
Initialization Matters: Privacy-Utility Analysis of Overparameterized Neural Networks

Anonymous Author(s)
Affiliation
Address
email

Abstract

1 We analytically investigate how overparameterization of models in randomized
2 machine learning algorithms impacts the information leakage about their training
3 data. Specifically, we prove a privacy bound for the KL divergence between model
4 distributions on worst-case neighboring datasets, and explore its dependence on
5 the initialization, width, and depth of fully connected neural networks. We find
6 that this KL privacy bound is largely determined by the expected squared gradient
7 norm relative to model parameters during training. Notably, for the special setting
8 of linearized network, our analysis indicates that the squared gradient norm (and
9 therefore the escalation of privacy loss) is tied directly to the per-layer variance of
10 the initialization distribution. By using this analysis, we demonstrate that privacy
11 bound improves with increasing depth under certain initializations (LeCun and
12 Xavier), while degrades with increasing depth under other initializations (He and
13 NTK). Our work reveals a complex interplay between privacy and depth that
14 depends on the chosen initialization distribution. We further prove excess empirical
15 risk bounds under a fixed KL privacy budget, and show that the interplay between
16 privacy utility trade-off and depth is similarly affected by the initialization.

17 1 Introduction

18 Deep neural networks (DNNs) in the over-parameterized regime (i.e., more parameters than data)
19 perform well in practice but the model predictions can easily leak private information about the
20 training data under inference attacks such as membership inference attacks [38] and reconstruction
21 attacks. [14, 5, 23] This leakage can be mathematically measured in terms of how much the algorithm’s
22 output distribution changes if it were trained on a neighboring dataset (that only differs in one record),
23 following the differential privacy (DP) framework [18].

24 To train differential private model, a typical way is randomize each gradient update in neural networks
25 training, e.g., stochastic gradient descent (SGD), which leads to the most widely applied differentially
26 private training algorithm in the literature – DP-SGD [1]. In each step, DP-SGD employs gradient
27 clipping and adds calibrated Gaussian noise, and thus it comes with a differential privacy guarantee
28 that scales with the noise multiplier (i.e., per-dimensional Gaussian noise standard deviation divided
29 by the clipping threshold) and number of training epochs. However, this privacy bound [1] is overly
30 general due to its independence on the network properties (e.g., width and depth) and training schemes
31 (e.g., initializations). Accordingly, a natural question arises in the community:

32 *How does overparameterization (e.g., increasing width and depth) of neural networks affect the*
33 *(worst-case) privacy bound of the training algorithm?*

Table 1: Our results for the privacy utility trade-off of training linearized network (3) via Langevin diffusion, under different width m , depth L and initializations. We set per-layer width $m_0 = d$, $m_1, \dots, m_{L-1} = m$ and $m_L = o$. We prove privacy bound in KL divergence, and obtain excess empirical risk bounds given KL privacy budget ε . For the excess risk bounds, we assume the network width $m = \Omega(n)$ is sufficiently large, and the data and network satisfy regularity assumption Assumption 2.1. For NTK, He and LeCun initialization, we observe that the privacy utility trade-off improves with overparameterization (increasing depth).

Init	Variance β_l for layer l	Gradient norm constant B (9)	Approximate lazy training distance \tilde{R} (12)	Excess Empirical risk under KL privacy bound ε (Corollary 5.4)
LeCun	$1/m_{l-1}$	$\frac{om(L-1+\frac{d}{m})}{2^{L-1}d}$	$\tilde{O}(\frac{n}{m(L-1)})$	$\tilde{O}\left(\frac{1}{n^2} + \sqrt{\frac{1}{2^L d \varepsilon}} \left(1 + \frac{d}{m(L-1)}\right)\right)$
He	$2/m_{l-1}$	$\frac{om(L-1+\frac{d}{m})}{d}$	$\tilde{O}(\frac{n}{2^L m(L-1)})$	$\tilde{O}\left(\frac{1}{n^2} + \sqrt{\frac{1}{2^L d \varepsilon}} \left(1 + \frac{d}{m(L-1)}\right)\right)$
NTK	$2/m_l, l < L$ $1/o, l = L$	$\frac{m(L-1)}{2} + o$	$\tilde{O}(\frac{n}{d \cdot 2^L \cdot (m(L-2)+1)})$	$\tilde{O}\left(\frac{1}{n^2} + \sqrt{\frac{1}{d \cdot 2^L \varepsilon}} \cdot \frac{m(L-1)+2}{m(L-2)+1}\right)$
Xavier	$\frac{2}{m_{l-1}+m_l}$	$\frac{L-1+\frac{d+o}{2m}}{2^{L-1}(1+\frac{d}{m})(1+\frac{o}{m})}$	-	-

- The Xavier initialization makes neural networks fall into non-lazy training regime [31], so we do not include the lazy training distance nor privacy-utility trade-off analysis here.

34 To answer this question, we would need new algorithmic framework and (or) new privacy analyses.
35 In this paper, we focus on analyzing privacy for the Langevin diffusion algorithm ¹. This is to avoid
36 artificially setting a sensitivity constraint on the gradient update and thus making the privacy bound
37 insensitive to the network overparameterization (as in DP-SGD analysis). Instead, we prove a KL
38 privacy bound for Langevin diffusion that scales with the expected gradient difference between the
39 training on any two worst-case neighboring datasets (Theorem 3.1). ² By proving precise upper
40 bounds on the expected ℓ_2 -norm of this gradient difference, we obtain KL privacy bounds for fully
41 connected neural network (Lemma 3.2) and its linearized variant (Corollary 4.2) that changes with
42 the network width, depth and per-layer variance for the initialization distribution. We summarized
43 the details of our KL privacy bounds in Table 1, and highlight our key observations below.

- 44 • Width always worsen privacy, under all the considered initialization distributions. Mean-
45 while, the interplay between network depth and privacy is much more complex and crucially
46 depends on what initialization distribution is used and how long the training time is.
- 47 • Specifically, when the initialization distribution has small per-layer variance (such as Le-
48 Cun and Xavier initialization), our KL privacy bound for training fully connected network
49 (with a small amount of time) and for training linearized network (with finite time) de-
50 cay exponentially with increasing depth, as long as the depth is large enough. To the
51 best of our knowledge, this is the first time that an improvement of privacy bound under
52 overparameterization is observed for randomized training algorithm.

53 To further understand how the privacy utility trade-off is affected by overparameterization, we also
54 analyze the excess empirical risk and excess population risk of training linearized network using
55 Langevin diffusion. Our risk bounds scale with the lazy training distance R (i.e., how close is the
56 initialization vector to an optimal solution for the empirical risk minimization problem), as well
57 as a constant B for expected gradient norm in Langevin diffusion. By analyzing these two terms
58 precisely under overparameterization, we prove that given any fixed KL privacy budget ε , our risk
59 bounds strictly improves with increasing depth and width for linearized network under LeCun and He
60 initialization. To our best knowledge, this is the first time that such a gain in privacy-utility trade-off

¹A key difference between this paper and existing privacy utility analysis of Langevin diffusion [21] is that we analyze in the absence of gradient clipping or Lipschitz assumption on loss function. Our results also readily extend to discretized noisy GD with constant step-size (as discussed in Appendix E).

²We focus on KL privacy loss because it is a more relaxed distinguishability notion than standard (ε, δ) -DP, and therefore could be upper bounded even without gradient clipping. Moreover, KL divergence enables upper bound for the advantage (relative success) of various inference attacks, as studied in recent works [32, 22].

61 due to overparameterization (increasing depth) is shown. Meanwhile, prior results only prove (nearly)
 62 dimension-independent privacy utility trade-off for such linear models in the literature [39, 26, 30].
 63 Our improvement demonstrates the unique benefits of our new KL privacy analysis in understanding
 64 the effect of overparameterization.

65 1.1 Related Works

66 **Overparameterization in DNNs and NTK.** Theoretical demonstration on the benefit of over-
 67 parameterization in DNNs occurs in global convergence [2, 17], generalization [3, 13]. Under
 68 proper initialization, the training dynamics of over-parameterized DNNs can be described by a kernel
 69 function, termed as neural tangent kernel (NTK) [25], which stimulates a series of analysis in DNNs.
 70 Accordingly, over-parameterization has been demonstrated to be beneficial/harmful to several topics
 71 in deep learning, e.g., robustness [12, 47], covariate shift [44]. However, the relationship between
 72 overparameterization and privacy (based on the differential privacy framework) remains largely an
 73 unsolved problem, as the training dynamics typically change [11] after adding new components in
 74 the privacy-preserving learning algorithm (such as DP-SGD [1]) to enforce privacy constraints.

75 **Membership inference privacy risk under overparameterization.** A recent line of works [42, 43]
 76 investigates how overparameterization affects the theoretical and empirical privacy in terms of
 77 membership inference advantage, and proves novel trade-off between privacy and generalization error.
 78 These are the closest works in the literature to our objective of investigating the interplay between
 79 privacy and overparameterization. However, Tan et al. [42, 43] focus on proving upper bounds for
 80 an average-case privacy risk defined by the advantage (relative success) of membership inference
 81 attack on models trained from randomly sampled training dataset from a population distribution. By
 82 contrast, our KL privacy bound is heavily based on the strongest adversary model in the differential
 83 privacy definition, and holds under an arbitrary *worst-case* pair of neighboring datasets that only
 84 differ in one record. Our setting for model (fully connected network) is also very different from that
 85 considered in Tan et al. [42, 43], thus requiring very different analysis tools.

86 **Differentially private learning in high dimension.** Standard results for private empirical risk
 87 minimization [6, 41] and private stochastic convex optimization [7, 9, 4] under ℓ_1 and ℓ_2 constraints
 88 suggest that there is an unavoidable factor d in the empirical risk and population risk that depends
 89 on the model dimension. However, for unconstrained optimization, it is possible to go across the
 90 dimension-dependency in proving risk bounds for certain class of problems (such as generalized
 91 linear model [39]). Recently, there is a growing line of works that prove dimension-independent
 92 excess risk bounds for differentially private learning, by utilizing the low-rank structure of data
 93 features [39] or gradient matrices [26, 30] in training. Several follow-up works [27, 10] further
 94 explore techniques to enforce the low-rank property (via random projection) and boost privacy utility
 95 trade-off. However, all the works focus on investigating a general high-dimensional problem for
 96 private learning, rather than separating the study for different network choices such as structure,
 97 width, depth and initializaiton. On the contrary, our study focus on the fully connected neural network
 98 and its linearized variant, which enables us to prove more precise privacy utility trade-off bounds for
 99 these particular networks under overparameterization.

100 2 Problem and Methodology

101 We consider the following standard multi-class supervised learning setting. Let $\mathcal{D} = (z_1, \dots, z_n)$ be
 102 a finite input dataset of size n , where each input data record $z_i = (x_i, y_i)$ contains a d -dimensional
 103 input feature vector $x_i \in \mathbb{R}^d$ and a label vector $y_i \in \mathcal{Y} = \{0, 1\}^o$ on o possible classes. The goal
 104 of learning is to learn a neural network output function $f_{\mathbf{W}}(\cdot) : \mathcal{X} \rightarrow \mathcal{Y}$ parameterized by \mathbf{W} that
 105 achieves high prediction performance on the training dataset \mathcal{D} . Formally, we consider the learning
 106 objective to be the empirical risk defined as follows.

$$\min_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathcal{D}) := \frac{1}{n} \sum_{i=1}^n \ell(f_{\mathbf{W}}(x_i); y_i), \quad (1)$$

107 where $\ell(f_{\mathbf{W}}(x_i); y_i)$ is a loss function that reflects the approximation quality of model predic-
 108 tion $f_{\mathbf{W}}(x_i)$ compared to the ground truth label y_i . For simplicity, we assume that the cross-

109 entropy loss is used, i.e., $\ell(\mathbf{f}_{\mathbf{W}}(\mathbf{x}); \mathbf{y}) = -\langle \mathbf{y}, \log \text{softmax}(\mathbf{f}_{\mathbf{W}}(\mathbf{x})) \rangle$ for multi-output network, and
 110 $\ell(\mathbf{f}_{\mathbf{W}}(\mathbf{x}); \mathbf{y}) = \log(1 + \exp(-\mathbf{y}\mathbf{f}_{\mathbf{W}}(\mathbf{x})))$ for single-output network.

111 **Network.** We focus on the *multi-output, fully connected, deep neural network (DNN) with ReLU*
 112 *activation* with depth L (i.e., $L - 1$ hidden layers). Denote the width of each hidden layer with
 113 m_1, \dots, m_{L-1} . The output function $\mathbf{f}_{\mathbf{W}}(\mathbf{x}) := \mathbf{h}_L(\mathbf{x})$ is defined recursively as follows.

$$\mathbf{h}_0(\mathbf{x}) = \mathbf{x}; \quad \mathbf{h}_l(\mathbf{x}) = \phi(\mathbf{W}_l \mathbf{x}) \text{ for } l = 1, \dots, L - 1; \quad \mathbf{h}_L(\mathbf{x}) = \mathbf{W}_L \mathbf{h}_{L-1}(\mathbf{x}) \quad (2)$$

114 where $\mathbf{h}_l(\mathbf{x})$ denotes the output after activation at the l -th layer, the parameter matrices of each layer of
 115 the neural network satisfy $\mathbf{W}_1 \in \mathbb{R}^{m_1 \times d}$, $\mathbf{W}_l \in \mathbb{R}^{m_l \times m_{l-1}}$, $l = 2, \dots, L - 1$ and $\mathbf{W}_L \in \mathbb{R}^{o \times m_{L-1}}$.
 116 The model parameter $\mathbf{W} := (\text{Vec}(\mathbf{W}_1), \dots, \text{Vec}(\mathbf{W}_L)) \in \mathbb{R}^{m_1 \cdot d + m_2 \cdot m_1 + \dots + o \cdot m_{L-1}}$ consists of
 117 concatenation of vectorizations for parameters of all the layers. For consistency, we also denote
 118 $m_0 = d$ and $m_L = o$.

119 We also analyze the following **linearized network**, which is used in prior works [28, 2, 34] as an
 120 important tool to (approximately and qualitatively) analyze the training dynamics of deep neural
 121 networks. More formally, the linearized network $\mathbf{f}_{\mathbf{W}}^{\text{lin},0}(\mathbf{x})$ is a first-order Taylor expansion of the
 122 fully connected ReLU network at initialization parameter $\mathbf{W}_0^{\text{lin}}$, as follows.

$$\mathbf{f}_{\mathbf{W}}^{\text{lin},0}(\mathbf{x}) \equiv \mathbf{f}_{\mathbf{W}_0^{\text{lin}}}(\mathbf{x}) + \left. \frac{\partial \mathbf{f}_{\mathbf{W}}(\mathbf{x})}{\partial \mathbf{W}} \right|_{\mathbf{W}=\mathbf{W}_0^{\text{lin}}} (\mathbf{W} - \mathbf{W}_0^{\text{lin}}), \quad (3)$$

123 where $\mathbf{f}_{\mathbf{W}_0^{\text{lin}}}(\mathbf{x})$ is the output function of the fully connected ReLU network (2) at initialization $\mathbf{W}_0^{\text{lin}}$.
 124 We denote $\mathcal{L}_0^{\text{lin}}(\mathbf{W}; \mathcal{D}) = \sum_{i=1}^n \ell \left(\mathbf{f}_{\mathbf{W}_0^{\text{lin}}}(\mathbf{x}_i) + \left. \frac{\partial \mathbf{f}_{\mathbf{W}}(\mathbf{x})}{\partial \mathbf{W}} \right|_{\mathbf{W}=\mathbf{W}_0^{\text{lin}}} (\mathbf{W} - \mathbf{W}_0^{\text{lin}}); \mathbf{y}_i \right)$ as the empir-
 125 ical loss function for training linearized network, by plugging (3) into (1).

126 **Langevin Diffusion.** In terms of optimization algorithm, we focus on the *Langevin diffusion*
 127 algorithm [29] with per-dimensional noise variance σ^2 . Note that we aim to *avoid gradient clipping*
 128 while still proving KL privacy bounds. After initializing the model parameters \mathbf{W}_0 at time zero, the
 129 model parameters \mathbf{W}_t at subsequent time t evolves as the below stochastic differential equation.

$$d\mathbf{W}_t = -\nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}) dt + \sqrt{2\sigma^2} d\mathbf{B}_t. \quad (4)$$

130 **Initialization Distribution.** The initialization of parameters \mathbf{W}_0 crucially affects the convergence of
 131 Langevin diffusion, as observed in prior literatures [46, 20, 19]. Moreover, when the network function
 132 depends on the initialization parameters (as in linearized network (3)), the stationary distribution of
 133 Langevin diffusion also depends on the initialization distribution (as discussed in Section 5). In this
 134 work, we investigate the following general class of Gaussian initialization distribution with (possibly
 135 depth-dependent) variance for the parameters in each layer. For any layer $l = 1, \dots, L$, we have that

$$[\mathbf{W}^l]_{ij} \sim \mathcal{N}(0, \beta_l), \text{ for } (i, j) \in [m_l] \times [m_{l-1}] \quad (5)$$

136 where $\beta_1, \dots, \beta_L > 0$ are the per-layer variance for Gaussian initialization. By choosing different
 137 variance, we recover many common initialization schemes in the literature, as summarized in Table 1.

138 2.1 Our objective and methodology

139 We aim to understand the relation between privacy, utility and over-parameterization (depth and width)
 140 for the Langevin diffusion algorithm (under different initialization distributions). To understand
 141 privacy, we prove a KL divergence upper bound for running Langevin diffusion on any two *worst-case*
 142 neighboring datasets \mathcal{D} and \mathcal{D}' of size n that only differ in one record, denoted as $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}$ and
 143 $(\mathbf{x}, \mathbf{y}') \in \mathcal{D}'$. For brevity, in later sections, we denote \mathbf{W}_t (with distribution p_t) and \mathbf{W}'_t (with
 144 distribution p'_t) as the trained model parameters after running Langevin diffusion (4) for time T on \mathcal{D}
 145 and \mathcal{D}' respectively. We make the following assumptions for privacy analysis in this paper.

146 **Assumption 2.1** (Bounded Data). We assume that all \mathbf{x} in the data domain is bounded s.t. $\|\mathbf{x}\|_2 \leq 1$.

147 To understand utility (under a given KL privacy budget), we aim to prove upper bounds for excess
 148 empirical risk and excess population risk given an arbitrarily fixed KL divergence privacy budget ε .
 149 Finally, we also investigate how trade-off between our KL privacy bound and risk bounds is affected
 150 by the network width and depth. For utility analysis, we additionally make the following fair and
 151 attainable assumptions on data and network regularity.

152 **Assumption 2.2** (Data and network regularity [33, Assumption 2.1]). For any training data $\mathbf{x}_i \in \mathcal{D}$,
 153 it satisfies that $\|\mathbf{x}_i\|_2 = 1$. Moreover, $\mathbf{x}_i \in \mathcal{D}$ are i.i.d. samples from a data distribution P_x that
 154 satisfies $\int \|\mathbf{x}\|_2^2 dP_x(\mathbf{x}) = 1$. We also assume that the network only has single output.

155 Note that Assumption 2.2 is only required for utility analysis, and is not need for our privacy bound.

156 3 KL Privacy for Training Fully Connected ReLU Neural Networks

157 In this section, we perform the composition-based privacy analysis in KL divergence for Langevin
 158 Diffusion on deep ReLU neural networks, under Gaussian initialization distribution specified by
 159 Eq. (5). More specifically, we prove upper bound for the KL divergence between distribution of output
 160 model parameters when running Langevin diffusion on an arbitrary pair of neighboring datasets \mathcal{D}
 161 and \mathcal{D}' .

162 Our first key insight is that by the joint convexity of KL divergence, it is possible to prove composition-
 163 based KL privacy bound under more relaxed condition regarding the sensitivity of gradient computa-
 164 tion (i.e., *without* gradient clipping).

165 **Theorem 3.1** (KL composition under possibly unbounded gradient difference). *The KL divergence*
 166 *between running Langevin diffusion (4) for DNN (2) on neighboring datasets \mathcal{D} and \mathcal{D}' satisfies*

$$KL(\mathbf{W}_T, \mathbf{W}'_T) \leq \frac{1}{2\sigma^2} \int_0^T \mathbb{E} \left[\|\nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}')\|_2^2 \right] dt. \quad (6)$$

167 **Proof sketch.** We compute the partial derivative of KL divergence with regard to time t , and then
 168 compute integral to bound the KL divergence. During computing the limit of differentiation, we
 169 upper bound KL divergence at time $t + \eta$ for small enough step-size with the divergence on path
 170 $[t, t + \eta]$. Then we use Girsanov's theorem to compute the KL divergence between the path of coupled
 171 Langevin diffusion processes. The complete proof is in Appendix B.1.

172 Theorem 3.1 is an extension of the standard additivity [45] of KL divergence (also known as chain
 173 rule [40]) for a finite sequence of distributions to continuous time processes with (possibly) unbounded
 174 drift difference. The key extension is that Theorem 3.1 does not require bounded sensitivity between
 175 Langevin Diffusion on neighboring datasets. Instead, it only requires finite second-order moment of
 176 drift difference (in the ℓ_2 -norm sense) between neighboring datasets $\mathcal{D}, \mathcal{D}'$. By using this extended
 177 KL composition Theorem 3.1, we prove KL privacy bound for running Langevin diffusion algorithm
 178 (without gradient clipping) on deep neural networks, by tracking the upper bound for ℓ_2 norm of the
 179 gradient difference throughout training (under mild assumptions) as follows.

180 **Lemma 3.2** (Drift Difference in Noisy Training). *Let M_T be the subspace spanned by gradients*
 181 *$\{\nabla \ell(f_{\mathbf{W}_t}(\mathbf{x}_i; \mathbf{y}_i) : (\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}, t \in [0, T]\}$ on each training data record throughout Langevin*
 182 *diffusion $(\mathbf{W}_t)_{t \in [0, T]}$. Denote $\|\cdot\|_{M_T}$ as the ℓ_2 norm of the projection of the input vector onto linear*
 183 *space M_T . Suppose that $\exists c, \beta > 0$ s.t. for any \mathbf{W}, \mathbf{W}' and \mathbf{x}, \mathbf{y} we have $\|\nabla \ell(f_{\mathbf{W}}(\mathbf{x}; \mathbf{y})) -$
 184 $\nabla \ell(f_{\mathbf{W}'}(\mathbf{x}; \mathbf{y}))\|_2 < \max\{c, \beta \|\mathbf{W} - \mathbf{W}'\|_{M_T}\}$. Then over the randomness of the Brownian motion
 185 \mathbf{B}_t and initialization distribution (5) in Langevin diffusion $(\mathbf{W}_t)_{t \in [0, T]}$, it satisfies that*

$$\int_0^T \mathbb{E}_{p_t} \left[\|\nabla \mathcal{L}(\mathbf{W}; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}; \mathcal{D}')\|_2^2 \right] dt \leq 2 \cdot T \cdot \underbrace{\mathbb{E}_{p_0} \left[\|\nabla \mathcal{L}(\mathbf{W}; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}; \mathcal{D}')\|_2^2 \right]}_{\text{gradient difference at initialization}} \\
+ 2 \left(\frac{e^{(2+\beta^2)T} - (2 + \beta^2)T}{2 + \beta^2} \right) \cdot \underbrace{\left(\mathbb{E}_{p_0} \left[\|\nabla \mathcal{L}(\mathbf{W}; \mathcal{D})\|_2^2 \right] + \sigma^2 \text{rank}(M_T) + c^2 \right)}_{\text{gradient difference fluctuation during training}} + \underbrace{2c^2 \cdot T}_{\text{non-smoothness cost}}. \quad (7)$$

186 **Remark 3.3.** The assumption $\|\nabla \ell(f_{\mathbf{W}}(\mathbf{x}; \mathbf{y})) - \nabla \ell(f_{\mathbf{W}'}(\mathbf{x}; \mathbf{y}))\|_2 < \max\{c, \beta \|\mathbf{W} - \mathbf{W}'\|_{M_T}\}$ is
 187 similar to smoothness condition for the loss function, but is more relaxed as it allows non-smoothness
 188 at places where the gradient is bounded. Therefore, the assumption holds under ReLU activation.

189 **Remark 3.4** (Gradient difference at initialization). The first term and in our upper bound linearly
 190 scales with the difference between gradients on neighboring datasets \mathcal{D} and \mathcal{D}' at initialization. Under
 191 different initializations, this gradient difference exhibits different dependency on the network depth
 192 and width, as we will prove theoretically in Lemma 4.1.

193 *Remark 3.5* (Gradient difference fluctuation during training). The second term in our upper bound is
 194 to bound the change of gradient difference norm during training, and is therefore proportional to the
 195 the rank of a subspace M_T spanned by gradients of all training data. Intuitively, this fluctuation is
 196 because Langevin diffusion adds per-dimensional noise with variance σ^2 , thus perturbing the training
 197 parameters away from the initialization at a rate of $O(\sigma\sqrt{\text{rank}(M_T)})$ in expected ℓ_2 distance.

198 ***Growth of KL privacy bound with regard to training time T .*** The first term in the gradient difference
 199 bound Lemma 3.2 grows linearly with the training time T , while the second term grows exponentially
 200 with regard to T . Consequently, for learning tasks that requires a long training time to converge, the
 201 second term will become the dominating term and the KL privacy bound suffers from exponential
 202 growth with regard to the training time. Nevertheless, if the total amount of required training time
 203 (for convergence) is small enough e.g. $T \leq \frac{1}{2(2+\beta^2)}$, then we have that $\frac{e^{(2+\beta^2)T} - (2+\beta^2)T}{2+\beta^2} < T$ and
 204 therefore the second term in the gradient difference upper bound accumulates at a lower than linear
 205 rate with increasing training time.

206 ***Dependence of KL privacy bound on network overparameterization.*** Under a fixed training time T
 207 and noise scale σ^2 , Lemma 3.2 predicts that the KL divergence upper bound in Theorem 3.1 is depen-
 208 dent on the gradient difference and gradient norm at initialization, and the rank of gradient subspace
 209 $\text{rank}(M_T)$ throughout training. We now discuss the how these two terms change under increasing
 210 width and depth, and whether there are possibilities to improve them under overparameterization.

- 211 1. The gradient norm at initialization crucially depend on how the per-layer variance in the
 212 Gaussian initialization distribution scales with the network width and depth. Therefore, it is
 213 possible to improve the KL privacy bound by using initialization distributions that enable
 214 smaller gradient difference at initialization, as we will theoretically show in Section 4.
- 215 2. Regarding the rank of gradient subspace $\text{rank}(M_T)$: when the gradients along the training
 216 trajectory span the whole optimization space, $\text{rank}(M)$ would equal the dimension of the
 217 learning problem. Consequently, the gradient fluctuation upper bound (and thus the KL
 218 privacy bound) worsens with increasing number of model parameters (overparameterization)
 219 in the worst-case. However, if the gradients are low-dimensional [39, 26, 37] or sparse [30],
 220 it is possible that $\text{rank}(M_T)$ will be dimension-independent and thus enable better bound for
 221 gradient fluctuation (and KL privacy bound). We leave this as an interesting open problem.

222 4 KL privacy bound for Linearized Network under overparameterization

223 In this section, we restrict ourselves to the training of linearized networks as described in (3), and
 224 investigate the interplay between KL privacy and overparameterization (increasing width and depth).
 225 The analysis of DNN via linearization is a commonly used technique in both theory and practice.
 226 In theory, DNN can work in the lazy training regime [16] (also called linear regime), under which
 227 the linearized network well approximates the training dynamics for deep neural network [28] and
 228 has been well studied by NTK. In practice, linearized network can still achieve decent performance,
 229 which provides a good justification of linearized networks. [37, 34]. We hope our analysis for
 230 linearized network serve as an initial attempt that would open a door to theoretically understanding
 231 the relationship between overparameterization and privacy.

232 To derive a composition-based KL privacy bound for training linearized network, we apply Theo-
 233 rem 3.1 which requires an upper bound for the norm of gradient difference between the training
 234 processes on neighboring datasets \mathcal{D} and \mathcal{D}' at any time t . Note that the empirical risk function
 235 for training linearized models enjoys convexity, and therefore requires a relatively short amount of
 236 training time for convergence. Therefore intuitively, the gradient difference between neighboring
 237 datasets does not change a lot during training, thus allowing us to prove tighter upper bound for the
 238 gradient difference norm for linearized networks (than Lemma 3.2).

239 In the following lemma, we prove that for linearized network, the gradient difference throughout
 240 training has a uniform upper bound that only depends on the network width, depth and initialization.

241 **Lemma 4.1** (Gradient Difference throughout training linearized network). *Under Assumption 2.1,*
 242 *taking over the randomness of the random initialization and the Brownian motion in Langevin*

243 diffusion, for any $t \in [0, T]$, it satisfies that

$$\mathbb{E} [\|\nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}) - \mathcal{L}(\mathbf{W}_t; \mathcal{D}')\|^2] \leq \frac{4B}{n^2}, \quad (8)$$

244 where n is the training dataset size, and B is a constant that only depends on the network width,
245 depth and initialization distribution as follows.

$$B := o \left(\prod_{i=1}^{L-1} \frac{\beta_i m_i}{2} \right) \sum_{l=1}^L \frac{\beta_L}{\beta_l}, \quad (9)$$

246 where o is the number of output classes, $\{m_i\}_{i=1}^L$ are the per-layer network widths, and $\{\beta_i\}_{i=1}^L$ are
247 the variances of Gaussian initialization at each layer.

248 Lemma 4.1 provides an precise analytical upper bound for the gradient difference during training
249 linearized network, by tracking the gradient distribution under fully connected feed-forward ReLU
250 network with Gaussian weight matrices. The full proof is in Appendix C.1 and is heavily based on
251 similar techniques for computing the gradient distribution in the NTK literature [2, 47]. By plugging
252 Eq. (8) into Theorem 3.1, we have the following KL privacy bound for training linearized network.

253 **Corollary 4.2** (KL privacy bound for training linearized network). *Under Assumption 2.1 and neural*
254 *networks (3) initialized by Gaussian distribution with per-layer variance $\{\beta_i\}_{i=1}^L$, running Langevin*
255 *diffusion for linearized network with time T on neighboring datasets satisfies that*

$$KL(\mathbf{W}_t^{lin} \| \mathbf{W}_T^{lin}) \leq \frac{2BT}{n^2 \sigma^2}. \quad (10)$$

256 where B is the constant that specifies the gradient norm upper bound, given by (9).

257 **Overparameterization affects privacy differently under different initialization.** Corollary 4.2 and
258 Lemma 4.1 suggest that the effect of network overparameterization on KL privacy bound crucially
259 relies on how the per-layer Gaussian initialization variance β_i is scaled with the per-layer network
260 width m_i and depth L . We summarize our KL privacy bound for linearized network under different
261 width, depth and initialization schemes in Table 1, and elaborate the comparison below.

262 **(1) LeCun initialization** uses small, width-independent variance for initializing the first layer $\beta_1 = \frac{1}{d}$
263 (where d is the number of input features), and width-dependent variance $\beta_2 = \dots = \beta_L = \frac{1}{m}$ for
264 initializing all the subsequent layers. Therefore, the second term $\sum_{l=1}^L \frac{\beta_L}{\beta_l}$ in the constant B (9)
265 increases linearly with the width m and depth L . However, due to $\frac{m_l \cdot \beta_l}{2} < 1$ for all $l = 2, \dots, L$,
266 the first product term $\prod_{l=1}^{L-1} \frac{\beta_l m_l}{2}$ in constant B decays with the increasing depth. Therefore, by
267 combining the two terms, we prove that the KL privacy bound worsens with increasing width, but
268 improves with increasing depth (as long as the depth is large enough). Similarly, under Xavier
269 initialization $\beta_l = \frac{2}{m_{l-1} + m_l}$, we prove that the KL privacy bound (especially the constant B (9))
270 improves with increasing depth as long as the depth is large enough.

271 **(2) NTK and He initializations** user large per-layer variance $\beta_l = \begin{cases} \frac{2}{m_l} & l = 1, \dots, L-1 \\ \frac{1}{\sigma} & l = L \end{cases}$ (for

272 NTK) and $\beta_l = \frac{2}{m_{l-1}}$ (for He). Consequently, the gradient difference under NTK or He initialization
273 is significantly larger than that under LeCun initialization. Specifically, the gradient norm constant
274 B (9) grows linearly with the width m and the depth L under He and NTK initializations, thus
275 suggesting a worsening of KL privacy bound under increasing width and depth.

276 5 Utility guarantees for Training Linearized Network

277 Our privacy analysis suggest that training linearized network under certain initialization schemes
278 (such as LeCun initialization) enable significantly better privacy bounds under overparameterization
279 by increasing depth. In this section, we further prove utility bounds for Langevin diffusion under
280 initialization schemes, and investigate the effect of overparameterization on the privacy utility trade-
281 off. In other words, we aim to understand whether there are any utility degradation for training
282 linearized network when using the more privacy-preserving initialization schemes.

283 **Convergence of training linearized network.** We now prove convergence of excess empirical risk in
 284 training linearized network via Langevin diffusion. This is well-studied problem in the literature for
 285 noisy gradient descent. We extend the convergence theorem to continuous-time Langevin diffusion
 286 below, and investigate factors that affect the convergence under overparameterization.

287 **Proposition 5.1** (Extension of [36, Theorem 2] and [39, Theorem 3.1]). *Let \mathbf{W}_0^{lin} be a randomly*
 288 *initialized parameter vector by (5). Let the empirical NTK feature mapping matrix for dataset*
 289 *training \mathcal{D} at initialization be $M_0 = (\nabla \mathbf{f}_{\mathbf{W}_0^{lin}}(\mathbf{x}_1) \cdots \nabla \mathbf{f}_{\mathbf{W}_0^{lin}}(\mathbf{x}_n))$. Let $\mathcal{L}_0^{lin}(\mathbf{W}; \mathcal{D})$ be the*
 290 *empirical loss for linearized network (3) expanded at initialization vector \mathbf{W}_0^{lin} . Then running*
 291 *Langevin diffusion (4) under empirical loss $\mathcal{L}_0^{lin}(\mathbf{W}; \mathcal{D})$ and initialization \mathbf{W}_0^{lin} for time T satisfies*
 292 *the following excess empirical risk bound*

$$\mathbb{E}[\mathcal{L}(\mathbf{W}_T^{lin})] - \mathbb{E}[\mathcal{L}(\mathbf{W}_0^*; \mathcal{D})] \leq \frac{2R}{T} + \frac{1}{2}\sigma^2\mathbb{E}[\text{rank}(M_0)] \left(1 + \log \frac{2BT^2}{R}\right) \quad (11)$$

293 where \mathbf{W}_0^* is an (exact or approximate) solution for the ERM problem on $\mathcal{L}_0^{lin}(\mathbf{W}; \mathcal{D})$, and $R =$
 294 $\mathbb{E}[\|\mathbf{W}_0^{lin} - \mathbf{W}_0^*\|_{M_0}^2]$ is the expected gap between initialization parameters \mathbf{W}_0 and solution \mathbf{W}_0^* .

295 **Remark 5.2.** The excess empirical risk bound Proposition 5.1 is smaller if data is low-rank, e.g.,
 296 image data, then $\text{rank}(M_0)$ is small. This is consistent with the prior dimension-independent private
 297 learning literature [26, 27, 30] and show benefit of low-dimensional gradients on private learning.

298 Proposition 5.1 highlights that the excess empirical risk scales with the expected gap R between
 299 initialization and optima (which we refer as the lazy training distance), the rank of the gradient
 300 subspace $\text{rank}(M_0)$, and the constant B that specifies upper bound for expected gradient norm
 301 during training. Specifically, the smaller is the lazy training distance R , the better is the excess risk
 302 bound for Langevin diffusion given fixed training time T and noise variance σ^2 . We have discussed
 303 how overparameterization affects the the gradient norm constant B and the gradient subspace rank
 304 $\text{rank}(M_0)$ in Section 3. Therefore, we only still need to investigate how the lazy training distance R
 305 changes with the network width, depth and initialization, as follows.

306 **Lazy training distance R decreases with increasing depth.** It is widely observed in the literature [16,
 307 48, 31] that under appropriate choices of initializations, gradient descent on fully connected neural
 308 network falls under a lazy training regime. That is, with high probability, there exists a (nearly)
 309 optimal solution for the ERM problem that is close to the initialization parameters in terms of l_2
 310 distance. Moreover, this lazy training distance R is closely related to the smallest eigenvalue of the
 311 NTK matrix. In the following proposition, we compute a near optimal solution via the pseudo inverse
 312 of the NTK matrix, and prove that it has small distance to the initialization parameters via existing
 313 lower bounds for the smallest eigenvalue of the NTK matrix [33].

314 **Proposition 5.3** (Bounding lazy training distance via smallest eigenvalue of the NTK matrix). *Under*
 315 *the data and network regularity Assumption 2.1, if the width $m_1 = \cdots = m_{L-1} = \Omega(n)$ is*
 316 *sufficiently large, then there exists an optimal solution $\mathbf{W}_0^{\frac{1}{n^2}}$ that satisfies $\mathcal{L}_0^{lin}(\mathbf{W}_0^{\frac{1}{n^2}}) \leq \frac{1}{n^2}$ and*
 317 *satisfies*

$$\tilde{R} = \mathbb{E}[\|\mathbf{W}_0^{\frac{1}{n}} - \mathbf{W}_0\|_2^2] \leq \begin{cases} \tilde{O}\left(\frac{n}{d \cdot 2^L \cdot (m(L-2)+1)}\right) & \text{for NTK initialization} \\ \tilde{O}\left(\frac{n}{2^L m(L-1)}\right) & \text{for He initialization} \\ \tilde{O}\left(\frac{n}{m(L-1)}\right) & \text{for LeCun initialization} \end{cases} \quad (12)$$

318

319 We refer to \tilde{R} as the approximate lazy training distance because $\mathbf{W}_0^{\frac{1}{n^2}}$ is only a nearly optimal
 320 solution for the ERM problem. Proposition 5.3 shows that this approximate lazy training distance
 321 improves with overparameterization (width and depth) under LeCun, He and NTK initializations.

322 **Privacy Excess empirical risk Tradeoff for training linearized network via Langevin diffusion.** We
 323 now use the approximate lazy training distance \tilde{R} to prove empirical risk bound, and combine it with
 324 our KL privacy bound Section 4 to show the privacy utility trade-off under overparameterization.

325 **Corollary 5.4** (Privacy utility trade-off for last iterate). *Assume that the data and network regularity*
 326 *Assumption 2.2 holds. Assume that all the conditions and definition for constants in Proposition 5.1*

327 holds. Then by setting $\sigma^2 = \frac{2BT}{\varepsilon n^2}$ and $T = \sqrt{\frac{2\varepsilon n \tilde{R}}{B}}$, we have that running Langevin diffusion for
 328 time T satisfies bound KL divergence ε , and has excess empirical risk upper bounded by

$$\mathbb{E}[\mathcal{L}(\mathbf{W}_T^{lin})] \leq \mathcal{O}\left(\frac{1}{n^2} + \sqrt{\frac{B\tilde{R}}{\varepsilon n}} \log(\varepsilon n)\right) \quad (13)$$

329 where B is the gradient norm constant Eq. (9), and \tilde{R} is the approximate lazy training distance in
 330 Eq. (12). A summary of B and \tilde{R} under different initializations is in Table 1.

331 Corollary 5.4 suggests that the excess empirical risk worsens in the presence of a stronger privacy
 332 constraint, i.e., under a small privacy budget ε , thus contributing to a trade-off between privacy and
 333 utility. However, the excess empirical risk also scales with constants such as the approximate lazy
 334 training distance \tilde{R} and the gradient norm constant B . These constants depend on network width,
 335 depth and initialization distributions, and therefore we prove privacy utility trade-offs for training
 336 linearized network that changes with overparameterization, as summarized in Table 1.

337 We would like to highlight that our privacy utility trade-off bound under LeCun and He initialization
 338 strictly improves with increasing width and depth as long as the data and network satisfy regularity
 339 Assumption 2.2 and the network width is large enough. To our best knowledge, this is the first time
 340 that a strictly improving privacy utility trade-off under overparameterization is shown in literature.
 341 This shows the benefits of our precise KL privacy analysis under overparameterization.

342 **Extension of privacy utility results to excess population risk.** Our privacy utility trade-off results
 343 can be generalized to to excess population risk, by additionally bounding the generalization error
 344 with standard stability-based arguments [35, 24, 8]. We elaborate the details below.

345 **Proposition 5.5** (Extended excess population risk bound for training linearized network). *Denote*
 346 $R_0(\mathbf{W}) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \in \text{pop}}[\ell(f_{\mathbf{W}_0}(\mathbf{x}) + \frac{\partial f_{\mathbf{W}_0}(\mathbf{x})}{\partial \mathbf{W}_0}(\mathbf{W} - \mathbf{W}_0); \mathbf{y})]$ *as the population risk of linearized*
 347 *network expanded at initialization vector \mathbf{W}_0 over population data distribution pop . Then under the*
 348 *conditions of Corollary 5.4, we have that for any dataset \mathcal{D} of size n , the following excess population*
 349 *risk upper bound holds.*

$$\mathbb{E}[R_0(\mathbf{W}_T)] - \mathbb{E}[\mathcal{L}(\mathbf{W}_{pop,0}^*; \mathcal{D})] \leq \mathcal{O}\left(\frac{1}{n^2} + \sqrt{\frac{B\tilde{R}}{\varepsilon n}} (\log(\varepsilon n) + \varepsilon)\right) \quad (14)$$

350 where the expectation is over the randomness of sampling the training dataset $\mathcal{D} \sim \text{pop}^n$ from the
 351 data population and the random coins for the Langevin diffusion training algorithm, and $\mathbf{W}_{pop,0}^* =$
 352 $\text{argmin}_{\mathbf{W}} R_0(\mathbf{W})$ is the optimal solution for the population risk minimization problem.

353 Proposition 5.5 shows that the excess population risk bound is almost the same as excess empirical
 354 risk, except that there is an additional generalization term that scales with the privacy budget ε .
 355 Intuitively, this is because the generalization error is proportional to the stability of model prediction
 356 function under different training dataset, which is smaller when the KL privacy loss ε is small.

357 6 Conclusion

358 We prove new KL privacy bound for training fully connected ReLU network (and its linearized variant)
 359 using the Langevin diffusion algorithm, and investigate how privacy is affected by the network width,
 360 depth and initialization. Our results suggest that there is a complex interplay between privacy
 361 and overparameterization (width and depth) that crucially relies on what initialization distribution
 362 is used and the how much the gradient fluctuates during training. To this end, we show that for
 363 training a linearized variant of fully connected network with finite time, it is possible to prove a KL
 364 privacy bound that improves with depth, as long as the initialization distribution is set appropriately
 365 (such as LeCun). We also study the excess empirical and population risk bounds for linearized
 366 network, and prove that the privacy-utility trade-off similarly improves as depth increases under
 367 LeCun initialization. This shows the gain of our new privacy analysis for capturing the effect of
 368 overparameterization. We leave it as an important open problem as to whether our privacy utility
 369 trade-off results for linearized network could be generalized to deep neural networks.

References

- 370
- 371 [1] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar,
372 and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC*
373 *conference on computer and communications security*, pages 308–318, 2016.
- 374 [2] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via
375 over-parameterization. In *International Conference on Machine Learning*, pages 242–252.
376 PMLR, 2019.
- 377 [3] Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of op-
378 timization and generalization for overparameterized two-layer neural networks. In *International*
379 *Conference on Machine Learning*, pages 322–332. PMLR, 2019.
- 380 [4] Hilal Asi, Vitaly Feldman, Tomer Koren, and Kunal Talwar. Private stochastic convex optimiza-
381 tion: Optimal rates in ℓ_1 geometry. In *International Conference on Machine Learning*, pages
382 393–403. PMLR, 2021.
- 383 [5] Borja Balle, Giovanni Cherubin, and Jamie Hayes. Reconstructing training data with informed
384 adversaries. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 1138–1156. IEEE,
385 2022.
- 386 [6] Raef Bassily, Adam Smith, and Abhradeep Thakurta. Private empirical risk minimization:
387 Efficient algorithms and tight error bounds. In *2014 IEEE 55th annual symposium on foundations*
388 *of computer science*, pages 464–473. IEEE, 2014.
- 389 [7] Raef Bassily, Vitaly Feldman, Kunal Talwar, and Abhradeep Guha Thakurta. Private stochastic
390 convex optimization with optimal rates. *Advances in neural information processing systems*, 32,
391 2019.
- 392 [8] Raef Bassily, Vitaly Feldman, Cristóbal Guzmán, and Kunal Talwar. Stability of stochastic
393 gradient descent on nonsmooth convex losses. *Advances in Neural Information Processing*
394 *Systems*, 33:4381–4391, 2020.
- 395 [9] Raef Bassily, Cristóbal Guzmán, and Michael Menart. Differentially private stochastic opti-
396 mization: New results in convex and non-convex settings. *Advances in Neural Information*
397 *Processing Systems*, 34:9317–9329, 2021.
- 398 [10] Raef Bassily, Mehryar Mohri, and Ananda Theertha Suresh. Differentially private learning with
399 margin guarantees. *arXiv preprint arXiv:2204.10376*, 2022.
- 400 [11] Zhiqi Bu, Hua Wang, and Qi Long. On the convergence and calibration of deep learning with
401 differential privacy. *arXiv preprint arXiv:2106.07830*, 2021.
- 402 [12] Sébastien Bubeck and Mark Sellke. A universal law of robustness via isoperimetry. *Journal of*
403 *the ACM*, 70(2):1–18, 2023.
- 404 [13] Yuan Cao and Quanquan Gu. Generalization bounds of stochastic gradient descent for wide
405 and deep neural networks. *Advances in neural information processing systems*, 32, 2019.
- 406 [14] Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine
407 Lee, Adam Roberts, Tom B Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training
408 data from large language models. In *USENIX Security Symposium*, volume 6, 2021.
- 409 [15] Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru R Zhang. Sampling is as
410 easy as learning the score: theory for diffusion models with minimal data assumptions. *arXiv*
411 *preprint arXiv:2209.11215*, 2022.
- 412 [16] Lenaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable program-
413 ming. *Advances in neural information processing systems*, 32, 2019.
- 414 [17] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global
415 minima of deep neural networks. In *International conference on machine learning*, pages
416 1675–1685. PMLR, 2019.

- 417 [18] Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Found.*
418 *Trends Theor. Comput. Sci.*, 9(3-4):211–407, 2014.
- 419 [19] Yoav Freund, Yi-An Ma, and Tong Zhang. When is the convergence time of langevin algorithms
420 dimension independent? a composite optimization viewpoint. *Journal of Machine Learning*
421 *Research*, 23(214):1–32, 2022.
- 422 [20] Arun Ganesh and Kunal Talwar. Faster differentially private samplers via rényi divergence
423 analysis of discretized langevin mcmc. *Advances in Neural Information Processing Systems*,
424 33:7222–7233, 2020.
- 425 [21] Arun Ganesh, Abhradeep Thakurta, and Jalaj Upadhyay. Langevin diffusion: An almost univer-
426 sal algorithm for private euclidean (convex) optimization. *arXiv preprint arXiv:2204.01585*,
427 2022.
- 428 [22] Chuan Guo, Alexandre Sablayrolles, and Maziar Sanjabi. Analyzing privacy leakage in machine
429 learning via multiple hypothesis testing: A lesson from fano. *arXiv preprint arXiv:2210.13662*,
430 2022.
- 431 [23] Niv Haim, Gal Vardi, Gilad Yehudai, Ohad Shamir, and Michal Irani. Reconstructing training
432 data from trained neural networks. *arXiv preprint arXiv:2206.07758*, 2022.
- 433 [24] Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of
434 stochastic gradient descent. In *International conference on machine learning*, pages 1225–1234.
435 PMLR, 2016.
- 436 [25] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and
437 generalization in neural networks. *Advances in neural information processing systems*, 31,
438 2018.
- 439 [26] Peter Kairouz, Monica Ribero Diaz, Keith Rush, and Abhradeep Thakurta. (nearly) dimension
440 independent private erm with adagrad rates via publicly estimated subspaces. In *Conference on*
441 *Learning Theory*, pages 2717–2746. PMLR, 2021.
- 442 [27] Shiva Prasad Kasiviswanathan. Sgd with low-dimensional gradients with applications to private
443 and distributed learning. In *Uncertainty in Artificial Intelligence*, pages 1905–1915. PMLR,
444 2021.
- 445 [28] Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-
446 Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models
447 under gradient descent. *Advances in neural information processing systems*, 32, 2019.
- 448 [29] DS Lemons, A Gythiel, and Paul Langevin’s. Sur la théorie du mouvement brownien [on the
449 theory of brownian motion]. *CR Acad. Sci.(Paris)*, 146:530–533, 1908.
- 450 [30] Xuechen Li, Daogao Liu, Tatsunori B Hashimoto, Huseyin A Inan, Janardhan Kulkarni, Yin-Tat
451 Lee, and Abhradeep Guha Thakurta. When does differentially private learning not suffer in high
452 dimensions? *Advances in Neural Information Processing Systems*, 35:28616–28630, 2022.
- 453 [31] Tao Luo, Zhi-Qin John Xu, Zheng Ma, and Yaoyu Zhang. Phase diagram for two-layer relu
454 neural networks at infinite-width limit. *The Journal of Machine Learning Research*, 22(1):
455 3327–3373, 2021.
- 456 [32] Saeed Mahloujifar, Alexandre Sablayrolles, Graham Cormode, and Somesh Jha. Optimal
457 membership inference bounds for adaptive composition of sampled gaussian mechanisms.
458 *arXiv preprint arXiv:2204.06106*, 2022.
- 459 [33] Quynh Nguyen, Marco Mondelli, and Guido F Montufar. Tight bounds on the smallest
460 eigenvalue of the neural tangent kernel for deep relu networks. In *International Conference on*
461 *Machine Learning*, pages 8119–8129. PMLR, 2021.
- 462 [34] Guillermo Ortiz-Jiménez, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. What can
463 linearized neural networks actually say about generalization? *Advances in Neural Information*
464 *Processing Systems*, 34:8998–9010, 2021.

- 465 [35] Shai Shalev-Shwartz, Ohad Shamir, Nathan Srebro, and Karthik Sridharan. Stochastic convex
466 optimization. In *COLT*, volume 2, page 5, 2009.
- 467 [36] Ohad Shamir and Tong Zhang. Stochastic gradient descent for non-smooth optimization:
468 Convergence results and optimal averaging schemes. In *International conference on machine
469 learning*, pages 71–79. PMLR, 2013.
- 470 [37] Vaishaal Shankar, Alex Fang, Wenshuo Guo, Sara Fridovich-Keil, Jonathan Ragan-Kelley,
471 Ludwig Schmidt, and Benjamin Recht. Neural kernels without tangents. In *International
472 Conference on Machine Learning*, pages 8614–8623. PMLR, 2020.
- 473 [38] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference
474 attacks against machine learning models. In *2017 IEEE symposium on security and privacy
475 (SP)*, pages 3–18. IEEE, 2017.
- 476 [39] Shuang Song, Thomas Steinke, Om Thakkar, and Abhradeep Thakurta. Evading the curse
477 of dimensionality in unconstrained private glms. In *International Conference on Artificial
478 Intelligence and Statistics*, pages 2638–2646. PMLR, 2021.
- 479 [40] Thomas Steinke. Adaptive data analysis. 2016.
- 480 [41] Kunal Talwar, Abhradeep Thakurta, and Li Zhang. Private empirical risk minimization beyond
481 the worst case: The effect of the constraint set geometry. *arXiv preprint arXiv:1411.5417*, 2014.
- 482 [42] Jasper Tan, Blake Mason, Hamid Javadi, and Richard Baraniuk. Parameters or privacy: A
483 provable tradeoff between overparameterization and membership inference. *Advances in Neural
484 Information Processing Systems*, 35:17488–17500, 2022.
- 485 [43] Jasper Tan, Daniel LeJeune, Blake Mason, Hamid Javadi, and Richard G Baraniuk. A blessing
486 of dimensionality in membership inference through regularization. In *International Conference
487 on Artificial Intelligence and Statistics*, pages 10968–10993. PMLR, 2023.
- 488 [44] Nilesh Tripuraneni, Ben Adlam, and Jeffrey Pennington. Covariate shift in high-dimensional
489 random feature regression. *arXiv preprint arXiv:2111.08234*, 2021.
- 490 [45] Tim Van Erven and Peter Harremoës. Rényi divergence and kullback-leibler divergence. *IEEE
491 Transactions on Information Theory*, 60(7):3797–3820, 2014.
- 492 [46] Santosh Vempala and Andre Wibisono. Rapid convergence of the unadjusted langevin algorithm:
493 Isoperimetry suffices. *Advances in neural information processing systems*, 32, 2019.
- 494 [47] Zhenyu Zhu, Fanghui Liu, Grigorios G Chrysos, and Volkan Cevher. Robustness in deep
495 learning: The good (width), the bad (depth), and the ugly (initialization). *arXiv preprint
496 arXiv:2209.07263*, 2022.
- 497 [48] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes over-
498 parameterized deep relu networks. *Machine learning*, 109:467–492, 2020.

499 **Contents**

500 **1 Introduction** **1**

501 1.1 Related Works 3

502 **2 Problem and Methodology** **3**

503 2.1 Our objective and methodology 4

504 **3 KL Privacy for Training Fully Connected ReLU Neural Networks** **5**

505 **4 KL privacy bound for Linearized Network under overparameterization** **6**

506 **5 Utility guarantees for Training Linearized Network** **7**

507 **6 Conclusion** **9**

508 **A Symbols and definitions** **13**

509 **B Deferred proofs for Section 3** **14**

510 B.1 Deferred proofs for Theorem 3.1 14

511 B.2 Deferred proofs for Lemma 3.2 15

512 **C Deferred proofs for Section 4** **17**

513 C.1 Bounding the gradient norm at initialization 17

514 C.2 Deferred proof for Lemma 4.1 19

515 **D Deferred proofs for Section 5** **20**

516 D.1 Excess empirical risk for training linearized network (average iterate) 20

517 D.2 Deferred proof for Proposition 5.1 21

518 D.3 Deferred proof for Proposition 5.3 23

519 D.4 Deferred proof for Corollary 5.4 24

520 D.5 Deferred proofs for Proposition 5.5 25

521 **E Discussion on extending our results to Noisy GD with constant step-size** **26**

522 **A Symbols and definitions**

523 Vectorization $\text{Vec}(\cdot)$ denotes the transformation that takes an input matrix $\mathbf{A} = (a_{ij})_{i \in [r], j \in [c]} \in$
524 $\mathbb{R}^{r \times c}$ (with r rows and c columns) and outputs a rc -dimensional column vector: $\text{Vec}(\mathbf{A}) =$
525 $(a_{1,1}, \dots, a_{r,1}, a_{1,2}, \dots, a_{r,2}, \dots, a_{1,c}, \dots, a_{r,c})^\top$.

526 Distribution p_t and p'_t : we denote p_t as the distribution of model parameters after running Langevin
527 diffusion on dataset \mathcal{D} with time t , and similarly denote p'_t as the distribution of model parameters
528 after running Langevin diffusion on dataset \mathcal{D}' with time t .

529 Softmax function: $\text{softmax}(\mathbf{y}) = \frac{e^{\mathbf{y}^{[j]}}}{\sum_{j=1}^o e^{\mathbf{y}^{[j]}}}$ where o is the number of output classes.

530 Neighboring datasets D and D' : two dataset with the same number of data records that differ in one
531 record. We also denote the differing records as $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}$ and $(\mathbf{x}', \mathbf{y}') \in \mathcal{D}'$.

532 o : number of output classes for the neural network.

533 B Deferred proofs for Section 3

534 B.1 Deferred proofs for Theorem 3.1

535 To prove the new composition theorem, we will use the Girsanov's Theorem. Here we follow the
536 presentation of [15, Theorem 6].

537 **Theorem B.1** (Implication of Girsanov's theorem [15, Theorem 6]). *Let $(\tilde{X}_t)_{t \in [0, \eta]}$ and $(X_t)_{t \in [0, \eta]}$
538 be two continuous-time processes over \mathbb{R}^r . Let P_T be the probability measure that corresponds to the
539 trajectory of $(\tilde{X}_t)_{t \in [0, \eta]}$, and let Q_T be the probability measure that corresponds to the trajectory of
540 $(X_t)_{t \in [0, \eta]}$. Suppose that the process $(\tilde{X}_t)_{t \in [0, \eta]}$ follows*

$$d\tilde{X}_t = \tilde{b}_t dt + \sigma_t d\tilde{B}_t,$$

541 where \tilde{B} is a Brownian motion, and the process $(X_t)_{t \in [0, \eta]}$ follows

$$dX_t = b_t dt + \sigma_t dB_t,$$

542 where B is a Brownian motion. We assume that for each $t > 0$, σ_t is a $r \times r$ symmetric positive
543 definite matrix. Then, provided that Novikov's condition holds,

$$\mathbb{E}_{Q_T} \exp \left(\frac{1}{2} \int_0^\eta \|\sigma_t^{-1}(\tilde{b}_t - b_t)\|_2^2 dt \right) < \infty, \quad (15)$$

544 we have that

$$\frac{dP_T}{dQ_T} = \exp \left(\int_0^\eta \sigma_t^{-1}(\tilde{b}_t - b_t) dB_t - \frac{1}{2} \int_0^\eta \|\sigma_t^{-1}(\tilde{b}_t - b_t)\|_2^2 dt \right).$$

545 Now we apply Girsanov's theorem on the coupled Langevin diffusion processes on neighboring
546 datasets, and obtain the following new composition theorem for KL divergence in the context of
547 privacy.

548 **Theorem 3.1** (KL composition under possibly unbounded gradient difference). The KL divergence
549 between running Langevin diffusion (4) for DNN (2) on neighboring datasets \mathcal{D} and \mathcal{D}' satisfies

$$KL(\mathbf{W}_T, \mathbf{W}'_T) \leq \frac{1}{2\sigma^2} \int_0^T \mathbb{E} \left[\|\nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}')\|_2^2 \right] dt. \quad (16)$$

550 *Proof.* Denote p_t as the distribution of model parameters after running Langevin diffusion on dataset
551 \mathcal{D} with time t , and similarly denote p'_t as the distribution of model parameters after running Langevin
552 diffusion on dataset \mathcal{D}' with time t . Then by definition,

$$\begin{aligned} \frac{\partial KL(p_t, p'_t)}{\partial t} &= \lim_{\eta \rightarrow 0} \frac{KL(p_{t+\eta}, p'_{t+\eta}) - KL(p_t, p'_t)}{\eta} \\ &\leq \lim_{\eta \rightarrow 0} \frac{KL(p_{t,t+\eta}, p'_{t,t+\eta}) - KL(p_t, p'_t)}{\eta}, \end{aligned} \quad (17)$$

553 where the last inequality is by the data processing inequality for KL divergence [45, Theorem 9]
554 (with the data processing operation given by $(\mathbf{W}_t, \mathbf{W}_{t+\eta}) \rightarrow \mathbf{W}_t$). Now we compute the term
555 $KL(p_{t,t+\eta}, p'_{t,t+\eta})$ as follows.

$$\begin{aligned} KL(p_{t,t+\eta}, p'_{t,t+\eta}) &= \mathbb{E}_{p_{t,t+\eta}(\mathbf{w}_t, \mathbf{w}_{t+\eta})} \left[\log \left(\frac{p_{t+\eta|t}(\mathbf{W}_{t+\eta} | \mathbf{W}_t) p_t(\mathbf{W}_t)}{p'_{t+\eta|t}(\mathbf{W}_{t+\eta} | \mathbf{W}_t) p'_t(\mathbf{W}_t)} \right) \right] \\ &= \mathbb{E}_{p_{t,t+\eta}(\mathbf{w}_t, \mathbf{w}_{t+\eta})} \left[\log \left(\frac{p_{t+\eta|t}(\mathbf{W}_{t+\eta} | \mathbf{W}_t)}{p'_{t+\eta|t}(\mathbf{W}_{t+\eta} | \mathbf{W}_t)} \right) \right] + \mathbb{E}_{p_t(\mathbf{w}_t)} \left[\log \left(\frac{p_t(\mathbf{W}_t)}{p'_t(\mathbf{W}_t)} \right) \right] \\ &= \mathbb{E}_{p_t(\mathbf{w}_t)} \left[KL(p_{t+\eta|t}, p'_{t+\eta|t}) \right] + KL(p_t, p'_t) \end{aligned} \quad (18)$$

556 Now we want to apply Girsanov's theorem to the following Langevin diffusion processes
 557 $(\mathbf{W}_{t+s|t})_{s \in [0, \eta]}$ and $(\mathbf{W}'_{t+s|t})_{s \in [0, \eta]}$.

$$\begin{aligned} d\mathbf{W}_{t+s|t} &= -\nabla \mathcal{L}(\mathbf{W}_{t+s}; \mathcal{D}) dt + \sqrt{2\sigma^2} dB_s \\ d\mathbf{W}'_{t+s|t} &= -\nabla \mathcal{L}(\mathbf{W}'_{t+s}; \mathcal{D}') dt + \sqrt{2\sigma^2} dB_s \end{aligned}$$

558 where we have the boundary condition that $\mathbf{W}_{t|t} = \mathbf{W}'_{t|t}$ due to the conditioning at time t . Note that
 559 when η is small enough, we have that the Novikov's condition in Eq. (15) holds because the exponent
 560 inside integration $\frac{1}{2} \int_0^\eta \|\sigma_t^{-1}(\tilde{b}_t - b_t)\|_2^2 dt$ scales linearly with η and is small when η is small enough.
 561 Therefore, by applying Girsanov's theorem, we have that

$$\begin{aligned} KL(p_{t+\eta|t}, p'_{t+\eta|t}) &\leq KL(p_{t:t+\eta|t}, p'_{t:t+\eta|t}) \\ &= \mathbb{E}_{p_{t:t+\eta|t}} \left[\int_0^\eta \sigma^{-1}(\tilde{b}_s - b_s) dB_s - \frac{1}{2} \int_0^\eta \|\sigma^{-1}(\tilde{b}_s - b_s)\|_2^2 ds \right] \end{aligned}$$

562 where $\tilde{b}_s - b_s = -\nabla \mathcal{L}(\mathbf{W}_{t+s}; \mathcal{D}) + \nabla \mathcal{L}(\mathbf{W}_{t+s}; \mathcal{D}')$. By Itô integration with regard to standard
 563 Brownian motion, we have that

$$KL(p_{t+\eta|t}, p'_{t+\eta|t}) \leq \frac{1}{2\sigma^2} \mathbb{E}_{p_{t:t+\eta|t}} \left[\int_0^\eta \|\nabla \mathcal{L}(\mathbf{W}_{t+s}; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}_{t+s}; \mathcal{D}')\|_2^2 ds \right] \quad (19)$$

564 By plugging Eq. (19) into Eq. (18), we have that

$$KL(p_{t,t+\eta}, p'_{t,t+\eta}) \leq \frac{1}{2\sigma^2} \mathbb{E}_{p_{t:t+\eta}} \left[\int_0^\eta \|\nabla \mathcal{L}(\mathbf{W}_{t+s}; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}_{t+s}; \mathcal{D}')\|_2^2 ds \right] + KL(p_t, p'_t) \quad (20)$$

565 By plugging Eq. (20) into Eq. (17), and by exchanging the order of expectation and integration, we
 566 have that

$$\begin{aligned} \frac{\partial KL(p_t, p'_t)}{\partial t} &\leq \frac{1}{2\sigma^2} \lim_{\eta \rightarrow 0} \frac{\int_0^\eta \mathbb{E}_{p_{t+s}} \left[\|\nabla \mathcal{L}(\mathbf{W}_{t+s}; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}_{t+s}; \mathcal{D}')\|_2^2 \right] ds}{\eta} \\ &= \frac{1}{2\sigma^2} \mathbb{E}_{p_t} \left[\|\nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}')\|_2^2 \right] \end{aligned} \quad (21)$$

567 Integrating Eq. (21) on $t \in [0, T]$ finishes the proof. \square

568 B.2 Deferred proofs for Lemma 3.2

569 **Lemma 3.2.** Let M_T be the subspace spanned by gradients $\{\nabla \ell(\mathbf{f}_{\mathbf{W}_t}(\mathbf{x}_i; \mathbf{y}_i) : (\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}, t \in [0, T]\}$
 570 on each training data record throughout Langevin diffusion $(\mathbf{W}_t)_{t \in [0, T]}$. Denote $\|\cdot\|_{M_T}$ as the
 571 ℓ_2 norm of the projection of the input vector onto linear space M_T . Suppose that $\exists c, \beta > 0$ s.t. for
 572 any \mathbf{W}, \mathbf{W}' and \mathbf{x}, \mathbf{y} we have $\|\nabla \ell(\mathbf{f}_{\mathbf{W}}(\mathbf{x}; \mathbf{y})) - \nabla \ell(\mathbf{f}_{\mathbf{W}'}(\mathbf{x}; \mathbf{y}))\|_2 < \max\{c, \beta \|\mathbf{W} - \mathbf{W}'\|_{M_T}\}$.
 573 Then over the randomness of the Brownian motion \mathbf{B}_t and initialization distribution (5) in Langevin
 574 diffusion $(\mathbf{W}_t)_{t \in [0, T]}$, it satisfies that

$$\begin{aligned} \int_0^T \mathbb{E}_{p_t} \left[\|\nabla \mathcal{L}(\mathbf{W}; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}; \mathcal{D}')\|_2^2 \right] dt &\leq 2 \cdot T \cdot \underbrace{\mathbb{E}_{p_0} \left[\|\nabla \mathcal{L}(\mathbf{W}; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}; \mathcal{D}')\|_2^2 \right]}_{\text{gradient difference at initialization}} \\ &+ 2 \left(\frac{e^{(2+\beta^2)T} - (2+\beta^2)T}{2+\beta^2} \right) \cdot \underbrace{\left(\mathbb{E}_{p_0} \left[\|\nabla \mathcal{L}(\mathbf{W}; \mathcal{D})\|_2^2 \right] + \sigma^2 \text{rank}(M_T) + c^2 \right)}_{\text{gradient difference fluctuation during training}} + \underbrace{2c^2 \cdot T}_{\text{non-smoothness cost}}. \end{aligned} \quad (22)$$

575 *Proof.* By definition of the neighboring datasets \mathcal{D} and \mathcal{D}' , we have that

$$\|\nabla \mathcal{L}(\mathbf{W}; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}; \mathcal{D}')\|_2^2 = \|\ell(\mathbf{f}_{\mathbf{W}}(\mathbf{x}; \mathbf{y})) - \nabla \ell(\mathbf{f}_{\mathbf{W}}(\mathbf{x}'; \mathbf{y}'))\|_2^2 \quad (23)$$

576 where (\mathbf{x}, \mathbf{y}) and $(\mathbf{x}', \mathbf{y}')$ are the differing records between two datasets. By the assumption that
 577 $\|\nabla\ell(f_{\mathbf{W}}(\mathbf{x}); \mathbf{y}) - \nabla\ell(f_{\mathbf{W}'}(\mathbf{x}); \mathbf{y})\|_2 < \max\{c, \beta\|\mathbf{W} - \mathbf{W}'\|_{M_T}\}$, and by the Cauchy-Schwarz
 578 inequality, we have that

$$\begin{aligned} \|\nabla\ell(f_{\mathbf{W}_t}(\mathbf{x}); \mathbf{y}) - \nabla\ell(f_{\mathbf{W}_t}(\mathbf{x}'); \mathbf{y}')\|_2^2 &\leq 2\|\nabla\ell(f_{\mathbf{W}_0}(\mathbf{x}); \mathbf{y}) - \nabla\ell(f_{\mathbf{W}_0}(\mathbf{x}'); \mathbf{y}')\|_2^2 \\ &\quad + 2\beta^2\|\mathbf{W}_t - \mathbf{W}_0\|_{M_T}^2 + 2c^2 \end{aligned} \quad (24)$$

579 The first term $\|\nabla\ell(f_{\mathbf{W}_0}(\mathbf{x}); \mathbf{y}) - \nabla\ell(f_{\mathbf{W}_0}(\mathbf{x}'); \mathbf{y}')\|_2^2$ is constant during training (as it only depends
 580 on the initialization). Therefore, we only need to bound the second term $\|\mathbf{W}_t - \mathbf{W}_0\|_{M_T}^2$. For brevity,
 581 we denote the term inside expectation as $d(\mathbf{W}) = \|\mathbf{W} - \mathbf{W}_0\|_{M_T}^2$. Then by definition we have that

$$\frac{\partial}{\partial t}\mathbb{E}_{p_t}[d(\mathbf{W})] = \lim_{\eta \rightarrow 0} \frac{\mathbb{E}_{p_{t+\eta}}[d(\mathbf{W})] - \mathbb{E}_{p_t}[d(\mathbf{W})]}{\eta}. \quad (25)$$

582 Denote Γ_s as the following random operator on model parameters θ .

$$\Gamma_s(\mathbf{W}) = \theta - s\nabla\mathcal{L}(\mathbf{W}; \mathcal{D}) + \sqrt{2\sigma^2}sZ$$

583 where $Z \sim \mathcal{N}(0, \mathbb{I})$. We first claim that the following equation holds.

$$\lim_{\eta \rightarrow 0} \frac{\mathbb{E}_{p_{t+\eta}}[d(\mathbf{W})] - \mathbb{E}_{p_t}[d(\Gamma_\eta(\mathbf{W}))]}{\eta} = 0 \quad (26)$$

584 This is by using Euler-Maruyama discretization method to approximate the solution \mathbf{W}_t of
 585 SDE Eq. (4). More specifically, the approximation error $\mathbb{E}_{p_{t+\eta}}[d(\mathbf{W})] - \mathbb{E}_{p_t}[d(\Gamma_\eta(\mathbf{W}))]$ is of
 586 size $O(r\eta^2)$, where r is the dimension of \mathbf{W} .

587 Therefore, by plugging Eq. (26) into Eq. (25), we have that

$$\frac{\partial}{\partial t}\mathbb{E}_{p_t}[d(\mathbf{W})] = \lim_{\eta \rightarrow 0} \frac{\mathbb{E}_{p_t}[d(\Gamma_\eta(\mathbf{W}))] - \mathbb{E}_{p_t}[d(\mathbf{W})]}{\eta}$$

588 Recall that $\nabla^2 d(\mathbf{W})$ exists almost everywhere with regard to $\mathbf{W} \sim p_t$. Therefore we could
 589 approximate the term $\mathbb{E}_{p_t}[d(\Gamma_\eta(\mathbf{W}); \mathcal{D}, \mathcal{D}')]$ via its second-order Taylor expansion at \mathbf{W} as follows.

$$\begin{aligned} \frac{\partial}{\partial t}\mathbb{E}_{p_t}[d(\mathbf{W})] &= \lim_{\eta \rightarrow 0} \frac{\mathbb{E}_{p_t}[\langle \nabla d(\mathbf{W}), -\eta\nabla\mathcal{L}(\mathbf{W}; \mathcal{D}) + \sqrt{2\sigma^2\eta}Z \rangle + \sigma^2\eta Z^\top \nabla^2 d(\mathbf{W})Z + o(\eta)]}{\eta} \\ &= -\mathbb{E}_{p_t}[\langle \nabla d(\mathbf{W}), \nabla\mathcal{L}(\mathbf{W}; \mathcal{D}) \rangle] + \sigma^2\mathbb{E}_{p_t}[\text{Tr}(\nabla^2 d(\mathbf{W}))] \end{aligned} \quad (27)$$

590 By plugging $d(\mathbf{W}) = \|\mathbf{W} - \mathbf{W}_0\|_{M_T}^2$ into the above equation, we have that

$$\frac{\partial}{\partial t}\mathbb{E}_{p_t}[\|\mathbf{W} - \mathbf{W}_0\|_{M_T}^2] \leq -2\mathbb{E}_{p_t}[\langle \mathbf{W} - \mathbf{W}_0, \nabla\mathcal{L}(\mathbf{W}; \mathcal{D}) \rangle] + \sigma^2\text{rank}(M_T) \quad (28)$$

$$\begin{aligned} &= -2\mathbb{E}_{p_t}[\langle \mathbf{W} - \mathbf{W}_0, \nabla\mathcal{L}(\mathbf{W}_0; \mathcal{D}) \rangle] + \sigma^2\text{rank}(M_T) \\ &\quad - 2\mathbb{E}_{p_t}[\langle \mathbf{W} - \mathbf{W}_0, \nabla\mathcal{L}(\mathbf{W}; \mathcal{D}) - \nabla\mathcal{L}(\mathbf{W}_0; \mathcal{D}) \rangle] \\ &\leq \mathbb{E}[\|\nabla\mathcal{L}(\mathbf{W}_0; \mathcal{D})\|_2^2] + \mathbb{E}_{p_t}[\|\mathbf{W} - \mathbf{W}_0\|_{M_T}^2] + \sigma^2\text{rank}(M_T) \end{aligned} \quad (29)$$

$$+ \mathbb{E}_{p_t}[\|\mathbf{W} - \mathbf{W}_0\|_{M_T}^2] + \mathbb{E}_{p_t}[\|\nabla\mathcal{L}(\mathbf{W}; \mathcal{D}) - \nabla\mathcal{L}(\mathbf{W}_0; \mathcal{D})\|_2^2] \quad (30)$$

591 where the last inequality is by using the Cauchy-schwarz inequality. By plugging the assumption
 592 that $\|\nabla\ell(f_{\mathbf{W}}(\mathbf{x}); \mathbf{y}) - \nabla\ell(f_{\mathbf{W}'}(\mathbf{x}); \mathbf{y})\|_2 < \max\{c, \beta\|\mathbf{W} - \mathbf{W}'\|_2\}$ into the above inequality, we
 593 have that

$$\frac{\partial}{\partial t}\mathbb{E}_{p_t}[\|\mathbf{W} - \mathbf{W}_0\|_{M_T}^2] \leq (2 + \beta^2)\mathbb{E}_{p_t}[\|\mathbf{W} - \mathbf{W}_0\|_{M_T}^2] + \mathbb{E}[\|\nabla\mathcal{L}(\mathbf{W}_0; \mathcal{D})\|_2^2] + c^2 + \sigma^2\text{rank}(M_T) \quad (31)$$

594 By solving the above ordinary differential inequality on $t \in [0, T]$, we have that

$$\mathbb{E}_{p_t}[\|\mathbf{W} - \mathbf{W}_0\|_{M_T}^2] \leq \frac{e^{(2+\beta^2)t} - 1}{\beta^2} (\mathbb{E}[\|\nabla\mathcal{L}(\mathbf{W}_0; \mathcal{D})\|_2^2] + \sigma^2\text{rank}(M_T) + c^2) \quad (32)$$

595 By plugging Eq. (32) into Eq. (24) and integrating over time $t \in [0, T]$, we have that

$$\int_0^T \mathbb{E}_{p_t} \left[\|\nabla \mathcal{L}(\mathbf{W}; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}; \mathcal{D}')\|_2^2 \right] dt \leq 2T \cdot \mathbb{E}_{p_0} \left[\|\nabla \mathcal{L}(\mathbf{W}; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}; \mathcal{D}')\|_2^2 \right] \quad (33)$$

$$+ 2 \left(\frac{e^{(2+\beta^2)T} - (2 + \beta^2)T}{2 + \beta^2} \right) \cdot (\mathbb{E}_{p_0} [\|\nabla \mathcal{L}(\mathbf{W}; \mathcal{D})\|_2^2] + \sigma^2 \text{rank}(M_T) + c^2) + 2c^2T. \quad (34)$$

596

□

597 C Deferred proofs for Section 4

598 C.1 Bounding the gradient norm at initialization

599 To bound the moment of ℓ_2 norm of the gradient $\frac{\partial f(\mathbf{x})}{\partial \mathbf{W}_i}$ of network output function, we need the
600 following (extended) lemmas from Zhu et al. [47].

601 **Lemma C.1** ([47, Lemma 1]). *Let $\mathbf{w} \sim \mathcal{N}(0, \sigma^2 \mathbb{I}_n)$, then for two fixed non-zero vectors $\mathbf{h}_1, \mathbf{h}_2 \in$
602 \mathbb{R}^n whose correlation is unknown, define two random variables $X = (\mathbf{w}^\top \mathbf{h}_1 1_{\{\mathbf{w}^\top \mathbf{h}_2 \geq 0\}})^2$ and
603 $Y = s(\mathbf{w}^\top \mathbf{h}_1)^2$, where $s \sim \text{Ber}(1, 1/2)$ follows a Bernoulli distribution with 1 trial and $\frac{1}{2}$ success
604 rate, and s and \mathbf{w} are independent random variables. Then X and Y have the same distribution.*

605 **Lemma C.2** (Extension of [47, Lemma 2]). *Given a fixed non-zero matrix $\mathbf{H}_1 \in \mathbb{R}^{p \times r}$ and a fixed
606 non-zero vector $\mathbf{h}_2 \in \mathbb{R}^p$ and, let $\mathbf{W} \in \mathbb{R}^{q \times p}$ be a random matrix with i.i.d. entries $W_{ij} \sim \mathcal{N}(0, \beta)$
607 and a matrix (or vector) $\mathbf{V} = \phi'(\mathbf{W}\mathbf{h}_2)\mathbf{W}\mathbf{H}_1 \in \mathbb{R}^{q \times r}$, then, we have $\mathbb{E} \left[\frac{\|\mathbf{V}\|_F^2}{\|\mathbf{H}_1\|_F^2} \right] = \frac{q\beta}{2}$.*

608 *Proof.* According to the definition of $\mathbf{V} = \phi'(\mathbf{W}\mathbf{h}_2)\mathbf{W}\mathbf{H}_1 \in \mathbb{R}^{q \times r}$, we have:

$$\|\mathbf{V}\|_F^2 = \sum_{i=1}^q \sum_{j=1}^r \left(D_{i,i} \langle \mathbf{W}^{[i]}, \mathbf{H}_1^{[j]} \rangle \right)^2,$$

609 where $D_{i,i} = 1_{\{\langle \mathbf{W}^{[i]}, \mathbf{h}_2 \rangle \geq 0\}}$, $\mathbf{W}^{[i]}$ is the i -th row of \mathbf{W} , and $\mathbf{H}_1^{[j]}$ is the j -th column vector of \mathbf{H}_1 .
610 Therefore by Lemma C.1, with i.i.d. Bernoulli random variable $\rho_1, \dots, \rho_q \sim \text{Ber}(1, 1/2)$, we have

$$\|\mathbf{V}\|_F^2 \stackrel{d}{=} \sum_{i=1}^q \sum_{j=1}^r \rho_i \langle \mathbf{W}^{[i]}, \mathbf{H}_1^{[j]} \rangle^2 = \sum_{i=1}^q \sum_{j=1}^r \rho_i \beta \|\mathbf{H}_1^{[j]}\|_2^2 \tilde{w}_{ij}^2.$$

611 where $\tilde{w}_{ij} = \langle \mathbf{W}^{[i]}, \mathbf{H}_1^{[j]} \rangle / \left(\sqrt{\beta \|\mathbf{H}_1^{[j]}\|_2^2} \right)$. By the fact that $\mathbf{W}^{[i]}$ has i.i.d. Gaussian entries, for
612 any fixed j , we have that $\tilde{w}_{ij} \sim \mathcal{N}(0, 1)$, $i = 1, \dots, q$ independently. Therefore, we have

$$\mathbb{E} [\|\mathbf{V}\|_F^2] = \sum_{i=1}^q \sum_{j=1}^r \mathbb{E}[\rho_i] \beta \|\mathbf{H}_1^{[j]}\|_2^2 \mathbb{E}[\tilde{w}_{ij}^2] = \frac{q\beta}{2} \mathbb{E}[\|\mathbf{H}_1\|_F^2].$$

613

□

614 Now, we are ready to prove output gradient expectation at random initialization as follows.

615 **Lemma C.3** (Output Gradient Expectation Bound at Random Initialization). *Fix any data record \mathbf{x} ,*
616 *then over the randomness of the initialization distributions for $\mathbf{W}_1, \dots, \mathbf{W}_L$, i.e., $\mathbf{W}_l \sim \mathcal{N}(0, \beta_l \mathbb{I})$*
617 *for $l = 1, \dots, L-1$, it satisfies that*

$$\mathbb{E}_{\mathbf{W}} \left[\left\| \frac{\partial f(\mathbf{x})}{\partial \text{Vec}(\mathbf{W})} \right\|_F^2 \right] = \|\mathbf{x}\|_{2^o}^2 \left(\prod_{i=1}^{L-1} \frac{\beta_i m_i}{2} \right) \sum_{l=1}^L \frac{\beta_L}{\beta_l}. \quad (35)$$

618 *Proof.* We use $\text{Vec}(\mathbf{W}_l)$ to denote the concatenation of all row vector of the parameter matrix \mathbf{W}_l .
619 By chain rule, for $l = 1, \dots, L-1$, we have that

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \text{Vec}(\mathbf{W}_l)} = \frac{\partial h_L(\mathbf{x})}{\partial h_{L-1}(\mathbf{x})} \left(\prod_{i=1}^{L-1-l} \frac{\partial h_{L-i}(\mathbf{x})}{\partial h_{L-1-i}(\mathbf{x})} \right) \frac{\partial h_l}{\text{Vec}(\mathbf{W}_l)} \quad (36)$$

$$= \mathbf{W}_L \left(\prod_{i=1}^{L-1-l} \sigma'_{L-i} \mathbf{W}_{L-i} \right) \sigma'_l \begin{pmatrix} h_{l-1}^\top & 0 & \cdots \\ \vdots & \vdots & \vdots \\ 0 & 0 & h_{l-1}^\top \end{pmatrix}_{m_l \times m_l m_{l-1}}. \quad (37)$$

620 Similarly, for the L -th layer, we have that

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \text{Vec}(\mathbf{W}_L)} = \begin{pmatrix} h_{L-1}^\top & 0 & \cdots \\ \vdots & \vdots & \vdots \\ 0 & 0 & h_{L-1}^\top \end{pmatrix}_{o \times o m_{L-1}}. \quad (38)$$

621 By properties of ReLU activation ϕ , we have $\phi'_{L-i} = \text{diag}[\text{sgn}(W_{L-i} h_{L-1-i})]$, where $\text{sgn}(x) =$
622 $\begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$ operates coordinate-wise with regard to the input matrix. Therefore, we have that for
623 $l = 1, \dots, L-1$

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \text{Vec}(\mathbf{W}_l)} = \mathbf{W}_L \left(\prod_{i=1}^{L-1-l} \text{diag}[\text{sgn}(W_{L-i} h_{L-1-i})] \mathbf{W}_{L-i} \right) \cdot \text{diag}[\text{sgn}(W_l h_{l-1})] \begin{pmatrix} h_{l-1}^\top & 0 & \cdots \\ \vdots & \vdots & \vdots \\ 0 & 0 & h_{l-1}^\top \end{pmatrix}_{m_l \times m_l m_{l-1}}.$$

624 For notational simplicity, we introduce the notation of $\mathbf{t}_l^{l'}$ for $l = 1, \dots, L-1$ and $l \leq l' < L$ as
625 follows.

$$\mathbf{t}_l^{l'} := \left(\prod_{i=L-l'}^{L-1-l} \text{diag}[\text{sgn}(W_{L-i} h_{L-1-i})] \mathbf{W}_{L-i} \right) \cdot \text{diag}[\text{sgn}(W_l h_{l-1})] \begin{pmatrix} h_{l-1}^\top & 0 & \cdots \\ \vdots & \vdots & \vdots \\ 0 & 0 & h_{l-1}^\top \end{pmatrix}_{m_l \times m_l m_{l-1}}.$$

626 Then by definition, we have that

$$\begin{aligned} \mathbb{E}_{\mathbf{W}} \left[\left\| \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \text{Vec}(\mathbf{W}_l)} \right\|_F^2 \right] &= \mathbb{E}_{\mathbf{W}} \left[\left\| \mathbf{W}_L \mathbf{t}_l^{L-1} \right\|_F^2 \right] = \mathbb{E}_{\mathbf{W}} \left[\frac{\left\| \mathbf{W}_L \mathbf{t}_l^{L-1} \right\|_F^2}{\left\| \mathbf{t}_l^{L-1} \right\|_F^2} \cdot \frac{\left\| \mathbf{t}_l^{L-1} \right\|_F^2}{\left\| \mathbf{t}_l^{L-2} \right\|_F^2} \cdots \frac{\left\| \mathbf{t}_l^{l+1} \right\|_F^2}{\left\| \mathbf{t}_l^l \right\|_F^2} \cdot \left\| \mathbf{t}_l^l \right\|_F^2 \right] \\ &= \mathbb{E}_{\mathbf{W}_1, \dots, \mathbf{W}_l} \left[\left\| \mathbf{t}_l^l \right\|_F^2 \cdot \mathbb{E}_{\mathbf{W}_{l+1}} \left[\frac{\left\| \mathbf{t}_l^{l+1} \right\|_F^2}{\left\| \mathbf{t}_l^l \right\|_F^2} \cdots \mathbb{E}_{\mathbf{W}_L} \left[\frac{\left\| \mathbf{W}_L \mathbf{t}_l^{L-1} \right\|_F^2}{\left\| \mathbf{t}_l^{L-1} \right\|_F^2} \right]} \right] \right]. \end{aligned}$$

627 By rotational invariance of Gaussian column vectors, we prove that for any possible value of \mathbf{t}_l^{L-1}
628 (which is completely determined by $\mathbf{W}_1, \dots, \mathbf{W}_{L-1}$ and \mathbf{x}), for any $l = 1, \dots, L-1$, we have that

$$\mathbb{E}_{\mathbf{W}_L} \left[\frac{\left\| \mathbf{W}_L \mathbf{t}_l^{L-1} \right\|_F^2}{\left\| \mathbf{t}_l^{L-1} \right\|_F^2} \right] = \mathbb{E}_{\mathbf{W}_L} \left[\frac{\left\| \mathbf{W}_L \mathbf{e}_1 \right\|_2^2}{\left\| \mathbf{e}_1 \right\|_2^2} \right] = \beta_L o. \quad (39)$$

629 By Lemma C.2, for any $l = 1, \dots, L-2$ and $l \leq l' \leq L-2$, we have that

$$\mathbb{E}_{\mathbf{W}_{l+1}} \left[\frac{\left\| \mathbf{t}_l^{l'+1} \right\|_F^2}{\left\| \mathbf{t}_l^{l'} \right\|_F^2} \right] = \frac{\beta_{l'+1}}{2} m_{l'+1}. \quad (40)$$

630 We now bound the last term $\mathbb{E}_{\mathbf{W}_1, \dots, \mathbf{W}_l} \left[\left\| \mathbf{t}_l^l \right\|_F^2 \right]$ for $l = 1, \dots, L-1$. By definition, we have that

$$\mathbb{E}_{\mathbf{W}_1, \dots, \mathbf{W}_l} \left[\left\| \mathbf{t}_l^l \right\|_F^2 \right] = \mathbb{E}_{\mathbf{W}_1, \dots, \mathbf{W}_l} \left[\sum_{i=1}^{m_l} \mathbf{1}_{\{\mathbf{W}_l^{(i)} h_{l-1} \leq 0\}} \cdot \left\| h_{l-1} \right\|_2^2 \right] = \frac{m_l}{2} \mathbb{E}_{\mathbf{W}_1, \dots, \mathbf{W}_{l-1}} \left[\left\| h_{l-1} \right\|_2^2 \right]. \quad (41)$$

631 To bound the term $\mathbb{E}_{\mathbf{W}_1, \dots, \mathbf{W}_l} [\|h_{l-1}\|_2^2]$, note that by Lemma C.2, we have that for any $l = 1, \dots, L-1$
 632 1

$$\mathbb{E}_{\mathbf{W}_{l-1}} \left[\frac{\|h_{l-1}(\mathbf{x})\|_2^2}{\|h_{l-2}(\mathbf{x})\|_2^2} \right] = \frac{\beta_{l-1}}{2} m_{l-1}. \quad (42)$$

633 Therefore, for any $l = 1, \dots, L$, we have that

$$\mathbb{E}_{\mathbf{W}_1, \dots, \mathbf{W}_{l-1}} [\|h_{l-1}(\mathbf{x})\|_2^2] = \mathbb{E}_{\mathbf{W}_1, \dots, \mathbf{W}_{l-1}} \left[\frac{\|h_{l-1}(\mathbf{x})\|_2^2}{\|h_{l-2}(\mathbf{x})\|_2^2} \dots \frac{\|h_1(\mathbf{x})\|_2^2}{\|\mathbf{x}\|_2^2} \right] \cdot \|\mathbf{x}\|_2^2 \quad (43)$$

$$= \left(\prod_{i=1}^{l-1} \frac{\beta_i}{2} m_i \right) \|\mathbf{x}\|_2^2. \quad (44)$$

634 By plugging (44) into (41), we have that

$$\mathbb{E}_{\mathbf{W}_1, \dots, \mathbf{W}_l} [\|t_l^l\|_2^2] = \frac{m_l}{2} \left(\prod_{i=1}^{l-1} \frac{\beta_i}{2} m_i \right) \|\mathbf{x}\|_2^2. \quad (45)$$

635 By combining (39), (40) and (45), we have for any $l = 1, \dots, L-1$

$$\begin{aligned} & \mathbb{E}_{\mathbf{W}} \left[\left\| \frac{\partial f(\mathbf{x})}{\partial \text{Vec}(\mathbf{W}_l)} \right\|_F^2 \right] \\ &= \frac{m_l}{2} \left(\prod_{i=1}^{l-1} \frac{\beta_i}{2} m_i \right) \cdot \left(\prod_{i=l+1}^{L-1} \frac{\beta_i}{2} m_i \right) \cdot \beta_L o \cdot \|\mathbf{x}\|_2^2 = \frac{\beta_L}{\beta_l} \|\mathbf{x}\|_2^2 o \left(\prod_{i=1}^{L-1} \frac{\beta_i m_i}{2} \right). \end{aligned}$$

636 On the other hand, by plugging Eq. (44) (under $\ell = L$) into Eq. (38), we have that

$$\mathbb{E}_{\mathbf{W}} \left[\left\| \frac{\partial f(\mathbf{x})}{\partial \text{Vec}(\mathbf{W}_L)} \right\|_2^2 \right] = o \left(\prod_{i=1}^{L-1} \frac{\beta_i}{2} m_i \right) \|\mathbf{x}\|_2^2$$

637 Therefore,

$$\mathbb{E}_{\mathbf{W}} \left[\left\| \frac{\partial f(\mathbf{x})}{\partial \text{Vec}(\mathbf{W})} \right\|_F^2 \right] = \sum_{l=1}^L \left\| \frac{\partial f(\mathbf{x})}{\partial \mathbf{W}_l} \right\|_F^2 = \|\mathbf{x}\|_2^2 o \left(\prod_{i=1}^{L-1} \frac{\beta_i m_i}{2} \right) \sum_{l=1}^L \frac{\beta_L}{\beta_l},$$

638 which suffices to prove Eq. (35). \square

639 C.2 Deferred proof for Lemma 4.1

640 Finally, we prove that the gradient difference between two training datasets under linearized network
 641 is bounded by a constant throughout training (which only depends on the network width, depth and
 642 initialization distribution).

643 **Lemma 4.1.** Under Assumption 2.1, taking over the randomness of the random initialization and the
 644 Brownian motion in Langevin diffusion, for any $t \in [0, T]$, it satisfies that

$$\mathbb{E} [\|\nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}')\|^2] \leq \frac{4B}{n^2}, \quad (46)$$

645 where n is the training dataset size, and B is a constant that only depends on the network width,
 646 depth and initialization distribution as follows.

$$B := o \left(\prod_{i=1}^{L-1} \frac{\beta_i m_i}{2} \right) \sum_{l=1}^L \frac{\beta_L}{\beta_l}, \quad (47)$$

647 where o is the number of output classes, $\{m_i\}_{i=1}^L$ are the per-layer network widths, and $\{\beta_i\}_{i=1}^L$ are
 648 the variances of Gaussian initialization at each layer.

649 *Proof.* Denote \mathbf{W} as the initialization parameters and denote \mathbf{W}_t^{lin} as the parameters for linearized
650 network after training time t . Then the gradient difference under linearized network and cross-entropy
651 loss function is as follows.

$$\begin{aligned} & \|\nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}')\|_F^2 \\ &= \left\| \frac{\nabla \mathbf{f}_{\mathbf{W}}(\mathbf{x})^\top (\text{softmax}(\mathbf{f}_{\mathbf{W}_t}(\mathbf{x})) - \mathbf{y})}{n} - \frac{\nabla \mathbf{f}_{\mathbf{W}}(\mathbf{x}')^\top (\text{softmax}(\mathbf{f}_{\mathbf{W}_t}(\mathbf{x}')) - \mathbf{y}')}{n} \right\|_F^2 \\ &\leq \frac{2}{n^2} (\|\nabla \mathbf{f}_{\mathbf{W}}(\mathbf{x})\|_F^2 + \|\nabla \mathbf{f}_{\mathbf{W}}(\mathbf{x}')\|_F^2). \end{aligned}$$

652 Plugging Lemma 4.1 into the above equation with data Assumption 2.1 suffice to prove the result. \square

653 D Deferred proofs for Section 5

654 D.1 Excess empirical risk for training linearized network (average iterate)

655 To prove empirical risk bound for the last iterate of training linearized network, we will first need to
656 prove the following intermediate result of excess empirical risk bound for average iterate.

657 **Lemma D.1** (Excess empirical risk for average iterate (Extension of [39, Theorem 3.1])). *Let \mathbf{W}_0^{lin}*
658 *be a randomly initialized parameter vector by (5). Let the empirical NTK feature matrix for dataset*
659 *training \mathcal{D} at initialization be $M_0 = (\nabla \mathbf{f}_{\mathbf{W}_0^{lin}}(\mathbf{x}_1) \cdots \nabla \mathbf{f}_{\mathbf{W}_0^{lin}}(\mathbf{x}_n))$. Let $\mathcal{L}_0^{lin}(\mathbf{W}; \mathcal{D})$ be the*
660 *empirical loss for linearized network (3) expanded at initialization vector \mathbf{W}_0^{lin} . Then running*
661 *Langevin diffusion (4) under empirical loss $\mathcal{L}_0^{lin}(\mathbf{W}; \mathcal{D})$ and initialization \mathbf{W}_0^{lin} for time T satisfies*
662 *the following excess empirical risk bound.*

$$\mathbb{E}[\mathcal{L}(\bar{\mathbf{W}}_T^{lin})] - \mathbb{E}[\mathcal{L}(\mathbf{W}_0^*; \mathcal{D})] \leq \frac{R}{2T} + \frac{1}{2}\sigma^2 \mathbb{E}[\text{rank}(M_0)]$$

663 where $\bar{\mathbf{W}}_T^{lin} = \frac{1}{T} \int \bar{\mathbf{W}}_t^{lin} dt$ is the average of all iterates, \mathbf{W}_0^* is an (exact or approximate) solution
664 for the ERM problem on $\mathcal{L}_0^{lin}(\mathbf{W}; \mathcal{D})$, and $R = \mathbb{E}[\|\mathbf{W}_0^{lin} - \mathbf{W}_0^*\|_{M_0}^2]$ is the expected gap between
665 initialization parameters \mathbf{W}_0 and solution \mathbf{W}_0^* .

666 *Proof.* Our proofs are heavily based on the idea in [39, Theorem 3.1] to work only in the parameter
667 space spanned by the input feature vectors. And our proof serves as an extension of their bound to
668 the continuous-time Langevin diffusion algorithm. We begin by using convexity of the empirical loss
669 function $\mathcal{L}^{lin}(\mathbf{W}; \mathcal{D})$ for linearized network to prove the following standard results

$$\mathcal{L}^{lin}(\bar{\mathbf{W}}_T^{lin}; \mathcal{D}) - \mathcal{L}^{lin}(\mathbf{W}_0^*; \mathcal{D}) \leq \langle \bar{\mathbf{W}}_T^{lin} - \mathbf{W}_0^*, \nabla \mathcal{L}^{lin}(\bar{\mathbf{W}}_T^{lin}; \mathcal{D}) \rangle \quad (48)$$

670 Denote $M_0 = (\nabla \mathbf{f}_{\mathbf{W}_0^{lin}}(\mathbf{x}_1) \cdots \nabla \mathbf{f}_{\mathbf{W}_0^{lin}}(\mathbf{x}_n))$. By computing the gradient under cross entropy
671 loss and linearized network, we have $\nabla \mathcal{L}^{lin}(\bar{\mathbf{W}}_T^{lin}; \mathcal{D})$ lies in the column space of M_0 . Denote Π_{M_0}
672 as the projection operator to the column space of M_0 , then (48) can be rewritten as

$$\mathcal{L}^{lin}(\bar{\mathbf{W}}_T^{lin}; \mathcal{D}) - \mathcal{L}^{lin}(\mathbf{W}_0^*; \mathcal{D}) \leq \langle \Pi_{M_0}(\bar{\mathbf{W}}_T^{lin} - \mathbf{W}_0^*), \nabla \mathcal{L}^{lin}(\bar{\mathbf{W}}_T^{lin}; \mathcal{D}) \rangle. \quad (49)$$

673 By taking expectation over training randomness and initialization distribution, we have

$$\mathbb{E}[\mathcal{L}(\bar{\mathbf{W}}_T^{lin})] - \mathbb{E}[\mathcal{L}(\mathbf{W}_0^*; \mathcal{D})] \leq \frac{1}{T} \int_0^T \mathbb{E} [\langle \Pi_{M_0}(\mathbf{W}_t^{lin} - \mathbf{W}_0^*), \nabla \mathcal{L}^{lin}(\mathbf{W}_t^{lin}; \mathcal{D}) \rangle] dt \quad (50)$$

674 We now rewrite $\mathbb{E} [\langle \Pi_{M_0}(\bar{\mathbf{W}}_t^{lin} - \mathbf{W}_0^*), \nabla \mathcal{L}^{lin}(\mathbf{W}_t^{lin}; \mathcal{D}) \rangle]$ by computing $\frac{\partial}{\partial t} \mathbb{E}[\|\mathbf{W}_t^{lin} - \mathbf{W}_0^*\|_{M_0}^2]$,
675 where $\|\mathbf{W}_t^{lin} - \mathbf{W}_0^*\|_{M_0}^2 = \Pi_{M_0}(\mathbf{W}_t^{lin} - \mathbf{W}_0^*)^\top \Pi_{M_0}(\mathbf{W}_t^{lin} - \mathbf{W}_0^*)$. By applying (27), we have

$$\frac{\partial}{\partial t} \mathbb{E}[\|\mathbf{W}_t^{lin} - \mathbf{W}_0^*\|_{M_0}^2] = -2\mathbb{E}[\langle \Pi_{M_0}(\mathbf{W}_t^{lin} - \mathbf{W}_0^*), \nabla \mathcal{L}^{lin}(\mathbf{W}_t^{lin}; \mathcal{D}) \rangle] + \sigma^2 \mathbb{E}[\text{rank}(M_0)] \quad (51)$$

676 Therefore by plugging (51) into (50), we have that

$$\mathbb{E}[\mathcal{L}^{lin}(\mathbf{W}; \mathcal{D}) - \mathcal{L}^{lin}(\mathbf{W}_0^*; \mathcal{D})] \leq -\frac{1}{2T} \int_0^T \frac{\partial}{\partial t} \mathbb{E}[\|\mathbf{W}_t - \mathbf{W}_0^*\|_{M_0}^2] dt + \frac{1}{2}\sigma^2 \mathbb{E}[\text{rank}(M_0)] \quad (52)$$

$$\leq \frac{1}{2T} \mathbb{E}[\|\mathbf{W}_0^{lin} - \mathbf{W}_0^*\|_{M_0}^2] + \frac{1}{2}\sigma^2 \mathbb{E}[\text{rank}(M_0)] \quad (53)$$

677 \square

678 **D.2 Deferred proof for Proposition 5.1**

679 We are now ready to prove the last iterate excess empirical risk bound for training linearized network.

680 **Proposition 5.1** (Excess empirical risk for training linearized network (last iterate)). Let \mathbf{W}_0^{lin} be
681 a randomly initialized parameter vector by (5). Let the empirical NTK feature mapping matrix for
682 dataset training \mathcal{D} at initialization be $M_0 = (\nabla \mathbf{f}_{\mathbf{W}_0^{lin}}(\mathbf{x}_1) \cdots \nabla \mathbf{f}_{\mathbf{W}_0^{lin}}(\mathbf{x}_n))$. Let $\mathcal{L}_0^{lin}(\mathbf{W}; \mathcal{D})$
683 be the empirical loss for linearized network (3) expanded at initialization vector \mathbf{W}_0^{lin} . Then running
684 Langevin diffusion (4) under empirical loss $\mathcal{L}_0^{lin}(\mathbf{W}; \mathcal{D})$ and initialization \mathbf{W}_0^{lin} for time T satisfies
685 the following excess empirical risk bound

$$\mathbb{E}[\mathcal{L}(\mathbf{W}_T^{lin})] - \mathbb{E}[\mathcal{L}(\mathbf{W}_0^*; \mathcal{D})] \leq \frac{2R}{T} + \frac{1}{2}\sigma^2\mathbb{E}[\text{rank}(M_0)] \left(1 + \log \frac{2BT^2}{R}\right) \quad (54)$$

686 where \mathbf{W}_0^* is an (exact or approximate) solution for the ERM problem on $\mathcal{L}_0^{lin}(\mathbf{W}; \mathcal{D})$, and $R =$
687 $\mathbb{E}[\|\mathbf{W}_0^{lin} - \mathbf{W}_0^*\|_{M_0}^2]$ is the expected gap between initialization parameters \mathbf{W}_0 and solution \mathbf{W}_0^* .

688 *Proof.* We first define the following potential function.

$$\Phi(t) = \frac{1}{T-t} \int_t^T \mathbb{E}[\mathcal{L}(\mathbf{W}_\tau^{lin}; \mathcal{D})] d\tau \quad (55)$$

689 By definition we have that the boundary values are $\Phi(0) = \frac{1}{T} \int_0^T \mathbb{E}[\mathcal{L}(\mathbf{W}_\tau^{lin}; \mathcal{D})] d\tau$ and $\Phi(T) =$
690 $\lim_{t \rightarrow T} \frac{1}{T-t} \int_t^T \mathbb{E}[\mathcal{L}(\mathbf{W}_\tau^{lin}; \mathcal{D})] d\tau = \mathbb{E}[\mathcal{L}(\mathbf{W}_T^{lin}; \mathcal{D})]$. Since we have proved upper bound for $\Phi(0)$
691 in the excess empirical risk bound for the average iterate Appendix D.1, to analyze $\Phi(T)$ the
692 loss difference between last iterate and average iterate, we only need to prove upper bound for
693 $\Phi(T) - \Phi(0)$.

694 By definition, we compute the partial derivative of the function $\Phi(t)$ with regard to time t as follows.

$$\begin{aligned} \frac{\partial}{\partial t} \Phi(t) &= \frac{1}{(T-t)^2} \int_t^T \mathbb{E}[\mathcal{L}(\mathbf{W}_\tau^{lin}; \mathcal{D})] d\tau - \frac{1}{T-t} \mathbb{E}[\mathcal{L}(\mathbf{W}_t^{lin}; \mathcal{D})] \\ &= \frac{1}{(T-t)^2} \int_t^T \mathbb{E}[\mathcal{L}(\mathbf{W}_\tau^{lin}; \mathcal{D}) - \mathcal{L}(\mathbf{W}_t^{lin}; \mathcal{D})] d\tau \\ &\leq \frac{1}{(T-t)^2} \int_t^T \mathbb{E}[\langle \nabla \mathcal{L}(\mathbf{W}_\tau^{lin}), \Pi_{M_0}(\mathbf{W}_\tau^{lin} - \mathbf{W}_t^{lin}) \rangle] d\tau \end{aligned} \quad (56)$$

695 where the last inequality is by the convexity of loss function for linearized network. Now to control
696 the integral in (56), we use a similar argument to (51) as follows.

$$\frac{\partial}{\partial \tau} \mathbb{E}[\|\mathbf{W}_\tau^{lin} - \mathbf{W}_t^{lin}\|_{M_0}^2] = -2\mathbb{E}[\langle \nabla \mathcal{L}(\mathbf{W}_\tau^{lin}), \Pi_{M_0}(\mathbf{W}_\tau^{lin} - \mathbf{W}_t^{lin}) \rangle] + \sigma^2 \mathbb{E}[\text{rank}(M_0)] \quad (57)$$

697 By plugging (57) into (56), we have that

$$\frac{\partial}{\partial t} \Phi(t) \leq \frac{\int_t^T -\frac{\partial}{\partial \tau} \mathbb{E}[\|\mathbf{W}_\tau^{lin} - \mathbf{W}_t^{lin}\|_{M_0}^2] + \sigma^2 \mathbb{E}[\text{rank}(M_0)] d\tau}{2(T-t)^2} \leq \frac{\sigma^2 \mathbb{E}[\text{rank}(M_0)]}{2(T-t)} \quad (58)$$

698 By intergrating the above equation over $t \in [0, T - \Delta T]$ where $\Delta \in (0, 1)$ is a tuning parameter that
699 we will determine later, we have that

$$\Phi(T - \Delta T) \leq \Phi(0) + \sigma^2 \mathbb{E}[\text{rank}(M_0)] \ln \frac{1}{\Delta} \quad (59)$$

700 Now we proceed to bound $\Phi(T) - \Phi(T - \Delta T)$. By definition of $\Phi(t)$, we have that

$$\begin{aligned} \Phi(T) - \Phi(T - \Delta T) &= \frac{1}{\Delta T} \int_{T-\Delta T}^T \mathbb{E}[\mathcal{L}(\mathbf{W}_\tau^{lin}; \mathcal{D}) - \mathcal{L}(\mathbf{W}_\tau^{lin}; \mathcal{D})] d\tau \\ &\leq \frac{1}{\Delta T} \int_{T-\Delta T}^T \mathbb{E}[\langle \nabla \mathcal{L}(\mathbf{W}_\tau^{lin}; \mathcal{D}), \Pi_{M_0}(\mathbf{W}_\tau^{lin} - \mathbf{W}_\tau^{lin}) \rangle] d\tau \end{aligned} \quad (60)$$

701 where (60) is by convexity of empirical loss function and by that $\nabla \mathcal{L}^{lin}(\mathbf{W}_T^{lin}; \mathcal{D})$ is in the linear
 702 space spanned by the network output gradient at initialization.

703 By applying the Cauchy-Schwartz inequality, we have that for any $\alpha \in (0, +\infty)$, the following
 704 inequality holds.

$$\Phi(T) - \Phi(T - \Delta T) \leq \frac{1}{\Delta T} \int_{T-\Delta T}^T \frac{1}{\alpha} \mathbb{E}[\|\nabla \mathcal{L}(\mathbf{W}_T^{lin}; \mathcal{D})\|_2^2] + \frac{\alpha}{4} \mathbb{E}[\|\mathbf{W}_T^{lin} - \mathbf{W}_\tau^{lin}\|_{M_0}^2] d\tau \quad (61)$$

$$\leq \frac{1}{\alpha} \max_{\tau \in [T-\Delta T, T]} \mathbb{E}[\|\nabla \mathcal{L}(\mathbf{W}_T^{lin}; \mathcal{D})\|_2^2] + \frac{\alpha}{4} \max_{\tau \in [T-\Delta T, T]} \mathbb{E}[\|\mathbf{W}_T^{lin} - \mathbf{W}_\tau^{lin}\|_{M_0}^2] \quad (62)$$

705 By Eq. (8), we have that when the data is normalized,

$$\mathbb{E}[\|\nabla \mathcal{L}(\mathbf{W}_T^{lin}; \mathcal{D})\|_2^2] \leq 2B \quad (63)$$

706 where $B = |\mathcal{Y}| \left(\prod_{i=1}^{L-1} \frac{\beta_i m_i}{2} \right) \sum_{l=1}^L \frac{\beta_L}{\beta_l}$. We now only need to bound $\mathbb{E}[\|\mathbf{W}_T^{lin} - \mathbf{W}_\tau^{lin}\|_{M_0}^2]$. By
 707 similar argument as (57), for any $t \geq \tau$, we have that

$$\begin{aligned} \frac{\partial}{\partial t} \mathbb{E} [\|\mathbf{W}_t^{lin} - \mathbf{W}_\tau^{lin}\|_{M_0}^2] &= -2\mathbb{E}[\langle \nabla \mathcal{L}(\mathbf{W}_t^{lin}), \Pi_{M_0}(\mathbf{W}_t^{lin} - \mathbf{W}_\tau^{lin}) \rangle] + \sigma^2 \mathbb{E}[\text{rank}(M_0)] \\ &\leq \frac{1}{\gamma} \mathbb{E}[\|\nabla \mathcal{L}(\mathbf{W}_t^{lin})\|_2^2] + \gamma \mathbb{E}[\|\mathbf{W}_t^{lin} - \mathbf{W}_\tau^{lin}\|_{M_0}^2] + \sigma^2 \mathbb{E}[\text{rank}(M_0)] \\ &\leq \gamma \mathbb{E}[\|\mathbf{W}_t^{lin} - \mathbf{W}_\tau^{lin}\|_{M_0}^2] + \frac{2B}{\gamma} + \sigma^2 \mathbb{E}[\text{rank}(M_0)] \end{aligned}$$

708 This is a linear ODE and can be solve closed formly as follows for $t \geq \tau$.

$$\mathbb{E} [\|\mathbf{W}_t^{lin} - \mathbf{W}_\tau^{lin}\|_{M_0}^2] \leq \left(\frac{2B}{\gamma^2} + \frac{\sigma^2 \mathbb{E}[\text{rank}(M_0)]}{\gamma} \right) (e^{\gamma(t-\tau)} - 1)$$

709 Since the above equation holds for any $\tau \leq t$, by setting $\tau = T - \Delta T$ and $t = T$ we have that

$$\max_{\tau \in [T-\Delta T, T]} \mathbb{E}[\|\mathbf{W}_T^{lin} - \mathbf{W}_\tau^{lin}\|_{M_0}^2] \leq \left(\frac{2B}{\gamma^2} + \frac{\sigma^2 \mathbb{E}[\text{rank}(M_0)]}{\gamma} \right) (e^{\gamma \Delta T} - 1) \quad (64)$$

$$\leq 4B(\Delta T)^2 + 2\sigma^2 \mathbb{E}[\text{rank}(M_0)] \Delta T \quad (65)$$

710 where (65) is by setting $\gamma = \frac{1}{\Delta T}$ and by $e^1 - 1 \leq 2$. By combining (62), (63) and (65), we have that

$$\begin{aligned} \Phi(T) - \Phi(T - \Delta T) &\leq \frac{2B}{\alpha} + \alpha B (\Delta T)^2 + \frac{\alpha \sigma^2 \mathbb{E}[\text{rank}(M_0)] \Delta T}{2} \\ &= \frac{3R}{2T} + \frac{\sigma^2 \mathbb{E}[\text{rank}(M_0)]}{2} \end{aligned} \quad (66)$$

711 where the last equation (66) is by setting $\alpha = \frac{2B}{R} T$ and $\Delta = \frac{R}{2BT^2}$ (note that here B and R are as
 712 given in the proposition statement). By combining (59) and (66) and using that $\Delta = \frac{R}{2BT^2}$, we have

$$\begin{aligned} \mathbb{E}[\mathcal{L}(\mathbf{W}_T^{lin})] &= \Phi(T) = \Phi(T) - \Phi(T - \Delta T) + \Phi(T - \Delta T) \\ &\leq \Phi(0) + \frac{\sigma^2 \mathbb{E}[\text{rank}(M_0)]}{2} \left(1 + \log \frac{2BT^2}{R} \right) + \frac{3R}{2T} \end{aligned}$$

713 Observe that $\Phi(0) = \mathbb{E}[\mathcal{L}(\bar{\mathbf{W}}_T^{lin})]$, therefore by using Lemma D.1 we have that

$$\mathbb{E}[\mathcal{L}(\mathbf{W}_T^{lin})] \leq \mathbb{E}[\mathcal{L}(\mathbf{W}_0^*; \mathcal{D})] + \frac{2R}{T} + \frac{\sigma^2 \mathbb{E}[\text{rank}(M_0)]}{2} \left(1 + \log \frac{2BT^2}{R} \right)$$

714 □

715 **D.3 Deferred proof for Proposition 5.3**

716 **Proposition 5.3** (Bounding lazy training distance via smallest eigenvalue of the NTK matrix). Under
 717 the data and network regularity Assumption 2.1, if the width $m_1 = \dots = m_{L-1} = \Omega(n)$ is
 718 sufficiently large, then there exists an optimal solution $\mathbf{W}_0^{\frac{1}{n^2}}$ that satisfies $\mathcal{L}_0^{lin}(\mathbf{W}_0^{\frac{1}{n^2}}) \leq \frac{1}{n^2}$ and
 719 satisfies

$$\tilde{R} = \mathbb{E}[\|\mathbf{W}_0^{\frac{1}{n}} - \mathbf{W}_0\|_2^2] \leq \begin{cases} \tilde{O}\left(\frac{n}{d \cdot 2^L \cdot (m(L-2)+1)}\right) & \text{for NTK initialization} \\ \tilde{O}\left(\frac{n}{2^L m(L-1)}\right) & \text{for He initialization} \\ \tilde{O}\left(\frac{n}{m(L-1)}\right) & \text{for LeCun initialization} \end{cases} \quad (67)$$

720 *Proof.* Given arbitrary initialization parameters \mathbf{W}_0 , we first construct an solution $\mathbf{W}_0^{\frac{1}{n^2}}$ that is nearly
 721 optimal for the ERM problem over $\mathcal{L}_0^{lin}(\mathbf{W})$. Specifically, let $\mathbf{W}_0^{\frac{1}{n^2}}$ have the following expression.

$$\mathbf{W}_0^{\frac{1}{n^2}} - \mathbf{W}_0 = M_0^\dagger \begin{pmatrix} 2 \ln n \cdot y_1 - \mathbf{f}_{\mathbf{W}_0}(\mathbf{x}_1) \\ \vdots \\ 2 \ln n \cdot y_n - \mathbf{f}_{\mathbf{W}_0}(\mathbf{x}_n) \end{pmatrix} \quad (68)$$

722 where $M_0 = \begin{pmatrix} \nabla \mathbf{f}_{\mathbf{W}_0}(\mathbf{x}_1)^\top \\ \vdots \\ \nabla \mathbf{f}_{\mathbf{W}_0}(\mathbf{x}_n)^\top \end{pmatrix}$ is the NTK feature matrix at initialization and \dagger denotes the
 723 pseudo-inverse. By random Gaussian initialization, and by the data regularity assumption As-

724 sumption 2.2, we have that $\text{rank}(M_0) = n$ with probability one, therefore $M_0^\dagger = M_0^\top (M_0 M_0^\top)^{-1}$

725 and $\begin{pmatrix} \mathbf{f}_{\mathbf{W}_0^{\frac{1}{n^2}}}(\mathbf{x}_1) \\ \vdots \\ \mathbf{f}_{\mathbf{W}_0^{\frac{1}{n^2}}}(\mathbf{x}_n) \end{pmatrix} = \begin{pmatrix} 2 \ln n \cdot y_1 \\ \vdots \\ 2 \ln n \cdot y_n \end{pmatrix}$ with probability one. By further using the definition of cross-

726 entropy loss for the single-output setting, we have that the solution $\mathbf{W}_0^{\frac{1}{n^2}}$ satisfies the following
 727 inequality.

$$\mathcal{L}_0^{lin}(\mathbf{W}_0^{\frac{1}{n^2}}) = \log(1 + \exp(-2 \ln n)) < \frac{1}{n^2} \quad (69)$$

728 We now only need to prove that the solution $\mathbf{W}_0^{\frac{1}{n^2}}$ is close to the initialization parameters \mathbf{W}_0 in
 729 expected ℓ_2 norm. By applying the holder inequality on (68), we have that

$$\tilde{R} = \mathbb{E}[\|\mathbf{W}_0^{\frac{1}{n^2}} - \mathbf{W}_0\|_2^2] \leq \mathbb{E} \left[\|M_0^\dagger\|_2^2 \cdot \left\| \begin{pmatrix} 2 \ln n \cdot y_1 - \mathbf{f}_{\mathbf{W}_0}(\mathbf{x}_1) \\ \vdots \\ 2 \ln n \cdot y_n - \mathbf{f}_{\mathbf{W}_0}(\mathbf{x}_n) \end{pmatrix} \right\|_2^2 \right] \quad (70)$$

$$\leq \mathbb{E} \left[\frac{1}{\lambda_{\min}(M_0 M_0^\top)} \cdot \left\| \begin{pmatrix} 2 \ln n \cdot y_1 - \mathbf{f}_{\mathbf{W}_0}(\mathbf{x}_1) \\ \vdots \\ 2 \ln n \cdot y_n - \mathbf{f}_{\mathbf{W}_0}(\mathbf{x}_n) \end{pmatrix} \right\|_2^2 \right] \quad (71)$$

730 We now prove bounds for the two terms on the right hand side separately. For the first term, when
 731 the width is sufficiently large $m = \Omega(n)$, by existing bound for the smallest eigenvalue of the NTK
 732 matrix $M_0 M_0^\top$ in [33, Theorem 4.1], we have that with high probability

$$\frac{1}{\lambda_{\min}(M_0 M_0^\top)} \leq O \left(\frac{1}{\left(d \prod_{l=1}^{L-1} m_l \right) \cdot \left(\prod_{l=1}^L \beta_l \right) \cdot \left(\sum_{l=2}^L \beta_l^{-1} \right)} \right) \quad (72)$$

733 For the second term, by Cauchy-Schwarz inequality, we have that

$$\begin{aligned} \mathbb{E} \left[\left\| \begin{pmatrix} 2 \ln n \cdot y_1 - \mathbf{f}_{\mathbf{W}_0}(\mathbf{x}_1) \\ \vdots \\ 2 \ln n \cdot y_n - \mathbf{f}_{\mathbf{W}_0}(\mathbf{x}_n) \end{pmatrix} \right\|_2^2 \right] &\leq 2 \cdot (2 \ln n)^2 \sum_{i=1}^n y_i^2 + 2 \sum_{i=1}^n \mathbb{E}[\|\mathbf{f}_{\mathbf{W}_0}(\mathbf{x}_i)\|_2^2] \\ &= 8n(\ln n)^2 + 2 \sum_{i=1}^n \mathbb{E}[\|\mathbf{f}_{\mathbf{W}_0}(\mathbf{x}_i)\|_2^2] \end{aligned}$$

734 By further using Eq. (44), we have that

$$\mathbb{E} \left[\left\| \begin{pmatrix} 2 \ln n \cdot y_1 - \mathbf{f}_{\mathbf{W}_0}(\mathbf{x}_1) \\ \vdots \\ 2 \ln n \cdot y_n - \mathbf{f}_{\mathbf{W}_0}(\mathbf{x}_n) \end{pmatrix} \right\|_2^2 \right] = O \left(n(\ln n)^2 + n \prod_{i=1}^L \frac{\beta_i m_i}{2} \right) \quad (73)$$

735 Therefore, by plugging Eq. (72) and Eq. (73) into Eq. (71), and by considering input dimension d as
736 const, we have that

$$\tilde{R} = O \left(\frac{n(\ln n)^2 + n \prod_{i=1}^L \frac{\beta_i m_i}{2}}{\left(d \prod_{l=1}^{L-1} \beta_l m_l \right) \cdot \beta_L \cdot \left(\sum_{l=2}^L \beta_l^{-1} \right)} \right)$$

737 By plugging the choice of initialization variance β_1, \dots, β_L for NTK, He and LeCun initialization
738 into the above equation, for single output network with $L \geq 2$, we have that

$$\tilde{R} = \begin{cases} \tilde{O}\left(\frac{n}{d \cdot 2^L \cdot (m(L-2)+1)}\right) & \text{for NTK initialization} \\ \tilde{O}\left(\frac{n}{2^L m(L-1)}\right) & \text{for He initialization} \\ \tilde{O}\left(\frac{n}{m(L-1)}\right) & \text{for LeCun initialization} \end{cases}$$

739 □

740 D.4 Deferred proof for Corollary 5.4

741 **Corollary 5.4** (Privacy utility trade-off for last iterate). Assume that the data and network regularity
742 Assumption 2.2 holds. Assume that all the conditions and definition for constants in Proposition 5.1
743 holds. Then by setting $\sigma^2 = \frac{2BT}{\varepsilon n^2}$ and $T = \sqrt{\frac{2\varepsilon n \tilde{R}}{B}}$, we have that running Langevin diffusion for
744 time T satisfies bound KL divergence ε , and has empirical excess risk upper bounded by

$$\mathbb{E}[\mathcal{L}(\mathbf{W}_T^{lin})] \leq \mathcal{O} \left(\frac{1}{n^2} + \sqrt{\frac{B \tilde{R}}{\varepsilon n}} \log(\varepsilon n) \right) \quad (74)$$

745 where B is the gradient norm constant Eq. (9), and \tilde{R} is the approximate lazy training distance in
746 Eq. (12). A summary of B and \tilde{R} under different initializations is in Table 1.

747 *Proof.* By setting $\mathbf{W}_0^* = \mathbf{W}_0^{\frac{1}{n^2}}$ in Proposition 5.3, we have that

$$\mathbb{E}[\mathcal{L}(\mathbf{W}_T^{lin})] - \mathbb{E}[\mathcal{L}(\mathbf{W}_0^{\frac{1}{n^2}}; \mathcal{D})] \leq \frac{2\tilde{R}}{T} + \frac{\sigma^2 \mathbb{E}[\text{rank}(M_0)]}{2} \left(1 + \log \frac{2BT^2}{\tilde{R}} \right)$$

748 By Proposition 5.3, we have that $\mathbb{E}[\mathcal{L}(\mathbf{W}_0^{\frac{1}{n^2}}; \mathcal{D})] \leq \frac{1}{n^2}$, therefore

$$\mathbb{E}[\mathcal{L}(\mathbf{W}_T^{lin})] \leq \frac{1}{n^2} + \frac{2\tilde{R}}{T} + \frac{\sigma^2 \mathbb{E}[\text{rank}(M_0)]}{2} \left(1 + \log \frac{2BT^2}{\tilde{R}} \right) \quad (75)$$

749 By plugging $\sigma^2 = \frac{2BT}{\varepsilon n^2}$ and $T = \sqrt{\frac{2\varepsilon n \tilde{R}}{B}}$ into (75), and by $\text{rank}(M_0) \leq n$, we have that

$$\mathbb{E}[\mathcal{L}(\mathbf{W}_T^{lin})] \leq \frac{1}{n^2} + \sqrt{\frac{2B \tilde{R}}{\varepsilon n}} (2 + \log(4\varepsilon n)) \quad (76)$$

750 □

751 **D.5 Deferred proofs for Proposition 5.5**

752 We will first need to prove following lemma to bound the uniform stability for the model parameters
753 during training linearized network.

754 **Lemma D.2** (Uniform stability for training linearized network). *Assume that the data and network
755 regularity Assumption 2.2 holds. Then it satisfies that*

$$\mathbb{E} [\|\mathbf{W}_T - \mathbf{W}'_T\|^2] \leq \frac{2B}{n^2} T^2 \quad (77)$$

756 where $B = |\mathcal{Y}| \left(\prod_{i=1}^{L-1} \frac{\beta_i m_i}{2} \right) \sum_{l=1}^L \frac{\beta_L}{\beta_l}$ is the constant defined in (37) for network output gradient
757 norm bound at initialization.

758 *Proof.* By definition, and by coupling the choice of Gaussian noise in the two Langevin diffusion
759 processes $(\mathbf{W}_t)_{t \in [0, T]}$ and $(\mathbf{W}'_t)_{t \in [0, T]}$, we have that

$$\begin{aligned} & \frac{\partial \mathbb{E} [\|\mathbf{W}_t - \mathbf{W}'_t\|^2]}{\partial t} = \lim_{\eta \rightarrow 0} \frac{\mathbb{E} [\|\mathbf{W}_{t+\eta} - \mathbf{W}'_{t+\eta}\|^2] - \mathbb{E} [\|\mathbf{W}_t - \mathbf{W}'_t\|^2]}{\eta} \\ &= \lim_{\eta \rightarrow 0} \frac{\mathbb{E} [\|\mathbf{W}_t - \eta \nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}) - \mathbf{W}'_t + \eta \nabla \mathcal{L}(\mathbf{W}'_t; \mathcal{D}')\|^2] - \mathbb{E} [\|\mathbf{W}_t - \mathbf{W}'_t\|^2]}{\eta} \\ &= \lim_{\eta \rightarrow 0} \frac{\eta^2 \mathbb{E} [\|\nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}'_t; \mathcal{D}')\|^2] - 2\eta \mathbb{E} [\langle \nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}'_t; \mathcal{D}'), \mathbf{W}_t - \mathbf{W}'_t \rangle]}{\eta} \\ &= -2\mathbb{E} [\langle \nabla \mathcal{L}(\mathbf{W}_t; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}'_t; \mathcal{D}'), \mathbf{W}_t - \mathbf{W}'_t \rangle] \leq 2\sqrt{\frac{2B}{n^2}} \cdot \mathbb{E} [\|\mathbf{W}_T - \mathbf{W}'_T\|^2] \end{aligned}$$

760 where the last inequality is by holder's inequality and by using Lemma 4.1. By solving the ordinary
761 differential equation with boundary condition $\mathbb{E} [\|\mathbf{W}_0 - \mathbf{W}'_0\|^2] = 0$, we have that

$$\mathbb{E} [\|\mathbf{W}_T - \mathbf{W}'_T\|^2] \leq \frac{2B}{n^2} T^2 \quad (78)$$

762

□

763 We are now ready to prove our excess population risk bound for training linearized network.

764 **Proposition 5.5.** Denote $R_0(\mathbf{W}) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \in \text{pop}} [\ell(f_{\mathbf{W}_0}(\mathbf{x}) + \frac{\partial f_{\mathbf{W}_0}(\mathbf{x})}{\partial \mathbf{W}_0}(\mathbf{W} - \mathbf{W}_0); \mathbf{y})]$ as the popula-
765 tion risk of linearized network expanded at initialization vector \mathbf{W}_0 over population data distribution
766 *pop*. Then under the conditions of Corollary 5.4, we have that for any dataset \mathcal{D} of size n , the
767 following excess population risk upper bound holds.

$$\mathbb{E}[R_0(\mathbf{W}_T)] - \mathbb{E}[\mathcal{L}(\mathbf{W}_{pop,0}^*; \mathcal{D})] \leq O \left(\frac{1}{n^2} + \sqrt{\frac{B\tilde{R}}{\varepsilon n}} (\log(\varepsilon n) + \varepsilon) \right) \quad (79)$$

768 where the expectation is over the randomness of sampling the training dataset $\mathcal{D} \sim \text{pop}^n$ from the
769 data population and the random coins for the Langevin diffusion training algorithm, and $\mathbf{W}_{pop,0}^* =$
770 $\text{argmin}_{\mathbf{W}} R_0(\mathbf{W})$ is the optimal solution for the population risk minimization problem.

771 *Proof.* By the uniform stability method [24, Theorem 2.2], we have the following generalization
772 error upper bound holds.

$$\alpha_{gen} = |\mathbb{E}[R_0(\mathbf{W}_T)] - \mathbb{E}_{\mathcal{D} \sim \text{pop}^n} [\mathcal{L}(\mathbf{W}_T; \mathcal{D})]| \leq \max_{z, \mathcal{D}, \mathcal{D}'} \mathbb{E} [\ell(f_{\mathbf{W}_T}(\mathbf{x}_z); \mathbf{y}_z) - \ell(f_{\mathbf{W}'_T}(\mathbf{x}_z); \mathbf{y}_z)] \quad (80)$$

773 where $\mathbf{z} = (\mathbf{x}_z, \mathbf{y}_z)$ is an arbitrary data point in the population data distribution *pop*. By convexity,
774 we further have the following bound for the uniform stability.

$$\ell(f_{\mathbf{W}_T}(\mathbf{x}_z); \mathbf{y}_z) - \ell(f_{\mathbf{W}'_T}(\mathbf{x}_z); \mathbf{y}_z) \leq \nabla \ell(f_{\mathbf{W}_T}(\mathbf{x}_z); \mathbf{y}_z)^\top (\mathbf{W}_T - \mathbf{W}'_T) \quad (81)$$

$$\leq \sqrt{\|\nabla \ell(f_{\mathbf{W}_T}(\mathbf{x}_z); \mathbf{y}_z)\|^2 \|\mathbf{W}_T - \mathbf{W}'_T\|^2} \quad (82)$$

775 By Eq. (37) and Lemma D.2, we further have that

$$\mathbb{E}[\ell(f_{\mathbf{W}_T}(\mathbf{x}_z); \mathbf{y}_z) - \ell(f_{\mathbf{W}'_T}(\mathbf{x}_z); \mathbf{y}_z)] \leq \sqrt{2B} \cdot \sqrt{\frac{2B}{n^2} T^2} = \frac{2B}{n} T \quad (83)$$

776 Therefore, by plugging the above equation into (80), we have that the generalization error satisfies

$$\alpha_{gen} = |\mathbb{E}_{\mathcal{D} \sim \text{pop}^n}[\mathcal{L}(\mathbf{W}_T; \mathcal{D})] - \mathbb{E}[R(\mathbf{W}_T)]| \leq \frac{2BT}{n} \leq 2\sqrt{\frac{2\varepsilon BR}{n}} \quad (84)$$

777 where the last inequality is by plugging our choice of $T = \sqrt{\frac{2\varepsilon n R}{B}}$ into the equation. On the other
778 hand, by Proposition 5.1, we have that the empirical risk is upper bounded as follows.

$$\begin{aligned} \mathbb{E}_{\mathcal{D} \sim \text{pop}^n}[\mathcal{L}(\mathbf{W}_T; \mathcal{D})] &\leq O\left(\sqrt{\frac{B\bar{R}}{\varepsilon n}} \log(\varepsilon n)\right) \\ &\leq \mathbb{E}[\mathcal{L}(\mathbf{W}_{\text{pop},0}^*; \mathcal{D})] + O\left(\sqrt{\frac{B\tilde{R}}{\varepsilon n}} \log(\varepsilon n)\right) \end{aligned} \quad (85)$$

779 Combining the generalization error term (84) and the excess empirical risk term (85) suffice to prove
780 the equation in the statement. \square

781 E Discussion on extending our results to Noisy GD with constant step-size

782 In this section, we discuss how to extend our privacy analyses to noisy GD with constant step-size.
783 Specifically, we only need to extend the KL composition theorem under possibly unbounded gradient
784 difference, i.e., Theorem 3.1, to the noisy GD algorithm.

785 **Theorem E.1** (KL composition for noisy GD under possibly unbounded gradient difference). *Let*
786 *the iterative update in noisy GD algorithm be defined by: $\mathbf{W}_{(k+1)} = \mathbf{W}_{(k)} - \eta \nabla \mathcal{L}(\mathbf{W}_{(k)}; \mathcal{D}) +$*
787 *$\sqrt{2\eta\sigma^2} Z_k$, where $Z_k \sim \mathcal{N}(0, \mathbb{I})$. Then the KL divergence between running noisy GD for DNN (2)*
788 *on neighboring datasets \mathcal{D} and \mathcal{D}' satisfies*

$$KL(\mathbf{W}_{(K)}, \mathbf{W}'_{(K)}) \leq \frac{1}{2\sigma^2} \sum_{k=0}^{K-1} \eta \cdot \mathbb{E} \left[\|\nabla \mathcal{L}(\mathbf{W}_{(k)}; \mathcal{D}) - \nabla \mathcal{L}(\mathbf{W}_{(k)}; \mathcal{D}')\|_2^2 \right]. \quad (86)$$

789 *Proof.* Denote $p_{(k)}$ as the distribution of model parameters after running noisy GD on dataset \mathcal{D} with
790 k steps, and similarly denote $p'_{(k)}$ as the distribution of model parameters after running noisy GD on
791 dataset \mathcal{D}' with k steps. Then by the data processing inequality for KL divergence [45, Theorem 9]
792 (with the data processing operation given by $(\mathbf{W}_{(k)}, \mathbf{W}_{(k+1)}) \rightarrow \mathbf{W}_{(k)}$), we have that

$$KL(p_{(k+1)}, p'_{(k+1)}) \leq KL(p_{(k), (k+1)}, p'_{(k), (k+1)}), \quad (87)$$

793 where $p_{(k), (k+1)}$ denotes the joint distribution of $(\mathbf{W}_{(k)}, \mathbf{W}_{(k+1)})$, and $p'_{(k), (k+1)}$ denotes the joint
794 distribution of $(\mathbf{W}'_{(k)}, \mathbf{W}'_{(k+1)})$. Now we expand the term $KL(p_{(k), (k+1)}, p'_{(k), (k+1)})$ by the Bayes
795 rule as follows.

$$\begin{aligned} &KL(p_{(k), (k+1)}, p'_{(k), (k+1)}) \quad (88) \\ &= \mathbb{E}_{p_{(k), (k+1)}}(\mathbf{W}_{(k)}, \mathbf{W}_{(k+1)}) \left[\log \left(\frac{p_{(k+1)|(k)}(\mathbf{W}_{(k+1)} | \mathbf{W}_{(k)}) p_{(k)}(\mathbf{W}_{(k)})}{p'_{(k+1)|(k)}(\mathbf{W}_{(k+1)} | \mathbf{W}_{(k)}) p'_{(k)}(\mathbf{W}_{(k)})} \right) \right] \\ &= \mathbb{E}_{p_{(k), (k+1)}}(\mathbf{W}_{(k)}, \mathbf{W}_{(k+1)}) \left[\log \left(\frac{p_{(k+1)|(k)}(\mathbf{W}_{(k+1)} | \mathbf{W}_{(k)})}{p'_{(k+1)|(k)}(\mathbf{W}_{(k+1)} | \mathbf{W}_{(k)})} \right) \right] + \mathbb{E}_{p_{(k)}(\mathbf{W}_{(k)})} \left[\log \left(\frac{p_{(k)}(\mathbf{W}_{(k)})}{p'_{(k)}(\mathbf{W}_{(k)})} \right) \right] \\ &= \mathbb{E}_{p_{(k)}(\mathbf{W}_{(k)})} \left[KL(p_{(k+1)|(k)}, p'_{(k+1)|(k)}) \right] + KL(p_{(k)}, p'_{(k)}) \end{aligned} \quad (89)$$

796 Observe that $p_{(k+1)|(k)}, p'_{(k+1)|(k)}$ are two Gaussian distributions with per-dimensional variance σ^2 ,
 797 due to the conditioning on the same model parameters $\mathbf{W}_{(k)}$ at iteration k . Therefore, by computing
 798 the KL divergence between two multivariate Gaussians, we have that

$$KL(p_{(k),(k+1)}, p'_{(k),(k+1)}) = \frac{1}{2\sigma^2} \cdot \eta \cdot \|\nabla\mathcal{L}(\mathbf{W}_{(k)}; \mathcal{D}) - \nabla\mathcal{L}(\mathbf{W}_{(k)}; \mathcal{D}')\|_2^2 \quad (90)$$

799 Therefore, by plugging Eq. (90) into Eq. (89) and Eq. (87), we have that

$$KL(p_{(k+1)}, p'_{(k+1)}) \leq \frac{\eta}{2\sigma^2} \mathbb{E} \left[\|\nabla\mathcal{L}(\mathbf{W}_{(k)}; \mathcal{D}) - \nabla\mathcal{L}(\mathbf{W}_{(k)}; \mathcal{D}')\|_2^2 \right] + KL(p_{(k)}, p'_{(k)}) \quad (91)$$

800 By summing (91) over $k = 0, \dots, K-1$ and observing that $KL(p_{(0)}, p'_{(0)}) = 0$ (as the initialization
 801 distribution is the same between noisy GD on \mathcal{D} and \mathcal{D}'), we finish the proof for Eq. (86). \square