
Theoretical Insights into In-context Learning with Unlabeled Data

Anonymous Author(s)

Affiliation

Address

email

Abstract

Recent research shows that in-context learning (ICL) can be effective even in settings where demonstrations have missing or incorrect labels. This motivates a deeper understanding of how sequence models leverage unlabeled data. We consider a canonical setting where the in-context demonstrations are drawn according to a binary Gaussian mixture model (GMM) and a certain fraction of the demonstrations have missing labels. We provide a comprehensive theoretical study to show that: (1) The loss landscape of one-layer linear attention learns the optimal fully-supervised learner but it completely fails to leverage the unlabeled data. (2) Multilayer as well as looped transformers can effectively leverage unlabeled data by implicitly constructing estimators of the form $\sum_{i \geq 0} a_i (X^\top X)^i X^\top \mathbf{y}$ with X and \mathbf{y} denoting features and visible labels. We shed light on the class of polynomials that can be expressed as a function of depth/looping and draw connections to iterative pseudo-labeling. Building on these insights and the importance of depth, we propose looping off-the-shelf tabular foundation models, such as TabPFN or TabICL, to enhance their semi-supervision capabilities. Extensive evaluations on real-world datasets reveal that our method significantly improves the semisupervised tabular learning performance over the standard single pass inference.

1 Introduction

In-Context Learning (ICL) is an intriguing capability of modern language models and has enjoyed remarkable empirical success (Brown et al., 2020; Min et al., 2022). This success is also being extended to multimodal scenarios (Zhou et al., 2024) as well as other modalities such as tabular data (Hollmann et al., 2022). The push toward test-time scaling and long-context models (Snell et al., 2024; Guo et al., 2025) has further boosted the benefits of ICL by allowing the model to ingest a large number of demonstrations. For instance, in “Many-shot in-context learning” paper, Agarwal et al. (2024) demonstrate that pushing more examples into context window can substantially boost the accuracy. The many-shot ICL setting naturally raises the question of when and how ICL can succeed with weaker supervision. As we can harness longer context models to boost predictive accuracy, we may indeed run out of high-quality demonstrations with verified answers/chain-of-thoughts and may want to utilize weaker data sources. This motivates our central question:

***Q:** How can transformers learn from unlabeled data? What are key architectural considerations?*

We primarily investigate this question under a semisupervised ICL (SS-ICL) setting with GMMs. Formally, given a prompt containing a dataset of feature-label pairs $(\mathbf{x}_i, y_i)_{i=1}^n \in \mathbb{R}^d \times \mathbb{R}$ as demonstrations and a query feature \mathbf{x} (see Eq. (2)), a model trained for ICL learns to predict the corresponding output y given prompt. For ICL with a supervised binary GMM model, we have $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}_{y_i}, \sigma^2 \mathbf{I})$ and $y_i \in \{-1, 1\}$, $i \in [n]$, and the component means $\boldsymbol{\mu}_{\pm 1}$ that parameterize the classification task are sampled from a prior task distribution. This prompt model is well studied under various fully-supervised

settings (Garg et al., 2022; Von Oswald et al., 2023; Ahn et al., 2023; Akyürek et al., 2023; Mahankali et al., 2024; Collins et al., 2024; Shen et al., 2024) where each demonstration includes a clearly labeled output. In our SS-ICL setting, only m out of n total samples have correct labels ($m \leq n$) either -1 or 1 , and remaining labels are unknown and fed to the model as $y_i = 0$.

In this work, we provide a comprehensive theoretical and empirical study of attention models with varying depths when trained with SS-ICL. Our analysis reveals **the importance of depth**: despite being able to implement the optimal supervised learner, single-layer linear attention completely fails to leverage unlabeled examples. In contrast, deeper or looped transformer architectures can emulate strong semi-supervision algorithms, approaching the performance of the Bayes-optimal classifier as depth increases. Informed by the importance of depth/looping, we also devise semisupervision strategies for tabular foundation models. Our specific contributions are:

- ◇ **Landscape of one-layer linear attention (§3)**: We study the optimization landscape of single-layer linear attention for the SS-ICL problem under an isotropic task prior. We prove that the global minimum of the loss function is the plug-in estimator, i.e., $\hat{y} = \text{sgn}(\mathbf{x}^\top \hat{\boldsymbol{\mu}})$ with $\hat{\boldsymbol{\mu}} = \mathbf{X}^\top \mathbf{y}$ where $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{y} \in \mathbb{R}^n$ are the feature-labels of the ICL demonstrations. This implies that 1-layer model learns Bayes-optimal classifier in the fully-supervised setting, but completely fails to make use of unlabeled data.
- ◇ **Depth is crucial but shallow can suffice (§4)**: We show that multilayer linear attention can emulate semisupervised learners by implementing polynomial estimators of the form

$$\hat{\boldsymbol{\mu}} = \sum_{i=0}^K a_i (\mathbf{X}^\top \mathbf{X})^i \mathbf{X}^\top \mathbf{y},$$

which can be interpreted as the model implicitly conducting *iterative pseudo-labeling*. We show that L -layer (or looped) attention can express up to $K = O(3^L)$ powers, highlighting exponentiation requires only logarithmic depth. We provide characterizations of the set of expressible polynomials through different constructions (where each layer updates features or labels of the previous layer). Corroborating these, experiments reveal that shallow transformers with $L \geq 2$ already achieve strong results and their performance can be approximately predicted through an eigen-estimator combining $i = 0$ and ∞ (see **SSPI- k**).

- ◇ **Applications to Tabular FMs (§5)**: Tabular foundation models such as TabPFN (Hollmann et al., 2022, 2025) and TabICL (Qu et al., 2025) represent a suitable application of theory as they also model the ICL examples with a single token. To harness unlabeled examples, we propose a novel strategy that iteratively creates soft pseudo-labels by *explicitly looping the tabular FM* while controlling validation risk. Focusing on the few-shot learning setting where TabPFN-v2 excels, we demonstrate that our approach can significantly improve predictive performance on various real-world datasets.

1.1 Related Work

Theoretical Analysis of In-Context Learning Recent work has developed theoretical frameworks for understanding in-context learning in transformers. Akyürek et al. (2023), Von Oswald et al. (2023) and Dai et al. (2023) demonstrated that transformers emulate gradient descent during ICL. Xie et al. (2022) offered a Bayesian perspective, while Zhang et al. (2023, 2024) showed transformers learn linear models in-context. Ahn et al. (2023) established they implement preconditioned gradient descent, and Mahankali et al. (2024) proved one-step gradient descent is optimal for single-layer linear attention. Works by Li et al. (2023) and Li et al. (2024) analyzed generalization capabilities of transformers. However, these frameworks exclusively focus on fully-supervised settings, leaving a critical gap in understanding how transformers handle partially labeled data—a common real-world scenario. Our work addresses this gap by providing the first theoretical characterization of semi-supervised in-context learning. Wang et al. considers a setting where the model observes demonstrations of the form (query, response _{i} , reward _{i}) and aims to correct its response based on the reward sequence. While our work has a different focus, it highlights that the model can correct/impute the missing labels in the context using implicit feedback from labeled demonstrations.

Semi-Supervised Learning Traditional semi-supervised learning (SSL) aims to leverage unlabeled data to improve classifier performance. For linear classifiers, Oymak & Gulcu (2020) characterized self-training iterations and demonstrated rejecting low-confidence samples; further theoretical

analyses of self-training/pseudo-labeling cover deep networks Wei et al. (2021) and models like gradient-boosted trees Kumar et al. (2020). For Gaussian Mixture Models (GMMs), Lelarge & Miolane (2019) quantified maximal improvement from unlabeled data, while Krishnapuram et al. (2004) developed graph-based priors. Learning GMMs via Expectation-Maximization (EM) or pseudo-labeling, especially with few labels, is well-studied. Ratsaby & Venkatesh (1995) provided early PAC-style bounds for GMMs learned from few labeled and many unlabeled points. Balakrishnan et al. (2017) offered further statistical guarantees for EM. Nigam et al. (2000) demonstrated empirically that EM (viewable as iterative pseudo-labeling Fan et al. (2023)) with pseudo-labels significantly reduces text classification error using unlabeled documents. These foundational works, with ongoing research in areas like agnostic learning Kwon & Caramanis (2020) and evolving theories Xu et al. (2021), underpin many SSL concepts. While these works established fundamental principles, they did not consider how these concepts apply to in-context learning with transformers. Our contribution bridges this gap by showing how transformer depth enables effective utilization of unlabeled examples within the prompt, essentially implementing semi-supervised learning without parameter updates.

2 Problem Setup and Preliminaries

We study ICL in the setting of semi-supervised classification, where the in-context demonstrations are drawn from a binary Gaussian mixture model (GMM). We begin by introducing core notation.

Notation: Denote the set $\{1, 2, \dots, n\}$ as $[n]$ and use bold letters, such as \mathbf{x} and \mathbf{X} , to represent vectors and matrices respectively. Let $\mathcal{Q}(\cdot)$ function return the right tail of the standard normal distribution.

We use $\text{sgn}(\cdot)$ denote the sign function which is defined as follows: $\text{sgn}(x) = \begin{cases} 1, & x \geq 0 \\ -1, & x < 0 \end{cases}$.

2.1 Semi-supervised Data Model

Consider a d -dimensional semi-supervised binary GMM with n examples $(\mathbf{x}_i, y_i)_{i=1}^n$, where $\mathbf{x}_i \in \mathbb{R}^d$ denotes the feature vector and $y_i \in \{-1, 0, 1\}$ represents the corresponding observed label, with $y_i = 0$ indicating a missing label, and each label is revealed independently with probability $p \in [0, 1]$. Specifically, the data is generated as follows (for each $i \in [n]$):

$$\mathbf{x}_i = y_i^c \cdot \boldsymbol{\mu} + \boldsymbol{\xi}_i, \quad y_i = \begin{cases} y_i^c, & \text{w.p. } p \\ 0, & \text{w.p. } 1 - p \end{cases} \quad \text{and} \quad y_i^c = \begin{cases} 1, & \text{w.p. } 1/2 \\ 0, & \text{w.p. } 1/2 \end{cases}. \quad (1)$$

Here $\boldsymbol{\mu} \sim \text{Unif}(\mathbb{S}^{d-1})$ denotes the task mean, which is sampled uniformly from the unit sphere, and $\boldsymbol{\xi}_i \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ is the random noise with $\sigma \geq 0$ being the noise level that controls the variability of \mathbf{x}_i around its mean. y_i^c denotes the true class label that is uniform over $\{-1, 1\}$. Observe that $p = 1$ corresponds to fully supervised learning and $p = 0$ corresponds to fully-unsupervised learning.

2.2 In-context Learning and Linear Attention

We build on the setting of (Garg et al., 2022; Mahankali et al., 2024; Zhang et al., 2023; Li et al., 2024) and construct the in-context prompts with examples drawn from the model (1) as follows.

Prompt Generation Given a task vector $\boldsymbol{\mu} \sim \text{Unif}(\mathbb{S}^{d-1})$, we sample $(n + 1)$ in-context demonstrations $(\mathbf{x}_i, y_i)_{i=1}^{n+1}$ according to (1) and construct the prompt

$$\mathbf{Z} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_n & \mathbf{x} \\ y_1 & y_2 & \cdots & y_n & 0 \end{bmatrix}^\top \in \mathbb{R}^{(n+1) \times (d+1)}. \quad (2)$$

We will investigate training a transformer such that given \mathbf{Z} as prompt, it correctly predicts the label $y := y_{n+1}^c$ of the query $\mathbf{x} := \mathbf{x}_{n+1}$ through ICL.

Model Architecture Our work primarily focuses on training of linear attention models. Given any prompt $\mathbf{Z} \in \mathbb{R}^{(n+1) \times (d+1)}$, which can be treated as a sequence of $(d + 1)$ -dimensional tokens, the linear attention mechanism outputs

$$\text{att}(\mathbf{Z}; \mathcal{W}) = (\mathbf{Z} \mathbf{W}_q \mathbf{W}_k^\top \mathbf{Z}^\top) \mathbf{M} \mathbf{Z} \mathbf{W}_v \quad (3)$$

where $\mathcal{W} := \{\mathbf{W}_k, \mathbf{W}_q, \mathbf{W}_v \in \mathbb{R}^{(d+1) \times (d+1)}\}$ denotes the key, query and value weight matrices, respectively. Therefore, given the prompt matrix $\mathbf{Z} \in \mathbb{R}^{(n+1) \times (d+1)}$ as input, the attention mechanism outputs a $(n+1)$ -length sequence (i.e., $\text{att}(\mathbf{Z}; \mathcal{W}) \in \mathbb{R}^{(n+1) \times (d+1)}$). Note that the label for the query \mathbf{x} is excluded from the prompt \mathbf{Z} . Similar to Ahn et al. (2023), we consider a training objective with a mask $\mathbf{M} = \begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$ to ensure inputs cannot attend to their own labels and training can be parallelized. To ensure that all in-context examples are treated equally and that the model remains invariant to their order/position, we do not apply a causal mask following Ahn et al. (2023). In contrast, Li et al. (2025) explores the use of causal masking in multi-layer linear attention and analyzes its impact on the final prediction.

Building upon the single-layer linear attention mechanism of (3), we can extend our model to multiple layers to capture more complex patterns. Consider optimizing an L -layer linear attention model and let \mathbf{Z}_ℓ be the input of ℓ th layer, $\ell \in [L]$. Additionally, let $\mathcal{W}_\ell := \{\mathbf{W}_{k\ell}, \mathbf{W}_{q\ell}, \mathbf{W}_{v\ell} \in \mathbb{R}^{(d+1) \times (d+1)}\}$ be the corresponding weight matrices of ℓ th layer. Then, the input prompt of ℓ th layer is defined by

$$\mathbf{Z}_\ell = \mathbf{Z}_{\ell-1} + \text{att}(\mathbf{Z}_{\ell-1}; \mathcal{W}_{\ell-1}) \quad \text{for} \quad \ell = 2, \dots, L,$$

and $\mathbf{Z}_1 = \mathbf{Z}$. We focus on the next-token prediction setting, where the model makes a prediction based on the final query token $[\mathbf{x}^\top \mathbf{0}]^\top$. Let $\mathbf{h} \in \mathbb{R}^{d+1}$ denote the linear prediction head. We define the output of the L -layer linear attention model at the last (query) token as

$$f_{\text{att-}L}(\mathbf{Z}) = \mathbf{h}^\top \text{att}(\mathbf{Z}_L; \mathcal{W}_L)_{[n+1]}. \quad (4)$$

Recalling the sign function, the predicted label for \mathbf{x} is given by $y_{\text{att-}L}(\mathbf{Z}) = \text{sgn}(f_{\text{att-}L}(\mathbf{Z}))$.

Model Training With our attention-based architecture established, we now turn to the training procedure and evaluation metrics. Consider the ICL setting where each input prompt \mathbf{Z} (cf. (2)) corresponds to a randomly sampled task vector $\boldsymbol{\mu} \sim \text{Unif}(\mathbb{S}^{d-1})$ and let $\ell(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ be the loss function. Additionally, define the set of attention weights $\mathcal{W}^{(L)} := \cup_{\ell=1}^L \mathcal{W}_\ell \in (\mathbb{R}^{(d+1) \times (d+1)})^{3L}$. The objective of L -layer linear attention takes the following form:

$$\min_{\mathcal{W}^{(L)}, \mathbf{h}} \mathcal{L}_{\text{att-}L}(\mathcal{W}^{(L)}, \mathbf{h}) \quad \text{where} \quad \mathcal{L}_{\text{att-}L}(\mathcal{W}^{(L)}, \mathbf{h}) = \mathbb{E}[\ell(y, f_{\text{att-}L}(\mathbf{Z}))]. \quad (5)$$

Here, $y := y_{n+1}^c$ and the expectation subsumes the randomness of $\boldsymbol{\mu}$ and $(\xi_i, y_i)_{i=1}^{n+1}$. The search space for $\mathcal{W}^{(L)}$ is $(\mathbb{R}^{(d+1) \times (d+1)})^{3L}$, and for \mathbf{h} is \mathbb{R}^{d+1} .

3 Loss Landscape of One-layer Linear Attention under SS-ICL

Previous work (Ahn et al., 2023; Li et al., 2024; Mahankali et al., 2024) has shown that an optimized single-layer linear attention implements a form of preconditioned gradient descent over the linear in-context demonstrations provided within the prompt. However, to the best of our knowledge, prior studies have not addressed the semi-supervised setting, where some in-context labels are missing. In this section, we analyze the optimization behavior of single-layer linear attention under the semi-supervised binary GMM setting described in Section 2, and demonstrate that the single-layer model learns the optimal fully-supervised learner, but fails to utilize the unlabeled data.

We begin with the following optimal supervised label estimator under our problem setting.

Supervised Plug-in (SPI) Estimator The plug-in method is a classical approach for supervised classification problems, aiming to find a linear combination of features that separates different categories. Under our problem setting, it also serves as the Bayes-optimal estimator given only labeled data (Mignacco et al., 2020; Lelarge & Miolane, 2019). Consider the binary semi-supervised GMM problem described in (1) with dataset $(\mathbf{x}_i, y_i)_{i=1}^n$, and let $\mathcal{I} \subset [n]$ represent the indices of labeled samples, e.g., $y_i \neq 0$ for $i \in \mathcal{I}$. The SPI estimator returns the task mean

$$\hat{\boldsymbol{\mu}}_s = \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} y_i \mathbf{x}_i. \quad (\text{SPI})$$

We next present the following theorem establishes that, under isotropic task prior, optimal single-layer linear attention is equivalent to the SPI estimation.

Theorem 1 Let the prompt (cf. (2)) be generated as described in Section 2.2. Consider the objective (cf. (5)) with $L = 1$ and squared loss function $\ell(y, \hat{y}) = (y - \hat{y})^2$, and denote the optimal prediction as y_{att-1}^* . Let $\hat{\mu}_s$ represent the SPI estimator defined in (SPI). Then, for any \mathbf{Z} from (2), we have

$$y_{att-1}^*(\mathbf{Z}) = \text{sgn}(\mathbf{x}^\top \hat{\mu}_s). \quad (6)$$

Additionally, its classification error obeys

$$\begin{aligned} \mathbb{P}(y_{att-1}^*(\mathbf{Z}) \neq y) &= \mathbb{E}_{g \sim \mathcal{N}(0,1), h \sim \chi_{d-1}^2} \left[\mathbb{Q} \left(\frac{1 + \varepsilon_\sigma g}{\sigma \sqrt{(1 + \varepsilon_\sigma g)^2 + \varepsilon_\sigma^2 h}} \right) \right] \\ &\leq \mathbb{Q} \left(\frac{1 - 10d\varepsilon_\sigma^2}{\sigma} \right) + e^{-d} + e^{-1/8\varepsilon_\sigma^2}. \end{aligned} \quad (\text{SPI-ERR})$$

where we define $\varepsilon_\sigma = \sigma / \sqrt{np}$ and χ_d^2 defines chi-squared distribution with d degrees of freedom.

The proof of Theorem 1 is deferred to the appendix. Eq. (6) shows that one-layer linear attention model indeed implements the optimal supervised predictor, assuming access to np labeled examples. Therefore, the classification error corresponds exactly to that of the SPI estimator. Such classification problem has been extensively studied. For example, Thrampoulidis et al. (2020); Wang & Thrampoulidis (2022). Most existing work focuses on a single classification task under asymptotic data regimes. In contrast, within the ICL framework considered in our setting, the task mean μ is randomly sampled, and the classification error is computed by averaging over random draws of \mathbf{Z} , y , and μ . Accordingly, in (SPI-ERR), we express the error in a simplified form as an expectation.

The experimental results in Figure 1 support Theorem 1, where dark blue circular markers represent the performance of the single-layer linear attention model, blue curves show the classification accuracy of the SPI estimator, and the red dotted curves depict $1 - \mathbb{P}(y_{att-1}^*(\mathbf{Z}) \neq y)$ as computed from (SPI-ERR). The alignments of these curves empirically validate Theorem 1. Implementation details and further discussion are provided in Section 5. Based on these results, we reach the following conclusion:

1-layer linear attention learns optimal supervised estimator but doesn't benefit from unlabeled data.

As shown in Figs 1b and 1c, when the number of labeled samples ($np = 10$) is fixed, increasing the number of unlabeled examples (even up to ~ 10000) has no effect on performance, as the dark blue markers remain at the same level.

At first glance, this may seem counterintuitive—while the data is unlabeled, it still contains information about the classification feature. For instance, the mean of the data points carries relevant information, and one might expect the model to extract and leverage this for better predictions. This expectation is particularly reasonable when a large amount of unlabeled data is available, as the sample covariance matrix approximates the population covariance, i.e., $\mathbb{E}[X^\top X/n] = \mu\mu^\top + \sigma^2\mathbf{I}$ where $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$. The key insight into why single-layer attention fails to leverage unlabeled data lies in the expectation structure. In our isotropic GMM setting where $\mu \sim \text{Unif}(\mathbb{S}^{d-1})$, the sample covariance matrix converges to $\mathbb{E}[X^\top X/n] = \mathbb{E}[\mu\mu^\top] + \sigma^2\mathbf{I} = (1/d + \sigma^2)\mathbf{I}$, which contains no task-specific information. The expectation across multiple tasks loses the signal from μ . This explains why single-layer attention, operating in a meta-learning framework across many tasks rather than optimizing for a single fixed task, cannot extract useful information from unlabeled data.

In the following section, we study multi-layer linear attention and demonstrate that it has the ability to propagate $X^\top X$ into deeper layers, thereby enabling the model to utilize the unlabeled data.

4 Multi-layer Attention and the Benefits of Depth

In this section, we explore how deeper attention models can effectively utilize the unlabeled data. Let

$$X = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_n \end{bmatrix}^\top \in \mathbb{R}^{n \times d} \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} y_1 & y_2 & \cdots & y_n \end{bmatrix}^\top \in \mathbb{R}^n. \quad (7)$$

We first present the following propositions to show that multi-layer as well as looped linear attention can be expressed as a polynomial function of $X^\top X$. This structure allows the models to leverage unlabeled data to improve the estimation of the task mean μ .

Proposition 1 Given an L -layer linear attention model described in Section 2.2 with input prompt \mathbf{Z} defined in (2), one can construct the key, query, value weight matrices and the linear prediction head such that the model outputs

$$f_{\text{att-L}}(\mathbf{Z}) = \mathbf{x}^\top \mathbf{A} \mathbf{X}^\top \mathbf{y}. \quad (8)$$

Then, the following \mathbf{A} matrices are achievable via label and feature updates:

- **Label propagation:** $\mathbf{A} = c \prod_{\ell=1}^{L-1} (\mathbf{I} + c_\ell \mathbf{X}^\top \mathbf{X})$ for arbitrary constants $\{c, c_1, \dots, c_{L-1}\}$;
- **Feature propagation:** $\mathbf{A} = c (\mathbf{X}^\top \mathbf{X})^{3^{L-1}-1}$ for an arbitrary constant c .

Proposition 2 Consider the same setting as in Proposition 1. There exists a single-layer linear attention model whose parameters can be constructed to reproduce the output in (8), with $c_\ell \equiv c'$ for some arbitrary constant c' .

The proofs of Proposition 1 and 2 are deferred to the appendix. In the following, we provide further clarification on the label and feature propagation.

1. The final prediction of the label propagation process can be rewritten as

$$f_{\text{att-L}}(\mathbf{Z}) = c \mathbf{x}^\top \mathbf{X}^\top \mathbf{y}_L \quad \text{where} \quad \mathbf{y}_{\ell+1} = (\mathbf{I} + c_\ell \mathbf{X} \mathbf{X}^\top) \mathbf{y}_\ell, \quad \text{for } \ell \in [L-1]$$

with $\mathbf{y}_1 = \mathbf{y}$. Here, \mathbf{y}_ℓ can be interpreted as the pseudo-labels input to the ℓ th layer, and each c_ℓ is parameterized by the attention mechanism in the corresponding layer. Although not exactly equivalent, the L -layer linear attention process shares similarities with the Expectation-Maximization (EM) algorithm for semi-supervised learning, with L iterations of pseudo-labeling and a different label update strategy.

2. In contrast, the feature propagation process yields the final prediction

$$f_{\text{all-L}}(\mathbf{Z}) = c \mathbf{x}_L^\top \mathbf{X}_L^\top \mathbf{y} \quad \text{where} \quad \mathbf{X}_{\ell+1} = (\mathbf{X}_\ell \mathbf{X}_\ell^\top) \mathbf{X}_\ell \quad \text{and} \quad \mathbf{x}_{\ell+1} = (\mathbf{X}_\ell^\top \mathbf{X}_\ell) \mathbf{x}_\ell, \quad \text{for } \ell \in [L-1]$$

with $\mathbf{X}_1 = \mathbf{X}$ and $\mathbf{x}_1 = \mathbf{x}$. Here, $(\mathbf{X}_\ell, \mathbf{x}_\ell)$ can be viewed as the input features at the ℓ th layer, encoding exponentially higher-order powers of $\mathbf{X}^\top \mathbf{X}$. This result highlights that a linear attention model requires only $O(\log k)$ layers to represent polynomial functions of degree k .

Our construction for *label propagation* is inherently related to the *gradient descent* emulation capability of linear attention Ahn et al. (2023). However, the *feature propagation* construction is fundamentally different and underscores the transformer's capability to implement rapid power iteration over the empirical covariance $\mathbf{X}^\top \mathbf{X}$. In the above constructions, each attention block with residual connections updates features or labels using one parameter, namely mappings of the form $\mathbf{X} \rightarrow \mathbf{X} + \alpha \mathbf{X} \mathbf{X}^\top \mathbf{X}$ or $\mathbf{y} \rightarrow \mathbf{y} + \beta \mathbf{X} \mathbf{X}^\top \mathbf{y}$. The lemma below shows that, even if the multilayer model can express polynomials of $\mathbf{X}^\top \mathbf{X}$ with exponential degrees in depth, the expressible manifold of polynomials has dimensionality linear in depth.

Lemma 1 (Label + Feature Propagation) For an L -layer linear attention model, the resulting eventual prediction corresponds to the matrix \mathbf{A} in Proposition 1 of the form

$$\mathbf{A} = \sum_{\ell=0}^{(3^L-3)/2} a_\ell (\mathbf{X}^\top \mathbf{X})^\ell. \quad (9)$$

The coefficients $\mathbf{a} := [a_0 \ a_1 \ \dots \ a_{(3^L-3)/2}]^\top$ lie on a manifold of dimension at most $2L$ as \mathbf{a} can be expressed as $\mathbf{a} = g(\mathbf{c})$ for some smooth function $g : \mathbb{R}^{2L} \rightarrow \mathbb{R}^{(3^L-3)/2}$ with \mathbf{c} representing the parameters of individual layers.

Recall the SPI estimator $\hat{\boldsymbol{\mu}}_s$ from (SPI), and that \mathbf{y} denotes the visible labels defined in Section 2.1 and (7). We have $\hat{\boldsymbol{\mu}}_s = \frac{1}{|\mathcal{I}|} \mathbf{X}^\top \mathbf{y}$. Motivated by Proposition 1 that multi-layer linear attention can implement higher-degree polynomials of $\mathbf{X}^\top \mathbf{X}$, we introduce the following SSPI estimator, which makes predictions based on the supervised estimate $\hat{\boldsymbol{\mu}}_s$ combined with higher-order debiased term of the form $(\mathbf{X}^\top \mathbf{X}/n - \sigma^2 \mathbf{I})^k$.

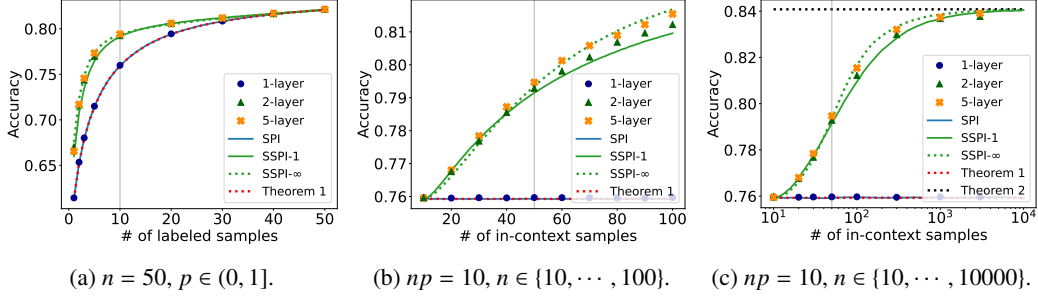


Figure 1: Experimental results support our theoretical findings presented in Sections 3 and 4. Details and discussion can be found in Sections 3, 4, and 5.

Semisupervised Plug-in (SSPI) Estimator Observe that the feature covariance satisfies $\mathbb{E}[X^\top X]/n = \mu\mu^\top + \sigma^2\mathbf{I}$, and the top eigenvector of the centered covariance matrix $(X^\top X/n - \sigma^2\mathbf{I})$ asymptotically aligns with either μ or $-\mu$. Therefore, with a substantial amount of unlabeled data, we propose the semisupervised plug-in (SSPI) estimator as follows:

$$\hat{\mu}_{ss-k} = \alpha \hat{\mu}_s + (1 - \alpha)(X^\top X/n - \sigma^2\mathbf{I})^k \hat{\mu}_s \quad (\text{SSPI-}k)$$

where $\hat{\mu}_s$ is the SPI estimator (c.f. (SPI)), and $\alpha \in [0, 1]$ controls the trade-off between the fully-supervised and semi-supervised estimator. The optimal choice of α depends on the problem parameters n, d and p . Note that as $k \rightarrow \infty$, the term $(X^\top X/n - \sigma^2\mathbf{I})^k$ converges (up to scaling) to a rank-one projection onto the top eigenvector of the debiased covariance matrix, effectively serving as an estimator for μ (up to sign).

In Figure 1, we present the prediction accuracies of 2-layer and 5-layer linear attention models, shown by green and orange markers, respectively. We also evaluate the SSPI algorithm with varying k values, where the green solid curve corresponds to SSPI-1, and the green dotted represents SSPI- ∞ , both using their respective optimal choices of α . The results reveal a close alignment between multi-layer linear attention and SSPI estimators. Notably, the 2-layer model outperforms SSPI-1, due to its ability to implement higher-degree polynomials of $X^\top X$ (cf. Proposition 1 and Equation (9)). When the sample size is sufficiently large (e.g., $n > 50$ in Figure 1b), the top eigenvector provides a more accurate estimate of the task mean, enabling SSPI- ∞ to achieve higher accuracy. Furthermore, since the 5-layer model is capable of representing higher-order functions than the 2-layer model, it can better estimate the top eigenvector, resulting in performance that closely matches that of SSPI- ∞ .

In the following, we analyze the optimal classifier of the form $\text{sgn}(\mathbf{x}^\top \mathbf{A} \hat{\mu}_s)$ for a GMM, and provide insights into its behavior in the asymptotic regime as $n \rightarrow \infty$.

Theorem 2 Consider a binary GMM defined in Section 2.1 and suppose that $(\mathbf{x}_i, y_i)_{i=1}^{n+1}$ is generated using a fixed μ following (1). Given matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$, define prediction

$$\hat{y}_A = \text{sgn}(\mathbf{x}^\top \mathbf{A} \hat{\mu}_s).$$

where $\hat{\mu}_s$ is the SPI estimator defined in (SPI). Let $\mathcal{A}^* := \min_{\mathbf{A} \in \mathbb{R}^{d \times d}} \mathbb{P}(\hat{y}_A \neq y)$ be its optimal solution set. Then, $\mu\mu^\top \in \mathcal{A}^*$. Additionally, it obeys

$$\mathbb{P}(\hat{y}_{\mu\mu^\top} \neq y) = Q(1/\sigma) + Q(\sqrt{np}/\sigma) - 2Q(1/\sigma)Q(\sqrt{np}/\sigma). \quad (10)$$

Note that, $\mathbb{P}(\hat{y}_{\mu\mu^\top} \neq y)$ depends on np and σ only, regardless of μ and d .

Theorem 3 Let the prompt \mathbf{Z} be generated as described in Section 2.2, and consider an L -layer linear attention model with $L \geq 2$ and $n = \infty$. Additionally, let $\hat{\mu}_s$ be the SPI estimator defined in (SPI). There exists model constructions such that for any \mathbf{Z} following (2), its prediction satisfies

$$y_{\text{att-}L}(\mathbf{Z}) = \text{sgn}(\mathbf{x}^\top \mu\mu^\top \hat{\mu}_s).$$

The proof follows directly from Proposition 1 (label propagation), which shows that multi-layer linear attention can output $\mathbf{x}^\top (X^\top X/n - \sigma^2\mathbf{I}) \hat{\mu}_s$. As $n \rightarrow \infty$, the empirical covariance converges to its expectation, i.e., $X^\top X/n - \sigma^2\mathbf{I} \rightarrow \mu\mu^\top$. The results in Figure 1c validate Theorem 3, showing that as n becomes large enough, (i.e., $n = 10000$) the predictions from both 2-layer and 5-layer linear attention models, as well as the SSPI-1 and SSPI- ∞ estimators, closely align with the classification error characterized in Theorem 2, depicted by the black dotted line.

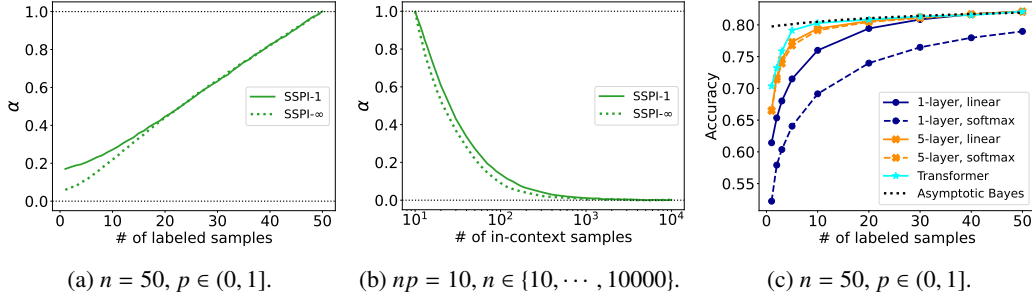


Figure 2: Additional experimental results. (a)&(b): Analysis of the optimal α values for the SSPI estimator (cf. (SSPI- k)) under varying (n, p, k) . (c): Comparison of different model architectures for the SS-ICL problem. Further details are provided in Section 5.

5 Experiments

In Sections 3 and 4, we introduced Figure 1 and demonstrated its consistency with our theoretical results. In this section, we describe the experimental setup and implementation details. Additionally, we present further empirical findings to investigate additional questions of interest. Motivated by Proposition 2, which suggests that looping can help leverage unlabeled data, Section 5.2 introduces an algorithm based on the TabPFN, showing how it can enhance prediction performance by incorporating a small amount of unlabeled data and iterative pseudo-labeling through model looping.

Experimental Setup Following Section 2, set $d = 10$ and noise level $\sigma = 1$. All models are trained using Adam optimizer with a learning rate of 10^{-3} for 40,000 epochs, with a batch size of 512. We use logistic loss in our experiments. Since our study focuses on the optimization landscape and model expressivity, and experiments are implemented via gradient descent, we repeat 10 trainings from random initialization and results are presented as the maximal test accuracy among those 10 trails.

5.1 Additional Observations

• **Exploration of Optimal α Values.** In Section 4, we introduced the SSPI- k estimator (cf. (SSPI- k)), but did not discuss the choice of the mixing parameter α , which plays a crucial role in balancing the contribution of the supervised estimator $\hat{\mu}_s$. Specifically, α controls how much weight is given to the purely supervised signal. In the fully supervised case, the optimal choice is $\alpha = 1$, as $\hat{\mu}_s$ corresponds to the Bayes-optimal estimator.

In Figures 2a and 2b, we empirically examine the optimal values of α . Given $\mu \sim \text{Unif}(\mathbb{S}^{d-1})$, we define the optimal α as the minimizer of the following cosine similarity-based objective:

$$\alpha^* := \min_{\alpha \in [0, 1]} \mathcal{L}(\alpha) \quad \text{where} \quad \mathcal{L}(\alpha) = 1 - \mathbb{E}[\text{cosine_similarity}(\mu_{ss-k}, \mu)].$$

For each setting, we optimize α using the Adam optimizer for 10,000 epochs with a batch size of 128 and a learning rate of 0.01. The results are shown in Figs 2a and 2b.

In Figure 2a, for both SSPI-1 and SSPI- ∞ , the optimal α starts near zero when the number of labeled examples is small, reflecting the limited utility of $\hat{\mu}_s$ in low-supervision regimes. As the number of labeled samples increases, α grows approximately linearly and approaches 1 when the problem becomes fully supervised. In Figure 2b, when $n = 10$ and $p = 1$ (i.e., all examples are labeled), the optimal α begins at 1. As n increases and the fraction of unlabeled data grows, α decreases significantly. This trend indicates that as the volume of unlabeled data increases, the SSPI estimator adaptively reduces reliance on the supervised component $\hat{\mu}_s$ and increases reliance on the semi-supervised component, which leverages the structure of the unlabeled data through $X^\top X$.

• **Comparison Across Different Model Architectures.** Beyond linear attention, we investigate additional model architectures under our SS-ICL setting. The comparison results are presented in Fig. 2c. The softmax attention model uses the same structure described in Section 2.2, with the only difference being the addition of a softmax operation in Eq. (3). The Transformer model introduces further nonlinearity and capacity by incorporating multi-layer perceptrons (MLPs) and layer normalization. The Transformer experiments are conducted with 5-layer models.

When comparing weaker models—such as 1-layer linear (dark blue solid) and softmax (dark blue dashed) attention—we observe that softmax attention consistently underperforms linear attention.

Table 1: Testing accuracy comparison between the baseline (trained on labeled samples only) and after 1 or 5 iterations of looping TabPFN-v2.

Dataset	OpenML1049		OpenML1464		OpenML1067		OpenML1494		OpenML1489		OpenML40981	
	10/10	20/20	10/10	20/20	10/10	20/20	10/10	20/20	10/10	20/20	10/10	20/20
Baseline	0.7497	0.8495	0.5707	0.7172	0.6883	0.7700	0.5914	0.6973	0.6216	0.6775	0.7420	0.7396
Loop-1	0.7929	0.8476	0.6066	0.7333	0.7160	0.8084	0.6091	0.6965	0.6555	0.7138	0.7445	0.7719
Loop-5	0.8287	0.8421	0.7178	0.7243	0.7710	0.8109	0.6394	0.7016	0.7016	0.7158	0.7578	0.7761

Notably, softmax attention fails to match the performance of the optimal supervised estimator, even when all labels are observed (i.e., when the number of labeled samples equals $n = 50$). Furthermore, increasing the depth of softmax attention (orange dashed curve for 5-layer softmax) still does not surpass the performance of 5-layer linear attention (orange solid curve). Among all architectures, the Transformer achieves the best performance due to its increased model capacity and expressiveness. Compared with Fig. 1a, where the orange and dark blue markers (linear attention) are identical, the Transformer significantly improves accuracy. This improvement highlights that SSPI, while effective, is not the optimal semi-supervised estimator. Although our semi-supervised setting assumes isotropic data, the characterization of its optimal algorithm remains an open and foundational problem for future exploration. In the figure, we also include the asymptotic Bayes-optimal curve (black dotted; derived from [Lelarge & Miolane \(2019\)](#)). As the number of samples increases, the results from linear attention, softmax attention, and Transformer all converge toward this optimal curve. We attribute the initial performance gap, particularly at low values of np (e.g., $np = 1$), to the scarcity of labeled data.

5.2 Tabular Experiments

To further investigate how model looping (Proposition 2) can improve label prediction, we introduce an algorithm that addresses unlabeled data by iteratively assigning pseudo-labels. We evaluate the algorithm on real-world datasets, with results presented in Table 1. We evaluated the effectiveness of our proposed looping strategy by iteratively applying TabPFN-v2 on real-world binary classification benchmarks from [Hollmann et al. \(2025\)](#). We tested two settings with equal numbers of labeled and unlabeled samples: (10, 10) and (20, 20). In each experiment, we first use the labeled data to assign soft pseudo-labels to the unlabeled data based on TabPFN-v2 predictions. These pseudo-labels are then updated iteratively through repeated looping, and each loop, we feed the model with both labeled and pseudo-labeled data. Full implementation details are provided in the supplementary material.

The results are summarized in Table 1, and each is averaged over 20 random training and test data splits. As a baseline, we use TabPFN-v2 feeding with only the labeled data to make one-shot predictions on the test set. We compare this to models after 1 iteration (Loop-1) and 5 iterations (Loop-5) of pseudo-label refinement. Our results show that the looping strategy significantly improves test accuracy. For instance, on the OpenML1464 dataset with 10 unlabeled examples, Loop-5 improves the baseline accuracy by approximately 25.8%. Notably, in some cases—such as OpenML1489—the (10 labeled, 10 unlabeled) setup with 5 loops outperforms the (20 labeled, 0 unlabeled) baseline, demonstrating the effectiveness of leveraging unlabeled data. These findings highlight that explicitly looping the tabular foundation model to iteratively refine soft pseudo-labels can substantially enhance performance by effectively incorporating information from unlabeled data.

6 Discussion and Limitations

Our paper introduces a theoretical study of semisupervised in-context learning and characterizes how transformer, specifically linear attention, models can harness unlabeled data in their context window to make inference. We show that depth is crucial to go beyond supervised estimation and utilize unlabeled data, and the latter is achieved by constructing estimators of the form $\hat{\mu} = \sum_{i \geq 0}^K a_i (X^\top X)^i X^\top y$. $\log K$ depth suffices to express a K th order polynomial which is in line with our synthetic and real experiments that corroborate that mild amount of depth/looping already achieves most of the benefit. Our core theoretical results are limited to linear attention models and it is important to understand the capabilities of the full transformer architecture. Indeed, transformer (MLP+softmax) empirically outperforms a linear attention model with equal number of layers, well approximating the Bayes optimal semisupervised estimator. It would also be exciting to go beyond the classification setting and examine how self-generated CoT rationales, as in ([Wu et al., 2023](#)), can enhance ICL capabilities for tasks that require reasoning/autoregression.

References

- Rishabh Agarwal, Avi Singh, Lei M Zhang, Bernd Bohnet, Stephanie Chan, Ankesh Anand, Zaheer Abbas, Azade Nova, John D Co-Reyes, Eric Chu, et al. Many-shot in-context learning. *arXiv preprint arXiv:2404.11018*, 2024.
- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing Systems*, 36, 2023.
- 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. URL <https://openreview.net/forum?id=0g0X4H8yN4I>.
- Sivaraman Balakrishnan, Martin J Wainwright, and Bin Yu. Statistical guarantees for the EM algorithm: From population to sample-based analysis. In *Annals of Statistics*, volume 45, pp. 77–120. Institute of Mathematical Statistics, 2017.
- 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.
- Liam Collins, Advait Parulekar, Aryan Mokhtari, Sujay Sanghavi, and Sanjay Shakkottai. In-context learning with transformers: Softmax attention adapts to function lipschitzness. *arXiv preprint arXiv:2402.11639*, 2024.
- Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. Why can GPT learn in-context? language models secretly perform gradient descent as meta-optimizers. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 4005–4019, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.247. URL <https://aclanthology.org/2023.findings-acl.247>.
- Chuang Fan, Shipeng Liu, Seyed Motamed, Shiyu Zhong, Silvio Savarese, Juan Carlos Niebles, Anima Anandkumar, Adrien Gaidon, and Stefan Scherer. Expectation maximization pseudo labels. *arXiv preprint arXiv:2305.01747*, 2023.
- Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- 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 llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Noah Hollmann, Samuel Müller, Katharina Eggenberger, and Frank Hutter. TabPFN: A transformer that solves small tabular classification problems in a second. *arXiv preprint arXiv:2207.01848*, 2022.
- Noah Hollmann, Samuel Müller, Lennart Purucker, Arjun Krishnakumar, Max Körfer, Shi Bin Hoo, Robin Tibor Schirmer, and Frank Hutter. Accurate predictions on small data with a tabular foundation model. *Nature*, 637(8045):319–326, 2025.
- Balaji Krishnapuram, David Williams, Ya Xue, Lawrence Carin, Mário Figueiredo, and Alexander Hartemink. On semi-supervised classification. *Advances in neural information processing systems*, 17, 2004.
- Ashish Kumar, Logan Engstrom, Andrew Ilyas, and Dimitris Tsipras. Understanding self-training for gradient-boosted trees. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pp. 1651–1662, 2020.

- 410 Jeongyeol Kwon and Constantine Caramanis. The em algorithm gives sample-optimality for learning
411 mixtures of well-separated gaussians. In *Conference on Learning Theory*, pp. 2425–2487. PMLR,
412 2020.
- 413 Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection.
414 *Annals of statistics*, pp. 1302–1338, 2000.
- 415 Marc Lelarge and Léo Miolane. Asymptotic bayes risk for gaussian mixture in a semi-supervised
416 setting. In *2019 IEEE 8th International Workshop on Computational Advances in Multi-Sensor
417 Adaptive Processing (CAMSAP)*, pp. 639–643. IEEE, 2019.
- 418 Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers
419 as algorithms: Generalization and stability in in-context learning. In *International Conference on
420 Machine Learning*, pp. 19565–19594. PMLR, 2023.
- 421 Yingcong Li, Ankit Singh Rawat, and Samet Oymak. Fine-grained analysis of in-context linear
422 estimation: Data, architecture, and beyond. *arXiv preprint arXiv:2407.10005*, 2024.
- 423 Yingcong Li, Davoud Ataee Tarzanagh, Ankit Singh Rawat, Maryam Fazel, and Samet Oymak.
424 Gating is weighting: Understanding gated linear attention through in-context learning. *arXiv
425 preprint arXiv:2504.04308*, 2025.
- 426 Arvind V. Mahankali, Tatsunori Hashimoto, and Tengyu Ma. One step of gradient descent is provably
427 the optimal in-context learner with one layer of linear self-attention. In *The Twelfth International
428 Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=8p3fu56lKc>.
- 430 Francesca Mignacco, Florent Krzakala, Yue Lu, Pierfrancesco Urbani, and Lenka Zdeborova. The
431 role of regularization in classification of high-dimensional noisy gaussian mixture. In *International
432 conference on machine learning*, pp. 6874–6883. PMLR, 2020.
- 433 Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke
434 Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv
435 preprint arXiv:2202.12837*, 2022.
- 436 Ojash Neopane. Lecture notes on high-dimensional statistics. [https://www.stat.cmu.edu/
437 ~arinaldo/Teaching/36709/S19/Scribed_Lectures/Feb26_Ojash.pdf](https://www.stat.cmu.edu/~arinaldo/Teaching/36709/S19/Scribed_Lectures/Feb26_Ojash.pdf), 2018.
- 438 Kamal Nigam, Andrew Kachites McCallum, Sebastian Thrun, and Tom Mitchell. Text classification
439 from labeled and unlabeled documents using em. *Machine Learning*, 39(2–3):103–134, 2000. doi:
440 10.1023/A:1007692713085.
- 441 Samet Oymak and Talha Cihad Gulcu. Statistical and algorithmic insights for semi-supervised
442 learning with self-training. *arXiv preprint arXiv:2006.11006*, 2020.
- 443 Jingang Qu, David Holzmüller, Gaël Varoquaux, and Marine Le Morvan. Tabicl: A tabular foundation
444 model for in-context learning on large data. *arXiv preprint arXiv:2502.05564*, 2025.
- 445 Joel Ratsaby and Santosh S. Venkatesh. Learning from a mixture of labeled and unlabeled examples
446 with parametric side information. In *Proceedings of the Eighth Annual Conference on Computa-
447 tional Learning Theory (COLT '95)*, pp. 412–417. ACM, 1995. doi: 10.1145/225298.225348.
- 448 Wei Shen, Ruida Zhou, Jing Yang, and Cong Shen. On the training convergence of transformers for
449 in-context classification. *arXiv preprint arXiv:2410.11778*, 2024.
- 450 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally
451 can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- 452 Christos Thrampoulidis, Samet Oymak, and Mahdi Soltanolkotabi. Theoretical insights into mul-
453 ticlass classification: A high-dimensional asymptotic view. *Advances in Neural Information
454 Processing Systems*, 33:8907–8920, 2020.
- 455 Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev,
456 Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In
457 *International Conference on Machine Learning*, pp. 35151–35174. PMLR, 2023.

- 458 Ke Wang and Christos Thrampoulidis. Binary classification of gaussian mixtures: Abundance of
 459 support vectors, benign overfitting, and regularization. *SIAM Journal on Mathematics of Data*
 460 *Science*, 4(1):260–284, 2022.
- 461 Yifei Wang, Yuyang Wu, Zeming Wei, Stefanie Jegelka, and Yisen Wang. A theoretical understanding
 462 of self-correction through in-context alignment. In *The Thirty-eighth Annual Conference on Neural*
 463 *Information Processing Systems*.
- 464 Colin Wei, Kendrick Shen, Yining Chen, and Tengyu Ma. Theoretical analysis of self-training with
 465 deep networks on unlabeled data. In *International Conference on Learning Representations (ICLR)*,
 466 2021.
- 467 Jingfeng Wu, Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, and Peter L Bartlett.
 468 How many pretraining tasks are needed for in-context learning of linear regression? *arXiv preprint*
 469 *arXiv:2310.08391*, 2023.
- 470 Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context
 471 learning as implicit bayesian inference. In *International Conference on Learning Representations*,
 472 2022. URL <https://openreview.net/forum?id=RdJVFCHjUMI>.
- 473 Jialin Xu, Zixuan Li, Sayan Mukherjee, and David Taylor. Towards understanding deep learning with
 474 persistent homology. *arXiv preprint arXiv:2106.06718*, 2021.
- 475 Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context.
 476 *arXiv preprint arXiv:2306.09927*, 2023.
- 477 Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context.
 478 *Journal of Machine Learning Research*, 25(49):1–55, 2024.
- 479 Yucheng Zhou, Xiang Li, Qianning Wang, and Jianbing Shen. Visual in-context learning for large
 480 vision-language models. In *Findings of the Association for Computational Linguistics ACL 2024*,
 481 pp. 15890–15902, 2024.

Appendix

Table of Contents

A Analysis of Single-layer Linear Attention	13
A.1 Supporting Lemmas	13
A.2 Proof of Theorem 1	16
B Analysis of Multi-layer Linear Attention	18
B.1 Proof of Proposition 1	18
B.2 Proof of Proposition 2	20
B.3 Proof of Lemma 1	20
B.4 Proof of Theorem 2	21
B.5 Non-asymptotic Analysis	23
C Additional Details on Tabular Experiments	25

A Analysis of Single-layer Linear Attention

A.1 Supporting Lemmas

Recap the SPI estimator from (SPI). Given a semi-supervised dataset $(\mathbf{x}_i, y_i)_{i=1}^n$ as described in Section 2.1, let \mathcal{I} denote the token indices set corresponding to the labeled demonstrations, that is, we have

$$y_i = \begin{cases} y_i^c, & i \in \mathcal{I} \\ 0, & \text{otherwise.} \end{cases} \quad (11)$$

Then, the SPI estimates the task mean via

$$\hat{\mu}_s = \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} y_i \mathbf{x}_i.$$

Let $\mathbf{W} \in \mathbb{R}^{d \times d}$ be the preconditioning matrix. We define the following objective:

$$\mathbf{W}^* := \arg \min_{\mathbf{W} \in \mathbb{R}^{d \times d}} \tilde{\mathcal{L}}(\mathbf{W}) \quad \text{where} \quad \tilde{\mathcal{L}}(\mathbf{W}) = \mathbb{E} \left[\left(\mathbf{x}^\top \mathbf{W} \sum_{i \in \mathcal{I}} y_i \mathbf{x}_i - y \right)^2 \right]. \quad (12)$$

Here, we set (\mathbf{x}, y) to be the query feature and its corresponding true label. The expectation subsumes the randomness in $(\mathbf{x}_i, y_i), (\mathbf{x}, y)$ as described in Section 2.1.

In this section, we provide a lemma that establishes equivalence between optimizing $\mathcal{L}_{\text{att-1}}(\mathbf{W}, \mathbf{h})$ (cf. (5) and choosing $L = 1$) and $\tilde{\mathcal{L}}(\mathbf{W})$.

Lemma 2 Consider ICL problem described in Section 2.2 with prompt defined in (2). Consider training with a single-layer linear attention with squared loss, that is, $L = 1$ and $\ell(y, \hat{y}) = (y - \hat{y})^2$. Recall the objectives from (5) and (12), and let $\mathcal{L}_{\text{att-1}}^*$ and $\tilde{\mathcal{L}}^* := \tilde{\mathcal{L}}(\mathbf{W}^*)$ be their corresponding optimal losses. Then, we have

$$\mathcal{L}_{\text{att-1}}^* = \tilde{\mathcal{L}}^*. \quad (13)$$

Additionally, let $f_{\text{att-1}}^* : \mathbb{R}^{(n+1) \times (d+1)} \rightarrow \mathbb{R}$ denote the optimal prediction (associated with the optimal loss $\mathcal{L}_{\text{att-1}}^*$). We have that $f_{\text{att-1}}^*$ is unique and for any prompt \mathbf{Z} (cf. (2))

$$f_{\text{att-1}}^*(\mathbf{Z}) = \mathbf{x}^\top \mathbf{W}^* \sum_{i \in \mathcal{I}} y_i \mathbf{x}_i. \quad (14)$$

516 **Proof.** Recap the single-layer linear attention model and its prediction from (3) and (4). We have

$$f_{\text{att-1}}(\mathbf{Z}) = \mathbf{h}^\top \text{att}(\mathbf{Z}; \mathcal{W})_{[n+1]} \quad \text{where} \quad \text{att}(\mathbf{Z}; \mathcal{W}) = (\mathbf{Z}\mathbf{W}_q\mathbf{W}_k^\top\mathbf{Z}^\top)\mathbf{M}\mathbf{Z}\mathbf{W}_v \quad (15)$$

517 with $\mathcal{W} := \{\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v\}$ being the set of the query, key and value matrices of the attention. Since \mathcal{W}
518 and \mathbf{h} are tunable parameters, without loss of generality and for simplicity, let

$$\mathbf{W} := \mathbf{W}_q\mathbf{W}_k^\top \quad \text{and} \quad \bar{\mathbf{h}} := \mathbf{W}_v\mathbf{h}.$$

519 Following the proof of Li et al., 2024, Proposition 1, similarly, we denote

$$\mathbf{W} = \begin{bmatrix} \bar{\mathbf{W}} & \mathbf{w}_1 \\ \mathbf{w}_2^\top & w \end{bmatrix} \quad \text{and} \quad \bar{\mathbf{h}} = \begin{bmatrix} \mathbf{h}_1 \\ h \end{bmatrix},$$

520 where $\bar{\mathbf{W}} \in \mathbb{R}^{d \times d}$, $\mathbf{w}_1, \mathbf{w}_2, \mathbf{h}_1 \in \mathbb{R}^d$, and $w, h \in \mathbb{R}$.

521 Additionally, let \mathcal{I} denote the token indices set corresponding to the labeled demonstrations (cf. 11).
522 Recall the prompt \mathbf{Z} from (2), and $\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$ and $\mathbf{y} = [y_1 \cdots y_n]^\top \in \mathbb{R}^n$ from (7).
523 Then we get

$$\mathbf{Z} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_n & \mathbf{x} \\ y_1 & y_2 & \cdots & y_n & 0 \end{bmatrix}^\top = \begin{bmatrix} \mathbf{X}^\top & \mathbf{x} \\ \mathbf{y}^\top & 0 \end{bmatrix}^\top \in \mathbb{R}^{(n+1) \times (d+1)}. \quad (16)$$

524 Combining (15) and (16) together, we can rewrite the one-layer linear prediction as

$$\begin{aligned} f_{\text{att-1}}(\mathbf{Z}) &= [\mathbf{x}^\top \ 0] \mathbf{W} \mathbf{Z}^\top \mathbf{M} \mathbf{Z} \bar{\mathbf{h}} \\ &= [\mathbf{x}^\top \ 0] \begin{bmatrix} \bar{\mathbf{W}} & \mathbf{w}_1 \\ \mathbf{w}_2^\top & w \end{bmatrix} \begin{bmatrix} \mathbf{X}^\top & \mathbf{x} \\ \mathbf{y}^\top & 0 \end{bmatrix} \begin{bmatrix} \mathbf{I}_n & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{X}^\top & \mathbf{x} \\ \mathbf{y}^\top & 0 \end{bmatrix}^\top \begin{bmatrix} \mathbf{h}_1 \\ h \end{bmatrix} \\ &= [\mathbf{x}^\top \bar{\mathbf{W}} \ \mathbf{x}^\top \mathbf{w}_1] \begin{bmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{y} \\ \mathbf{y}^\top \mathbf{X} & \mathbf{y}^\top \mathbf{y} \end{bmatrix} \begin{bmatrix} \mathbf{h}_1 \\ h \end{bmatrix} \\ &= [\mathbf{x}^\top \bar{\mathbf{W}} \ \mathbf{x}^\top \mathbf{w}_1] \begin{bmatrix} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1 + h \mathbf{X}^\top \mathbf{y} \\ \mathbf{y}^\top \mathbf{X} \mathbf{h}_1 + h \mathbf{y}^\top \mathbf{y} \end{bmatrix} \\ &= \mathbf{x}^\top \bar{\mathbf{W}} (\mathbf{X}^\top \mathbf{X} \mathbf{h}_1 + h \mathbf{X}^\top \mathbf{y}) + \mathbf{x}^\top \mathbf{w}_1 (\mathbf{y}^\top \mathbf{X} \mathbf{h}_1 + h \mathbf{y}^\top \mathbf{y}) \\ &= \mathbf{x}^\top (h \bar{\mathbf{W}} + \mathbf{w}_1 \mathbf{h}_1^\top) \mathbf{X}^\top \mathbf{y} + \mathbf{x}^\top (\bar{\mathbf{W}} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1 + h \mathbf{y}^\top \mathbf{y} \mathbf{w}_1) \\ &= \mathbf{x}^\top \tilde{\mathbf{W}} \mathbf{X}^\top \mathbf{y} + \mathbf{x}^\top (\bar{\mathbf{W}} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1 + m h \mathbf{w}_1) \end{aligned}$$

525 where $\tilde{\mathbf{W}} := h \bar{\mathbf{W}} + \mathbf{w}_1 \mathbf{h}_1^\top$ and we define $m := |\mathcal{I}|$.

526 Next, recall the loss from (5) and consider the squared loss function, $\ell(y, \hat{y}) = (y - \hat{y})^2$. We have

$$\begin{aligned} \mathcal{L}_{\text{att-1}}(\mathcal{W}, \mathbf{h}) &= \mathbb{E} \left[(f_{\text{att-1}}(\mathbf{Z}) - y)^2 \right] \\ &= \mathbb{E} \left[\left(\mathbf{x}^\top \tilde{\mathbf{W}} \mathbf{X} \mathbf{y} + \mathbf{x}^\top (\bar{\mathbf{W}} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1 + m h \mathbf{w}_1) - y \right)^2 \right] \\ &= \mathbb{E} \left[\left(y \mathbf{x}^\top \tilde{\mathbf{W}} \mathbf{X} \mathbf{y} + y \mathbf{x}^\top (\bar{\mathbf{W}} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1 + m h \mathbf{w}_1) - 1 \right)^2 \right]. \end{aligned}$$

527 For simplicity and without loss of generality, we omit y and use \mathbf{x} to represent $y\mathbf{x}$. Note that
528 the distribution of (updated) \mathbf{x} is not conditioned on its class and given mean vector $\boldsymbol{\mu}$, it follows
529 $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$. Similarly, let \mathbf{x}_i represent $y_i^c \mathbf{x}_i$. We can then write

$$\begin{aligned} \mathcal{L}_{\text{att-1}}(\mathcal{W}, \mathbf{h}) &= \mathbb{E} \left[\left(\mathbf{x}^\top \tilde{\mathbf{W}} \sum_{i \in \mathcal{I}} \mathbf{x}_i + \mathbf{x}^\top (\bar{\mathbf{W}} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1 + m h \mathbf{w}_1) - 1 \right)^2 \right] \\ &= \mathbb{E} \left[\left(\mathbf{x}^\top \tilde{\mathbf{W}} \sum_{i \in \mathcal{I}} \mathbf{x}_i - 1 \right)^2 \right] + \mathbb{E} \left[\left(\mathbf{x}^\top (\bar{\mathbf{W}} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1 + m h \mathbf{w}_1) \right)^2 \right] \\ &\quad + 2 \mathbb{E} \left[\left(\mathbf{x}^\top \tilde{\mathbf{W}} \sum_{i \in \mathcal{I}} \mathbf{x}_i - 1 \right) \left(\mathbf{x}^\top (\bar{\mathbf{W}} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1 + m h \mathbf{w}_1) \right) \right]. \end{aligned} \quad (17)$$

530 We start with showing that for any given parameters $\mathbf{W} \in \mathbb{R}^{(d+1) \times (d+1)}$, $\mathbf{h} \in \mathbb{R}^{d+1}$,
 531 $\mathbb{E}\left[\left(\mathbf{x}^\top \tilde{\mathbf{W}} \sum_{i \in \mathcal{I}} \mathbf{x}_i - 1\right) \left(\mathbf{x}^\top (\tilde{\mathbf{W}} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1 + m \mathbf{h} \mathbf{w}_1)\right)\right] = 0$. To prove it, we first expand

$$\begin{aligned} & (\mathbf{x}^\top \tilde{\mathbf{W}} \sum_{i \in \mathcal{I}} \mathbf{x}_i - 1)(\mathbf{x}^\top (\tilde{\mathbf{W}} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1 + m \mathbf{h} \mathbf{w}_1)) \\ &= \underbrace{(\mathbf{x}^\top \tilde{\mathbf{W}} \sum_{i \in \mathcal{I}} \mathbf{x}_i)(\mathbf{x}^\top \tilde{\mathbf{W}} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1)}_{(a)} - \underbrace{\mathbf{x}^\top \tilde{\mathbf{W}} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1}_{(b)} + \underbrace{(\mathbf{x}^\top \tilde{\mathbf{W}} \sum_{i \in \mathcal{I}} \mathbf{x}_i)(m \mathbf{h} \mathbf{x}^\top \mathbf{w}_1)}_{(c)} - \underbrace{m \mathbf{h} \mathbf{x}^\top \mathbf{w}_1}_{(d)}. \end{aligned}$$

532 In the following, we consider the expectations of (a), (b), (c), (d) sequentially, all of which take the
 533 value zero. First note that since $\boldsymbol{\mu} \sim \text{Unif}(\mathbb{S}^{d-1})$ and $(\boldsymbol{\xi}_i)_{i=1}^n, \boldsymbol{\xi} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$, the odd moments of $\boldsymbol{\mu}, \boldsymbol{\xi}$
 534 and $\boldsymbol{\xi}_i, i \in [n]$ are all zeros.

$$\begin{aligned} (a) : \quad & \mathbb{E}\left[\left(\mathbf{x}^\top \tilde{\mathbf{W}} \sum_{i \in \mathcal{I}} \mathbf{x}_i\right) (\mathbf{x}^\top \tilde{\mathbf{W}} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1)\right] \\ &= \mathbb{E}\left[\left(\boldsymbol{\mu} + \boldsymbol{\xi}\right)^\top \tilde{\mathbf{W}} \sum_{i \in \mathcal{I}} (\boldsymbol{\mu} + \boldsymbol{\xi}_i) (\boldsymbol{\mu} + \boldsymbol{\xi})^\top \tilde{\mathbf{W}} \sum_{i \in [n]} (\boldsymbol{\mu} + \boldsymbol{\xi}_i) (\boldsymbol{\mu} + \boldsymbol{\xi}_i)^\top \mathbf{h}_1\right] \\ &= \sum_{i \in \mathcal{I}} \sum_{j \in [n]} \mathbb{E}\left[(\boldsymbol{\mu} + \boldsymbol{\xi})^\top \tilde{\mathbf{W}} (\boldsymbol{\mu} + \boldsymbol{\xi}_i) (\boldsymbol{\mu} + \boldsymbol{\xi})^\top \tilde{\mathbf{W}} (\boldsymbol{\mu} + \boldsymbol{\xi}_j) (\boldsymbol{\mu} + \boldsymbol{\xi}_j)^\top \mathbf{h}_1\right] \\ &= \sum_{i \in \mathcal{I}} \sum_{j \in [n]} \mathbb{E}\left[\boldsymbol{\mu}^\top \tilde{\mathbf{W}} \boldsymbol{\mu} \boldsymbol{\mu}^\top \tilde{\mathbf{W}} (\boldsymbol{\mu} \boldsymbol{\mu}^\top + \boldsymbol{\xi}_j \boldsymbol{\xi}_j^\top) \mathbf{h}_1 + \boldsymbol{\xi}^\top \tilde{\mathbf{W}} \boldsymbol{\mu} \boldsymbol{\xi}^\top \tilde{\mathbf{W}} (\boldsymbol{\mu} \boldsymbol{\mu}^\top + \boldsymbol{\xi}_j \boldsymbol{\xi}_j^\top) \mathbf{h}_1\right] \\ &= 0, \end{aligned}$$

535

$$\begin{aligned} (b) : \quad & \mathbb{E}\left[\mathbf{x}^\top \tilde{\mathbf{W}} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1\right] \\ &= \mathbb{E}\left[(\boldsymbol{\mu} + \boldsymbol{\xi})^\top \tilde{\mathbf{W}} \sum_{i \in [n]} (\boldsymbol{\mu} + \boldsymbol{\xi}_i) (\boldsymbol{\mu} + \boldsymbol{\xi}_i)^\top \mathbf{h}_1\right] \\ &= \mathbb{E}\left[\boldsymbol{\mu}^\top \tilde{\mathbf{W}} \sum_{i \in [n]} (\boldsymbol{\mu} \boldsymbol{\mu}^\top + \boldsymbol{\xi}_i \boldsymbol{\xi}_i^\top) \mathbf{h}_1\right] \\ &= 0, \end{aligned}$$

536

$$\begin{aligned} (c) : \quad & \mathbb{E}\left[\left(\mathbf{x}^\top \tilde{\mathbf{W}} \sum_{i \in \mathcal{I}} \mathbf{x}_i\right) (m \mathbf{h} \mathbf{x}^\top \mathbf{w}_1)\right] \\ &= m \mathbf{h} \mathbb{E}\left[(\boldsymbol{\mu} + \boldsymbol{\xi})^\top \tilde{\mathbf{W}} \sum_{i \in \mathcal{I}} (\boldsymbol{\mu} + \boldsymbol{\xi}_i) (\boldsymbol{\mu} + \boldsymbol{\xi})^\top \mathbf{w}_1\right] \\ &= m \mathbf{h} \sum_{i \in \mathcal{I}} \mathbb{E}\left[(\boldsymbol{\mu} + \boldsymbol{\xi})^\top \tilde{\mathbf{W}} \boldsymbol{\mu} (\boldsymbol{\mu} + \boldsymbol{\xi})^\top \mathbf{w}_1\right] \\ &= m \mathbf{h} \sum_{i \in \mathcal{I}} \mathbb{E}\left[\boldsymbol{\mu}^\top \tilde{\mathbf{W}} \boldsymbol{\mu} \boldsymbol{\mu}^\top \mathbf{w}_1 + \boldsymbol{\xi}^\top \tilde{\mathbf{W}} \boldsymbol{\mu} \boldsymbol{\xi}^\top \mathbf{w}_1\right] \\ &= 0, \end{aligned}$$

537

$$(d) : \quad \mathbb{E}\left[m \mathbf{h} \mathbf{x}^\top \mathbf{w}_1\right] = 0.$$

538 Therefore, loss in (17) returns

$$\mathcal{L}_{\text{att-1}}(\mathcal{W}, \mathbf{h}) = \mathbb{E}\left[\underbrace{\left(\mathbf{x}^\top \tilde{\mathbf{W}} \sum_{i \in \mathcal{I}} \mathbf{x}_i - 1\right)^2}_{\tilde{\mathcal{L}}(\tilde{\mathbf{W}})}\right] + \mathbb{E}\left[\left(\mathbf{x}^\top (\tilde{\mathbf{W}} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1 + m \mathbf{h} \mathbf{w}_1)\right)^2\right].$$

Here, the first term $\mathbb{E}[(\mathbf{x}^\top \tilde{\mathbf{W}} \sum_{i \in \mathcal{I}} \mathbf{x}_i - 1)^2] = \tilde{\mathcal{L}}(\tilde{\mathbf{W}})$ where $\tilde{\mathcal{L}}(\tilde{\mathbf{W}})$ is defined in (12). Recall that $\tilde{\mathbf{W}} = h\bar{\mathbf{W}} + \mathbf{w}_1 \mathbf{h}_1^\top$. Then for any $\tilde{\mathbf{W}} \in \mathbb{R}^{d \times d}$, setting $\mathbf{h}_1 = \mathbf{w}_1 = \mathbf{0}_d$ and $h = 1$ returns $\mathbb{E}[(\mathbf{x}^\top (\tilde{\mathbf{W}} \mathbf{X}^\top \mathbf{X} \mathbf{h}_1 + m \mathbf{h} \mathbf{w}_1))^2] = 0$, and then

$$\mathcal{L}_{\text{att-1}}(\mathbf{W}, \mathbf{h}) = \mathbb{E} \left[\left(\mathbf{x}^\top \tilde{\mathbf{W}} \sum_{i \in \mathcal{I}} \mathbf{x}_i - 1 \right)^2 \right]$$

Therefore, optimizing $\mathcal{L}_{\text{att-1}}(\mathbf{W}, \mathbf{h})$ returns the same minima as optimizing $\tilde{\mathcal{L}}(\mathbf{W})$, which completes the proof of (13). Note that optimal loss $\mathcal{L}_{\text{att-1}}^*$ depends on the labeled data $i \in \mathcal{I}$ only.

Furthermore, since $\tilde{\mathcal{L}}(\mathbf{W})$ is strongly convex (see (18)), \mathbf{W}^* exists and is unique. Therefore, (13) and uniqueness of \mathbf{W}^* leads to the conclusion (14). ■

Lemma 3 Consider the objective defined in (12) with semi-supervised data following Section 2. Then the optimal solution \mathbf{W}^* satisfies

$$\mathbf{W}^* = c\mathbf{I}$$

for some $c > 0$.

Proof. Recap the Objective (12) and its optimal solution \mathbf{W}^* . Let \mathcal{I} be the index set corresponding the labeled in-context examples, and $|\mathcal{I}| = m$. Note that, m is also a random variable, independent of $\mathbf{x}_i, y_i^c, \mathbf{x}, y$.

As in the proof of Lemma 2, we use \mathbf{x} to represent $y\mathbf{x}$ and \mathbf{x}_i to represent $y_i^c \mathbf{x}_i$ for simplicity, where (updated) $\mathbf{x}_i, \mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$. Letting $\boldsymbol{\xi}', \boldsymbol{\xi}, \boldsymbol{\xi}_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ be independent, we obtain

$$\begin{aligned} \tilde{\mathcal{L}}(\mathbf{W}) &= \mathbb{E} \left[\left(\mathbf{x}^\top \mathbf{W} \sum_{i \in \mathcal{I}} \mathbf{x}_i - 1 \right)^2 \right] \\ &= \mathbb{E} \left[\left((\boldsymbol{\mu} + \boldsymbol{\xi})^\top \mathbf{W} \sum_{i \in \mathcal{I}} (\boldsymbol{\mu} + \boldsymbol{\xi}_i) - 1 \right)^2 \right] \\ &= \mathbb{E} \left[\left((\boldsymbol{\mu} + \boldsymbol{\xi})^\top \mathbf{W} (m\boldsymbol{\mu} + \sqrt{m}\boldsymbol{\xi}') - 1 \right)^2 \right] \\ &= \mathbb{E} \left[m^2 (\boldsymbol{\mu}^\top \mathbf{W} \boldsymbol{\mu})^2 + m (\boldsymbol{\mu}^\top \mathbf{W} \boldsymbol{\xi}')^2 + m^2 (\boldsymbol{\xi}'^\top \mathbf{W} \boldsymbol{\mu})^2 + m (\boldsymbol{\xi}'^\top \mathbf{W} \boldsymbol{\xi}')^2 + 1 \right] - 2 \mathbb{E} [m \boldsymbol{\mu}^\top \mathbf{W} \boldsymbol{\mu}] \\ &= \frac{\mathbb{E}[m^2]}{d(d+2)} (\text{tr}(\mathbf{W})^2 + \text{tr}(\mathbf{W}\mathbf{W}^\top) + \text{tr}(\mathbf{W}^2)) + \frac{\mathbb{E}[m+m^2]}{d} \sigma^2 \text{tr}(\mathbf{W}\mathbf{W}^\top) \\ &\quad + \mathbb{E}[m] \sigma^4 \text{tr}(\mathbf{W}\mathbf{W}^\top) + 1 - \frac{2 \mathbb{E}[m]}{d} \text{tr}(\mathbf{W}). \end{aligned} \tag{18}$$

Differentiating it results in

$$\nabla_{\mathbf{W}} \tilde{\mathcal{L}}(\mathbf{W}) = \frac{2 \mathbb{E}[m^2]}{d(d+2)} (\text{tr}(\mathbf{W})\mathbf{I} + \mathbf{W} + \mathbf{W}^\top) + \frac{2 \mathbb{E}[m+m^2]\sigma^2}{d} \mathbf{W} + 2 \mathbb{E}[m] \sigma^4 \mathbf{W} - \frac{2 \mathbb{E}[m]}{d} \mathbf{I}.$$

Setting $\nabla_{\mathbf{W}} \tilde{\mathcal{L}}(\mathbf{W}) = 0$, we obtain the unique optimal \mathbf{W}^*

$$\mathbf{W}^* = \frac{1}{(1 + \sigma^2) \mathbb{E}[m^2] / \mathbb{E}[m] + \sigma^2 + \sigma^4 d} \mathbf{I},$$

which leads to the conclusion that $\mathbf{W}^* = c\mathbf{I}$, for $c = \frac{1}{(1 + \sigma^2) \mathbb{E}[m^2] / \mathbb{E}[m] + \sigma^2 + \sigma^4 d} > 0$. It completes the proof. ■

A.2 Proof of Theorem 1

Note that Theorem 1 has been slightly modified by adding a constant factor of 8.

Proof. Note that (6) can be easily proven using Lemmas 2 and 3. Then, we focus on proving (SPI-ERR).

562 Given that (6) holds, we can rewrite its classification error as

$$\mathbb{P}(y_{\text{att-1}}^*(\mathbf{Z}) \neq y) = \mathbb{P}(\text{sgn}(\mathbf{x}^\top \hat{\boldsymbol{\mu}}_s) \neq y) = \mathbb{P}(\text{sgn}(y\mathbf{x}^\top \hat{\boldsymbol{\mu}}_s) \neq 1) \quad (19)$$

563 where $\hat{\boldsymbol{\mu}}_s = \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} y_i \mathbf{x}_i$ defined in (SPI) and \mathcal{I} is the index set of labeled samples. Let $m = |\mathcal{I}|$.

564 Recall from Section 2.1 where $\mathbf{x} \sim \mathcal{N}(y \cdot \boldsymbol{\mu}, \sigma^2 \mathbf{I})$. We can rewrite

$$y\mathbf{x} = \boldsymbol{\mu} + \sigma \mathbf{g}_1 \quad \text{where} \quad \mathbf{g}_1 \sim \mathcal{N}(0, \mathbf{I}).$$

565 Then for any given $\boldsymbol{\mu}, \hat{\boldsymbol{\mu}}_s$, we get

$$\begin{aligned} \mathbb{P}(\text{sgn}(y\mathbf{x}^\top \hat{\boldsymbol{\mu}}_s) \neq 1 \mid \boldsymbol{\mu}, \hat{\boldsymbol{\mu}}_s) &= \mathbb{P}((\boldsymbol{\mu} + \sigma \mathbf{g}_1)^\top \hat{\boldsymbol{\mu}}_s < 0 \mid \boldsymbol{\mu}, \hat{\boldsymbol{\mu}}_s) \\ &= \mathbb{P}(\boldsymbol{\mu}^\top \hat{\boldsymbol{\mu}}_s < \sigma \mathbf{g}_1^\top \hat{\boldsymbol{\mu}}_s \mid \boldsymbol{\mu}, \hat{\boldsymbol{\mu}}_s) \\ &= Q\left(\frac{\boldsymbol{\mu}^\top \hat{\boldsymbol{\mu}}_s}{\sigma \|\hat{\boldsymbol{\mu}}_s\|_{\ell_2}}\right). \end{aligned} \quad (20)$$

566 Here Q -function is the tail distribution function of the standard normal distribution.

567 Next, similarly, given that $\mathbf{x}_i \sim \mathcal{N}(y_i \cdot \boldsymbol{\mu}, \sigma^2 \mathbf{I})$ for $i \in \mathcal{I}$, we can rewrite

$$\hat{\boldsymbol{\mu}}_s = \frac{1}{m} \sum_{i \in \mathcal{I}} y_i \mathbf{x}_i = \boldsymbol{\mu} + \frac{\sigma}{\sqrt{m}} \mathbf{g}_2 \quad \text{where} \quad \mathbf{g}_2 \sim \mathcal{N}(0, \mathbf{I}).$$

568 Then combining (19) and (20), we have

$$\begin{aligned} \mathbb{P}(y_{\text{att-1}}^*(\mathbf{Z}) \neq y) &= \mathbb{E}_{\boldsymbol{\mu}, \mathbf{g}_2} \left[Q\left(\frac{\boldsymbol{\mu}^\top \hat{\boldsymbol{\mu}}_s}{\sigma \|\hat{\boldsymbol{\mu}}_s\|_{\ell_2}}\right) \right] \\ &= \mathbb{E}_{\boldsymbol{\mu}, \mathbf{g}_2} \left[Q\left(\frac{\boldsymbol{\mu}^\top (\boldsymbol{\mu} + \frac{\sigma}{\sqrt{m}} \mathbf{g}_2)}{\sigma \left\| \boldsymbol{\mu} + \frac{\sigma}{\sqrt{m}} \mathbf{g}_2 \right\|_{\ell_2}}\right) \right] \\ &= \mathbb{E}_{\boldsymbol{\mu}, \mathbf{g}_2} \left[Q\left(\frac{1 + \frac{\sigma}{\sqrt{m}} \boldsymbol{\mu}^\top \mathbf{g}_2}{\sigma \sqrt{1 + 2 \frac{\sigma}{\sqrt{m}} \boldsymbol{\mu}^\top \mathbf{g}_2 + \frac{\sigma^2}{m} \|\mathbf{g}_2\|_{\ell_2}^2}}\right) \right]. \end{aligned}$$

569 Note that for any $\boldsymbol{\mu}$ with $\|\boldsymbol{\mu}\|_{\ell_2} = 1$, we have $\boldsymbol{\mu}^\top \mathbf{g}_2 \sim \mathcal{N}(0, 1)$. Therefore, we can write

$$\boldsymbol{\mu}^\top \mathbf{g}_2 = g \quad \text{where} \quad g \sim \mathcal{N}(0, 1),$$

570 and let $\mathbf{U} \in \mathbb{R}^{d \times d}$ be a unitary matrix with first row being $\boldsymbol{\mu}$. We can write

$$\|\mathbf{g}_2\|_{\ell_2}^2 = \|\mathbf{U} \mathbf{g}_2\|_{\ell_2}^2 = g^2 + h \quad \text{where} \quad h \sim \chi_{d-1}^2.$$

571 Here, χ_{d-1}^2 denotes chi-squared distribution with $(d-1)$ degrees of freedom. Then, we get

$$\begin{aligned} \mathbb{P}(y_{\text{att-1}}^*(\mathbf{Z}) \neq y) &= \mathbb{E}_{g, h} \left[Q\left(\frac{1 + \frac{\sigma}{\sqrt{m}} g}{\sigma \sqrt{1 + 2 \frac{\sigma}{\sqrt{m}} g + \frac{\sigma^2}{m} (g^2 + h)}}\right) \right] \\ &= \mathbb{E}_{g, h} \left[Q\left(\frac{1 + \frac{\sigma}{\sqrt{m}} g}{\sigma \sqrt{(1 + \frac{\sigma}{\sqrt{m}} g)^2 + \frac{\sigma^2}{m} h}}\right) \right], \\ &= \mathbb{E}_{g, h} \left[Q\left(\frac{1 + \varepsilon_\sigma g}{\sigma \sqrt{(1 + \varepsilon_\sigma g)^2 + \varepsilon_\sigma^2 h}}\right) \right], \end{aligned}$$

572 where $\varepsilon_\sigma := \sigma / \sqrt{m}$. It completes the proof of (SPI-ERR).

573 Next, we derive an upper bound for $\mathbb{P}(y_{\text{att-1}}^*(\mathbf{Z}) \neq y)$. Let $c := \varepsilon_\sigma^{-1}$. Then we have

$$\begin{aligned}
\mathbb{P}(y_{\text{att-1}}^*(\mathbf{Z}) \neq y) &= \mathbb{E}_{g,h} \left[Q \left(\frac{c+g}{\sigma \sqrt{(c+g)^2 + h}} \right) \right] \\
&= \mathbb{E}_{g \geq -\frac{c}{2}, h} \left[Q \left(\frac{c+g}{\sigma \sqrt{(c+g)^2 + h}} \right) \right] + \mathbb{E}_{g < -\frac{c}{2}, h} \left[Q \left(\frac{c+g}{\sigma \sqrt{(c+g)^2 + h}} \right) \right] \\
&\leq \mathbb{E}_{g \geq -\frac{c}{2}, h} \left[Q \left(\frac{c+g}{\sigma \sqrt{(c+g)^2 + h}} \right) \right] + Q(c/2) \\
&= \mathbb{E}_{g \geq -\frac{c}{2}, h} \left[Q \left(\frac{1}{\sigma \sqrt{1 + h/(c+g)^2}} \right) \right] + Q(c/2), \tag{21}
\end{aligned}$$

574 where the inequality comes from the fact that $\mathbb{P}(g \leq -c/2) = Q(c/2)$ and $Q(x) \leq 1$ for any $x \in \mathbb{R}$, and
575 we have

$$\frac{1}{\sqrt{1 + h/(c+g)^2}} \geq 1 - \frac{1}{2} \frac{h}{(c+g)^2} \geq 1 - \frac{2h}{c^2}.$$

576 Here the first inequality comes from that $\frac{1}{\sqrt{1+x}} \geq 1 - \frac{1}{2}x$ and the second utilizes that $g \geq -\frac{c}{2}$.

577 Since $h \sim \chi_{d-1}^2$, from the Laurent-Massart inequality (Laurent & Massart, 2000), we have that

$$\mathbb{P}(h \geq d - 1 + 2\sqrt{(d-1)t_1} + 2t_1) \leq e^{-t_1}.$$

578 Therefore, we have that with probability at least $1 - e^{-t_1}$

$$\frac{1}{\sqrt{1 + h/(c+g)^2}} \geq 1 - \frac{2(d-1 + 2\sqrt{(d-1)t_1} + 2t_1)}{c^2}.$$

579 Setting $t_1 = d$, we get

$$\frac{1}{\sqrt{1 + h/(c+g)^2}} \geq 1 - \frac{10d}{c^2}.$$

580 Combining the result with (21), since $Q(x) \leq 1$ for $x \in \mathbb{R}$ and $Q(x) \leq e^{-x^2/2}$ for $x > 1$, we get that

$$\begin{aligned}
\mathbb{P}(y_{\text{att-1}}^*(\mathbf{Z}) \neq y) &\leq e^{-d} + Q(c/2) + Q \left(\frac{1}{\sigma} \left(1 - \frac{10d}{c^2} \right) \right) \\
&\leq e^{-d} + e^{-1/8\varepsilon_\sigma^2} + Q \left(\frac{1}{\sigma} \left(1 - 10d\varepsilon_\sigma^2 \right) \right).
\end{aligned}$$

581 It completes the proof.

582 ■

583 B Analysis of Multi-layer Linear Attention

584 B.1 Proof of Proposition 1

585 **Proof.** We consider the following model constructions for the attention matrices in the ℓ 'th layer,
586 $\ell \in [L]$ and the final linear prediction head:

$$\begin{aligned}
\ell\text{'th layer: } \mathbf{W}_{q\ell} \mathbf{W}_{k\ell}^\top &= \begin{bmatrix} \mathbf{I}_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \quad \text{and} \quad \mathbf{W}_{v\ell} = \begin{bmatrix} a_\ell \mathbf{I}_d & \mathbf{0} \\ \mathbf{0} & b_\ell \end{bmatrix}; \\
\text{Prediction head: } \mathbf{h} &= \begin{bmatrix} \mathbf{0}_d \\ c \end{bmatrix}. \tag{22}
\end{aligned}$$

587 Suppose the input to ℓ 'th layer is

$$\mathbf{Z}_\ell = \begin{bmatrix} \mathbf{X}_\ell & \mathbf{y}_\ell \\ \mathbf{x}_\ell^\top & y_\ell \end{bmatrix} \in \mathbb{R}^{(n+1) \times (d+1)} \quad \text{where} \quad \mathbf{Z}_1 = \mathbf{Z} = \begin{bmatrix} \mathbf{X} & \mathbf{y} \\ \mathbf{x}^\top & 0 \end{bmatrix}.$$

Recapping the model construction from (22), the ℓ 'th layer output returns

$$\begin{aligned}
(\mathbf{Z}_\ell \mathbf{W}_{q\ell} \mathbf{W}_{k\ell}^\top \mathbf{Z}_\ell^\top \mathbf{M}) \mathbf{Z}_\ell \mathbf{W}_{v\ell} &= \begin{bmatrix} \mathbf{X}_\ell & \mathbf{y}_\ell \\ \mathbf{x}_\ell^\top & y_\ell \end{bmatrix} \begin{bmatrix} \mathbf{I}_d & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{X}_\ell^\top & \mathbf{x}_\ell \\ \mathbf{y}_\ell^\top & y_\ell \end{bmatrix} \mathbf{M} \begin{bmatrix} \mathbf{X}_\ell & \mathbf{y}_\ell \\ \mathbf{x}_\ell^\top & y_\ell \end{bmatrix} \begin{bmatrix} a_\ell \mathbf{I}_d & 0 \\ 0 & b_\ell \end{bmatrix} \\
&= \begin{bmatrix} \mathbf{X}_\ell \mathbf{X}_\ell^\top & \mathbf{X}_\ell \mathbf{x}_\ell \\ \mathbf{x}_\ell^\top \mathbf{X}_\ell^\top & \mathbf{x}_\ell^\top \mathbf{x}_\ell \end{bmatrix} \begin{bmatrix} \mathbf{I}_n & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} a_\ell \mathbf{X}_\ell & b_\ell \mathbf{y}_\ell \\ a_\ell \mathbf{x}_\ell^\top & b_\ell y_\ell \end{bmatrix} \\
&= \begin{bmatrix} a_\ell \mathbf{X}_\ell \mathbf{X}_\ell^\top \mathbf{X}_\ell & b_\ell \mathbf{X}_\ell \mathbf{X}_\ell^\top \mathbf{y}_\ell \\ a_\ell \mathbf{x}_\ell^\top \mathbf{X}_\ell^\top \mathbf{X}_\ell & b_\ell \mathbf{x}_\ell^\top \mathbf{X}_\ell^\top \mathbf{y}_\ell \end{bmatrix}.
\end{aligned} \tag{23}$$

Therefore, the input of $(\ell + 1)$ 'th layer is

$$\begin{aligned}
\mathbf{Z}_{\ell+1} &= \mathbf{Z}_\ell + \begin{bmatrix} a_\ell \mathbf{X}_\ell \mathbf{X}_\ell^\top \mathbf{X}_\ell & b_\ell \mathbf{X}_\ell \mathbf{X}_\ell^\top \mathbf{y}_\ell \\ a_\ell \mathbf{x}_\ell^\top \mathbf{X}_\ell^\top \mathbf{X}_\ell & b_\ell \mathbf{x}_\ell^\top \mathbf{X}_\ell^\top \mathbf{y}_\ell \end{bmatrix} \\
&= \begin{bmatrix} \mathbf{X}_\ell + a_\ell \mathbf{X}_\ell \mathbf{X}_\ell^\top \mathbf{X}_\ell & \mathbf{y}_\ell + b_\ell \mathbf{X}_\ell \mathbf{X}_\ell^\top \mathbf{y}_\ell \\ \mathbf{x}_\ell^\top + a_\ell \mathbf{x}_\ell^\top \mathbf{X}_\ell^\top \mathbf{X}_\ell & y_\ell + b_\ell \mathbf{x}_\ell^\top \mathbf{X}_\ell^\top \mathbf{y}_\ell \end{bmatrix} \in \mathbb{R}^{(n+1) \times (d+1)}.
\end{aligned} \tag{24}$$

• **Label propagation:** We first focus on deriving label propagation results. Suppose that we have

$$a_\ell = 0 \quad \text{for } \ell \in [L].$$

Then following (23), the output of ℓ 'th layer takes the following form:

$$(\mathbf{Z}_\ell \mathbf{W}_{q\ell} \mathbf{W}_{k\ell}^\top \mathbf{Z}_\ell^\top \mathbf{M}) \mathbf{Z}_\ell \mathbf{W}_{v\ell} = \begin{bmatrix} 0 & b_\ell \mathbf{X}_\ell \mathbf{X}_\ell^\top \mathbf{y}_\ell \\ 0 & b_\ell \mathbf{x}_\ell^\top \mathbf{X}_\ell^\top \mathbf{y}_\ell \end{bmatrix}.$$

Here, the first d coordinates of each token's output are zeros, and therefore, the corresponding input coordinates remain unchanged, and we have

$$\mathbf{X}_\ell \equiv \mathbf{X} \quad \text{and} \quad \mathbf{x}_\ell \equiv \mathbf{x} \quad \text{for } \ell \in [L].$$

The prediction (based on the last token output and after applying prediction head) is given by

$$f_{\text{all-L}}(\mathbf{Z}) = c b_L \mathbf{x}^\top \mathbf{X}^\top \mathbf{y}_L. \tag{25}$$

We next focus on obtaining \mathbf{y}_L . From (24), we have

$$\mathbf{y}_{\ell+1} = \mathbf{y}_\ell + b_\ell \mathbf{X} \mathbf{X}^\top \mathbf{y}_\ell = (\mathbf{I} + b_\ell \mathbf{X} \mathbf{X}^\top) \mathbf{y}_\ell.$$

Therefore,

$$\mathbf{y}_L = \prod_{\ell=1}^{L-1} (\mathbf{I} + b_\ell \mathbf{X} \mathbf{X}^\top) \mathbf{y}.$$

Combining with (25) results in

$$f_{\text{all-L}}(\mathbf{Z}) = c b_L \mathbf{x}^\top \mathbf{X}^\top \prod_{\ell=1}^{L-1} (\mathbf{I} + b_\ell \mathbf{X} \mathbf{X}^\top) \mathbf{y} = c b_L \mathbf{x}^\top \prod_{\ell=1}^{L-1} (\mathbf{I} + b_\ell \mathbf{X}^\top \mathbf{X}) \mathbf{X}^\top \mathbf{y}.$$

It completes the proof.

• **Feature propagation:** We now focus on the feature propagation setting. In contrast to the label propagation, let us assume that

$$a_\ell \rightarrow \infty \quad \text{and} \quad b_\ell \rightarrow 0^+ \quad \text{for } \ell \in [L].$$

The prediction (following (23), based on the last token output and after applying prediction head) is given by

$$f_{\text{all-L}}(\mathbf{Z}) = c b_L \mathbf{x}_L^\top \mathbf{X}_L^\top \mathbf{y}_L. \tag{26}$$

We first obtain \mathbf{y}_L . From (24) (since $b_\ell \rightarrow 0$), we have

$$\mathbf{y}_{\ell+1} = \mathbf{y}_\ell + b_\ell \mathbf{X} \mathbf{X}^\top \mathbf{y}_\ell = \mathbf{y}_\ell.$$

604 Therefore,

$$\mathbf{y}_L = \mathbf{y}.$$

605 Next, we focus on $\mathbf{X}_L, \mathbf{x}_L$. From (24), as $a_\ell \rightarrow \infty$, we have

$$\begin{aligned}\mathbf{X}_{\ell+1} &= \mathbf{X}_\ell + a_\ell \mathbf{X}_\ell \mathbf{X}_\ell^\top \mathbf{X}_\ell = \mathbf{X}_\ell (\mathbf{I} + a_\ell \mathbf{X}_\ell^\top \mathbf{X}_\ell) = a_\ell \mathbf{X}_\ell \mathbf{X}_\ell^\top \mathbf{X}_\ell; \\ \mathbf{x}_{\ell+1}^\top &= \mathbf{x}_\ell^\top + a_\ell \mathbf{x}_\ell^\top \mathbf{X}_\ell^\top \mathbf{X}_\ell = \mathbf{x}_\ell^\top (\mathbf{I} + a_\ell \mathbf{X}_\ell^\top \mathbf{X}_\ell) = a_\ell \mathbf{x}_\ell^\top \mathbf{X}_\ell^\top \mathbf{X}_\ell.\end{aligned}$$

606 Therefore,

$$\begin{aligned}\mathbf{X}_L &= a_{L-1} \mathbf{X}_{L-1} (\mathbf{X}_{L-1}^\top \mathbf{X}_{L-1}) \\ &= a_{L-1} a_{L-2}^3 \mathbf{X}_{L-2} (\mathbf{X}_{L-2}^\top \mathbf{X}_{L-2})^{\frac{3^2-1}{2}} \\ &= a_{L-1} a_{L-2}^3 a_{L-3}^{3^2} \mathbf{X}_{L-3} (\mathbf{X}_{L-3}^\top \mathbf{X}_{L-3})^{\frac{3^3-1}{2}} \\ &= \dots \\ &= a_{L-1} a_{L-2}^3 a_{L-3}^{3^2} \dots a_1^{3^{L-2}} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{\frac{3^{L-1}-1}{2}},\end{aligned}$$

607 and

$$\begin{aligned}\mathbf{x}_L^\top &= a_{L-1} \mathbf{x}_{L-1}^\top (\mathbf{X}_{L-1}^\top \mathbf{X}_{L-1}) \\ &= a_{L-1} a_{L-2}^3 \mathbf{x}_{L-2}^\top (\mathbf{X}_{L-2}^\top \mathbf{X}_{L-2})^{\frac{3^2-1}{2}} \\ &= a_{L-1} a_{L-2}^3 a_{L-3}^{3^2} \mathbf{x}_{L-3}^\top (\mathbf{X}_{L-3}^\top \mathbf{X}_{L-3})^{\frac{3^3-1}{2}} \\ &= \dots \\ &= a_{L-1} a_{L-2}^3 a_{L-3}^{3^2} \dots a_1^{3^{L-2}} \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{\frac{3^{L-1}-1}{2}}.\end{aligned}$$

608 Combining all together with (26), we have that

$$\begin{aligned}f_{\text{all-}L}(\mathbf{Z}) &= c b_L \mathbf{x}_L^\top \mathbf{X}_L^\top \mathbf{y}_L \\ &= c b_L \left(\prod_{\ell=1}^{L-1} a_\ell^{3^{L-1-\ell}} \right)^2 \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{3^{L-1}-1} \mathbf{X}^\top \mathbf{y}.\end{aligned}$$

609 It completes the proof. ■

610 B.2 Proof of Proposition 2

611 **Proof.** The proof follows directly by adopting the same model construction and proof strategy as in
612 Proposition 1, under the additional assumption that

$$a_\ell = a \quad \text{and} \quad b_\ell = b \quad \text{for} \quad \ell \in [L].$$

613 ■

614 B.3 Proof of Lemma 1

615 **Proof.** In the proof of Proposition 1, we showed how to derive the label and feature propagation
616 results by restricting the construction to either $a_\ell \equiv 0$ (for label propagation) or $(a_\ell \rightarrow \infty, b_\ell \rightarrow 0)$
617 (for feature propagation). Here, we consider a propagation process without imposing restrictions on
618 the choices of (a_ℓ, b_ℓ) , and study the form of the final prediction returned by the model.

619 To avoid the notation conflict, we express the matrix \mathbf{A} in (9) as

$$\mathbf{A} = \sum_{k=0}^K e_k (\mathbf{X}^\top \mathbf{X})^k$$

620 and let $\mathbf{e} = [e_0 \ e_2 \ \dots \ e_{(3^{L-3}-1)/2}]^\top \in \mathbb{R}^{K+1}$.

621 Recall the same model construction used in the proof of Proposition 1, defined in (22). From (23),
 622 we have that

$$f_{\text{att-}L}(\mathbf{Z}) = cb_L \mathbf{x}_L^\top \mathbf{X}_L^\top \mathbf{y}_L$$

623 where following (24), we have

$$\begin{aligned} \mathbf{X}_{\ell+1} &= (\mathbf{I} + a_\ell \mathbf{X}_\ell \mathbf{X}_\ell^\top) \mathbf{X}_\ell, \\ \mathbf{x}_{\ell+1}^\top &= \mathbf{x}_\ell^\top (\mathbf{I} + a_\ell \mathbf{X}_\ell^\top \mathbf{X}_\ell), \\ \mathbf{y}_{\ell+1} &= (\mathbf{I} + b_\ell \mathbf{X}_\ell \mathbf{X}_\ell^\top) \mathbf{y}_\ell. \end{aligned}$$

624 At each layer, the operations performed are linear combinations and multiplications involving $\mathbf{X}_\ell^\top \mathbf{X}_\ell$
 625 and identity matrices scaled by the parameters (a_ℓ, b_ℓ) . Thus, each coefficient e_k of $(\mathbf{X}^\top \mathbf{X})^k$ depends
 626 smoothly on the scalar parameters (a_ℓ, b_ℓ) .

627 From (23) and (24), we have that

$$\begin{aligned} f_{\text{att-}L}(\mathbf{Z}) &= cb_L \mathbf{x}_L^\top \mathbf{X}_L^\top \mathbf{y}_L \\ &= cb_L \cdot \mathbf{x}_{L-1}^\top (\mathbf{I} + a_{L-1} \mathbf{X}_{L-1}^\top \mathbf{X}_{L-1}) (\mathbf{I} + a_{L-1} \mathbf{X}_{L-1}^\top \mathbf{X}_{L-1}) \mathbf{X}_{L-1}^\top \cdot (\mathbf{I} + b_{L-1} \mathbf{X}_{L-1} \mathbf{X}_{L-1}^\top) \mathbf{y}_{L-1} \\ &= \dots \end{aligned} \tag{27}$$

628 That is, in the final $f_{\text{att-}L}(\mathbf{Z})$ expression, the coefficients corresponding to different degrees of $(\mathbf{X}^\top \mathbf{X})^k$
 629 depend on the model parameters cb_L and $(a_\ell, b_\ell)_{\ell=1}^{L-1}$, which together have at most $2L - 1$ degrees
 630 of freedom. Let $\mathbf{c} = [cb_L \ a_1 \ \dots \ a_{L-1} \ b_1 \ \dots \ b_{L-1}]^\top$. This means there exists a smooth function
 631 $g : \mathbb{R}^{2L-1} \rightarrow \mathbb{R}^K$ such that: $\mathbf{e} = g(\mathbf{c})$.

632 It remains to show that an L -layer linear attention model can produce terms involving powers of $\mathbf{X}^\top \mathbf{X}$
 633 up to degree $(3^L - 3)/2$.

634 Let $f(\mathbf{Z})$ be a function that contains terms of the form $\mathbf{X}^\top (\mathbf{X}^\top \mathbf{X})^k \mathbf{X}^\top \mathbf{y}$ for various powers k . Define
 635 $\mathcal{P}(f(\mathbf{Z}))$ as the projection that extracts the highest degree k present in $f(\mathbf{Z})$. For example, $\mathcal{P}(\mathbf{x}^\top (\mathbf{I} +$
 636 $(\mathbf{X}^\top \mathbf{X})^2) \mathbf{X}^\top \mathbf{y}) = 2$. Then from (27), we have

$$\begin{aligned} \mathcal{P}(f_{\text{att-}L}(\mathbf{Z})) &= \mathcal{P}(\mathbf{x}_L^\top \mathbf{X}_L^\top \mathbf{y}_L) \\ &= \mathcal{P}(\mathbf{x}_{L-1}^\top (\mathbf{X}_{L-1}^\top \mathbf{X}_{L-1})^3 \mathbf{X}_{L-1}^\top \mathbf{y}_{L-1}) \\ &= \mathcal{P}(\mathbf{x}_{L-2}^\top (\mathbf{X}_{L-2}^\top \mathbf{X}_{L-2}) (\mathbf{X}_{L-2}^\top \mathbf{X}_{L-2})^3 (\mathbf{X}_{L-2}^\top \mathbf{X}_{L-2})^2 \mathbf{X}_{L-2}^\top \mathbf{y}_{L-2}) \\ &= \mathcal{P}(\mathbf{x}_{L-2}^\top (\mathbf{X}_{L-2}^\top \mathbf{X}_{L-2})^{3^2+3} \mathbf{X}_{L-2}^\top \mathbf{y}_{L-2}) \\ &= \mathcal{P}(\mathbf{x}_{L-3}^\top (\mathbf{X}_{L-3}^\top \mathbf{X}_{L-3}) (\mathbf{X}_{L-3}^\top \mathbf{X}_{L-3})^{3^3+3^2} (\mathbf{X}_{L-3}^\top \mathbf{X}_{L-3})^2 \mathbf{X}_{L-3}^\top \mathbf{y}_{L-3}) \\ &= \mathcal{P}(\mathbf{x}_{L-3}^\top (\mathbf{X}_{L-3}^\top \mathbf{X}_{L-3})^{3^3+3^2+3} \mathbf{X}_{L-3}^\top \mathbf{y}_{L-3}) \\ &= \dots \\ &= \mathcal{P}(\mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{3^{L-1}+\dots+3^2+3} \mathbf{X}^\top \mathbf{y}) \\ &= 3^{L-1} + \dots + 3^2 + 3 = \frac{3^L - 3}{2}. \end{aligned}$$

637 It completes the proof.

638 ■

639 B.4 Proof of Theorem 2

640 **Proof.** Let $\boldsymbol{\xi} \sim \mathcal{N}(0, \mathbf{I})$ and rewrite $\mathbf{y}\mathbf{x} = \boldsymbol{\mu} + \sigma\boldsymbol{\xi}$. For any matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$, the prediction error of
 641 $\hat{\mathbf{y}}_A = \text{sgn}(\mathbf{x}^\top \mathbf{A} \hat{\boldsymbol{\mu}}_s)$ given $\hat{\boldsymbol{\mu}}_s$ returns

$$\begin{aligned} \mathbb{P}(\hat{\mathbf{y}}_A \neq y \mid \hat{\boldsymbol{\mu}}_s) &= \mathbb{P}(\mathbf{y}\mathbf{x}^\top \mathbf{A} \hat{\boldsymbol{\mu}}_s < 0 \mid \hat{\boldsymbol{\mu}}_s) \\ &= \mathbb{P}((\boldsymbol{\mu} + \sigma\boldsymbol{\xi})^\top \mathbf{A} \hat{\boldsymbol{\mu}}_s < 0 \mid \hat{\boldsymbol{\mu}}_s) \\ &= \mathcal{Q}\left(\frac{\boldsymbol{\mu}^\top \mathbf{A} \hat{\boldsymbol{\mu}}_s}{\sigma \|\mathbf{A} \hat{\boldsymbol{\mu}}_s\|_{\ell_2}}\right). \end{aligned} \tag{28}$$

642 For any $\mathbf{A} \in \mathbb{R}^{d \times d}$, we can decompose it as

$$\mathbf{A} = \sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{v}_i^\top$$

643 where $\mathbf{u}_1 = \boldsymbol{\mu}$, $\|\mathbf{u}_i\|_{\ell_2} = 1$ and $\mathbf{u}_i^\top \mathbf{u}_j = 0$ for any $i \neq j$. Let $\lambda_1 > 0$. Then, we get

$$\begin{aligned} \boldsymbol{\mu}^\top \mathbf{A} \hat{\boldsymbol{\mu}}_s &= \boldsymbol{\mu}^\top \left(\sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{v}_i^\top \right) \hat{\boldsymbol{\mu}}_s \\ &= \sum_{i=1}^d \lambda_i \boldsymbol{\mu}^\top \mathbf{u}_i \mathbf{v}_i^\top \hat{\boldsymbol{\mu}}_s \\ &= \lambda_1 \boldsymbol{\mu}^\top \mathbf{u}_1 \mathbf{v}_1^\top \hat{\boldsymbol{\mu}}_s \\ &= \lambda_1 \mathbf{v}_1^\top \hat{\boldsymbol{\mu}}_s. \end{aligned} \quad (29)$$

644 Now consider $\|\mathbf{A} \hat{\boldsymbol{\mu}}_s\|_{\ell_2}$ where we have

$$\begin{aligned} \mathbf{A} \hat{\boldsymbol{\mu}}_s &= \sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{v}_i^\top \hat{\boldsymbol{\mu}}_s \\ &= \lambda_1 \boldsymbol{\mu} \mathbf{v}_1^\top \hat{\boldsymbol{\mu}}_s + \sum_{i=2}^d \lambda_i \mathbf{u}_i \mathbf{v}_i^\top \hat{\boldsymbol{\mu}}_s. \end{aligned}$$

645 Since $\mathbf{u}_i, i \neq 1$ is orthogonal to $\boldsymbol{\mu}$, $\lambda_1 \boldsymbol{\mu} \mathbf{v}_1^\top \hat{\boldsymbol{\mu}}_s$ is orthogonal to $\sum_{i=2}^d \lambda_i \mathbf{u}_i \mathbf{v}_i^\top \hat{\boldsymbol{\mu}}_s$. Therefore, given $\|\mathbf{u}_i\|_{\ell_2} = 1$
646 for all $i \in [d]$, it obeys

$$\|\mathbf{A} \hat{\boldsymbol{\mu}}_s\|_{\ell_2}^2 = \|\lambda_1 \boldsymbol{\mu} \mathbf{v}_1^\top \hat{\boldsymbol{\mu}}_s\|_{\ell_2}^2 + \sum_{i=2}^d \|\lambda_i \mathbf{u}_i \mathbf{v}_i^\top \hat{\boldsymbol{\mu}}_s\|_{\ell_2}^2 = (\lambda_1 \mathbf{v}_1^\top \hat{\boldsymbol{\mu}}_s)^2 + \lambda_1^2 \sum_{i=2}^d (\lambda_1^{-1} \lambda_i \mathbf{v}_i^\top \hat{\boldsymbol{\mu}}_s)^2. \quad (30)$$

647 For simplicity, define

$$\Delta(\hat{\boldsymbol{\mu}}_s) = \sum_{i=2}^d (\lambda_1^{-1} \lambda_i \mathbf{v}_i^\top \hat{\boldsymbol{\mu}}_s)^2$$

648 where we have

$$\Delta(\hat{\boldsymbol{\mu}}_s) \geq 0 \quad \text{and} \quad \Delta(-\hat{\boldsymbol{\mu}}_s) = \Delta(\hat{\boldsymbol{\mu}}_s).$$

649 Recall that $\hat{\boldsymbol{\mu}}_s$ is the SPI estimator (cf. (SPI)). Let $|I| = m$. We can write $\hat{\boldsymbol{\mu}}_s = \boldsymbol{\mu} + \frac{\sigma}{\sqrt{m}} \boldsymbol{\xi}'$ where
650 $\boldsymbol{\xi}' \sim \mathcal{N}(0, \mathbf{I})$.

651 Using (28), (29) and (30), the classification error becomes

$$\begin{aligned} \mathbb{P}(\hat{\mathbf{y}}_A \neq y) &= \mathbb{E}_{\hat{\boldsymbol{\mu}}_s} \left[\mathcal{Q} \left(\frac{\boldsymbol{\mu}^\top \mathbf{A} \hat{\boldsymbol{\mu}}_s}{\sigma \|\mathbf{A} \hat{\boldsymbol{\mu}}_s\|_{\ell_2}} \right) \right] \\ &= \mathbb{E}_{\hat{\boldsymbol{\mu}}_s} \left[\mathcal{Q} \left(\frac{\mathbf{v}_1^\top \hat{\boldsymbol{\mu}}_s}{\sigma \sqrt{(\mathbf{v}_1^\top \hat{\boldsymbol{\mu}}_s)^2 + \Delta(\hat{\boldsymbol{\mu}}_s)}} \right) \right]. \end{aligned}$$

652 First, note that for any $x > 0$, $\mathcal{Q}(x) < 0.5 < \mathcal{Q}(-x)$. Therefore, the optimal $\mathbf{v}_1 \in \mathbb{R}^d$ maximizes
653 $\mathbb{P}(\mathbf{v}_1^\top \hat{\boldsymbol{\mu}}_s > 0)$. Let $\mathbf{v}_1^* := \arg \max_{\mathbf{v}_1 \in \mathbb{R}^d} \mathbb{P}(\mathbf{v}_1^\top \hat{\boldsymbol{\mu}}_s > 0)$. Given that $\hat{\boldsymbol{\mu}}_s \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2/m\mathbf{I})$, we have that
654 $\mathbf{v}_1^* = c\boldsymbol{\mu}$ for $c > 0$. Let $c = 1$ and therefore, $\mathbf{v}_1^* = \boldsymbol{\mu}$ without loss of generality. Then we obtain

$$\mathbb{P}(\hat{\mathbf{y}}_A \neq y) = \mathbb{E}_{\hat{\boldsymbol{\mu}}_s} \left[\mathcal{Q} \left(\frac{\boldsymbol{\mu}^\top \hat{\boldsymbol{\mu}}_s}{\sigma \sqrt{(\boldsymbol{\mu}^\top \hat{\boldsymbol{\mu}}_s)^2 + \Delta(\hat{\boldsymbol{\mu}}_s)}} \right) \right].$$

655 Let $f(\hat{\boldsymbol{\mu}}_s)$ be the probability density function of $\hat{\boldsymbol{\mu}}_s$. Since $\hat{\boldsymbol{\mu}}_s \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2/m\mathbf{I})$, then it satisfies

$$f(\hat{\boldsymbol{\mu}}_s) \geq f(-\hat{\boldsymbol{\mu}}_s) \quad \text{for any } \boldsymbol{\mu}^\top \hat{\boldsymbol{\mu}}_s > 0. \quad (31)$$

Therefore, the classification error becomes

$$\begin{aligned}
\mathbb{P}(\hat{y}_A \neq y) &= \int_{\hat{\mu}_s} f(\hat{\mu}_s) Q\left(\frac{\mu^\top \hat{\mu}_s}{\sigma \sqrt{(\mu^\top \hat{\mu}_s)^2 + \Delta(\hat{\mu}_s)}}\right) d\hat{\mu}_s \\
&= \int_{\mu^\top \hat{\mu}_s > 0} f(\hat{\mu}_s) Q\left(\frac{\mu^\top \hat{\mu}_s}{\sigma \sqrt{(\mu^\top \hat{\mu}_s)^2 + \Delta(\hat{\mu}_s)}}\right) + f(-\hat{\mu}_s) Q\left(\frac{-\mu^\top \hat{\mu}_s}{\sigma \sqrt{(\mu^\top \hat{\mu}_s)^2 + \Delta(\hat{\mu}_s)}}\right) d\hat{\mu}_s \\
&= \int_{\mu^\top \hat{\mu}_s > 0} (f(\hat{\mu}_s) - f(-\hat{\mu}_s)) Q\left(\frac{\mu^\top \hat{\mu}_s}{\sigma \sqrt{(\mu^\top \hat{\mu}_s)^2 + \Delta(\hat{\mu}_s)}}\right) + f(-\hat{\mu}_s) d\hat{\mu}_s.
\end{aligned}$$

Following (31), to minimize the error, we need minimize $Q\left(\frac{\mu^\top \hat{\mu}_s}{\sigma \sqrt{(\mu^\top \hat{\mu}_s)^2 + \Delta(\hat{\mu}_s)}}\right)$ for $\mu^\top \hat{\mu}_s > 0$, which can be easily done by choosing $\lambda_i = 0$ for $i \geq 2$. Then we get $\Delta(\hat{\mu}_s) \equiv 0$. Therefore, the optimal solution set \mathcal{A}^* defined in Theorem 2 satisfies:

$$\mathcal{A}^* = \{\lambda_1 \mu \mu^\top \mid \lambda_1 > 0\}.$$

Combining all together, we obtain

$$\begin{aligned}
\mathbb{P}(\hat{y}_A \neq y) &= \int_{\mu^\top \hat{\mu}_s > 0} (f(\hat{\mu}_s) - f(-\hat{\mu}_s)) Q\left(\frac{1}{\sigma}\right) + f(-\hat{\mu}_s) d\hat{\mu}_s \\
&= \int_{\mu^\top \hat{\mu}_s > 0} f(\hat{\mu}_s) d\hat{\mu}_s \cdot Q\left(\frac{1}{\sigma}\right) + \int_{\mu^\top \hat{\mu}_s < 0} f(\hat{\mu}_s) d\hat{\mu}_s \cdot \left(1 - Q\left(\frac{1}{\sigma}\right)\right) \\
&= Q\left(-\frac{\sqrt{m}}{\sigma}\right) Q\left(\frac{1}{\sigma}\right) + Q\left(\frac{\sqrt{m}}{\sigma}\right) \left(1 - Q\left(\frac{1}{\sigma}\right)\right) \\
&= \left(1 - Q\left(\frac{\sqrt{m}}{\sigma}\right)\right) Q\left(\frac{1}{\sigma}\right) + Q\left(\frac{\sqrt{m}}{\sigma}\right) \left(1 - Q\left(\frac{1}{\sigma}\right)\right) \\
&= Q\left(\frac{1}{\sigma}\right) + Q\left(\frac{\sqrt{m}}{\sigma}\right) - 2Q\left(\frac{\sqrt{m}}{\sigma}\right) Q\left(\frac{1}{\sigma}\right).
\end{aligned}$$

It completes the proof. ■

B.5 Non-asymptotic Analysis

In Section 4 and Theorem 3, we showed that with infinitely many unlabeled samples, an L -layer linear attention model (for $L \geq 2$) can implement the predictor described in Theorem 2 with optimal \mathbf{A} choice, achieving the classification error given by (10). In this section, we turn to the non-asymptotic setting where n is finite, and analyze the model's performance under this regime.

Theorem 4 *Let the prompt \mathbf{Z} be generated as described in Section 2.2, and consider an L -layer linear attention model with $L \geq 2$. Let $\hat{\mu}_s$ be the SPI estimator defined in (SPI), and denote the optimal prediction as $y_{\text{att-}L}^*(\mathbf{Z})$. Additionally, suppose that the number of labeled samples satisfies $m := np \geq d\sigma^2$. Then, there exists a universal constant $C > 0$ such that the classification error satisfies*

$$\mathbb{P}(y^*(\mathbf{Z}) \neq y) \leq e^{C\sqrt{d/n}} \cdot Q\left(\frac{1}{\sigma}\right) + Q\left(\frac{\sqrt{m}}{\sigma}\right) + e^{-d}.$$

Proof. Recap from Proposition 1. For any L -layer attention model with $L \geq 2$, it can output

$$f_{\text{att-}L}(\mathbf{Z}) = \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X} / n - \sigma^2 \mathbf{I}) \hat{\mu}_s. \quad (32)$$

Let

$$\hat{y} = \text{sgn}(f_{\text{att-}L}(\mathbf{Z}))$$

with $f_{\text{att-}L}(\mathbf{Z})$ defined in (32). Then we have

$$\mathbb{P}(y^*(\mathbf{Z}) \neq y) \leq \mathbb{P}(\hat{y} \neq y).$$

Therefore, in the following, we focus on upper-bounding the classification error $\mathbb{P}(\hat{y} \neq y)$ corresponding to (32). Given that the optimal prediction under the form $\text{sgn}(\mathbf{x}^\top \mathbf{A} \hat{\boldsymbol{\mu}}_s)$ is given by $\hat{y}_{\boldsymbol{\mu}\boldsymbol{\mu}^\top} := \text{sgn}(\mathbf{x}^\top \boldsymbol{\mu} \boldsymbol{\mu}^\top \hat{\boldsymbol{\mu}}_s)$ (cf. Theorem 2), with its corresponding error presented in (10). To analyze the performance of \hat{y} , we study its difference from the prediction $\hat{y}_{\boldsymbol{\mu}\boldsymbol{\mu}^\top}$.

To begin with, let $\mathbf{g}_i = \boldsymbol{\xi}_i / \sigma \sim \mathcal{N}(0, \mathbf{I})$ and $\mathbf{g} = \sum_{i=1}^n \boldsymbol{\xi}_i / \sigma \sqrt{n} \sim \mathcal{N}(0, \mathbf{I})$. For simplicity, let $\mathbf{A} := \mathbf{X}^\top \mathbf{X} / n - \sigma^2 \mathbf{I}$. We get

$$\begin{aligned} \mathbf{A} &= \frac{1}{n} \mathbf{X}^\top \mathbf{X} - \sigma^2 \mathbf{I} \\ &= \frac{1}{n} \left(\sum_{i=1}^n \boldsymbol{\mu} \boldsymbol{\mu}^\top + \boldsymbol{\mu} \boldsymbol{\xi}_i^\top + \boldsymbol{\xi}_i \boldsymbol{\mu}^\top + \boldsymbol{\xi}_i \boldsymbol{\xi}_i^\top \right) - \sigma^2 \mathbf{I} \\ &= \boldsymbol{\mu} \boldsymbol{\mu}^\top + \frac{\sigma}{\sqrt{n}} (\boldsymbol{\mu} \mathbf{g}^\top + \mathbf{g} \boldsymbol{\mu}^\top) + \sigma^2 \left(\frac{\sum_{i=1}^n \mathbf{g}_i \mathbf{g}_i^\top}{n} - \mathbf{I} \right). \end{aligned}$$

From the Laurent-Massart inequality (Laurent & Massart, 2000), we have that with probability at least $1 - e^{-t_1}$ (assuming $t_1 \geq d$),

$$\frac{1}{\sqrt{n}} \|\boldsymbol{\mu} \mathbf{g}^\top + \mathbf{g} \boldsymbol{\mu}^\top\| \leq \frac{2\|\mathbf{g}\|}{\sqrt{n}} \leq 6 \sqrt{\frac{t_1}{n}}. \quad (33)$$

Additionally, from Neopane (2018), we have that with probability at least $1 - e^{-t_2}$ (assuming $t_2 \geq d$)

$$\left\| \frac{\sum_{i=1}^n \mathbf{g}_i \mathbf{g}_i^\top}{n} - \mathbf{I} \right\| \leq C_2 \cdot \sqrt{\frac{t_2}{n}}. \quad (34)$$

Define

$$\boldsymbol{\Delta} := \mathbf{A} - \boldsymbol{\mu} \boldsymbol{\mu}^\top = \frac{\sigma}{\sqrt{n}} (\boldsymbol{\mu} \mathbf{g}^\top + \mathbf{g} \boldsymbol{\mu}^\top) + \sigma^2 \left(\frac{\sum_{i=1}^n \mathbf{g}_i \mathbf{g}_i^\top}{n} - \mathbf{I} \right).$$

Combining (33) and (34), we get with probability at least $1 - e^{-t}$ (for $t \geq d$)

$$\|\boldsymbol{\Delta}\| \leq C \sqrt{\frac{t}{n}}.$$

and therefore, with probability at least $1 - e^{-t}$

$$\boldsymbol{\mu}^\top \mathbf{A} \hat{\boldsymbol{\mu}}_s - \boldsymbol{\mu}^\top \boldsymbol{\mu} \boldsymbol{\mu}^\top \hat{\boldsymbol{\mu}}_s = \boldsymbol{\mu}^\top \boldsymbol{\Delta} \hat{\boldsymbol{\mu}}_s \leq \|\boldsymbol{\Delta}\| \cdot \|\hat{\boldsymbol{\mu}}_s\|.$$

Since $\hat{\boldsymbol{\mu}}_s \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{m} \mathbf{I})$, similar to (33), with probability at least $1 - e^{-t_3}$ (assuming $d \leq t_3 \leq m/\sigma^2$), we can bound

$$\|\hat{\boldsymbol{\mu}}_s\| \leq 1 + \frac{\sigma}{\sqrt{m}} \|\mathbf{g}'\| \leq 1 + 3\sigma \sqrt{\frac{t_3}{m}} \leq 4.$$

Then consider a significantly large n (to ensure that $\|\boldsymbol{\Delta}\| \leq 1/8$). With probability at least $1 - e^{-\min(t, t_3)}$, we can bound

$$\left| \frac{\boldsymbol{\mu}^\top \mathbf{A} \hat{\boldsymbol{\mu}}_s}{\|\mathbf{A} \hat{\boldsymbol{\mu}}_s\|_{\ell_2}} - \frac{\boldsymbol{\mu}^\top \boldsymbol{\mu} \boldsymbol{\mu}^\top \hat{\boldsymbol{\mu}}_s}{\|\boldsymbol{\mu} \boldsymbol{\mu}^\top \hat{\boldsymbol{\mu}}_s\|_{\ell_2}} \right| \leq \frac{\|\boldsymbol{\Delta}\| \cdot \|\hat{\boldsymbol{\mu}}_s\|}{1 - \|\boldsymbol{\Delta}\| \cdot \|\hat{\boldsymbol{\mu}}_s\|} \leq C' \sqrt{\frac{t}{n}}.$$

Recall (28) from the proof of Theorem 2. The error for any given $\boldsymbol{\mu}, \hat{\boldsymbol{\mu}}_s$ is presented by

$$\mathbb{P}(\hat{y} \neq y \mid \boldsymbol{\mu}, \hat{\boldsymbol{\mu}}_s) = \mathcal{Q} \left(\frac{\boldsymbol{\mu}^\top \mathbf{A} \hat{\boldsymbol{\mu}}_s}{\sigma \|\mathbf{A} \hat{\boldsymbol{\mu}}_s\|_{\ell_2}} \right).$$

Note that we have $\mathcal{Q}(x - \delta) \leq e^{x\delta} \mathcal{Q}(x)$ for $0 \leq \delta \ll x$. Then, for any $\boldsymbol{\mu}, \hat{\boldsymbol{\mu}}_s$ satisfying $\boldsymbol{\mu}^\top \hat{\boldsymbol{\mu}}_s > 0$, with probability at least $1 - e^{-\min(t, t_3)}$, we have

$$\begin{aligned} \mathbb{P}(\hat{y} \neq y \mid \boldsymbol{\mu}, \hat{\boldsymbol{\mu}}_s) &= \mathcal{Q} \left(\frac{\boldsymbol{\mu}^\top \mathbf{A} \hat{\boldsymbol{\mu}}_s}{\sigma \|\mathbf{A} \hat{\boldsymbol{\mu}}_s\|_{\ell_2}} \right) \\ &\leq e^{C'' \sqrt{t/n}} \cdot \mathcal{Q} \left(\frac{1}{\sigma} \right). \end{aligned}$$

694 Combining all together, we obtain

$$\begin{aligned}\mathbb{P}(\hat{y} \neq y) &\leq \mathbb{P}(\mu^\top \hat{\mu}_s > 0) \left(e^{C'' \sqrt{t/n}} \cdot \mathcal{Q}\left(\frac{1}{\sigma}\right) + e^{-t} \right) + \mathbb{P}(\mu^\top \hat{\mu}_s < 0) \\ &\leq \left(e^{C'' \sqrt{t/n}} \cdot \mathcal{Q}\left(\frac{1}{\sigma}\right) + e^{-t} \right) + \mathcal{Q}\left(\frac{\sqrt{m}}{\sigma}\right).\end{aligned}$$

695 Choosing $t = d$ completes the proof. ■

696 C Additional Details on Tabular Experiments

Algorithm 1 LoopTabFM: Looping Tabular FM with Soft Pseudo-labels and Risk-aware Updates

Require: Dataset $\mathcal{D}_{\text{lab}}, \mathcal{D}_{\text{unlab}}$, looping iterations K

```

1: procedure LOOPING( $\mathcal{D}_{\text{lab}}, \mathcal{D}_{\text{unlab}}, K$ )
2:   Base model/ $\text{FM}_0 \leftarrow \text{TabPFN-v2}(\mathcal{D}_{\text{lab}})$ 
3:    $\mathcal{D}_{\text{unlab}} \leftarrow \text{FM}_0(\mathcal{D}_{\text{unlab}})$  ▷ Assign pseudo labels via  $\hat{y}^{\text{soft}} \leftarrow \text{FM}_0(x \in \mathcal{D}_{\text{unlab}})$ .
4:    $\text{FM}_{\text{best}} \leftarrow \text{FM}_0$ 
5:    $\mathcal{R}_{\text{val}} = \text{Val\_Risk}(\mathcal{D}_{\text{unlab}})$ 
6:   for Looping iteration  $k = 1, \dots, K$  do
7:      $\text{FM}_k \leftarrow \text{TabPFN-v2}(\mathcal{D}_{\text{lab}}, \mathcal{D}_{\text{unlab}})$ 
8:      $\mathcal{D}_{\text{unlab}} \leftarrow \text{FM}_k(\mathcal{D}_{\text{unlab}})$  ▷ Update pseudo labels via  $\hat{y}^{\text{soft}} \leftarrow \text{FM}_k(x \in \mathcal{D}_{\text{unlab}})$ .
9:     if  $\text{Val\_Risk}(\mathcal{D}_{\text{unlab}}) < \mathcal{R}_{\text{val}}$  then
10:       $\text{FM}_{\text{best}} \leftarrow \text{FM}_k$ 
11:       $\mathcal{R}_{\text{val}} = \text{Val\_Risk}(\mathcal{D}_{\text{unlab}})$ 
12:     end if
13:   end for
14:   return  $\text{FM}_{\text{best}}$ 
15: end procedure
16: procedure VAL_RISK( $\mathcal{D}_{\text{unlab}}$ )
17:   return  $\frac{1}{|\mathcal{D}_{\text{unlab}}|} \sum_i \min(|\hat{y}_i^{\text{soft}} - 1|, |\hat{y}_i^{\text{soft}} + 1|)$ 
18:   ▷  $\hat{y}^{\text{soft}}$  corresponds to the assigned soft label for feature in  $\mathcal{D}_{\text{unlab}}$ .
19: end procedure

```

697 In this section, we provide additional details regarding the tabular experiments discussed in Section 5.2.
698 We propose the LoopTabFM algorithm with its details outlined in Algorithm 1. Suppose that we
699 are given labeled \mathcal{D}_{lab} and unlabeled $\mathcal{D}_{\text{unlab}}$ datasets during training. The overall workflow of the
700 algorithm proceeds as follows:

- 701 1. **Base Model:** Train TabPFN on the labeled dataset \mathcal{D}_{lab} and treat the resulting model as the
702 base model (Loop-0). Its test accuracy is reported in Table 1.
- 703 2. **Pseudo-Label Assignment:** Using the current model (e.g., Loop- k), generate predictions
704 for the unlabeled data $\mathcal{D}_{\text{unlab}}$. Assign soft pseudo-labels based on these predictions. Note
705 that the model outputs are scalars (i.e., elements of \mathbb{R}) and can be interpreted as soft labels.
- 706 3. **Model Update:** Construct a new prompt that includes both labeled examples with their true
707 labels and unlabeled examples with their assigned soft pseudo-labels. Fit the TabPFN to this
708 combined prompt to obtain the updated model (Loop- $(k + 1)$). Repeat from Step 2 until the
709 maximum number of looping iterations is reached.
- 710 ★ **Model Validation:** To improve the stability of the looping process, we introduce an
711 additional validation step and retain the model with the lowest validation risk as the final
712 (best) model. Specifically, suppose that the unlabeled data has been assigned soft pseudo-
713 labels, i.e., $\mathcal{D}_{\text{unlab}} = \{(x_i, \hat{y}_i^{\text{soft}})_{i=1}^n\}$. The validation risk is then computed over the pseudo
714 labels as:

$$\text{Val_Risk}(\mathcal{D}_{\text{unlab}}) = \frac{1}{n} \sum_{i \in [n]} \min(|\hat{y}_i^{\text{soft}} - 1|, |\hat{y}_i^{\text{soft}} + 1|),$$

715 which penalizes predictions that deviate from confident binary labels ± 1 .

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: The abstract and introduction accurately reflect the paper's contributions and scope.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: Our assumptions underlying the algorithm and their necessity are discussed. Also, the limitation is discussed in the supplementary material.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: Our assumptions underlying the analysis and algorithm are discussed. Detailed proofs can be found in the supplementary material.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: All the information needed to reproduce the main experimental results of the paper are provided, either in the main paper or in the supplementary material.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We have attached the code for implementing the algorithm and reproducing the experiments in the supplementary material.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: All the information needed to reproduce the main experimental results of the paper are provided, either in the main paper or in the supplementary material.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Detailed experiment results with errors is included in the supplementary.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Details can be found in the supplementary material.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: We have reviewed and confirmed that the research conducted in the paper conform with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: There is no societal impact of the work performed.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We cited the original paper that produced the code package or dataset.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

977 • For existing datasets that are re-packaged, both the original license and the license of
978 the derived asset (if it has changed) should be provided.
979 • If this information is not available online, the authors are encouraged to reach out to
980 the asset’s creators.

981 **13. New assets**

982 Question: Are new assets introduced in the paper well documented and is the documentation
983 provided alongside the assets?

984 Answer: [NA]

985 Justification: The paper does not release new assets.

986 Guidelines:

987 • The answer NA means that the paper does not release new assets.
988 • Researchers should communicate the details of the dataset/code/model as part of their
989 submissions via structured templates. This includes details about training, license,
990 limitations, etc.
991 • The paper should discuss whether and how consent was obtained from people whose
992 asset is used.
993 • At submission time, remember to anonymize your assets (if applicable). You can either
994 create an anonymized URL or include an anonymized zip file.

995 **14. Crowdsourcing and research with human subjects**

996 Question: For crowdsourcing experiments and research with human subjects, does the paper
997 include the full text of instructions given to participants and screenshots, if applicable, as
998 well as details about compensation (if any)?

999 Answer: [NA]

1000 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1001 Guidelines:

1002 • The answer NA means that the paper does not involve crowdsourcing nor research with
1003 human subjects.
1004 • Including this information in the supplemental material is fine, but if the main contribu-
1005 tion of the paper involves human subjects, then as much detail as possible should be
1006 included in the main paper.
1007 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1008 or other labor should be paid at least the minimum wage in the country of the data
1009 collector.

1010 **15. Institutional review board (IRB) approvals or equivalent for research with human**
1011 **subjects**

1012 Question: Does the paper describe potential risks incurred by study participants, whether
1013 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1014 approvals (or an equivalent approval/review based on the requirements of your country or
1015 institution) were obtained?

1016 Answer: [NA]

1017 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1018 Guidelines:

1019 • The answer NA means that the paper does not involve crowdsourcing nor research with
1020 human subjects.
1021 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1022 may be required for any human subjects research. If you obtained IRB approval, you
1023 should clearly state this in the paper.
1024 • We recognize that the procedures for this may vary significantly between institutions
1025 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1026 guidelines for their institution.
1027 • For initial submissions, do not include any information that would break anonymity (if
1028 applicable), such as the institution conducting the review.

1029 **16. Declaration of LLM usage**
1030 Question: Does the paper describe the usage of LLMs if it is an important, original, or
1031 non-standard component of the core methods in this research? Note that if the LLM is used
1032 only for writing, editing, or formatting purposes and does not impact the core methodology,
1033 scientific rigorousness, or originality of the research, declaration is not required.
1034 Answer: [NA]
1035 Justification: This research does not involve LLMs as any important components.
1036 Guidelines:
1037 • The answer NA means that the core method development in this research does not
1038 involve LLMs as any important, original, or non-standard components.
1039 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
1040 for what should or should not be described.
1041