# Understanding Mamba in In-Context Learning with Outliers: A Theoretical Generalization Analysis

Hongkang Li Rensselaer Polytechnic Institute, USA

**Songtao Lu** *The Chinese University of Hong Kong, Hong Kong* 

Xiaodong Cui IBM Research, USA

Pin-Yu Chen IBM Research, USA

Meng Wang Rensselaer Polytechnic Institute, USA LOHEK330@GMAIL.COM STLU@CSE.CUHK.EDU.HK

CUIX@US.IBM.COM

PINYU.CHEN@IBM.COM

WANGM7@RPI.EDU

### Abstract

The Mamba model has gained significant attention for its computational advantages over Transformerbased models, while achieving comparable performance across a wide range of language tasks. Like Transformers, Mamba exhibits in-context learning (ICL) capabilities, i.e., making predictions for new tasks based on a prompt containing input-label pairs and a query, without requiring fine-tuning. Despite its empirical success, the theoretical understanding of Mamba remains limited, largely due to the nonlinearity introduced by its gating mechanism. To the best of our knowledge, this paper presents the first theoretical analysis of the training dynamics of a one-layer Mamba model, which consists of a linear attention component followed by a nonlinear gating layer, and its ICL generalization on unseen binary classification tasks, even when the prompt includes additive outliers. Our analysis shows that Mamba leverages the linear attention layer to select informative context examples and uses the nonlinear gating layer to suppress the influence of outliers. By establishing and comparing to the analysis of linear Transformers under the same setting, we show that although Mamba may require more training iterations to converge, it maintains accurate predictions even when the proportion of outliers exceeds the threshold that a linear Transformer can tolerate. These theoretical findings are supported by empirical experiments.

# 1. Introduction

Transformer-based large language models (LLMs) [10, 29, 65, 72] have demonstrated remarkable capabilities across a wide range of language, vision, and reasoning tasks. However, they face efficiency challenges when processing long sequences due to the quadratic time and memory complexity of the self-attention mechanism with respect to sequence length [20, 26]. To address this, many efficient alternative architectures have been proposed, including state space models (SSMs) such as S4 [27, 28] and H3 [21]. Among them, Mamba [26] has attracted significant attention for its strong empirical performance, linear computational complexity, and hardware-friendly properties that enable efficient parallelization. These advantages have sparked growing interest in understanding the mechanism of Mamba and whether it can match or surpass the capabilities of Transformer models.

One particularly intriguing property of LLMs is *in-context learning (ICL)* [10, 23], which allows a pre-trained model to generalize to new tasks without any parameter updates. By simply augmenting the input with a prompt containing a few labeled examples from the new task, the model can produce accurate predictions for unseen tasks. While LLMs have demonstrated impressive ICL generalization, their performance is sensitive to the quality of the context examples [58, 78]. In particular, ICL performance can degrade significantly in the presence of outliers or adversarial attacks on prompts, such as data poisoning, resulting in incorrect predictions [6, 34, 42, 67, 76, 85].

Recent empirical work [7, 24, 30, 39, 66, 75] has demonstrated that Mamba can also perform ICL on function learning and natural language processing tasks. [24, 66] show that Mamba is competitive with Transformers of similar size in some ICL tasks and outperforms them in settings with many outliers, such as regression with corrupted examples. On the other hand, studies such as [7, 39, 66, 75] identify limitations of Mamba in retrieval-based and long-context reasoning tasks. Despite these empirical insights, several fundamental questions remain open:

# Why and how can a Mamba model be trained to perform in-context generalization to new tasks? How robust is it to outliers? Under what conditions can Mamba outperform Transformers for ICL?

[56] and [57] analyze Mamba-like architectures, e.g., H3 and gated linear attention, and show that the global minima of the loss landscapes correspond to models whose outputs, when given a prompt, are equivalent to those of a model performing a weighted preconditioned gradient descent using the context examples. This serves as the counterpart to the preconditioned gradient descent interpretation of ICL in Transformers [1]. [41] shows that continuous SSMs can learn dynamic systems in context. [9] proves that Mamba is expressive enough to represent optimal Laplacian smoothing. However, these studies do not address whether practical training methods can reliably yield Mamba models with ICL capabilities, nor do they provide theoretical guarantees for generalization or robustness in the presence of outliers.

Theoretical Works	Standard Mamba	Mechanism Analysis	Convergence Analysis	Generalization Guarantee	Outliers in Context
[56]		$\checkmark$			
[57]		$\checkmark$			
[41]		$\checkmark$			
[9]	$\checkmark$				
Ours	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$

Table 1: Comparison with existing works on theoretical analysis of Mamba-like models.

### 1.1. Major Contributions

This paper presents the first theoretical analysis of the training dynamics of Mamba models and their resulting ICL performance, including scenarios where context examples in the prompt contain outliers. We focus on training Mamba on binary classification tasks in which input data consist of both relevant patterns, which determine the label, and irrelevant patterns, which do not. Additionally, context inputs may include additive outliers that perturb the labels. While our analysis is based on one-layer Mamba architectures, this setting aligns with the scope of state-of-the-art theoretical studies on the training dynamics and generalization of Transformers and other neural networks, which also typically focus on one-hidden-layer models [50, 56, 57, 82]. Our main contributions are as follows:

1. Quantitative analysis of ICL emergence and robustness to outliers in Mamba. We characterize the number of context examples and training iterations required for a Mamba model to acquire ICL capabilities for new tasks that were not present during training. We prove that when trained with prompts that may contain a finite number of outlier patterns, Mamba can generalize in-context on new tasks when the context examples contain unseen outliers that are linear combinations of the training-time outliers. Furthermore, Mamba can maintain accurate ICL generalization even when the fraction of outlier-containing context examples approaches 1, demonstrating strong robustness.

2. Theoretical comparison between Mamba and linear Transformers. We provide a theoretical characterization of the convergence and generalization properties of linear Transformers trained on the same tasks. While linear Transformers may converge faster with smaller batch sizes, they can only in-context generalize effectively when the fraction of outlier-containing context examples is less than 1/2, much less than that for Mamba. Moreover, linear Transformers require significantly more context examples than Mamba to achieve comparable generalization performance. This highlights Mamba's superior robustness to a high density of outliers in ICL.

3. Theoretical characterization of the mechanism by which Mamba implements ICL. We show that the equivalent linear attention mechanism in Mamba selects context examples that share the same relevant pattern as the query, while the nonlinear gating mechanism suppresses corrupted examples and applies an exponential decay in importance based on index distance, emphasizing examples closer to the query. Together, these mechanisms enable Mamba to suppress irrelevant or corrupted context examples and focus on informative and nearby ones, achieving effective and robust ICL.

#### 1.2. Related Works

**Theoretical Analysis of ICL.** Existing theoretical works of ICL primarily focus on Transformerbased models. [1, 2, 8, 23, 74] illustrate that Transformers can implement many machine learning algorithms, such as gradient-based methods, via ICL. [14, 36, 50, 77, 82] provably investigate the training dynamics and generalization of ICL on single/multi-head Transformers. [44, 63, 81] extend the analysis to learning complicated nonlinear functions by ICL.

**Connections Between Mamba and Transformers.** [3] finds that Mamba exhibits explainability metrics comparable to those of Transformers. [20] shows that SSMs and variants of attention mechanisms share a large intersection and can be viewed as duals of each other. [34] notes a similarity between the forget gate in Mamba and the positional encodings in Transformers. The complementary strengths, Mamba's computational efficiency and Transformers' ability to capture global dependencies, have motivated the development of hybrid architectures [32, 45, 79].

**Optimization and Generalization of the Attention Architecture.** Some other works focus on the optimization and generalization of attention-based models without nonlinear gating beyond the ICL setting. [38, 40, 48, 49, 55, 80] study the generalization of one-layer Transformers in classification tasks by formulating spatial association, key features, or the semantic structure of the input. [35, 62, 69] investigate the problem in next-token prediction based on the partial order, bigram, or semantic association assumption. [14, 33] extend the analysis to multi-head attention networks.

### 2. Problem Formulation

The learning model, Mamba, is proposed in [25]. Given the input  $U = (u_1, \dots, u_m) \in \mathbb{R}^{d_0 \times m}$ , the model outputs  $o_i$  recursively through the hidden states  $h_i$ ,  $i \in [m]$ . Starting from  $h_0 = U$ , a

one-layer Mamba can be formulated as

$$\boldsymbol{h}_{i} = \boldsymbol{h}_{i-1} \odot \tilde{\boldsymbol{A}}_{i} + (\boldsymbol{u}_{i} \boldsymbol{1}_{m}^{\top}) \odot \tilde{\boldsymbol{B}}_{i} \quad \in \mathbb{R}^{d_{0} \times m}, \quad \forall i \in [m]$$

$$\boldsymbol{o}_{i} = \boldsymbol{h}_{i} \boldsymbol{C}_{i} \quad \in \mathbb{R}^{d+1},$$

$$(1)$$

where  $\mathbf{1}_m$  is an all-ones vector in  $\mathbb{R}^m$ ,  $\tilde{\mathbf{B}}_i = \mathbf{1}_{d_0}(\Delta_i \mathbf{B}_i)(\exp(\Delta_i \mathbf{A}) - \mathbf{I}_m)(\Delta_i \mathbf{A})^{-1} \in \mathbb{R}^{d_0 \times m}$ ,  $\mathbf{B}_i = \mathbf{u}_i^\top \mathbf{W}_B^\top \in \mathbb{R}^{1 \times m}$  with  $\mathbf{W}_B \in \mathbb{R}^{m \times d_0}$ ,  $\tilde{\mathbf{A}}_i = \mathbf{1}_{d_0} \operatorname{diag}(\exp(\Delta_i \mathbf{A}))^\top \in \mathbb{R}^{d_0 \times m}$ ,  $\mathbf{C}_i = \mathbf{W}_C \mathbf{u}_i \in \mathbb{R}^m$  with  $\mathbf{W}_C \in \mathbb{R}^{m \times d_0}$ .  $\odot$  and  $\exp(\cdot)$  are element-wise product and exponential operations, respectively.  $\operatorname{diag}(\cdot) : \mathbb{R}^{d_0 \times d_0} \to \mathbb{R}^{d_0}$  outputs the diagonal of the input as a vector.  $\sigma(\cdot) : z \in \mathbb{R} \mapsto (1 + \exp(-z))^{-1} \in \mathbb{R}$  is the sigmoid function. According to [25, 31], we select  $\mathbf{A} = -\mathbf{I}_m \in \mathbb{R}^{m \times m}$ ,  $\Delta_i = \operatorname{softplus}(\mathbf{w}^\top \mathbf{u}_i) = \log(1 + \exp(\mathbf{w}^\top \mathbf{u}_i)) \in \mathbb{R}$  with  $\mathbf{w} \in \mathbb{R}^{d_0}$ . for simplicity of analysis.

Following the theoretical setup used in recent in-context learning (ICL) analyses [23, 36, 50, 56, 57], we consider training a model on prompts from a subset of tasks to endow it with ICL capabilities on unseen tasks. This framework is motivated by the observation [15] that although LLMs are typically trained without supervised labels, natural text often contains implicit inputoutput pairs, i.e., phrases following similar templates, that resemble the prompt-query format used in our setup. Specifically, we consider a set of binary classification tasks  $\mathcal{T}$ , where for a certain task  $f \in \mathcal{T}$ , the label  $z \in \{+1, -1\}$  of a given input query  $\mathbf{x}_{query} \in \mathbb{R}^d$  is determined by  $z = f(\mathbf{x}_{query}) \in \{+1, -1\}$ . Then, the prompt  $\mathbf{P}$  for  $\mathbf{x}_{query}$  is constructed as

$$P = \begin{pmatrix} \boldsymbol{x}_1 & \boldsymbol{x}_2 & \cdots & \boldsymbol{x}_l & \boldsymbol{x}_{query} \\ \boldsymbol{y}_1 & \boldsymbol{y}_2 & \cdots & \boldsymbol{y}_l & 0 \end{pmatrix}$$
  
$$:= (\boldsymbol{p}_1, \boldsymbol{p}_2, \cdots, \boldsymbol{p}_{query}) \in \mathbb{R}^{(d+1) \times (l+1)}, \qquad (2)$$

where  $y_i = f(x_i)$ ,  $i \in [l]$ . With the prompt P in (2) as the input to the Mamba model in (1) with m = l + 1 and  $d_0 = d + 1$ , the output of one-layer Mamba can be rewritten as

$$F(\Psi; \boldsymbol{P}) = \boldsymbol{e}_{d+1}^{\top} \boldsymbol{o}_{l+1} = \sum_{i=1}^{l+1} G_{i,l+1}(\boldsymbol{w}) y_i \boldsymbol{p}_i^{\top} \boldsymbol{W}_B^{\top} \boldsymbol{W}_C \boldsymbol{p}_{query},$$
where  $G_{i,l+1}(\boldsymbol{w}) = \begin{cases} \Pi_{j=i+1}^{l+1} (1 - \sigma(\boldsymbol{w}^{\top} \boldsymbol{p}_j)) \sigma(\boldsymbol{w}^{\top} \boldsymbol{p}_i), & i < l+1, \\ \sigma(\boldsymbol{w}^{\top} \boldsymbol{p}_{query}), & i = l+1, \end{cases}$ 
(3)

where  $e_{d+1} = (0, \dots, 0, 1)^{\top} \in \mathbb{R}^{d+1}$  and  $\Psi = \{W_B, W_C, w\}$  is the set of trainable parameters. The derivation of (3) can be found in Appendix F.1. From (3), one can observe that a one-layer Mamba is equivalent to a **linear attention** layer parameterized by  $W_B$  and  $W_C$  followed by a **nonlinear gating** layer  $G_{i,l+1}(w)$  for  $i \in [l+1]$ . Specifically,  $W_B$  and  $W_C$  can be respectively interpreted as the key and query parameters in a Transformer model. Therefore, a Transformer with linear attention, commonly studied in the context of ICL [82], can be viewed as a special case of the formulation in (3) by removing the nonlinear gating, i.e., setting  $G_{i,l+1}(w) = 1$  for all  $i \in [l+1]$ . We adopt this simplified formulation when comparing Mamba and Transformers in Section A.3.

Given N training examples consisting of prompt-label pairs  $(\mathbf{P}^n, z^n)_{n=1}^N$ , the model is trained by solving the empirical risk minimization problem using the hinge loss:

$$\min_{\Psi} R_N(\Psi) := \frac{1}{N} \sum_{n=1}^N \ell(\Psi; \boldsymbol{P}^n, z^n) = \frac{1}{N} \sum_{n=1}^N \max\{0, 1 - z^n \cdot F(\Psi; \boldsymbol{P}^n)\}.$$
 (4)

Each prompt  $P^n$  is generated from a distribution D, where the query  $x_{query}^n$  and all context inputs  $x_i^n$  are sampled independently, and the associated task  $f^n$  is drawn from a set of training tasks  $\mathcal{T}_{tr} \subset \mathcal{T}$ .

**Training Algorithm**: The model is trained using stochastic gradient descent (SGD) with step size  $\eta$  with batch size B, summarized in Algorithm 1.  $W_B^{(0)}$  and  $W_C^{(0)}$  are initialized such that the first d diagonal entries of  $W_B^{(0)}$  and  $W_C^{(0)}$  are set as  $\delta \in (0, 0.2]$ .  $w^{(0)}$  follows Gaussian  $\mathcal{N}(0, \mathbf{I}_{d+1}/(d+1))$ .

ICL Generalization in the Presence of Outliers: The testing prompt P' follows an unknown distribution  $\mathcal{D}'$ , which is different from the training prompt P and may contain outliers. Then, the ICL generalization of the model  $\Psi$  is computed as the classification error across all tasks in  $\mathcal{T}$ , including those never appear during the training stage, i.e.,

$$L_{f\in\mathcal{T},\mathbf{P}'\sim\mathcal{D}'}^{0-1}(\Psi;\mathbf{P}',z) = \mathbb{E}_{f\in\mathcal{T},\mathbf{P}'\sim\mathcal{D}'} \big[ \mathbb{1}[z\cdot F(\Psi;\mathbf{P}')<0] \big].$$
(5)

#### 3. Main Theoretical Insights

We formulate a class of binary classification tasks where the labels in each task are determined by two selected relevant patterns. The model is trained on a subset of these tasks using prompts that may include context examples corrupted by additive outliers. We then evaluate the model's performance on unseen tasks, where the prompts can contain outliers not observed during training.

P1. Theoretical Characterization of Learning Dynamics, ICL Generalization, and Robustness to Outliers in Mamba Models. We provide quantitative guarantees that training with prompts can lead to favorable ICL generalization on unseen tasks, and these results hold even in the presence of outliers (Theorems 1 and 2). Specifically, if a fraction  $p_a \in [0, 1)$  of the context examples in the training prompts contain additive outliers, we prove the learned model still generalizes accurately at test time, as long as the fraction of outliers in the testing prompt, denoted by  $\alpha$ , is less than  $\min\{1, p_a \cdot l_{tr}/l_{ts}\}$  where  $l_{tr}$  and  $l_{ts}$  are the number of examples in the training and testing prompts, respectively. Notably, the outliers in the test prompt may be previously unseen and can be formed as almost arbitrary positive linear combinations of a finite set of outlier patterns seen during training.

P2. A Comparison Between Mamba and Linear Transformer Models. We theoretically analyze the convergence and ICL generalization of a one-layer linear Transformer (Theorems 3 and 4) for comparison. Our results show that linear Transformers require smaller batch sizes, fewer iterations, and can tolerate larger-magnitude outliers for successful training convergence compared to Mamba. However, linear Transformers can only generalize well when the test prompt has an outlier fraction  $\alpha < 1/2$ , whereas Mamba could maintain accurate generalization even if  $\alpha$  goes to 1. Moreover, even when both models can achieve ICL, e.g., when  $\alpha$  is close to 1/2, linear Transformers require significantly more context examples to achieve comparable performance. Thus, despite requiring more effort during training, Mamba models demonstrate superior robustness to outliers during ICL.

**P3. Mechanism of Mamba Models in Implementing ICL.** Our analysis shows that the linear attention layer in Mamba selectively emphasizes context examples that share the same relevant pattern as the query, while the nonlinear gating layer promotes examples that are both close to the query and free of additive outliers. This dual mechanism enables the trained Mamba to suppress irrelevant or corrupted context examples and focus on informative examples close to the query, thus achieving successful and robust ICL.

The details of main theoretical results and experiments can be found in Appendices A and B.

# 4. Conclusion, Limitations, and Future Works

This paper theoretically studies the learning dynamics, ICL generalization, and the robustness to outliers of Mamba models, together with a characterization of how different components of Mamba contribute to the ICL mechanism. Our analysis also provides a theoretical comparison between Mamba and linear Transformer models. Although based on a one-layer Mamba structure on binary classification tasks, this work provides a deeper theoretical understanding and provable advantages of Mamba. Future directions include designing general Mamba-based language/multi-modal models.

# References

- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. *arXiv preprint arXiv:2306.00297*, 2023.
- [2] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. In *The Eleventh International Conference on Learning Representations*, 2023.
- [3] Ameen Ali, Itamar Zimerman, and Lior Wolf. The hidden attention of mamba models. *arXiv* preprint arXiv:2403.01590, 2024.
- [4] Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. In Advances in neural information processing systems, pages 6155–6166, 2019.
- [5] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In *International Conference on Machine Learning*, pages 242–252. PMLR, 2019.
- [6] Usman Anwar, Johannes Von Oswald, Louis Kirsch, David Krueger, and Spencer Frei. Adversarial robustness of in-context learning in transformers for linear regression. arXiv preprint arXiv:2411.05189, 2024.
- [7] Simran Arora, Sabri Eyuboglu, Michael Zhang, Aman Timalsina, Silas Alberti, James Zou, Atri Rudra, and Christopher Re. Simple linear attention language models balance the recallthroughput tradeoff. In *International Conference on Machine Learning*, pages 1763–1840. PMLR, 2024.
- [8] Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. *arXiv preprint arXiv:2306.04637*, 2023.
- [9] Marco Bondaschi, Nived Rajaraman, Xiuying Wei, Kannan Ramchandran, Razvan Pascanu, Caglar Gulcehre, Michael Gastpar, and Ashok Vardhan Makkuva. From markov to laplace: How mamba in-context learns markov chains. arXiv preprint arXiv:2502.10178, 2025.
- [10] 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.

- [11] Alon Brutzkus and Amir Globerson. An optimization and generalization analysis for maxpooling networks. In *Uncertainty in Artificial Intelligence*, pages 1650–1660. PMLR, 2021.
- [12] Yuan Cao and Quanquan Gu. Generalization bounds of stochastic gradient descent for wide and deep neural networks. In Advances in Neural Information Processing Systems, pages 10836–10846, 2019.
- [13] Stephanie Chan, Adam Santoro, Andrew Lampinen, Jane Wang, Aaditya Singh, Pierre Richemond, James McClelland, and Felix Hill. Data distributional properties drive emergent in-context learning in transformers. *Advances in neural information processing systems*, 35:18878–18891, 2022.
- [14] Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Training dynamics of multi-head softmax attention for in-context learning: Emergence, convergence, and optimality. arXiv preprint arXiv:2402.19442, 2024.
- [15] Yanda Chen, Chen Zhao, Zhou Yu, Kathleen McKeown, and He He. Parallel structures in pre-training data yield in-context learning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8582–8592, 2024.
- [16] Zixiang Chen, Yuan Cao, Quanquan Gu, and Tong Zhang. A generalized neural tangent kernel analysis for two-layer neural networks. *Advances in Neural Information Processing Systems*, 33, 2020.
- [17] Mohammed Nowaz Rabbani Chowdhury, Shuai Zhang, Meng Wang, Sijia Liu, and Pin-Yu Chen. Patch-level routing in mixture-of-experts is provably sample-efficient for convolutional neural networks. In *International Conference on Machine Learning*, 2023.
- [18] Mohammed Nowaz Rabbani Chowdhury, Meng Wang, Kaoutar El Maghraoui, Naigang Wang, Pin-Yu Chen, and Christopher Carothers. A provably effective method for pruning experts in fine-tuned sparse mixture-of-experts. arXiv preprint arXiv:2405.16646, 2024.
- [19] Amit Daniely and Eran Malach. Learning parities with neural networks. *Advances in Neural Information Processing Systems*, 33:20356–20365, 2020.
- [20] Tri Dao and Albert Gu. Transformers are ssms: generalized models and efficient algorithms through structured state space duality. In *Proceedings of the 41st International Conference on Machine Learning*, pages 10041–10071, 2024.
- [21] Daniel Y Fu, Tri Dao, Khaled Kamal Saab, Armin W Thomas, Atri Rudra, and Christopher Re. Hungry hungry hippos: Towards language modeling with state space models. In *The Eleventh International Conference on Learning Representations*, 2023.
- [22] Haoyu Fu, Yuejie Chi, and Yingbin Liang. Guaranteed recovery of one-hidden-layer neural networks via cross entropy. *IEEE Transactions on Signal Processing*, 68:3225–3235, 2020.
- [23] 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.

- [24] Riccardo Grazzi, Julien Niklas Siems, Simon Schrodi, Thomas Brox, and Frank Hutter. Is mamba capable of in-context learning? In *International Conference on Automated Machine Learning*, pages 1–1. PMLR, 2024.
- [25] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- [26] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. In First Conference on Language Modeling, 2024.
- [27] Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, and Christopher Ré. Combining recurrent, convolutional, and continuous-time models with linear state space layers. *Advances in neural information processing systems*, 34:572–585, 2021.
- [28] Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured state spaces. In *International Conference on Learning Representations*, 2022.
- [29] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. arXiv preprint arXiv:2501.12948, 2025.
- [30] John T Halloran, Manbir Gulati, and Paul F Roysdon. Mamba state-space models can be strong downstream learners. arXiv e-prints, pages arXiv–2406, 2024.
- [31] Dongchen Han, Ziyi Wang, Zhuofan Xia, Yizeng Han, Yifan Pu, Chunjiang Ge, Jun Song, Shiji Song, Bo Zheng, and Gao Huang. Demystify mamba in vision: A linear attention perspective. *arXiv preprint arXiv:2405.16605*, 2024.
- [32] Ali Hatamizadeh and Jan Kautz. Mambavision: A hybrid mamba-transformer vision backbone. *arXiv preprint arXiv:2407.08083*, 2024.
- [33] Jianliang He, Xintian Pan, Siyu Chen, and Zhuoran Yang. In-context linear regression demystified: Training dynamics and mechanistic interpretability of multi-head softmax attention. arXiv preprint arXiv:2503.12734, 2025.
- [34] Pengfei He, Han Xu, Yue Xing, Hui Liu, Makoto Yamada, and Jiliang Tang. Data poisoning for in-context learning. arXiv preprint arXiv:2402.02160, 2024.
- [35] Ruiquan Huang, Yingbin Liang, and Jing Yang. Non-asymptotic convergence of training transformers for next-token prediction. Advances in Neural Information Processing Systems, 37:80634–80673, 2024.
- [36] Yu Huang, Yuan Cheng, and Yingbin Liang. In-context convergence of transformers. In NeurIPS 2023 Workshop on Mathematics of Modern Machine Learning, 2023.
- [37] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In Advances in neural information processing systems, pages 8571–8580, 2018.
- [38] Samy Jelassi, Michael Sander, and Yuanzhi Li. Vision transformers provably learn spatial structure. Advances in Neural Information Processing Systems, 35:37822–37836, 2022.

- [39] Samy Jelassi, David Brandfonbrener, Sham M Kakade, et al. Repeat after me: Transformers are better than state space models at copying. In *Forty-first International Conference on Machine Learning*, 2024.
- [40] Jiarui Jiang, Wei Huang, Miao Zhang, Taiji Suzuki, and Liqiang Nie. Unveil benign overfitting for transformer in vision: Training dynamics, convergence, and generalization. Advances in Neural Information Processing Systems, 37:135464–135625, 2024.
- [41] Federico Arangath Joseph, Kilian Konstantin Haefeli, Noah Liniger, and Caglar Gulcehre. Hippo-prophecy: State-space models can provably learn dynamical systems in context. arXiv preprint arXiv:2407.09375, 2024.
- [42] Nikhil Kandpal, Matthew Jagielski, Florian Tramèr, and Nicholas Carlini. Backdoor attacks for in-context learning with language models. In *The Second Workshop on New Frontiers in Adversarial Machine Learning*, 2023.
- [43] Stefani Karp, Ezra Winston, Yuanzhi Li, and Aarti Singh. Local signal adaptivity: Provable feature learning in neural networks beyond kernels. *Advances in Neural Information Processing Systems*, 34:24883–24897, 2021.
- [44] Juno Kim and Taiji Suzuki. Transformers learn nonlinear features in context: Nonconvex mean-field dynamics on the attention landscape. In *International Conference on Machine Learning*, pages 24527–24561. PMLR, 2024.
- [45] Barak Lenz, Opher Lieber, Alan Arazi, Amir Bergman, Avshalom Manevich, Barak Peleg, Ben Aviram, Chen Almagor, Clara Fridman, Dan Padnos, et al. Jamba: Hybrid transformer-mamba language models. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [46] Hongkang Li, Meng Wang, Sijia Liu, Pin-Yu Chen, and Jinjun Xiong. Generalization guarantee of training graph convolutional networks with graph topology sampling. In *International Conference on Machine Learning*, pages 13014–13051. PMLR, 2022.
- [47] Hongkang Li, Shuai Zhang, and Meng Wang. Learning and generalization of one-hidden-layer neural networks, going beyond standard gaussian data. In 2022 56th Annual Conference on Information Sciences and Systems (CISS), pages 37–42. IEEE, 2022.
- [48] Hongkang Li, Meng Wang, Sijia Liu, and Pin-Yu Chen. A theoretical understanding of shallow vision transformers: Learning, generalization, and sample complexity. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview. net/forum?id=jClGv3Qjhb.
- [49] Hongkang Li, Meng Wang, Tengfei Ma, Sijia Liu, ZAIXI ZHANG, and Pin-Yu Chen. What improves the generalization of graph transformer? a theoretical dive into self-attention and positional encoding. In *NeurIPS 2023 Workshop: New Frontiers in Graph Learning*, 2023. URL https://openreview.net/forum?id=BaxFC3z9R6.

- [50] Hongkang Li, Meng Wang, Songtao Lu, Xiaodong Cui, and Pin-Yu Chen. How do nonlinear transformers learn and generalize in in-context learning? In *Forty-first International Conference on Machine Learning*, 2024. URL https://openreview.net/forum?id= 14HTPws9P6.
- [51] Hongkang Li, Meng Wang, Songtao Lu, Xiaodong Cui, and Pin-Yu Chen. How do nonlinear transformers acquire generalization-guaranteed cot ability? In *High-dimensional Learning Dynamics 2024: The Emergence of Structure and Reasoning*, 2024.
- [52] Hongkang Li, Meng Wang, Songtao Lu, Xiaodong Cui, and Pin-Yu Chen. Training nonlinear transformers for chain-of-thought inference: A theoretical generalization analysis. *arXiv* preprint arXiv:2410.02167, 2024.
- [53] Hongkang Li, Meng Wang, Shuai Zhang, Sijia Liu, and Pin-Yu Chen. Learning on transformers is provable low-rank and sparse: A one-layer analysis. In 2024 IEEE 13rd Sensor Array and Multichannel Signal Processing Workshop (SAM), pages 1–5. IEEE, 2024.
- [54] Hongkang Li, Shuai Zhang, Yihua Zhang, Meng Wang, Sijia Liu, and Pin-Yu Chen. How does promoting the minority fraction affect generalization? a theoretical study of one-hidden-layer neural network on group imbalance. *IEEE Journal of Selected Topics in Signal Processing*, 2024.
- [55] Hongkang Li, Yihua Zhang, Shuai Zhang, Meng Wang, Sijia Liu, and Pin-Yu Chen. When is task vector provably effective for model editing? a generalization analysis of nonlinear transformers. *arXiv preprint arXiv:2504.10957*, 2025.
- [56] Yingcong Li, Ankit S Rawat, and Samet Oymak. Fine-grained analysis of in-context linear estimation: Data, architecture, and beyond. *Advances in Neural Information Processing Systems*, 37:138324–138364, 2024.
- [57] Yingcong Li, Davoud Ataee Tarzanagh, Ankit Singh Rawat, Maryam Fazel, and Samet Oymak. Gating is weighting: Understanding gated linear attention through in-context learning. arXiv preprint arXiv:2504.04308, 2025.
- [58] Jiachang Liu, Dinghan Shen, Yizhe Zhang, William B Dolan, Lawrence Carin, and Weizhu Chen. What makes good in-context examples for gpt-3? In Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures, pages 100–114, 2022.
- [59] Yuankai Luo, Hongkang Li, Qijiong Liu, Lei Shi, and Xiao-Ming Wu. Node identifiers: Compact, discrete representations for efficient graph learning. *arXiv preprint arXiv:2405.16435*, 2024.
- [60] Yuankai Luo, Hongkang Li, Lei Shi, and Xiao-Ming Wu. Enhancing graph transformers with hierarchical distance structural encoding. *Advances in Neural Information Processing Systems*, 37:57150–57182, 2024.
- [61] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. Foundations of machine learning. MIT press, 2018.

- [62] Eshaan Nichani, Jason D. Lee, and Alberto Bietti. Understanding factual recall in transformers via associative memories. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=hwSmPOAmhk.
- [63] Kazusato Oko, Yujin Song, Taiji Suzuki, and Denny Wu. Pretrained transformer efficiently learns low-dimensional target functions in-context. In *The Thirty-eighth Annual Conference* on Neural Information Processing Systems, 2024. URL https://openreview.net/ forum?id=uHcG5Y6fdB.
- [64] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction heads. arXiv preprint arXiv:2209.11895, 2022.
- [65] OpenAI. Gpt-4 technical report. OpenAI, 2023.
- [66] Jongho Park, Jaeseung Park, Zheyang Xiong, Nayoung Lee, Jaewoong Cho, Samet Oymak, Kangwook Lee, and Dimitris Papailiopoulos. Can mamba learn how to learn? a comparative study on in-context learning tasks. *arXiv preprint arXiv:2402.04248*, 2024.
- [67] Yao Qiang, Xiangyu Zhou, and Dongxiao Zhu. Hijacking large language models via adversarial in-context learning. *arXiv preprint arXiv:2311.09948*, 2023.
- [68] Gautam Reddy. The mechanistic basis of data dependence and abrupt learning in an in-context classification task. In *The Twelfth International Conference on Learning Representations*, 2024.
- [69] Yunwei Ren, Zixuan Wang, and Jason D Lee. Learning and transferring sparse contextual bigrams with linear transformers. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [70] Zhenmei Shi, Junyi Wei, and Yingyu Liang. A theoretical analysis on feature learning in neural networks: Emergence from inputs and advantage over fixed features. In *International Conference on Learning Representations*, 2021.
- [71] Jiawei Sun, Hongkang Li, and Meng Wang. How do skip connections affect graph convolutional networks with graph sampling? a theoretical analysis on generalization, 2024. URL https: //openreview.net/forum?id=J2pMoN2pon.
- [72] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971, 2023.
- [73] Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. arXiv preprint arXiv:1011.3027, 2010.
- [74] 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.

- [75] Roger Waleffe, Wonmin Byeon, Duncan Riach, Brandon Norick, Vijay Korthikanti, Tri Dao, Albert Gu, Ali Hatamizadeh, Sudhakar Singh, Deepak Narayanan, et al. An empirical study of mamba-based language models. arXiv preprint arXiv:2406.07887, 2024.
- [76] Alexander Wan, Eric Wallace, Sheng Shen, and Dan Klein. Poisoning language models during instruction tuning. In *International Conference on Machine Learning*, pages 35413–35425. PMLR, 2023.
- [77] Jingfeng Wu, Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, and Peter L Bartlett. How many pretraining tasks are needed for in-context learning of linear regression? arXiv preprint arXiv:2310.08391, 2023.
- [78] Zhiyong Wu, Yaoxiang Wang, Jiacheng Ye, and Lingpeng Kong. Self-adaptive in-context learning: An information compression perspective for in-context example selection and ordering. *ACL*, 2023.
- [79] Qianxiong Xu, Xuanyi Liu, Lanyun Zhu, Guosheng Lin, Cheng Long, Ziyue Li, and Rui Zhao. Hybrid mamba for few-shot segmentation. *Advances in Neural Information Processing Systems*, 37:73858–73883, 2024.
- [80] Hongru Yang, Bhavya Kailkhura, Zhangyang Wang, and Yingbin Liang. Training dynamics of transformers to recognize word co-occurrence via gradient flow analysis. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [81] Tong Yang, Yu Huang, Yingbin Liang, and Yuejie Chi. In-context learning with representations: Contextual generalization of trained transformers. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [82] Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. arXiv preprint arXiv:2306.09927, 2023.
- [83] Shuai Zhang, Hongkang Li, Meng Wang, Miao Liu, Pin-Yu Chen, Songtao Lu, Sijia Liu, Keerthiram Murugesan, and Subhajit Chaudhury. On the convergence and sample complexity analysis of deep q-networks with  $\epsilon$ -greedy exploration. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [84] Yihua Zhang, Hongkang Li, Yuguang Yao, Aochuan Chen, Shuai Zhang, Pin-Yu Chen, Meng Wang, and Sijia Liu. Visual prompting reimagined: The power of activation prompts, 2024. URL https://openreview.net/forum?id=0b328CMwn1.
- [85] Shuai Zhao, Meihuizi Jia, Luu Anh Tuan, Fengjun Pan, and Jinming Wen. Universal vulnerabilities in large language models: Backdoor attacks for in-context learning. In *Proceedings of the* 2024 Conference on Empirical Methods in Natural Language Processing, pages 11507–11522, 2024.
- [86] Kai Zhong, Zhao Song, Prateek Jain, Peter L Bartlett, and Inderjit S Dhillon. Recovery guarantees for one-hidden-layer neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 4140–4149, 2017. URL https:// arxiv.org/pdf/1706.03175.pdf.

### **Appendix A. Main Theoretical Results**

We first summarize insights of our theoretical results in Section 3. Then, we introduce our formulation for analysis in Section A.1. Section A.2 presents the theoretical results of learning for ICL generalization with Mamba. Section A.3 analyzes linear Transformers for a comparison with Mamba models. We finally characterize the ICL mechanism by the trained Mamba in Section A.4.

# A.1. Data and Tasks Modeling

We follow the definition of tasks in [11, 40, 48]. Specifically, there are  $M_1$  relevant patterns  $\{\mu_j\}_{j=1}^{M_1} \cup \{\mu_k\}_{k=1}^{M_2}$  with  $M_1 + M_2 < d$ . All the patterns from  $\{\mu_j\}_{j=1}^{M_1} \cup \{\nu_k\}_{k=1}^{M_2}$  are orthogonal to each other, with  $\|\mu_j\| = \|\nu_k\| = \beta$  for  $j \in [M_1]$ ,  $k \in [M_2]$ , and the constant  $\beta \ge 1$ . Each input x contains one relevant pattern that determines the label, and one irrelevant pattern that does not affect the label. We consider a set of binary classification tasks in  $\mathcal{T}$  where the binary labels are determined by the relevant patterns. For instance, for a task f that is determined by  $(\mu_a, \mu_b), a, b \in [M_1]$ , the label of  $x_{query}$  is z = 1 (or z = -1) if the input  $x_{query}$  contains  $\mu_a$  (or  $\mu_b$ ), respectively.

**Training Stage:** For a given task f, we consider learning with a  $p_a \in [0, 1)$  fraction of examples containing additive outliers  $\{\boldsymbol{v}_r^*\}_{r=1}^V$  that can affect the label of corresponding examples in each prompt, where  $\boldsymbol{v}_s^* \perp \boldsymbol{\mu}_j$ ,  $\boldsymbol{v}_s^* \perp \boldsymbol{\nu}_k$  for any  $j \in [M_1]$ ,  $k \in [M_2]$ , and  $s \in [V]$ . The input of each context example satisfies

$$\boldsymbol{x} = \begin{cases} \boldsymbol{\mu}_j + \kappa \boldsymbol{\nu}_k + \kappa_a \boldsymbol{v}_s^*, & \text{with a probability of } p_a \\ \boldsymbol{\mu}_j + \kappa \boldsymbol{\nu}_k, & \text{with a probability of } 1 - p_a, \end{cases}$$
(6)

for some  $s \in [V]$ , where  $j \in [M_1]$  and  $k \in [M_2]$  are arbitrarily selected.  $\kappa$  follows a uniform distribution U(-K, K) with  $K \leq 1/2$ .  $v_s^*$  is uniformly sampled from  $\{v_r^*\}_{r=1}^V$ . No additive outliers exist in  $x_{query}$ . We then present the definition of training prompts.

**Definition 1** (Training prompts) Given a task  $f \in \mathcal{T}$  with  $\mu_a$  and  $\mu_b$  as the two different decisive patterns, a training prompt  $P \sim D$  with  $l_{tr}$  context examples is constructed as follows.

- $x_{query}$  follows the second line of (6) with j equally selected from  $\{a, b\}$  and contains no  $v_s^*$ .
- Each  $x_i$  contains  $\mu_a$  or  $\mu_b$  with equal probability  $i \in [l_{tr}]$ , following (6).
- $y_i = +1$  (or  $y_i = -1$ ) if the relevant pattern of  $x_i$  is  $\mu_a$  (or  $\mu_b$ ), and  $x_i$  does not contain any  $v_s^*$ .  $y_i$  is selected from  $\{+1, -1\}$  with equal probability if  $x_i$  contains a certain  $v_s^*$  for  $s \in [V]$ .

We include outliers in the training prompt ( $p_a \ge 0$ ) to encourage the model to learn to ignore examples containing outliers. This improves robustness during inference when prompts may also include such outliers. Our motivation stems from noise-aware training to mitigate data poisoning or hijacking attacks in ICL [34, 67, 76], where prompts are corrupted with noisy or random labels. When  $p_a = 0$ , the setup reduces to the case where training prompts contain no outliers, aligning with the theoretical setup in [36, 50, 82].

**Inference Stage**: During inference, we consider that the outliers in the testing prompt can differ from those in the training prompt in several ways, including their direction, magnitude, and the

fraction of examples affected. Specifically, the data input during the testing follow

$$\boldsymbol{x} = \begin{cases} \boldsymbol{\mu}_{j} + \kappa' \boldsymbol{\nu}_{k} + \kappa'_{a} \boldsymbol{v}_{s}^{*\prime}, & \text{with a probability of } \alpha \\ \boldsymbol{\mu}_{j} + \kappa' \boldsymbol{\nu}_{k}, & \text{with a probability of } 1 - \alpha, \end{cases}$$
(7)

for some  $v_s^{*'} \in \mathcal{V}'$ ,  $\kappa_a' > 0$ , and  $\kappa' \sim U(-K', K')$  with K' > 1.  $\alpha \in [0, 1)$  is the probability of examples containing the testing additive outliers in  $\mathcal{V}'$ .

**Definition 2** (*Testing prompts*) Given a task  $f \in \mathcal{T}$  with  $\mu_a$  and  $\mu_b$  as the relevant patterns, a testing  $\mathbf{P}' \sim \mathcal{D}'$  with  $l_{ts}$  context examples is constructed as follows. each testing query  $\mathbf{x}_{query}$  only follows the second line of (7) without outliers. Each context input  $\mathbf{x}_i$ ,  $i \in [l_{ts}]$ , follows (7). If  $\mathbf{x}_i$  does not contain any  $\mathbf{v}_s^* \in \mathcal{V}'$ , then  $y_i = +1$  (or  $y_i = -1$ ) if the relevant pattern of  $\mathbf{x}_i$  is  $\mu_a$  (or  $\mu_b$ ). If  $\mathbf{x}_i$  contains a certain  $\mathbf{v}_s^* \in \mathcal{V}'$ , then  $y_i$  can be an arbitrary function that maps  $\mathbf{x}_i$  to  $\{+1, -1\}$ .

The testing prompt P' differs from the training prompt P in two key aspects. First, the outlier patterns, the magnitude of the outliers, and the magnitude of the irrelevant patterns can differ from those in P. While the training prompts include V distinct outlier patterns, the testing prompts may contain an unbounded number of outlier variations. Second, the labels associated with examples containing outliers can be generated by any deterministic or probabilistic function. This flexibility allows our framework to model a wide range of noisy testing prompts in practice. For instance,

**Example 1** Consider a data poisoning attack on a text sentiment classification task in [34, 76]. In one such attack, whenever the phrase "James Bond" is inserted into the example, the label is always set to positive, regardless of the original sentiment of the input. This illustrates a case where all examples containing the outlier are deterministically mapped to a targeted label +1.

### A.2. Learning, Generalization, and Sample Complexity Analysis of Mamba Models

To enable the model learned from data in training tasks  $\mathcal{T}_{tr}$  to generalize well across all tasks in  $\mathcal{T}$ , we require Condition 3.2 from [50] for  $\mathcal{T}_{tr}$ . We restate this condition as Condition 1, along with a construction of a training task set that satisfies it in the Appendix. The high-level idea is that the training tasks  $\mathcal{T}_{tr}$  should uniformly cover all of the relevant patterns and labels appearing in  $\mathcal{T}$  such that no bias from the training tasks is introduced to the learning process.

Following [48, 70], we assume the training labels are balanced, i.e.,  $||\{n : z^n = +1\}| - |\{n : z^n = -1\}|| = O(\sqrt{N})$ . Let  $B_T := \max\{\epsilon^{-2}, M_1(1 - p_a)^{-1}\} \cdot \log \epsilon^{-1}$ . We have the following result.

**Theorem 1** (Convergence and Sample Complexity of Mamba) For any  $\epsilon > 0$ , of (i)  $B \gtrsim B_M := \max\{B_T, \beta^{-4}V^2\kappa_a^{-2}(1-p_a)^{-2}\log\epsilon^{-1}\}$ , (ii)  $V\beta^{-4} \lesssim \kappa_a \lesssim V\beta(1-p_a)p_a^{-1}\epsilon^{-1}$ , and (iii)

$$p_a^{-1} poly(M_1^{\kappa_a}) \gtrsim l_{tr} \gtrsim (1 - p_a)^{-1} \log M_1,$$
(8)

then (iv) after

$$T \ge T_M = \Theta(\eta^{-1}(1-p_a)^{-1}\beta^{-2}M_1)$$
(9)

iterations with  $\eta \leq 1$  and using N = BT samples, we have that

$$L_{f\in\mathcal{T},\boldsymbol{P}\sim\mathcal{D}}^{0-1}(\Psi^{(T)};\boldsymbol{P},z) \le \epsilon.$$
(10)

**Remark 1** Theorem 1 provides the convergence and sample complexity analysis of training a onelayer Mamba model to enhance its ICL ability. We characterize the sufficient conditions on the batch size, the magnitude of additive outliers, the prompt length, and the required number of iterations. The convergent model has desirable generalization on all tasks in T, including those not appearing in the training data, when the prompt is constructed in the same way as the training data.

Condition (ii) requires that the magnitude of additive outliers be moderate and scale with V. This ensures that outliers are neither too small to be easily detectable by the model nor excessively large (i.e., less than  $\Theta(\epsilon^{-1})$ ), which would diminish the influence of relevant patterns. Conditions (iii) and (iv) show that the required number of context examples in the prompt and the number of iterations scale as  $(1 - p_a)^{-1}$ . This implies that a higher fraction of outlier-containing context examples slows convergence and requires more context examples.

**Remark 2** When  $p_a = 0$ , Theorem 1 reduces to the case where Mamba is trained with prompts that contain no outliers and serves as the Mamba counterpart to Theorem 3.3 in [50], which addresses Transformers. Although [36, 50] analyze ICL training without outliers for Transformers, their analyses do not directly extend to Mamba due to the significant structural differences between the two architectures. To the best of our knowledge, we are the first to analyze the training dynamics of Mamba in the ICL setting, under a more general scenario where prompts may contain outliers.

We then study the generalization performance on testing prompts with distribution-shifted additive outliers using the trained Mamba.

**Theorem 2** (ICL Generalization on Distribution-shifted Prompts with Outliers) During the inference, if (a) the outlier pattern  $v_s^{*'}$  belongs to

$$\mathcal{V}' = \Big\{ \boldsymbol{v} \Big| \boldsymbol{v} = \sum_{i=1}^{V} \lambda_i \boldsymbol{v}_i^*, \sum_{i=1}^{V} \lambda_i \ge L > 0, \|\boldsymbol{v}\| = 1 \Big\},$$
(11)

(b) the outlier magnitude  $\kappa'_a \in [\kappa_a, \Theta(V\beta p_a^{-1}\kappa_a^{-1}L^{-1}(1-p_a)\epsilon^{-1})]$ , (c)  $\alpha < \min(1, p_a l_{tr}/l_{ts})$ , and (d) the number of context examples

$$\alpha^{-1} \operatorname{poly}(M_1^{\kappa_a}) \gtrsim l_{ts} \gtrsim (1-\alpha)^{-1} \log M_1, \tag{12}$$

then for testing prompt P' defined by Definition 2, the trained model  $\Psi^{(T)}$  satisfies

$$L_{f\in\mathcal{T},\mathbf{P}'\sim\mathcal{D}'}^{0-1}(\Psi^{(T)};\mathbf{P}',z)\leq\epsilon.$$
(13)

#### Remark 3

Theorem 2 shows that the model trained under Theorem 1 generalizes well and remains robust when tested on prompts containing a signification fraction of unseen distribution-shifted outliers. Each additive outlier in the test prompt can be expressed as a linear combination of the V training outlier patterns, with coefficients summing to a positive value (Condition (a)). This formulation captures a wide range of possible outlier patterns at test time. Notably, the fraction of examples with outliers  $\alpha$  in the test prompt is less than  $\min(1, p_a l_{tr}/l_{ts})$ , which can be close to 1 if the prompt length is selected in a way such that  $p_a l_{tr}/l_{ts} \ge 1$  (Condition (c)). Thus, Mamba can be trained to maintain ICL generalization in the presence of a large fraction of outlier examples. Conditions (b) and (d) impose mild requirements on the outlier magnitude and the context length, respectively. Condition (b) requires that the magnitude of test-time outliers be at least as large as that of the training outliers. Condition (d) ensures that the context prompt is sufficiently long to include enough clean examples for correct prediction, while also imposing an upper bound on the total number of outliers.

#### A.3. Theoretical Results for Linear Transformers and A Comparison with Mamba Models

In this section, we compare Mamba with linear Transformer, where the Transformer model is formulated by setting the nonlinear gating  $G_{i,l+1}(w) = 1$  in (3) for  $i \in [l+1]$ , as discussed in Section 2. Such a comparison can help understand the impact of the nonlinear gating on model training, in-context generalization, and robustness.<sup>1</sup>.

**Theorem 3** (Convergence and Sample Complexity for Transformer) As long as (i)  $B \gtrsim B_T$ , (ii)  $\kappa_a \lesssim VB(1-p_a)p_a^{-1}\epsilon^{-1}$ , (iii)

$$l_{tr} \gtrsim (1 - p_a)^{-2} p_a \log M_1,$$
 (14)

then (iv) after

$$T \ge T_T = \Theta(\eta^{-1}(1-p_a)^{-1}\beta^{-2}l_{tr}^{-1}M_1)$$
(15)

*iterations with*  $\eta \leq 1$  *and* N = BT *samples, we have that* 

$$L^{0-1}_{f\in\mathcal{T},\boldsymbol{P}\sim\mathcal{D}}(\Psi^{(T)};\boldsymbol{P},z)\leq\epsilon.$$
(16)

### **Remark 4**

Theorem 3 characterizes the sufficient conditions for the convergence and generalization of training a one-layer Transformer with linear attention using prompts containing outliers as formulated by Definition 1. Comparing conditions (i)-(iv) with those in Theorem 1 on Mamba models, one can see that, to achieve a  $\epsilon$  generalization error, linear Transformers need a smaller batch size, a smaller number of training iterations, and a less restrictive requirement for the magnitude of additive outliers. To see this, Theorem 1 indicates that the required batch size for Mamba models is at least  $B_M$ , which is defined as the larger of value  $B_T$  and another constant, while the required batch size for linear Transformers is  $B_T$ . The required number of training iterations for Mamba is  $T_M$ , which equals  $\Theta(l_{tr}) \cdot T_T$ , and that is larger than that for linear Transformers,  $T_T$ , by a scaling of  $\Theta(l_{tr}) > 1$ . The conditions for  $\kappa_a$  for Mamba and Transformer models share the same upper bound, but  $\kappa_a$  has an extra lower bound for Mamba.

**Theorem 4** (Generalization using Transformers) During the inference, if (a) in Theorem 2, (b)  $\kappa'_a \leq \Theta(V\beta p_a^{-1}(1-p_a)\kappa_a^{-1}L^{-1}l_{tr}\epsilon^{-1})$ , (c)  $\alpha \in [0, 1/2)$ , and (d) the number of context examples

$$l_{ts} \gtrsim \max\{\Theta((1-\alpha)^{-1}), \Theta((\frac{1}{2}-\alpha)^{-2}\alpha)\}\log M_1,$$
 (17)

then the trained model  $\Psi^{(T)}$  satisfies

$$\underline{L}_{f\in\mathcal{T},\boldsymbol{P}'\sim\mathcal{D}'}^{0-1}(\Psi^{(T)};\boldsymbol{P},z)\leq\epsilon.$$
(18)

<sup>1.</sup> The comparison is made between sufficient conditions for the desired generalization. Given the consistent training setup and analytical tools used, we still believe this is a fair comparison with the main insights validated empirically in Section B.1. Although providing necessary conditions would lead to a more rigorous comparison, we leave this as future work due to its technical difficulty.

#### **Remark 5**

Theorem 4 establishes the conditions under which a Transformer model, trained according to Theorem 3, can generalize effectively on testing prompts with possible outliers, as defined in Definition 2. In contrast to Theorem 2 for Mamba, the Transformer guarantees generalization only when the outlier fraction satisfies  $\alpha < 1/2$ , whereas Mamba can remain robust when  $\alpha$  goes to 1 (Condition (c)). This highlights that Mamba achieves better in-context generalization performance in the presence of distribution-shifted additive outliers, particularly when outlier-containing context examples are in the majority. This conclusion is consistent with the empirical findings of [66], which observed that Mamba outperforms Transformers in many-outlier regression tasks.

### A.4. The Mechanism of Mamba in implementing ICL

We next examine the mechanism by which the trained Mamba model from Theorem 1 performs ICL on prompts containing additive outliers. This analysis provides deeper insights into the differences between Mamba and Transformer models. We begin by showing, in Corollary 1, that the linear attention of the learned Mamba model assigns greater weight to context examples that share the same relevant pattern as the query.

#### **Corollary 1**

Let  $\mathcal{N}_1 \subseteq [l_{ts}]$  denote the index sets of context examples that share the same relevant pattern as the query  $\mathbf{x}_{query}$ . Then, for the model trained by Theorem 1 after  $T \ge T_M$  iterations in (9), we have with a high probability, for  $\mathbf{P}'$  defined by Definition 2,

$$\sum_{i \in \mathcal{N}_1} \tilde{\boldsymbol{p}}_i^{\top} \boldsymbol{W}_B^{(T)^{\top}} \boldsymbol{W}_C^{(T)} \tilde{\boldsymbol{p}}_{query} \ge \Theta(1); \sum_{i \in [l_{ts}] \setminus \mathcal{N}_1} \tilde{\boldsymbol{p}}_i^{\top} \boldsymbol{W}_B^{(T)^{\top}} \boldsymbol{W}_C^{(T)} \tilde{\boldsymbol{p}}_{query} \le \Theta((1-p_a)^{-1}\epsilon).$$
(19)

#### Remark 6

Corollary 1 illustrates that for the testing prompt  $\mathcal{P}'$ , the learned Mamba model will let the attention scores be concentrated on examples with the same relevant pattern as the query, i.e., the sum of these attention scores will increase to be larger than  $\Theta(1)$ , while the sum of attention score on examples with other different relevant pattern from the query is upper bounded by a small order of  $(1 - p_a)^{-1}\epsilon$ . This enforces the model to focus on examples with the same relevant pattern as the query when making the prediction.

Corollary 1 reveals an insight similar to the "induction head" mechanism [13, 64, 68] observed in softmax attention layers for ICL. However, our result is established in the context of linear attention, suggesting that different attention variants may share fundamentally similar internal mechanisms.

We then show that the nonlinear gating mechanism in Mamba models enables ICL by effectively ignoring context examples containing outliers and focusing on those that are closer to the query.

#### **Corollary 2**

For the trained model by Theorem 1 after  $T \ge T_M$  iterations in (9), we have that with a high probability, for  $\tilde{p}_i$  that contain a  $v_s^{*'} \in \mathcal{V}'$ ,

$$G_{i,l_{ts}+1}(\boldsymbol{w}^{(T)}) \le O(poly(M_1)^{-1}).$$
 (20)

Denote  $h(j) \in [l_{ts}]$   $(j \leq l_{ts})$  as the index of context example that is the *j*-th closest to the query and does not contain any  $\mathbf{v}_s^{*'} \in \mathcal{V}'$ . Then, with a high probability, we have

$$G_{h(j),l_{ts}+1}(\boldsymbol{w}^{(T)}) \ge \Theta(1/2^{j-1}).$$
 (21)

### **Remark 7**

Corollary 2 indicates that the nonlinear gating function  $G_{i,l_{ts}+1}(\boldsymbol{w}^{(T)})$  serves two main purposes: (i) filtering out examples containing additive outliers and (ii) inducing a local bias, as observed in [31], that focuses on examples near the query. Specifically, (20) unveils that on examples with outliers,  $G_{i,l_{ts}+1}(\boldsymbol{w}^{(T)})$  is close to 0, effectively suppressing their influence. (21) shows that for clean examples, the nonlinear gating values decay exponentially with the distance (in index) from the query. Hence, combing Corollaries 1 and 2, one can see that the model primarily relies on examples that are close to the query, do not contain outliers, and share the same relevant pattern as the query for prediction, resulting in desirable ICL performance even in the presence of outliers.

Corollary 2 characterizes the role of the nonlinear gating layer, Mamba's key structural difference from the Transformer. This distinction explains their performance gap: while nonlinear gating makes Mamba more challenging to optimize, it also enables Mamba to suppress outlier-containing examples more effectively, resulting in superior robustness when handling prompts with many outliers.

# **Appendix B. Experiment**

We generate synthetic data following Section A.1. Let d = 30,  $M_1 = 6$ ,  $M_2 = 10$ , V = 3. For generalization with unseen outliers, let  $v_1^{*'} = 0.7v_1^* + 0.6v_2^* - 0.4v_3^*$ ,  $v_2^{*'} = 0.4v_1^* + 0.7v_2^* - 0.6v_3^*$ ,  $v_3^{*'} = -0.7v_1^* + 0.5v_2^* + 0.5v_3^*$ , with L = 0.3.  $l_{ts} = l_{tr} = 20$ . Let  $\delta = 0.2$ ,  $\beta = 3$ ,  $\kappa_a = 2$ . The experiments are conducted on a single NVIDIA RTX A5000 GPU.

#### B.1. Comparison of One-Layer Mamba and Linear Transformer on ICL with Outliers

The learning model is a one-layer Mamba defined in (3) and a one-layer single-head Transformer by making  $G_{i,l+1}(w) = 1$  for  $i \in [l+1]$ . We set  $p_a = 0.6$ . We consider three types of outlier-relevant labeling functions during inference. If the context examples in a given prompt P' contains any additive outlier, the corresponding context label will be (A) flipped, (B) mapping to one targeted label out of  $\{+1, -1\}$ , or (C) randomly chosen from  $\{+1, -1\}$  with equal probability. Figure 1 shows that under three different forms of outliers, the classification error of Mamba is smaller than 0.01 even when  $\alpha$  is close to 0.8. In contrast, the classification error of linear Transformers is large as long as  $\alpha > 1/2$ . This is consistent with Remark 5: the linear transformer can tolerate at most a 1/2 fraction of outliers in the prompt, whereas Mamba can tolerate a fraction of outliers close to that seen during training, which can be close to 1.



Figure 1: ICL classification error of Mamba and linear Transformer against  $\alpha$  with different prompt outliers. (A) Label flipping. (B) Targeted labeling. (C) Random labeling.

### **B.2.** The ICL Mechanism of Multi-Layer Mamba

The learning model is a three-layer Mamba and a three-layer single-head linear Transformer.  $p_a = 0.4$ . Figure 2 shows the first-layer attention scores in the testing prompt. The sum of attention scores on the examples that share the same pattern as the query is significantly larger than that on examples with other patterns, and this gap increases during training. This verifies Corollary 1. Figure 3 shows that the first-layer gating values with  $\alpha = 0.3$  of outlier-containing examples are very small (red bars), while those of clean examples are relatively large and exhibit an approximately exponential decay with increasing distance from the query (green bars). This is consistent with (20) and (21) in Corollary 2. The results of attention scores and gating values in the other two layers exhibit the same trend as the first layer and are shown in Section C in Appendix due to the space limit.

Next, we study the impact of the positions of context examples with  $\alpha = 0.5$ . Table 2 presents the ICL performance under three different placements of outlier examples: all positioned farthest from the query (FQ), closest to the query (CQ), or at random positions (R). We find that Mamba is highly sensitive to the position of outliers, whereas the linear Transformer (LT) is much less affected. This is because, when outliers are placed close to the query, the clean examples that share the same pattern as the query are pushed farther away, and the gating values on these examples decay exponentially according to (21), thereby degrading ICL performance.



Table 2: ICL accuracy of 3-layer

attention scores on examples with the Figure 3: The 1st-layer gating value Mamba and linear Transformers (LT) same or a different relevant pattern as of examples with (red) or without with different example arrangement. the query. (green) additive outliers.

# Appendix C. Additional Experiments, Related Works, and the Algorithm

Figure 2: The summation of 1st-layer

We first show the visualization result of the second and the third linear attention and nonlinear gating layers of the three-layer Mamba analyzed in Section B.2. The conclusions in Figures 4 and 5 are aligned with Figures 2 and 3, respectively.



Figure 4: The summation of attention scores in the 2nd and 3rd layers.

We then introduce other related theoretical works on optimization and generalization of neural networks in this section. Some works [22, 47, 54, 83, 86] study the generalization of neural networks using the model recovery framework by investigating the local convexity around a ground truth parameter of the problem. The neural-tangent-kernel (NTK) analyses [4, 5, 12, 16, 37, 46, 71] study



Figure 5: The gating values of examples with or without outliers in the 2nd and 3rd layers.

this problem in the overparameterized setting to linearize the neural network around the initialization, with the resulting generalization performance irrelevant to the feature distribution. Another line of works [11, 17–19, 43, 48, 51–53, 59, 60, 70, 84] studies the generalization of neural networks by formulating data that contains discriminative and unimportant features. Our analysis in this work is aligned with the last framework to probe the generalization of Mamba and Transformers.

We next present the training algorithm introduced in Section 2.

Algorithm 1 Training with Stochastic Gradient Descent (SGD)

- 1: Hyperparameters: The step size  $\eta$ , the number of iterations T, batch size B.
- 2: Initialization:  $W_B^{(0)}$  and  $W_C^{(0)}$  are initialized such that the first d diagonal entries of  $W_B^{(0)}$  and  $W_C^{(0)}$  are set as  $\delta \in (0, 0.2]$ .  $w^{(0)} \sim \mathcal{N}(0, I_{d+1}/(d+1))$ .
- 3: Training by SGD: For each iteration, we independently sample P ~ D, f ∈ T<sub>tr</sub> to form a batch of training prompt and labels {P<sup>n</sup>, z<sup>n</sup>}<sub>n∈B<sub>t</sub></sub> as introduced in Section A.1. Each relevant pattern is sampled equally likely in each batch. For each t = 0, 1, · · · , T − 1 and W<sup>(t)</sup> ∈ Ψ<sup>(t)</sup>,

$$\boldsymbol{W}^{(t+1)} = \boldsymbol{W}^{(t)} - \eta \cdot \frac{1}{B} \sum_{n \in \mathcal{B}_t} \nabla_{\boldsymbol{W}^{(t)}} \ell(\boldsymbol{\Psi}^{(t)}; \boldsymbol{P}^n, \boldsymbol{z}^n).$$
(22)

4: **Output:**  $W_B^{(T)}, W_C^{(T)}, w^{(T)}$ .

# Appendix D. Key Lemmas

We first present Table 3 for a summary of notations used in the proof.

**Lemma 1** (Multiplicative Chernoff bounds, Theorem D.4 of [61]) Let  $X_1, \dots, X_m$  be independent random variables drawn according to some distribution  $\mathcal{D}$  with mean p and support included in [0, 1]. Then, for any  $\gamma \in [0, \frac{1}{p} - 1]$ , the following inequality holds for  $\hat{p} = \frac{1}{m} \sum_{i=1}^{m} X_i$ :

$$\Pr(\hat{p} \ge (1+\gamma)p) \le e^{-\frac{mp\gamma^2}{3}},\tag{23}$$

$$\Pr(\hat{p} \le (1 - \gamma)p) \le e^{-\frac{mp\gamma^2}{2}}.$$
(24)

**Definition 1** [73] We say X is a sub-Gaussian random variable with sub-Gaussian norm K > 0, if  $(\mathbb{E}|X|^p)^{\frac{1}{p}} \leq K\sqrt{p}$  for all  $p \geq 1$ . In addition, the sub-Gaussian norm of X, denoted  $||X||_{\psi_2}$ , is defined as  $||X||_{\psi_2} = \sup_{p\geq 1} p^{-\frac{1}{2}} (\mathbb{E}|X|^p)^{\frac{1}{p}}$ .

Notations	Annotation
$ ilde{A}_i,  ilde{B}_i, C_i$	Parameters in Mamba.
$\sigma(\cdot)$	sigmoid function.
$x_s^n, y_s^n$	$x_s^n$ is the input data for classification. $y_s^n$ is the label for $x_s^n$ .
$P^n, z^n$	$P^n$ is a prompt that consists of the query and $l$ pairs of
	examples of $x_s^n$ and $y_s^n$ , $s \in [l]$ . $z^n \in \{+1, -1\}$ is the
	binary label of $p_{query}^n$ .
$F(\Psi; \mathbf{P}^n), \ell(\Psi; \mathbf{P}^n, z^n)$	$F(\Psi; \mathbf{P}^n)$ is the model output for $\mathbf{P}^n$ with $\Psi$ as the parame-
	ter. $\ell(\Psi; \mathbf{P}^n, z^n)$ is the loss function given the input $\mathbf{P}^n$ and
	the corresponding label $z^n$ .
$L^{0-1}_{f\in\mathcal{T},\boldsymbol{P}'\sim\mathcal{D}'}(\Psi;\boldsymbol{P}',z)$	The classification error of $\Psi$ given $P' \sim D'$ as the input and
	$f\in\mathcal{T}.$
$\mu_j,  u_k$	$\mu_j$ and $\nu_k$ are the relevant and irrelevant patterns in the data
	formulation.
$M_1, M_2$	$M_1$ is the number of relevant patterns. $M_2$ is the number of
	irrelevant patterns.
$oldsymbol{v}_s^*,oldsymbol{v}_s^{*\prime},\kappa_a,\kappa_a^\prime$	$v_s^*, s \in [V]$ is the additive outlier for training. $v_s^{*'}$ is the
	additive outlier for testing. $\kappa_a$ and $\kappa'_a$ are the magnitudes of
	outliers in training and testing.
$p_a, \alpha$	$p_a$ is the probability of examples containing additive outliers
	in training prompts. $\alpha$ is the probability of examples contain-
	ing outliers in testing prompts.
$ \mathcal{B}_b $	$\mathcal{B}_b$ is the SGD batch at the <i>b</i> -th iteration. $l_{ts}$ is the prompt
	length of the testing data.
$l_{tr}, l_{ts}$	$l_{tr}$ is the prompt length of the training data. $l_{ts}$ is the prompt
	length of the testing data.
$\mathcal{O}(), \Omega(), \Theta()$	We follow the convention that $f(x) = O(g(x))$ (or $\Omega(g(x))$ ,
	$\Theta(g(x)))$ means that $f(x)$ increases at most, at least, or in
	the order of $g(x)$ , respectively.
$ \gtrsim,\lesssim$	$f(x) \gtrsim g(x)$ (or $f(x) \lesssim g(x)$ ) means that $f(x) \ge \Omega(g(x))$
	$ $ (or $f(x) \lesssim \mathcal{O}(g(x))).$

Table 3: Summary of Notations

**Lemma 2** ([73] Proposition 5.1, Hoeffding's inequality) Let  $X_1, X_2, \dots, X_N$  be independent centered sub-gaussian random variables, and let  $K = \max_i ||\mathbf{X}_i||_{\psi_2}$ . Then for every  $\mathbf{a} = (a_1, \dots, a_N) \in \mathbb{R}^N$  and every  $t \ge 0$ , we have

$$\Pr\left(\left|\sum_{i=1}^{N} a_i X_i\right| \ge t\right) \le e \cdot \exp\left(-\frac{ct^2}{K^2 \|\boldsymbol{a}\|^2}\right),\tag{25}$$

where c > 0 is an absolute constant.

**Lemma 3** For any  $j \neq j', j'' \in [M_1]$ ,  $k \neq k' \in [M_2]$ , and  $s \in [V]$ , j'' where  $\mu_j$  and  $\mu_{j''}$ form a training task, and j' where  $\mu_j$  and  $\mu_{j'}$  does not form a training task, we have that for  $W \in \{W_B, W_C\}$ , if  $B \gtrsim \max\{(1 - p_a)^{-1}M_1 \log \epsilon^{-1}, (1 - p_a)^{-2} \log \epsilon^{-1}\}$ ,

$$-(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})\eta \cdot \sum_{b=1}^{t+1} \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{W}^{(b)}} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top} \gtrsim \eta(t+1) \frac{1}{M_{1}} (1-p_{a})\beta, \quad (26)$$

$$\left| (\boldsymbol{v}_s^{*\top}, \boldsymbol{0}^{\top}) \boldsymbol{\eta} \cdot \sum_{b=1}^{t+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}^{(b)}} (\boldsymbol{\mu}_j^{\top}, \boldsymbol{0}^{\top})^{\top} \right| \le \frac{\eta \beta (t+1) p_a \kappa_a}{M_1 V} \cdot \sqrt{\frac{\log B}{B}}, \quad (27)$$

$$-(\boldsymbol{\mu}_{j'}^{\top}, 0^{\top})\eta \cdot \sum_{b=1}^{t_0+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_B^{(b)}} (\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top} = 0,$$
(28)

$$-(\boldsymbol{\mu}_{j''}{}^{\top}, 0^{\top})\eta \cdot \sum_{b=1}^{t+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}^{(t)}} (\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top} \le -\eta(t+1) \frac{1}{M_1} (1-p_a)\beta, \quad (29)$$

$$\left| - (\boldsymbol{\nu}_k^{\top}, 0^{\top})\eta \cdot \sum_{b=1}^{t+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}^{(b)}} (\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top} \right| \le \frac{\eta(t+1)\beta}{M_1 M_2} \sqrt{\frac{\log B}{B}}, \quad (30)$$

$$-(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})\boldsymbol{\eta} \cdot \sum_{b=1}^{t+1} \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{W}^{(b)}} (\boldsymbol{\nu}_{k}^{\top}, 0^{\top})^{\top} \Big| \leq \frac{\eta(t+1)\beta}{M_{1}M_{2}} \sqrt{\frac{\log B}{B}}, \quad (31)$$

$$-(\boldsymbol{\nu}_{k}^{\top}, 0^{\top})\eta \cdot \sum_{b=1}^{t+1} \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{W}^{(b)}} (\boldsymbol{\nu}_{k}^{\top}, 0^{\top})^{\top} \Big| \leq \frac{\eta(t+1)\beta}{M_{2}} \sqrt{\frac{\log B}{B}}, \quad (32)$$

$$\left| - (\boldsymbol{\nu}_{k'}^{\top}, 0^{\top})\eta \cdot \sum_{b=1}^{t+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}^{(b)}} (\boldsymbol{\nu}_k^{\top}, 0^{\top})^{\top} \right| \le \frac{\eta(t+1)\beta}{M_2^2} \sqrt{\frac{\log B}{B}}.$$
 (33)

Lemma 4 When  $t \lesssim \min\{\eta^{-1}\beta^{-2}\kappa_a^{-1}(1-p_a)^{-1}V, \eta^{-1}M_1^{\frac{2}{3}}\beta^{-\frac{2}{3}}\kappa_a^{-\frac{1}{3}}(1-p_a)^{-1}V^{\frac{1}{3}}\}$ , as long as  $l \gtrsim (1-p_a)^{-1}\log\epsilon^{-1}$ , (34)

$$B \gtrsim \beta^{-4} \kappa_a^{-2} (1 - p_a)^{-2} V^2 \log \epsilon^{-1},$$
(35)

we have that for any  $s \in [V]$ ,

$$\mathbf{v}_{s}^{*\top}\mathbf{w}^{(t)} \lesssim -\frac{\eta\beta^{2}t\kappa_{a}(1-p_{a})}{V} - \eta\sum_{i=1}^{t}i^{2}(\frac{\eta^{2}(1-p_{a})^{3}\beta^{2}}{M_{1}^{2}})\frac{\kappa_{a}}{V},$$
(36)

$$(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})\boldsymbol{w}^{(t)} = \Theta\left(-\frac{\eta(1-p_{a})\beta^{2}(t)}{M_{1}} - \sum_{i=1}^{t-1}i^{2}\cdot\left(\frac{\eta^{3}(1-p_{a})^{3}\beta^{2}}{M_{1}^{3}}\right)\right).$$
(37)

For  $p_s$  that does not contain any  $v_o^*$ ,  $o \in [V]$ , and  $p_r$  that contains a  $v_o^*$ ,  $o \in [V]$ ,  $r \neq s$ , we have

$$-\frac{\eta(1-p_a)\beta^2 t}{M_1} - \sum_{i=1}^t i^2 \cdot \left(\frac{\eta^3(1-p_a)^3\beta^2}{M_1^3}\right) \lesssim \boldsymbol{w}^{(t)^{\top}} \boldsymbol{p}_s < 0,$$
(38)

$$\boldsymbol{w}^{(t)^{\top}}\boldsymbol{p}_{r} \lesssim -\eta t \beta^{2} \kappa_{a} (1-p_{a}) < \boldsymbol{w}^{(t)^{\top}}\boldsymbol{p}_{s} < 0.$$
(39)

**Lemma 5** When  $t \gtrsim \eta^{-1}(1-p_a)^{-1}\beta^{-2}M_1$  and  $\kappa_a \gtrsim V\beta^{-4}$ , we have

$$\boldsymbol{w}^{(t)} \boldsymbol{p}_i \lesssim -\log M_1, \tag{40}$$

for  $p_i$  that contains a  $v_s^*$ ,  $s \in [V]$ , and

$$\boldsymbol{w}^{(t)^{\top}}\boldsymbol{p}_i \gtrsim -\Theta(1). \tag{41}$$

for  $p_i$  that does not contain any  $v_s^*$ ,  $s \in [V]$ .

**Lemma 6** When  $t \lesssim \min\{\eta^{-1}\beta^{-2}\kappa_a^{-1}(1-p_a)^{-1}V, \eta^{-1}M_1^{\frac{2}{3}}((1-p_a)\beta)^{-\frac{2}{3}}(\kappa_a(1-p_a))^{-\frac{1}{3}}V^{\frac{1}{3}}\},$  we have

$$\sum_{i=1}^{l} G_{i,l+1}(\boldsymbol{w}^{(t)})(l-i+1) \le \Theta(1).$$
(42)

**Condition 1** (Condition 3.2 of [50]) For any given  $j \in [M_1]$  and either label +1 or -1, the number of tasks in  $\mathcal{T}_{tr}$  that map  $\mu_j$  to that label is  $|\mathcal{T}_{tr}|/M_1(\geq 1)$ .

We introduce a construction of  $\mathcal{T}_{tr}$  that satisfies Condition 1 as follows. Let the *i*-th task function  $(i \in [M_1 - 1])$  in  $\mathcal{T}_{tr}$  map the queries with  $\mu_i$  and  $\mu_{i+1}$  as the relevant patterns to +1 and -1, respectively. The  $M_1$ -th task function maps  $\mu_{M_1}$  and  $\mu_1$  to +1 and -1, respectively. We can easily verify that such a  $\mathcal{T}_{tr}$  satisfies Condition 1 in this case.

# **Appendix E. Proof of Main Theorems**

### E.1. Proof of Theorem 1

**Proof** We know that there exists gradient noise caused by imbalanced patterns in each batchTherefore, by Hoeffding's inequality (25), for any  $W \in \Psi$ ,

$$\Pr\left(\left\|\frac{1}{|\mathcal{B}_b|}\sum_{n\in\mathcal{B}_b}\frac{\partial\ell(\Psi;\boldsymbol{P}^n,z^n)}{\partial\boldsymbol{W}} - \mathbb{E}\left[\frac{\partial\ell(\Psi;\boldsymbol{P}^n,z^n)}{\partial\boldsymbol{W}}\right]\right\| \ge \left|\mathbb{E}\left[\frac{\partial\ell(\Psi;\boldsymbol{P}^n,z^n)}{\partial\boldsymbol{W}}\right]\epsilon\right) \le e^{-B\epsilon^2} \le \epsilon,$$
(43)

if  $B \gtrsim \epsilon^{-2} \log \epsilon^{-1}$ . Combining (35), we require

$$B \gtrsim \max\{\beta^{-4} \kappa_a^{-2} (1-p_a)^{-2}, \epsilon^{-2}, M_1 (1-p_a)^{-1}\} \cdot \log \epsilon^{-1}.$$
(44)

When  $t \ge T = \Theta(\eta^{-1}(1-p_a)^{-1}\beta^{-2}M_1)$ , we have that for  $\boldsymbol{W} \in \{\boldsymbol{W}_B, \boldsymbol{W}_C\}$  and any  $j \in [M_1]$ ,

$$(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})\boldsymbol{W}^{(T)}(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top}$$

$$= (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})(\boldsymbol{W}^{(0)} - \eta \cdot \sum_{b=1}^{T} \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{W}^{(b)}})(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top}$$

$$\gtrsim 1,$$
(45)

where the last step comes from (26) in Lemma 3. Then, for  $p_i$  that shares the same pattern as the query, we have

$$\boldsymbol{p}_{i}^{\top} \boldsymbol{W}_{B}^{(T)^{\top}} \boldsymbol{W}_{C}^{(T)} \boldsymbol{p}_{query} \gtrsim \beta^{2} (1 + \kappa_{a} \mathbb{1}[\boldsymbol{p}_{i} \text{ contains any } \boldsymbol{v}_{s}^{*}]) + 1 - (1 - p_{a})^{-1} \epsilon \beta^{-1} / M_{2}$$

$$- (1 - p_{a})^{-1} p_{a} \kappa_{a} V^{-1} \beta^{-1} \epsilon \mathbb{1}[\boldsymbol{p}_{i} \text{ contains any } \boldsymbol{v}_{s}^{*}],$$

$$(46)$$

as long as  $\epsilon \in (0,1)$ .  $(1-p_a)^{-1}\epsilon/M_2$  comes from the correlation between  $\mu_j$  and  $\nu_k$ ,  $\nu_*$  and between  $\nu_k$  and  $\nu_*$ , and  $B \gtrsim \epsilon^{-2} \log \epsilon^{-1}$ . For  $p_i$  that shares a different pattern that does not form a training task from the query, with a high probability, we have

$$\boldsymbol{p}_{i}^{\top}\boldsymbol{W}_{B}^{(T)}{}^{\top}\boldsymbol{W}_{C}^{(T)}\boldsymbol{p}_{query} \leq (1-p_{a})^{-1}\epsilon\beta^{-1}/M_{2} + (1-p_{a})^{-1}p_{a}\kappa_{a}V^{-1}\beta^{-1}\epsilon\mathbb{1}[\boldsymbol{p}_{i} \text{ contains any } \boldsymbol{v}_{s}^{*}].$$
(47)

Meanwhile, for  $p_i$  that contains a  $v_s^*$ ,  $s \in [V]$ , we have

$$G_{i,l+1}(\boldsymbol{w}^{(T)}) \le \sigma(\boldsymbol{w}^{(T)^{\top}}\boldsymbol{p}_i) \lesssim O(\operatorname{poly}(M_1^{\kappa_a})^{-1}), \tag{48}$$

by Lemma 5. We have that for the  $p_{i^*}$  that does not contain any  $v_s^*$ ,  $s \in [V]$  and is the closest to the query, by Lemma 5,

$$G_{i^*,l+1}(\boldsymbol{w}^{(T)}) \gtrsim (1 - \frac{1}{\operatorname{poly}(M_1^{\kappa_a})})^{lp_a} \sigma(\boldsymbol{w}^{(T)^{\top}} \boldsymbol{p}_{i^*})$$
  
$$\gtrsim (1 - \frac{lp_a}{\operatorname{poly}(M_1^{\kappa_a})}) \sigma(\boldsymbol{w}^{(T)^{\top}} \boldsymbol{p}_{i^*})$$
  
$$\gtrsim (1 - \frac{lp_a}{\operatorname{poly}(M_1^{\kappa_a})}).$$
(49)

Hence, for  $\boldsymbol{P}$  with z = +1, with a high probability, we have

 $(\mathbf{T})$ 

$$F(\Psi^{(T)}, \mathbf{P})$$

$$\gtrsim (1 - (1 - p_a)^{-1} \epsilon / M_2 - (1 - p_a)^{-1} p_a \kappa_a V^{-1} \beta^{-1} \epsilon) \cdot \sum_{i=1}^{l_{tr}(1 - p_a)^{-1}} (1$$

$$- \max_{\mathbf{p}_i \text{ contains no } \mathbf{v}_s^*} \{\sigma(\mathbf{w}^{(T)^{\top}} \mathbf{p}_i)\})^{i-1} \cdot \min_{\mathbf{p}_i \text{ contains no } \mathbf{v}_s^*} \{\sigma(\mathbf{w}^{(T)^{\top}} \mathbf{p}_i)\}$$

$$\gtrsim \frac{(1 - (1 - \max_{\mathbf{p}_i \text{ contains no } \mathbf{v}_s^*} \{\sigma(\mathbf{w}^{(T)^{\top}} \mathbf{p}_i)\})^{l_{tr}(1 - p_a)}) \cdot \min_{\mathbf{p}_i \text{ contains no } \mathbf{v}_s^*} \{\sigma(\mathbf{w}^{(T)^{\top}} \mathbf{p}_i)\}}{\max_{\mathbf{p}_i \text{ contains no } \mathbf{v}_s^*} \{\sigma(\mathbf{w}^{(T)^{\top}} \mathbf{p}_i)\}}$$

$$>\Theta(1) \cdot (1 - \frac{1}{M_1})$$

$$>1,$$
(50)

where the second to last step holds if  $p_a^{-1}$  poly $(M_1^{\kappa_a}) \gtrsim l_{tr} \gtrsim (1 - p_a)^{-1} \log M_1$  and for  $p_i$  that contains no  $\boldsymbol{v}_s^*$ ,  $\sigma(\boldsymbol{w}^{(T)^{\top}}\boldsymbol{p}_i) \in (0, 1/2)$ . Similarly, we can also derive that for  $\boldsymbol{P}$  with z = -1, we have

$$F(\Psi^{(T)}, \mathbf{P}) < -1.$$
 (51)

Then, we study in-domain generalization. By (43), for any given testing prompt embedding P with z = +1, we have that with a high probability of  $1 - \epsilon$ ,

$$F(\Psi^{(T)}; \mathbf{P}) \ge 1 - \epsilon, \tag{52}$$

and if z = -1,

$$F(\Psi^{(T)}; \mathbf{P}) \le -1 + \epsilon.$$
(53)

Therefore,

$$L^{0-1}_{\boldsymbol{x}_{query} \sim \mathcal{D}, f \in \mathcal{T}}(\Psi^{(T)}; \boldsymbol{P}, z) \le \epsilon.$$
(54)

# E.2. Proof of Theorem 2

#### Proof

By Lemma 3, we have that for any  $j \in [M_1]$  and  $k \neq k' \in [M_2]$ ,

$$(\boldsymbol{\nu}_k^{\top}, 0^{\top}) \boldsymbol{W}^{(T)}(\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top} \lesssim \frac{\epsilon (1 - p_a)^{-1} \beta^{-1}}{M_2},$$
(55)

$$(\boldsymbol{\mu}_{j}^{\top}, \boldsymbol{0}^{\top})\boldsymbol{W}^{(T)}(\boldsymbol{\nu}_{k}^{\top}, \boldsymbol{0}^{\top})^{\top} \lesssim \frac{\epsilon(1 - p_{a})^{-1}\beta^{-1}}{M_{2}},$$
(56)

$$(\boldsymbol{\nu}_{k}^{\top}, 0^{\top}) \boldsymbol{W}^{(T)}(\boldsymbol{\nu}_{k}^{\top}, 0^{\top})^{\top} \lesssim \frac{\epsilon (1 - p_{a})^{-1} \beta^{-1} M_{1}}{M_{2}}.$$
 (57)

$$(\boldsymbol{\nu}_{k}^{\top}, 0^{\top}) \boldsymbol{W}^{(T)}(\boldsymbol{\nu}_{k'}^{\top}, 0^{\top})^{\top} \lesssim \frac{\epsilon (1 - p_{a})^{-1} \beta^{-1} M_{1}}{M_{2}^{2}}.$$
 (58)

Meanwhile, we have that for  $v_s^{*\prime} \in \mathcal{V}'$  with  $v_s^{*\prime} = \sum_{i=1}^V \lambda_i v_s^*$ ,

$$(\boldsymbol{v}_{s}^{\prime *^{\top}}, \boldsymbol{0}^{\top})\boldsymbol{W}^{(T)}(\boldsymbol{\mu}_{j}^{\top}, \boldsymbol{0}^{\top})^{\top} \lesssim \epsilon(1 - p_{a})^{-1}p_{a}\kappa_{a}V^{-1}\beta^{-1} \cdot L.$$
(59)

Therefore, we have that for  $p_i$  that shares the same pattern as the query,

$$\boldsymbol{p}_{i}^{\top} \boldsymbol{W}_{B}^{(T)} \boldsymbol{W}_{C}^{(T)} \boldsymbol{p}_{query} \gtrsim 1 - \epsilon (1 - p_{a})^{-1} \cdot \frac{1}{M_{2}} - \epsilon (1 - p_{a})^{-1} p_{a} V^{-1} \kappa_{a} \beta^{-1} \cdot \kappa_{a}^{\prime} L.$$
(60)

For  $p_i$  that shares a different pattern from the query, we have

$$|\boldsymbol{p}_{i}^{\top} \boldsymbol{W}_{B}^{(T)}^{\top} \boldsymbol{W}_{C}^{(T)} \boldsymbol{p}_{query}| \lesssim \epsilon (1 + (1 - p_{a})^{-1} / M_{2} + (1 - p_{a})^{-1} p_{a} V^{-1} \kappa_{a} \beta^{-1} \cdot \kappa_{a}' L).$$
(61)

Meanwhile, for  $oldsymbol{p}_i$  that contains a  $oldsymbol{v}_s^{*\prime} \in \mathcal{V}'$ , we have

$$G_{i,l+1}(\boldsymbol{w}^{(T)}) \le \sigma(\boldsymbol{w}^{(T)^{\top}}\boldsymbol{p}_i) \lesssim O(\operatorname{poly}(M_1^{\kappa_a'})^{-1}), \tag{62}$$

by Lemma 5. We have that for the  $p_{i^*}$  that does not contain any  $v_s^{*'} \in \mathcal{V}'$  and is the closest to the query, by Lemma 5,

$$G_{i^*,l+1}(\boldsymbol{w}^{(T)}) \gtrsim (1 - \frac{1}{\operatorname{poly}(M_1^{\kappa_a'})})^{l_{ts}\alpha} \sigma(\boldsymbol{w}^{(T)\top} \boldsymbol{p}_{i^*})$$
  
$$\gtrsim (1 - \frac{l_{ts}\alpha}{\operatorname{poly}(M_1^{\kappa_a'})}).$$
(63)

Hence, for P' with z = +1, with a high probability, we have

$$F(\Psi^{(T)}, g(\mathbf{P})) \\ \geq (1 - (1 - p_{a})^{-1} \epsilon / M_{2} - \epsilon (1 - p_{a})^{-1} p_{a} V^{-1} \kappa_{a} \beta^{-1} \cdot \kappa_{a}' L) \cdot \sum_{i=1}^{l_{ts}(1 - \alpha) - 1} (1 \\ - \max_{p_{i} \text{ contains no } v_{s}^{*} \in \mathcal{V}'} \{\sigma(\boldsymbol{w}^{(T)^{\top}} p_{i})\})^{i-1} \cdot \min_{p_{i} \text{ contains no } v_{s}^{*} \in \mathcal{V}'} \{\sigma(\boldsymbol{w}^{(T)^{\top}} p_{i})\} \\ \geq \Theta((1 - (1 - p_{a})^{-1} \epsilon / M_{2} - \epsilon (1 - p_{a})^{-1} p_{a} V^{-1} \kappa_{a} \beta^{-1} \cdot (\kappa_{a} + \kappa_{a}' L - \kappa_{a})) \\ \cdot (1 - \frac{l_{ts} \alpha}{\text{poly}(M_{1}^{\kappa_{a}'})})) \\ = \Theta((1 - \epsilon (1 - p_{a})^{-1} p_{a} V^{-1} \kappa_{a} \beta^{-1} \cdot (\kappa_{a}' L - \kappa_{a}))(1 - \frac{l_{tr} p_{a}}{\text{poly}(M_{1}^{\kappa_{a}})}) \\ \cdot (1 - \frac{\frac{l_{ts} \alpha}{\text{poly}(M_{1}^{\kappa_{a}})} - \frac{l_{tr} p_{a}}{\text{poly}(M_{1}^{\kappa_{a}})}})) \\ \geq \Theta(1 - \epsilon (1 - p_{a})^{-1} p_{a} V^{-1} \kappa_{a} \beta^{-1} \cdot (\kappa_{a}' L - \kappa_{a}) - (\frac{l_{ts} \alpha}{\text{poly}(M_{1}^{\kappa_{a}})} - \frac{l_{tr} p_{a}}{\text{poly}(M_{1}^{\kappa_{a}})})) \\ \geq 1 - (\epsilon (1 - p_{a})^{-1} p_{a} V^{-1} \kappa_{a} \beta^{-1} \cdot (\kappa_{a}' L - \kappa_{a}) + \frac{l_{ts} \alpha}{\text{poly}(M_{1}^{\kappa_{a}})} - \frac{l_{tr} p_{a}}{\text{poly}(M_{1}^{\kappa_{a}})}), \end{cases}$$

where we consider the worst-case order that makes all examples that contain  $v_s^{*'} \in \mathcal{V}'$  right before the query, such that there is a scaling of  $1 - \frac{l_{ts}\alpha}{\operatorname{poly}(M_1^{\kappa'_a})}$  in the second step. The trained model still selects examples with the same pattern as the query no matter whether there is a certain  $v_s'^*$ added to the token if  $\kappa'_a \leq V\beta p_a^{-1}(1-p_a)\kappa_a^{-1}L^{-1}\epsilon^{-1}$ . Then, flipping the labels of examples with any of  $v_s'^*$  can change the model output the most. If  $l_{ts} \leq \alpha^{-1}\operatorname{poly}(M_1^{\kappa_a})$ ,  $\kappa_a \leq \kappa'_a \leq \Theta(L^{-1}(\kappa_a+V\beta p_a^{-1}(1-p_a)\kappa_a^{-1}\epsilon^{-1}))$ ,  $\alpha \leq \min\{1, p_a \cdot l_{tr}/l_{ts}\}$ , we have that that with a probability of  $1 - \log M_1$ ,

$$F(\Psi^{(T)}, g(\boldsymbol{P})) > 0 \tag{65}$$

Therefore, we can derive that

$$L^{0-1}_{\boldsymbol{x}_{query}\sim\mathcal{D},f\in\mathcal{T}}(\Psi^{(T)};\boldsymbol{P},z)\leq\epsilon.$$
(66)

# E.3. Proof of Theorem 3

# Proof

By the Chernoff bound of Bernoulli distribution in Lemma 1, we can obtain that for any n and  $s \in [V]$ ,

$$\Pr\left(\frac{1}{l}\sum_{i=1}^{l}\mathbb{1}[\boldsymbol{p}_{i}^{n} \text{ contains } \boldsymbol{\mu}_{a} \text{ and no any } \boldsymbol{v}_{s}^{*}] \leq (1-c)(1-p_{a})\frac{1}{2}\right) \leq e^{-lc^{2}\frac{(1-p_{a})}{2}} = \epsilon, \quad (67)$$

for some  $c \in (0, 1)$ . Hence, with a high probability,

$$l \gtrsim (1 - p_a)^{-1} \log \epsilon^{-1}.$$
 (68)

We know that there exists gradient noise caused by imbalanced patterns in each batchTherefore, by Hoeffding's inequality (25), for any  $W \in \{W_Q, W_K\}$ ,

$$\Pr\left(\left\|\frac{1}{|\mathcal{B}_b|}\sum_{n\in\mathcal{B}_b}\frac{\partial\ell(\Psi;\boldsymbol{P}^n,z^n)}{\partial\boldsymbol{W}} - \mathbb{E}\left[\frac{\partial\ell(\Psi;\boldsymbol{P}^n,z^n)}{\partial\boldsymbol{W}}\right]\right\| \ge \left|\mathbb{E}\left[\frac{\partial\ell(\Psi;\boldsymbol{P}^n,z^n)}{\partial\boldsymbol{W}}\right]\epsilon\right) \le e^{-B\epsilon^2} \le \epsilon,$$
(69)

if  $B \gtrsim \epsilon^{-2} \log \epsilon^{-1}$ . Therefore, we require

$$B \gtrsim \max\{\epsilon^{-2}, (1-p_a)^{-1}M_1\}\log\epsilon^{-1}.$$
 (70)

Let  $G_{i,l+1}(\boldsymbol{w}^{(T)}) = 1$  for any  $i \leq l+1$ . Following the proof in Theorem 1, we have that when

$$T \ge \Theta(\eta^{-1}(1-p_a)^{-1}l_{tr}^{-1}\beta^{-1}M_1), \tag{71}$$

we have

$$F(\Psi^{(T)}, \mathbf{P}) \gtrsim (1 - (1 - p_a)^{-1} \epsilon / M_2 - (1 - p_a)^{-1} p_a \kappa_a V^{-1} \beta^{-1} \epsilon)$$
  
>1. (72)

Therefore, by Lemma 1,

$$L_{\boldsymbol{x}_{query}\sim\mathcal{D},f\in\mathcal{T}}^{0-1}(\Psi^{(T)};\boldsymbol{P},z)$$

$$\lesssim \Pr(\frac{1}{l}\sum_{i=1}^{l}\mathbb{1}[\boldsymbol{p}_{i} \text{ with the same pattern as } \boldsymbol{p}_{query} \text{ but a flipped label}] > \frac{1}{2})$$

$$= \Pr(\frac{1}{l}\sum_{i=1}^{l}\mathbb{1}[\boldsymbol{p}_{i} \text{ with the same pattern as } \boldsymbol{p}_{query} \text{ but a flipped label}] - \frac{p_{a}}{2} > \frac{p_{a}}{2} \cdot \frac{1-p_{a}}{p_{a}})$$

$$\leq e^{-l(1-p_{a})^{2}p_{a}}$$

$$\leq \epsilon,$$
(73)

 $\text{if } l \ge (1 - p_a)^{-2} p_a \log \epsilon^{-1}.$ 

# E.4. Proof of Theorem 4

**Proof** By setting  $G_{i,l+1}(\boldsymbol{w}^{(T)}) = 1$  for any  $i \leq l+1$ , we have for any  $j \in [M_1], k' \neq k \in [M_2]$ 

$$(\boldsymbol{\nu}_{k}^{\top}, 0^{\top}) \boldsymbol{W}^{(T)}(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top} \lesssim \frac{\epsilon \beta^{-1} (1 - p_{a})^{-1} l_{tr}^{-1}}{M_{2}},$$
(74)

$$(\boldsymbol{\mu}_{j}^{\top}, 0^{\top}) \boldsymbol{W}^{(T)}(\boldsymbol{\nu}_{k}^{\top}, 0^{\top})^{\top} \lesssim \frac{\epsilon \beta^{-1} (1 - p_{a})^{-1} l_{tr}^{-1}}{M_{2}}.$$
 (75)

$$(\boldsymbol{\nu}_{k}^{\top}, 0^{\top}) \boldsymbol{W}^{(T)}(\boldsymbol{\nu}_{k}^{\top}, 0^{\top})^{\top} \lesssim \frac{\epsilon \beta^{-1} (1 - p_{a})^{-1} l_{tr}^{-1} M_{1}}{M_{2}}.$$
 (76)

$$(\boldsymbol{\nu}_{k'}^{\top}, 0^{\top}) \boldsymbol{W}^{(T)}(\boldsymbol{\nu}_{k}^{\top}, 0^{\top})^{\top} \lesssim \frac{\epsilon \beta^{-1} (1 - p_a)^{-1} l_{tr}^{-1} M_1}{M_2^2}.$$
 (77)

Meanwhile, we have that for  $\boldsymbol{v}_s^{*\prime} \in \mathcal{V}'$  with  $\boldsymbol{v}_s^{*\prime} = \sum_{i=1}^V \lambda_i \boldsymbol{v}_s^*$ ,

$$(\boldsymbol{v}_{s}^{\prime *^{\top}}, 0^{\top})\boldsymbol{W}^{(T)}(\boldsymbol{\mu}_{j}^{\prime \top}, 0^{\top})^{\top} \lesssim \epsilon \beta^{-1} (1 - p_{a})^{-1} p_{a} \kappa_{a} V^{-1} l_{tr}^{-1} \kappa_{a}^{\prime} L.$$
(78)

Therefore, we have that for  $p_i$  that shares the same pattern as the query,

$$\boldsymbol{p}_{i}^{\top} \boldsymbol{W}_{B}^{(T)} \boldsymbol{W}_{C}^{(T)} \boldsymbol{p}_{query} \gtrsim 1 - \epsilon \cdot \frac{\beta^{-1} (1 - p_{a})^{-1} l_{tr}^{-1}}{M_{2}} - \epsilon (1 - p_{a})^{-1} \beta^{-1} p_{a} \kappa_{a} V^{-1} l_{tr}^{-1} L \kappa_{a}^{\prime}.$$
 (79)

For  $p_i$  that shares a different pattern from the query, we have

$$|\boldsymbol{p}_{i}^{\top}\boldsymbol{W}_{B}^{(T)}|^{\top}\boldsymbol{W}_{C}^{(T)}\boldsymbol{p}_{query}| \lesssim \epsilon(1+\beta^{-1}(1-p_{a})^{-1}l_{tr}^{-1}/M_{2}+(1-p_{a})^{-1}\beta^{-1}p_{a}\kappa_{a}V^{-1}l_{tr}^{-1}\kappa_{a}'L).$$
(80)

Therefore, the trained model still selects examples with the same pattern as the query no matter whether there is a certain  $v'_s$  added to the token if  $\kappa'_a \leq V\beta p_a^{-1}(1-p_a)\kappa_a^{-1}L^{-1}l_{tr}\epsilon^{-1}$ . Then, flipping the labels of examples with any of  $v'_s$  can change the model output the most. With  $\alpha < 1/2$ , we can derive that

$$L_{\boldsymbol{x}_{query}\sim\mathcal{D},f\in\mathcal{T}}^{0-1}(\Psi^{(T)};\boldsymbol{P},z) = \Pr(\frac{1}{l_{ts}}\sum_{i=1}^{l_{ts}}\mathbbm{1}[\boldsymbol{p}_{i} \text{ with the same pattern as } \boldsymbol{p}_{query} \text{ but a flipped label}] - \frac{\alpha}{2} > \frac{\alpha}{2} \cdot \frac{\frac{1}{2} - \alpha}{\alpha}) \quad (81)$$
$$\leq e^{-l_{ts}(\frac{1}{2} - \alpha)^{2}\alpha} \leq \epsilon,$$

as long as

$$l_{ts} \ge \max\{\Theta((1-\alpha)^{-1}), \Theta((\frac{1}{2}-\alpha)^{-2}\alpha)\}\log\epsilon^{-1}.$$
 (82)

# E.4.1. PROOF OF COROLLARY 1

**Proof** The first part of (19) comes from (46) since  $\beta \ge 1$  is a constant. The second part of (19) comes from (47) plus  $\kappa_a V^{-1} \beta^{-1} \ge \beta^{-5}$  with  $\beta \ge 1$  as a constant order.

# E.4.2. PROOF OF COROLLARY 2

**Proof** (20) comes from (62) plus  $\kappa'_a \ge \Theta(1)$ . (21) is derived as follows. By (63), we have

$$G_{h(1),l_{ts}+1}(\boldsymbol{w}^{(T)}) \ge \Theta(1).$$
(83)

Then, combining (39) and (20), we have that if  $p_s$  does not contain any outliers,

$$1 - \sigma(\boldsymbol{w}^{(T)^{\top}}\boldsymbol{p}_s) \ge \frac{1}{2}.$$
(84)

Then, with a high probability

$$G_{h(j),l_{ts}+1}(\boldsymbol{w}^{(T)}) \ge G_{h(j),l_{ts}+1}(\boldsymbol{w}^{(T)}) \cdot \frac{1}{2^{j-1}} \cdot (1 - \Theta(\operatorname{poly}(M_1)^{-1}))^{l_{ts}\alpha} \cdot \Theta(1) \ge \Theta(\frac{1}{2^{j-1}}).$$
(85)

# **Appendix F. Proof of Supportive Lemmas**

### F.1. Derivation of (3)

Proof

By formulation in Section 2, we have

$$\tilde{\boldsymbol{A}}_{i} = \mathbf{1}_{d_{0}} \operatorname{diag}(\exp(\Delta_{i}\boldsymbol{A}))^{\top} = \mathbf{1}_{d_{0}} \operatorname{diag}(e^{-\boldsymbol{I}_{l+1}\Delta_{i}})^{\top} = \mathbf{1}_{d_{0}} \operatorname{diag}(e^{-\boldsymbol{I}_{l+1}\log(1+e^{\boldsymbol{w}^{\top}\boldsymbol{x}_{i}})})^{\top} = \mathbf{1}_{d_{0}} \operatorname{diag}(e^{-\boldsymbol{I}_{l+1}\log(1+e^{\boldsymbol{w}^{\top}\boldsymbol{x}_{i}})})^{\top} = \mathbf{1}_{d_{0}} \mathbf{1}_{l+1}(\frac{1}{1+e^{\boldsymbol{w}^{\top}\boldsymbol{x}_{i}}})^{\top} = \mathbf{1}_{d_{0}} \mathbf{1}_{l+1}^{\top}(1-\sigma(\boldsymbol{w}^{\top}\boldsymbol{x}_{i})) \in \mathbb{R}^{d_{0} \times (l+1)}, \quad \sigma(\cdot) : \text{sigmoid function}$$

$$\tilde{\boldsymbol{B}}_{i} = \mathbf{1}_{d_{0}}(\Delta_{i}\boldsymbol{B}_{i})(\exp(\Delta_{i}\boldsymbol{A}) - \boldsymbol{I})(\Delta_{i}\boldsymbol{A})^{-1} = \mathbf{1}_{d_{0}}\boldsymbol{B}_{i}(\boldsymbol{I}_{l+1}\frac{1}{1+e^{\boldsymbol{w}^{\top}\boldsymbol{x}_{i}}} - \boldsymbol{I}_{l+1})(-\boldsymbol{I}_{l+1}) = \mathbf{1}_{d_{0}}\sigma(\boldsymbol{w}^{\top}\boldsymbol{x}_{i})\boldsymbol{B}_{i} := \boldsymbol{s}_{i}\boldsymbol{B}_{i} \in \mathbb{R}^{d_{0} \times (l+1)},$$
(87)

with  $s_i = \mathbf{1}_{d_0} \sigma(\boldsymbol{w}^{\top} \boldsymbol{x}_i)$ . Therefore,

$$h_{i} = h_{i-1} \odot \tilde{A}_{i} + (p_{i} \mathbf{1}_{l+1}^{\top}) \tilde{B}_{i}$$

$$= h_{i-1} \odot \tilde{A}_{i} + (p_{i} \odot s_{i}) B_{i}$$

$$= (h_{i-2} \odot \tilde{A}_{i-1} + (p_{i-1} \odot s_{i}) B_{i-1}) \odot \tilde{A}_{i} + p_{i} B_{i}$$

$$= h_{i-2} \odot \tilde{A}_{i-1} \odot \tilde{A}_{i} + (p_{i-1} \odot s_{i}) B_{i-1} \odot \tilde{A}_{i} + (p_{i} \odot s_{i}) B_{i}$$

$$= \cdots$$

$$= h_{0} \odot \tilde{A}_{1} \odot \cdots \odot \tilde{A}_{i} + \sum_{j=1}^{i} (p_{j} \odot s_{j}) B_{j} \odot \tilde{A}_{j+1} \cdots \odot \tilde{A}_{i} + (p_{i} \odot s_{i}) B_{i}$$

$$= \sum_{j=1}^{i} (p_{j} \odot s_{j}) B_{j} \odot (\tilde{A}_{i} \odot \cdots \odot \tilde{A}_{j+1}) + (p_{i} \odot s_{i}) B_{i},$$
(88)

Then, given  $oldsymbol{W}_C \in \mathbb{R}^{(l+1) imes d_0}$ , we have

$$o_{i} = h_{i}C_{i}$$

$$= h_{i}W_{C}p_{i}$$

$$= \sum_{j=1}^{i} (p_{j} \odot s_{j})B_{j}(\tilde{A}_{i} \odot \cdots \odot \tilde{A}_{j+1})W_{C}p_{i} + (p_{i} \odot s_{i})B_{i}W_{C}p_{i}$$

$$= \sum_{j=1}^{i} p_{j}p_{j}^{\top}W_{B}^{\top}\prod_{k=j+1}^{i} (1 - \sigma(w^{\top}p_{k})) \cdot \sigma(w^{\top}p_{j})W_{C}p_{i} + p_{j}p_{i}^{\top}W_{B}^{\top}\sigma(w^{\top}p_{i})W_{C}p_{i}$$

$$:= \sum_{j=1}^{i} G_{j,i}(w)p_{j}p_{j}^{\top}W_{B}^{\top}W_{C}p_{i},$$
(89)

where

$$G_{j,i}(\boldsymbol{w}) := \begin{cases} \prod_{k=j+1}^{i} (1 - \sigma(\boldsymbol{w}^{\top} \boldsymbol{p}_j)) \sigma(\boldsymbol{w}^{\top} \boldsymbol{p}_j), & \text{if } j < i \\ \sigma(\boldsymbol{w}^{\top} \boldsymbol{p}_i), & \text{if } j = i, \end{cases}$$
(90)

with  $\sigma(\cdot)$  as the sigmoid function. Therefore, we can obtain (3), i.e.,

$$F(\Psi; \boldsymbol{P}) = \boldsymbol{e}_{d+1}^{\top} \boldsymbol{o}_{l+1} = \sum_{i=1}^{l+1} G_{i,l+1}(\boldsymbol{w}) y_i \boldsymbol{p}_i^{\top} \boldsymbol{W}_B^{\top} \boldsymbol{W}_C \boldsymbol{p}_{query}.$$
(91)

# F.2. Proof of Lemma 3

**Proof** (a) When  $F(\Psi; \mathbf{P}^n) \in (-1, 1)$  for some  $n \in [N]$ , we have

$$\frac{\partial \ell(\Psi; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_C} = -z^n \sum_{i=1}^l G^n_{i,l+1}(\boldsymbol{w}) y^n_i \boldsymbol{W}_B \boldsymbol{p}^n_i \boldsymbol{p}^n_{query}^{\mathsf{T}}.$$
(92)

When t = 0, we know that with high probability,

$$|\boldsymbol{w}^{(0)^{\top}}\boldsymbol{x}_{j}| \lesssim \xi = \frac{1}{d+1},\tag{93}$$

$$|\sigma(\boldsymbol{w}^{(0)}^{\top}\boldsymbol{x}_j) - \frac{1}{2}| \lesssim \frac{|1 - e^{\pm\xi}|}{2(1 + e^{\pm\xi})} \lesssim \xi.$$
(94)

Then,

$$\frac{1}{2^{l+2-i}}(1-\xi(l+2-i)) \le G_{i,l+1}^{n(0)}(\boldsymbol{w}) \lesssim \frac{1}{2^{l+2-i}}(1+\xi(l+2-i)).$$
(95)

Let the IDR pattern of  $\mu_{query}^n$  be  $\mu_j$ ,  $j \in [M_1]$ . Note that  $\frac{1}{2} \cdot p_a$  fraction of examples correspond to  $\mu_j$  with poisoned labels. For different f,  $y_*^f = 1$  or -1 with 1/2 probability. By Lemma 1, we have for any  $i \in l$ ,

$$\Pr\left(\frac{1}{|\mathcal{B}_b|}\sum_{i\in\mathcal{B}_b}\mathbb{1}[\boldsymbol{x}_i^n \text{ contains } \boldsymbol{\mu}_j \text{ and no } \boldsymbol{v}_s^*] - (1-p_a) \le -\frac{c}{M_1}(1-p_a)\right) \le e^{-\frac{B(1-p_a)}{M_1}} \le \epsilon,$$
(96)

for some  $c \in (0,1)$  and  $\epsilon > 0$  if

$$B \gtrsim (1 - p_a)^{-1} M_1 \log \epsilon^{-1}.$$
 (97)

By (25), let  $\mathcal{B}_b' = \{i : i \in \mathcal{B}_b, x_i^n \text{ contains } \mu_j \text{ and } \nu_s^*, s \in [V]\}$  we have

$$\Pr\left(\left|\frac{1}{|\mathcal{B}_{b}'|}\sum_{i\in\mathcal{B}_{b}'}(\mathbb{1}[y_{i}^{n}=z^{n}]-\mathbb{1}[y_{i}^{n}=-z^{n}])\right| \geq \sqrt{\frac{\log B}{B}}\right) \leq M_{1}^{-C},\tag{98}$$

for some  $c \in (0, 1)$  and C > 1. Therefore, we have

$$-(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})\eta \cdot \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\Psi^{(0)}; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{W}_{C}^{(0)}} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top}$$

$$= (\boldsymbol{\mu}_{j}^{\top}, 0^{\top}) \frac{\eta}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} z^{n} \sum_{i=1}^{l} G_{i,l+1}^{n}(\boldsymbol{w}^{(0)}) y_{i}^{n} \boldsymbol{W}_{B}^{(0)} \boldsymbol{p}_{i}^{n} \boldsymbol{p}_{query}^{n}^{\top} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top}$$

$$\cdot \mathbb{1}[\boldsymbol{x}_{i}^{n} \text{ does not contain any } \boldsymbol{v}_{s}^{*}] + (\boldsymbol{\mu}_{j}^{\top}, 0^{\top}) \frac{\eta}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} z^{n} \sum_{i=1}^{l} G_{i,l+1}^{n}(\boldsymbol{w}^{(0)})$$

$$\cdot y_{i}^{n} \boldsymbol{W}_{B}^{(0)} \boldsymbol{p}_{i}^{n} \boldsymbol{p}_{query}^{n}^{\top} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top} \mathbb{1}[\boldsymbol{x}_{i}^{n} \text{ contains any } \boldsymbol{v}_{s}^{*}]$$

$$\geq \eta \cdot \frac{1}{2M_{1}} (1 - p_{a}) \sum_{i=1}^{l} G_{i,l+1}^{n}(\boldsymbol{w}^{(0)}) \beta - \eta \cdot \frac{1}{2M_{1}} \sum_{i=1}^{l} G_{i,l+1}^{n}(\boldsymbol{w}^{(0)}) \beta p_{a} \sqrt{\frac{\log B}{B}}$$

$$\geq \eta \frac{1}{4M_{1}} (1 - p_{a}) \beta (1 - \xi l),$$

$$(99)$$

where the last step holds if

$$B \gtrsim (1 - p_a)^{-2} \log \epsilon^{-1}.$$
 (100)

For  $\mu_{j'}, j' \neq j$ , that does not form a task in the training set, we have

$$-(\boldsymbol{\mu}_{j'}^{\top}, \boldsymbol{0}^{\top})\boldsymbol{\eta} \cdot \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\boldsymbol{\Psi}^{(0)}; \boldsymbol{P}^n, \boldsymbol{z}^n)}{\partial \boldsymbol{W}_C^{(0)}} (\boldsymbol{\mu}_j^{\top}, \boldsymbol{0}^{\top})^{\top} = 0$$
(101)

For  $\mu_{j''}, j'' \neq j$ , that forms a task in the training set, we have

$$- (\boldsymbol{\mu}_{j''}^{\top}, 0^{\top}) \eta \cdot \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\Psi^{(0)}; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{W}_{C}^{(0)}} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top}$$

$$= (\boldsymbol{\mu}_{j''}^{\top}, 0^{\top}) \frac{\eta}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} z^{n} \sum_{i=1}^{l} G_{i,l+1}^{n}(\boldsymbol{w}^{(0)}) y_{i}^{n} \boldsymbol{W}_{B}^{(0)} \boldsymbol{p}_{i}^{n} \boldsymbol{p}_{query}^{n}^{\top} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top} \qquad (102)$$

$$\lesssim -\eta \cdot \frac{1}{4M_{1}} (1 - p_{a}) \beta (1 - \xi l).$$

For  $\nu_k, \nu_{k'}$  with  $k, k' \in [M_2]$ , we have

$$-(\boldsymbol{\nu}_{k}^{\top},0^{\top})\boldsymbol{\eta}\cdot\frac{1}{|\mathcal{B}_{b}|}\sum_{n\in\mathcal{B}_{b}}\frac{\ell(\Psi^{(0)};\boldsymbol{P}^{n},z^{n})}{\partial\boldsymbol{W}_{C}^{(0)}}(\boldsymbol{\mu}_{j}^{\top},0^{\top})^{\top}\Big|\leq\frac{\eta\beta}{M_{1}M_{2}}\sqrt{\frac{\log B}{B}},$$
(103)

$$-(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})\boldsymbol{\eta} \cdot \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\Psi^{(0)}; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{W}_{C}^{(0)}} (\boldsymbol{\nu}_{k}^{\top}, 0^{\top})^{\top} \Big| \leq \frac{\eta \beta}{M_{2}M_{1}} \sqrt{\frac{\log B}{B}}.$$
 (104)

$$-(\boldsymbol{\nu}_{k'}{}^{\mathsf{T}},0{}^{\mathsf{T}})\eta \cdot \frac{1}{|\mathcal{B}_b|} \sum_{n\in\mathcal{B}_b} \frac{\ell(\Psi^{(0)};\boldsymbol{P}^n,z^n)}{\partial \boldsymbol{W}_C^{(0)}} (\boldsymbol{\nu}_k^{\mathsf{T}},0^{\mathsf{T}})^{\mathsf{T}} \Big| \le \frac{\eta\beta}{M_2^2} \sqrt{\frac{\log B}{B}}.$$
 (105)

$$-(\boldsymbol{\nu}_{k}^{\top},0^{\top})\eta\cdot\frac{1}{|\mathcal{B}_{b}|}\sum_{n\in\mathcal{B}_{b}}\frac{\ell(\Psi^{(0)};\boldsymbol{P}^{n},z^{n})}{\partial\boldsymbol{W}_{C}^{(0)}}(\boldsymbol{\nu}_{k}^{\top},0^{\top})^{\top}\Big|\leq\frac{\eta\beta}{M_{2}}\sqrt{\frac{\log B}{B}}.$$
(106)

Since that for  $oldsymbol{x}_i^n$  that contains  $oldsymbol{
u}_s^*$  for a certain  $s\in [V]$ ,

$$\Pr(y_i^n = z^n) = \Pr(y_i^n = -z^n) = \frac{1}{2},$$
(107)

we have

$$\left| (\boldsymbol{\nu}_{s}^{*\top}, 0^{\top}) \eta \cdot \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\Psi^{(0)}; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{W}_{C}^{(0)}} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top} \right|$$

$$= \left| (\boldsymbol{\nu}_{s}^{*\top}, 0^{\top}) \frac{\eta}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} z^{n} \sum_{i=1}^{l} G_{i,l+1}^{n} (\boldsymbol{w}^{(0)}) y_{i}^{n} \boldsymbol{W}_{B}^{(0)} \boldsymbol{p}_{i}^{n} \boldsymbol{p}_{query}^{n}^{\top} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top} \right|$$

$$\leq \frac{\eta \beta p_{a} \kappa_{*}}{M_{1} V} \cdot \sqrt{\frac{\log B}{B}},$$

$$(108)$$

Suppose that the conclusion holds when  $t = t_0$ . Then, when  $t = t_0 + 1$ , we have

$$- (\boldsymbol{\mu}_{j}^{\top}, 0^{\top}) \eta \cdot \sum_{b=1}^{t_{0}+1} \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{W}_{C}^{(b)}} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top} \\ = (\boldsymbol{\mu}_{j}^{\top}, 0^{\top}) \sum_{b=1}^{t_{0}+1} \frac{\eta}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} z^{n} \sum_{i=1}^{l} G_{i,l+1}^{n}(\boldsymbol{w}^{(b)}) y_{i}^{n} \boldsymbol{W}_{B}^{(b)} \boldsymbol{p}_{i}^{n} \boldsymbol{p}_{query}^{n}^{\top} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top} \\ \gtrsim \eta \cdot \sum_{b=1}^{t_{0}+1} \frac{1}{2M_{1}} (1 - p_{a}) \sum_{i=1}^{l} G_{i,l+1}^{n}(\boldsymbol{w}^{(t_{0})}) \beta \\ \gtrsim \eta (t_{0} + 1) \frac{1}{M_{1}} (1 - p_{a}) \beta. \end{cases}$$
(109)

The last step holds since  $\sum_{i=1}^{l} G_{i,l+1}^{n}(\boldsymbol{w}^{(t_0)}) \gtrsim 1$ . Similarly, we have that for any  $s \in [V]$ ,

$$\left| (\boldsymbol{\nu}_s^{*\top}, 0^{\top}) \eta \cdot \sum_{b=1}^{t_0+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_C^{(b)}} (\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top} \right| \le \frac{\eta \beta(t_0+1) p_a \kappa_*}{M_1} \cdot \sqrt{\frac{\log B}{B}}, \quad (110)$$

For  $\mu_{j'}, j' \neq j$ , that forms a task in the training set, we have

$$-(\boldsymbol{\mu}_{j'}^{\top}, 0^{\top})\eta \cdot \sum_{b=1}^{t_0+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_C^{(b)}} (\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top} = 0$$
(111)

For  $\mu_{j''}, j'' \neq j$ , that forms a task in the training set, we have

$$- (\boldsymbol{\mu}_{j''}^{\top}, 0^{\top}) \eta \cdot \sum_{b=1}^{t_0+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_C^{(b)}} (\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top}$$

$$\leq (\boldsymbol{\mu}_j^{\top}, 0^{\top}) \eta \cdot \sum_{b=1}^{t_0+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_C^{(b)}} (\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top}.$$

$$(112)$$

For  $\nu_k$ ,  $\nu_{k'}$  with  $k \neq k' \in [M_2]$ , we have

$$\left| - (\boldsymbol{\nu}_k^{\top}, 0^{\top}) \eta \cdot \sum_{b=1}^{t_0+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_C^{(b)}} (\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top} \right| \le \frac{\eta(t_0+1)\beta}{M_1 M_2} \sqrt{\frac{\log B}{B}}, \quad (113)$$

$$\left| - (\boldsymbol{\mu}_{j}^{\top}, \boldsymbol{0}^{\top})\boldsymbol{\eta} \cdot \sum_{b=1}^{t_{0}+1} \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\boldsymbol{\Psi}^{(b)}; \boldsymbol{P}^{n}, \boldsymbol{z}^{n})}{\partial \boldsymbol{W}_{C}^{(b)}} (\boldsymbol{\nu}_{k}^{\top}, \boldsymbol{0}^{\top})^{\top} \right| \leq \frac{\eta(t_{0}+1)\beta}{M_{1}M_{2}} \sqrt{\frac{\log B}{B}}, \quad (114)$$

$$\left| - (\boldsymbol{\nu}_k^{\top}, 0^{\top})\eta \cdot \sum_{b=1}^{t_0+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_C^{(b)}} (\boldsymbol{\nu}_k^{\top}, 0^{\top})^{\top} \right| \le \frac{\eta(t_0+1)\beta}{M_2} \sqrt{\frac{\log B}{B}}, \quad (115)$$

$$\left| - (\boldsymbol{\nu}_{k'}^{\top}, 0^{\top})\eta \cdot \sum_{b=1}^{t_0+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_C^{(b)}} (\boldsymbol{\nu}_k^{\top}, 0^{\top})^{\top} \right| \le \frac{\eta(t_0+1)\beta}{M_2^2} \sqrt{\frac{\log B}{B}}, \quad (116)$$

Then, we complete the induction.

(b) We then characterize the gradient updates of  $W_B$ . We have that when  $F(\Psi; \mathbf{P}^n) \in (-1, 1)$  for some  $n \in [N]$ ,

$$\frac{\partial \ell(\Psi; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_B} = -z^n \sum_{i=1}^{l+1} G_{i,l+1}^n(\boldsymbol{w}) y_i \boldsymbol{W}_C \boldsymbol{p}_{query} \boldsymbol{p}_i^\top.$$
(117)

We also use induction to complete the proof. Similar to the analysis of  $W_C$ , we have that when t = 0,

$$-(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})\eta \cdot \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\Psi^{(0)}; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{W}_{B}^{(0)}} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top}$$

$$= (\boldsymbol{\mu}_{j}^{\top}, 0^{\top}) \frac{\eta}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} z^{n} \sum_{i=1}^{l} G_{i,l+1}^{n}(\boldsymbol{w}^{(0)}) y_{i}^{n} \boldsymbol{W}_{C}^{(0)} \boldsymbol{p}_{query}^{n} \boldsymbol{p}_{i}^{n\top} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top}$$

$$\geq \eta \cdot \frac{1}{2M_{1}} (1 - p_{a}) \sum_{i=1}^{l} G_{i,l+1}^{n}(\boldsymbol{w}^{(0)}) \beta - \eta \cdot \frac{1}{2M_{1}} \sum_{i=1}^{l} G_{i,l+1}^{n}(\boldsymbol{w}^{(0)}) \beta p_{a} \sqrt{\frac{\log B}{B}}$$

$$\geq \eta \frac{1}{4M_{1}} (1 - p_{a}) \beta (1 - \xi l).$$
(118)

For  $\mu_{j'}, j' \neq j$ , that does not form a task in the training stage, we have

$$-(\boldsymbol{\mu}_{j'}^{\top}, 0^{\top})\eta \cdot \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(0)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_B^{(0)}} (\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top} = 0.$$
(119)

For  $\mu_{j''}, j'' \neq j$ , that forms a task in the training stage, we have

$$-(\boldsymbol{\mu}_{j''}^{\top}, 0^{\top})\eta \cdot \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(0)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_B^{(0)}} (\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top} \le -\eta \cdot \frac{1}{4M_1} (1 - p_a)\beta (1 - \xi l).$$
(120)

For  $\boldsymbol{\nu}_k, \boldsymbol{\nu}_{k'}$  with  $k \neq k' \in [M_2]$ , we have

$$-(\boldsymbol{\nu}_{k}^{\top}, 0^{\top})\eta \cdot \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\Psi^{(0)}; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{W}_{B}^{(0)}} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top} \Big| \leq \frac{\eta \beta}{M_{1}M_{2}} \sqrt{\frac{\log B}{B}},$$
(121)

$$-(\boldsymbol{\mu}_{j}^{\top}, \boldsymbol{0}^{\top})\boldsymbol{\eta} \cdot \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\boldsymbol{\Psi}^{(0)}; \boldsymbol{P}^{n}, \boldsymbol{z}^{n})}{\partial \boldsymbol{W}_{B}^{(0)}} (\boldsymbol{\nu}_{k}^{\top}, \boldsymbol{0}^{\top})^{\top} \Big| \leq \frac{\eta \beta}{M_{1}M_{2}} \sqrt{\frac{\log B}{B}}.$$
 (122)

$$\left| - (\boldsymbol{\nu}_{k}^{\top}, 0^{\top})\eta \cdot \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\Psi^{(0)}; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{W}_{B}^{(0)}} (\boldsymbol{\nu}_{k}^{\top}, 0^{\top})^{\top} \right| \leq \frac{\eta\beta}{M_{2}} \sqrt{\frac{\log B}{B}}.$$
 (123)

$$\left| - (\boldsymbol{\nu}_{k'}^{\top}, 0^{\top})\eta \cdot \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(0)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_B^{(0)}} (\boldsymbol{\nu}_k^{\top}, 0^{\top})^{\top} \right| \le \frac{\eta\beta}{M_2^2} \sqrt{\frac{\log B}{B}}.$$
 (124)

We also have that for any  $s \in [V]$ ,

$$\left| (\boldsymbol{\nu}_s^{*\top}, 0^{\top}) \eta \cdot \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(0)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_B^{(0)}} (\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top} \right| \le \frac{\eta \beta p_a \kappa_*}{M_1 V} \cdot \sqrt{\frac{\log B}{B}}, \quad (125)$$

Therefore, the conclusions hold when t = 0. Suppose that the conclusions also hold when  $t = t_0$ . Then, when  $t = t_0 + 1$ , we have

$$-(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})\eta \cdot \sum_{b=1}^{t_{0}+1} \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{W}_{B}^{(b)}} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})^{\top}$$

$$\gtrsim \eta \cdot \sum_{c=1}^{t_{0}+1} \frac{1}{2M_{1}} (1-p_{a}) \sum_{i=1}^{l} G_{i,l+1}^{n}(\boldsymbol{w}^{(t_{0})})\beta$$

$$\gtrsim \eta (t_{0}+1) \frac{1}{M_{1}} (1-p_{a})\beta.$$
(126)

For  $\mu_{j'}$ ,  $j' \neq j$ , that does not form a task in the training set, we have

$$-(\boldsymbol{\mu}_{j'}^{\top}, 0^{\top})\boldsymbol{\eta} \cdot \sum_{b=1}^{t_0+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_B^{(b)}} (\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top} = 0$$
(127)

For  $\mu_{j''}, j'' \neq j$ , that forms a task in the training set, we have

$$-(\boldsymbol{\mu}_{j''}^{\top}, 0^{\top})\eta \cdot \sum_{b=1}^{t_0+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_B^{(b)}} (\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top}$$

$$\leq -\eta (t_0+1) \frac{1}{M_1} (1-p_a)\beta.$$
(128)

For  $\boldsymbol{\nu}_k, \boldsymbol{\nu}_{k'}$  with  $k \neq k' \in [M_2]$ , we have

$$\left| - (\boldsymbol{\nu}_k^{\top}, 0^{\top})\eta \cdot \sum_{b=1}^{t_0+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_B^{(b)}} (\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top} \right| \le \frac{\eta(t_0+1)\beta}{M_1 M_2} \sqrt{\frac{\log B}{B}}, \quad (129)$$

$$\left| - (\boldsymbol{\mu}_{j}^{\top}, \boldsymbol{0}^{\top})\boldsymbol{\eta} \cdot \sum_{b=1}^{t_{0}+1} \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\boldsymbol{\Psi}^{(b)}; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{W}_{B}^{(b)}} (\boldsymbol{\nu}_{k}^{\top}, \boldsymbol{0}^{\top})^{\top} \right| \leq \frac{\eta(t_{0}+1)\beta}{M_{1}M_{2}} \sqrt{\frac{\log B}{B}}.$$
(130)

$$\left| - (\boldsymbol{\nu}_{k}^{\top}, 0^{\top})\eta \cdot \sum_{b=1}^{t_{0}+1} \frac{1}{|\mathcal{B}_{b}|} \sum_{n \in \mathcal{B}_{b}} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{W}_{B}^{(b)}} (\boldsymbol{\nu}_{k}^{\top}, 0^{\top})^{\top} \right| \leq \frac{\eta(t_{0}+1)\beta}{M_{2}} \sqrt{\frac{\log B}{B}}.$$
 (131)

$$\left| - (\boldsymbol{\nu}_{k'}^{\top}, 0^{\top})\eta \cdot \sum_{b=1}^{t_0+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_B^{(b)}} (\boldsymbol{\nu}_k^{\top}, 0^{\top})^{\top} \right| \le \frac{\eta(t_0+1)\beta}{M_2^2} \sqrt{\frac{\log B}{B}}.$$
 (132)

We also have that for any  $s \in [V]$ ,

$$\left| (\boldsymbol{\nu}_s^{*\top}, 0^{\top}) \eta \cdot \sum_{b=1}^{t_0+1} \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\ell(\Psi^{(b)}; \boldsymbol{P}^n, z^n)}{\partial \boldsymbol{W}_B^{(b)}} (\boldsymbol{\mu}_j^{\top}, 0^{\top})^{\top} \right| \le \frac{\eta \beta(t_0+1) p_a \kappa_*}{M_1 V} \cdot \sqrt{\frac{\log B}{B}}, \quad (133)$$

# F.3. Proof of Lemma 4

**Proof** When  $F(\Psi; \mathbf{P}^n) \in (-1, 1)$  for some  $n \in [N]$ ,

$$\begin{split} &\frac{\partial \ell(\Psi; \mathbf{P}^{n}, z^{n})}{\partial w} \\ &= -z^{n} \sum_{i=1}^{l} y_{i}^{n} p_{i}^{n^{\top}} \mathbf{W}_{B}^{\top} \mathbf{W}_{C} p_{query}^{n} \frac{\partial G_{i,l+1}^{n}(w)}{\partial w} \\ &= -z^{n} \sum_{i=1}^{l} y_{i}^{n} p_{i}^{n^{\top}} \mathbf{W}_{B}^{\top} \mathbf{W}_{C} p_{query}^{n} \frac{\partial \prod_{j=i+1}^{l+1} (1 - \sigma(w^{\top} p_{j}^{n})) \sigma(w^{\top} p_{i}^{n})}{\partial w} \\ &= -z^{n} \sum_{i=1}^{l} y_{i}^{n} p_{i}^{n^{\top}} \mathbf{W}_{B}^{\top} \mathbf{W}_{C} p_{query}^{n} (\sum_{s=i+1}^{l+1} \prod_{j=i+1, j \neq s}^{l+1} (1 - \sigma(w^{\top} p_{j}^{n})) \mathbb{1}[j < l+1]) \sigma(w^{\top} p_{i}^{n}) \\ &\cdot \frac{\partial (1 - \sigma(w^{\top} p_{s}^{n}))}{\partial w} + \prod_{j=i+1}^{l+1} (1 - \sigma(w^{\top} p_{j}^{n})) \frac{\partial \sigma(w^{\top} p_{i}^{n})}{\partial w}) \\ &= -z^{n} \sum_{i=1}^{l} y_{i}^{n} p_{i}^{n^{\top}} \mathbf{W}_{B}^{\top} \mathbf{W}_{C} p_{query}^{n} (\sum_{s=i+1}^{l+1} \prod_{j=i+1, j \neq s}^{l+1} (1 - \sigma(w^{\top} p_{j}^{n})) \mathbb{1}[j < l+1]) \sigma(w^{\top} p_{i}^{n}) \\ &\cdot (1 - \sigma(w^{\top} p_{s}^{n})) \sigma(w^{\top} p_{s}^{n}) (-p_{s}^{n}) + \prod_{j=i+1}^{l+1} (1 - \sigma(w^{\top} p_{j}^{n})) (1 - \sigma(w^{\top} p_{i}^{n})) \sigma(w^{\top} p_{i}^{n}) p_{i}^{n}) \\ &= z^{n} \sum_{i=1}^{l} y_{i}^{n} p_{i}^{n^{\top}} \mathbf{W}_{B}^{\top} \mathbf{W}_{C} p_{query}^{n} (\sum_{s=i+1}^{l+1} \prod_{j=i+1}^{l+1} (1 - \sigma(w^{\top} p_{j}^{n})) \mathbb{1}[j < l+1]) \cdot \sigma(w^{\top} p_{s}^{n}) \\ &\cdot \sigma(w^{\top} p_{i}^{n}) p_{s}^{n} - \prod_{j=i}^{l+1} (1 - \sigma(w^{\top} p_{j}^{n})) \sigma(w^{\top} p_{i}^{n}) p_{i}^{n}) \\ &= z^{n} \sum_{i=1}^{l} y_{i}^{n} p_{i}^{n^{\top}} \mathbf{W}_{B}^{\top} \mathbf{W}_{C} p_{query}^{n} G_{i,l+1}^{n}(w) (\sum_{s=i+1}^{l+1} \sigma(w^{\top} p_{s}^{n}) p_{s}^{n} - (1 - \sigma(w^{\top} p_{i}^{n})) p_{i}^{n}). \end{split}$$

$$(134)$$

When t = 1, we have

$$\boldsymbol{w}^{(1)} = \boldsymbol{w}^{(0)} - \frac{\eta}{|\mathcal{B}_{1}|} \sum_{n \in \mathcal{B}_{1}} \frac{\partial \ell(\Psi; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{w}^{(0)}}$$
  
$$= \boldsymbol{w}^{(0)} - \frac{\eta}{|\mathcal{B}_{1}|} \sum_{n \in \mathcal{B}_{1}} z^{n} \sum_{i=1}^{l} y^{n}_{i} \boldsymbol{p}^{n^{\top}}_{i} \boldsymbol{W}^{(0)^{\top}}_{B} \boldsymbol{W}^{(0)}_{C} \boldsymbol{p}^{n}_{query} G^{n}_{i,l+1}(\boldsymbol{w}^{(0)}) \qquad (135)$$
  
$$\cdot (\sum_{s=i+1}^{l+1} \sigma(\boldsymbol{w}^{(0)^{\top}} \boldsymbol{p}^{n}_{s}) \boldsymbol{p}^{n}_{s} - (1 - \sigma(\boldsymbol{w}^{(0)^{\top}} \boldsymbol{p}^{n}_{i})) \boldsymbol{p}^{n}_{i})$$

For  $p_i^n$  that contains a  $v_s^*$ , the corresponding  $y_i^n$  is consistent with  $z^n$  with a probability of 1/2. Given Hoeffding's bound (25), this part generates a gradient update as

$$\left\|\frac{\eta}{|\mathcal{B}_{1}|}\sum_{n\in\mathcal{B}_{1}}z^{n}\sum_{1\leq i\leq l,\boldsymbol{p}_{i}^{n}\text{ does not contain any }\boldsymbol{v}_{s}^{*}}y_{i}^{n}\boldsymbol{p}_{i}^{n\top}\boldsymbol{W}_{B}^{(0)^{\top}}\boldsymbol{W}_{C}^{(0)}\boldsymbol{p}_{query}^{n}G_{i,l+1}^{n}(\boldsymbol{w}^{(0)})\right.$$

$$\left.\cdot\left(\sum_{s=i+1}^{l+1}\sigma(\boldsymbol{w}^{(0)^{\top}}\boldsymbol{p}_{s}^{n})\boldsymbol{p}_{s}^{n}-(1-\sigma(\boldsymbol{w}^{(0)^{\top}}\boldsymbol{p}_{i}^{n}))\boldsymbol{p}_{i}^{n})\right\|$$

$$\leq\eta\sqrt{\frac{\log B}{B}}$$

$$(136)$$

by (95) and  $\sum_{i=1}^{l} \frac{l}{2^{l}} \leq 2$ . Then, with a high probability, for  $s \in [V]$ ,  $\xi = 1/(d+1)$ ,  $\boldsymbol{n}^{*\top} \boldsymbol{w}^{(1)}$ 

$$\leq \xi + \eta \sqrt{\frac{\log B}{B}} - \eta \beta^2 \frac{1}{|\mathcal{B}_1|} \sum_{n \in \mathcal{B}_b} \sum_{1 \leq i \leq l, p_i^n \text{ does not contain any } v_s^*} G_{i,l+1}^n(w^{(0)})$$

$$\cdot (\sum_{s=i+1}^{l+1} \sigma(w^{(0)^\top} p_s^n) v_s^{*\top} p_s^n - (1 - \sigma(w^{(0)^\top} p_i^n)) v_s^{*\top} p_i^n)$$

$$\leq \xi + \eta \sqrt{\frac{\log B}{B}} - \eta \beta^2 \sum_{i=1}^l \frac{1}{2^{l+2-i}V} \cdot \kappa_a \sum_{s=i+1}^{l+1} \frac{1}{2}(1 - p_a)$$

$$= \xi + \eta \sqrt{\frac{\log B}{B}} - \eta \beta^2 \sum_{i=1}^l \frac{\kappa_a}{2^{l+2-i}V} \cdot \frac{(1 - p_a)(l - i + 1)}{2}$$

$$= \xi + \eta \sqrt{\frac{\log B}{B}} - \eta \beta^2 \cdot \sum_{i=1}^l \frac{\kappa_a i}{2^{2+i}V} \cdot \frac{1 - p_a}{2}$$

$$\leq \xi + \eta \sqrt{\frac{\log B}{B}} - \frac{\eta \beta^2 \kappa_a (1 - p_a)}{V}$$

$$\leq -\frac{\eta \beta^2 \kappa_a (1 - p_a)}{V}.$$
(137)

The second step comes from (95) and the fact that

$$\Pr\left(\left|\frac{1}{l|\mathcal{B}_{1}|}\sum_{n\in\mathcal{B}_{1}}\sum_{i=1}^{l}\mathbb{1}[\boldsymbol{p}_{i}^{n} \text{ does not contain any } \boldsymbol{v}_{s}^{*}]G_{i,l+1}^{n}(\boldsymbol{w}^{(0)})\sum_{s=i+1}^{l+1}\sigma(\boldsymbol{w}^{(0)^{\top}}\boldsymbol{p}_{s}^{n})\right.\\\left.\cdot\boldsymbol{v}_{s}^{*^{\top}}\boldsymbol{p}_{s}^{n}-(1-p_{a})\mathbb{E}[\frac{1}{l|\mathcal{B}_{1}|}\sum_{n\in\mathcal{B}_{1}}\sum_{i=1}^{l}G_{i,l+1}^{n}(\boldsymbol{w}^{(0)})\sum_{s=i+1}^{l+1}\sigma(\boldsymbol{w}^{(0)^{\top}}\boldsymbol{p}_{s}^{n})\boldsymbol{v}_{s}^{*^{\top}}\boldsymbol{p}_{s}^{n}]\right|\\\geq c\cdot(1-p_{a})\mathbb{E}[\frac{1}{l|\mathcal{B}_{1}|}\sum_{n\in\mathcal{B}_{1}}\sum_{i=1}^{l}G_{i,l+1}^{n}(\boldsymbol{w}^{(0)})\sum_{s=i+1}^{l+1}\sigma(\boldsymbol{w}^{(0)^{\top}}\boldsymbol{p}_{s}^{n})\boldsymbol{v}_{s}^{*^{\top}}\boldsymbol{p}_{s}^{n}]\right)\\\lesssim e^{-lB(1-p_{a})^{2}c^{2}}\leq\epsilon$$
(138)

for some  $c \in (0, 1)$ , and

$$Bl \ge (1 - p_a)^{-2} \log \epsilon^{-1} \tag{139}$$

by Lemma 2 since  $p_i^n$  contains  $v_s^*$  with a probability of  $p_a/V$ . The last step holds with a high probability if

$$B \gtrsim \beta^{-4} \kappa_a^{-2} (1 - p_a)^{-2} V^2 \log \epsilon^{-1}.$$
 (140)

We can also derive that for any  $j \in [M_1]$ ,

$$(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})\boldsymbol{w}^{(1)}$$

$$\leq \xi + \frac{\eta}{M_{1}}\sqrt{\frac{\log B}{B}} - \frac{\eta\beta^{2}}{|\mathcal{B}_{1}|} \sum_{n\in\mathcal{B}_{b}} \sum_{1\leq i\leq l, \boldsymbol{p}_{i}^{n} \text{ does not contain any } \boldsymbol{v}_{s}^{*}} G_{i,l+1}^{n}(\boldsymbol{w}^{(0)})(\sum_{s=i+1}^{l+1} \sigma(\boldsymbol{w}^{(0)^{\top}}\boldsymbol{p}_{s}^{n}))$$

$$\cdot (\boldsymbol{\mu}_{j}^{\top}, 0^{\top})\boldsymbol{p}_{s}^{n} - (1 - \sigma(\boldsymbol{w}^{(0)^{\top}}\boldsymbol{p}_{i}^{n}))(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})\boldsymbol{p}_{i}^{n})$$

$$\lesssim \xi + \frac{\eta}{M_{1}}\sqrt{\frac{\log B}{B}} - \eta\beta^{2}\sum_{i=1}^{l} \frac{1}{2^{l+2-i}} \cdot \frac{(1 - p_{a})}{2M_{1}}(l - i + 1)$$

$$\lesssim \xi + \frac{\eta}{M_{1}}\sqrt{\frac{\log B}{B}} - \frac{\eta(1 - p_{a})\beta^{2}}{M_{1}}$$

$$\lesssim - \frac{\eta(1 - p_{a})\beta^{2}}{M_{1}}.$$
(141)

The second step of (141) comes from the fact that

$$\Pr\left(\left|\frac{1}{l|\mathcal{B}_{1}|}\sum_{n\in\mathcal{B}_{1}}\sum_{i=1}^{l}\mathbb{1}[\boldsymbol{p}_{i}^{n} \text{ does not contain any } \boldsymbol{v}_{s}^{*}]G_{i,l+1}^{n}(\boldsymbol{w}^{(0)})\sum_{s=i+1}^{l+1}\sigma(\boldsymbol{w}^{(0)^{\top}}\boldsymbol{p}_{s}^{n})\right.-(1-p_{a})\mathbb{E}[\frac{1}{l|\mathcal{B}_{1}|}\sum_{n\in\mathcal{B}_{1}}\sum_{i=1}^{l}G_{i,l+1}^{n}(\boldsymbol{w}^{(0)})\sum_{s=i+1}^{l+1}\sigma(\boldsymbol{w}^{(0)^{\top}}\boldsymbol{p}_{s}^{n})]\right]\geq c\cdot(1-p_{a})\mathbb{E}[\frac{1}{l|\mathcal{B}_{1}|}\sum_{n\in\mathcal{B}_{1}}\sum_{i=1}^{l}G_{i,l+1}^{n}(\boldsymbol{w}^{(0)})\sum_{s=i+1}^{l+1}\sigma(\boldsymbol{w}^{(0)^{\top}}\boldsymbol{p}_{s}^{n})]\right)\\\lesssim e^{-lB(1-p_{a})^{2}c^{2}}\leq M_{1}^{-C}$$

$$(142)$$

for some  $c \in (0, 1)$ , C > 1, and

$$Bl \ge (1 - p_a)^{-2} \log \epsilon^{-1} \tag{143}$$

by Lemma 2 since  $p_i^n$  does not contain any  $v_s^*$  with a probability of  $1 - p_a$ .

The last step of (141) holds if  $B \gtrsim \beta^{-4}$  and  $\xi \lesssim \frac{1}{M_1}$ . Similarly, we also have

$$(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})\boldsymbol{w}^{(1)}$$

$$\geq -\xi - \frac{\eta}{M_{1}}\sqrt{\frac{\log B}{B}} - \frac{\eta\beta^{2}}{|\mathcal{B}_{1}|} \sum_{n \in \mathcal{B}_{b}} \sum_{1 \leq i \leq l, \boldsymbol{p}_{i}^{n} \text{ does not contain any } \boldsymbol{v}_{s}^{*}} G_{i,l+1}^{n}(\boldsymbol{w}^{(0)})$$

$$\cdot (\sum_{s=i+1}^{l+1} \sigma(\boldsymbol{w}^{(0)^{\top}}\boldsymbol{p}_{s}^{n})(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})\boldsymbol{p}_{s}^{n}$$

$$\gtrsim - \frac{\eta(1-p_{a})\beta^{2}}{M_{1}}.$$
(144)

Hence, the conclusion holds when t = 1. Meanwhile, for any  $k \in [M_2]$ ,

$$(\boldsymbol{\nu}_k^{\top}, 0^{\top})\boldsymbol{w}^{(1)} \leq \xi + \frac{\eta}{M_2}\sqrt{\frac{\log B}{B}}.$$
(145)

Suppose that the conclusion holds when  $t = t_0$  for  $t_0 \leq \min\{\eta^{-1}\beta^{-2}\kappa_a^{-1}(1-p_a)^{-1}V, \eta^{-1}M_1^{\frac{2}{3}}\beta^{-\frac{2}{3}}\kappa_a^{-\frac{1}{3}}(1-p_a)^{-1}V^{\frac{1}{3}}\}$ . Then, when  $t = t_0 + 1$ , we have that for  $\boldsymbol{p}_s^n$  that does not contain any  $\boldsymbol{v}_s^*, s \in [V]$ 

$$-\frac{\eta(1-p_a)\beta^2 t_0}{M_1} - \sum_{i=1}^{t_0} i^2 \cdot \left(\frac{\eta^3(1-p_a)^3\beta^2}{M_1^3}\right) \lesssim \boldsymbol{w}^{(t_0)^{\top}} \boldsymbol{p}_s^n \lesssim t_0 \cdot \left(-\frac{\eta\beta^2}{M_1} + \frac{\eta}{M_2}\sqrt{\frac{\log B}{B}} + \xi\right) < 0.$$
(146)

For another  $p_r^n$ ,  $r \neq s$ , that contains a  $v_s^*$ ,  $s \in [V]$ ,

$$\boldsymbol{w}^{(t_0)^{\top}}\boldsymbol{p}_r^n \lesssim t_0 \cdot (0 - \eta \beta^2 \kappa_a (1 - p_a)) < \boldsymbol{w}^{(t_0)^{\top}} \boldsymbol{p}_s^n < 0.$$
(147)

Then, with a high probability, we have for any  $s \in [V]$ ,

$$\boldsymbol{v}_{s}^{*^{\top}} \boldsymbol{w}^{(t)} = \boldsymbol{v}_{s}^{*^{\top}} (\boldsymbol{w}^{(t-1)} - \eta \frac{\partial \ell(\Psi; \boldsymbol{P}^{n}, z^{n})}{\partial \boldsymbol{w}})$$

$$\leq -\eta \beta^{2} t_{0} \kappa_{a} (1 - p_{a}) - \eta \sum_{i=1}^{t_{0}-1} i^{2} (\frac{\eta^{2} (1 - p_{a})^{3} \beta^{2}}{M_{1}^{2}}) \kappa_{a} - \eta \frac{z^{n}}{|\mathcal{B}_{1}|} \sum_{n \in \mathcal{B}_{b}} \sum_{i=1}^{l} y_{i}^{n} (\beta^{2} + \frac{\eta^{2} t_{0}^{2} (1 - p_{a})^{2} \beta^{2}}{M_{1}^{2}}) G_{i,l+1}^{n} (\boldsymbol{w}^{(t_{0})}) (\sum_{s=i+1}^{l+1} \sigma (\boldsymbol{w}^{(t_{0})^{\top}} \boldsymbol{p}_{s}^{n}) \boldsymbol{v}_{s}^{*^{\top}} \boldsymbol{p}_{s}^{n} - (1 - \sigma (\boldsymbol{w}^{(t_{0})^{\top}} \boldsymbol{p}_{i}^{n})) \boldsymbol{v}_{s}^{*^{\top}} \boldsymbol{p}_{i}^{n}),$$

$$(148)$$

where the last step is by (109) and (126). Following our proof idea in the case of t = 1, we have that for  $p_i^n$  that contains a  $v_s^*$ ,  $s \in [V]$ , the corresponding  $y_i^n$  has a probability of 1/2 to be both binary

labels. Then, by Hoeffding' bound (25), we have

$$\left\|\frac{\eta}{|\mathcal{B}_{1}|}\sum_{n\in\mathcal{B}_{1}}z^{n}\sum_{1\leq i\leq l,\boldsymbol{p}_{i}^{n} \text{ contains } \boldsymbol{v}_{s}^{*}}y_{i}^{n}\boldsymbol{p}_{i}^{n\top}\boldsymbol{W}_{B}^{(t_{0})^{\top}}\boldsymbol{W}_{C}^{(t_{0})}\boldsymbol{p}_{query}^{n}G_{i,l+1}^{n}(\boldsymbol{w}^{(t_{0})})\right)$$

$$\cdot\left(\sum_{s=i+1}^{l+1}\sigma(\boldsymbol{w}^{(t_{0})^{\top}}\boldsymbol{p}_{s}^{n})\boldsymbol{p}_{s}^{n}-(1-\sigma(\boldsymbol{w}^{(t_{0})^{\top}}\boldsymbol{p}_{i}^{n}))\boldsymbol{p}_{i}^{n})\right\|$$

$$\leq\eta\sqrt{\frac{\log B}{B}}.$$
(149)

Then, with a high probability,

$$\begin{split} &\eta \frac{z^{n}}{|\mathcal{B}_{1}|} \sum_{n \in \mathcal{B}_{b}} \sum_{i=1}^{l} y_{i}^{n} (\beta^{2} + \frac{\eta^{2} t_{0}^{2} (1-p_{a})^{2} \beta^{2}}{M_{1}^{2}}) G_{i,l+1}^{n} (\boldsymbol{w}^{(t_{0})}) (\sum_{s=i+1}^{l+1} \sigma (\boldsymbol{w}^{(t_{0})^{\top}} \boldsymbol{p}_{s}^{n}) \boldsymbol{v}_{s}^{*^{\top}} \boldsymbol{p}_{s}^{n} \\ & \cdot -(1 - \sigma (\boldsymbol{w}^{(t_{0})^{\top}} \boldsymbol{p}_{i}^{n})) \boldsymbol{v}_{s}^{*^{\top}} \boldsymbol{p}_{i}^{n}) \\ & \gtrsim -\eta \sqrt{\frac{\log B}{B}} + \eta \frac{z^{n}}{|\mathcal{B}_{1}|} \sum_{n \in \mathcal{B}_{b}} \sum_{\boldsymbol{p}_{i}^{n} \operatorname{does not contain} \boldsymbol{v}_{s}^{*} z^{n} \boldsymbol{y}_{i}^{n=1}} y_{i}^{n} (\beta^{2} + \frac{\eta^{2} t_{0}^{2} (1-p_{a})^{2} \beta^{2}}{M_{1}^{2}}) \\ & \cdot G_{i,l+1}^{n} (\boldsymbol{w}^{(t_{0})}) (\sum_{s=i+1}^{l+1} \sigma (\boldsymbol{w}^{(t_{0})^{\top}} \boldsymbol{p}_{s}^{n}) \boldsymbol{v}_{s}^{*^{\top}} \boldsymbol{p}_{s}^{n} - (1 - \sigma (\boldsymbol{w}^{(t_{0})^{\top}} \boldsymbol{p}_{i}^{n})) \boldsymbol{v}_{s}^{*^{\top}} \boldsymbol{p}_{i}^{n}) \\ & = -\eta \sqrt{\frac{\log B}{B}} + \eta \frac{1}{|\mathcal{B}_{1}|} \sum_{n \in \mathcal{B}_{b}} \sum_{\boldsymbol{p}_{i}^{n} \operatorname{does not contain} \boldsymbol{v}_{s}^{*}} (\beta^{2} + \frac{\eta^{2} t_{0}^{2} (1-p_{a})^{2} \beta^{2}}{M_{1}^{2}}) G_{i,l+1}^{n} (\boldsymbol{w}^{(t_{0})}) \\ & \cdot \sum_{s=i+1}^{l+1} \sigma (\boldsymbol{w}^{(t_{0})^{\top}} \boldsymbol{p}_{s}^{n}) \boldsymbol{v}_{s}^{*^{\top}} \boldsymbol{p}_{s}^{n} \qquad (150) \\ & \gtrsim -\eta \sqrt{\frac{\log B}{B}} + \eta \frac{1}{|\mathcal{B}_{1}|} \sum_{n \in \mathcal{B}_{b}} \sum_{\boldsymbol{p}_{i}^{n} \operatorname{does not contain} \boldsymbol{v}_{s}^{*}} (\beta^{2} + \frac{\eta^{2} t_{0}^{2} (1-p_{a})^{2} \beta^{2}}{M_{1}^{2}}) G_{i,l+1}^{n} (\boldsymbol{w}^{(t_{0})}) \\ & \cdot (l - i + 1) \frac{\kappa_{a}}{V} \\ & \geq -\eta \sqrt{\frac{\log B}{B}} + \eta (\beta^{2} + \frac{\eta^{2} t_{0}^{2} (1-p_{a})^{2} \beta^{2}}{M_{1}^{2}}) \mathbb{E} \left[ \sum_{i=1}^{l} G_{i,l+1}^{n} (\boldsymbol{w}^{(t_{0})}) (l - i + 1) \frac{\kappa_{a} (1-p_{a})}{V} \right] \\ & \geq -\eta \sqrt{\frac{\log B}{B}} + \eta (\beta^{2} + \frac{\eta^{2} t_{0}^{2} (1-p_{a})^{2} \beta^{2}}{M_{1}^{2}}) \mathbb{E} \left[ \sum_{i=1}^{l} G_{i,l+1}^{n} (\boldsymbol{w}^{(t_{0})}) \frac{\kappa_{a} (1-p_{a})}{V} \right] \\ & \geq -\eta \sqrt{\frac{\log B}{B}} + \eta (\beta^{2} + \frac{\eta^{2} t_{0}^{2} (1-p_{a})^{2} \beta^{2}}{M_{1}^{2}}) \frac{\kappa_{a} (1-p_{a})}{V}, \end{aligned}$$

where the fourth step follows the idea of (138) since

$$G_{i,l+1}^{n}(\boldsymbol{w}^{(t_{0})})(l-i+1) \le \Theta(1),$$
(151)

for any  $i \in [l]$  and  $n \in \mathcal{B}_b$ . The last step of (150) follows from

$$\sum_{i=1}^{l} G_{i,l+1}^{n}(\boldsymbol{w}^{(t_{0})}) = 1 - \sigma(\boldsymbol{w}^{(t_{0})^{\top}}\boldsymbol{p}_{query}^{n}) - \prod_{i=1}^{l+1} (1 - \sigma(\boldsymbol{w}^{(t_{0})^{\top}}\boldsymbol{p}_{i}^{n})) \ge \frac{1}{4}, \quad (152)$$

since

$$\sigma(\boldsymbol{w}^{(t_0)^{\top}}\boldsymbol{p}_{query}^n) < \sigma(0) = \frac{1}{2},$$
(153)

by (146), and with a high probability,

$$\prod_{i=1}^{l+1} (1 - \sigma(\boldsymbol{w}^{(t_0)^{\top}} \boldsymbol{p}_i^n)) \leq \prod_{\boldsymbol{p}_i^n \text{ does not contain any } \boldsymbol{v}_s^*} (1 - \sigma(\boldsymbol{w}^{(t_0)^{\top}} \boldsymbol{p}_i^n)) \\
\leq (1 - \frac{1}{1 + e^{-\frac{V}{\kappa_a M_1}}})^{l(1-p_a)} \\
\leq \frac{1}{4},$$
(154)

where the last step holds if

$$l \gtrsim (1 - p_a)^{-1} \log M_1.$$
 (155)

The second step of (154) comes from (146) and

$$\Pr\left(\left|\frac{1}{l}\sum_{i=1}^{l}\mathbb{I}[\boldsymbol{p}_{i}^{n} \text{ does not contain } \boldsymbol{v}_{s}^{*}] - (1-p_{a})\right| \geq c \cdot (1-p_{a})\right) \lesssim e^{-l(1-p_{a})c^{2}} \leq M_{1}^{-C} \quad (156)$$

by Lemma 1 for some  $c \in (0, 1), C > 1$ , and

$$l \ge (1 - p_a)^{-1} \log M_1. \tag{157}$$

Then, by plugging (150) into (148), we have

$$\begin{aligned} \boldsymbol{v}_{s}^{*^{\top}} \boldsymbol{w}^{(t_{0}+1)} \\ &\leq -\frac{\eta \beta^{2} t_{0} \kappa_{a}(1-p_{a})}{V} - \eta \sum_{i=1}^{t_{0}-1} i^{2} (\frac{\eta^{2}(1-p_{a})^{3}\beta^{2}}{M_{1}^{2}}) \frac{\kappa_{a}}{V} + \eta \sqrt{\frac{\log B}{B}} - \eta (\beta^{2} \\ &+ \frac{\eta^{2} t_{0}^{2}(1-p_{a})^{2}\beta^{2}}{M_{1}^{2}}) \cdot \frac{\kappa_{a}(1-p_{a})}{V} \\ &= -\frac{\eta \beta^{2}(t_{0}+1)\kappa_{a}(1-p_{a})}{V} - \eta \sum_{i=1}^{t_{0}} i^{2} (\frac{\eta^{2}(1-p_{a})^{3}\beta^{2}}{M_{1}^{2}}) \frac{\kappa_{a}}{V} + \eta \sqrt{\frac{\log B}{B}} \\ &\lesssim -\frac{\eta \beta^{2}(t_{0}+1)\kappa_{a}(1-p_{a})}{V} - \eta \sum_{i=1}^{t_{0}} i^{2} (\frac{\eta^{2}(1-p_{a})^{3}\beta^{2}}{M_{1}^{2}}) \frac{\kappa_{a}}{V}, \end{aligned}$$
(158)

where the last step holds given (140) and  $t_0 \lesssim \min\{\eta^{-1}\beta^{-2}\kappa_a^{-1}(1-p_a)^{-1}V, \eta^{-1}M_1^{\frac{2}{3}}\beta^{-\frac{2}{3}}\kappa_a^{-\frac{1}{3}}(1-p_a)^{-1}V^{\frac{1}{3}}\}$ . We can also derive that for any  $j \in [M_1]$ ,

$$\begin{split} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top}) \boldsymbol{w}^{(l)} \\ \leq \xi + \frac{\eta}{M_{1}} \sqrt{\frac{\log B}{B}} - \frac{\eta(1 - p_{a})\beta^{2}t_{0}}{M_{1}} - \sum_{i=1}^{t_{0}-1} i^{2} \cdot \left(\frac{\eta^{3}(1 - p_{a})^{3}\beta^{2}}{M_{1}^{3}}\right) - \frac{\eta}{|\mathcal{B}_{1}|} \sum_{n \in \mathcal{B}_{b}} \\ \sum_{p_{i}^{n} \text{ does not contain any } \boldsymbol{v}_{s}^{*}} \left(\beta^{2} + \frac{\eta^{2}t_{0}^{2}(1 - p_{a})^{2}\beta^{2}}{M_{1}^{2}}\right) G_{i,l+1}^{n}(\boldsymbol{w}^{(t_{0})}) \cdot \left(\sum_{s=i+1}^{l+1} \sigma(\boldsymbol{w}^{(t_{0})^{\top}}\boldsymbol{p}_{s}^{n}\right) \\ \cdot (\boldsymbol{\mu}_{j}^{\top}, 0^{\top}) \boldsymbol{p}_{s}^{n} - (1 - \sigma(\boldsymbol{w}^{(t_{0})^{\top}}\boldsymbol{p}_{i}^{n}))(\boldsymbol{\mu}_{j}^{\top}, 0^{\top}) \boldsymbol{p}_{i}^{n}) \\ \lesssim \xi + \frac{\eta}{M_{1}} \sqrt{\frac{\log B}{B}} - \frac{\eta(1 - p_{a})\beta^{2}t_{0}}{M_{1}} - \sum_{i=1}^{t_{0}-1} i^{2} \cdot \left(\frac{\eta^{3}(1 - p_{a})^{3}\beta^{2}}{M_{1}^{3}}\right) - \frac{\eta}{|\mathcal{B}_{1}|} \sum_{n \in \mathcal{B}_{b}} \sum_{i=1}^{l} (\beta^{2} + \frac{\eta^{2}t_{0}^{2}(1 - p_{a})^{2}\beta^{2}}{M_{1}^{2}}) \cdot G_{i,l+1}^{n}(\boldsymbol{w}^{(t_{0})})(l - i + 1) \cdot \frac{(1 - p_{a})}{M_{1}} \\ \lesssim \xi + \frac{\eta}{M_{1}} \sqrt{\frac{\log B}{B}} - \frac{\eta(1 - p_{a})\beta^{2}t_{0}}{M_{1}} - \sum_{i=1}^{t_{0}-1} i^{2} \cdot \left(\frac{\eta^{3}(1 - p_{a})^{3}\beta^{2}}{M_{1}^{3}}\right) - \eta \frac{(1 - p_{a})}{M_{1}} (\beta^{2} + \frac{\eta^{2}t_{0}^{2}(1 - p_{a})^{2}\beta^{2}}{M_{1}^{2}}) \\ \lesssim \xi + \frac{\eta}{M_{1}} \sqrt{\frac{\log B}{B}} - \frac{\eta(1 - p_{a})\beta^{2}(t_{0} + 1)}{M_{1}} - \sum_{i=1}^{t_{0}-1} i^{2} \cdot \left(\frac{\eta^{3}(1 - p_{a})^{3}\beta^{2}}{M_{1}^{3}}\right) \\ - \frac{\eta(1 - p_{a})}{M_{1}} \left(\frac{\eta^{2}t_{0}^{2}(1 - p_{a})^{2}\beta^{2}}{M_{1}^{2}}\right) \\ \lesssim - \frac{\eta(1 - p_{a})\beta^{2}(t_{0} + 1)}{M_{1}} - \sum_{i=1}^{t_{0}} i^{2} \cdot \left(\frac{\eta^{3}(1 - p_{a})^{3}\beta^{2}}{M_{1}^{3}}\right), \end{split}$$

where the second step of (159) follows the second step in (141) using Lemma 2. Meanwhile,

$$\begin{aligned} (\boldsymbol{\mu}_{j}^{+}, 0^{+}) \boldsymbol{w}^{(t)} \\ \gtrsim &- \xi - \frac{\eta}{M_{1}} \sqrt{\frac{\log B}{B}} - \frac{\eta (1 - p_{a}) \beta^{2} t_{0}}{M_{1}} - \sum_{i=1}^{t_{0}-1} i^{2} \cdot \left(\frac{\eta^{3} (1 - p_{a})^{3} \beta^{2}}{M_{1}^{3}}\right) - \frac{\eta}{|\mathcal{B}_{1}|} \sum_{n \in \mathcal{B}_{b}} \sum_{i=1}^{l} (\beta^{2} + \frac{\eta^{2} t_{0}^{2} (1 - p_{a})^{2} \beta^{2}}{M_{1}^{2}}) \cdot G_{i,l+1}^{n} (\boldsymbol{w}^{(t_{0})}) (l - i + 1) \cdot \frac{(1 - p_{a})}{M_{1}} \\ \gtrsim &- \xi - \frac{\eta}{M_{1}} \sqrt{\frac{\log B}{B}} - \frac{\eta (1 - p_{a}) \beta^{2} t_{0}}{M_{1}} - \sum_{i=1}^{t_{0}-1} i^{2} \cdot \left(\frac{\eta^{3} (1 - p_{a})^{3} \beta^{2}}{M_{1}^{3}}\right) - \eta \frac{(1 - p_{a})}{M_{1}} (\beta^{2} + \frac{\eta^{2} t_{0}^{2} (1 - p_{a})^{2} \beta^{2}}{M_{1}^{2}}) \\ \approx &- \frac{\eta (1 - p_{a}) \beta^{2} (t_{0} + 1)}{M_{1}} - \sum_{i=1}^{t_{0}} i^{2} \cdot \left(\frac{\eta^{3} (1 - p_{a})^{3} \beta^{2}}{M_{1}^{3}}\right), \end{aligned}$$

where the second step is by Lemma 6. Therefore, we complete the induction.

# F.4. Proof of Lemma 5

# Proof

Let

$$t_0 = \Theta(\eta^{-1}(1-p_a)^{-1}\beta^{-2}M_1).$$
(161)

(a) We first prove that for any  $s \in [V]$ ,

$$(\boldsymbol{v}_s^{*\top}, \boldsymbol{0}^{\top})\boldsymbol{w}^{(t)} \le \Theta(-\log(2+t\gamma_1))$$
(162)

for some  $\gamma_1 > 0$  by induction. When  $t = \min\{\eta^{-1}\beta^{-2}\kappa_a^{-1}(1-p_a)^{-1}V, \eta^{-1}M_1^{\frac{2}{3}}\beta^{-\frac{2}{3}}\kappa_a^{-\frac{1}{3}}(1-p_a)^{-1}V^{\frac{1}{3}}\}$ , we have

$$(\boldsymbol{v}_{s}^{*\top}, \boldsymbol{0}^{\top})\boldsymbol{w}^{(t)} \lesssim -\Theta(1) \le \Theta(-\log(2+\eta^{-1}\beta^{-\frac{2}{3}}\kappa_{a}^{-\frac{1}{3}}M_{1}^{\frac{2}{3}}(1-p_{a})^{-1}V^{\frac{1}{3}}\gamma_{1}))$$
(163)

by Lemma 4 for any  $\gamma_1 > 0$ , since that  $1 + \eta^{-1}\beta^{-\frac{2}{3}}\kappa_a^{-\frac{1}{3}}M_1^{\frac{2}{3}}(1-p_a)^{-1}V^{\frac{1}{3}}\gamma_1 \ge \Theta(1)$  and  $\gamma_1 > 0$ . Therefore, (162) holds when

$$t = \min\{\eta^{-1}\beta^{-2}\kappa_a^{-1}(1-p_a)^{-1}V, \eta^{-1}M_1^{\frac{2}{3}}\beta^{-\frac{2}{3}}\kappa_a^{-\frac{1}{3}}(1-p_a)^{-1}V^{\frac{1}{3}}\}.$$
 (164)

Suppose that when  $t \le t_2$  with  $t_2 > \min\{\eta^{-1}\beta^{-2}\kappa_a^{-1}(1-p_a)^{-1}V, \eta^{-1}M_1^{\frac{2}{3}}\beta^{-\frac{2}{3}}\kappa_a^{-\frac{1}{3}}(1-p_a)^{-1}V^{\frac{1}{3}}\}$ and  $t_2 \le t_0$ , the conclusion still holds. Then, when  $t = t_2 + 1$ , we have

$$(\boldsymbol{v}_{s}^{*\top}, 0^{\top})\boldsymbol{w}^{(t)} \lesssim -\log(2+t_{2}\gamma_{1}) - \frac{\eta(1-p_{a})\kappa_{a}}{V}(\beta^{2} + \frac{\eta^{2}t_{2}^{2}(1-p_{a})^{2}\beta^{2}}{M_{1}^{2}}) \cdot \frac{1}{1+e^{\log(2+t_{2}\gamma_{1})}}$$

$$= -\log(2+t_{2}\gamma_{1}) - \frac{\eta(1-p_{a})\kappa_{a}}{V}(\beta^{2} + \frac{\eta^{2}t_{2}^{2}(1-p_{a})^{2}\beta^{2}}{M_{1}^{2}}) \cdot (3+t_{2}\gamma_{1})^{-1}$$

$$\lesssim -\log(2+(t_{2}+1)\gamma_{1}),$$

$$(165)$$

where the last step comes from the following.

(i)

$$\frac{\eta(1-p_a)\beta^2\kappa_a}{V}(3+t_2\gamma_1)^{-1} \gtrsim \log(1+\frac{\gamma_1}{2+t_2\gamma_1}) = \log(2+(t_2+1)\gamma_1) - \log(2+t_2\gamma_1),$$
(166)

where the first step is from

$$\gamma_1 \le \eta (1 - p_a) \beta^2. \tag{167}$$

(ii)

$$\eta^{3} \frac{(1-p_{a})^{3} \kappa_{a}}{M_{1}^{2} V} \beta^{2} t_{2}^{2} (3+t_{2} \gamma_{1})^{-1} \gtrsim \log(2+(t_{2}+1)\gamma_{1}) - \log(2+t_{2} \gamma_{1}), \tag{168}$$

which comes from

$$\gamma_1 \le \frac{\eta(1-p_a)\beta^{-2}\kappa_a}{V}.$$
(169)

Therefore, (162) can be rewritten as

$$(\boldsymbol{v}_s^{*\top}, \boldsymbol{0}^{\top})\boldsymbol{w}^{(t)} \le \Theta(-\log(2 + t \cdot \eta(1 - p_a)\beta^2)),$$
(170)

when  $\kappa_a \ge V\beta^{-4}$ , so that the conclusion holds when  $t = t_2 + 1$ . Thus, the induction can be completed. We can then derive that when  $t = t_0$ , we have

$$(\boldsymbol{v}_s^{*\top}, 0^{\top}) \boldsymbol{w}^{(t_0)} \le \Theta(-\log(2 + t_0 \cdot \eta(1 - p_a)\beta^2)) \lesssim -\log(M_1),$$
 (171)

and for  $p_i$  that contains  $\nu_*$ ,

$$\sigma(\boldsymbol{p}_i^{\top} \boldsymbol{w}^{(t)}) \lesssim \frac{1}{\operatorname{poly}(M_1)}.$$
(172)

(b) We then prove that

$$(\boldsymbol{\mu}_{j}^{\top}, \boldsymbol{0}^{\top})\boldsymbol{w}^{(t)} \ge \Theta(-\log(2 + \frac{t\gamma_{2}}{M_{1}}))$$
(173)

for  $j \in [M_1]$  and some  $\gamma_2 > 0$  by induction. When  $t = \min\{\eta^{-1}\beta^{-2}\kappa_a^{-1}(1-p_a)^{-1}V, \eta^{-1}M_1^{\frac{2}{3}}\beta^{-\frac{2}{3}}\kappa_a^{-\frac{1}{3}}(1-p_a)^{-1}V^{\frac{1}{3}}\}$ , we have

$$(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})\boldsymbol{w}^{(t)} \gtrsim -\frac{1}{M_{1}} \ge \Theta(-\log(2+\eta^{-1}\beta^{-\frac{2}{3}}\kappa_{a}^{-\frac{1}{3}}M_{1}^{-\frac{1}{3}}(1-p_{a})^{-1}V^{\frac{1}{3}}\gamma_{2}))$$
(174)

by Lemma 4 for any  $\gamma_2 > 0$ , since that  $1 + \eta^{-1} \beta^{-\frac{2}{3}} \kappa_a^{-\frac{1}{3}} M_1^{-\frac{1}{3}} (1 - p_a)^{-1} V^{\frac{1}{3}} \gamma_2 \gg M_1^{-1}$  and  $\gamma_2 \ge 1$ . Therefore, (173) holds when

$$t = \min\{\eta^{-1}\beta^{-2}\kappa_a^{-1}(1-p_a)^{-1}V, \eta^{-1}M_1^{\frac{2}{3}}\beta^{-\frac{2}{3}}\kappa_a^{-\frac{1}{3}}(1-p_a)^{-1}V^{\frac{1}{3}}\}.$$
 (175)

Suppose that when  $t \le t_2$  with  $t_2 > \min\{\eta^{-1}\beta^{-2}\kappa_a^{-1}(1-p_a)^{-1}V, \eta^{-1}M_1^{\frac{2}{3}}\beta^{-\frac{2}{3}}\kappa_a^{-\frac{1}{3}}(1-p_a)^{-1}V^{\frac{1}{3}}\}$  and  $t_2 \le t_0$ , the conclusion still holds. Then, when  $t = t_2 + 1$ , we have

$$\begin{aligned} (\boldsymbol{\mu}_{j}^{\top}, 0^{\top}) \boldsymbol{w}^{(t)} \\ \gtrsim &- \log(2 + \frac{t_{2} \gamma_{2}}{M_{1}}) - \eta \frac{(1 - p_{a})}{M_{1}} (\beta^{2} + \frac{\eta^{2} t_{2}^{2} (1 - p_{a})^{2} \beta^{2}}{M_{1}^{2}}) \cdot \frac{1}{1 + e^{\log(2 + \frac{t_{2} \gamma_{2}}{M_{1}})}} \\ &= -\log(2 + \frac{t_{2} \gamma_{2}}{M_{1}}) - \eta \frac{(1 - p_{a})}{M_{1}} (\beta^{2} + \frac{\eta^{2} t_{2}^{2} (1 - p_{a})^{2} \beta^{2}}{M_{1}^{2}}) \cdot (3 + \frac{t_{2} \gamma_{2}}{M_{1}})^{-1} \\ \gtrsim &- \log(2 + \frac{(t_{2} + 1) \gamma_{2}}{M_{1}}), \end{aligned}$$
(176)

where the last step comes from the following.

(i)

$$\eta \frac{(1-p_a)}{M_1} \beta^2 (3 + \frac{t_2 \gamma_2}{M_1})^{-1} \lesssim \log(1 + \frac{\frac{\gamma_2}{M_1}}{2 + \frac{t_2 \gamma_2}{M_1}}) = \log(2 + \frac{(t_2+1)\gamma_2}{M_1}) - \log(2 + \frac{t_2 \gamma_2}{M_1}),$$
(177)

where the first step is from

$$\gamma_2 \ge \eta (1 - p_a) \beta^2. \tag{178}$$

(ii)

$$\eta^{3} \frac{(1-p_{a})^{3}}{M_{1}^{3}} \beta^{2} t_{2}^{2} (3 + \frac{t_{2} \gamma_{2}}{M_{1}})^{-1} \lesssim \log(2 + \frac{(t_{2}+1)\gamma_{2}}{M_{1}}) - \log(2 + \frac{t_{2} \gamma_{2}}{M_{1}}), \tag{179}$$

which comes from

$$\gamma_2 \ge \eta (1 - p_a) \beta^{-2}. \tag{180}$$

Therefore, (173) can be rewritten as

$$(\boldsymbol{\mu}_{j}^{\top}, \boldsymbol{0}^{\top})\boldsymbol{w}^{(t)} \ge \Theta(-\log(2 + t \cdot \frac{\eta(1 - p_{a})\beta^{2}}{M_{1}})),$$
(181)

so that the conclusion holds when  $t = t_2 + 1$ . Thus, the induction can be completed. We can then derive that when  $t = t_0$ , we have

$$(\boldsymbol{\mu}_{j}^{\top}, 0^{\top})\boldsymbol{w}^{(t_{0})} \ge \Theta(-\log(2 + t_{0} \cdot \frac{\eta(1 - p_{a})\beta^{2}}{M_{1}})) \ge -\log(3) \ge -\Theta(1),$$
(182)

and for  $p_i$  that does not contain  $\nu_*$ ,

$$\sigma(\boldsymbol{p}_i^{\top} \boldsymbol{w}^{(t)}) \gtrsim \Theta(1).$$
(183)

# F.5. Proof of Lemma 6

**Proof** Given a prompt P defined in (2) with  $(x_1, x_2, \dots, x_l, x_{query})$ , let  $x_{l+1} = x_{query}$ . Define

$$\hat{P}^{i} = \begin{pmatrix}
x_{i+1} & x_{i+2} & \cdots & x_{l} & x_{l+1} & x_{1} & x_{2} & \cdots & x_{i} \\
y_{i+1} & y_{i+2} & \cdots & y_{l} & y_{l+1} & y_{1} & y_{2} & \cdots & y_{i}
\end{pmatrix}$$

$$\coloneqq = \begin{pmatrix}
\hat{x}_{1}^{i} & \hat{x}_{2}^{i} & \cdots & \hat{x}_{l}^{i} & \hat{x}_{l+1}^{i} \\
\hat{y}_{1}^{i} & \hat{y}_{2}^{i} & \cdots & \hat{y}_{l}^{i} & \hat{y}_{l+1}^{i}
\end{pmatrix}$$

$$\coloneqq = (\hat{p}_{1}^{i}, \hat{p}_{2}^{i}, \cdots, \hat{p}_{l}^{i}, \hat{p}_{l+1}^{i}),$$
(184)

which is a rotation of in-context examples for  $i \in [l] \cup \{0\}$ . Therefore, we have

$$\begin{split} &\sum_{i=1}^{l} G_{i,l+1}(\boldsymbol{w}^{(t)})(l-i+1) \\ &= \sum_{i=1}^{l} G_{i,l+1}^{0}(\boldsymbol{w}^{(t)})(l-i+1) \\ &\leq \sum_{i=1}^{l} G_{i,l+1}^{0}(\boldsymbol{w}^{(t)}) + \sum_{i=1}^{l} G_{i,l+1}^{l}(\boldsymbol{w}^{(t)})(1 - \sigma(\boldsymbol{w}^{(t)^{\top}}\hat{p}_{1}^{l})) + \sum_{i=1}^{l} G_{i,l+1}^{l-1}(\boldsymbol{w}^{(t)})(1 \\ &- \sigma(\boldsymbol{w}^{(t)^{\top}}\hat{p}_{1}^{l-1}))(1 - \sigma(\boldsymbol{w}^{(t)^{\top}}\hat{p}_{2}^{l-1})) + \dots + \sum_{i=1}^{l} G_{i,l+1}^{2}(\boldsymbol{w}^{(t)}) \prod_{j=1}^{l-1} (1 - \sigma(\boldsymbol{w}^{(t)^{\top}}\hat{p}_{j}^{2})) \\ &\leq \max_{j \in [l]} \left\{ \sum_{i=1}^{l} G_{i,l+1}^{j}(\boldsymbol{w}^{(t)}) \right\} \cdot (1 + (1 - \sigma(\boldsymbol{w}^{(t)^{\top}}\hat{p}_{1}^{l})) + (1 - \sigma(\boldsymbol{w}^{(t)^{\top}}\hat{p}_{1}^{l-1}))(1 \\ &- \sigma(\boldsymbol{w}^{(t)^{\top}}\hat{p}_{2}^{l-1})) + \dots + \prod_{j=1}^{l-1} (1 - \sigma(\boldsymbol{w}^{(t)^{\top}}\hat{p}_{j}^{2}))) \\ &\leq 1 + (1 - \sigma(\boldsymbol{w}^{(t)^{\top}}\hat{p}_{1}^{l})) + (1 - \sigma(\boldsymbol{w}^{(t)^{\top}}\hat{p}_{1}^{l-1}))(1 - \sigma(\boldsymbol{w}^{(t)^{\top}}\hat{p}_{2}^{l-1})) + \dots \\ &+ \prod_{j=1}^{l-1} (1 - \sigma(\boldsymbol{w}^{(t)^{\top}}\hat{p}_{j}^{2})) \\ &\leq 1 + 1 - c + (1 - c)^{2} + \dots + (1 - c)^{l-1} \\ &\leq \frac{1}{c} \\ &\leq \Theta(1), \end{split}$$

where the third to last step holds since that when  $t \leq \min\{\eta^{-1}\beta^{-2}\kappa_a^{-1}(1-p_a)^{-1}V, \eta^{-1}M_1^{\frac{2}{3}}((1-p_a)\beta)^{-\frac{2}{3}}(\kappa_a(1-p_a))^{-\frac{1}{3}}V^{\frac{1}{3}}\}$ , there exists  $c \in (0,1)$  and  $C \in (0,1)$ , C > c, such that  $c \leq \sigma(\boldsymbol{w}^{(t)}^{\top}\boldsymbol{p}_j) \leq C$  for any  $j \in [l]$ .