactive and heterogeneous structure of Transformer models: the attention map depends on the hidden representations, which in turn evolve through both self-attention and feedforward layers. This bidirectional dependency makes theoretical analysis extremely difficult, as noted by Yang et al. (2022). In this paper, to overcome this difficulty, we adopt a

more pragmatic approach: Rather than attempting to derive the structure of attention maps from first principles, we posit a structured form for the attention map, that is strong enough to enable meaningful theoretical results, yet general enough to be empirically validated and extendable to broader settings. Our goal is not to explain why attention maps take this form, but to demonstrate that, if they do, they can give rise to the observed structure in hidden representations, and empirically verify the validity of these assumptions in practice.

Appendix B. Detailed Discussion About the Theoretical Results

In Theorem 7, we proved that the representations converge to the top eigenspace of $M = D^{-1/2}WD^{-1/2}$. This matrix has been extensively studied in the literature on graph learning and spectral methods (Kipf and Welling, 2016; Li et al., 2018; Wu et al., 2019; Yang et al., 2021; Spielman, 2019; Trevisan, 2013). Specifically, since W is a symmetric matrix with non-negative entries, it defines a (weighted) undirected graph. Let \hat{L} be the symmetrically normalized Laplacian matrix of this graph (Spielman, 2019), it holds that $M = I - \hat{L}$, i.e. M and \hat{L} share the same set of eigenvectors, with the order of eigenvalues being reversed. Thus, top eigenvectors of M corresponds to low eigenvalues of \hat{L} , which is known to encode the low-frequency (smooth) and low-energy signals on the graph, as they tends to assign similar values to adjacent nodes. These low eigenvectors of \hat{L} are often used as coordinates for graph visualization, as it is know that they form figures that match human intuition (Tutte, 1963)¹, and exactly explains why the ICLR phenomenon, where the hidden representations encode global graph structure, emerges in such settings.

We confirm this theoretical prediction by reproducing the experiments in Park et al. (2024) and compare the principal components of the actual hidden representations and the analytical prediction, i.e. top eigenvectors of \mathbf{W} . The result is shown in Figure 5. It is clear that the analytical predictions align closely with the empirical principle components.

B.1. How Does ICL Suppress Original Semantic Meaning?

One of the most surprising observations in Park et al. (2024) is that, under their proposed DGP, ICL can produce word embeddings that no longer reflect the original semantic meaning of each word, but instead align solely with the structure imposed by the DGP. Our theory provides a natural explanation for this phenomenon: the "meaning" encoded in a word representation can be seen as a combination of multiple frequency components. However, as both the model depth and sequence length increase, higher-frequency components are progressively suppressed through the double convergence process. As a result, the semantic features associated with these higher-frequency components are effectively erased, and the representation becomes increasingly dominated by the low-frequency structure induced by the DGP.

B.2. Why Start From the 2nd Eigenvector?

In both Theorem 7 and Figure 5, we intentionally omit the first eigenvector of M. This is due to the coincidental alignment between the Laplacian and PCA. In short, the 1st eigenvector of M corresponds to a constant vector added to each z'_x ; however, PCA involves a **centralization step** that removes the mean component from the data. Specifically,

^{1.} Likely because humans also have a low-frequency bias in visual processing.

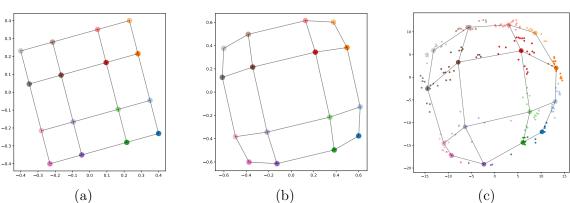


Figure 5: Comparison between empirical observations and theoretical predictions. (a): The 2nd and 3rd eigenvectors of $\widehat{\boldsymbol{W}}$, illustrating the idealized grid structure hypothesized in Park et al. (2024); (b): The 2nd and 3rd eigenvectors of \boldsymbol{M} , as predicted by Theorem 7; (c): The first two principal components of the actual hidden representations from a pre-trained Transformer (Llama-3.1-8B), collected at context positions 2360-2560. In all subfigures, each point represents the latent representation of a word in the vocabulary. The x- and y-axes represent the values of the 2nd and 3rd eigenvectors (or the 1st and 2nd principal components), respectively. A small rotation was applied to panels (a) and (b) to better visually align them with (c). This is valid since the 2nd and 3rd eigenvalues of \boldsymbol{W} are equal, and their eigenspace is invariant under rotation.

it can proved that the first eigenvector of M is exactly $d^{-1/2}$ (Spielman, 2019), and the centralization operation is to projecting $V_n^{(L)}$ onto the space orthogonal to $\mathbf{1}$, which is equivalent to projecting $V_n^{(L)} \mathbf{D}^{1/2}$ onto the space orthogonal to $d^{-1/2}$, effectively removing the component aligned with the first eigenvector of M.

B.3. Why Are Peripheral Nodes Compressed?

In Park et al. (2024), the authors keenly observed that the empirical PCA results (say, for the grid graph as in Figure 5(c)), despite roughly showing the underlying grid structure, appears slightly compressed near the periphery, compared to the actual grid formed by the eigenvectors of the original graph (as in Figure 5(a), or Figure 7 in Park et al. (2024)). The authors attributed this distortion to uneven visitation frequencies in the random walk:

... due to lack of periodic boundary conditions, concepts that are present in the inner 2×2 region of the grid are visited more frequently during a random walk on the graph, while the periphery of the graph has a lower visitation frequency.

While there is indeed a difference in visitation frequency, we argue that it is not the most fundamental explanation. The true cause lies in the context-wise process. As shown in Theorem 4, the transformation applied by the attention map to the latent representations is modulated by the stationary distribution π . As a result, the actual graph the model is aware of is the reweighted graph W instead of the original one \widetilde{W} . The top eigenvectors are twisted a little bit according to node degrees since it is reweighted by π .

B.4. Why Energy Decreases but Doesn't Vanish?

In Park et al. (2024), the authors hypothesized that the structure of the representation is a consequence of energy minimization. While this observation aligns with empirical trends, we argue that energy decay is not the fundamental cause, since 1) it doesn't explain why the model follows the principle of energy decaying, and 2) the energy does not actually converges to 0, despite the 0-energy solutions are actually easy to find. Instead, both the energy decay and the structured representations are consequences of the double convergence.

Formally, let Π_i be the projection operator onto the *i*-th eigenspace of M, the energy of the latent representation $\{z_x\}_{x\in[c]}$ under the graph W can be decomposed as follows:

$$\sum_{x,y\in[c]} w_{x,y} \|\mathbf{z}_x - \mathbf{z}_y\|^2 = \sum_{i=1}^c \sum_{x,y\in[c]} w_{x,y} \|\Pi_i(\mathbf{z}_x - \mathbf{z}_y)\|^2.$$
 (2)

As shown in Theorem 7, the projections of the latent representations onto the low eigenspaces converges to $\mathbf{0}^2$, and thus $\|\Pi_i(\mathbf{z}_x - \mathbf{z}_y)\|$ converges to 0 for large i. The energy decay is thus a consequence of the representation leaving corresponding eigenspace.

On the other hand, for the top eigenspaces of M (i.e. small i), the the projection components persist or even grow. This explains why the total energy does not decay to zero: the representation is leaving high-frequency eigenspaces, but accumulating energy in low-frequency ones.

We validate this explanation empirically in Figure 6. While the overall energy decreases across layers, the energy in low-frequency directions (e.g., Component 1 and 2) increases, confirming our theoretical prediction: energy decay arises from the attenuation of high-frequency components, whereas the persistence of low-frequency components prevents the total energy from converging to 0.

B.5. Predicted Robustness Against Noise

Our theoretical framework predicts that high-frequency components in the hidden representations will naturally decay over context and layers. This implies that LLMs performing ICL should be inherently robust to high-frequency noise. Since natural signals are typically dominated by low-frequency structure (Field, 1987), this suggests that ICL should be able to tolerate and even correct a moderate amount of errors in the input.

To test this prediction, we conduct an additional experiment under a noisy data-generating process. Specifically, during the random walk over the graph \mathcal{G} , we inject noise by allowing the sequence to transition to a uniformly random token in [c] with 1% probability at each step, rather than to a graph neighbor. This corruption can be viewed as temporarily replacing the original graph with a complete graph, which introduces purely high-frequency components into the sequence.

In Figure 7, we measure the number of non-neighbor transitions (i.e. token pairs that do not correspond to valid edges in \mathcal{G}) within a sliding window of the last 500 tokens. We compare this quantity in both the input sequence and the output predicted by the ICL model. While the input maintains a constant error rate due to the injected corruption, the ICL

^{2.} In principle, Theorem 7 is a relative result. However, with a similar proof one can show that the numerator also actually converge to 0 as long as the corresponding eigenvalues are significantly smaller than 1.

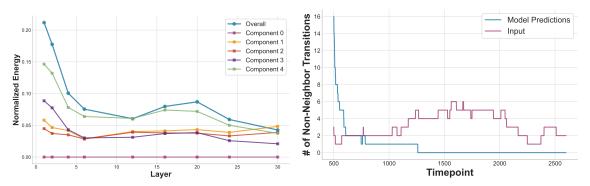


Figure 6: The Normalized energy evo-Figure 7: Predicted Robustness Against

lution across layers. Each "Component" Input Noise. We inject 1% random noise curve shows the energy along a direction de- into the input sequence and plot the number fined by the k-th eigenvector of M. We only of non-neighbor transitions within a sliding show the first 5 components as an illustration, window of the last 500 tokens. The "Overall" curve represents the average energy across all directions. To eliminate scale effects, the matrix $\mathbf{Z}^{(\ell)}$ is normalized to have unit Frobenius norm at each layer.

output gradually eliminates these errors. Once the context becomes sufficiently long, the model consistently produces transitions that respect the original graph structure, effectively achieving 100% accuracy despite the noisy input.

This result provides further evidence that ICL dynamics favor low-frequency structure and can suppress high-frequency perturbations. This result also explains previous observations that ICL can implicitly denoise input data (Cheng et al., 2025; Alazraki et al., 2025).

Appendix C. Preliminaries

Throughout this paper, we use bold upper-case letters to represent matrices (e.g. A), bold lower-case letters to represent vectors (e.g. x) and calligraphic upper-case letters to represent sequences of vectors (e.g. $\mathcal{V} = \{v_k\}_{k=1}^n$). For a matrix or a vector, we use plain lower-case letters to represent their entries (e.g. $a_{i,j}$ represents the i, j-th entry of A). For a number $n \in \mathbb{N}$, we denote $\{1, 2, \dots, n\}$ by [n]. For a set S, we use 2^S to represent its power set. We use $\mathbf{1}$ to represent a vector with all-entries being 1, whose dimensionality is inferred

from context. For a logical statement ϕ , we define $\mathbb{1}_{\{\phi\}} = \begin{cases} 1 & \phi \text{ is true} \\ 0 & \phi \text{ is false} \end{cases}$ to be its indicator

function. Given a sequence of vectors $\mathcal{Z} = \{z_x\}_{x=1}^c \in (\mathbb{R}^d)^c$, we define $\mathbf{Z} = \text{mat } \mathcal{Z} \in \mathbb{R}^{d \times c}$ to be the matrix formed by column-wise stacking of the vectors in \mathcal{Z} , i.e. the *i*-th column of \mathbf{Z} is \mathbf{z}_i . Given a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ and a sequence of vectors $\mathcal{V} = \{v_k\}_{k=1}^n \in (\mathbb{R}^d)^n$, let $\mathbf{A}\mathcal{V}$ to be a sequence defined as $\mathbf{A}\mathcal{V} = \left\{\sum_{j=1}^n a_{k,j} v_j\right\}_{k=1}^n$. For a vector function $\sigma : \mathbb{R}^d \to \mathbb{R}^d$ and a matrix $\mathbf{Z} \in \mathbb{R}^{c \times d}$, we use $\sigma(\mathbf{Z})$ to represent applying σ column-wisely to \mathbf{Z} .

C.1. Data Generation

Throughout this paper, we assume the input is a sequence of tokens $\mathcal{X} = \{x_k\}_{k=1}^n \in [c]^n$, where n is the sequence length and c is the vocabulary size (each token is simply a number in [c]). We assume that the sequence is sufficiently long, i.e., n > 10c.

Let $\mathcal{G} = ([c], \mathcal{E})$ be a connected undirected graph defined over the vocabulary (with each node being a word). Let $\widetilde{\boldsymbol{W}} = \{\widetilde{w}_{x,y}\}_{x,y\in[c]}$ denote the adjacency matrix of \mathcal{G} , and let $\boldsymbol{\pi} = \{\pi_x\}_{x\in[c]} \in \mathbb{R}^c$ be the stationary distribution \mathcal{G} (i.e. $\boldsymbol{\pi}$ is the L_1 normalized Perron vector of $\widetilde{\boldsymbol{W}}$).

We define the **reweighted adjacency** matrix $\boldsymbol{W} = \{w_{x,y}\}_{x,y \in [c]}$ as $w_{x,y} = \widetilde{w}_{x,y} \pi_x \pi_y$. Note that \boldsymbol{W} is also a non-negative symmetric matrix, and can therefore be viewed as the adjacency matrix of a reweighted version of \mathcal{G} . Define the (reweighted) degree vector $\boldsymbol{d} = \{d_x\}_{x \in [c]} \in \mathbb{R}^c$ as $d_x = \sum_{y \in [c]} w_{x,y}$. Let $\boldsymbol{D} = \operatorname{diag}(\boldsymbol{d})$ be the degree matrix of \boldsymbol{W} .

We assume the input sequence \mathcal{X} satisfies the following data-generating process: the first c tokens are fixed as a traversal of the vocabulary (i.e. $x_k = k$ for $k \in [c]$), and starting from x_{c+1} , the remaining tokens are generated by a random walk on the graph \mathcal{G} , with the initial token x_{c+1} being sampled from the stationary distribution π (i.e. the probability of $x_{c+1} = y$ is π_y)³.

Finally, for any word $x \in [c]$ and position $k \in [n]$, let $F_{x,k}$ be the frequency of token x in the first k elements of the sequence, i.e. $F_{x,k} = \sum_{j=1}^{k} \mathbb{1}_{\{x_j = x\}}$.

C.2. Model Architecture

Throughout the paper, we consider a simplified yet deep and non-linear Transformer model, as described in Algorithm 1, where $d \in \mathbb{N}$ is the hidden and input dimension, $L \in \mathbb{N}$ is the number of layers, $\mathbf{A}^{(\ell)} = \left\{a_{k,j}^{(\ell)}\right\}_{k,j\in[n]}$ is the (single-head) attention map at layer ℓ , and $\sigma^{(\ell)}: \mathbb{R}^d \to \mathbb{R}^d$ is the neuron-wise transformations at layer ℓ (which may include feedforward networks (FFNs), normalization layers, and other non-linearities).

As discussed in Appendix A, the most critical simplification in our model is that we treat the attention maps as given, instead of generated from hidden representations. In exchange for this simplification, we are able to explicitly characterize the structure of the hidden representations in a deep and non-linear model, enabling us to rigorously explain multiple in-context learning behaviors. This contrasts with prior works that remain entangled in the complexity of layer-wise interactions in full Transformer architectures. Moreover, we validate our assumption empirically on real models in Appendix K, finding that our assumptions on self-attention effectively explains more than 70% of actual attention connections, indicating its practical justifiability.

This model described in Algorithm 1 also omits several other standard components such as normalizations and residual connections. We note that this is to not over-complicating the theoretical results while still capturing the core mechanisms and challenges of the model. There is flexibility in our analysis to include other components, but we choose to focus on a

^{3.} This DGP is essentially the same as Park et al. (2024). We fix the first c tokens only to avoid trivial but complicated edge cases in the analysis, and since we study the asymptotic behavior of the model, the effect of the first a few tokens is negligible.

minimal version in the main paper to clearly highlight the key ideas behind our theoretical results. See Appendix J for further discussion on integrating additional components.

C.3. Methodology Outline

As noted before, we treat the attention maps in the Transformer as given, rather than dynamically generated from the hidden representations. Specifically, we assume the attention maps in the model are composed of a class of structured maps we refer to as **nice** attentions, formally defined in Theorem 3. These are attention maps whose connectivity patterns are determined by a function of the input tokens. This concept can be viewed as a generalization of the notion of induction heads (see eq. (1)), as illustrated by the structural similarity between eq. (1) and eq. (5).

Algorithm 1 Transformer forward process input:
$$\left\{ oldsymbol{v}_k^{(1)} \right\}_{k=1}^n \in \left(\mathbb{R}^d \right)^n$$
 as $oldsymbol{v}_k^{(1)} = oldsymbol{b}_{x_k}$. for $\ell = 1, 2 \cdots L - 1$ do for $k = 1, 2 \cdots n$ do $oldsymbol{u}_k^{(\ell)} = \sum_{j=1}^k a_{k,j}^{(\ell)} oldsymbol{v}_j^{(\ell)};$ $oldsymbol{v}_k^{(\ell+1)} = \sigma^{(\ell)} \left(oldsymbol{u}_k^{(\ell)} \right).$ end for end for output: $\left\{ oldsymbol{v}_k^{(L)} \right\}_{k=1}^n$.

Given this assumption, we are able to prove that the there is a **double convergence** process in the inference dynamics: 1) within each layer, the hidden representations converge to a set of "latent representations" that only encodes token identity and does not have position information; 2) then, each layer of nice attentions operates as a transformation on these latent representations. Consequently, the latent representations evolve and converge progressively across layers.

Furthermore, as the hidden representations can be viewed as having two axes, which we call the neuron axis and token axis, and the double convergence happens in the token axis, it can be shown that any transformations in the neuron axis, as long as being well conditioned, will not affect this double convergence process. This observation allows us to "insert" any neuron-wise transformations (such as such as FFNs and normalizations) between the self-attention layers without affecting the double convergence behavior, enabling a modular and robust theoretical framework.

Appendix D. Context-wise Convergence

We begin by establishing a general result: if the attention map reflects a specific function of the underlying tokens, and the representations converge to a set of latent representations, then the output of the attention layer also converges, towards a transformed set of latent representations. We start the presentation of this result by defining latent representations and nice attention maps.

Definition 1 For a sequence of d-dimensional vectors $\mathcal{U} = \{u_k\}_{k=1}^n \in (\mathbb{R}^d)^n$, if there exists a number $\gamma > 0$ and another sequence $\mathcal{Z} = \{z_x\}_{x \in [c]}$, such that

$$\forall k \in [n], \|\boldsymbol{u}_k - \boldsymbol{z}_{x_k}\| \le \frac{\gamma}{\sqrt{k}},\tag{3}$$

Extended Abstract Track

then we say \mathcal{U} is a **good sequence** converging to \mathcal{Z} with parameter γ , and \mathcal{Z} is the **latent** representation of \mathcal{U} .

The concept of latent representations captures the idea that the hidden representation for each token converges to a vector determined solely by the token identity and independent of position.

Definition 2 If a matrix $\mathbf{A} = \{a_{k,j}\}_{k,j \in [n]} \in \mathbb{R}_+^{n \times n}$ is lower-triangular and row-stochastic, i.e. it satisfies $a_{k,j} = 0$ for all j > k, and $\sum_{j=1}^k a_{k,j} = 1$ for all $k \in [n]$, then we say \mathbf{A} is an attention map. Moreover, for an attention map \mathbf{A} , if there exists a scalar $\psi > 0$, such that

$$\forall k \in [n], j \in [k], \sum_{i=1}^{j} a_{k,i} \le \frac{\psi j}{k}, \tag{4}$$

then we say A is a nice attention map with parameter ψ .

Nice attention maps are attention maps with a soft uniformity and locality: they prevent the attention from overly concentrating disproportionately on early tokens.

Definition 3 If $A \in \mathbb{R}^{n \times n}$ is a nice attention map with parameter ψ , and there is a function $f: [c] \to 2^{[c]}$, such that for all k > c and $j \in [k]$,

$$a_{k,j} > 0 \implies x_j \in f(x_k),$$
 (5)

and moreover,

$$\forall k \in [n], \forall y \in f(x), \left| \sum_{j \in [k]} a_{k,j} \mathbb{1}_{\{x_j = y\}} - \frac{F_{y,k}}{\sum_{y' \in f(x)} F_{y',k}} \right| \le \frac{\psi}{\sqrt{k}}$$
 (6)

then we say A reflects the function f.

Intuitively, nice attentions that reflects a function matches a functionally defined neighborhood of the current token, which can be viewed as a generalized notion of induction heads. Moreover, a nice attention that reflects a function also requires attention weights to distribute roughly proportionally among all words.

Notice that in Theorem 3, the k-th row of the attention map $\{a_{k,j}\}_{j=1}^n$ is only well defined only there exists $j \in [k]$ such that $x_j \in f(x_k)$ (otherwise all attention weights in this row are 0, violating the assumption that each row of A sums up to 1). This is guaranteed under our setup because we have explicitly set the first c tokens to be a traversal over the entire vocabulary.

We are now ready to state the main theorem of this section. Notice that this theorem stated here depends on the DGP, as it relies on the distribution of $F_{x,k}$. However, it is possible to prove a weaker but more general version of this theorem that is independent of the DGP. See Appendix I for more details.

Theorem 4 There exists a constant C > 0 that only depends on \mathcal{G} , and an event with probability at least 0.999, such that within this event, the following statement holds. Suppose $\mathcal{V} = \{\boldsymbol{v}_k\}_{k=1}^n \in (\mathbb{R}^d)^n$ is a good sequence converging to $\mathcal{Z} = \{\boldsymbol{z}_x\}_{x \in [c]}$ with parameter γ , and $\boldsymbol{A} = \{a_{k,j}\}_{k,j \in [n]} \in \mathbb{R}^{n \times n}$ is a nice attention map with parameter ψ that reflects a function $f: [c] \to 2^{[c]}$. Then $\boldsymbol{A}\mathcal{V}$ is a good sequence converging to

$$\mathcal{Z}' = \left\{ \frac{\sum_{y \in f(x)} \pi_y \mathbf{z}_y}{\sum_{y \in f(x)} \pi_y} \right\}_{x \in [c]}$$

$$(7)$$

with parameter $C\psi(\gamma + N) + C\log n$, where $N = \max_{y \in [c]} \|\boldsymbol{z}_y\|$.

Theorem 4 shows that applying a nice attention map that reflects a function to a good sequence yields another good sequence, and that the attention operation effectively acts on the latent representations. In other words, attention maps of this kind preserve convergence and transform latent representations in a token-consistent way.

Appendix E. Layer-wise Convergence

In Theorem 4, we showed how latent representations evolve under a single attention layer. In this section, we study how these latent representations change across layers.

Attention Maps. To analyze layer-wise convergence, we must have a more specific assumption on what exactly are the functions that the attention maps reflect. Specifically, we assume each attention map is a weighted combination of four basic types of maps: $A^{(\ell,A)}$, $A^{(\ell,B)}$, $A^{(\ell,O)}$ and $A^{(T)}$, that satisfies the following assumptions respectively:

- 1. A-type (self connections): $\mathbf{A}^{(\ell,A)}$ is a nice attention map with parameter $\psi_A^{(\ell)}$ that reflects $f_A: x \mapsto \{x\}$;
- 2. B-type (neighbor connections): $A^{(\ell,B)}$ is a nice attention map with parameter $\psi_B^{(\ell)}$ that reflects $f_B: x \mapsto \{y \in [c] | \widetilde{w}_{x,y} > 0\}$;
- 3. O-type (other connections): $A^{(\ell,O)}$ is a nice attention map with parameter $\psi_O^{(\ell)}$ that reflects $f_O: x \mapsto [c]$;
- 4. T-type (trivial connections, i.e. attention sinks): $\mathbf{A}^{(T)}$ satisfies $a_{i,j}^{(T)} = 1$ only when j = 1.

We assume the attention map at layer ℓ , i.e. $\boldsymbol{A}^{(\ell)}$, takes the form

$$\mathbf{A}^{(\ell)} = \rho_A^{(\ell)} \mathbf{A}^{(\ell,A)} + \rho_B^{(\ell)} \mathbf{A}^{(\ell,B)} + \rho_O^{(\ell)} \mathbf{A}^{(\ell,O)} + \rho_T^{(\ell)} \mathbf{A}^{(T)}, \tag{8}$$

where $\rho_{\tau}^{(\ell)} \geq 0$ ($\tau \in \{A, B, O, T\}$) is the weight of the τ -th type connections, and $\sum_{\tau \in \{A, B, O, T\}} \rho_{\tau}^{(\ell)} = 1$. Empirically, we find that these four types explain over 70% of attention connections in real models (see Appendix K for details), highlighting the empirical soundness of this classification.

Extended Abstract Track

The Role of FFN. To enable meaningful layer-wise convergence results with neuron-wise transformations involved, we also need assumptions on the nonlinearity $\sigma^{(\ell)}$ applied at each layer. We introduce the following concept.

Definition 5 For a set $\mathscr{Z} \subseteq \mathbb{R}^{d \times c}$, an orthonormal matrix $\mathbf{U} \in \mathbb{R}^{d \times p}$ and scalars $\gamma_1, \gamma_2 > 0$, if a function $\sigma : \mathbb{R}^d \to \mathbb{R}^d$ satisfies that for any matrix $\mathbf{Z} \in \mathscr{Z}$,

$$\gamma_1 \| \mathbf{Z} \mathbf{D}^{1/2} \mathbf{U} \| \le \| \sigma(\mathbf{Z}) \mathbf{D}^{1/2} \mathbf{U} \| \le \gamma_2 \| \mathbf{Z} \mathbf{D}^{1/2} \mathbf{U} \|,$$
 (9)

then we say σ is a great mapping w.r.t. $(\gamma_1, \gamma_2, \mathcal{Z}, U)$.

Intuitively, a great mapping is well-behaved (e.g. smooth) along a given subspace at the point of the latent representations. This allows us to insert FFNs or other neuron-wise transformations without disrupting convergence.

E.1. Main Results

Before stating the main results, we first define the concept of the spectral gap of a matrix.

Definition 6 For a symmetric matrix $\mathbf{M} \in \mathbb{R}^{c \times c}$ and index $q \in [c]$, let its eigenvalues be $\{\lambda_k\}_{k=1}^c$, arranging in a non-decreasing order of absolute values, define $\delta_q(\mathbf{M}) = \left|\frac{\lambda_q}{\lambda_{q+1}}\right|$ be the spectral gap of \mathbf{M} .

With the above assumptions and definitions, we are ready to present our main theorem.

Theorem 7 There exists an event with probability at least 0.99 such that the following statement holds. Let $n_{[x]}$ be the largest $k \in [n]$ such that $x_k = x$. Let $\mathbf{M} = \mathbf{D}^{-1/2}\mathbf{W}\mathbf{D}^{-1/2}$. Let the eigen-decomposition of \mathbf{M} be

$$\boldsymbol{M} = \begin{bmatrix} \boldsymbol{f}_1 & \boldsymbol{X} & \boldsymbol{Y} \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \boldsymbol{\Lambda} & \\ & & \boldsymbol{\Lambda}' \end{bmatrix} \begin{bmatrix} \boldsymbol{f}_1^\top \\ \boldsymbol{X}^\top \\ \boldsymbol{Y}^\top, \end{bmatrix}$$
(10)

where the eigenvalues are arranged in a non-increasing order of absolute values; $\boldsymbol{\Lambda}$ contains q-1 eigenvalues and $\boldsymbol{\Lambda}'$ contains n-q eigenvalues. Let $\boldsymbol{A}^{(\ell)}$ be defined as in eq. (8), and $\left\{\boldsymbol{v}_k^{(\ell)}\right\}_{k=1}^n$ be defined as in Algorithm 1. Suppose each $\sigma^{(\ell)}$ is a great mapping w.r.t. $\left(\gamma_1^{(\ell)}, \gamma_2^{(\ell)}, \mathbb{R}^d, \boldsymbol{X}\right)$ and $\left(\gamma_1'^{(\ell)}, \gamma_2'^{(\ell)}, \mathbb{R}^d, \boldsymbol{Y}\right)$, and have finite Lipschitz constant. Suppose there exists $\epsilon > 0$ satisfying that $\delta_q\left(\rho_A^{(\ell)}\boldsymbol{I} + \rho_B^{(\ell)}\boldsymbol{M}\right)\frac{\gamma_1^{(\ell)}}{\gamma_2'^{(\ell)}} \leq 1 - \epsilon$ for all $\ell \in [L]$. Then

$$\lim_{L \to \infty} \lim_{n \to \infty} \frac{\left\| \boldsymbol{V}_n^{(L)} \boldsymbol{D}^{1/2} \boldsymbol{Y} \right\|}{\left\| \boldsymbol{V}_n^{(L)} \boldsymbol{D}^{1/2} \boldsymbol{X} \right\|} = 0, \tag{11}$$

where $oldsymbol{V}_n^{(L)} = \max \left\{ oldsymbol{v}_{n_{[x]}}^{(L)}
ight\}_{x \in [c]}.$

Appendix F. Proof Theoretical Results w.r.t. Context-wise Convergence

In this section, we prove Theorem 4. The proof start by identifying a high-probability event in the random walk sequence that ensures it is "regular" enough.

F.1. Events in a Random Walk Sequence

Theorem 8 (Theorem 1 in (Fan et al., 2021)) Given graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with stationary distribution π , there exists constant C > 0 satisfies the following statement. Let $\{x_i\}_{i=1}^{\infty} \in \mathcal{V}^{\mathbb{N}}$ be a random walk sequence on \mathcal{G} starting from the stationary distribution, and $\{f_i\}_{i=1}^{\infty}$ be a sequence of bounded functions satisfying $f_i(\mathcal{V}) \subseteq [-\alpha_i, \alpha_i]$, then for any $k \in \mathbb{N}$ and $\epsilon > 0$,

$$\mathbb{P}\left\{ \left| \sum_{i=1}^{k} f_i(x_i) - \sum_{i=1}^{k} \pi(f_i) \right| > \epsilon \right\} \le 2 \exp\left(-\frac{C\epsilon^2}{k \sum_{k=1}^{k} \alpha_i^2} \right), \tag{12}$$

where $\pi(f_i) = \sum_{x \in \mathcal{V}} \pi_x f_i(x)$ is the expectation of f_i under the distribution defined by π .

Notice that in theorem 8, we view the spectral property of the graph as constant and absorb it into C. Below is a direct corollary of Theorem 8.

Corollary 9 There exists constants C, C' (that probably depends on G) such that the following inequality holds.

$$\forall S \subseteq [c], \mathbb{P}\left\{ \left| \frac{\sum_{y \in S} F_{y,k}}{k} - \sum_{y \in S} \pi_y \right| > \epsilon \right\} \le 2 \exp\left(-C\epsilon^2 k\right), \tag{13}$$

in other words, with probability at least 0.999, we have

$$\forall k \in [n], \forall S \subseteq [c], \left| \frac{\sum_{y \in S} F_{y,k}}{k} - \sum_{y \in S} \pi_y \right| \le \frac{C' \log(n)}{\sqrt{k}}. \tag{14}$$

Notice that we assumed n > 10c, therefore the following statement is also a direct corollary of Theorem 8.

Corollary 10 The following statement holds with probability at least 0.999: for any $x \in [c]$, $F_{x,n} \geq F_{x,\lceil n/2 \rceil} + 1$.

F.2. Proof of Theorem 4

Now we prove Theorem 4. We start by a lemma that is easy to verify.

Lemma 11 If a, b, r, s > 0 satisfies $|a - r| \le \epsilon$ and $|b - s| \le \epsilon$ and $\epsilon \le s/2$, then $\left|\frac{a}{b} - \frac{r}{s}\right| \le 2\epsilon \frac{r+s}{s^2}$.

Proof Let
$$\{u_k\}_{k=1}^n = A\mathcal{V}$$
. Let $\mathbf{z}'_x = \frac{\sum_{y \in f(x)} \pi_y \mathbf{z}_y}{\sum_{y \in f(x)} \pi_y}$. For $k \leq c$, we have

$$\|\boldsymbol{u}_{k} - \boldsymbol{z}_{k}'\| \le \sum_{j=1}^{k} a_{k,j} \|\boldsymbol{v}_{j}\| + \|\boldsymbol{z}_{k}'\| \le \sum_{j=1}^{k} \psi(\gamma + N) + N \le \frac{\sqrt{c(c+1)}(\psi\gamma + \psi N)}{\sqrt{k}}.$$
 (15)

Therefore, the error can be absorbed into the $C\psi(\gamma+N)$ terms since C is allowed to be dependent on c. Below we only consider k>c. Moreover, if $f(x_k)=\emptyset$, then it is obvious that $u_k=0=z'_{x_k}$. Therefore, in the following we only consider the case where $f(x)\neq\emptyset$.

Let
$$R_k = \sum_{y \in f(x)} F_{y,k}$$
. Let $\widetilde{\boldsymbol{z}}_x^{(k)} = \sum_{y \in f(x)} \frac{F_{y,k} \boldsymbol{z}_y}{R_k}$. We have

$$\left\| \boldsymbol{u}_{k} - \widetilde{\boldsymbol{z}}_{x}^{(k)} \right\| = \left\| \sum_{\boldsymbol{y} \in f(x)} \left(\sum_{\substack{j \in [k] \\ x_{j} = \boldsymbol{y}}} a_{k,j} \boldsymbol{v}_{j} - \frac{\boldsymbol{z}_{\boldsymbol{y}} F_{\boldsymbol{y},j}}{R_{k}} \right) \right\|$$

$$(16)$$

$$= \left\| \sum_{y \in f(x)} \sum_{\substack{j \in [k] \\ x_j = y}} \left(a_{k,j} \boldsymbol{v}_j - \frac{\boldsymbol{z}_{x_j}}{R_k} \right) \right\|$$
 (17)

$$= \left\| \sum_{y \in f(x)} \left(\sum_{\substack{j \in [k] \\ x_j = y}} \left(a_{k,j} \boldsymbol{v}_j - \boldsymbol{z}_{x_j} \right) \right) + \sum_{y \in f(x)} \boldsymbol{z}_y \sum_{\substack{j \in [k] \\ x_j = k}} \left(a_{k,j} - \frac{1}{R_k} \right) \right\|$$
(18)

$$\leq \sum_{\substack{j \in [k] \\ x_j \in f(x)}} a_{k,j} \| \boldsymbol{v}_j - \boldsymbol{z}_{x_j} \| + N \sum_{y \in f(x)} \left| \frac{F_{y,k}}{R_k} - \sum_{\substack{j \in [k] \\ x_j = y}} a_{k,j} \right|$$
(19)

$$\stackrel{\text{(i)}}{\leq} \gamma \left(\sum_{\substack{j \in [k] \\ x_j \in f(x)}} \frac{a_{k,j}}{\sqrt{j}} \right) + \frac{Nc\psi}{\sqrt{k}} \tag{20}$$

$$\leq \gamma \left(\sum_{j=1}^{k} \frac{a_{k,j}}{\sqrt{j}} \right) + \frac{Nc\psi}{\sqrt{k}}, \tag{21}$$

where in (i) we use the condition that \mathcal{V} is a good sequence converging to \mathcal{Z} and that \mathcal{A} reflects f.

Now, define $S_{k,j} = \sum_{i=1}^{j} a_{k,j}$ and $S_{k,0} = 0$. Since \mathbf{A} is a nice attention map with parameter ψ , we have $S_j \leq \frac{\psi_j}{k}$. Notice that $a_{k,j} = S_{k,j} - S_{k,j-1}$. Thus we have

$$\sum_{j=1}^{k} \frac{a_{k,j}}{\sqrt{j}} = \sum_{j=1}^{k} \frac{1}{\sqrt{j}} \left(S_{k,j} - S_{k,j-1} \right)$$
 (22)

$$= \frac{S_{k,k}}{\sqrt{k}} + \sum_{j=0}^{k-1} S_{k,j} \left(\frac{1}{\sqrt{j}} - \frac{1}{\sqrt{j+1}} \right)$$
 (23)

$$\leq \frac{\psi}{\sqrt{k}} + \frac{\psi}{k} \sum_{j=1}^{k-1} j \left(\frac{1}{\sqrt{j}} - \frac{1}{\sqrt{j+1}} \right) \tag{24}$$

$$= \frac{\psi}{\sqrt{k}} + \frac{\psi}{k} \sum_{i=1}^{k-1} \frac{1}{\sqrt{j}} - \frac{\psi(k-1)}{k\sqrt{k}}$$
 (25)

$$=\frac{\psi}{k}\sum_{j=1}^{k}\frac{1}{\sqrt{j}}.$$
(26)

Notice that,

$$\sum_{j=1}^{k} \frac{1}{\sqrt{j}} = 1 + \sum_{j=2}^{k} \int_{j-1}^{j} \frac{1}{\sqrt{j}} dx \le 1 + \sum_{j=2}^{k} \int_{j-1}^{j} \frac{1}{\sqrt{x}} dx = 1 + \int_{1}^{k} \frac{1}{\sqrt{x}} dx = 2\sqrt{k}.$$
 (27)

Subtracting eq. (27) into the argument above, we obtain

$$\left\| \boldsymbol{u}_k - \widetilde{\boldsymbol{z}}_x^{(k)} \right\| \le \frac{2\psi\gamma + Nc\psi}{\sqrt{k}}.$$
 (28)

Moreover,

$$\left\| \widetilde{\boldsymbol{z}}_{x}^{(k)} - \boldsymbol{z}_{x}' \right\| = \left\| \sum_{y \in f(x)} \left(\frac{F_{y,k}}{R} - \frac{\pi_{y}}{\sum_{y' \in f(x)} \pi_{y'}} \right) \boldsymbol{z}_{y} \right\|$$
(29)

$$\leq N \sum_{y \in [c]} \left| \frac{F_{y,k}/k}{R/k} - \frac{\pi_y}{\sum_{y' \in f(x)} \pi_{y'}} \right|.$$
 (30)

Define $a_{k,y} = F_{y,k}/k$, $b_{k,y} = R/k$, $r_y = \pi_y$ and $s_y = \sum_{y' \in f(x)} \pi_{y'}$. From Theorem 9 we have there exists a constant number C > 0 that only depends on \mathcal{G} , and an event whose probability is at least 0.999, such that for all $k \in [n]$ and $y \in [c]$ we have $\max\{|a_{y,k} - r_y|, |b_{k,y} - s_y|\} \leq \frac{C \log n}{\sqrt{k}}$ (notice that this event is only related to the random walk, and does not depend on the specific values of \mathcal{V} , \mathbf{A} , etc.).

Let $\rho = \min_{y \in [c]} \pi_y \in (0,1)$. Let $C' = 2C/\rho > C$ that also only depends on \mathcal{G} . If $k \ge \frac{4C^2(\log n)^2}{\rho^2}$, then $\frac{C\log n}{\sqrt{k}} \le \frac{\rho}{2} \le \frac{s_y}{2}$, thus from Theorem 11 we have $\left|\frac{a_{y,k}}{r_y} - \frac{b_{y,k}}{s_y}\right| \le \frac{C\log n}{\sqrt{k}} \le \frac{C'\log n}{\sqrt{k}}$. On the other hand, if $k < \frac{4C^2(\log n)^2}{\rho^2}$, we have $\frac{C'\log n}{\sqrt{k}} \ge 1 \ge \left|\frac{a_{y,k}}{r_y} - \frac{b_{y,k}}{s_y}\right|$. Therefore,

Extended Abstract Track

we conclude that, in an event at happens with probability at least 0.999, for all $k \in [n]$ and $y \in [c]$ we have $\left| \frac{F_{y,k}}{R} - \frac{\pi_y}{\sum_{y' \in f(x)} \pi_{y'}} \right| \leq \frac{C' \log n}{\sqrt{k}}$.

Combining the above arguments, we conclude that with probability at least 0.999, it holds that for all k,

$$\|\boldsymbol{u}_k - \boldsymbol{z}'_{x_k}\| \le \frac{Nc\psi + 2\psi\gamma + C'c\log n}{\sqrt{k}},$$
 (31)

which is the desired conclusion.

Appendix G. Proof of Theoretical Results w.r.t. Layer-wise Convergence

We first prove that under the specific conditions given in Appendix E, how does the latent representations evolves.

Lemma 12 Suppose $\mathcal{V} = \{\boldsymbol{v}_k\}_{k=1}^n \in (\mathbb{R}^d)^n$ is a good sequence converging to $\mathcal{Z} = \{\boldsymbol{z}_x\}_{x \in [c]}$ with parameter γ , and $\mathcal{V}' = \{\boldsymbol{v}_k'\}_{k=1}^n \in (\mathbb{R}^d)^n$ is a good sequence converging to $\mathcal{Z}' = \{\boldsymbol{z}_x'\}_{x \in [c]}$ with parameter γ' . Moreover, suppose $T, G : \mathbb{R}^d \to \mathbb{R}^d$ are Lipschitz continuous functions with Lipschitz constants L_T , L_G respectively. Then, we have $\{T\boldsymbol{v}_k + G\boldsymbol{v}_k'\}$ is a good sequence converging to $\{T\boldsymbol{z}_x + G\boldsymbol{z}_x'\}_{x \in \mathcal{V}}$ with parameter $L_T\gamma + L_G\gamma'$.

Proof Only need to notice that for any $k \in [n]$,

$$\left\| \left(T \boldsymbol{v}_{k} + G \boldsymbol{v}_{k}^{\prime} \right) - \left(T \boldsymbol{z}_{x_{k}} + G \boldsymbol{z}_{x_{k}}^{\prime} \right) \right\| \leq L_{T} \left\| \boldsymbol{v}_{k} - \boldsymbol{z}_{x_{k}} \right\| + L_{G} \left\| \boldsymbol{v}_{k}^{\prime} - \boldsymbol{z}_{x_{k}}^{\prime} \right\| \leq \frac{L_{T} \gamma + L_{G} \gamma^{\prime}}{\sqrt{k}}. \tag{32}$$

Lemma 13 There exists a scalar number C > 0 that possibly depends on the graph \mathcal{G} , and an event with probability at least 0.999, such that the following statement holds. For any layer ℓ , if $\mathcal{V} = \{\boldsymbol{v}_k\}_{k=1}^n$ is a good sequence converging to $\mathcal{Z} = \{\boldsymbol{z}_x\}_{x \in [c]}$ with parameter γ , and let $\boldsymbol{A}^{(\ell)}$ be defined as in eq. (8), then $\mathcal{U} = \boldsymbol{A}^{(\ell)}\mathcal{V}$ is a good sequence converging to $\boldsymbol{Z}' = \{\boldsymbol{z}'_x\}_{x \in [c]}$, where

$$\boldsymbol{z}_{x}' = \rho_{A}^{(\ell)} \boldsymbol{z}_{x} + \rho_{B}^{(\ell)} \sum_{y \in [c]} \frac{w_{x,y}}{d_{x}} \boldsymbol{z}_{y} + \rho_{O}^{(\ell)} \sum_{y \in [c]} \pi_{\mathcal{G}}(y) \boldsymbol{z}_{y} + \rho_{T}^{(\ell)} \boldsymbol{v}_{1}, \tag{33}$$

with parameter

$$\kappa = C(\gamma + N) \left(\sum_{\tau \in \{A, B, O\}} \rho_{\tau}^{(\ell)} \psi_{\tau}^{(\ell)} \right) + C \log n \sum_{\tau \in \{A, B, O\}} \rho_{\tau}^{(\ell)}, \tag{34}$$

where $N = \max_{y \in [c]} \|\boldsymbol{z}_y\|$.

Proof

Let $\widehat{\boldsymbol{A}}^{(\ell)} = \boldsymbol{A}^{(\ell)} - \rho_T^{(\ell)} \boldsymbol{A}^{(T)}$, and let $\widehat{\mathcal{U}} = \left(\widehat{\boldsymbol{A}}^{(\ell)} \mathcal{V}; \widehat{\mathcal{Z}}'\right)$, where $\widehat{\mathcal{Z}}' = \left\{\widehat{\boldsymbol{z}}_x'\right\}_{x \in [c]}$ is defined as

$$\widehat{\boldsymbol{z}}_{x}' = \boldsymbol{z}_{x}' - \rho_{T}^{(\ell)} \boldsymbol{v}_{1}. \tag{35}$$

For a token $x \in [c]$, let $\mathcal{N}(x)$ be the set of all neighbors of x. Notice that for any $y \in \mathcal{N}(x)$, we have

$$\frac{\pi_y}{\sum_{y' \in \mathcal{N}(x)} \pi_{y'}} = \frac{\pi_x \pi_y}{\pi_x \sum_{y' \in [c]} \widetilde{w}_{x,y'} \pi_{y'}} = \frac{w_{x,y}}{d_x}.$$
 (36)

Therefore, from Theorem 12 and Theorem 4, we have $\widehat{\mathcal{U}}$ converges to $\widehat{\mathcal{Z}}'$ with parameter κ . Since for k > c, $\mathbf{u}_k = \widehat{\mathbf{u}}_k + \rho_T^{(\ell)} \sum_{j=1}^k a_{k,j}^{(T)} \mathbf{v}_j = \widehat{\mathbf{u}}_k + \rho_T^{(\ell)} \mathbf{v}_1$, and $\mathbf{z}'_{x_k} = \widehat{\mathbf{z}}'_x + \rho_T^{(\ell)} \mathbf{v}_1$, we have

$$\|\boldsymbol{u}_k - \boldsymbol{z}_k'\| = \|\widehat{\boldsymbol{u}}_k - \widehat{\boldsymbol{z}}_k'\| \le \frac{\kappa}{\sqrt{k}},$$
 (37)

we have \mathcal{U} is also a good sequence converging to \mathcal{Z}' with parameter κ .

G.1. Evolution of the Latent Representation

From this sub-section, we focus on the evolution of the latent representation across layers and show where do they converge.

Lemma 14 Let $\rho_A, \rho_B, \rho_O, \rho_T > 0$. Let $\mathbf{D} = \text{diag}(\mathbf{W1})$ is the degree matrix. Let $\mathbf{z} \in \mathbb{R}^c$ be a vector, and let \mathbf{z}' be defined as

$$z' = \rho_A z + \rho_B D^{-1} W z + \rho_O \langle \alpha, z \rangle \mathbf{1} + \rho_T \mathbf{1}. \tag{38}$$

Then, for any U^{\top} be a projection on to a subspace orthogonal to $D^{1/2}\mathbf{1}$, we have

$$\boldsymbol{U}^{\top} \boldsymbol{D}^{1/2} \boldsymbol{z}' = \boldsymbol{U}^{\top} \boldsymbol{M} \boldsymbol{D}^{1/2} \boldsymbol{z}, \tag{39}$$

where $\mathbf{M} = \rho_A \mathbf{I} + \rho_B \mathbf{D}^{-1/2} \mathbf{W} \mathbf{D}^{-1/2}$.

Proof

Let $\boldsymbol{E} = \boldsymbol{D}^{1/2} \mathbf{1} \boldsymbol{\alpha}^{\top} \boldsymbol{D}^{-1/2}$. We have

$$z' = \left(\rho_A \mathbf{I} + \rho_B \mathbf{D}^{-1} \mathbf{W} + \rho_O \mathbf{1} \boldsymbol{\alpha}^{\top}\right) z + \rho_T \mathbf{1}$$
(40)

$$= \boldsymbol{D}^{-1/2} \left(\rho_A \boldsymbol{I} + \rho_B \boldsymbol{D}^{-1} \boldsymbol{W} + \rho_O \mathbf{1} \boldsymbol{\alpha}^{\top} \right) \boldsymbol{D}^{1/2} \boldsymbol{z} + \rho_T \mathbf{1}. \tag{41}$$

$$= D^{-1/2} (M + \rho_O E) D^{1/2} z + \rho_T 1.$$
(42)

Let $\widetilde{\boldsymbol{z}}' = \boldsymbol{D}^{1/2} \boldsymbol{z}', \ \widetilde{\boldsymbol{z}} = \boldsymbol{D}^{1/2} \boldsymbol{z}$, and $\widetilde{\boldsymbol{1}} = \boldsymbol{D}^{1/2} \boldsymbol{1}$. Thus, we have

$$\widetilde{z}' = M\widetilde{z} + \rho_O E\widetilde{z} + \rho_T \widetilde{1} \tag{43}$$

Extended Abstract Track

Notice that, since E is a rank-1 matrix, its image space is span $\tilde{\mathbf{1}}$: for any vector $\mathbf{x} \in \mathbb{R}^c$,

$$Ex = D^{1/2} \mathbf{1} \boldsymbol{\alpha}^{\top} D^{-1/2} x = \left\langle \boldsymbol{\alpha}, D^{-1/2} x \right\rangle D^{1/2} \mathbf{1} = \left\langle \boldsymbol{\alpha}, D^{-1/2} x \right\rangle \tilde{\mathbf{1}}.$$
(44)

Therefore, we have

$$\boldsymbol{U}^{\top} \widetilde{\boldsymbol{z}}' = \boldsymbol{U}^{\top} \boldsymbol{M} \widetilde{\boldsymbol{z}} + \rho_O \boldsymbol{U}^{\top} (\boldsymbol{E} \widetilde{\boldsymbol{z}}) + \rho_T \boldsymbol{U}^{\top} \widetilde{\boldsymbol{1}}$$
(45)

$$= \boldsymbol{U}^{\top} \boldsymbol{M} \widetilde{\boldsymbol{z}} + \left(\rho_O \left\langle \boldsymbol{\alpha}, \boldsymbol{D}^{-1/2} \widetilde{\boldsymbol{z}} \right\rangle + \rho_T \right) \boldsymbol{U}^{\top} \widetilde{\boldsymbol{1}}$$
 (46)

$$= \boldsymbol{U}^{\top} \boldsymbol{M} \widetilde{\boldsymbol{z}}. \tag{47}$$

Corollary 15 Let $\rho_A, \rho_B, \rho_O, \rho_T > 0$. Let $\mathcal{Z} = \{z_x\}_{x \in [c]} \in (\mathbb{R}^d)^c$ be a sequence, and define sequence $\mathcal{Z}' = \{z_x'\}_{x \in [c]}$ as follows:

$$\boldsymbol{z}_{x}' = \rho_{A}\boldsymbol{z}_{x} + \frac{\rho_{B}}{d_{x}} \sum_{y \in [c]} w_{x,y} \boldsymbol{z}_{y} + \rho_{O} \sum_{y \in [c]} \alpha_{y} \boldsymbol{z}_{y} + \rho_{T} \boldsymbol{v}_{1}. \tag{48}$$

Let $\mathbf{Z} = \text{mat } \mathcal{Z} \in \mathbb{R}^{d \times c}$ and $\mathbf{Z}' = \text{mat } \mathcal{Z}' \in \mathbb{R}^{d \times c}$. Let $\mathbf{M} = \rho_A \mathbf{I} + \rho_B \mathbf{D}^{-1/2} \mathbf{W} \mathbf{D}^{-1/2}$, and $\{\lambda_k\}_{k=1}^n$ be its eigenvalues, arranging in a non-increasing order of absolute values. Let the eigen-decomposition of \mathbf{M} be

$$\boldsymbol{M} = \begin{bmatrix} \boldsymbol{f} & \boldsymbol{X} & \boldsymbol{Y} \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \boldsymbol{\Lambda} & \\ & & \boldsymbol{\Lambda}' \end{bmatrix} \begin{bmatrix} \boldsymbol{f}^{\top} \\ \boldsymbol{X}^{\top} \\ \boldsymbol{Y}^{\top}, \end{bmatrix}$$
(49)

where $\mathbf{\Lambda} = \{\lambda_k\}_{k=2}^q$ and $\mathbf{\Lambda}' = \{\lambda_k\}_{k=q+1}^c$. Then, we have

$$\frac{\left\| \mathbf{Z}' \mathbf{D}^{1/2} \mathbf{X} \right\|}{\left\| \mathbf{Z} \mathbf{D}^{1/2} \mathbf{X} \right\|} \ge \delta_{\mathbf{M}} \frac{\left\| \mathbf{Z}' \mathbf{D}^{1/2} \mathbf{Y} \right\|}{\left\| \mathbf{Z} \mathbf{D}^{1/2} \mathbf{Y} \right\|}$$

$$(50)$$

Proof From Theorem 14, we have

$$Z'D^{1/2}X = ZD^{1/2}MX = ZD^{1/2}X\Lambda,$$
(51)

therefore

$$\left\| \mathbf{Z}' \mathbf{D}^{1/2} \mathbf{X} \right\|_{F} = \left\| \mathbf{Z} \mathbf{D}^{1/2} \mathbf{X} \mathbf{\Lambda} \right\| \ge \left\| \mathbf{Z} \mathbf{D}^{1/2} \mathbf{X} \right\|_{F} \left\| \mathbf{\Lambda} \right\| \ge \left| \lambda_{q} \right| \left\| \mathbf{Z} \mathbf{D}^{1/2} \mathbf{X} \right\|_{F}. \tag{52}$$

Similarly, we have

$$\left\| \mathbf{Z}' \mathbf{D}^{1/2} \mathbf{Y} \right\|_{F} \le |\lambda_{q+1}| \left\| \mathbf{Z} \mathbf{D}^{1/2} \mathbf{Y} \right\|_{F}. \tag{53}$$

The proposition directly follows.

Corollary 16 Under the same condition as in Theorem 15, let $\sigma : \mathbb{R}^d \to \mathbb{R}^d$ be a great mapping w.r.t. $(\gamma_1, \gamma_2, \{Z'\}, X)$ and $(\gamma'_1, \gamma'_2, \{Z'\}, Y)$. Then, we have

$$\frac{\left\|\sigma\left(\mathbf{Z}'\right)\mathbf{D}^{1/2}\mathbf{X}\right\|}{\left\|\mathbf{Z}\mathbf{D}^{1/2}\mathbf{X}\right\|} \ge \delta_{\mathbf{M}} \frac{\gamma_{1}}{\gamma_{2}'} \frac{\left\|\sigma\left(\mathbf{Z}'\right)\mathbf{D}^{1/2}\mathbf{Y}\right\|}{\left\|\mathbf{Z}\mathbf{D}^{1/2}\mathbf{Y}\right\|}.$$
(54)

Proof Only need to repeat the proof of Theorem 15 and use the definition of great mappings. Notice that

$$\left\|\sigma\left(\mathbf{Z}'\right)\mathbf{D}^{1/2}\mathbf{X}\right\|_{F} \ge \gamma_{1} \left\|\mathbf{Z}'\mathbf{D}^{1/2}\mathbf{X}\right\|_{F} \ge \gamma_{1} |\lambda_{q}| \left\|\mathbf{Z}'\mathbf{D}^{1/2}\mathbf{X}\right\|_{F}.$$
 (55)

Similarly,

$$\left\|\sigma\left(\mathbf{Z}'\right)\mathbf{D}^{1/2}\mathbf{Y}\right\|_{F} \leq \gamma_{2}' \left\|\mathbf{Z}'\mathbf{D}^{1/2}\mathbf{Y}\right\|_{F} \leq \gamma_{2}' |\lambda_{q+1}| \left\|\mathbf{Z}'\mathbf{D}^{1/2}\mathbf{Y}\right\|_{F}.$$
 (56)

The proposition directly follows.

G.2. Proof of Theorem 7

Let E_1 be an event with probability at least 0.999 defined in Theorem 13 holds. Let $E_2 = \{ \forall x \in [c], n_{[x]} > n/2 \}$, Theorem 10 shows that E_2 also holds with probability at least 0.999. In the following, we condition on the event $E_1 \cap E_2$, which holds with probability at least 0.99.

Let $\mathcal{Z}^{(\ell)} = \left\{ \boldsymbol{z}_x^{(\ell)} \right\}_{x=1}^c \in (\mathbb{R}^d)^c$ be defined as follows: $\boldsymbol{z}_x^{(1)} = \boldsymbol{b}_x$, and

$$\boldsymbol{z}_{x}^{(\ell+1)} = \rho_{A}^{(\ell)} \boldsymbol{z}_{x}^{(\ell)} + \frac{\rho_{B}^{(\ell)}}{d_{x}} \sum_{y \in [c]} w_{x,y} \boldsymbol{z}_{y}^{(\ell)} + \rho_{O} \sum_{y \in [c]} \alpha_{y} \boldsymbol{z}_{y}^{(\ell)} + \rho_{T}^{(\ell)} \boldsymbol{v}_{1}^{(\ell)}.$$
 (57)

for any $\ell \in [L-1]$, and $\mathbf{z}_x'^{(\ell)} = \sigma^{(\ell)}\left(\mathbf{z}_k^{(\ell)}\right)$. From the definition, it is obvious that $\left\{\mathbf{v}_k^{(1)}\right\}_{k=1}^n$ converges to $\mathcal{Z}^{(1)}$ with parameter 0.

• We first fix an $L \in \mathbb{N}$ and analyze the context-wise convergence. First consider an arbitrary $x \in [c]$. Using Theorem 13 with an induction, it is not hard to prove that for each $\ell \in L$ there exists a $\gamma^{(\ell)} = \text{poly} \log n$ (since we are going to take limit w.r.t. n, we view all other values independent of n as constants; notice that ℓ is a fixed index here), such that $\left\{ \boldsymbol{v}_k^{(\ell)} \right\}_{k=1}^n$ is a good sequence converging to $\left\{ \boldsymbol{z}_y'^{(\ell)} \right\}_{y \in [c]}$ with parameter $\gamma^{(\ell)}$. Therefore, we have

$$\left\| \boldsymbol{v}_{n_{[x]}}^{(\ell)} - \sigma \left(\boldsymbol{z}_{x}^{(\ell)} \right) \right\| \leq \frac{\gamma^{(\ell)}}{\sqrt{n_{[x]}}} < \frac{\sqrt{2} \operatorname{poly} \log n}{\sqrt{n}}.$$
 (58)

Therefore, taking $\ell = L$ and sending $n \to \infty$, we have

$$\lim_{n \to \infty} \left\| \boldsymbol{v}_{n_{[x]}}^{(L)} - \boldsymbol{z}_x^{\prime(L)} \right\| = 0, \tag{59}$$

Extended Abstract Track

which is equivalent to

$$\lim_{n \to \infty} \mathbf{v}_{n_{[x]}}^{(L)} = \mathbf{z}_x^{\prime(\ell)}.$$
 (60)

Notice that this holds for any $x \in [c]$. Therefore, let $\mathbf{Z}'^{(\ell)} = \max \left\{ \mathbf{z}_x'^{(\ell)} \right\}_{x \in [c]} \in \mathbb{R}^{d \times c}$, we have when $\left\| \mathbf{V}_n'^{(L)} \mathbf{D}^{1/2} \mathbf{X} \right\| > 0$, we have

$$\lim_{n\to\infty} \frac{\left\| \boldsymbol{V}_n^{(L)} \boldsymbol{D}^{1/2} \boldsymbol{Y} \right\|}{\left\| \boldsymbol{V}_n^{(L)} \boldsymbol{D}^{1/2} \boldsymbol{X} \right\|} = \frac{\left\| \boldsymbol{Z}'^{(L)} \boldsymbol{D}^{1/2} \boldsymbol{Y} \right\|}{\left\| \boldsymbol{Z}'^{(L)} \boldsymbol{D}^{1/2} \boldsymbol{X} \right\|}.$$
 (61)

• Next, we consider the layer-wise evolution. Theorem 16 shows that for any $\ell \in [L]$,

$$\frac{\left\| \mathbf{Z}^{\prime(\ell)} \mathbf{D}^{1/2} \mathbf{Y} \right\|}{\left\| \mathbf{Z}^{\prime(\ell)} \mathbf{D}^{1/2} \mathbf{X} \right\|} \leq \delta_q \left(\rho_A \mathbf{I} + \rho_B \mathbf{M} \right) \frac{\gamma_1^{(\ell)}}{\gamma_2^{\prime(\ell)}} \frac{\left\| \mathbf{Z}^{\prime(\ell-1)} \mathbf{D}^{1/2} \mathbf{Y} \right\|}{\left\| \mathbf{Z}^{\prime(\ell-1)} \mathbf{D}^{1/2} \mathbf{X} \right\|} \leq (1 - \epsilon) \frac{\left\| \mathbf{Z}^{\prime(\ell-1)} \mathbf{D}^{1/2} \mathbf{Y} \right\|}{\left\| \mathbf{Z}^{\prime(\ell-1)} \mathbf{D}^{1/2} \mathbf{X} \right\|}.$$
(62)

Theorem 16 also confirms that $\left\| \boldsymbol{Z}'^{(L)} \boldsymbol{D}^{1/2} \boldsymbol{X} \right\| > 0$. Using an easy induction, we have

$$\frac{\left\| \mathbf{Z}^{\prime(L)} \mathbf{D}^{1/2} \mathbf{Y} \right\|}{\left\| \mathbf{Z}^{\prime(L)} \mathbf{D}^{1/2} \mathbf{X} \right\|} \le (1 - \epsilon)^{L-1} \frac{\left\| \mathbf{Z}^{\prime(1)} \mathbf{D}^{1/2} \mathbf{Y} \right\|}{\left\| \mathbf{Z}^{\prime(1)} \mathbf{D}^{1/2} \mathbf{X} \right\|}.$$
 (63)

Therefore,

$$\lim_{L \to \infty} \frac{\left\| \mathbf{Z}^{\prime(L)} \mathbf{D}^{1/2} \mathbf{Y} \right\|}{\left\| \mathbf{Z}^{\prime(L)} \mathbf{D}^{1/2} \mathbf{X} \right\|} = 0.$$
 (64)

Combining above arguments, we obtain that

$$\lim_{L \to \infty} \lim_{n \to \infty} \frac{\left\| \boldsymbol{V}_n^{(L)} \boldsymbol{D}^{1/2} \boldsymbol{Y} \right\|}{\left\| \boldsymbol{V}_n^{(L)} \boldsymbol{D}^{1/2} \boldsymbol{X} \right\|} = \lim_{L \to \infty} \frac{\left\| \boldsymbol{Z}'^{(L)} \boldsymbol{D}^{1/2} \boldsymbol{Y} \right\|}{\left\| \boldsymbol{Z}'^{(L)} \boldsymbol{D}^{1/2} \boldsymbol{X} \right\|} = 0.$$
 (65)

Appendix H. The Role of FFN: What Mappings are Great Mappings

Lemma 17 Let $\sigma : \mathbb{R}^d \to \mathbb{R}^d$ be defined as $\sigma(z) = Wz$, where W is a non-singular matrix, then σ is a great mapping w.r.t. $(\gamma_{\min}, \gamma_{\max}, \mathbb{R}^d, U)$ for any orthonormal matrix U, where γ_{\min} and γ_{\max} are the smallest and largest singular values of matrix W, respectively.

The proof of Theorem 17 is obvious.

Lemma 18 Let $\sigma : \mathbb{R}^d \to \mathbb{R}^d$ be a great mapping w.r.t. $(\gamma_1, \gamma_2, \mathscr{Z}, U)$ and $\sigma' : \mathbb{R}^d \to \mathbb{R}^d$ be a great mapping w.r.t. $(\gamma'_1, \gamma'_2, \mathscr{Z}, U)$, then $\sigma_1 \circ \sigma_2$ is a great mapping w.r.t. $(\gamma_1 \gamma'_1, \gamma_2 \gamma'_2, \mathscr{Z}, U)$.

The proof of Theorem 18 is obvious.

Appendix I. A Generalized Framework That is Independent of Data Generating

In the main paper, Theorem 4 conditions on the specific data generating process used in Park et al. (2024). This is because we need to match the distribution of the attention connection to each token with their stationary distribution in the data. In this section, we show that, with a slightly stronger assumption on the attention map, a general result of the context-wise convergence can be derived.

Definition 19 If a matrix $\mathbf{A} = \{a_{k,j}\}_{k,j \in [n]} \in \mathbb{R}^{n \times n}$ is an nice attention map with parameter ψ , and there is a mapping $f : [c] \to [c]$, such that for all $k, j \in [n]$,

$$a_{k,j} > 0 \implies x_j = f(x_k),$$
 (66)

then we say A reflects f.

Notice that Theorem 19 is basically Theorem 3 but limits the function f maps each node to only one node, instead of a set of nodes as in Theorem 3 (and in this case eq. (6) automatically holds). Although this condition seems stronger, we note that following the same idea used in the main paper, that we can compose multiple attention maps into one, this still represents a large family of allowed attention maps.

Next, we prove a similar result as Theorem 4 under Theorem 19 that is independent of input distribution.

Theorem 20 Suppose $\mathcal{V} = \{\boldsymbol{v}_k\}_{k=1}^n \in (\mathbb{R}^d)^n$ is a good sequence converging to $\mathcal{Z} = \{\boldsymbol{z}_k\}_{k \in [c]}$ with parameter γ , and $\boldsymbol{A} = \{a_{k,j}\}_{k,j \in [n]} \in \mathbb{R}^{n \times n}$ is a nice attention map with parameter ψ that reflects a function $f: [c] \to [c]$. Then $\boldsymbol{A}\mathcal{V}$ is a good sequence converging to

$$\mathcal{Z}' = \left\{ z_{f(x)} \right\}_{x \in [c]} \tag{67}$$

with parameter $2\gamma\psi + c(\gamma + 2N)$, where $N = \max_{x \in [c]} \|z_x\|$.

Proof Let $AV = \{u_k\}_{k=1}^n$.

If $k \leq c$, we have

$$\|\boldsymbol{u}_k - \boldsymbol{z}_{f(x_k)}\| \le \|\boldsymbol{u}_k - z_{x_k}\| + \|\boldsymbol{z}_{f(x_k)}\| + \|\boldsymbol{z}_{x_k}\| \le \gamma + 2N \le \frac{c(\gamma + 2N)}{\sqrt{k}}.$$
 (68)

Below, we only need to consider k > c.

Extended Abstract Track

We have

$$\|\boldsymbol{u}_{k} - \boldsymbol{z}_{f(x_{k})}\| = \left\| \sum_{\substack{j \in [k] \\ x_{j} = f(x_{k})}} a_{k,j} \left(\boldsymbol{v}_{j} - \boldsymbol{z}_{f(x_{k})}\right) \right\|$$

$$(69)$$

$$\leq \sum_{\substack{j \in [k] \\ x_j = f(x_k)}} a_{k,j} \left\| \boldsymbol{v}_j - \boldsymbol{z}_{x_j} \right\| \tag{70}$$

$$\leq \sum_{\substack{j \in [k] \\ x_j = f(x_k)}} a_{k,j} \| \boldsymbol{v}_j - \boldsymbol{z}_{x_j} \|$$

$$\leq \sum_{\substack{j \in [k] \\ x_j = f(x_k)}} \frac{a_{k,j} \gamma}{\sqrt{j}}$$

$$(70)$$

$$\leq \gamma \sum_{j=1}^{k} \frac{a_{k,j}}{\sqrt{j}}.$$
(72)

Now, define $S_{k,j} = \sum_{i=1}^{j} a_{k,j}$ and $S_{k,0} = 0$. Since \boldsymbol{A} is nice, we have $S_j \leq \frac{\psi j}{k}$. Notice that $a_{k,j} = S_{k,j} - S_{k,j-1}$. Thus we have

$$\frac{1}{\gamma} \left\| \boldsymbol{u}_k - \boldsymbol{z}_{f(x_k)} \right\| \le \sum_{j=1}^k \frac{a_{k,j}}{\sqrt{j}} \tag{73}$$

$$= \sum_{j=1}^{k} \frac{1}{\sqrt{j}} \left(S_{k,j} - S_{k,j-1} \right) \tag{74}$$

$$= \frac{S_{k,k}}{\sqrt{k}} + \sum_{j=0}^{k-1} S_{k,j} \left(\frac{1}{\sqrt{j}} - \frac{1}{\sqrt{j+1}} \right)$$
 (75)

$$\leq \frac{\psi}{\sqrt{k}} + \frac{\psi}{k} \sum_{i=1}^{k-1} j \left(\frac{1}{\sqrt{j}} - \frac{1}{\sqrt{j+1}} \right) \tag{76}$$

$$= \frac{\psi}{\sqrt{k}} + \frac{\psi}{k} \sum_{i=1}^{k-1} \frac{1}{\sqrt{j}} - \frac{\psi(k-1)}{k\sqrt{k}}$$
 (77)

$$=\frac{\psi}{k}\sum_{j=1}^{k}\frac{1}{\sqrt{j}}\tag{78}$$

$$\leq \frac{2\psi}{\sqrt{k}}.\tag{79}$$

Thus we conclude that $\|\boldsymbol{u}_k - \boldsymbol{z}_{f(x_k)}\| \leq \frac{2\gamma\psi}{\sqrt{k}}$ for any $k \in [n]$, which means \boldsymbol{U} converges to $\{z_{f(x_k)}\}_{k=1}^n$ with parameter $2\gamma\psi$.

Using, Theorem 12 which is also independent of DGP, we can prove the generalized results for a large family of attention maps by combining multiple attention maps that satisfies Theorem 19.

Appendix J. Integrating Other Transformer Components

In this section, we discuss how our theoretical results can extend to more complex and realistic Transformer architectures beyond the simplified model described in Algorithm 1. We first note that in our framework, any neuron-wise transformation (i.e., operations that apply independently to the representation of each token) can be absorbed into the definition of the great mapping σ . This includes FFNs as well as normalizations such as LayerNorm or RMSNorm. Therefore, here we focus here two architectural components not yet discussed: residual connections and multi-head attention.

Residual Connection The most critical step in the proof of Theorem 7 is Theorem 13, which establishes that applying a nice attention map to a good sequence results in another good sequence, and that the attention operates implicitly on the latent representations to which the sequence converges. Introducing residual connections here is straightforward: with residual connection, eq. (33) would become

$$\boldsymbol{z}_{x}' = \left(1 + \rho_{A}^{(\ell)}\right) \boldsymbol{z}_{x} + \rho_{B}^{(\ell)} \sum_{y \in [c]} \frac{w_{x,y}}{d_{x}} \boldsymbol{z}_{y} + \rho_{O}^{(\ell)} \sum_{y \in [c]} \pi_{\mathcal{G}}(y) \boldsymbol{z}_{y} + \rho_{T}^{(\ell)} \boldsymbol{v}_{1}, \tag{80}$$

which is simply adding 1 to the $\rho_A^{(\ell)}$ coefficient. This modification only affects the δ_q term in Theorem 7, which becomes

$$\delta_q \left[(\rho_A + 1) \mathbf{I} + \rho_B \mathbf{M} \right], \tag{81}$$

which makes the spectral gap smaller. This can slower the convergence, but will not prevent it as long as the spectral gap of M is large enough.

Multi-head Attention Theorem 12 shows that any Lipschitz combination of good sequences remains a good sequence. Since multi-head attention can be viewed as a Lipschitz combination of multiple single-head attentions, it follows that a multi-head attention mechanism also satisfies Theorem 13, as long as each individual head does. All other parts of the theoretical framework extend accordingly. Notice that the coefficients in Theorem 13 may differ by constant factors in the multi-head case, but this does not affect the asymptotic conclusions in Theorem 7, which only concern limiting behavior.

Appendix K. Empirical Verification

As discussed in the main text, Theorem 7 relies on a relatively strong structural assumption about the attention maps. It is therefore essential to verify whether these assumptions hold in practice. In this section, we empirically examine this question.

Specifically, in Figure 8, we compute the proportion of A,B,T type attention connections defined in Appendix E. Note that Type O connections (i.e., connecting to arbitrary tokens) are excluded from this analysis, as they cover all positions. For each head, we compute the fraction of attention weights (across all layers) that fall into types A, B, or T. Overlapping cases (e.g., the connection from the second token in the sequence to the first one can be considered both as B and as T type) are counted only once. The figure shows that a large proportion of attention weights (>72% in total) indeed falls into these structured types, lending empirical support to our theoretical assumptions.

Extended Abstract Track 1.0 0.8 0.850 0.819 0.809 0.772 0.799 0.746 0.774 0.795 0.794 0.795 0.795 0.796 0.795 0.666 0.894 0.675 0.676 0.677 0.6784 0.677 0.6784 0.6784 0.6784 0.688 0.6882 0

Figure 8: **Proportion of structured attention connections.** For each attention head, we sum attention weights that falls into type A,B and T across all layers, and divide them by total attention weights (which is equal to the number of tokens per layer, since the attention weights are normalized). The dotted horizontal line is the average proportion over all heads.