# How Private is Your Attention?
# Bridging Privacy with In-Context Learning

**Soham Bonnerjee**
sohambonnerjee@uchicago.edu

**Yeon Zhen Wei (Kingsley)**
yeon@uchicago.edu

**Anna Asch**
aasch@uchicago.edu

**Sagnik Nandy**
sagnik@uchicago.edu

**Promit Ghosal**
promit@uchicago.edu

## Abstract

In-context learning (ICL)—the ability of transformer-based models to perform new tasks from examples provided at inference time—has emerged as a hallmark of modern language models. While recent works have investigated the mechanisms underlying ICL, its feasibility under formal privacy constraints remains largely unexplored. In this paper, we propose a differentially private pretraining algorithm for linear attention heads and present the first theoretical analysis of the privacy–accuracy trade-off for ICL in linear regression. Our results characterize the fundamental tension between optimization and privacy-induced noise, formally capturing behaviors observed in private training via iterative methods. Additionally, we show that our method is robust to adversarial perturbations of training prompts, unlike standard ridge regression. All theoretical findings are supported by extensive simulations across diverse settings.

## 1 Introduction

Attention-based models, particularly large language models (LLMs), have demonstrated remarkable capabilities in performing *in-context learning* [Brown et al., 2020, Lieber et al., 2021, Rae et al., 2021, Black et al., 2022, Bubeck et al., 2023]. This paradigm has transformed human-AI interaction, enabling AI models to tackle complex tasks without explicit parameter updates. A growing body of theoretical work has aimed to explain this emergent behavior [Dong et al., 2022, Akyürek et al., 2022, Garg et al., 2022, Wang et al., 2023, Xie et al., 2022], often using simplified settings. These studies suggest that transformers can implicitly infer patterns or rules from training examples in the prompt and apply them to new, related inputs during inference.

The growing use of LLM-based agents in sensitive domains such as medicine [Li et al., 2025, Dennstädt et al., 2025] and psychology [Ke et al., 2024] underscores the urgent need for robust privacy safeguards. In particular, model providers must prevent adversaries from extracting sensitive training data, a risk highlighted by recent work demonstrating that LLMs can memorize and reveal specific examples when prompted adversarially [Carlini et al., 2021, 2022, Tirumala et al., 2022]. A principled approach to mitigating such leakage is *differential privacy* (DP) [Dwork et al., 2006], which ensures that an algorithm's output remains nearly unchanged when a single training point is modified. This is typically achieved by injecting calibrated noise to limit individual influence.

However, integrating privacy-preserving mechanisms into the pretraining process of a transformer inevitably degrades the downstream performance of in-context learning on test prompts. This trade-off motivates a rigorous study of the cost of privacy of *in-context differentially-private* algorithms: what additional error is incurred at test time?

## 1.1 Main Results

We study the effect of differentially-private pretraining on in-context learning (ICL) for linear regression, where each data point is a noisy linear response to input features. We propose a differentially-private pretraining algorithm for a linear attention head that performs ICL—predicting the response for a query input by attending to a sequence of labeled input-output examples. The model is trained on $N$ prompts, each containing $L$ feature-response pairs sampled from a noisy linear model, and optimized to minimize squared prediction error on the query token. To enforce privacy, we apply the Gaussian mechanism—gradient clipping followed by additive noise—commonly used in private empirical risk minimization [Dwork et al., 2006, Chaudhuri et al., 2011, Abadi et al., 2016, Cai et al., 2021]. Our method, `NoisyHead` (Algorithm 1, Section 3), formalizes this approach.

We define the *cost of privacy* as the difference, between attention heads trained with and without privacy constraints, in average prediction error of the response to a query token from a held-out test prompt. Our main theoretical result characterizes how the *cost of privacy* scales with the number of training prompts $N$, the prompt length $L$, the token dimension $D$, and the privacy parameters $(\varepsilon, \delta)$. We state it informally below:

**Theorem 1.1** (Informal). *In the low dimensional regime, when $L$ and $\sqrt{N}$ are asymptotically of same order and $D = O(1)$, the cost of privacy satisfies*

$$Cost \ of \ Privacy \lesssim \frac{1}{N^{3/2}L^2} \frac{\log(1/\delta)}{\varepsilon^2}.$$

*In the high dimensional regime, when $N/D^2 = O(1)$ and $L/D = O(1)$, the cost of privacy scales as*

$$Cost \ of \ Privacy \lesssim \frac{D^2}{N^2 L^2} \frac{\log(1/\delta)}{\varepsilon^2},$$

*up to* polylog *factors.*

A formal version of this result is presented in Theorem 4.2, followed by a detailed discussion of its implications. The theorem highlights that the cost of privacy exhibits fundamentally different behavior in the low- and high-dimensional regimes. In the *low-dimensional* setting, the minimax cost of privacy for learning a linear model from $L$ labeled data points is known to scale as $(\varepsilon L)^{-2} \cdot \log(1/\delta)$, as established in Cai et al. [2021]. The result above shows that leveraging contextual data reduces this cost to $N^{-3/2}(\varepsilon L)^{-2} \cdot \log(1/\delta)$. However, because test-time prediction requires learning an unseen coefficient vector $w$, we do not achieve the rate $N^{-2}(\varepsilon L)^{-2} \cdot \log(1/\delta)$, which would be expected if the coefficient was identical across all training and test prompts. In contrast, in the *high-dimensional* regime, where the feature dimension scales with the number of prompts $N$, we incur an additional multiplicative factor of $\sqrt{N}$ in the denominator due to the increased complexity of the learning problem.

We also show that our private pretraining procedure is more robust to adversarial perturbations of training prompts than its non-private counterpart. When a fraction of prompts are corrupted, the prediction risk on test instances remains significantly more stable under our method — a property especially relevant given recent concerns about adversarial attacks in LLMs [Anwar et al., 2024].

Our key contributions are as follows:

(**1**) We propose a differentially-private pretraining algorithm (`NoisyHead`) based on the *Gaussian mechanism* for training linear attention heads to perform in-context learning in linear regression (see Algorithm 1). Our method is motivated by the differentially-private stochastic gradient descent algorithm [Abadi et al., 2016], containing a tuned noise-injection at the gradient steps.

(**2**) We provide a detailed theoretical analysis of the excess risk incurred by enforcing differential privacy during pretraining in Theorem 4.1. In particular, it characterizes the privacy–utility trade-off, quantifying the impact of privacy constraints on the prediction error of `NoisyHead` across any number of iterations $T$ of the algorithm. This trade-off exhibits dichotomous behavior depending on how the feature dimension $D$ scales with the number of training samples $N$. We identify two distinct regimes: one where $D = O(\log N)$ and another where $N/D^2 = O(1)$. These lead to qualitatively different error decay rates with respect to $N$, $L$, and $D$, as formalized in Theorem 4.2. In the over-parametrized setting when $N, L^2, D^2$ are asymptotically of the same order, we show that there is a delicate interplay between the number of training iterations and the generalization error on unseen prompts. Due to the

injection of noise at each iteration, longer training can degrade generalization, necessitating careful selection of the number of optimization steps. This highlights the importance of "early stopping" for the algorithm. See Proposition 4.1 and the following remark for related discussion.

($\mathbf{4}$) We establish that `NoisyHead` exhibits a notable robustness property under adversarial perturbations to the training data, particularly during the pretraining stage. Compared to the baseline method proposed in Lu et al. [2024], our approach shows significantly less degradation in generalization error in the presence of such perturbations. In the baseline setting, where the linear attention module is pretrained using ridge regression, even moderately large perturbations can induce a distributional shift in the training data, leading to inaccurate estimation of model weights and consequently poor generalization. In contrast, `NoisyHead` incorporates a truncation mechanism that clips responses, predictors, and weights within prescribed compact sets. This simple yet effective step restricts the influence of corrupted or outlying data points, enhancing robustness to adversarial noise introduced during training. Theoretical support for this robustness is provided in Theorem 5.1.

($\mathbf{5}$) We conduct a comprehensive empirical study to validate the theoretical predictions of our analysis. In both low- and high-dimensional regimes (Section 6.1), we demonstrate that the excess prediction risk of `NoisyHead` decays with increasing sample size and privacy parameter, consistent with the rates derived in Theorem 4.2. Moreover, in the overparameterized regime (Section 6.2), our experiments reveal a distinct phase transition in the generalization error: initially decreasing due to optimization, but eventually increasing due to cumulative noise from differential privacy. This phenomenon, visualized in Figure 2, substantiates the theoretical trade-off outlined in Proposition 4.1 and underscores the critical role of early stopping. Finally, robustness experiments (Section 6.3) confirm that `NoisyHead` maintains stable performance under adversarial perturbations, while ridge-based pretraining degrades significantly. These results highlight the practical utility of our method and affirm the relevance of our theoretical contributions in realistic settings.

## 1.2  Related literature and notations

Since its introduction by Dwork et al. [2006], differential privacy has become a cornerstone of privacy-preserving machine learning, inspiring a wide range of algorithms across classical and deep learning tasks [Cai et al., 2021, Wang and Xu, 2019, Gu et al., 2024, Jain and Thakurta, 2013, Ni et al., 2016, Ji et al., 2019, Abadi et al., 2016, Feldman et al., 2018]. In parallel, recent work has explored the in-context learning (ICL) capabilities of transformers, demonstrating that pretraining enables them to emulate diverse algorithms—including ridge regression, generalized linear models, Lasso, and neural networks—purely from contextual examples [Dai et al., 2023], with theoretical insights provided for linear attention models by Zhang et al. [2024] and Lu et al. [2024]. Despite significant advances in both areas, their intersection remains underexplored: while prior work has investigated differentially-private pretraining for transformers [Majmudar et al., 2022, Yu et al., 2023, Li et al., 2022] and evaluated the privacy properties of language models [Hoory et al., 2021, Anil et al., 2021], the impact of privacy on downstream ICL performance has not been theoretically analyzed. This paper bridges this gap by providing the first rigorous analysis of how imposing differential privacy during pretraining influences the in-context learning capabilities of attention-based models.

### 1.2.1  Notation

In this paper, we denote the set $\{1, \ldots, n\}$ by $[n]$. $d$-dimensional Euclidean space is $\mathbb{R}^d$, with $\mathbb{R}^d_{>0}$ the positive orthant. The set of $m \times n$ real matrices is $\mathbb{R}^{m \times n}$, and $\mathbb{S}^{d-1}$ denotes the $d$-dimensional unit sphere. The Frobenius norm of a matrix $A$ is $\|A\|_F$, and $\langle \cdot, \cdot \rangle$ denotes the standard inner product. We write $a_n \lesssim b_n$ if $a_n \leq C b_n$ for some constant $C > 0$, and $a_n \asymp b_n$ if $C_1 b_n \leq a_n \leq C_2 b_n$ for some constants $C_1, C_2 > 0$. We also write $a_n \asymp b_n$ as $a_n = \Theta(b_n)$.

## 2  Problem Formulation

We consider a set-up where we observe a sequence of labeled tokens $\{(y_i, x_i) : i \in \{1, \ldots, L\}\}$, for $x_i \overset{i.i.d}{\sim} \mathcal{U}(\mathbb{S}^{D-1})$ and $y_i = w^\top x_i + \epsilon_i$, with $w \sim \mathcal{N}_D(0, \mathbb{I}_D)$ and $\epsilon_i \overset{i.i.d}{\sim} \mathcal{N}(0, \tau^2)$. Here $\mathcal{U}(\mathbb{S}^{D-1})$ denotes the uniform distribution on the $D$-dimensional hypersphere and $\mathcal{N}_k(\mu, \Sigma)$ denotes the $k$ dimensional normal distribution with mean $\mu$ and covariance $\Sigma$. For a test token $(y_{L+1}, x_{L+1})$

generated independently from the same distribution as the training tokens, we want to predict $y_{L+1}$ based on $x_{L+1}$.

This setting was used by Zhang et al. [2024] and Lu et al. [2024], both of whom considered the noiseless case of $\tau^2 = 0$. As proposed therein, we embed the prompt as

$$E = \begin{pmatrix} x_1 & x_2 & \cdots & x_L & x_{L+1} \\ y_1 & y_2 & \cdots & y_L & 0 \end{pmatrix} \in \mathbb{R}^{(D+1) \times (L+1)}. \tag{2.1}$$

This matrix is passed through a single linear attention head as follows:

$$f(E; \theta) = E + W^{PV} E \cdot \frac{E^\top W^{KQ} E}{L}, \tag{2.2}$$

where $\theta = (W_{PV}, W_{KQ})$ with $W_{PV} \in \mathbb{R}^{(D+1) \times (D+1)}$ and $W_{KQ} \in \mathbb{R}^{(D+1) \times (D+1)}$. The prediction of the query response is given by the $(D+1, L+1)$-th entry of $f(E; \theta)$; that is, $\widehat{y}_{L+1}(E) = (f(E; \theta))_{(D+1, L+1)}$. We aim to learn the parameters of the model $f(E; \theta)$ by pretraining the model based on $N$ training prompts $\{(y_{k,1}, x_{k,1}), \ldots, (y_{k,L}, x_{k,L}), (y_{k,L+1}, x_{k,L+1})\}_{k=1}^N$, where the $L+1$-th token is the query token. Putting the prompts into matrices $E_1, \ldots, E_N$, we have

$$E_k := \begin{pmatrix} x_{k,1} & x_{k,2} & \cdots & x_{k,L} & x_{k,L+1} \\ y_{k,1} & y_{k,2} & \cdots & y_{k,L} & 0 \end{pmatrix} \in \mathbb{R}^{(D+1) \times (L+1)}.$$

Now we minimize the standard loss function $\mathcal{L}(\theta) = \frac{1}{2N} \sum_{i=1}^N (\widehat{y}_{L+1}(E_k) - y_{k,L+1})^2$. The predictor $(f(E; \theta))_{(D+1, L+1)}$ can be simplified by linear algebra to

$$\hat{y}_{L+1} := [(f(E; \theta)_{[D+1, L+1]}] = \begin{bmatrix} (w_{21}^{PV})^\top & w_{22}^{PV} \end{bmatrix} \left( \frac{EE^\top}{L} \right) \begin{bmatrix} W_{11}^{KQ} \\ (w_{21}^{KQ})^\top \end{bmatrix} x_{L+1}, \tag{2.3}$$

where we have used the matrices $W^{PV}$ and $W^{KQ}$, partitioned as follows:

$$W^{PV} = \begin{bmatrix} W_{11}^{PV} & w_{12}^{PV} \\ (w_{12}^{PV})^\top & w_{22}^{PV} \end{bmatrix}, \quad W^{KQ} = \begin{bmatrix} W_{11}^{KQ} & w_{12}^{KQ} \\ (w_{12}^{KQ})^\top & w_{22}^{KQ} \end{bmatrix},$$

with $W_{11}^{PV}, W_{11}^{KQ} \in \mathbb{R}^{D \times D}$, $w_{21}^{PV}, w_{21}^{KQ} \in \mathbb{R}^D$, and $w_{22}^{PV}, w_{22}^{KQ} \in \mathbb{R}$. The quadratic form (2.3) can be expanded to yield

$$\hat{y}_{L+1} = \frac{1}{L} \langle x_{L+1}, Q_W^{(1)} + Q_W^{(2)} \rangle, \tag{2.4}$$

where $Q_W^{(1)} := w_{22}^{PV} W_{11}^{KQ} \sum_{i=1}^L y_i x_i + w_{22}^{PV} w_{12}^{KQ} \sum_{i=1}^L y_i^2$ and $Q_W^{(2)} := W_{11}^{KQ} \sum_{i=1}^{\ell+1} x_i x_i^\top w_{12}^{PV} + w_{12}^{KQ} \sum_{i=1}^\ell y_i x_i^\top w_{12}^{PV}$. Following Yu et al. [2023] and Zhang et al. [2024], we adopt the assumption that $w_{12}^{KQ} = 0$ and $w_{12}^{PV} = 0$ throughout this paper. This particular choice is also explained in Section A.1. Let us define

$$\Gamma = w_{22}^{PV} W_{11}^{KQ} \in \mathbb{R}^{D \times D}, \quad \text{and} \quad Z = \frac{1}{L} x_{L+1} \sum_{i=1}^L y_i x_i^\top \in \mathbb{R}^{D \times D}. \tag{2.5}$$

With this definition of $\Gamma$ and $Z$, the predictor $\widehat{y}$ simplifies to the inner product $\widehat{y} = \langle \Gamma, Z \rangle$, and we train the model using the following regularized squared error loss:

$$\mathcal{L}_\lambda(\Gamma) := \frac{1}{N} \sum_{i=1}^N (y_i - \langle \Gamma, Z_i \rangle)^2 + \lambda \|\Gamma\|_F^2. \tag{2.6}$$

The solution to this optimization problem is denoted by $\Gamma^\star \in \mathbb{R}^{D \times D}$, whose vectorized form is given by

$$\text{vec}(\Gamma^\star; E_1, \ldots, E_N) = \left( \lambda N I + \sum_{k=1}^N \text{vec}(Z_k) \text{vec}(Z_k)^\top \right)^{-1} \sum_{k=1}^N y_{k,L+1} \text{vec}(Z_k). \tag{2.7}$$

4

**Algorithm 1** In-Context Differentially private pretraining of linear attention head (`NoisyHead`)

---

**Input:** Training prompts $(E_k)_{k \in [N]} \in \mathbb{R}^{(D+1) \times (L+1)}$; noise scale $\sigma$; privacy parameters $\varepsilon, \delta$; clipping parameter $\mathcal{C} \geq 0$; projection parameters $R, G \geq 0$; regularization parameter $\lambda := \lambda(n, d) \geq c > 0$; number of iterations $T$; step-size $\eta_0$; and initialization $\Gamma^0 \in \mathbb{R}^{D \times D}$ with $\|\Gamma^0\|_F \leq R$.

- For $k \in [N]$, $\widetilde{Z}_k := \Pi_G \left( L^{-1} x_{k,L+1} \sum_{i=1}^L \mathtt{clip}_{\mathcal{C}}(y_{k,i}) x_{k,i}^\top \right)$.

- For $t$ in $0, 1, \ldots, T-1$:

    - Generate $z_t \in \mathbb{R}^{D \times D}$ such that $\mathrm{vec}(z_t) \sim \mathcal{N}_{D^2} \left( 0, 2\eta_0^2 \frac{T^2 \sigma^2}{\varepsilon^2 N^2} \log \frac{1.25T}{\delta} \mathbb{I}_{D^2} \right)$.

    - Do $\Gamma^{t+1} = \Pi_R \left( (1 - 2\lambda\eta_0)\Gamma^t - \eta_0 N^{-1} \sum_{k=1}^N \left( \langle \Gamma^t, \widetilde{Z}_k \rangle - \mathtt{clip}(y_{k,L+1}) \right) \widetilde{Z}_k + z_t \right)$.

    **Output:** $\hat{\Gamma} := \Gamma^T$.

---

## 3 Differentially Private Pretraining

In this section, we present our differentially-private pretraining program of a linear attention network. Before proceeding to the main algorithm, we recall the definition of differential privacy.

**Definition 3.1.** *A randomized algorithm $\mathcal{M}(\cdot)$ over a set of prompts is said to be in-context $(\varepsilon, \delta)$-differentially private if for any two sequences of prompts $\mathcal{D} = (E_1, \ldots, E_N)$ and $\mathcal{D}' = (E_1', \ldots, E_N')$ differing in at most one entry, and for all measurable subsets $\mathcal{W}$ of outputs,*

$$\mathbb{P}[\mathcal{M}(\mathcal{D}) \in \mathcal{W}] \leq e^\varepsilon \mathbb{P}[\mathcal{M}(\mathcal{D}') \in \mathcal{W}] + \delta.$$

The probability is taken over the internal randomness of the mechanism $\mathcal{M}$, while the prompt sequences $\mathcal{D}$ and $\mathcal{D}'$ are treated as fixed. This definition ensures that the inclusion or exclusion of any individual data point has a limited effect on the algorithm's output, thereby preserving privacy. A standard approach to enforce DP in the iterative training of machine learning models (e.g. gradient descent) is to inject noise at each update step. The cumulative effect of this noise is carefully calibrated to satisfy user-specified $(\varepsilon, \delta)$-differential privacy guarantees but minimize degradation in model performance. This technique, introduced as *differentially private stochastic gradient descent*, has been echoed in recent works [Abadi et al., 2016, Cai et al., 2021, Zhang et al., 2021, Gopi et al., 2021, Majmudar et al., 2022, Bombari and Mondelli, 2025]. In what follows, we improvise the aforementioned differentially-private training strategy while using the gradient descent to minimize the regularized loss $\mathcal{L}_\lambda(\Gamma)$ over a sequence of prompts:

$$\Gamma^{t+1} = (1 - 2\lambda\eta_0)\Gamma^t - \frac{\eta_0}{N} \sum_{k=1}^N \left( \langle \Gamma^t, Z_k \rangle - y_{k,L+1} \right) Z_k,$$

where $\eta_0$ is the learning rate, and $\lambda$ is the regularization parameter.

To ensure privacy, we inject carefully calibrated Gaussian noise into each update step. The variance of this noise is set proportional to the $\ell_2$-*sensitivity* of the update, which measures the maximum change in the update (in Frobenius norm) resulting from the change of a single training example. Formally, the $\ell_2$-sensitivity at iteration $t$ is defined as:

$$\Delta(\hat{\Gamma}) = \left\| \hat{\Gamma}(E_1, \ldots, E_N) - \hat{\Gamma}(E_1', \ldots, E_N') \right\|_F, \tag{3.1}$$

where the datasets $(E_1, \ldots, E_N)$ and $(E_1', \ldots, E_N')$ differ in exactly one training prompt. Intuitively, privacy is preserved because an adversary observing the output of the algorithm (i.e., the final parameters) cannot reliably distinguish whether a change in the result is due to the presence or absence of a particular training prompt or due to the added random noise. However, in the problem setup considered in this paper, the $\ell_2$-sensitivity of the gradient updates may not be uniformly bounded across all possible sequences of training prompts due to the unbounded nature of the weights $w$ and the noise $\epsilon$. To mitigate this, we clip the responses $y_k$ and project the gradient updates $\Gamma^t$ onto

compact sets. With these modifications, our differentially-private pretraining algorithm is presented in Algorithm 1, where the clipping and projection operators are defined as follows:

$$\texttt{clip}_{\mathcal{C}}(x) := \arg\min_{y \in [-\mathcal{C}, \mathcal{C}]} \|x - y\|_2, \quad \Pi_R(X) := \arg\min_{\substack{Y \in \mathbb{R}^{(D+1) \times (D+1)} \\ \|Y\|_F \leq R}} \|X - Y\|_F.$$

**Theorem 3.2.** *Given the set of hyperparameters* $(\mathcal{C}, R, G) \in \mathbb{R}^3_{>0}$, *Algorithm 1 is* $(\varepsilon, \delta)$-*differentially private if the noise scale* $\sigma \geq 2G(\mathcal{C} + RG)$.

Theorem 3.2 (which we prove in Section A.2) hints at the minimum amount of noise to be injected in the gradient descent step to achieve differential privacy. In particular, the amount of noise depends crucially on the projection parameters $\mathcal{C}$ and $R$. On the other hand, the higher the noise variance $\sigma^2$, the more we expect the predictive performance of the differentially-private estimate $\hat{\Gamma}$ to degrade compared to the ridge estimate $\Gamma^\star$ (as defined in (2.6)). However, it can still be argued that performing an appropriate number of iterations, governed by an "early stopping criterion", can improve accuracy. In fact, the additional error from noise injection can be made much smaller than the overall gradient descent error by properly tuning the hyper-parameters. This angle is explored in detail in Section 4.

## 4   Cost of In-Context Differential Privacy

In this section, we rigorously characterize the additional error incurred due to enforcing privacy constraints in Algorithm 1. Let $E^{\texttt{test}}$ be a test prompt and $y^{\texttt{test}}_{L+1}$ be the corresponding query response. Let us consider the prediction error in the test prompt given by

$$\mathcal{L}_{\texttt{test}}(\Gamma) = (y^{\texttt{test}}_{L+1} - \langle \Gamma, Z(E^{\texttt{test}}) \rangle)^2,$$

where $Z(E^{\texttt{test}})$ is constructed from $E^{\texttt{test}}$ as described in (2.5). We bound $\mathcal{L}_{\texttt{test}}(\Gamma)$ by the following two types of error terms:

$$\mathcal{L}_{\texttt{test}}(\hat{\Gamma}) \leq 2\mathcal{L}_{\texttt{test}}(\Gamma^\star) + 2(\langle \hat{\Gamma}, Z(E^{\texttt{test}}) \rangle - \langle \Gamma^\star, Z(E^{\texttt{test}}) \rangle)^2.$$

While $\mathcal{L}_{\texttt{test}}(\Gamma^\star)$ is the prediction error of the non-private procedure, the extra error is proportional to $(\langle \hat{\Gamma}, Z(E^{\texttt{test}}) \rangle - \langle \Gamma^\star, Z(E^{\texttt{test}}) \rangle)^2$. The following theorem characterizes this extra error.

**Theorem 4.1.** *Consider the pretrained weights* $\hat{\Gamma}$ *generated by running* `NoisyHead` *(Algorithm 1) on prompt set* $(E_1, \ldots, E_N)$, *ensuring* $(\varepsilon, \delta)$ *differential privacy for* $T$ *iterations with a fixed stepsize* $\eta_0 \in \left(\frac{\lambda}{c(2\lambda + G^2)^2}, \frac{\lambda}{(2\lambda + G^2)^2}\right)$ *for some large* $c > 1$, *and* $\Gamma^\star$ *generated by solving the ridge regression described in* (2.6). *If the clipping and projection parameters are set as:*

$$\nu = 1 + \tau^2, \quad \mathcal{C} = \sqrt{2\nu \log(NL/\kappa)}, \quad G = \frac{\mathcal{C}}{\sqrt{L}}\left(1 + \frac{(\log(N/\kappa))^{1/2}}{D}\right),$$

$$G_0 = \frac{\mathcal{C}}{\sqrt{L}}\left(1 + \frac{(\log(1/\kappa))^{1/2}}{D}\right), \quad \text{and } R \asymp \lambda^{-1}\mathcal{C}^2\sqrt{\frac{N}{L}}\left(1 + \frac{(\log(1/\kappa))^{1/2}}{D}\right), \quad (4.1)$$

*then for a test prompt* $E$ *independent of* $(E_k)_{k \in [N]}$,

$$(\langle \hat{\Gamma}, Z \rangle - \langle \Gamma^\star, Z \rangle)^2 \leq G_0^2\left((1 - \eta_0\lambda)^T R^2 + \sigma^2\eta_0^2 D \frac{T^2 \log(2T/\delta)}{N^2\varepsilon^2}\right). \quad (4.2)$$

*with probability greater than* $1 - c_1 \exp(-c_2 D) - 4\kappa$, *where* $Z$ *is formed via* $E$ *as in* (2.5).

The above theorem is proved in Section A.3.

*Remark* 4.1. The "*cost of privacy*" on the right hand side of (4.2) naturally decomposes into two components. The first arises from the optimization error of gradient descent, hereby referred to as the "*cost of descent*", and is given by $(1 - \eta_0\lambda)^T$, where $\eta_0$ is the step size, $\lambda$ is the strong convexity parameter, and $T$ is the number of gradient steps. The second component stems from the noise injected at each iteration to ensure DP, and takes the form $\sigma^2\eta_0^2 D \cdot \frac{T^2 \log(2T/\delta)}{N^2\varepsilon^2}$, where $\sigma^2$ denotes the variance of the added noise, $D$ is the feature dimension, $N$ is the number of training samples, and $(\varepsilon, \delta)$ are the privacy parameters. This term will henceforth be referred to as the "*cost of noise injection*". The trade-off between these two terms plays a crucial role in determining the

generalization error. While the optimization error decays exponentially with $T$, the privacy-induced error increases quadratically. Therefore, it is essential to choose an optimal stopping time for the gradient descent iterations. This optimal stopping time depends on the problem hyper-parameters $\eta_0$, $\lambda$, and the feature dimension $D$. In the following theorem (proved in Section A.4), we characterize how the interplay between the dimensionality and the stopping time governs the behavior of the generalization error in different settings of interest.

**Theorem 4.2.** *Assume $N \gtrsim D^2 L^{-2}$, and suppose the noise scale $\sigma$ in Algorithm 1 satisfies $\sigma \asymp 2G(\mathcal{C} + RG)$ where $(\mathcal{C}, R, G)$ is the set of hyper-parameters. Then, under the assumptions and hyper-parameter specifications of Theorem 4.1, the following assertions hold:*

(i) *(Low-dimensional setting) If $\kappa \gtrsim \exp(-D^2)$, and $D^2 \lesssim \log(NL)$, then, after $T = \frac{\log(N^2 L D^3)}{\log((1-\eta_0\lambda)^{-1})}$ many iterations of Algorithm 1 with $\eta_0 \asymp \lambda \asymp 1$ such that $\eta_0\lambda \in (0, 1)$, the cost of privacy of $\hat{\Gamma}(= \Gamma^T)$ behaves as follows:*

$$(\langle \hat{\Gamma}, Z \rangle - \langle \Gamma^\star, Z \rangle)^2 \lesssim \nu^5 \frac{\log^{10}(NL)}{NL^3} \left( 1 + \frac{\log(1/\delta)}{\varepsilon^2} \right), \tag{4.3}$$

*with probability at least $1 - c_1 \exp(-c_2 D) - 4\kappa$.*

(ii) *(High-dimensional setting) If $\kappa \gtrsim (NL)^{-1}$, and $D^2 \gtrsim \log(NL)$, then for some $r$ (possibly depending on $N$, $L$ and $D$), let $T = \frac{\log r}{\log((1-\eta_0\lambda)^{-1})}$. If $\eta_0 < \frac{\lambda}{(2\lambda+G^2)^2}$, then*

$$(\langle \hat{\Gamma}, Z \rangle - \langle \Gamma^\star, Z \rangle)^2$$
$$\lesssim \frac{N\nu^5 \log^3(NL)}{L^2\lambda^2 r} \left( 1 + \frac{D\, r \log^3 r}{N^3} \left( 1 + \log^2(NL)\frac{N}{L^2\lambda^2} \right) \frac{\log(1/\delta)}{\varepsilon^2} \right), \tag{4.4}$$

*with probability at least $1 - c_1 \exp(-c_2 D) - 4\kappa$.*

*Remark* 4.2. In the low-dimensional setting with a specific choice of $T$, the cost of gradient descent is negligible, and the cost of noise injection becomes the dominant contributor to the overall cost of privacy. Notably, the restriction on $D$ renders it irrelevant in determining the cost of privacy in this regime. On the other hand, among the many possible high-dimensional scenarios, a particularly interesting case is the over-parameterized regime where $N \asymp L^2 \asymp D^2$.

**Proposition 4.1.** *If $N \asymp L^2 \asymp D^2$, then with $\lambda \asymp \frac{N}{D}$, it holds that*

$$\left( \langle \hat{\Gamma}, Z \rangle - \langle \Gamma^\star, Z \rangle \right)^2 \lesssim \nu^5 \log^3(NL) \frac{D^2}{NL^2 r} \left( 1 + r \log^3 r \cdot \frac{D}{N^3} \cdot \frac{\log(1/\delta)}{\varepsilon^2} \right), \tag{4.5}$$

*with probability at least $1 - c_1 \exp(-c_2 D) - 4\kappa$. In particular, when $r \asymp N$, or, equivalently, $T \asymp \log N$, it holds that*

$$\left( \langle \hat{\Gamma}, Z \rangle - \langle \Gamma^\star, Z \rangle \right)^2 \lesssim \nu^5 \log^3(NL) \frac{D^2}{N^2 L^2} \tag{4.6}$$

*with probability at least $1 - c_1 \exp(-c_2 D) - 4\kappa$.*

This result is proved in Section A.5. The choice of $\lambda$ in Proposition 4.1 is standard and also appears in the ridge regression analysis of Lu et al. [2024]. Equation (4.5) highlights the trade-off between the two components of the cost of privacy, as previously discussed in Remark 4.1. For fixed values of $N$, $L$, and $D$, the test risk of the estimates generated by `NoisyHead` decreases with the number of iterations $T$ at a rate of $\Theta(e^{-T})$ whereas the test error increases at a rate of $\Theta(T^3)$. As a result, in this high-dimensional regime, the *optimal stopping point* for pretraining is $T = \Theta(\log N)$ iterations. This phenomena is explored numerically in Section 6.1.

# 5    Robustness properties of `NoisyHead`

In this section, we demonstrate that `NoisyHead` is inherently robust to adversarial perturbations to the training data. Specifically, we show that such perturbations during the pretraining stage affect the generalization error of our method significantly less than the baseline approach proposed in Lu et al. [2024].

Consider a set of training prompts $E_1, \dots, E_N$, and suppose a malicious attacker aims to degrade performance on an independent test prompt $E$ by perturbing the training data, thereby inducing inaccurate estimation of the weights in the attention module. To disrupt the training process, the attacker selects a prompt uniformly at random from the training set, say $E_i$, and replaces it with a perturbed version,

$$E_{\text{bad},i}(\mu, \alpha) = \begin{pmatrix} x'_{i,1} & x'_{i,2} & \cdots & x'_{i,L} & x'_{i,L+1} \\ y'_{i,1} & y'_{i,2} & \cdots & y'_{i,L} & 0 \end{pmatrix} \in \mathbb{R}^{(D+1) \times (L+1)}, \quad (5.1)$$

where the perturbed components are given by $x'_{i,k} = x_{i,k} + \mu$ for all $k \in [L+1]$ and $y'_{i,\ell} = y_{i,\ell} + \alpha$ for all $\ell \in [L]$. Let the parameter trained by the `NoisyHead` algorithm acting on the perturbed set of prompts $(E_1, \dots, E_{i-1}, E_{\text{bad},i}, E_{i+1}, \dots, E_N)$ be $\hat{\Gamma}_{\text{bad}}$. Correspondingly, let the parameter trained on the original, unperturbed prompts be $\hat{\Gamma}$. Let the ridge regression solutions of (2.7) on the "perturbed" and "original" set of prompts, be denoted by $\Gamma^\star_{\text{bad}}$ and $\Gamma^\star$, respectively. Then the following theorem characterizes the robustness properties of the estimates generated by `NoisyHead`.

**Theorem 5.1.** *Consider the* `NoisyHead` *algorithm with the hyper-parameter specifications as in Theorem 4.1. Further, consider an adversarial prompt perturbation as in (5.1), with $\mu, \alpha$ satisfying*

$$\alpha^2 \mu^4 \le c_u N L \lambda \quad \text{and} \quad \alpha^2 \mu^2 \ge c_\ell \mathcal{C}^2 L^{-1/2} (1 \vee \lambda N R^2 L^{-1/2}), \quad (5.2)$$

*for large enough constant $c_u > 0$ and small enough constant $c_\ell > 0$. If $\kappa > N e^{-D^2}$ and $\lambda > \mathcal{C}^2 L^{-1}$, then for an "unperturbed" test prompt $E$ and the corresponding $Z$ from (2.5), it holds that*

$$(\langle \hat{\Gamma}, Z \rangle - \langle \hat{\Gamma}_{bad}, Z \rangle)^2 \lesssim \frac{N}{L^2} \log^2(NL/\kappa) < \frac{\alpha^2 \mu^2}{N\lambda} \le (\langle \Gamma^\star, Z \rangle - \langle \Gamma^\star_{bad}, Z \rangle)^2, \quad (5.3)$$

*with probability at least $1 - c_1 \exp(-c_2 D) - 5\kappa$ for constants $c_1, c_2 > 0$.*

*Remark* 5.1. Theorem 5.1 (proved in Section A.6) shows that under bounded perturbations, pretraining with `NoisyHead` yields generalization error closer to that from the unperturbed setup than does ridge regression. If $\lambda \asymp 1$ and $D^2 \gtrsim \log N$, the bounds in (5.2) simplify to $\frac{N^2}{L^2} \lesssim \alpha^2 \mu^2 \le \alpha^2 \mu^4 \lesssim NL$. In the regime $\frac{N}{L^2} \log^2(NL) \to 0$, an adversary can choose $\alpha, \mu$ such that $\alpha^2 \mu^2 \to \infty$ while still satisfying $\alpha^2 \mu^4 \lesssim NL$, leading to $(\langle \Gamma^\star, Z \rangle - \langle \Gamma^\star_{\text{bad}}, Z \rangle)^2 \xrightarrow{\mathbb{P}} \infty$. In contrast, `NoisyHead` ensures $(\langle \hat{\Gamma}, Z \rangle - \langle \hat{\Gamma}_{\text{bad}}, Z \rangle)^2 \xrightarrow{\mathbb{P}} 0$ even under such adversarial conditions, as confirmed by experiments in Section 6.3.

# 6 Numerical experiments

We evaluate the empirical behavior of the `NoisyHead` algorithm. Section 6.1 examines how prediction risk changes under different privacy constraint strengths. Section 6.2 explores the trade-off between optimization and noise under different iteration counts. Section 6.3 validates the robustness of `NoisyHead` to adversarial perturbations. All code to reproduce the figures can be found at https://github.com/kingsleyyeon/DP.
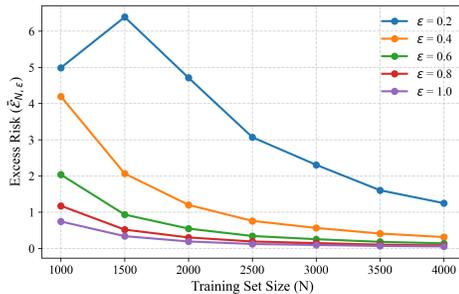


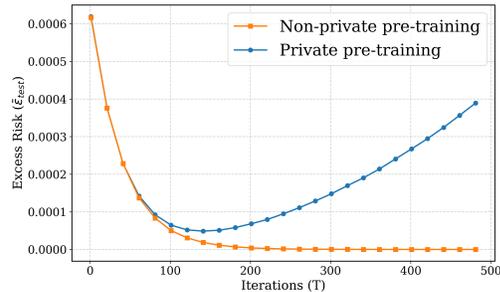Figure 1: Excess risk of `NoisyHead` for the low-dimensional set-up with $D = 5$.



Figure 2: Interplay between the cost of descent and the cost of privacy in the overparameterized setting with $N = 1000$ and $\varepsilon = 0.8$.

## 6.1 Effect of privacy on prediction risk

We evaluate the impact of privacy on ICL in a low-dimensional setting with $D = 5$. Full simulation details, as well as high-dimensional experiments, are provided in Appendix 6.1. In this experiment, we vary the number of prompts $N$ and privacy level $\varepsilon$. We use $T = \log N^{5/2} / \log(1 - \lambda\eta_0)$ and set other parameters according to Theorem 4.1. The excess test risk, averaged over $B = 500$ trials, is measured relative to that of ridge regression as $\mathcal{E}_{\text{test}} = \frac{1}{n_{\text{test}}} \sum_{k=1}^{n_{\text{test}}} \left( \langle \hat{\Gamma} - \Gamma^\star, Z_{k,\text{test}} \rangle \right)^2$, where $\hat{\Gamma}$ is the DP estimate and $\Gamma^\star$ is the ridge regression solution. As shown in Figure 1, the excess risk decreases with $N$ and increases under stricter privacy, aligning with Theorem 4.2.

## 6.2 Effect of early stopping in over-parametrized setting

Next, we fix $N = 1000$ and study how the test error evolves with the number of iterations $T$, in an over-parameterized regime with $L \asymp D \asymp \sqrt{N}$ and $\varepsilon = 0.8$. For each $T$, we average test error over 500 trials, comparing that of the differentially-private training with that of non-private gradient descent. Figure 2 illustrates a phase transition in DP training: first the error decreases with $T$ as optimization improves the solution, but then it passes a critical point and the error rises as injected noise accumulates. Early stopping around $T \approx 140$ optimally balances under-optimization and noise accumulation in this setting. This validates the need for early stopping under privacy constraints.
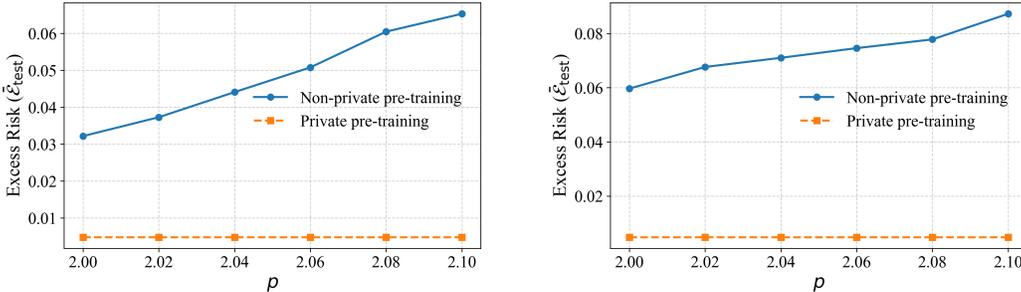


Figure 3: Comparison of prediction error under adversarial perturbations for different values of $c$. Left: $c = 2$; Right: $c = 4$. The differentially private estimator (`NoisyHead`) consistently outperforms the ridge estimator ($\Gamma^\star$) as the perturbation magnitude $\alpha = cN^p$ increases.

## 6.3 Robustness of `NoisyHead`

We test robustness by changing one training prompt with the additive perturbation $\alpha = cN^p$ (as described in Section 5) for $c \in \{2, 4\}$ and $p \in [2, 2.1]$, while fixing $N = 5000$, $L = 500$, and $D = 5$. We compare the prediction error of `NoisyHead` and ridge regression on 500 test prompts across 500 trials. Figure 3 shows that while ridge regression is increasingly affected by larger perturbations, `NoisyHead` remains robust, demonstrating its resilience to adversarial training examples.

## 7 Conclusion

Maintaining privacy during the pretraining of large language models is an increasingly important challenge as such architectures become ubiquitous. To the best of our knowledge, this work provides the first systematic theoretical characterization of differentially-private in-context learning. We quantify the *cost of privacy* on the performance of linear attention heads and formally justify the widely observed phenomenon of *early stopping* [Zhang et al., 2023, Majmudar et al., 2022, Bu et al., 2024, Bombari and Mondelli, 2025] in the context of training attention-based models under differential privacy—previously unexplored even for simplified architectures using the attention mechanism. Recent studies [Dai et al., 2023, Vladymyrov et al., 2024, Liang et al., 2025] show that multi-layered transformers can emulate gradient-based learning. Our framework offers a pathway toward understanding the theoretical behavior of such models when executing privacy-preserving pretraining, with potential implications for mitigating the "regurgitation" Carlini et al. [2021] behavior observed in large language models.

# References

M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pages 308–318. ACM, 2016.

E. Akyürek, J. Andreas, and K. G. Lin. What learning algorithm is in-context learning? investigations with linear models. In *International Conference on Learning Representations (ICLR)*, 2022. URL https://openreview.net/forum?id=rd4DBk3whU.

R. Anil, B. Ghazi, V. Gupta, R. Kumar, and P. Manurangsi. Large-scale differentially private bert. *arXiv preprint arXiv:2108.01624*, 2021.

U. Anwar, J. Von Oswald, L. Kirsch, D. Krueger, and S. Frei. Adversarial robustness of in-context learning in transformers for linear regression. *arXiv preprint arXiv:2411.05189*, 2024.

S. Black, S. Biderman, E. Hallahan, Q. Anthony, L. Gao, L. Golding, H. He, C. Leahy, K. McDonell, J. Phang, et al. Gpt-neox-20b: An open-source autoregressive language model. https://arxiv.org/abs/2204.06745, 2022.

S. Bombari and M. Mondelli. Privacy for free in the overparameterized regime. *Proceedings of the National Academy of Sciences*, 122(15):e2423072122, 2025.

T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901, 2020.

Z. Bu, X. Zhang, S. Zha, M. Hong, and G. Karypis. Pre-training differentially private models with limited public data. *Advances in Neural Information Processing Systems*, 37:94652–94683, 2024.

S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. Lee, Y. T. Li, S. Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. https://arxiv.org/abs/2303.12712, 2023.

T. T. Cai, Y. Wang, and L. Zhang. The cost of privacy: optimal rates of convergence for parameter estimation with differential privacy. *Ann. Statist.*, 49(5):2825–2850, 2021. ISSN 0090-5364,2168-8966. doi: 10.1214/21-aos2058. URL https://doi.org/10.1214/21-aos2058.

N. Carlini, F. Tramer, E. Wallace, M. Jagielski, A. Herbert-Voss, et al. Extracting training data from large language models. In *USENIX Security Symposium*, 2021. URL https://arxiv.org/abs/2012.07805.

N. Carlini, F. Tramer, E. Wallace, M. Jagielski, A. Herbert-Voss, and K. Lee. Membership inference attacks from first principles. In *Proceedings of the IEEE Symposium on Security and Privacy*, 2022. URL https://arxiv.org/abs/2112.03570.

K. Chaudhuri, C. Monteleoni, and A. D. Sarwate. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12:1069–1109, 2011.

Z. Dai, N. Golowich, C. Zhang, J. Sohl-Dickstein, and B. Neyshabur. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. URL https://openreview.net/forum?id=ceqN1LULfV.

F. Dennstädt, J. Hastings, P. M. Putora, M. Schmerder, and N. Cihoric. Implementing large language models in healthcare while balancing control, collaboration, costs and security. *npj Digital Medicine*, 8(1):143, 2025. doi: 10.1038/s41746-025-01476-7. URL https://www.nature.com/articles/s41746-025-01476-7. Available in PMC: 2025 Mar 6.

X. Dong, Y. Xu, and D. Radev. General-purpose in-context learning by meta-learning transformers. *arXiv preprint arXiv:2212.10559*, 2022. URL https://arxiv.org/abs/2212.10559.

C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *Proceedings of the 3rd Theory of Cryptography Conference (TCC)*, volume 3876 of *Lecture Notes in Computer Science*, pages 265–284. Springer, 2006. doi: 10.1007/11681878_14. URL https://link.springer.com/chapter/10.1007/11681878_14.

V. Feldman, T. Zrnic, R. Bassily, and K. Talwar. Privacy amplification by iteration. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 31, 2018.

S. Garg, X. Li, T. Zhou, and S. Ermon. What can transformers learn in-context? a case study of simple function classes. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/hash/fd09bd8a5e61ae58b07b9052d9e4a6e4-Abstract-Conference.html.

S. Gopi, P. Jain, P. K. Kothari, and A. G. Thakurta. Dp-agd: Private adaptive gradient descent with optimal utility. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.

X. Gu, W. Li, X. Zhang, Q. Wang, B. Ding, and Y. Zhang. Differential privacy mechanisms in neural tangent kernel regression. *arXiv preprint arXiv:2407.13621*, 2024. URL https://arxiv.org/abs/2407.13621.

S. Hoory, A. Feder, A. Tendler, S. Erell, A. Peled-Cohen, I. Laish, H. Nakhost, U. Stemmer, A. Benjamini, A. Hassidim, et al. Learning and evaluating a differentially private pre-trained language model. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1178–1189, 2021.

P. Jain and A. Thakurta. Differentially private learning with kernels. In *Proceedings of the 30th International Conference on Machine Learning*, pages 118–126. PMLR, 2013. URL https://proceedings.mlr.press/v28/jain13.html.

T. Ji, C. Luo, Y. Guo, J. Ji, W. Liao, and P. Li. Differentially private community detection in attributed social networks. In *Proceedings of The Eleventh Asian Conference on Machine Learning*, pages 16–31. PMLR, 2019. URL https://proceedings.mlr.press/v101/ji19a.html.

G. Kamath and J. Ullman. A primer on private statistics. *arXiv preprint arXiv:2005.00010*, 2020. URL https://arxiv.org/abs/2005.00010.

L. Ke, S. Tong, P. Cheng, and K. Peng. Exploring the frontiers of llms in psychological applications: A comprehensive review. *arXiv preprint arXiv:2401.01519*, 2024. URL https://arxiv.org/abs/2401.01519.

X. Li, T. Ma, R. Zhang, M. Ganjali, and O. Mir. Large language models can be strong differentially private learners. In *International Conference on Learning Representations (ICLR)*, 2022. URL https://arxiv.org/abs/2110.05679.

X. Li, V. Patel, S. Kumar, et al. Privacy preserving strategies for electronic health records in the era of large language models. *Journal of Biomedical Informatics*, 137:104567, Jan 2025. doi: 10.1016/j.jbi.2024.104567. URL https://pubmed.ncbi.nlm.nih.gov/39820020/.

H. Liang, K. Balasubramanian, and L. Lai. Transformers handle endogeneity in in-context linear regression. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=QfhU3ZC2g1.

D. Lieber, T. Wolf, I. Golan, A. Shmidman, O. Sharir, and Y. Shoham. Jurassic-1: Technical details and evaluation. https://www.ai21.com/blog/jurassic-1-open-access, 2021. AI21 Labs Whitepaper.

Y. Lu, M. Letey, J. A. Zavatone-Veth, A. Maiti, and C. Pehlevan. In-context learning by linear attention: Exact asymptotics and experiments. In *NeurIPS 2024 Workshop on Mathematics of Modern Machine Learning*, 2024.

J. Majmudar, C. Dupuy, C. Peris, S. Smaili, R. Gupta, and R. Zemel. Differentially private decoding in large language models. *arXiv preprint arXiv:2205.13621*, 2022.

B. Ni, N. Li, and W. Li. Detecting communities under differential privacy. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*, pages 597–606. ACM, 2016. doi: 10.1145/2994620.2994624. URL https://dl.acm.org/doi/10.1145/2994620.2994624.

J. W. Rae, S. Borgeaud, T. Cai, K. Millican, J. Hoffmann, F. Song, J. Aslanides, S. Henderson, R. Ring, S. Young, et al. Scaling language models: Methods, analysis & insights from training gopher. https://arxiv.org/abs/2112.11446, 2021.

P. Rigollet and J.-C. Hütter. High-dimensional statistics. *arXiv preprint arXiv:2310.19244*, 2023.

A. Tirumala, Adithya Renduchintala et al. Memorization in large language models: Quantifying, understanding, and reducing. In *arXiv preprint arXiv:2202.07646*, 2022. URL https://arxiv.org/abs/2202.07646.

M. Vladymyrov, J. V. Oswald, M. Sandler, and R. Ge. Linear transformers are versatile in-context learners. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=p1ft33Mu3J.

D. Wang and J. Xu. On sparse linear regression in the local differential privacy model. In *Proceedings of the 36th International Conference on Machine Learning*, pages 6628–6637. PMLR, 2019. URL https://proceedings.mlr.press/v97/wang19m.html.

Y. Wang, H. Zhang, and Y. Liu. In-context learning for attention scheme: from single softmax regression to multiple softmax regression via a tensor trick. In *International Conference on Learning Representations (ICLR)*, 2023. URL https://openreview.net/forum?id=j-D4N9T4U3.

Y. Xie, Y. Lu, Z. Lin, J. Wei, E. Zhang, C. Raffel, L. Kong, T. Hashimoto, and C. D. Manning. An explanation of in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2211.15661*, 2022. URL https://arxiv.org/abs/2211.15661.

X. Yu, S. Gopi, et al. Differentially private fine-tuning of language models. *arXiv preprint arXiv:2303.08774*, 2023. URL https://arxiv.org/abs/2303.08774.

R. Zhang, S. Frei, and P. L. Bartlett. Trained transformers learn linear models in-context. *J. Mach. Learn. Res.*, 25:Paper No. [49], 55, 2024. ISSN 1532-4435,1533-7928.

T. Zhang, T. Zhu, K. Gao, W. Zhou, and P. S. Yu. Balancing learning model privacy, fairness, and accuracy with early stopping criteria. *IEEE Trans. Neural Netw. Learn. Syst.*, 34(9):5557–5569, 2023. ISSN 2162-237X,2162-2388. doi: 10.1109/tnnls.2021.3129592. URL https://doi.org/10.1109/tnnls.2021.3129592.

X. Zhang, K. Talwar, and D. Zhang. Understanding the difficulty of training transformers with differential privacy. *arXiv preprint arXiv:2104.05800*, 2021.

# A  Theoretical Details

## A.1  Choice of parameters

Yu et al. [2023] provides an intuitive explanation for the behavior of the predictor (2.4). The term $w_{22}^{PV} W_{11}^{KQ}$ is approximately equal to $\mathbb{E}[(X^\top X)^{-1}]$, capturing the inverse second-moment structure of the features. The second term does not depend on the features, and the third term is independent of the labels $y$. They act as an extra additive term, which can be assumed to have no significant impact on the final prediction. The fourth term represents the effect of projecting the input features $x_i$ onto the direction $w_{12}^{PV}$ in the final prediction. However, since the features are assumed to be isotropic, it is reasonable to expect that projections onto any particular direction carry no special predictive value. Consequently, it is justified to assume that $w_{12}^{KQ} = 0$ and $w_{12}^{PV} = 0$, which simplifies the predictor to:

$$\hat{y}_{L+1} = \frac{1}{L} \left\langle x_{L+1}, w_{22}^{PV} W_{11}^{KQ} \sum_{i=1}^{L} y_i x_i \right\rangle.$$

This assumption is further supported by the observation in Zhang et al. [2024], where the authors show that when the parameters of $W^{PV}$ and $W^{KQ}$ are learned via gradient flow on the average reconstruction loss $\mathbb{E}[(\widehat{y} - y)^2]$, initializing with $w_{12}^{KQ} = 0$ and $w_{12}^{PV} = 0$ ensures that the parameter remains zero throughout training.

## A.2    Proof of Theorem 3.2

Consider two datasets of prompts $(E_k)_{k \in [N]}$ and $(E'_k)_{k \in [N]}$ that differ in exactly one prompt. Without loss of generality, assume $E_1 \neq E'_1$ and $E_k = E'_k$ for all $k \geq 2$. The change in the gradient update due to this single difference is bounded by:

$$\frac{\eta_0}{N}\bigg( \|\langle \hat{\Gamma}, \widetilde{Z}_1 \rangle \widetilde{Z}_1\|_F + \|\langle \hat{\Gamma}, \widetilde{Z}'_1 \rangle \widetilde{Z}'_1\|_F + \|\texttt{clip}_{\mathcal{C}}(y_{1,L+1})\widetilde{Z}_1\|_F + \|\texttt{clip}_{\mathcal{C}}(y'_{1,L+1})\widetilde{Z}'_1\|_F \bigg)$$

$$\leq \frac{2\eta_0(RG^2 + \mathcal{C}G)}{N} \leq \frac{\eta_0 \sigma}{N}, \tag{A.1}$$

where the final inequality follows from the assumption on $\sigma$.

By Lemma 2.5 of Kamath and Ullman [2020], each gradient step in Algorithm 1 is $(\varepsilon/T, \delta/T)$-differentially private. The overall guarantee then follows by composition, using Fact 2.2 of Cai et al. [2021].

## A.3    Proof of Theorem 4.1

Consider the set $\mathcal{D}_1 := \{\|\Gamma^\star\|_F \leq R\}$. Moreover, denote

$$\widetilde{\Gamma}^{t+1} = (1 - 2\lambda)\Gamma^t - \eta_0 N^{-1} \sum_{k=1}^{N} \Big( \langle \Gamma^t, \widetilde{Z}_k \rangle - \texttt{clip}(y_{L+1}) \Big) \widetilde{Z}_k.$$

Clearly, $\Gamma^{t+1} = \Pi_R(\widetilde{\Gamma}^{t+1} + \boldsymbol{z}_t)$. Under $\mathcal{D}_1$, it is easy to see that

$$\|\hat{\Gamma} - \Gamma^\star\|_F^2 \leq \|\widetilde{\Gamma}^T + \boldsymbol{z}_{T-1} - \Gamma^\star\|_F^2 \leq (1 + C_0^{-1})\|\widetilde{\Gamma}^T - \Gamma^\star\|_F^2 + (1 + C_0)\|\boldsymbol{z}_{T-1}\|_F^2, \tag{A.2}$$

where, the choice of the constant $C_0$ ensures

$$(1 + C_0^{-1})\kappa < 1 - \eta_0 \lambda \,, \text{ with } \kappa := 1 - 2\eta_0 \lambda + \eta_0^2 (G^2 + 2\lambda)^2. \tag{A.3}$$

Further consider the sets $\mathcal{D}_2 := \Big\{ \max_{k \in [N]} \big\| \sum_{i=1}^{L} x_{k,i} \big\|_2 \leq GL\mathcal{C}^{-1} \Big\}$, and $\mathcal{D}_3 := \big\{ \max_{k \in [N], i \in [L+1]} |y_{k,i}| \leq \mathcal{C} \big\}$. Since $\|x_{k,L+1}\|_2 = 1$, under the events $\mathcal{D}_2$ and $\mathcal{D}_3$, it follows that

$$\max_{k \in [N]} \bigg\| L^{-1} x_{k,L+1} \sum_{i=1}^{L} y_{k,i} x_{k,i}^\top \bigg\|_F \leq G, \tag{A.4}$$

which implies $\widetilde{Z}_k = Z_k$ for all $k \in [N]$ by the definition of $Z_k$ in (2.5). The sets $\mathcal{D}_i, i = 1, 2, 3$ allow us to bear down the classical theory of convex minimization, and our choice of the parameters $R, G$ and $\mathcal{C}$ will emphasize that these events occur with high probability. In particular, under $\mathcal{D} := \cap_{i=1}^{3} \mathcal{D}_i$, we note the $\mathcal{L}$ is $\lambda$-strongly convex:

$$\langle \nabla_\Gamma \mathcal{L}(\Gamma, (Z_k)_{k \in [N]}), \Gamma - \Gamma^\star \rangle \geq \lambda \|\Gamma - \Gamma^\star\|_F^2, \tag{A.5}$$

and the $(G^2 + 2\lambda)$-smooth:

$$\big\| \nabla_\Gamma \mathcal{L}(\Gamma, (Z_k)_{k \in [N]}) - \nabla_\Gamma \mathcal{L}(\Gamma', (Z_k)_{k \in [N]}) \big\|_F \leq (G^2 + 2\lambda) \|\Gamma - \Gamma'\|_F. \tag{A.6}$$

Therefore, for the term $\|\widetilde{\Gamma}^T - \Gamma^\star\|_F$ in (A.2),

$$\|\widetilde{\Gamma}^T - \Gamma^\star\|_F^2 = \|\Gamma^{T-1} - \eta_0 \nabla_{\Gamma^{T-1}} \mathcal{L}(\Gamma^{T-1}, (Z_k)_{k \in [N]}) - \Gamma^\star\|_F^2 \leq \kappa \|\Gamma^{T-1} - \Gamma^\star\|_F^2, \tag{A.7}$$

where we recall $\kappa$ from (A.3), and (A.7) employs (A.5) and (A.6). Note that we must require $\kappa < 1$, which makes use of $\eta_0 < \frac{\lambda}{(G^2 + 2\lambda)^2}$. Putting (A.7) back into (A.2), one obtains under $\mathcal{D}$ that

$$\|\hat{\Gamma} - \Gamma^\star\|_F^2 \leq (1 + C_0^{-1})\kappa \|\Gamma^{T-1} - \Gamma^\star\|_F^2 + (1 + C_0)\|\boldsymbol{z}_{T-1}\|_F^2.$$

13

Proceeding recursively, we can show that for all $T > 1$, we have

$$\|\hat{\Gamma} - \Gamma^\star\|_F^2 \leq (1 - \eta_0\lambda)^T R^2 + (1 + C_0) \sum_{i=0}^{T-1} (1 - \eta_0\lambda)^{T-i-1} \|z_i\|_F^2. \tag{A.8}$$

Since the errors $(z_i)_{i=1}^T$ are independent of the prompts $(E_k)_{k\in[B]}$, an application of Lemma A.2. of Cai et al. [2021] implies

$$\|\hat{\Gamma} - \Gamma^\star\|_F^2 \lesssim (1 - \eta_0\lambda)^T R^2 + \sigma^2\eta_0^2 D \frac{T^2 \log(2T/\delta)}{N^2\varepsilon^2},$$

with probability at least $1 - c_1\exp(-c_2 D)$ under $\mathcal{D}$. An application of Cauchy-Schwarz inequality entails

$$(\langle\hat{\Gamma}, Z\rangle - \langle\Gamma^\star, Z\rangle)^2 \leq G_0^2\left((1 - \eta_0\lambda)^T R^2 + \sigma^2\eta_0^2 D \frac{T^2 \log(2T/\delta)}{N^2\varepsilon^2}\right) \tag{A.9}$$

with probability at least $1 - c_1\exp(-c_2 D)$ under $\mathcal{D} \cap \{\|Z\|_F \leq G_0\}$. Now we turn to tackling the individual events $\mathcal{D}_i$, $i = 1, 2, 3$. For $\mathcal{D}_3$, note that if $(x_i)_{i\in[L]} \overset{i.i.d.}{\sim} \mathcal{U}(\mathbb{S}^{D-1})$ and $w \sim N(0, \mathbb{I}_D)$ independently of $x_i$'s, then $(w^\top x_i) \sim N(0, 1)$ marginally. Therefore, Lemma A.1 implies

$$\mathbb{P}(\mathcal{D}_3) \geq 1 - \kappa, \quad \text{for } \mathcal{C} = \sqrt{2\nu\log(NL/\kappa)}. \tag{A.10}$$

Furthermore, with $G \asymp \frac{\mathcal{C}}{\sqrt{L}}\left(1 + \left(\frac{\log(N/\kappa)}{D^2}\right)^{1/4}\right)$, from Lemma A.2 we get that

$$\mathbb{P}(\mathcal{D}_2) \geq 1 - \kappa. \tag{A.11}$$

Finally, noting that

$$\sum_{k=1}^N y_{k,L+1} \operatorname{vec}(Z_k) \leq \max_{k\in[N]} |y_{k,L+1}| \left(\max_{k\in[N],i\in[L]} |y_{k,i}|\right) \frac{1}{L}\sum_{k,i} \operatorname{vec}(x_{k,L+1}x_{k,i}^\top),$$

an application of Lemma A.1 on (2.7), in conjunction with Lemma A.3, yields

$$\mathbb{P}(\mathcal{D}_1) \geq 1 - \kappa, \text{ with } R = \lambda^{-1}\mathcal{C}^2\sqrt{\frac{N}{L}}\left(1 + \left(\frac{\log(1/\kappa)}{D^2}\right)^{1/4}\right). \tag{A.12}$$

Finally, similar to Lemma A.2 it can be argued that

$$\mathbb{P}(\|Z\|_F \leq G_0) \geq 1 - \kappa, \quad \text{for } G_0 \asymp \frac{\mathcal{C}}{\sqrt{L}}\left(1 + \left(\frac{\log(1/\kappa)}{D^2}\right)^{1/4}\right). \tag{A.13}$$

Summarizing (A.10)-(A.12), it holds that

$$\mathbb{P}(\mathcal{D} \cap \{\|Z\|_F \leq G_0\}) \geq 1 - 4\kappa. \tag{A.14}$$

Putting these bounds back into (A.9), we invoke (A.14) to conclude (4.2).

## A.4 Proof of Theorem 4.2

(i) *Low-dimensional setting.* Recall $T = \frac{\log(N^2 L D^3)}{\log((1-\eta_0\lambda)^{-1})}$. Note that, with $\kappa > e^{-D^2}$, we have $G_0 \lesssim \mathcal{C}/\sqrt{L}$, $R \lesssim \lambda^{-1}\mathcal{C}^2\sqrt{N/L}$, $\eta_0 \asymp \frac{\lambda}{(\lambda+G^2)^2} \lesssim 1/\lambda$ and $\lambda \asymp 1 \asymp \eta_0$ from (4.1). Hence, the first term of (4.3) can be bounded as

$$G_0^2(1 - \eta_0\lambda)^T R^2 \lesssim \frac{\mathcal{C}^6}{NL^3 D^3} \lesssim \nu^3\log^3\frac{NL}{\kappa}\frac{1}{NL^3 D^3}.$$

Moreover, from $\log(\frac{1}{\kappa}) \lesssim D^2 \lesssim \log(NL)$, we have that

$$G_0^2(1 - \eta_0\lambda)^T R^2 \lesssim \nu^3(\log^3 NL) \cdot \frac{1}{NL^3}. \tag{A.15}$$

14

On the other hand, write the second term as

$$G_0^2 \sigma^2 \eta_0^2 D \frac{T^2}{N^2} \frac{\log(2T/\delta)}{\varepsilon^2} \lesssim \frac{\mathcal{C}^2}{L} (\mathcal{C}G + RG^2)^2 D \frac{T^3}{N^2} \frac{\log(1/\delta)}{\varepsilon^2}. \qquad \text{(A.16)}$$

Clearly, for $\sigma$, one obtains,

$$\mathcal{C}G + RG^2 \lesssim \frac{\mathcal{C}^2}{\sqrt{L}} \left( 1 + \left( \frac{\log(NL)}{D^2} \right)^{1/2} \right) + \mathcal{C}^2 \sqrt{\frac{N}{L}} \frac{\mathcal{C}^2}{L} \left( 1 + \left( \frac{\log(NL)}{D^2} \right) \right)$$

$$\lesssim \nu^2 \log^3(NL) \frac{\sqrt{N}}{LD^2},$$

where the second inequality is attained by using $\log(N/\kappa) = \log N + \log 1/\kappa \lesssim \log N + D^2 \lesssim \log NL$, and the final assertion follows from $(\log NL)/D^2 >> (\sqrt{\log NL})/D$. Therefore, from (A.16), the second term is bounded by,

$$\lesssim \frac{\nu^5 \log^7(NL)}{L} \cdot \frac{N}{L^2 D^4} \cdot D \cdot \frac{\log^3(N^2 LD^3)}{N^2} \cdot \frac{\log(1/\delta)}{\varepsilon^2}$$

$$\lesssim \frac{\nu^5 \log^7(NL) \log^3(N^2 LD^3)}{NL^3 D^3} \cdot \frac{\log(1/\delta)}{\varepsilon^2}$$

$$\lesssim \frac{\nu^5 \log^{10}(NL)}{NL^3} \cdot \frac{\log(1/\delta)}{\varepsilon^2}. \qquad \text{(A.17)}$$

Combining (A.15) and (A.17) yields the proof for the low-dimensional case.

(ii) *High-dimensional setting.* Here, $\kappa \gtrsim (NL)^{-1}$, and $D^2 \gtrsim \log(NL)$ implies that $\frac{\log(NL/\kappa)}{D^2} \lesssim 1$. We also have $G \lesssim \mathcal{C}/\sqrt{L}$, $G_0 \lesssim \mathcal{C}/\sqrt{L}$ and, $R \lesssim \lambda^{-1} \mathcal{C}^2 \sqrt{N/L}$. The first term of (4.4) can be bounded as

$$G_0^2 R^2 (1 - \eta_0 \lambda)^T \lesssim \frac{\mathcal{C}^2}{L} \cdot \mathcal{C}^4 \frac{N}{L} \cdot \frac{1}{\lambda^2} \cdot \frac{1}{r} \lesssim \nu^3 \log^3(NL) \cdot \frac{N}{L^2 \lambda^2 r}. \qquad \text{(A.18)}$$

Furthermore, for the second term, observe that $\eta_0 \lesssim 1/\lambda$, and

$$\sigma = G(\mathcal{C} + RG) \asymp \frac{\mathcal{C}^2}{\sqrt{L}} + \mathcal{C}^2 \sqrt{\frac{N}{L}} \cdot \frac{1}{\lambda} \cdot \frac{\mathcal{C}^2}{L} = \frac{\mathcal{C}^2}{\sqrt{L}} \left( 1 + \frac{\sqrt{N}}{\lambda L} \mathcal{C}^2 \right).$$

Therefore, the second term can be bounded as

$$G_0^2 \sigma^2 \eta_0^2 D \frac{T^3}{N^2} \leq \frac{\mathcal{C}^2}{L} \cdot \sigma^2 \frac{1}{\lambda} D \cdot \frac{T^3}{N^2} \leq \frac{\mathcal{C}^2}{L} \cdot \log^3 r \cdot \frac{D}{N^2 \lambda^2} \cdot \sigma^2$$

$$\leq \frac{\mathcal{C}^2}{L} \log^3 r \frac{D}{N^2 \lambda^2} \cdot \frac{\mathcal{C}^4}{L} \left( 1 + \frac{N}{\lambda^2 L^2} \mathcal{C}^4 \right)$$

$$\leq \frac{\mathcal{C}^6 D}{N^2 L^2 \lambda^2} \left( 1 + \frac{N}{\lambda^2 L^2} \log^2(NL) \right) \log^3 r$$

$$\lesssim \nu^3 \frac{N \log^3 r}{L^2 \lambda^2} \log^3 NL \frac{D}{N^3} \left( 1 + \frac{N}{\lambda^2 L^2} \log^2 NL \right). \qquad \text{(A.19)}$$

Assertions (A.18) and (A.19) conclude the proof.

## A.5  Proof of Proposition 4.1

From $\lambda \asymp N/D \asymp \sqrt{N}$, it follows that $\frac{D^2}{NL^2} \asymp \frac{N}{L^2 \lambda^2} \asymp N^{-1}$, and hence, $\frac{\log^2(NL)}{\lambda^2} \lesssim 1$. Therefore, from (4.4), (4.5) follows trivially. Further, the first term in (4.5) dominates as long as $r \log^3 r \ll \frac{N^3}{D}$, yielding (4.6) when $r \asymp N$.

## A.6 Proof of Theorem 5.1

Recall the definition of $G_0$ and $R$ from Theorem 4.1. In view of $\kappa > N \exp(-D^2)$, (A.4) and $\|\hat{\Gamma}\|_F \vee \|\Gamma^{T^{\text{bad}}}\|_F \leq R$, using (A.13), it holds that

$$(\langle \hat{\Gamma}, Z \rangle - \langle \hat{\Gamma}_{\text{bad}}, Z \rangle)^2 \leq \frac{\mathcal{C}^2}{L} R^2, \tag{A.20}$$

with probability at least $1 - \kappa$. On the other hand, for the analysis of the ridge estimates, recall (2.7). Clearly, from Lemma A.3, it holds with probability $\geq 1 - 2\kappa$ that

$$\left\| \sum_{k=1}^{N} \text{vec}(Z_k)\, \text{vec}(Z_k)^\top \right\|_F \leq \mathcal{C}^2 \frac{N}{L} < N\lambda, \tag{A.21}$$

where the final equality follows from $\lambda > \mathcal{C}^2 L^{-1}$. Moreover, from Lemma A.2, it holds with probability $\geq 1 - 2\kappa$ that

$$\left\| \text{vec}(Z_{\text{bad},i})\, \text{vec}(Z_{\text{bad},i})^\top \right\|_F \lesssim \frac{\mathcal{C}^2}{L} + \frac{\alpha^2 \mu^4}{L} \asymp \frac{\alpha^2 \mu^4}{L} \lesssim N\lambda, \tag{A.22}$$

where the first part of the inequality follows from the lower bound on $\alpha^2 \mu^2$ and the second inequality follows from the upper bound on $\alpha^2 \mu^4$ as stated in (5.2). Consequently, combining (2.7) with (A.21) and (A.22) jointly yields,

$$\| \text{vec}(\Gamma_{\text{bad}}^\star - \Gamma^\star) \| \geq \frac{\| y'_{i,L+1} \text{Vec}(Z'_i) - y_{i,L+1} \text{Vec}(Z_i) \|}{N\lambda} \tag{A.23}$$

with probability at least $1 - 4\kappa$. Since $\| y_{i,L+1} \text{vec}(Z_i) \| \leq \frac{\mathcal{C}^2}{\sqrt{L}}$ with probability at least $1 - \kappa$, invoking (5.2), yet another application of Lemma A.2 yields

$$\| y'_{i,L+1} \text{Vec}(Z_{\text{bad},i}) - y_{i,L+1} \text{Vec}(Z_i) \| \geq \alpha^2 \mu^2 \tag{A.24}$$

with probability at least $1 - \kappa$. Since (5.2) also implies $\frac{\alpha^2 \mu^2}{N\lambda} > R^2 \frac{\mathcal{C}^2}{L}$, from (A.20), (A.22), (A.23), and (A.24), we obtain (5.3).

## A.7 Auxiliary Lemmas

The following lemmas are instrumental to proving our theorems 4.1 and 5.1, and hereby are listed. In particular, Lemma A.1 and A.3 follows using Hoeffding's inequality and a union bound argument.

**Lemma A.1.** *If $z_{kj} \sim N(0, 1+\tau^2)$ $k \in [N], j \in [L]$ are not necessarily independent, then*

$$\mathbb{P}\left( \max_{k,j} |z_{ij}| \lesssim \sqrt{(1+\tau^2)\log\left(\frac{4NL}{\kappa}\right)} \right) \geq 1 - \kappa.$$

For the Lemmas A.2 and A.3, note that for any vector $x \in \mathbb{R}^D$, the Euclidean norm $\|x\|_2 = \sup_{a \in \mathbb{S}^{D-1}} a^\top x$. For any fixed $a \in \mathbb{S}^{D-1}$ and $k \in [N]$, $a^\top \sum_{i=1}^{L} x_{k,i}$ is a sub-Gaussian random variable with variance proxy $L/D$. Therefore

$$\mathbb{P}\left[ \left| a^\top \sum_{i=1}^{L} x_{k,i} \right| > \sqrt{L/D}\, t \right] \lesssim \exp\left(-t^2\right).$$

Therefore, using a covering number argument similar to Theorem 1.19 of Rigollet and Hütter [2023] one can show the following.

**Lemma A.2.** *Suppose $(x_{k,i})_{k \in [N], i \in [L]} \overset{i.i.d.}{\sim} \mathcal{U}(\mathbb{S}^{D-1})$. Then,*

$$\mathbb{P}\left( \max_{k \in [N]} \| \sum_{i=1}^{L} x_{k,i} \|_2 \lesssim \sqrt{L}(1 + D^{-1/2}(\log(\frac{N}{\kappa}))^{1/2}) \right) \geq 1 - \kappa.$$

**Lemma A.3.** *Suppose $(x_{k,i})_{k \in [N], i \in [L]} \overset{i.i.d.}{\sim} \mathcal{U}(\mathbb{S}^{D-1})$. Then,*

$$\mathbb{P}\left( \| \sum_{k=1}^{N} \sum_{i=1}^{L} \text{vec}(x_{k,L+1} x_{k,i}^\top) \|_2 \lesssim \sqrt{NL}(1 + D^{-1}(\log(\frac{1}{\kappa}))^{1/2}) \right) \geq 1 - \kappa.$$
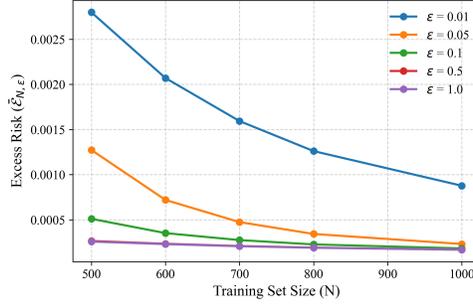
Figure 4: Excess risk of `NoisyHead` as a function of training set size $N$ for different values of the privacy parameter $\varepsilon$ with $D = \lfloor\sqrt{N}\rfloor$.

## B  Numerical Experiments Details

This section details and extends upon the numerical examples presented in Section 6.

### B.1  Effect of privacy on prediction risk: low- vs. high-dimensional regimes

In Section 6.1, we empirically investigate how the level of privacy, parameterized by $\varepsilon$, affects the prediction accuracy of `NoisyHead` through its impact on the excess risk.

**Low-dimensional regime.** We first consider the low-dimensional regime with feature dimension fixed at $D = 5$. Training set sizes are varied over $N \in \{1000, 1500, 2000, 2500, 3000, 3500, 4000\}$, with prompt length set as $L = \lfloor\sqrt{N}\rfloor$ and privacy levels $\varepsilon \in \{0.2, 0.4, 0.6, 0.8, 1.0\}$. The hyperparameters $\mathcal{C}, \mathcal{G}, \mathcal{R}$ are chosen according to Theorem 4.1, with $\kappa = 1$ and $\delta = 10^{-5}$. The step size is set as $\eta_0 = 3.17/(5 + G^2)^2$, where $G$ denotes an upper bound on the norm of the projected features $\widetilde{Z}$, and the ridge regularization parameter is fixed at $\lambda = 5$. We work in a noiseless setting with $\tau^2 = 0$.

For each $(N, \varepsilon)$ pair, we generate $N$ prompts according to (2.1) and run $T$ iterations of `NoisyHead`, where $T = \log N^{5/2}/\log(1 - \lambda\eta_0)$, as prescribed by Theorem 4.2. Test performance is evaluated on $n_{\text{test}} = 500$ held-out prompts. Each test prediction is computed as $\langle\hat{\Gamma}, Z_{\text{test}}\rangle$, where $Z_{\text{test}}$ is constructed using (2.5). The excess risk is defined as

$$\mathcal{E}_{\text{test}} = \frac{1}{n_{\text{test}}} \sum_{k=1}^{n_{\text{test}}} \left(\langle\hat{\Gamma} - \Gamma^\star, Z_{k,\text{test}}\rangle\right)^2, \tag{B.1}$$

where $\Gamma^\star$ denotes the non-private ridge estimator trained on the same data. We repeat the entire procedure $B = 500$ times and report the average excess risk $\bar{\mathcal{E}}_{N,\varepsilon}$.

The left panel of Figure 4 illustrates that for each fixed $\varepsilon$, the excess risk decreases with $N$ at a super-quadratic rate, in agreement with Theorem 4.2(i). For fixed $N$, the excess risk also decreases with increasing $\varepsilon$, highlighting the trade-off between privacy and predictive accuracy in this regime.

**High-dimensional regime.** We also consider the high-dimensional regime where $D \asymp L \asymp \sqrt{N}$. We vary $N \in \{500, 600, 700, 800, 900, 1000\}$ and $\varepsilon \in \{0.01, 0.05, 0.1, 0.5, 1.0\}$. The ridge regularization parameter is set as $\lambda = N/D$, and we use a fixed number of iterations $T = 5$ with step size $\eta = 0.07ND/(N + DG^2)^2$. All other parameters mirror those used in the low-dimensional setting. The average excess risk, computed over $B = 500$ repetitions, is reported in the right panel of Figure 4. The excess risk decreases with both $N$ and $\varepsilon$, though at a slower rate than in the low-dimensional case, consistent with Theorem 4.2 and reflecting the increased challenge of private learning in high dimensions.

### B.2  Effect of early stopping in over-parametrized setting

In Section 6.2, we investigate how the number of gradient descent steps $T$ affects the test performance of the linear attention head trained using `NoisyHead`. We fix $N = 1000$ and consider the overparam-

eterized setting with $L \asymp D \asymp \lfloor \sqrt{N} \rfloor$, under fixed privacy parameters $\varepsilon = 0.8$ and $\delta = 10^{-5}$. The step size is set as $\eta_0 = 0.007\lambda/(\lambda + G^2)^2$, and remaining hyperparameters follow the setup from the previous experiment. We vary $T$ over $\{1, 20, 40, \ldots, 480\}$ and compute the average test error over $500$ held-out prompts, repeated over $B = 500$ independent trials.

Figure 2 plots the evolution of two components of the prediction error: the *cost of descent* (blue) incurred by underoptimization, and the *cost of privacy* (orange) due to noise injection. For small $T$, the descent cost dominates and the error decreases with additional optimization. However, beyond a critical number of iterations, the cost of privacy dominates, causing error to increase as more noise accumulates. This trade-off, predicted theoretically in Remark 4.2, yields a phase transition in the test error under privacy constraints. In contrast, in the noiseless setting (approximating ridge regression), the error decreases monotonically with $T$.

### B.3  Robustness of `NoisyHead`

In Section 6.3 we evaluate the robustness of `NoisyHead` under adversarial perturbations, following the setup in Section 5. A single training prompt is perturbed by adding 1 to all features and $\alpha = cN^p$ to all responses, with $c \in \{2, 4\}$ and $p \in \{2, 2.02, 2.04, 2.06, 2.08, 2.1\}$. We fix $N = 5000$, $L = 500$, $D = 5$, $\varepsilon = 0.5$, and $\delta = 10^{-2}$. Generalization error is measured via (5.3). We compare the ridge estimator $\Gamma^\star$ with the output of `NoisyHead` after $T = \log N$ iterations, using $\lambda = 0.01$ and step size $\eta_0 = 0.007/(0.01 + G^2)^2$, with all other parameters unchanged. Figure 3 reports the average prediction error over 500 test prompts, averaged over 500 trials. As $p$ increases, ridge regression becomes increasingly sensitive to the perturbation, while differentially private pretraining with `NoisyHead` remains substantially more robust.