
Deep Networks as Denoising Algorithms: Sample-Efficient Learning of Diffusion Models in High-Dimensional Graphical Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We investigate the efficiency of deep neural networks for approximating scoring
2 functions in diffusion-based generative modeling. While existing approximation
3 theories leverage the smoothness of score functions, they suffer from the curse of di-
4 mensionality for intrinsically high-dimensional data. This limitation is pronounced
5 in graphical models such as Markov random fields, where the approximation
6 efficiency of score functions remains unestablished.

7 To address this, we note score functions can often be well-approximated in graphical
8 models through variational inference denoising algorithms. Furthermore, these
9 algorithms can be efficiently represented by neural networks. We demonstrate this
10 through examples, including Ising models, conditional Ising models, restricted
11 Boltzmann machines, and sparse encoding models. Combined with off-the-shelf
12 discretization error bounds for diffusion-based sampling, we provide an efficient
13 sample complexity bound for diffusion-based generative modeling when the score
14 function is learned by deep neural networks.

15 1 Introduction

16 In recent years, diffusion models [Sohl-Dickstein et al., 2015, Ho et al., 2020, Song and Ermon,
17 2019, Song et al., 2020] have emerged as a leading approach for generative modeling, achieving
18 state-of-the-art results across diverse domains. Given a dataset of n independent and identically
19 distributed samples $\{\mathbf{x}_i\}_{i=1}^n$ drawn from an unknown distribution $\mu \in \mathcal{P}(\mathbb{R}^d)$, diffusion models aim
20 to learn a generative model that produces new samples $\hat{\mathbf{x}} \sim \hat{\mu}$ that match this distribution. Popular
21 diffusion models such as DDPM [Ho et al., 2020] achieve this through a two-step procedure:

- 22 • **Step 1.** Fit approximate score functions $\hat{\mathbf{s}}_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$ for $t \in [0, T]$ by minimizing the
23 following empirical risk over a neural network class \mathcal{F} :

$$\hat{\mathbf{s}}_t = \arg \min_{\text{NN} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \|\sigma_t^{-1} \mathbf{g}_i + \text{NN}(\lambda_t \mathbf{x}_i + \sigma_t \mathbf{g}_i)\|_2^2. \quad (\text{ERM})$$

24 In the above display, $\mathbf{g}_i \sim_{iid} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and $(\lambda_t, \sigma_t^2) = (e^{-t}, 1 - e^{-2t})$.

- 25 • **Step 2.** Discretize the following stochastic differential equation (SDE) from Gaussian
26 initialization, whose drift term is given by the fitted approximate score functions:

$$d\mathbf{Y}_t = (\mathbf{Y}_t + 2\hat{\mathbf{s}}_{T-t}(\mathbf{Y}_t))dt + \sqrt{2}d\mathbf{B}_t, \quad t \in [0, T], \quad \mathbf{Y}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d), \quad (\text{SDE})$$

27 and take the approximate sample $\hat{\mathbf{x}} = \mathbf{Y}_T \in \mathbb{R}^d$.

28 *Score functions* $s_t(\mathbf{z})$ are central to the diffusion model framework. Given infinite data and model
 29 capacity, the minimizer of the empirical risk in Eq. (ERM) yields the score function,

$$s_t(\mathbf{z}) = \nabla_{\mathbf{z}} \log \mu_t(\mathbf{z}), \quad \mu_t(\mathbf{z}) : \text{density of } \mathbf{z}, \mathbf{x} \sim \mu \text{ and } [\mathbf{z}|\mathbf{x}] \sim \mathcal{N}(\lambda_t \mathbf{x}, \sigma_t^2 \mathbf{I}_d). \quad (\text{Score})$$

30 The sample quality from diffusion models relies on two key factors: (1) how well \hat{s}_t approximates s_t ;
 31 and (2) how accurately the SDE discretization scheme approximates process (SDE). Recent work has
 32 made substantial progress on controlling the SDE discretization error in diffusion models, assuming
 33 access to a good score function estimator [Chen et al., 2022a, 2023a, Lee et al., 2023, Li et al., 2023a,
 34 Benton et al., 2023]. However, understanding when neural networks can accurately estimate the score
 35 function itself remains less explored. Some analyses rely on strong distributional assumptions for
 36 score function realizability [Shah et al., 2023, Yuan et al., 2023], while others exploit the smoothness
 37 of score functions, incurring the curse of dimensionality [Oko et al., 2023, Chen et al., 2023b]. These
 38 results do not cover many common high-dimensional graphical models for images and text, such as
 39 Markov random fields or restricted Boltzmann machines [Geman and Graffigne, 1986, Ranzato et al.,
 40 2010, Conroy and O’leary, 2001].

41 **A new perspective on score function approximation.** We provide a new perspective on approxi-
 42 mating diffusion model score functions with neural networks. First, we observe that by Tweedie’s
 43 formula, score functions s_t are related to denoising functions \mathbf{m}_t :

$$s_t(\mathbf{z}) = (\lambda_t \cdot \mathbf{m}_t(\mathbf{z}) - \mathbf{z}) / \sigma_t^2, \quad \mathbf{m}_t(\mathbf{z}) = \mathbb{E}_{(\mathbf{x}, \mathbf{g}) \sim \mu \otimes \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[\mathbf{x} | \lambda_t \mathbf{x} + \sigma_t \mathbf{g} = \mathbf{z}]. \quad (\text{Denoiser})$$

44 Our key insight is that if the data distribution μ arises from a graphical model, these denoisers $\mathbf{m}_t(\mathbf{z})$
 45 can often be approximated by variational inference (VI) algorithms, which takes the form

$$\mathbf{m}_t(\mathbf{z}) \approx \mathbf{f}_{\text{out}}(\mathbf{u}^{(L)}), \quad \mathbf{u}^{(\ell)} = \mathbf{f}_{\ell}(\mathbf{u}^{(\ell-1)}), \quad \ell \in \{1, \dots, L\}, \quad \mathbf{u}^{(0)} = \mathbf{f}_{\text{in}}(\mathbf{z}). \quad (\text{VI})$$

46 For instance, when μ is an Ising model, \mathbf{m}_t can be approximated by an iterative algorithm that
 47 minimizes a VI objective [Jordan et al., 1999, Wainwright et al., 2008]. Each update step \mathbf{f}_{ℓ} is
 48 composed of simple operations, including matrix-vector multiplication and pointwise nonlinearity,
 49 which can be captured by a two-layer neural network $\mathbf{f}_{\ell}(\mathbf{u}) \approx \mathbf{u} + \mathbf{W}_1 \cdot \text{ReLU}(\mathbf{W}_2 \mathbf{u})$. By comparing
 50 updates (VI) and residual network forms (ResNet), we can see how the iterative variational inference
 51 steps directly translate to residual block approximations. This establishes a clear connection between
 52 variational inference in graphical models and score approximation in diffusion models.

53 2 Preliminaries: the DDPM sampling scheme

Algorithm 1 The DDPM sampling scheme

Require: $\{\mathbf{x}_i\}_{i \in [n]}$, (d, D, L, M, B) , $(N, T, \delta, \{t_k\}_{0 \leq k \leq N})$ with $0 = t_0 < \dots < t_N = T - \delta$.

- 1: // Computing the approximate score function
- 2: Sample $\{\mathbf{g}_i\}_{i \in [n]} \sim \text{iid } \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$.
- 3: **for** $t \in \{T - t_k\}_{0 \leq k \leq N-1}$ **do**
- 4: Solve the ERM problem below for $t = T - t_k$:

$$\widehat{\mathbf{W}}_t = \arg \min_{\mathbf{W} \in \mathcal{W}_{d, D, L, M, B}} \frac{1}{n} \sum_{i=1}^n \left\| \sigma_t^{-1} \mathbf{g}_i + \text{P}_t[\text{ResN}_{\mathbf{W}}](\lambda_t \mathbf{x}_i + \sigma_t \mathbf{g}_i) \right\|_2^2. \quad (1)$$

- 5: Take the approximate score function to be $\hat{s}_t(\mathbf{z}) = \text{P}_t[\text{ResN}_{\widehat{\mathbf{W}}_t}](\mathbf{z})$.
- 6: // Sampling by discretizing the stochastic differential equation
- 7: Sample $\widehat{\mathbf{Y}}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$.
- 8: **for** $k = 0, \dots, N - 1$ **do**
- 9: Sample $\mathbf{G}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. Calculate $\widehat{\mathbf{Y}}_{k+1}$ using the exponential integrator scheme:

$$\widehat{\mathbf{Y}}_{k+1} = e^{\gamma_k} \cdot \widehat{\mathbf{Y}}_k + 2(e^{\gamma_k} - 1) \cdot \hat{s}_{T-t_k}(\widehat{\mathbf{Y}}_k) + \sqrt{e^{2\gamma_k} - 1} \cdot \mathbf{G}_k, \quad \gamma_k = t_{k+1} - t_k. \quad (2)$$

Return: $\hat{\mathbf{x}} = \widehat{\mathbf{Y}}_N$.

54 This section provides details on the two-step DDPM sampling scheme in Algorithm 1. The inputs
 55 of the algorithm are n IID samples $\{\mathbf{x}_i\}_{i \in [n]}$ from μ . The algorithm also receives parameters

56 (d, D, L, M, B) for specifying the ResNet class, and $(N, T, \delta, \{t_k\}_{0 \leq k \leq N})$ for specifying the time
 57 discretization scheme. The first step of the algorithm performs empirical risk minimization to compute
 58 the approximate score functions \hat{s}_t (lines 2-5). The second step generates a sample by discretizing
 59 the reverse-time SDE using the fitted score functions (lines 7-9). We discuss the score learning and
 60 SDE discretization steps in more detail below.

61 **ERM and the ResNet class.** The first step of Algorithm 1 solves an ERM problem (1) to fit the score
 62 functions. This regresses standard Gaussian noises $\{\mathbf{g}_i\}_{i \in [n]}$ on the noisy samples $\{\lambda_t \mathbf{x}_i + \sigma_t \mathbf{g}_i\}_{i \in [n]}$,
 63 using a standard ResNet architecture $\text{ResN}_{\mathbf{W}} : \mathbb{R}^d \rightarrow \mathbb{R}^d$. The ResNet is parameterized by a set
 64 of weight matrices $\mathbf{W} = \{\mathbf{W}_1^{(\ell)} \in \mathbb{R}^{D \times M}, \mathbf{W}_2^{(\ell)} \in \mathbb{R}^{M \times D}\}_{\ell \in [L]} \cup \{\mathbf{W}_{\text{in}} \in \mathbb{R}^{(d+1) \times D}, \mathbf{W}_{\text{out}} \in$
 65 $\mathbb{R}^{D \times d}\}$ with embedding dimension D , number of layers L , and hidden-layer width M . It applies
 66 iterative residual blocks with ReLU nonlinearities to map an input \mathbf{z} to an output in \mathbb{R}^d :

$$\text{ResN}_{\mathbf{W}}(\mathbf{z}) = \mathbf{W}_{\text{out}} \mathbf{u}^{(L)}, \quad \mathbf{u}^{(\ell)} = \mathbf{u}^{(\ell-1)} + \mathbf{W}_1^{(\ell)} \text{ReLU}(\mathbf{W}_2^{(\ell)} \mathbf{u}^{(\ell-1)}), \quad \mathbf{u}^{(0)} = \mathbf{W}_{\text{in}}[\mathbf{z}; 1].$$

(ResNet)

67 The minimization in (1) is over the ResNets whose weights are contained in a B -bounded set,
 68 specified by parameters (d, D, L, M, B)

$$\mathcal{W}_{d,D,L,M,B} := \left\{ \mathbf{W} = \{\mathbf{W}_1^{(\ell)}, \mathbf{W}_2^{(\ell)}\}_{\ell \in [L]} \cup \{\mathbf{W}_{\text{in}}, \mathbf{W}_{\text{out}}\} : \|\mathbf{W}\| \leq B \right\}. \quad (3)$$

69 Here the norm of ResNet weights is defined as

$$\|\mathbf{W}\| := \max_{\ell \in [L]} \{\|\mathbf{W}_1^{(\ell)}\|_{\text{op}} + \|\mathbf{W}_2^{(\ell)}\|_{\text{op}}\} \vee \max \{\|\mathbf{W}_{\text{in}}\|_{\text{op}}, \|\mathbf{W}_{\text{out}}\|_{\text{op}}\}. \quad (4)$$

70 For technical reasons, we truncate the ResNet output using P_t . Given a function $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$, we
 71 define $\text{P}_t[f](\mathbf{z}) = \text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(f(\mathbf{z}) + \sigma_t^{-2} \mathbf{z}) - \sigma_t^{-2} \mathbf{z}$, where $\text{proj}_R(\mathbf{z})$ is the projector of $\mathbf{z} \in \mathbb{R}^d$
 72 into the R -Euclidean ball. Note that when $f(\mathbf{z})$ is a score function, $f(\mathbf{z}) + \sigma_t^{-2} \mathbf{z}$ is a rescaled
 73 denoising function and should be bounded for data distribution with compact support. This operator
 74 is a technical detail that could be eliminated in practice — it is only used to control the generalization
 75 error of the empirical risk minimization problem.

76 **Choice of the discretization scheme.** We choose a particular scheme that uses a uniform grid
 77 in the first phase and an exponential decaying grid in the second phase. As shown in Benton et al.
 78 [2023], such a scheme provides a sharp sampling error control. We delay the detailed description of
 79 our discretization scheme to Appendix A.1.

80 **The conditional diffusion model.** In conditional generative modeling tasks, we observe IID
 81 samples $\{(\mathbf{x}_i, \boldsymbol{\theta}_i)\}_{i \in [n]} \sim_{\text{iid}} \mu$, and our goal is to learn a model to generate new samples $\hat{\mathbf{x}}$ from the
 82 conditional distribution $\mu(\mathbf{x}|\boldsymbol{\theta})$ for a given $\boldsymbol{\theta}$.

83 The DDPM sampling scheme can be simply adapted to solve conditional generative modeling tasks,
 84 as per Algorithm 2. Specifically, we modify the ResNet in empirical risk minimization to take the
 85 form (ResNet-Conditional), admitting inputs $(\lambda_t \mathbf{x}_i + \sigma_t \mathbf{g}_i, \boldsymbol{\theta}_i) \in \mathbb{R}^d \times \mathbb{R}^m$. The approximated score
 86 functions $\hat{s}_t(\mathbf{z})$ become conditional $\hat{s}_t(\mathbf{z}; \boldsymbol{\theta}) = \text{P}_t[\text{ResN}_{\widehat{\mathbf{W}}_t}](\mathbf{z}, \boldsymbol{\theta})$, estimating the conditional score
 87 functions $\mathbf{s}_t(\mathbf{z}; \boldsymbol{\theta}) = \nabla_{\mathbf{z}} \log \mu_t(\mathbf{z}, \boldsymbol{\theta})$, where μ_t is the joint density of $(\mathbf{z}, \boldsymbol{\theta})$ when $(\mathbf{x}, \boldsymbol{\theta}) \sim \mu$ and
 88 $[\mathbf{z}|\mathbf{x}] \sim \mathcal{N}(\lambda_t \mathbf{x}, \sigma_t^2 \mathbf{I}_d)$. Details of the conditional algorithm are provided in Appendix D.1.

89 3 Diffusion models for Ising models

90 The Ising model $\mu \in \mathcal{P}(\{\pm 1\}^d)$ is a distribution over the discrete hypercube, with probability mass
 91 function characterized by an energy function of spin configurations. Specifically,

$$\mu(\mathbf{x}) = Z^{-1} \exp\{\langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle / 2\}, \quad \mathbf{x} \in \{\pm 1\}^d, \quad Z = \sum_{\mathbf{x} \in \{\pm 1\}^d} \exp\{\langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle / 2\}. \quad (\text{Ising})$$

92 The Ising model stands as one of the most fundamental graphical models; it belongs to the exponential
 93 family, yet its normalizing constant, Z , does not possess an analytic expression.

94 Consider the task of generative modeling where the input consists of IID samples $\{\mathbf{x}_i\}_{i \in [n]} \sim \mu$
 95 derived from the Ising model. To demonstrate that Algorithm 1 outputs valid samples, we need to

96 control the estimation error $\mathbb{E}[\|\hat{\mathbf{s}}_t(\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2]$. Recall that \mathbf{s}_t relates to \mathbf{m}_t . To calculate $\mathbf{m}_t(\mathbf{z})$,
 97 one often seeks to minimize certain type of free energy, for instance, the naive variational Bayes free
 98 energy [Wainwright et al., 2008]. To establish our main result, we will assume the consistency of a
 99 free energy minimizer with the denoiser.

100 **Assumption 1.** Let $\mathbf{x} \sim \mu(\boldsymbol{\sigma}) \propto \exp\{\langle \boldsymbol{\sigma}, \mathbf{A}\boldsymbol{\sigma} \rangle / 2\}$ and $\mathbf{z} \sim \mathcal{N}(\lambda_t \mathbf{x}, \sigma_t^2 \mathbf{I}_d)$. Denote the marginal
 101 distribution of \mathbf{z} by μ_t . For any fixed t , assume that there exists $\varepsilon_{\text{VI},t}^2(\mathbf{A}) < \infty$ and $\mathbf{K} = \mathbf{K}(\mathbf{A}, t)$
 102 with $\|\mathbf{K} - \mathbf{A}\|_{\text{op}} \leq A < 1$, such that

$$\begin{aligned} \mathbb{E}_{\mathbf{z} \sim \mu_t} [\|\hat{\mathbf{m}}_t(\mathbf{z}) - \mathbf{m}_t(\mathbf{z})\|_2^2] / d &\leq \varepsilon_{\text{VI},t}^2(\mathbf{A}), \\ \hat{\mathbf{m}}_t(\mathbf{z}) &= \operatorname{argmin}_{\mathbf{m} \in [-1,1]^d} \left\{ \sum_{i=1}^d -\mathfrak{h}_{\text{bin}}(m_i) - \frac{1}{2} \langle \mathbf{m}, \mathbf{A}\mathbf{m} \rangle - \frac{\lambda_t}{\sigma_t^2} \langle \mathbf{z}, \mathbf{m} \rangle + \frac{1}{2} \langle \mathbf{m}, \mathbf{K}\mathbf{m} \rangle \right\}. \end{aligned}$$

103 In Appendix A.2, we will discuss cases in which the VI approximation error $\varepsilon_{\text{VI},t}^2(\mathbf{A})$ can be well-
 104 controlled. Given Assumption 1 holds, we are ready to provide a bound on the estimation error
 105 of the approximate score function. We give a proof outline in Appendix A.4 and the full proof in
 106 Appendix E.

107 **Theorem 1.** Let Assumption 1 hold. Let $\{\hat{\mathbf{s}}_{T-t_k}\}_{0 \leq k \leq N-1}$ be the approximate score function given
 108 by Algorithm 1 in which we take

$$D = 3d, \quad M \geq 4d, \quad B \geq 7 \cdot (M/d) \cdot \log(M) + 1/\min_k \{T - t_k\} + \sqrt{d}.$$

109 Then with probability at least $1 - \eta$, for any $t \in \{T - t_k\}_{0 \leq k \leq N-1}$, we have

$$\mathbb{E}_{\mathbf{z} \sim \mu_t} [\|\hat{\mathbf{s}}_t(\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2] / d \lesssim \lambda_t^2 \sigma_t^{-4} \cdot \left(\varepsilon_{\text{VI},t}^2(\mathbf{A}) + \varepsilon_{\text{ResN}}^2 + \varepsilon_{\text{gen}}^2 \right), \quad (5)$$

110 where

$$\varepsilon_{\text{ResN}}^2 = \frac{d^2}{M^2(1-A)^2} + A^{2L}, \quad \varepsilon_{\text{gen}}^2 = \sqrt{\frac{(MdL + d^2)[T + L \log(BL)] + \log(N/\eta)}{n}}. \quad (6)$$

111 Combining Theorem 1 with off-the-shelf results on the DDPM discretization error [Benton et al.,
 112 2023], we obtain the following bound on the sampling error in terms of KL divergence:

113 **Corollary 1.** Let Assumption 1 hold. Consider the two-phase discretization scheme as in Definition
 114 1. Denote the distribution of the output of Algorithm 1 as $\hat{\mu}$. Then, with probability at least $1 - \eta$, we
 115 have

$$\text{KL}(\mu_\delta, \hat{\mu}) / d \lesssim \varepsilon_{\text{score}}^2 + \varepsilon_{\text{disc}}^2, \quad (7)$$

116 where

$$\varepsilon_{\text{score}}^2 \leq \delta^{-1} \cdot \left(\sup_{0 \leq k \leq N-1} \varepsilon_{\text{VI},T-t_k}^2 + \varepsilon_{\text{ResN}}^2 + \varepsilon_{\text{gen}}^2 \right), \quad \varepsilon_{\text{disc}}^2 \leq \kappa^2 N + \kappa T + e^{-2T}. \quad (8)$$

117 Equation (7) provides control on the KL divergence between μ_δ and $\hat{\mu}$ normalized by dimension
 118 d . If the right-hand side is small, this guarantees the two distributions are close in an average per-
 119 coordinate sense: for two d -dimensional product distributions $\mu = \mathcal{N}(0, 1)^{\otimes d}$ and $\nu = \mathcal{N}(\varepsilon, 1)^{\otimes d}$
 120 that are close per coordinate, their KL divergence scales as $\text{KL}(\mu, \nu) \asymp d \cdot \varepsilon^2$, growing linearly with
 121 d . Furthermore, it is possible to derive bounds on the distance between the original distribution μ
 122 (instead of μ_δ) and the learned distribution $\hat{\mu}$ using other DDPM discretization analyses such as
 123 Chen et al. [2022a, 2023a], Li et al. [2023a]. We provide additional discussions of our results in
 124 Appendix A.3. We discuss other related works and future directions in Appendix C.

125 **Generalization to other high-dimensional graphical models** To demonstrate the flexibility of our
 126 proposed framework, we also generalize the results in this section to other high-dimensional graphical
 127 models in Appendix B. Specifically, we consider latent variable Ising models (Appendix B.1), the
 128 conditional Ising models for the conditional generative modeling task (Appendix B.2), and the sparse
 129 coding models (Appendix B.3).

References

- Francis Bach. Breaking the curse of dimensionality with convex neural networks. *The Journal of Machine Learning Research*, 18(1):629–681, 2017.
- Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. *arXiv preprint arXiv:2306.04637*, 2023.
- Jean Barbier, Mohamad Dia, Nicolas Macris, and Florent Krzakala. The mutual information in random linear estimation. In *2016 54th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 625–632. IEEE, 2016.
- Jean Barbier, Florent Krzakala, Nicolas Macris, Léo Miolane, and Lenka Zdeborová. Optimal errors and phase transitions in high-dimensional generalized linear models. *Proceedings of the National Academy of Sciences*, 116(12):5451–5460, 2019.
- Andrew R Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information theory*, 39(3):930–945, 1993.
- Joe Benton, Valentin De Bortoli, Arnaud Doucet, and George Deligiannidis. Linear convergence bounds for diffusion models via stochastic localization. *arXiv preprint arXiv:2308.03686*, 2023.
- David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
- David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- Mark Borgerding, Philip Schniter, and Sundeep Rangan. Amp-inspired deep networks for sparse linear inverse problems. *IEEE Transactions on Signal Processing*, 65(16):4293–4308, 2017.
- Philippe Carmona and Yueyun Hu. Universality in sherrington–kirkpatrick’s spin glass model. In *Annales de l’Institut Henri Poincaré (B) Probability and Statistics*, volume 42, pages 215–222. Elsevier, 2006.
- Michael Celentano. Sudakov-fernique post-amp, and a new proof of the local convexity of the tap free energy. *arXiv preprint arXiv:2208.09550*, 2022.
- Michael Celentano, Zhou Fan, and Song Mei. Local convexity of the tap free energy and amp convergence for z2-synchronization. *arXiv preprint arXiv:2106.11428*, 2021.
- Michael Celentano, Zhou Fan, Licong Lin, and Song Mei. Mean-field variational inference with the tap free energy: Geometric and statistical properties in linear models. In *in preparation*, 2023+.
- Sourav Chatterjee. Spin glasses and stein’s method. *Probability theory and related fields*, 148(3-4): 567–600, 2010.
- Sourav Chatterjee and Amir Dembo. Nonlinear large deviations. *Advances in Mathematics*, 299: 396–450, 2016.
- Hongrui Chen, Holden Lee, and Jianfeng Lu. Improved analysis of score-based generative modeling: User-friendly bounds under minimal smoothness assumptions. In *International Conference on Machine Learning*, pages 4735–4763. PMLR, 2023a.
- Minshuo Chen, Kaixuan Huang, Tuo Zhao, and Mengdi Wang. Score approximation, estimation and distribution recovery of diffusion models on low-dimensional data. *arXiv preprint arXiv:2302.07194*, 2023b.
- Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru R Zhang. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. *arXiv preprint arXiv:2209.11215*, 2022a.
- Sitan Chen, Sinho Chewi, Holden Lee, Yuanzhi Li, Jianfeng Lu, and Adil Salim. The probability flow ode is provably fast. *arXiv preprint arXiv:2305.11798*, 2023c.

- 176 Sitan Chen, Giannis Daras, and Alex Dimakis. Restoration-degradation beyond linear diffusions:
177 A non-asymptotic analysis for ddim-type samplers. In *International Conference on Machine*
178 *Learning*, pages 4462–4484. PMLR, 2023d.
- 179 Xiaohan Chen, Jialin Liu, Zhangyang Wang, and Wotao Yin. Theoretical linear convergence of
180 unfolded ista and its practical weights and thresholds. *Advances in Neural Information Processing*
181 *Systems*, 31, 2018.
- 182 Yongxin Chen, Sinho Chewi, Adil Salim, and Andre Wibisono. Improved analysis for a proximal
183 algorithm for sampling. In *Conference on Learning Theory*, pages 2984–3014. PMLR, 2022b.
- 184 John M Conroy and Dianne P O’leary. Text summarization via hidden markov models. In *Proceed-*
185 *ings of the 24th annual international ACM SIGIR conference on Research and development in*
186 *information retrieval*, pages 406–407, 2001.
- 187 George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control,*
188 *signals and systems*, 2(4):303–314, 1989.
- 189 Ronald DeVore, Guergana Petrova, and Przemyslaw Wojtaszczyk. Approximation of functions of
190 few variables in high dimensions. *Constructive Approximation*, 33(1):125–143, 2011.
- 191 Ronald DeVore, Boris Hanin, and Guergana Petrova. Neural network approximation. *Acta Numerica*,
192 30:327–444, 2021.
- 193 David L Donoho, Arian Maleki, and Andrea Montanari. Message-passing algorithms for compressed
194 sensing. *Proceedings of the National Academy of Sciences*, 106(45):18914–18919, 2009.
- 195 Ahmed El Alaoui, Andrea Montanari, and Mark Sellke. Sampling from the sherrington-kirkpatrick
196 gibbs measure via algorithmic stochastic localization. In *2022 IEEE 63rd Annual Symposium on*
197 *Foundations of Computer Science (FOCS)*, pages 323–334. IEEE, 2022.
- 198 Ronen Eldan. Thin shell implies spectral gap up to polylog via a stochastic localization scheme.
199 *Geometric and Functional Analysis*, 23(2):532–569, 2013.
- 200 Ronen Eldan. Gaussian-width gradient complexity, reverse log-sobolev inequalities and nonlinear
201 large deviations. *Geometric and Functional Analysis*, 28(6):1548–1596, 2018.
- 202 Ronen Eldan. Analysis of high-dimensional distributions using pathwise methods. *Proceedings of*
203 *ICM, to appear*, 2022.
- 204 Ronen Eldan, Frederic Koehler, and Ofer Zeitouni. A spectral condition for spectral gap: fast mixing
205 in high-temperature ising models. *Probability theory and related fields*, 182(3-4):1035–1051,
206 2022.
- 207 Zhou Fan, Song Mei, and Andrea Montanari. Tap free energy, spin glasses and variational inference.
208 2021.
- 209 Zhou Fan, Yufan Li, and Subhabrata Sen. Tap equations for orthogonally invariant spin glasses at
210 high temperature. *arXiv preprint arXiv:2202.09325*, 2022.
- 211 Oliver Y Feng, Ramji Venkataramanan, Cynthia Rush, Richard J Samworth, et al. A unifying tutorial
212 on approximate message passing. *Foundations and Trends® in Machine Learning*, 15(4):335–536,
213 2022.
- 214 Stuart Geman and Christine Graffigne. Markov random field image models and their applications
215 to computer vision. In *Proceedings of the international congress of mathematicians*, volume 1,
216 page 2. Berkeley, CA, 1986.
- 217 Davide Ghio, Yatin Dandi, Florent Krzakala, and Lenka Zdeborová. Sampling with flows, diffusion
218 and autoregressive neural networks: A spin-glass perspective. *arXiv preprint arXiv:2308.14085*,
219 2023.
- 220 Behrooz Ghorbani, Hamid Javadi, and Andrea Montanari. An instability in variational inference for
221 topic models. In *International conference on machine learning*, pages 2221–2231. PMLR, 2019.

- 222 Angeliki Giannou, Shashank Rajput, Jy-yong Sohn, Kangwook Lee, Jason D Lee, and Dimitris Pa-
 223 pailiopoulos. Looped transformers as programmable computers. *arXiv preprint arXiv:2301.13196*,
 224 2023.
- 225 Karol Gregor and Yann LeCun. Learning fast approximations of sparse coding. In *Proceedings of the*
 226 *27th international conference on international conference on machine learning*, pages 399–406,
 227 2010.
- 228 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
 229 *Neural Information Processing Systems*, 33:6840–6851, 2020.
- 230 Kurt Hornik. Some new results on neural network approximation. *Neural networks*, 6(8):1069–1072,
 231 1993.
- 232 Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are
 233 universal approximators. *Neural networks*, 2(5):359–366, 1989.
- 234 Vishesh Jain, Frederic Koehler, and Elchanan Mossel. The mean-field approximation: Information
 235 inequalities, algorithms, and complexity. In *Conference On Learning Theory*, pages 1326–1347.
 236 PMLR, 2018.
- 237 Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduction to
 238 variational methods for graphical models. *Machine learning*, 37:183–233, 1999.
- 239 Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence of score-based generative modeling for general
 240 data distributions. In *International Conference on Algorithmic Learning Theory*, pages 946–985.
 241 PMLR, 2023.
- 242 Marc Lelarge and Leo Miolane. Fundamental limits of symmetric low-rank matrix estimation.
 243 *Probability Theory and Related Fields*, 173:859–929, 2019.
- 244 Gen Li, Yuting Wei, Yuxin Chen, and Yuejie Chi. Towards faster non-asymptotic convergence for
 245 diffusion-based generative models. *arXiv preprint arXiv:2306.09251*, 2023a.
- 246 Yufan Li, Zhou Fan, Subhabrata Sen, and Yihong Wu. Random linear estimation with rotationally-
 247 invariant designs: Asymptotics at high temperature. In *2023 IEEE International Symposium on*
 248 *Information Theory (ISIT)*, pages 156–161. IEEE, 2023b.
- 249 Bingbin Liu, Jordan T Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. Transformers
 250 learn shortcuts to automata. *arXiv preprint arXiv:2210.10749*, 2022a.
- 251 Xingchao Liu, Lemeng Wu, Mao Ye, and Qiang Liu. Let us build bridges: Understanding and
 252 extending diffusion generative models. *arXiv preprint arXiv:2208.14699*, 2022b.
- 253 Yao Ma, Xiaorui Liu, Tong Zhao, Yozen Liu, Jiliang Tang, and Neil Shah. A unified view on
 254 graph neural networks as graph signal denoising. In *Proceedings of the 30th ACM International*
 255 *Conference on Information & Knowledge Management*, pages 1202–1211, 2021.
- 256 Antoine Maillard, Laura Foini, Alejandro Lage Castellanos, Florent Krzakala, Marc Mézard, and
 257 Lenka Zdeborová. High-temperature expansions and message passing algorithms. *Journal of*
 258 *Statistical Mechanics: Theory and Experiment*, 2019(11):113301, 2019.
- 259 Tanya Marwah, Zachary Lipton, and Andrej Risteski. Parametric complexity bounds for approx-
 260 imating pdes with neural networks. *Advances in Neural Information Processing Systems*, 34:
 261 15044–15055, 2021.
- 262 Tanya Marwah, Zachary Chase Lipton, Jianfeng Lu, and Andrej Risteski. Neural network approxi-
 263 mations of pdes beyond linearity: A representational perspective. In *International Conference on*
 264 *Machine Learning*, pages 24139–24172. PMLR, 2023.
- 265 Marc Mezard and Andrea Montanari. *Information, physics, and computation*. Oxford University
 266 Press, 2009.
- 267 Thomas P Minka. Expectation propagation for approximate bayesian inference. *arXiv preprint*
 268 *arXiv:1301.2294*, 2013.

- 269 Andriy Mnih and Geoffrey Hinton. Three new graphical models for statistical language modelling.
270 In *Proceedings of the 24th international conference on Machine learning*, pages 641–648, 2007.
- 271 Vishal Monga, Yuelong Li, and Yonina C Eldar. Algorithm unrolling: Interpretable, efficient deep
272 learning for signal and image processing. *IEEE Signal Processing Magazine*, 38(2):18–44, 2021.
- 273 Andrea Montanari. Sampling, diffusions, and stochastic localization. *arXiv preprint*
274 *arXiv:2305.10690*, 2023.
- 275 Andrea Montanari and Yuchen Wu. Posterior sampling from the spiked models via diffusion processes.
276 *arXiv preprint arXiv:2304.11449*, 2023.
- 277 Sumit Mukherjee and Subhabrata Sen. Variational inference in high-dimensional linear regression.
278 *The Journal of Machine Learning Research*, 23(1):13703–13758, 2022.
- 279 Kazusato Oko, Shunta Akiyama, and Taiji Suzuki. Diffusion models are minimax optimal distribution
280 estimators. *arXiv preprint arXiv:2303.01861*, 2023.
- 281 Vardan Papayan, Yaniv Romano, and Michael Elad. Convolutional neural networks analyzed via
282 convolutional sparse coding. *The Journal of Machine Learning Research*, 18(1):2887–2938, 2017.
- 283 Judea Pearl. Reverend bayes on inference engines: A distributed hierarchical approach. In *Proba-
284 bilistic and Causal Inference: The Works of Judea Pearl*, pages 129–138. 1982.
- 285 Allan Pinkus. Approximation theory of the mlp model in neural networks. *Acta numerica*, 8:143–195,
286 1999.
- 287 Timm Plefka. Convergence condition of the tap equation for the infinite-ranged ising spin glass
288 model. *Journal of Physics A: Mathematical and general*, 15(6):1971, 1982.
- 289 Marc’ Aurelio Ranzato, Alex Krizhevsky, and Geoffrey Hinton. Factored 3-way restricted boltzmann
290 machines for modeling natural images. In *Proceedings of the thirteenth international conference on
291 artificial intelligence and statistics*, pages 621–628. JMLR Workshop and Conference Proceedings,
292 2010.
- 293 Kulin Shah, Sitan Chen, and Adam Klivans. Learning mixtures of gaussians using the ddpm objective.
294 *arXiv preprint arXiv:2307.01178*, 2023.
- 295 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
296 learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*,
297 pages 2256–2265. PMLR, 2015.
- 298 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.
299 *Advances in neural information processing systems*, 32, 2019.
- 300 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
301 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint*
302 *arXiv:2011.13456*, 2020.
- 303 Michel Talagrand. *Spin glasses: a challenge for mathematicians: cavity and mean field models*,
304 volume 46. Springer Science & Business Media, 2003.
- 305 David J Thouless, Philip W Anderson, and Robert G Palmer. Solution of ’solvable model of a spin
306 glass’. *Philosophical Magazine*, 35(3):593–601, 1977.
- 307 Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge
308 Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019. doi:
309 10.1017/9781108627771.
- 310 Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational
311 inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- 312 Colin Wei, Yining Chen, and Tengyu Ma. Statistically meaningful approximation: a case study on
313 approximating turing machines with transformers. *Advances in Neural Information Processing*
314 *Systems*, 35:12071–12083, 2022.

315 E Weinan, Chao Ma, and Lei Wu. Barron spaces and the compositional function spaces for neural
316 network models. *arXiv preprint arXiv:1906.08039*, 2019.

317 Qiang Wu. Thouless-anderson-palmer equations for the multi-species sherrington-kirkpatrick model.
318 *arXiv preprint arXiv:2308.09099*, 2023.

319 Dmitry Yarotsky. Error bounds for approximations with deep relu networks. *Neural Networks*, 94:
320 103–114, 2017.

321 Yaodong Yu, Sam Buchanan, Druv Pai, Tianzhe Chu, Ziyang Wu, Shengbang Tong, Benjamin D
322 Haeffele, and Yi Ma. White-box transformers via sparse rate reduction. *arXiv preprint*
323 *arXiv:2306.01129*, 2023a.

324 Yaodong Yu, Tianzhe Chu, Shengbang Tong, Ziyang Wu, Druv Pai, Sam Buchanan, and Yi Ma. Emer-
325 gence of segmentation with minimalistic white-box transformers. *arXiv preprint arXiv:2308.16271*,
326 2023b.

327 Hui Yuan, Kaixuan Huang, Chengzhuo Ni, Minshuo Chen, and Mengdi Wang. Reward-directed
328 conditional diffusion: Provable distribution estimation and reward improvement. *arXiv preprint*
329 *arXiv:2307.07055*, 2023.

330 Jian Zhang and Bernard Ghanem. Ista-net: Interpretable optimization-inspired deep network for
331 image compressive sensing. In *Proceedings of the IEEE conference on computer vision and pattern*
332 *recognition*, pages 1828–1837, 2018.

333 Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong
334 Du, Chang Huang, and Philip HS Torr. Conditional random fields as recurrent neural networks. In
335 *Proceedings of the IEEE international conference on computer vision*, pages 1529–1537, 2015.

336 Contents

337	1 Introduction	1
338	2 Preliminaries: the DDPM sampling scheme	2
339	3 Diffusion models for Ising models	3
340	A More details for Section 3	10
341	A.1 Discretization scheme	10
342	A.2 Verifying the assumption in examples	10
343	A.3 Discussions	11
344	A.4 Proof outline of Theorem 1	12
345	B Generalization to other high-dimensional graphical models	13
346	B.1 Diffusion models for latent variable Ising models	13
347	B.2 Conditional diffusion models for Ising models	14
348	B.3 Diffusion models for sparse coding	15
349	C Other related work and future directions	16
350	D Technical preliminaries	18
351	D.1 DDPM conditional sampling scheme	18
352	D.2 Sampling error bound of the DDPM scheme	18

353	D.3	Generalization error of empirical risk minimization over ResNets	21
354	D.3.1	Result for Ising models	21
355	D.3.2	Result for Sparse coding	23
356	D.4	Uniform approximation of the denoiser	25
357	D.5	Approximation error of fixed point iteration	26
358	D.6	Properties of two-phase time discretization scheme	26
359	E	Proofs for Section 3: Ising models	27
360	E.1	Proof of Theorem 1	27
361	E.2	Proof of Corollary 1	29
362	E.3	Proofs for Section A.2	29
363	E.3.1	Proof of Lemma 1	29
364	E.3.2	Proof of Lemma 2	30
365	F	Proofs for Section B: Generalization to other models	31
366	F.1	Proof of Theorem 2	31
367	F.2	Proof of Theorem 3	33
368	F.3	Proof of Theorem 4	35
369	F.4	Proof of Lemma 3	38

370 A More details for Section 3

371 A.1 Discretization scheme

372 **Definition 1** (Two-phase discretization scheme [Benton et al., 2023]). *The two-phase discretization*
373 *scheme has parameters $(\kappa, N_0, N, T, \delta) \in (0, 1) \times \mathbb{N} \times \mathbb{N} \times \mathbb{R} \times (0, 1)$, where (κ, N_0, N) are free*
374 *parameters and (T, δ) are fully determined by (κ, N_0, N) . In the first uniform phase, the N_0 time*
375 *steps have equal length κ . In the second exponential phase, the $N - N_0$ steps decay with rate*
376 *$1/(1 + \kappa) \in (0, 1)$. The last time step t_N has a gap $\delta = (1 + \kappa)^{N_0 - N} \in (0, 1)$ to T .*

377 *Specifically, we take $t_0 = 0$, $t_k = k\kappa$ for $k \leq N_0$, $t_{N_0} = N_0\kappa = T - 1$, $t_{N_0+k} = T - (1 + \kappa)^{-k}$*
378 *for $0 \leq k \leq N - N_0$, and $t_N = T - (1 + \kappa)^{N_0 - N} = T - \delta$. Defining $\gamma_k = t_{k+1} - t_k$, we have*
379 *$\gamma_k = \kappa$ for $k \leq N_0 - 1$, and $\gamma_{N_0+k} = \kappa/(1 + \kappa)^{k+1}$ for $0 \leq k \leq N - N_0 - 1$. See [Benton et al.,*
380 *2023, Figure 1] for a pictorial illustration of this scheme.*

381 A.2 Verifying the assumption in examples

382 This section provides examples that admit controlled VI approximation error $\varepsilon_{\text{VI},t}^2$. The results in this
383 section are proved in Appendix E.3.

384 **Ising model in the VB consistency regime.** There is a line of work studying the consistency of
385 the naive mean-field variational Bayes (VB) free energy in Ising models under high-temperature
386 conditions [Chatterjee and Dembo, 2016, Eldan, 2018, Jain et al., 2018, Mukherjee and Sen, 2022].
387 We build on this by providing a quantitative bound on the variational inference approximation error
388 for a general coupling matrix \mathbf{A} in this regime.

389 **Lemma 1.** *Assume $\|\mathbf{A}\|_{\text{op}} < 1/2$. Then for any t , Assumption 1 is satisfied for $\mathbf{K} = \mathbf{0}$, and*

$$\varepsilon_{\text{VI},t}^2(\mathbf{A}) \leq \frac{4}{1 - 2\|\mathbf{A}\|_{\text{op}}} \frac{\|\mathbf{A}\|_F^2}{d}. \quad (9)$$

390 As an example, for the ferromagnetic Ising model we have $\mathbf{A} = \beta \mathbf{1}\mathbf{1}^\top/d$, giving $\varepsilon_{\text{VI},t}^2(\mathbf{A}) \leq$
 391 $[4\beta^2/(1 - 2\beta)]/d$. This shows the VI approximation error vanishes as $\beta < 1/2$ and $d \rightarrow \infty$.
 392 However, this is not a particularly interesting regime for Ising models, since they can be well-
 393 approximated by a product distribution when β is small [Chatterjee and Dembo, 2016, Eldan, 2018].

394 **The Sherrington-Kirkpatrick model.** The Sherrington-Kirkpatrick model assumes $\mathbf{A} = \beta \mathbf{J}$,
 395 where $\mathbf{J} \sim \text{GOE}(d)$ is a symmetric Gaussian random matrix with off-diagonal entries that are
 396 IID Gaussian with variance $1/d$. Prior work has shown that the VB free energy does not provide
 397 consistent estimation in this model [Ghorbani et al., 2019, Fan et al., 2021]. Instead, the variational
 398 objective that yields a consistent estimator of the Gibbs mean is the Thouless-Anderson-Palmer (TAP)
 399 free energy [Thouless et al., 1977, Fan et al., 2021, El Alaoui et al., 2022]. Using results on the TAP
 400 free energy, the variational inference (VI) approximation error can be controlled for this model when
 401 $\beta < 1/4$.

402 **Lemma 2.** [Corollary of Lemma 4.10 of El Alaoui et al. [2022]] Assume $\mathbf{A} = \beta \mathbf{J}$ where $\mathbf{J} \sim$
 403 $\text{GOE}(d)$ and $\beta < 1/4$. Then for any t , there exists matrices $\mathbf{K} = c_t \mathbf{I}_d$ for some c_t , such that with
 404 high probability, $\|\mathbf{A} - c_t \mathbf{I}_d\|_{\text{op}} \leq A < 1$ and

$$\varepsilon_{\text{VI},t}(\beta \mathbf{J}) \xrightarrow{P} 0, \quad \text{as } d \rightarrow \infty.$$

405 Lemma 2 provides a qualitative result on the consistency of variational inference (VI) for the
 406 Sherrington-Kirkpatrick model, but does not give a non-asymptotic error bound. To establish a
 407 non-asymptotic guarantee, one could potentially leverage tools like the smart path method [Talagrand,
 408 2003, Theorem 2.4.20] or Stein’s method [Chatterjee, 2010]. We conjecture it is possible to prove a
 409 quantitative error bound of order $C(\beta)/d$ using these techniques, as illustrated in [Talagrand, 2003,
 410 Theorem 2.4.20].

411 **Other Ising models.** We conjecture that Lemma 2 could extend to a variety of other models
 412 including non-Gaussian Wigner matrices [Carmona and Hu, 2006], heterogeneous variances [Wu,
 413 2023], orthogonally invariant spin glasses [Fan et al., 2022], and spiked matrix models with non-
 414 Rademacher priors [Fan et al., 2021, Lelarge and Miolane, 2019]:

- 415 • *Non-Gaussian Wigner matrices.* We have $\mathbf{A} = \beta \mathbf{J}$ where \mathbf{J} is a symmetric random matrix
 416 whose off-diagonal elements are independent with variance $1/d$ and satisfy some moment
 417 condition. This generalizes GOE matrices to non-Gaussian distributions. Since these
 418 matrices have similar properties to GOE matrices [Carmona and Hu, 2006], we conjecture
 419 Lemma 2 should hold.
- 420 • *Heterogeneous variance: multi-species Sherrington-Kirkpatrick models.* We have $\mathbf{A} = \beta \mathbf{J}$
 421 where \mathbf{J} is a random matrix with independent entries but heterogeneous variance. An
 422 example is the bipartite Sherrington-Kirkpatrick model specified by a set $S \subseteq [d]$, with
 423 $J_{ij} = J_{ji} \sim \mathcal{N}(0, 1/d)$ for $i \in S$ and $j \in S^c$, and $J_{ij} = 0$ for $i, j \in S$ or $i, j \in S^c$. The
 424 TAP equations verifying Assumption 1 has been shown to hold in similar models [Wu, 2023]
 425 in the high-temperature regime $\beta \leq \beta_0$.
- 426 • *Orthogonally invariant spin glass models.* We have $\mathbf{A} = \beta \mathbf{J}$, where $\mathbf{J} = \mathbf{O}\mathbf{E}\mathbf{O}^\top \in$
 427 $\mathbb{R}^{d \times d}$. Here, $\mathbf{O} \sim \text{Haar}(\text{SO}(d))$ is a uniform random orthogonal matrix and $\mathbf{E} =$
 428 $\text{diag}(e_1, \dots, e_d) \in \mathbb{R}^{d \times d}$ is a diagonal matrix. The TAP equations have been shown
 429 for related models [Fan et al., 2022] in the high-temperature regime.
- 430 • *Spiked matrix models.* Suppose we observe $\mathbf{Y} = \mathbf{u}\mathbf{u}^\top + \mathbf{J}$ where $\mathbf{J} \sim \text{GOE}(d)$ and
 431 $\mathbf{u} \in \mathbb{R}^d$ with $u_i \sim_{iid} \pi_0$ for some distribution $\pi_0 \in \mathcal{P}(\mathbb{R})$. The posterior distribution of
 432 \mathbf{u} given observation \mathbf{Y} is given by $\mu(\mathbf{x}) \propto \exp\{\langle \mathbf{x}, \mathbf{Y}\mathbf{x} \rangle / 2 - \|\mathbf{x}\|_2^4 / (4n)\} \pi_0^d(\mathbf{x})$. Taking
 433 this μ as the sample distribution, we conjecture that Assumption 1 can be verified for this
 434 model Fan et al. [2021], Lelarge and Miolane [2019].

435 A.3 Discussions

436 **More explicit sample complexity bounds.** Corollary 1 provides a sampling error bound in
 437 terms of the KL divergence of μ_δ and $\hat{\mu}$. To interpret this bound, assume $\hat{\mu}$ satisfies a
 438 dimension-free transportation-information inequality, i.e., $W_1^2(\mu_\delta, \hat{\mu}) \lesssim \text{KL}(\mu_\delta, \hat{\mu})$. Further assume

439 $\sup_t \varepsilon_{\text{VI},t}^2(\mathbf{A}) \lesssim 1/d$ (conjectured to hold for the SK model when $\beta < 1$). Since $W_1^2(\mu_\delta, \mu)/d \lesssim \delta$,
 440 this implies

$$W_1^2(\mu, \hat{\mu})/d \lesssim W_1^2(\mu, \mu_\delta)/d + \text{KL}(\mu_\delta, \hat{\mu})/d \lesssim \delta + \varepsilon_{\text{score}}^2 + \varepsilon_{\text{disc}}^2.$$

441 By the formulation of $\varepsilon_{\text{score}}^2$ and $\varepsilon_{\text{disc}}^2$ in Eq. (6) and (8) and by $\sup_t \varepsilon_{\text{VI},t}^2(\mathbf{A}) \lesssim 1/d$, to ensure
 442 $W_1^2(\mu, \hat{\mu})/d \lesssim \varepsilon^2$, it suffices to take

$$\begin{aligned} \delta &\asymp \varepsilon^2, & T &\asymp \log(1/\varepsilon), & \kappa &\asymp \varepsilon^2/\log(1/\varepsilon), & N &\asymp \log^2(1/\varepsilon)/\varepsilon^2, \\ d &\asymp 1/\varepsilon^4, & M &\asymp 1/\varepsilon^6, & L &\asymp \log(1/\varepsilon), & n &\asymp \log^3(1/\varepsilon)/\varepsilon^{18}. \end{aligned}$$

443 **The role of dimensionality.** In contrast to existing results [Oko et al., 2023, Chen et al., 2023b]
 444 in which the score estimation error bounds exhibit a curse of dimensionality, our result seems to
 445 demonstrate a ‘‘blessing of dimensionality’’. Specifically, the term $\varepsilon_{\text{VI},t}^2$ in Theorem 1 is independent
 446 of ResNet size, sample size, and will typically vanish as dimension d goes to infinity. However, we
 447 cannot conclude that score estimation actually becomes easier for higher-dimensional Ising models,
 448 since our result only provides an upper bound on the estimation error. Whether score approximation
 449 truly simplifies with increasing dimensions is an open question deserving further investigation.

450 **Generalizing Assumption 1.** While Assumption 1 provides a sufficient condition for efficient
 451 score approximation, it is stronger than necessary. For example, in the Sherrington-Kirkpatrick
 452 model when $\mathbf{A} = \beta \mathbf{J}$ where $\mathbf{J} \sim \text{GOE}(d)$, an efficient sampling algorithm is known when $\beta < 1$
 453 [Celentano, 2022]. However, we can only verify Assumption 1 for $\beta \leq \beta_0$ for some $1/4 < \beta_0 < 1/2$.
 454 Nevertheless, we believe one can weaken our assumption to show score estimation is efficient for any
 455 $\beta < 1$ by leveraging local convexity of the TAP free energy of the SK model, proved in El Alaoui
 456 et al. [2022], Celentano [2022].

457 **The choice of sampling scheme and discretization scheme.** Importantly, our score estimation
 458 error bound in Theorem 1 can combine with sampling schemes beyond DDPM, as it does not rely
 459 on a specific diffusion model. For instance, stochastic localization schemes [Eldan, 2013, El Alaoui
 460 et al., 2022, Montanari and Wu, 2023, Montanari, 2023] estimate the denoiser rather than the score,
 461 and our analysis can be adapted to bound the denoiser estimation error, enabling sampling guarantees
 462 for stochastic localization. Additionally, the discretization scheme and sampling error bound in
 463 Corollary 1 may not be optimal. The analysis could likely be sharpened, or the discretization
 464 improved, to provide tighter error guarantees.

465 A.4 Proof outline of Theorem 1

466 Here, we outline the proof of Theorem 1, with full details in Appendix E.

467 Recall that we have $\hat{s}_t(\mathbf{z}) = \text{P}_t[\text{ResN}_{\widehat{\mathbf{W}}}] (\mathbf{z})$, where $\widehat{\mathbf{W}} = \text{argmin}_{\mathbf{W} \in \mathcal{W}} \widehat{\mathbb{E}}[\|\text{P}_t \text{ResN}_{\mathbf{W}}(\mathbf{z}) +$
 468 $\sigma_t^{-1} \mathbf{g}\|_2^2]$ for $\mathcal{W} = \mathcal{W}_{d,D,L,M,B}$. Here, $\widehat{\mathbb{E}}$ denotes averaging over the empirical data distribution. By
 469 standard error decomposition analysis in empirical risk minimization theory, we have:

$$\begin{aligned} \mathbb{E}[\|\text{P}_t[\text{ResN}_{\widehat{\mathbf{W}}}] (\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d &\leq \inf_{\mathbf{W} \in \mathcal{W}} \mathbb{E}[\|\text{P}_t[\text{ResN}_{\mathbf{W}}] (\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d \\ &+ 2 \sup_{\mathbf{W} \in \mathcal{W}} \left| \widehat{\mathbb{E}}[\|\text{P}_t[\text{ResN}_{\mathbf{W}}] (\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d - \mathbb{E}[\|\text{P}_t[\text{ResN}_{\mathbf{W}}] (\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d \right|. \end{aligned}$$

470 Furthermore, a standard identity in diffusion model theory shows:

$$\mathbb{E}[\|\hat{s}_t(\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2]/d = \mathbb{E}[\|\hat{s}_t(\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d + C, \quad C = \mathbb{E}[\|\mathbf{s}_t(\mathbf{z})\|_2^2]/d - \mathbb{E}[\|\sigma_t^{-1} \mathbf{g}\|_2^2]/d.$$

471 Combining the above yields:

$$\mathbb{E}[\|\hat{s}_t(\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2]/d \leq \bar{\varepsilon}_{\text{app}}^2 + \bar{\varepsilon}_{\text{gen}}^2,$$

472 where $\bar{\varepsilon}_{\text{app}}^2$ is the approximation error and $\bar{\varepsilon}_{\text{gen}}^2$ is the generalization error,

$$\begin{aligned} \bar{\varepsilon}_{\text{app}}^2 &= \inf_{\mathbf{W} \in \mathcal{W}} \mathbb{E}[\|\text{P}_t[\text{ResN}_{\mathbf{W}}] (\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2]/d, \\ \bar{\varepsilon}_{\text{gen}}^2 &= 2 \sup_{\mathbf{W} \in \mathcal{W}} \left| \widehat{\mathbb{E}}[\|\text{P}_t[\text{ResN}_{\mathbf{W}}] (\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d - \mathbb{E}[\|\text{P}_t[\text{ResN}_{\mathbf{W}}] (\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d \right|. \end{aligned}$$

473 The generalization error $\bar{\varepsilon}_{\text{gen}}^2$ can be controlled by a standard empirical process analysis. We simply
 474 use a parameter counting argument to control this term, which can be found in Proposition 6. This
 475 gives rise to the term $\varepsilon_{\text{gen}}^2$ in (6).

476 To control the approximation error $\bar{\varepsilon}_{\text{app}}^2$, we note that $\mathbf{s}_t(\mathbf{z}) = (\lambda_t \cdot \mathbf{m}_t(\mathbf{z}) - \mathbf{z})/\sigma_t^2$, where $\mathbf{m}_t(\mathbf{z}) =$
 477 $\mathbb{E}_{(\mathbf{x}, \mathbf{g}) \sim \mu \otimes \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[\mathbf{x} | \lambda_t \mathbf{x} + \sigma_t \mathbf{g} = \mathbf{z}]$ is the denoiser. Thus, approximating the score function reduces
 478 to approximating $\mathbf{m}_t(\mathbf{z})$ using a ResNet. By Assumption 1, the denoiser \mathbf{m}_t can be approximated
 479 by the minimizer of a variational free energy $\mathcal{F}_t^{\text{VI}}$. This minimizer can be found by a fixed point
 480 iteration, which can further be approximated by a ResNet.

481 More specifically, simple calculus shows that the minimizer $\hat{\mathbf{m}} = \hat{\mathbf{m}}_t$ of the variational free energy
 482 $\mathcal{F}_t^{\text{VI}}$ satisfies the fixed point equation

$$\hat{\mathbf{m}} = \tanh(\mathbf{U}\hat{\mathbf{m}} + \mathbf{h}), \quad \mathbf{U} = \mathbf{A} - \mathbf{K}, \quad \mathbf{h} = \lambda_t \sigma_t^{-2} \mathbf{z}.$$

483 When $\|\mathbf{U}\|_{\text{op}} < 1$, this can be efficiently solved by fixed point iteration

$$\hat{\mathbf{m}} \approx \mathbf{m}^L, \quad \mathbf{m}^{\ell+1} = \tanh(\mathbf{U}\mathbf{m}^\ell + \mathbf{h}), \quad \mathbf{m}^0 = \mathbf{0}.$$

484 This fixed point iteration can further be approximated by the ResNet structure (ResNet), where
 485 \tanh is approximated by a linear combination of ReLU activations. Lemma 5 and 6 analyze this
 486 approximation error $\varepsilon_{\text{ResN}}^2$. Our analysis shows that the total approximation error $\bar{\varepsilon}_{\text{app}}^2$ is controlled
 487 by $\varepsilon_{\text{VI}}^2 + \varepsilon_{\text{ResN}}^2$. Adding the generalization error yields the overall score estimation error bound in
 488 Eq. (5).

489 B Generalization to other high-dimensional graphical models

490 B.1 Diffusion models for latent variable Ising models

491 In the latent variable Ising model μ , we have a coupling matrix $\mathbf{A} = [\mathbf{A}_{11}, \mathbf{A}_{12}; \mathbf{A}_{12}^\top, \mathbf{A}_{22}] \in$
 492 $\mathbb{R}^{(d+m) \times (d+m)}$ (where $\mathbf{A}_{11} \in \mathbb{R}^{d \times d}$, $\mathbf{A}_{12} \in \mathbb{R}^{d \times m}$, and $\mathbf{A}_{22} \in \mathbb{R}^{m \times m}$), specifying a joint distribu-
 493 tion over $(\mathbf{x}, \boldsymbol{\theta}) \in \{\pm 1\}^{d+m}$,

$$\mu(\mathbf{x}, \boldsymbol{\theta}) \propto \exp\{\langle \mathbf{x}, \mathbf{A}_{11} \mathbf{x} \rangle / 2 + \langle \mathbf{x}, \mathbf{A}_{12} \boldsymbol{\theta} \rangle + \langle \boldsymbol{\theta}, \mathbf{A}_{22} \boldsymbol{\theta} \rangle / 2\}, \quad \mathbf{x} \in \{\pm 1\}^d, \boldsymbol{\theta} \in \{\pm 1\}^m. \quad (10)$$

494 Note that the joint distribution over $(\mathbf{x}, \boldsymbol{\theta})$ is still an Ising model. However, here we will treat $\boldsymbol{\theta}$ as a
 495 latent variable and consider generative modeling for the marginal distribution $\mu(\mathbf{x}) = \sum_{\boldsymbol{\theta}} \mu(\mathbf{x}, \boldsymbol{\theta})$
 496 when $\boldsymbol{\theta}$ is unobserved. When $\mathbf{A}_{11} = \mathbf{0}$ and $\mathbf{A}_{22} = \mathbf{0}$, this model reduces to a restricted Boltzmann
 497 machine, which is often used to model natural image distributions [Ranzato et al., 2010].

498 We still consider the generative modeling task where we observe $\{\mathbf{x}_i\}_{i \in [m]} \sim_{\text{iid}} \mu$, and our goal is to
 499 sample a new $\hat{\mathbf{x}} \sim \hat{\mu}$ with $\hat{\mu} \approx \mu$. To show the DDPM scheme (Algorithm 1) provides a controlled
 500 error bound, we need to bound the score estimation error [Benton et al., 2023]. This estimation error
 501 can be controlled if we assume the denoiser minimizes a VI objective.

502 **Assumption 2** (Consistency of the free energy minimizer in marginal Ising models). *Let $\boldsymbol{\sigma} =$*
 503 *$(\mathbf{x}, \boldsymbol{\theta}) \sim \mu(\mathbf{x}, \boldsymbol{\theta}) \propto \exp\{\langle \boldsymbol{\sigma}, \mathbf{A} \boldsymbol{\sigma} \rangle / 2\}$ and $\mathbf{z} \sim \mathcal{N}(\lambda_t \mathbf{x}, \sigma_t^2 \mathbf{I}_d)$. For any fixed t , assume that there*
 504 *exists $\varepsilon_{\text{VI},t}^2(\mathbf{A}) < \infty$ and $\mathbf{K} = \mathbf{K}(\mathbf{A}, t) \in \mathbb{R}^{(d+m) \times (d+m)}$ with $\|\mathbf{K} - \mathbf{A}\|_{\text{op}} \leq A < 1$, such that*

$$\mathbb{E}_{\mathbf{z} \sim \mu_t} [\|\hat{\mathbf{m}}_t(\mathbf{z}) - \mathbf{m}_t(\mathbf{z})\|_2^2] / d \leq \varepsilon_{\text{VI},t}^2(\mathbf{A}), \quad \hat{\mathbf{m}}_t(\mathbf{z}) = [\hat{\boldsymbol{\omega}}_t(\mathbf{z})]_{1:d},$$

$$\hat{\boldsymbol{\omega}}_t(\mathbf{z}) = \operatorname{argmin}_{\boldsymbol{\omega} \in [-1, 1]^{d+m}} \left\{ \sum_{i=1}^{d+m} -\text{h}_{\text{bin}}(\omega_i) - \frac{1}{2} \langle \boldsymbol{\omega}, \mathbf{A} \boldsymbol{\omega} \rangle - \frac{\lambda_t}{\sigma_t^2} \langle \mathbf{z}, \boldsymbol{\omega}_{1:d} \rangle + \frac{1}{2} \langle \boldsymbol{\omega}, \mathbf{K} \boldsymbol{\omega} \rangle \right\}.$$

505 Assumption 2 can be verified in concrete examples. Lemma 1 still applies in this model: when
 506 $\|\mathbf{A}\|_{\text{op}} < 1/2$, taking $\mathbf{K} = \mathbf{0}$ gives $\varepsilon_{\text{VI},t}^2(\mathbf{A}) \leq 4d^{-1}(1 - 2\|\mathbf{A}\|_{\text{op}})^{-1} \|\mathbf{A}\|_F^2$. We conjecture that for
 507 \mathbf{A} being spin glass models like the Sherrington-Kirkpatrick model at high temperature, there exists
 508 \mathbf{K} such that $\mathbb{E}[\varepsilon_{\text{VI},t}^2(\mathbf{A})] \rightarrow 0$ as $d, m \rightarrow \infty$. Given Assumption 2, the following theorem provides a
 509 score estimation error bound and a sampling error bound in latent variable Ising models, proved in
 510 Appendix F.1.

511 **Theorem 2.** Let Assumption 2 hold. Let $\{\hat{\mathbf{s}}_{T-t_k}\}_{0 \leq k \leq N-1}$ be the approximate score function given
 512 by Algorithm 1 in which we take

$$D = 3(d+m), \quad M \geq 4(d+m), \quad B \geq 7 \cdot (M/(d+m)) \cdot \log(M) + \sqrt{d+m} + 1/\min_k\{T-t_k\}.$$

513 Then with probability at least $1 - \eta$, for any $t \in \{T-t_k\}_{0 \leq k \leq N-1}$, we have

$$\mathbb{E}_{\mathbf{z} \sim \mu_t} [\|\hat{\mathbf{s}}_t(\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2]/d \lesssim \lambda_t^2 \sigma_t^{-4} \cdot \left(\varepsilon_{\text{VI},t}^2(\mathbf{A}) + \varepsilon_{\text{ResN}}^2 + \varepsilon_{\text{gen}}^2 \right), \quad (11)$$

514 where $\varepsilon_{\text{VI},t}^2$ is given in Assumption 2, and

$$\begin{aligned} \varepsilon_{\text{ResN}}^2 &= \frac{d+m}{d} \left(\frac{(d+m)^2}{M^2(1-A)^2} + A^{2L} \right), \\ \varepsilon_{\text{gen}}^2 &= \sqrt{\frac{(ML+d)(d+m)[T+L \log(BL)] + \log(N/\eta)}{n}}. \end{aligned} \quad (12)$$

515 Furthermore, consider the two-phase discretization scheme as in Definition 1, we have with probability
 516 $1 - \eta$ that

$$\text{KL}(\mu_\delta, \hat{\mu})/d \lesssim \delta^{-1} \cdot \left(\sup_{0 \leq k \leq N-1} \varepsilon_{\text{VI},T-t_k}^2 + \varepsilon_{\text{ResN}}^2 + \varepsilon_{\text{gen}}^2 \right) + \kappa^2 N + \kappa T + e^{-2T}. \quad (13)$$

517 B.2 Conditional diffusion models for Ising models

518 In the conditional Ising model, we also have a coupling matrix $\mathbf{A} = [\mathbf{A}_{11}, \mathbf{A}_{12}; \mathbf{A}_{12}^\top, \mathbf{A}_{22}] \in$
 519 $\mathbb{R}^{(d+m) \times (d+m)}$, specifying a joint distribution over $(\mathbf{x}, \boldsymbol{\theta}) \in \{\pm 1\}^{d+m}$ as in Eq. (10). However, we
 520 now consider the conditional generative modeling task where we observe $\{(\mathbf{x}_i, \boldsymbol{\theta}_i)\}_{i \in [n]} \sim_{\text{iid}} \mu$.
 521 The goal is to sample $\hat{\mathbf{x}} \sim \hat{\mu}(\cdot | \boldsymbol{\theta}) \approx \mu(\cdot | \boldsymbol{\theta})$ for a given $\boldsymbol{\theta}$. Such problems naturally arise in image
 522 imputation tasks, where $(\mathbf{x}, \boldsymbol{\theta})$ represents a full image, $\boldsymbol{\theta}$ is the observed part, and \mathbf{x} is the missing
 523 part to impute.

524 The conditional generative modeling task can be solved using the conditional DDPM scheme (Al-
 525 gorithm 2 as described in Appendix D.1). To bound the error, we need to control the estima-
 526 tion error of the conditional score $\mathbf{s}_t(\mathbf{z}; \boldsymbol{\theta}) = \nabla_{\mathbf{z}} \log \mu_t(\mathbf{z}, \boldsymbol{\theta})$. By Tweedie’s formula, we have
 527 $\mathbf{s}_t(\mathbf{z}; \boldsymbol{\theta}) = (\lambda_t \mathbf{m}_t(\mathbf{z}; \boldsymbol{\theta}) - \mathbf{z})/\sigma_t^2$, where $\mathbf{m}_t(\mathbf{z}; \boldsymbol{\theta}) := \mathbb{E}_{(\mathbf{x}, \boldsymbol{g}) \sim \mu \otimes \mathcal{N}(0,1)}[\mathbf{x} | \boldsymbol{\theta}, \mathbf{z} = \lambda_t \mathbf{x} + \sigma_t \boldsymbol{g}]$ is
 528 the conditional denoiser. We assume the following about $\mathbf{m}_t(\mathbf{z}; \boldsymbol{\theta})$.

529 **Assumption 3** (Consistency of the free energy minimizer in conditional Ising models). Let $(\mathbf{x}, \boldsymbol{\theta}) \sim$
 530 $\mu(\mathbf{x}, \boldsymbol{\theta}) \propto \exp\{\langle \boldsymbol{\sigma}, \mathbf{A}\boldsymbol{\sigma} \rangle/2\}$ and $\mathbf{z} \sim \mathcal{N}(\lambda_t \mathbf{x}, \sigma_t^2 \mathbf{I}_d)$. For any fixed t , assume that there exists
 531 $\varepsilon_{\text{VI},t}^2(\mathbf{A}) < \infty$ and $\mathbf{K} = \mathbf{K}(\mathbf{A}, t) \in \mathbb{R}^{d \times d}$ with $\|\mathbf{K} - \mathbf{A}_{11}\|_{\text{op}} \leq A < 1$, such that

$$\begin{aligned} \mathbb{E}_{(\boldsymbol{\theta}, \mathbf{z})} [\|\hat{\mathbf{m}}_t(\mathbf{z}; \boldsymbol{\theta}) - \mathbf{m}_t(\mathbf{z}; \boldsymbol{\theta})\|_2^2]/d &\leq \varepsilon_{\text{VI},t}^2(\mathbf{A}), \\ \hat{\mathbf{m}}_t(\mathbf{z}; \boldsymbol{\theta}) &= \operatorname{argmin}_{\mathbf{m} \in [-1,1]^d} \left\{ \sum_{i=1}^d -\text{h}_{\text{bin}}(m_i) - \frac{1}{2} \langle \mathbf{m}, \mathbf{A}_{11} \mathbf{m} \rangle - \langle \mathbf{m}, \mathbf{A}_{12} \boldsymbol{\theta} \rangle - \frac{\lambda_t}{\sigma_t^2} \langle \mathbf{z}, \mathbf{m} \rangle + \frac{1}{2} \langle \mathbf{m}, \mathbf{K} \mathbf{m} \rangle \right\}. \end{aligned}$$

532 Assumption 3 can be verified in concrete examples. Lemma 1 still applies in this model: when
 533 $\|\mathbf{A}_{11}\|_{\text{op}} < 1/2$, taking $\mathbf{K} = \mathbf{0}$ gives $\varepsilon_{\text{VI},t}^2(\mathbf{A}) \leq 4d^{-1}(1-2\|\mathbf{A}_{11}\|_{\text{op}})^{-1}\|\mathbf{A}_{11}\|_F^2$. We conjecture
 534 that $\mathbb{E}[\varepsilon_{\text{VI},t}^2(\mathbf{A})] \rightarrow 0$ as $d, m \rightarrow \infty$ for \mathbf{A} being spin glass models at high temperature. Given
 535 Assumption 3, the following theorem provides a conditional score estimation error bound and a
 536 conditional sampling error bound in conditional Ising models, proved in Appendix F.2.

537 **Theorem 3.** Let Assumption 3 hold. Let $\{\hat{\mathbf{s}}_{T-t_k}\}_{0 \leq k \leq N-1}$ be the approximate score function given
 538 by Algorithm 2 in which we take

$$D = 4d, \quad M \geq 4d, \quad B \geq 7 \cdot (M/d) \cdot \log(M) + \sqrt{d} + 1/\min_k\{T-t_k\} + \|\mathbf{A}_{12}\|_{\text{op}} \cdot (M/d + 1).$$

539 Then with probability at least $1 - \eta$, for any $t \in \{T-t_k\}_{0 \leq k \leq N-1}$, we have

$$\mathbb{E}_{(\boldsymbol{\theta}, \mathbf{z})} [\|\hat{\mathbf{s}}_t(\mathbf{z}; \boldsymbol{\theta}) - \mathbf{s}_t(\mathbf{z}; \boldsymbol{\theta})\|_2^2]/d \lesssim \lambda_t^2 \sigma_t^{-4} \cdot \left(\varepsilon_{\text{VI},t}^2(\mathbf{A}) + \varepsilon_{\text{ResN}}^2 + \varepsilon_{\text{gen}}^2 \right),$$

540 where $\varepsilon_{\text{VI},t}^2$ is given in Assumption 3, and

$$\begin{aligned}\varepsilon_{\text{ResN}}^2 &= \frac{d^2}{M^2(1-A)^2} + A^{2L}, \\ \varepsilon_{\text{gen}}^2 &= \sqrt{\frac{(MdL + d(d+m))[T + L \log(BLd^{-1}(m+d))] + \log(N/\eta)}{n}}.\end{aligned}\tag{14}$$

541 Furthermore, consider the two-phase discretization scheme as in Definition 1, we have with probability
542 $1 - \eta$ that

$$\mathbb{E}_{\boldsymbol{\theta} \sim \mu}[\text{KL}(\mu_\delta(\cdot|\boldsymbol{\theta}), \hat{\mu}(\cdot|\boldsymbol{\theta}))/d] \lesssim \delta^{-1} \cdot \left(\sup_{0 \leq k \leq N-1} \varepsilon_{\text{VI},T-t_k}^2 + \varepsilon_{\text{ResN}}^2 + \varepsilon_{\text{gen}}^2 \right) + \kappa^2 N + \kappa T + e^{-2T}.$$

543 We note the score estimation and sampling error bounds in Theorem 3 are averaged over $\boldsymbol{\theta} \sim \mu(\boldsymbol{\theta}) =$
544 $\sum_{\mathbf{x} \in \{\pm 1\}^d} \mu(\mathbf{x}, \boldsymbol{\theta})$, the marginal of $\boldsymbol{\theta}$. These do not ensure error bounds for any fixed $\boldsymbol{\theta}$.

545 B.3 Diffusion models for sparse coding

546 In sparse coding, there is a fixed dictionary $\mathbf{A} \in \mathbb{R}^{d \times m}$. Our observations are noisy, sparse linear
547 combinations of the columns of the dictionary: $\mathbf{x}_i = \mathbf{A}\boldsymbol{\theta}_i + \boldsymbol{\varepsilon}_i$ for $i \in [n]$. Here $\boldsymbol{\varepsilon}_i \sim_{\text{iid}} \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I}_d)$
548 are noise vectors, and $\boldsymbol{\theta}_i \sim_{\text{iid}} \pi_0^{\otimes m}$ are sparse coefficient vectors, with $\pi_0 \in \mathcal{P}(\mathbb{R})$ having a Dirac
549 delta mass at 0. Given observations $\{\mathbf{x}_i\}_{i \in [n]}$, sparse coding typically aims to recover \mathbf{A} and estimate
550 $\{\boldsymbol{\theta}_i\}_{i \in [n]}$. Instead, we consider the generative modeling problem — learning a model to generate
551 new samples $\hat{\mathbf{x}}$ resembling the observations $\{\mathbf{x}_i\}_{i \in [n]}$.

552 The generative modeling task for sparse coding can be solved by the DDPM sampling scheme
553 (Algorithm 1). To control the score estimation error, we make the following assumption on the
554 following denoising function e_t , which requires a little modification in the sparse coding setting:

$$e_t(\mathbf{z}_*) := \mathbb{E}_{(\mathbf{z}_*, \boldsymbol{\theta})}[\boldsymbol{\theta} | \mathbf{z}_*], \quad \mathbf{z}_* = \mathbf{A}\boldsymbol{\theta} + \bar{\boldsymbol{\varepsilon}}, \quad \bar{\boldsymbol{\varepsilon}}_j \sim_{\text{iid}} \mathcal{N}(0, \tau^2 + \sigma_t^2/\lambda_t^2).\tag{15}$$

555 **Assumption 4** (Consistency of the free energy minimizer in sparse coding). Fix $\mathbf{A} \in \mathbb{R}^{d \times m}$.
556 Consider the Bayesian linear model $\mathbf{z}_* = \mathbf{A}\boldsymbol{\theta} + \bar{\boldsymbol{\varepsilon}} \in \mathbb{R}^d$, $\bar{\boldsymbol{\varepsilon}}_j \sim_{\text{iid}} \mathcal{N}(0, \bar{\tau}_t^2)$ where $\bar{\tau}_t^2 = \tau^2 + \sigma_t^2/\lambda_t^2$
557 and $\theta_i \sim_{\text{iid}} \pi_0$ where $\pi_0 \in \mathcal{P}([-\Pi, \Pi])$. Assume that for any $t > 0$, there exist $(\nu_t, \mathbf{K}_t, \varepsilon_{\text{VI},t}^2)$ that
558 depend on $(\pi_0, \mathbf{A}, \tau, t)$ with $\|\mathbf{A}^\top \mathbf{A}/\bar{\tau}_t^2 - \mathbf{K}_t\|_{\text{op}} \leq A < 1/\Pi^2$, such that

$$\begin{aligned}\mathbb{E}_{\mathbf{z}_* \sim \mu_t}[\|\hat{e}_t(\mathbf{z}_*) - e_t(\mathbf{z}_*)\|_2^2]/m &\leq \varepsilon_{\text{VI},t}^2(\mathbf{A}), \\ \hat{e}_t(\mathbf{z}_*) &= \operatorname{argmin}_{\mathbf{e} \in [-\Pi, \Pi]^m} \left\{ \sum_{i=1}^m \max_{\lambda} \left[\lambda m_i - \log \mathbb{E}_{\beta \sim \pi_0} [e^{\lambda \beta - \beta^2 \nu_t / 2}] \right] + \frac{1}{2\bar{\tau}_t^2} \|\mathbf{z}_* - \mathbf{A}\mathbf{e}\|_2^2 - \frac{1}{2} \langle \mathbf{e}, \mathbf{K}_t \mathbf{e} \rangle \right\}.\end{aligned}$$

We also use a different truncation operator in Algorithm 1, replacing P_t by \bar{P}_t :

$$\bar{P}_t[f](\mathbf{z}) = \operatorname{proj}_{\sqrt{m} \|\mathbf{A}\|_{\text{op}} \Pi \lambda_t (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1}}(f(\mathbf{z}) + (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \mathbf{z}) - (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \mathbf{z}.$$

559 Given Assumption 4, the following theorem provides a score estimation error bound in sparse coding
560 models, proved in Appendix F.3.

561 **Theorem 4.** Let Assumption 4 hold. Let $\{\hat{\mathbf{s}}_{T-t_k}\}_{0 \leq k \leq N-1}$ be the approximate score function given
562 by Algorithm 1 in which we take

$$\begin{aligned}D &= 3m + d, \quad M \geq 4m, \\ B &\geq (M/m) \cdot (A + 1 + 2\Pi^2 + w_\star) + 2\Pi + 6 + (\|\mathbf{A}\|_{\text{op}} + 1) / \min_k \{T - t_k\} + \tau^{-2} \|\mathbf{A}\|_{\text{op}} + \sqrt{m},\end{aligned}$$

563 where w_\star is defined in Eq. (65). Then with probability at least $1 - \eta$, when $n \geq \log(2/\eta)$, for any
564 $t \in \{T - t_k\}_{0 \leq k \leq N-1}$, we have the following score estimation error bound

$$\begin{aligned}\mathbb{E}_{(\boldsymbol{\theta}, \mathbf{z})}[\|\hat{\mathbf{s}}_t(\mathbf{z}; \boldsymbol{\theta}) - \mathbf{s}_t(\mathbf{z}; \boldsymbol{\theta})\|_2^2]/d \\ \lesssim \lambda_t^2 \|\mathbf{A}\|_{\text{op}}^2 (1 + \tau^{-4}) \cdot \frac{m}{d} \cdot \left(\varepsilon_{\text{VI},t}^2(\mathbf{A}) + \varepsilon_{\text{ResN}}^2 \right) + \left(\lambda_t^2 \|\mathbf{A}\|_{\text{op}}^2 (1 + \tau^{-4}) \Pi^2 \cdot \frac{m}{d} + \frac{\lambda_t^2}{\sigma_t^2} (1 + \tau^2) \right) \varepsilon_{\text{gen}}^2,\end{aligned}$$

565 for $\varepsilon_{\text{VI},t}^2$ as given in Assumption 4, and

$$\varepsilon_{\text{ResN}}^2 = \Pi^2 \cdot (\Pi^2 A)^{2L} + \frac{m^2 \Pi^2}{(1 - \Pi^2 A)^2 M^2}, \quad \varepsilon_{\text{gen}}^2 = \sqrt{\frac{(dD + LDM) \cdot (T + L) \cdot \iota}{n}}. \quad (16)$$

566 where $\iota = \log(LBnmT(1 + \tau)(1 + \|\mathbf{A}\|_{\text{op}}\Pi)\tau^{-1}N\eta^{-1})$.

567 Theorem 4 can be further combined with an off-the-shelf discretization bound as in Theorem 5 to
568 derive a sampling error bound.

569 **Verifying the assumption.** The VI approximation error $\varepsilon_{\text{VI}}^2$ in Assumption 4 converges to 0 as
570 $d, m \rightarrow \infty$ when \mathbf{A} is a rotationally invariant design matrix, by choosing the variational objective
571 to be the TAP free energy [Thouless et al., 1977]. Specifically, assume the SVD decomposition
572 $\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{O}^\top$ where $\mathbf{Q} \in \mathbb{R}^{d \times d}$ and $\mathbf{O} \in \mathbb{R}^{m \times m}$ are orthonormal, and $\mathbf{D} \in \mathbb{R}^{d \times m}$ is diagonal.
573 Assume that $\mathbf{O} \sim \text{Haar}(\text{SO}(m))$ is independent of everything else, and the diagonal elements of
574 \mathbf{D} have certain empirical distribution converging to a bounded distribution D . As an example, \mathbf{A}
575 with IID Gaussian entries of variance $1/m$ is rotationally invariant. Under the assumption that \mathbf{A} is
576 rotationally invariant, a corollary of [Li et al., 2023b, Theorem 1.11] gives the following lemma, with
577 proof contained in Appendix F.4.

578 **Lemma 3** (Corollary of Li et al. [2023b] Theorem 1.11). *Let $\mathbf{A} \in \mathbb{R}^{d \times m}$ be a rotationally invariant
579 design matrix and let Assumption 7 hold. Then for any $\pi_0, \alpha = d/m$, and limiting distribution D ,
580 there exists $\tau^2 > 0$, such that for any t , there exists matrices $\mathbf{K} = c_t \mathbf{I}_d$ for some c_t , such that*

$$\varepsilon_{\text{VI},t}(\mathbf{A}) \xrightarrow{a.s.} 0, \quad d, m \rightarrow \infty, \quad d/m \rightarrow \alpha.$$

581 Although Lemma 3 does not provide non-asymptotic control of the VI approximation error, we
582 believe this could be obtained through more refined analysis.

583 C Other related work and future directions

584 **Score function approximation in diffusion models.** Neural network-based score function ap-
585 proximation has been recently studied in Oko et al. [2023], Chen et al. [2023b], Yuan et al. [2023],
586 Shah et al. [2023]. Oko et al. [2023] assumes that the data distribution $\mu \in \mathcal{P}(\mathbb{R}^d)$ has a density
587 with s -order bounded derivatives and shows that estimating the score to precision ε requires network
588 size and sample complexity at least $\varepsilon^{-d/s}$. This suffers from the curse of dimensionality unless the
589 data distribution is very smooth ($s \asymp d$). Oko et al. [2023], Chen et al. [2023b] avoid the curse
590 of dimensionality by assuming that the data distribution has a low-dimensional structure, but this
591 assumption does not apply to high-dimensional graphical models. Shah et al. [2023] considers
592 Gaussian mixture models where the score function has a closed form, enabling parameterized by a
593 small shallow network.

594 In contrast, we assume the data distribution is a graphical model, common for images and text [Blei
595 et al., 2003, Mnih and Hinton, 2007, Geman and Graffigne, 1986]. Assuming the efficiency of
596 variational inference approximation, we show that the score can be well-approximated by a network
597 polynomial in dimension, enabling efficient learning from polynomial samples. Our graphical model
598 assumption and algorithm unrolling of variational inference perspective circumvent dimensionality
599 issues faced by prior work.

600 **Discretizing the diffusion process.** Recent work has studied the convergence rates of the discretized
601 reverse SDEs/ODEs for diffusion models [Liu et al., 2022b, Li et al., 2023a, Lee et al., 2023, Chen
602 et al., 2022b, 2023d, 2022a, 2023c,a, Benton et al., 2023]. In particular, Chen et al. [2023a], Benton
603 et al. [2023] provide minimal assumptions to quantitatively control the KL divergence between the
604 perturbed and data distributions. These assumptions include the second moment bound and the
605 controlled score estimation error. Our work focuses on controlling the score estimation error, a goal
606 that is orthogonal to analyzing discretization schemes. Specifically, we directly leverage the result of
607 Benton et al. [2023] to provide an end-to-end error bound.

608 **Stochastic localization.** Stochastic localization, proposed by Eldan [2013, 2022], is another sam-
609 pling scheme similar to diffusion models. Recent works have developed algorithmic sampling

610 techniques based on stochastic localization [El Alaoui et al., 2022, Montanari and Wu, 2023, Ce-
611 lentano, 2022]. Montanari [2023] shows the equivalence of stochastic localization to the DDPM
612 sampling scheme in the Gaussian setting and proposes various ways of generalizing stochastic local-
613 ization schemes. While we present our results in the diffusion model framework, our methods can
614 also provide sampling error bound for stochastic localization schemes.

615 **Neural network approximation theory.** Classical neural network approximation theory typically
616 relies on assumptions that the target function is smooth or hierarchically smooth [Cybenko, 1989,
617 Hornik et al., 1989, Hornik, 1993, Pinkus, 1999, DeVore et al., 2011, Weinan et al., 2019, Yarotsky,
618 2017, Barron, 1993, Bach, 2017, DeVore et al., 2021]. These enable overcoming the curse of
619 dimensionality for higher-order smooth or low-dimensional target functions [Barron, 1993, Weinan
620 et al., 2019, Bach, 2017]. However, when applying them to score function approximation in diffusion
621 models, it is unclear whether such assumptions hold for the score function of high-dimensional
622 graphical models.

623 A recent line of work investigated the expressiveness of neural networks through an algorithm
624 approximation viewpoint [Wei et al., 2022, Bai et al., 2023, Giannou et al., 2023, Liu et al., 2022a,
625 Marwah et al., 2021, 2023]. Wei et al. [2022], Bai et al. [2023], Giannou et al. [2023], Liu et al.
626 [2022a] show that transformers can efficiently approximate several algorithm classes, such as gradient
627 descent and Turing machines. Marwah et al. [2021, 2023] demonstrate that deep networks can
628 efficiently approximate PDE solutions by approximating the gradient dynamics. We also adopt this
629 algorithmic perspective for neural network approximation but apply it to score function approximation
630 for diffusion models.

631 **Variational inference in graphical models.** Variational inference is commonly used to approximate
632 the marginal statistics of graphical models [Pearl, 1982, Jordan et al., 1999, Minka, 2013, Mezard
633 and Montanari, 2009, Wainwright et al., 2008, Blei et al., 2017]. In certain regimes, such as graphical
634 models in the high temperature, naive variational Bayes has been shown to yield consistent posterior
635 estimates [Chatterjee and Dembo, 2016, Eldan, 2018, Jain et al., 2018, Mukherjee and Sen, 2022].
636 For high dimensional statistical models in the low signal-to-noise ratio regime, approximate message
637 passing [Donoho et al., 2009, Feng et al., 2022] and equivalently TAP variational inference [Thouless
638 et al., 1977, Ghorbani et al., 2019, Fan et al., 2021, Celentano et al., 2021, Celentano, 2022, Celentano
639 et al., 2023+], can achieve consistent estimation of the Bayes posterior. Our paper directly adopts
640 results developed for variational inference methods in spin glass models and Bayesian linear models
641 [Talagrand, 2003, Chatterjee, 2010, Barbier et al., 2019, 2016, Fan et al., 2021, 2022, Li et al., 2023b,
642 Celentano et al., 2021, Celentano, 2022, Celentano et al., 2023+].

643 **Algorithm unrolling.** A line of work has focused on neural network denoising by unrolling iterative
644 denoising algorithms into deep networks [Gregor and LeCun, 2010, Zheng et al., 2015, Zhang and
645 Ghanem, 2018, Pappayan et al., 2017, Ma et al., 2021, Chen et al., 2018, Borgerding et al., 2017, Monga
646 et al., 2021, Yu et al., 2023a,b]. These approaches include unrolling ISTA for LASSO into recurrent
647 nets [Gregor and LeCun, 2010, Zhang and Ghanem, 2018, Pappayan et al., 2017, Borgerding et al.,
648 2017], unrolling belief propagation for Markov random fields into recurrent nets [Zheng et al., 2015],
649 and unrolling graph denoising algorithms into graph neural nets [Ma et al., 2021]. Our work also
650 adopts this algorithm unrolling viewpoint, but with a different goal: while the prior literature has
651 mainly focused on devising better denoising algorithms, our work uses this perspective to provide
652 neural network approximation theories for diffusion-based generative models.

653 **Algorithmic hard phase.** The algorithm unrolling perspective can also shed light on the failure
654 mode of score approximation, namely when score functions cannot be efficiently represented by
655 neural networks. For example, we can conclude that the score function of the Sherrington-Kirkpatrick
656 model with $\beta > 1$ cannot be efficiently represented by a neural network, as it was proven in El Alaoui
657 et al. [2022] that there is no stable algorithm to sample the SK Gibbs measure for $\beta > 1$. More
658 generally, the relationship between hardness of sampling, hardness of diffusion-based sampling, and
659 hardness of score approximation deserves further investigation. Recent work such as Ghio et al.
660 [2023] provides a valuable discussion on this important topic.

661 **Future directions.** Our work leaves open several interesting questions. One issue is that for fixed
662 dimension d , our score approximation error does not decay as the network size and sample size

663 increase, and is lower bounded by the variational inference approximation error $\varepsilon_{\text{VI}}^2$. To resolve
 664 this, one approach could consider a hierarchy of variational inference algorithms, such as Plefka’s
 665 expansion [Plefka, 1982, Maillard et al., 2019], which provide increasingly accurate approximations.
 666 Using these hierarchical approximations within our framework could potentially reduce the score
 667 approximation error.

668 Another open question is understanding the algorithms that diffusion neural networks like U-nets
 669 and transformers implement in diffusion models for image tasks. One hypothesis is that U-nets
 670 with convolution layers are implementing some form of variational inference denoising on graphical
 671 models with certain locality and invariance structures. It would be interesting to test this hypothesis
 672 on real image datasets.

673 Finally, an exciting direction is leveraging the algorithmic unrolling perspective to design improved
 674 neural network architectures for diffusion models. The resulting architectures could potentially be
 675 more interpretable and achieve better emergent capabilities, as illustrated by recent works like Yu
 676 et al. [2023a,b].

677 D Technical preliminaries

678 D.1 DDPM conditional sampling scheme

679 We provide the details of the DDPM conditional sampling scheme (Algorithm 2) as mentioned in
 680 Section 2. The algorithm still has two steps, with minor modifications from unconditional DDPM
 681 (Algorithm 1). In the first step, empirical risk minimization (Eq. (18)) fits manually-generated noises
 682 $\{\mathbf{g}_i\}_{i \in [n]}$ using the noisy samples and conditioning variables $\{(\lambda_t \mathbf{x}_i + \sigma_t \mathbf{g}_i; \boldsymbol{\theta}_i)\}_{i \in [n]}$. The ResNet
 683 $\text{ResN}_{\mathbf{W}} : \mathbb{R}^d \times \mathbb{R}^m \rightarrow \mathbb{R}^d$ is parameterized by $\mathbf{W} = \{\mathbf{W}_1^{(\ell)} \in \mathbb{R}^{D \times M}, \mathbf{W}_2^{(\ell)} \in \mathbb{R}^{M \times D}\}_{\ell \in [L]} \cup$
 684 $\{\mathbf{W}_{\text{in}} \in \mathbb{R}^{(d+m+1) \times D}, \mathbf{W}_{\text{out}} \in \mathbb{R}^{D \times d}\}$ and is defined iteratively as

$$\text{ResN}_{\mathbf{W}}(\mathbf{z}, \boldsymbol{\theta}) = \mathbf{W}_{\text{out}} \mathbf{u}^{(L)}, \quad \mathbf{u}^{(\ell)} = \mathbf{u}^{(\ell-1)} + \mathbf{W}_1^{(\ell)} \text{ReLU}(\mathbf{W}_2^{(\ell)} \mathbf{u}^{(\ell-1)}), \quad \mathbf{u}^{(0)} = \mathbf{W}_{\text{in}}[\mathbf{z}; \boldsymbol{\theta}; 1].$$

(ResNet-Conditional)

685 The only difference between (ResNet) and (ResNet-Conditional) is the input dimension. Minimization
 686 is over the ResNets with weights in the set (for parameters d, m, D, L, M, B):

$$\mathcal{W}_{d,m,D,L,M,B} := \left\{ \mathbf{W} = \{\mathbf{W}_1^{(\ell)}, \mathbf{W}_2^{(\ell)}\}_{\ell \in [L]} \cup \{\mathbf{W}_{\text{in}}, \mathbf{W}_{\text{out}}\} : \|\mathbf{W}\| \leq B \right\},$$

$$\|\mathbf{W}\| := \max_{\ell \in [L]} \{\|\mathbf{W}_1^{(\ell)}\|_{\text{op}} + \|\mathbf{W}_2^{(\ell)}\|_{\text{op}}\} \vee \max\{\|\mathbf{W}_{\text{in}}\|_{\text{op}}, \|\mathbf{W}_{\text{out}}\|_{\text{op}}\}.$$

(17)

687 We still truncate the ResNet output using P_t : for $f : \mathbb{R}^d \times \mathbb{R}^m \rightarrow \mathbb{R}^d$, we define $P_t[f](\mathbf{z}, \boldsymbol{\theta}) =$
 688 $\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}(f(\mathbf{z}, \boldsymbol{\theta}) + \sigma_t^{-2} \mathbf{z}) - \sigma_t^{-2} \mathbf{z}}$, where proj_R projects $\mathbf{z} \in \mathbb{R}^d$ into the R -Euclidean ball.

689 The second step of Algorithm 2 still discretizes the backward SDE through the exponential integrator
 690 scheme (19) and the two-phase discretization scheme (Definition 1). However, we replace the score
 691 function $\hat{\mathbf{s}}_t(\hat{\mathbf{Y}}_k)$ with the conditional score function $\hat{\mathbf{s}}_t(\hat{\mathbf{Y}}_k; \boldsymbol{\theta}) = P_t[\text{ResN}_{\hat{\mathbf{W}}_t}](\hat{\mathbf{Y}}_k, \boldsymbol{\theta})$.

692 D.2 Sampling error bound of the DDPM scheme

693 In this section, we state a result from Benton et al. [2023], which establishes the convergence of
 694 the DDPM discretization scheme, when evaluated using Kullback-Leibler (KL) divergence, with
 695 only minimal assumptions required. A slight generalization of the result in Benton et al. [2023] is
 696 necessary, generalizing the identity covariance assumption to a general covariance matrix. The proof
 697 requires little modification, but we present a proof sketch here for completeness.

698 Suppose we are interested in drawing samples from μ in \mathbb{R}^d . The forward process that evolves
 699 according to the Ornstein-Uhlenbeck (OU) process is defined as the following SDE:

$$dX_t = -X_t dt + \sqrt{2} dB_t, \quad X_0 \sim \mu, \quad 0 \leq t \leq T. \quad (20)$$

700 In the above display, $(B_t)_{0 \leq t \leq T}$ is a standard Brownian motion in \mathbb{R}^d . We denote by μ_t the
 701 distribution of X_t . One can check that $X_t \stackrel{d}{=} e^{-t} X_0 + \sqrt{1 - e^{-2t}} \mathbf{g}$ for $\mathbf{g} \sim \mathcal{N}(\cdot, \mathbf{I}_d)$ that is

Algorithm 2 The DDPM conditional sampling scheme

Require: Samples $\{(\mathbf{x}_i, \boldsymbol{\theta}_i)\}_{i \in [n]} \subseteq \mathbb{R}^d \times \mathbb{R}^m$. Conditional latent variable $\boldsymbol{\theta}$. ResNet parameters (d, m, D, L, M, B) . Discretization scheme parameters $(N, T, \delta, \{t_k\}_{0 \leq k \leq N})$ with $0 = t_0 < \dots < t_N = T - \delta$. Denote $\gamma_k = t_{k+1} - t_k$.

1: // Computing the approximate conditional score function

2: Sample $\{\mathbf{g}_i\}_{i \in [n]} \sim \text{iid } \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$.

3: **for** $t \in \{T - t_k\}_{0 \leq k \leq N}$ **do**

4: Solve the ERM problem below for $t = T - t_k$:

$$\widehat{\mathbf{W}}_t = \arg \min_{\mathbf{W} \in \mathcal{W}_{d,m,D,L,M,B}} \frac{1}{n} \sum_{i=1}^n \left\| \sigma_t^{-1} \mathbf{g}_i + \text{P}_t[\text{ResN}_{\mathbf{W}}](\lambda_t \mathbf{x}_i + \sigma_t \mathbf{g}_i, \boldsymbol{\theta}_i) \right\|_2^2. \quad (18)$$

5: Take the approximate score function to be $\hat{s}_t(\mathbf{z}; \boldsymbol{\theta}) = \text{P}_t[\text{ResN}_{\widehat{\mathbf{W}}_t}](\mathbf{z}, \boldsymbol{\theta})$.

6: // Sampling by discretizing the stochastic differential equation

7: Sample $\widehat{\mathbf{Y}}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$.

8: **for** $k = 0, \dots, N - 1$ **do**

9: Sample $\mathbf{G}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. Calculate $\widehat{\mathbf{Y}}_{k+1}$ using the exponential integrator scheme: (here $\boldsymbol{\theta}$ is provided as an input)

$$\widehat{\mathbf{Y}}_{k+1} = e^{\gamma_k} \cdot \widehat{\mathbf{Y}}_k + 2(e^{\gamma_k} - 1) \cdot \hat{s}_{T-t_k}(\widehat{\mathbf{Y}}_k; \boldsymbol{\theta}) + \sqrt{e^{2\gamma_k} - 1} \cdot \mathbf{G}_k. \quad (19)$$

Return: $\hat{\mathbf{x}} = \widehat{\mathbf{Y}}_N$.

702 independent of X_0 . The reverse process that corresponds to process (20) is defined via the SDE

$$dY_t = \{Y_t + 2\nabla \mu_{T-t}(Y_t)\} dt + \sqrt{2} dB'_t, \quad Y_0 \sim \mu_T. \quad (21)$$

703 An approximation to continuous-time process (21) is obtained via performing time discretization,
704 which directly leads to a sampling algorithm. More precisely, for $0 = t_0 < t_1 < \dots < t_N = T - \delta$,
705 we let

$$d\widehat{Y}_t = \{\widehat{Y}_t + 2\hat{s}_{T-t_k}(\widehat{Y}_t)\} dt + d\widehat{B}_t \quad \text{for } t_k \leq t \leq t_{k+1}, \quad \widehat{Y}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d), \quad (22)$$

706 where $\hat{s}_{T-t}(\cdot)$ is an estimate of the true score function $s_{T-t}(\cdot) = \nabla \log \mu_{T-t}(\cdot)$. We denote by p_t
707 the marginal distribution of \widehat{Y}_t , and set $\gamma_k = t_{k+1} - t_k$. In addition, we assume there exists $\kappa > 0$,
708 such that $\gamma_k \leq \kappa \cdot \min\{1, T - t_{k+1}\}$.

709 Next, we state the assumptions required to establish the discretization error bound of the DDPM
710 sampling scheme.

711 **Assumption 5** (Rescaled version of Benton et al. [2023] Assumption 1). *The score function estimator*
712 *\hat{s}_t satisfies*

$$\sum_{k=0}^{N-1} \gamma_k \mathbb{E}_{\mathbf{x} \sim \mu_{T-t_k}} \left[\|\nabla \log \mu_{T-t_k}(\mathbf{x}) - \hat{s}_{T-t_k}(\mathbf{x})\|_2^2 \right] \leq d \cdot \varepsilon_{\text{score}}^2.$$

713 **Assumption 6.** *The data distribution μ has finite second moment: $\mathbb{E}_{\mathbf{x}_0 \sim \mu} [\|\mathbf{x}_0\|_2^2] \leq d \cdot B$, where*
714 *$B \geq 1$ is a fixed constant.*

715 With Assumptions 5 and 6, we are ready to state the main theorem for this part.

716 **Theorem 5.** [Theorem 1 of Benton et al. [2023]] *Let Assumptions 5 and 6 hold. Then there exists a*
717 *numerical constant $C_0 > 0$, such that*

$$\text{KL}(\mu_\delta, p_{t_N}) \leq C_0 \cdot d \cdot (\varepsilon_{\text{score}}^2 + \kappa^2 NB + \kappa TB + e^{-2T} B).$$

718 *Proof sketch of Theorem 5.*

719 **Part 1.** We first control the quantity

$$E_{s,t} = \mathbb{E} \left[\|\nabla \log \mu_{T-t}(Y_t) - \nabla \log \mu_{T-s}(Y_s)\|_2^2 \right],$$

720 where $0 \leq s \leq t \leq T$. According to Lemma 2 of Benton et al. [2023], we have

$$\begin{aligned} d \left(\|\nabla \log \mu_{T-t}(Y_t) - \nabla \log \mu_{T-s}(Y_s)\|_2^2 \right) &= -2 \|\nabla \log \mu_{T-t}(Y_t) - \nabla \log \mu_{T-s}(Y_s)\|_2^2 dt \\ &\quad - 2 \{ \nabla \log \mu_{T-t}(Y_t) - \nabla \log \mu_{T-s}(Y_s) \} \cdot \nabla \log \mu_{T-s}(Y_s) dt + 2 \|\nabla^2 \log \mu_{T-t}(Y_t)\|_F^2 dt \\ &\quad + 2\sqrt{2} \{ \nabla \log \mu_{T-t}(Y_t) - \nabla \log \mu_{T-s}(Y_s) \} \cdot \nabla^2 \log \mu_{T-t}(Y_t) \cdot dB'_t. \end{aligned} \quad (23)$$

721 In the above display, s is fixed and t varies. Taking expectation and integrate over $[s, t]$, we obtain

$$\begin{aligned} \mathbb{E} \left[\|\nabla \log \mu_{T-t}(Y_t) - \nabla \log \mu_{T-s}(Y_s)\|_2^2 \right] &= \mathbb{E} \int_s^t -2 \|\nabla \log \mu_{T-r}(Y_r) - \nabla \log \mu_{T-s}(Y_s)\|_2^2 dr \\ &\quad - \mathbb{E} \int_s^t 2 \{ \nabla \log \mu_{T-r}(Y_r) - \nabla \log \mu_{T-s}(Y_s) \} \cdot \nabla \log \mu_{T-s}(Y_s) dr \\ &\quad + \mathbb{E} \int_s^t 2 \|\nabla^2 \log \mu_{T-r}(Y_r)\|_F^2 dr. \end{aligned}$$

722 Observe that all terms above are integrable. Hence, we may apply Fubini's theorem and interchange
723 integration and expectation, which gives

$$\begin{aligned} \frac{dE_{s,t}}{dt} &= -2 \mathbb{E} \left[\|\nabla \log \mu_{T-t}(Y_t) - \nabla \log \mu_{T-s}(Y_s)\|_2^2 \right] \\ &\quad + 2 \mathbb{E} \left[\{ \nabla \log \mu_{T-s}(Y_s) - \nabla \log \mu_{T-t}(Y_t) \} \cdot \nabla \log \mu_{T-s}(Y_s) \right] + 2 \mathbb{E} \left[\|\nabla^2 \log \mu_{T-t}(Y_t)\|_F^2 \right]. \end{aligned}$$

724 Invoking Cauchy-Schwartz inequality, we have

$$\frac{dE_{s,t}}{dt} \leq \mathbb{E} \left[\|\nabla \log \mu_{T-s}(Y_s)\|_2^2 \right] + 2 \mathbb{E} \left[\|\nabla^2 \log \mu_{T-t}(Y_t)\|_F^2 \right]. \quad (24)$$

725 Next, we upper bound $\mathbb{E} \left[\|\nabla \log \mu_{T-s}(Y_s)\|_2^2 \right]$ and $\mathbb{E} \left[\|\nabla^2 \log \mu_{T-t}(Y_t)\|_F^2 \right]$, respectively.

726 Lemma 3 of Benton et al. [2023] gives

$$\begin{aligned} \nabla \log \mu_t(\mathbf{x}_t) &= -\sigma_t^{-2} \mathbf{x}_t + e^{-t} \sigma_t^{-2} \mathbf{m}_t(\mathbf{x}_t), \\ \nabla^2 \log \mu_t(\mathbf{x}_t) &= -\sigma_t^{-2} \mathbf{I} + e^{-2t} \sigma_t^{-4} \boldsymbol{\Sigma}_t(\mathbf{x}_t), \end{aligned} \quad (25)$$

727 where $\sigma_t^2 = 1 - e^{-2t}$, $\mathbf{m}_t(\mathbf{x}_t) = \mathbb{E}_{\mu_0 | \mu_t(\mathbf{x}_0 | \mathbf{x}_t)}[\mathbf{x}_0]$, and $\boldsymbol{\Sigma}_t(\mathbf{x}_t) = \text{Cov}_{\mu_0 | \mu_t(\mathbf{x}_0 | \mathbf{x}_t)}[\mathbf{x}_0]$. By Eq. (25),
728 we see that

$$\begin{aligned} &\mathbb{E}_{\mathbf{x}_t \sim \mu_t} \left[\|\nabla \log \mu_t(\mathbf{x}_t)\|_2^2 \right] \\ &= \sigma_t^{-4} \mathbb{E}_{\mathbf{x}_t \sim \mu_t} \left[\|\mathbf{x}_t\|_2^2 \right] - 2e^{-t} \sigma_t^{-4} \mathbb{E}_{\mathbf{x}_t \sim \mu_t} [\mathbf{x}_t \cdot \mathbf{m}_t(\mathbf{x}_t)] + e^{-2t} \sigma_t^{-4} \mathbb{E}_{\mathbf{x}_t \sim \mu_t} \left[\|\mathbf{m}_t(\mathbf{x}_t)\|_2^2 \right]. \end{aligned}$$

729 Note that

$$\begin{aligned} \mathbb{E}_{\mathbf{x}_t \sim \mu_t} [\mathbf{x}_t \cdot \mathbf{m}_t(\mathbf{x}_t)] &= \mathbb{E}_{\mathbf{x}_t \sim \mu_t} [\mathbf{x}_t \cdot \mathbf{x}_0] = e^{-t} \mathbb{E}_{\mathbf{x}_0 \sim \mu_0} [\|\mathbf{x}_0\|^2] \leq dB e^{-t}, \\ \text{Tr}(\boldsymbol{\Sigma}_t(\mathbf{x}_t)) &= \mathbb{E}[\|\mathbf{x}_0\|^2 | \mathbf{x}_t] - \|\mathbf{m}_t(\mathbf{x}_t)\|_2^2, \end{aligned}$$

730 hence

$$\begin{aligned} &\mathbb{E}_{\mathbf{x}_t \sim \mu_t} \left[\|\nabla \log \mu_t(\mathbf{x}_t)\|_2^2 \right] \\ &= \sigma_t^{-4} \cdot (e^{-2t} \mathbb{E}[\|\mathbf{x}_0\|^2] + \sigma_t^2 d) - 2e^{-2t} \sigma_t^{-4} \mathbb{E}[\|\mathbf{x}_0\|^2] + e^{-2t} \sigma_t^{-4} \cdot (\mathbb{E}[\|\mathbf{x}_0\|^2] - \mathbb{E}[\text{Tr}(\boldsymbol{\Sigma}_t(\mathbf{x}_t))]) \\ &= \sigma_t^{-2} d - e^{-2t} \sigma_t^{-4} \mathbb{E}[\text{Tr}(\boldsymbol{\Sigma}_t(\mathbf{x}_t))] \leq d \sigma_t^{-2}. \end{aligned} \quad (26)$$

731 That is to say, we have $\mathbb{E} \left[\|\nabla \log \mu_{T-s}(Y_s)\|_2^2 \right] \leq d \sigma_{T-s}^{-2}$. We write $\boldsymbol{\Sigma}_t = \boldsymbol{\Sigma}_t(\mathbf{x}_t)$ for short. The
732 second part of Eq. (25) implies that

$$\mathbb{E}_{\mathbf{x}_t \sim \mu_t} \left[\|\nabla^2 \log \mu_t(\mathbf{x}_t)\|_F^2 \right] = \sigma_t^{-4} d - 2\sigma_t^{-6} e^{-2t} \mathbb{E}[\text{Tr}(\boldsymbol{\Sigma}_t)] + e^{-4t} \sigma_t^{-8} \mathbb{E}[\text{Tr}(\boldsymbol{\Sigma}_t^2)]. \quad (27)$$

733 Lemma 1 of Benton et al. [2023] gives

$$\frac{e^{2t} \sigma_t^4}{2} \frac{d}{dt} \mathbb{E}[\boldsymbol{\Sigma}_t] = \mathbb{E}[\boldsymbol{\Sigma}_t^2]. \quad (28)$$

734 Putting together Eq. (27) and (28), we obtain

$$\begin{aligned} \mathbb{E}_{\mathbf{x}_t \sim \mu_t} [\|\nabla^2 \log \mu_t(\mathbf{x}_t)\|_F^2] &= d\sigma_t^{-4} - 2\sigma_t^{-6} e^{-2t} \mathbb{E}[\text{Tr}(\boldsymbol{\Sigma}_t)] + \frac{e^{-2t} \sigma_t^{-4}}{2} \frac{d}{dt} \mathbb{E}[\text{Tr}(\boldsymbol{\Sigma}_t)] \\ &\leq d\sigma_t^{-4} + \frac{1}{2} \frac{d}{dt} (\sigma_t^{-4} \mathbb{E}[\text{Tr}(\boldsymbol{\Sigma}_t)]). \end{aligned} \quad (29)$$

735 Putting together Eq. (26) and (29), we get

$$\begin{aligned} &\mathbb{E} [\|\nabla \log \mu_{T-s}(Y_s)\|_2^2] + 2\mathbb{E} [\|\nabla^2 \log \mu_{T-t}(Y_t)\|_F^2] \\ &\leq \sigma_{T-s}^{-2} d + 2d\sigma_{T-t}^{-4} - \frac{d}{dr} (\sigma_{T-r}^{-4} \mathbb{E}[\text{Tr}(\boldsymbol{\Sigma}_{T-r})]) \Big|_{r=t}. \end{aligned}$$

736 We define

$$E_{s,t}^{(1)} := d\sigma_{T-s}^{-2} + 2d\sigma_{T-t}^{-4}, \quad E_{s,t}^{(2)} := -\frac{d}{dr} (\sigma_{T-r}^{-4} \mathbb{E}[\text{Tr}(\boldsymbol{\Sigma}_{T-r})]) \Big|_{r=t}.$$

737 According to Eq. (24) and notice that $E_{t_k, t_k} = 0$, we have

$$E_{t_k, t} \leq \int_{t_k}^t \{ \mathbb{E} [\|\nabla \log \mu_{T-t_k}(Y_{t_k})\|_2^2] + 2\mathbb{E} [\|\nabla^2 \log \mu_{T-s}(Y_s)\|_F^2] \} ds \leq \int_{t_k}^t (E_{t_k, s}^{(1)} + E_{t_k, s}^{(2)}) ds.$$

738 Following exactly the same procedure as in Benton et al. [2023], we conclude that there exists a
739 positive numerical constant C_0 , such that

$$\sum_{k=0}^{N-1} \int_{t_k}^{t_{k+1}} \mathbb{E} [\|\nabla \log \mu_{T-t}(Y_t) - \nabla \log \mu_{T-t_k}(Y_{t_k})\|_2^2] \leq C_0(\kappa^2 dNB + \kappa dTB).$$

740 **Part 2.** We denote by Q the distribution of Y_{t_N} derived from process (21), and P^{μ_T} the distribution
741 of process (22) at time t_N initialized at μ_T . By Proposition 3 of Benton et al. [2023], we obtain

$$\text{KL}(Q \| P^{\mu_T}) \leq \sum_{k=0}^{N-1} \int_{t_k}^{t_{k+1}} \mathbb{E} [\|\nabla \log \mu_{T-t}(Y_t) - \hat{s}_{T-t_k}(Y_{t_k})\|_2^2] dt,$$

742 which by triangle inequality is no smaller than

$$\begin{aligned} &2 \sum_{k=0}^{N-1} \gamma_k \mathbb{E} [\|\nabla \log \mu_{T-t_k}(Y_{t_k}) - \hat{s}_{T-t_k}(Y_{t_k})\|_2^2] dt \\ &+ 2 \sum_{k=0}^{N-1} \int_{t_k}^{t_{k+1}} \mathbb{E} [\|\nabla \log \mu_{T-t_k}(Y_{t_k}) - \nabla \log \mu_{T-t}(Y_t)\|_2^2] \\ &\leq 2d \cdot \varepsilon_{\text{score}}^2 + 2C_0(\kappa^2 dNB + \kappa dTB). \end{aligned}$$

743 We denote by P the distribution of process (22) at time t_N initialized at $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. By Eq. (19) of
744 Benton et al. [2023], we have

$$\text{KL}(Q \| P) = \text{KL}(Q \| P^{\mu_T}) + \text{KL}(\mu_T \| \mathcal{N}(\mathbf{0}, \mathbf{I}_d)).$$

745 Proposition 4 of Benton et al. [2023] gives $\text{KL}(\mu_T \| \mathcal{N}(\mathbf{0}, \mathbf{I}_d)) \lesssim dB e^{-2T}$. Putting together the
746 above upper bounds, we arrive at the following conclusion:

$$\text{KL}(\mu_\delta \| p_{t_N}) \leq C_0 \cdot d \cdot (B e^{-2T} + \kappa^2 NB + \kappa TB + \varepsilon_{\text{score}}^2),$$

747 thus concluding the proof of Theorem 5. \square

748 D.3 Generalization error of empirical risk minimization over ResNets

749 D.3.1 Result for Ising models

750 Note that the conditional (and unconditional) DDPM methods estimate the score function $\hat{s}_t =$
751 $P_t \text{ResN}_{\widehat{\mathbf{W}}_t}$ by solving the following ERM problem:

$$\begin{aligned} \widehat{\mathbf{W}}_t &= \text{argmin}_{\mathbf{W} \in \mathcal{W}_{d,m,D,L,M,B}} \widehat{R}_n(\mathbf{W}), \\ \widehat{R}_n(\mathbf{W}) &= \frac{1}{nd} \sum_{i=1}^n \|\sigma_t^{-1} \mathbf{g}_i + P_t(\text{ResN}_{\mathbf{W}}(\lambda_t \mathbf{x}_i + \sigma_t \mathbf{g}_i, \boldsymbol{\theta}_i))\|_2^2. \end{aligned} \quad (30)$$

752 Here, $\mathbf{x}_i, \mathbf{g}_i \in \mathbb{R}^d$, and $\boldsymbol{\theta}_i \in \mathbb{R}^m$ follow $\{(\mathbf{x}_i, \boldsymbol{\theta}_i, \mathbf{z}_i)\}_{i \in [n]} \sim_{iid} \mu \otimes \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. Recall that the
753 truncation operator gives $\text{P}_t[f](\mathbf{z}, \boldsymbol{\theta}) = \text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(f(\mathbf{z}, \boldsymbol{\theta}) + \sigma_t^{-2} \mathbf{z}) - \sigma_t^{-2} \mathbf{z}$. In cases where $\boldsymbol{\theta}_i$
754 does not exist (unconditional DDPM), we simply set $m = 0$. The population risk gives

$$R(\mathbf{W}) := \frac{1}{d} \mathbb{E}_{(\mathbf{x}, \boldsymbol{\theta}, \mathbf{g}) \sim \mu \otimes \mathcal{N}(\mathbf{0}, \mathbf{I}_d)} \left[\left\| \sigma_t^{-1} \mathbf{g}_0 + \text{P}_t(\text{ResN}_{\mathbf{W}}([\lambda_t \mathbf{x} + \sigma_t \mathbf{g}, \boldsymbol{\theta}]) \right\|_2^2 \right].$$

755 In the proposition below, we provide a uniform upper bound for $|\widehat{R}(\mathbf{W}) - R(\mathbf{W})|$ over
756 $\mathcal{W}_{d,m,D,L,M,B}$, where the ResNet class is given by Eq. (17).

757 **Proposition 6.** Assume that $\mu \in \mathcal{P}([-1, 1]^{d+m})$. There exists a numerical constant $C > 0$, such
758 that with probability at least $1 - \eta$,

$$\begin{aligned} & \sup_{\mathbf{W} \in \mathcal{W}_{d,m,D,L,M,B}} \left| \widehat{R}(\mathbf{W}) - R(\mathbf{W}) \right| \\ & \leq C \cdot \frac{\lambda_t^2}{\sigma_t^4} \cdot \sqrt{\frac{[(d+m)D + LDM] \cdot [L \cdot \log(LB(m+d)/d) + \log(\lambda_t^{-1})] + \log(1/\eta)}{n}}. \end{aligned}$$

759 *Proof of Proposition 6.* The proof of this proposition uses the following lemma.

760 **Lemma 4** (Proposition A.4 of Bai et al. [2023]). Suppose that $\{X_w\}_{w \in \Theta}$ is a zero-mean random
761 process given by

$$X_w \equiv \frac{1}{n} \sum_{i=1}^n f(z_i; w) - \mathbb{E}_z[f(z; w)],$$

762 where z_1, \dots, z_n are i.i.d samples from a distribution \mathbb{P}_z such that the following assumption holds:

763 (a) The index set Θ is equipped with a distance ρ and diameter B . Further, assume that for
764 some constant A , for any ball Θ' of radius r in Θ , the covering number admits upper bound
765 $\log N(\Delta; \Theta', \rho) \leq d \log(2Ar/\Delta)$ for all $0 < \Delta \leq 2r$.

766 (b) For any fixed $w \in \Theta$ and z sampled from \mathbb{P}_z , the random variable $f(z; w) - \mathbb{E}_z[f(z; w)]$ is
767 a σ -sub-Gaussian random variable ($\mathbb{E}[e^{\lambda[f(z; w) - \mathbb{E}_z[f(z; w)]]}] \leq e^{\lambda^2 \sigma^2 / 2}$ for any $\lambda \in \mathbb{R}$).

768 (c) For any $w, w' \in \Theta$ and z sampled from \mathbb{P}_z , the random variable $f(z; w) - f(z; w')$ is a
769 $\sigma' \rho(w, w')$ -sub-Gaussian random variable ($\mathbb{E}[e^{\lambda[f(z; w) - f(z; w')]}] \leq e^{\lambda^2 (\sigma')^2 \rho^2(w, w') / 2}$ for
770 any $\lambda \in \mathbb{R}$).

771 Then with probability at least $1 - \eta$, it holds that

$$\sup_{w \in \Theta} |X_w| \leq C \sigma \sqrt{\frac{d \cdot \log(2A(1 + B\sigma'/\sigma)) + \log(1/\eta)}{n}},$$

772 where C is a universal constant.

773 In Lemma 4, we can take $z = (\mathbf{g}, \mathbf{x}, \boldsymbol{\theta})$, $w = \mathbf{W}$, $\Theta = \mathcal{W}_{d,m,D,L,M,B}$, $\rho(w, w') = \|\mathbf{W} - \mathbf{W}'\|$,
774 and $f(z_i; w) = d^{-1} \|\sigma_t^{-1} \mathbf{g}_i + \text{P}_t(\text{ResN}_{\mathbf{W}}(\lambda_t \mathbf{x}_i + \sigma_t \mathbf{g}_i, \boldsymbol{\theta}_i))\|_2^2$. Therefore, to show Proposition 6,
775 we just need to apply Lemma 4 by checking (a), (b), (c).

776 **Check (a).** We note that the index set $\Theta = \mathcal{W}_{d,m,D,L,M,B}$ equipped with $\rho(w, w') = \|\mathbf{W} - \mathbf{W}'\|$
777 has diameter $2B$. Further note that $\mathcal{W}_{d,m,D,L,M,B}$ has dimension bounded by $4(d+m)D + 2LDM$.
778 According to Example 5.8 of Wainwright [2019], it holds that $\log N(\Delta; \mathcal{W}_{d,m,D,L,M,r}, \|\cdot\|) \leq$
779 $[4(d+m)D + 2LDM] \cdot \log(1 + 2r/\Delta)$ for any $0 < \Delta \leq 2r$. This verifies (a).

780 **Check (b).** By the definition of the projection operator that $\text{P}_t[f](\mathbf{z}) = \text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(f(\mathbf{z}) + \sigma_t^{-2} \mathbf{z}) -$
781 $\sigma_t^{-2} \mathbf{z}$ and that $\mathbf{z} = \lambda_t \mathbf{x} + \sigma_t \mathbf{g}$, we have

$$\begin{aligned} 0 \leq f(z; w) &= d^{-1} \|\sigma_t^{-1} \mathbf{g} + \text{P}_t(\text{ResN}_{\mathbf{W}}(\lambda_t \mathbf{x} + \sigma_t \mathbf{g}, \boldsymbol{\theta}))\|_2^2 \\ &= d^{-1} \|\lambda_t \sigma_t^{-2} \mathbf{x} + \text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\text{ResN}_{\mathbf{W}_1}(\lambda_t \mathbf{x} + \sigma_t \mathbf{g}, \boldsymbol{\theta}) + \sigma_t^{-2} \mathbf{z})\|_2^2 \\ &\leq 4\lambda_t^2 \sigma_t^{-4}. \end{aligned}$$

782 As a consequence, $f(z, w) - \mathbb{E}_z[f(z, w)]$ is a $\sigma = 4\lambda_t^2\sigma_t^{-4}$ sub-Gaussian random variable.

783 **Check (c).** Direct calculation yields

$$\begin{aligned}
& |f(z; w_1) - f(z; w_2)| \\
&= \frac{1}{d} \left| \|\sigma_t^{-1}\mathbf{g} + \bar{\text{P}}_t(\text{ResN}_{\mathbf{W}_1}(\lambda_t\mathbf{x} + \sigma_t\mathbf{g}, \boldsymbol{\theta}))\|_2^2 - \|\sigma_t^{-1}\mathbf{g} + \bar{\text{P}}_t(\text{ResN}_{\mathbf{W}_2}(\lambda_t\mathbf{x} + \sigma_t\mathbf{g}, \boldsymbol{\theta}))\|_2^2 \right| \\
&= \frac{1}{d} \left| \left\| -\lambda_t\sigma_t^{-2}\mathbf{x} + \text{proj}_{\lambda_t\sigma_t^{-2}\sqrt{d}}(\text{ResN}_{\mathbf{W}_1}(\lambda_t\mathbf{x} + \sigma_t\mathbf{g}, \boldsymbol{\theta}) + \sigma_t^{-2}\mathbf{z}) \right\|_2^2 \right. \\
&\quad \left. - \left\| -\lambda_t\sigma_t^{-2}\mathbf{x} + \text{proj}_{\lambda_t\sigma_t^{-2}\sqrt{d}}(\text{ResN}_{\mathbf{W}_2}(\lambda_t\mathbf{x} + \sigma_t\mathbf{g}, \boldsymbol{\theta}) + \sigma_t^{-2}\mathbf{z}) \right\|_2^2 \right| \\
&\leq \frac{8\lambda_t}{\sigma_t^2\sqrt{d}} \cdot \left\| \text{proj}_{\lambda_t\sigma_t^{-2}\sqrt{d}}(\text{ResN}_{\mathbf{W}_1}(\lambda_t\mathbf{x} + \sigma_t\mathbf{g}, \boldsymbol{\theta}) + \sigma_t^{-2}\mathbf{z}) \right. \\
&\quad \left. - \text{proj}_{\lambda_t\sigma_t^{-2}\sqrt{d}}(\text{ResN}_{\mathbf{W}_2}(\lambda_t\mathbf{x} + \sigma_t\mathbf{g}, \boldsymbol{\theta}) + \sigma_t^{-2}\mathbf{z}) \right\|_2 \\
&\lesssim \frac{2\lambda_t L(B^2 + 1)^L}{\sigma_t^2} \cdot \frac{1}{\sqrt{d}} \left(\lambda_t \|\mathbf{x}\|_2 + \sigma_t \|\mathbf{g}\|_2 + \|\boldsymbol{\theta}\|_2 \right) \cdot \|\mathbf{W}_1 - \mathbf{W}_2\|.
\end{aligned}$$

784 Notice that $(\mathbf{x}, \boldsymbol{\theta}, \mathbf{g}) \sim \mu \otimes \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and note that $\mu \in \mathcal{P}([-1, 1]^{d+m})$, we have that $\|\boldsymbol{\theta}\|_2/\sqrt{d}$ is
785 $\sqrt{m/d}$ -bounded and is thus $\mathcal{O}(\sqrt{m/d})$ -sub-Gaussian, $\|\mathbf{x}\|_2/\sqrt{d}$ is 1-bounded and is thus $\mathcal{O}(1)$ -
786 sub-Gaussian, and $\|\mathbf{g}\|_2/\sqrt{d}$ is $\mathcal{O}(1)$ -sub-Gaussian. As a consequence, $f(z; w_1) - f(z; w_2)$ is
787 $\sigma'\rho(w_1, w_2) = C \cdot \lambda_t\sigma_t^{-2}L(B^2 + 1)^L\sqrt{(m+d)/d} \cdot \|\mathbf{W}_1 - \mathbf{W}_2\|$ sub-Gaussian.

788 Therefore, we apply Lemma 4, and use the fact that

$$\begin{aligned}
& \log(2(1 + B\sigma'/\sigma)) \\
&= \log(2(1 + (C/2)B\lambda_t^{-1}\sigma_t^2L(B^2 + 1)^L\sqrt{(m+d)/d})) \lesssim L\log(LB(m+d)/d) + \log(\lambda_t^{-1}).
\end{aligned}$$

789 This concludes the proof of Proposition 6. \square

790 D.3.2 Result for Sparse coding

791 In the setting of sparse coding, we assume a fixed dictionary $\mathbf{A} \in \mathbb{R}^{d \times m}$. The model $\mathbf{x} \sim \mu$ is given
792 by $\mathbf{x} = \mathbf{A}\boldsymbol{\theta} + \boldsymbol{\varepsilon}$, where $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \tau^2\mathbf{I}_d)$ is independent of anything else and $\theta_i \sim_{iid} \pi_\theta \in \mathcal{P}([-1, 1])$
793 for $i \in [m]$. Assume that we have $\{(\mathbf{x}_i, \mathbf{g}_i)\}_{i \in [n]} \sim_{iid} \mu \otimes \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. We are interested in estimating
794 the score function $\hat{\mathbf{s}}_t = \bar{\text{P}}_t \text{ResN}_{\widehat{\mathbf{W}}_t}$ by solving the following ERM problem:

$$\begin{aligned}
\widehat{\mathbf{W}}_t &= \text{argmin}_{\mathbf{W} \in \mathcal{W}_{d,D,L,M,B}} \widehat{R}_n(\mathbf{W}), \\
\widehat{R}_n(\mathbf{W}) &= \frac{1}{nd} \sum_{i=1}^n \|\sigma_t^{-1}\mathbf{g}_i + \bar{\text{P}}_t(\text{ResN}_{\mathbf{W}}(\lambda_t\mathbf{x}_i + \sigma_t\mathbf{g}_i))\|_2^2.
\end{aligned} \tag{31}$$

795 Here, the truncation operator gives $\bar{\text{P}}_t[f](\mathbf{z}) = \text{proj}_{\sqrt{m}\|\mathbf{A}\|_{\text{op}}\Pi \cdot \lambda_t(\sigma_t^2 + \tau^2\lambda_t^2)^{-1}}(f(\mathbf{z}) + (\sigma_t^2 + \tau^2\lambda_t^2)^{-1}\mathbf{z}) - (\sigma_t^2 + \tau^2\lambda_t^2)^{-1}\mathbf{z}$. The corresponding population risk gives

$$R(\mathbf{W}) := \frac{1}{d} \mathbb{E}_{(\mathbf{x}, \mathbf{g}) \sim \mu \otimes \mathcal{N}(\mathbf{0}, \mathbf{I}_d)} \left[\|\sigma_t^{-1}\mathbf{g} + \bar{\text{P}}_t(\text{ResN}_{\mathbf{W}}(\lambda_t\mathbf{x} + \sigma_t\mathbf{g}))\|_2^2 \right].$$

797 In the proposition below, we provide a uniform upper bound for $|\widehat{R}(\mathbf{W}) - R(\mathbf{W})|$ over $\mathcal{W}_{d,D,L,M,B}$
798 in the sparse coding setup, where the ResNet class is given by Eq. (3).

799 **Proposition 7.** *Under the setting of sparse coding stated above, there exists a numerical constant*
800 $C > 0$, *such that with probability at least $1 - \eta$, for $n \geq \log(2/\eta)$, we have*

$$\begin{aligned}
& \sup_{\mathbf{W} \in \mathcal{W}_{d,D,L,M,B}} \left| \widehat{R}(\mathbf{W}) - R(\mathbf{W}) \right| \lesssim \left(\lambda_t^2 \|\mathbf{A}\|_{\text{op}}^2 \Pi^2 (\tau^{-4} + 1) \frac{m}{d} + \frac{\lambda_t^2}{\sigma_t^2} (1 + \tau^2) \right) \\
& \times \sqrt{\frac{(dD + LDM) \cdot [T + L\log(LB) + \log(nmT(\tau + 1)(\|\mathbf{A}\|_{\text{op}}\Pi + 1)\tau^{-1})] + \log(2/\eta)}{n}}.
\end{aligned}$$

801 *Proof of Proposition 7.* Note that $\{(\mathbf{x}_i, \mathbf{g}_i)\}_{i \in [n]} \sim_{iid} \mu \times \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ where μ is the sparse coding
802 model. Then we must have $\mathbf{x}_i = \mathbf{A}\boldsymbol{\theta}_i + \boldsymbol{\varepsilon}_i$ for some $(\boldsymbol{\theta}_i, \boldsymbol{\varepsilon}_i) \sim_{iid} \pi_0^m \times \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I}_d)$. Denote
803 $\mathbf{z} = (\mathbf{g}, \mathbf{x}, \boldsymbol{\varepsilon})$, $\mathbf{w} = \mathbf{W}$, and

$$f(\mathbf{z}; \mathbf{w}) = d^{-1} (\|\sigma_t^{-1} \mathbf{g} + \bar{\mathbf{P}}_t(\text{ResN}_{\mathbf{W}}(\lambda_t \mathbf{x} + \sigma_t \mathbf{g}))\|_2^2 - \|(\sigma_t^{-1} - \sigma_t(\sigma_t^2 + \tau^2 \lambda_t^2)^{-1}) \mathbf{g} - \lambda_t(\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \boldsymbol{\varepsilon}\|_2^2).$$

804 We further denote $\mathbf{z} = \lambda_t \mathbf{x} + \sigma_t \mathbf{g}$. Note that we have

$$\begin{aligned} & |f(\mathbf{z}; \mathbf{w}_1) - f(\mathbf{z}; \mathbf{w}_2)| \\ &= \frac{1}{d} \left| \|\sigma_t^{-1} \mathbf{g} + \bar{\mathbf{P}}_t(\text{ResN}_{\mathbf{W}_1}(\lambda_t \mathbf{x} + \sigma_t \mathbf{g}))\|_2^2 - \|\sigma_t^{-1} \mathbf{g} + \bar{\mathbf{P}}_t(\text{ResN}_{\mathbf{W}_2}(\lambda_t \mathbf{x} + \sigma_t \mathbf{g}))\|_2^2 \right| \\ &\leq \frac{1}{d} \left| \|\text{proj}_{\sqrt{m} \|\mathbf{A}\|_{\text{op}} \Pi \cdot \lambda_t (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1}}(\text{ResN}_{\mathbf{W}_1}(\mathbf{z}) + (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \mathbf{z}) - (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \mathbf{z} + \sigma_t^{-1} \mathbf{g}\|_2^2 \right. \\ &\quad \left. - \|\text{proj}_{\sqrt{m} \|\mathbf{A}\|_{\text{op}} \Pi \cdot \lambda_t (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1}}(\text{ResN}_{\mathbf{W}_2}(\mathbf{z}) + (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \mathbf{z}) - (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \mathbf{z} + \sigma_t^{-1} \mathbf{g}\|_2^2 \right| \\ &\lesssim \left(\frac{\sqrt{m} \lambda_t \Pi \|\mathbf{A}\|_{\text{op}}}{d(\sigma_t^2 + \tau^2 \lambda_t^2)} + \frac{\lambda_t}{d(\sigma_t^2 + \tau^2 \lambda_t^2)} \|\boldsymbol{\varepsilon}\|_2 + \frac{\tau^2 \lambda_t^2}{d \sigma_t (\sigma_t^2 + \tau^2 \lambda_t^2)} \|\mathbf{g}\|_2 \right) \\ &\quad \times \left\| \text{proj}_{\sqrt{m} \|\mathbf{A}\|_{\text{op}} \Pi \cdot \lambda_t (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1}}(\text{ResN}_{\mathbf{W}_1}(\mathbf{z}) + (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \mathbf{z}) \right. \\ &\quad \left. - \text{proj}_{\sqrt{m} \|\mathbf{A}\|_{\text{op}} \Pi \cdot \lambda_t (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1}}(\text{ResN}_{\mathbf{W}_2}(\mathbf{z}) + (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \mathbf{z}) \right\|_2 \\ &\lesssim \frac{L(B^2 + 1)^L \lambda_t}{(\sigma_t^2 + \tau^2 \lambda_t^2)} \cdot \left(\frac{\sqrt{m} \Pi \|\mathbf{A}\|_{\text{op}}}{\sqrt{d}} + \frac{\|\boldsymbol{\varepsilon}\|_2}{\sqrt{d}} + \frac{\tau^2 \lambda_t \|\mathbf{g}\|_2}{\sqrt{d} \sigma_t} \right) \times \\ &\quad \frac{1}{\sqrt{d}} \left(\lambda_t \|\mathbf{A}\boldsymbol{\theta}\|_2 + \lambda_t \|\boldsymbol{\varepsilon}\|_2 + \sigma_t \|\mathbf{g}\|_2 \right) \cdot \|\mathbf{W}_1 - \mathbf{W}_2\|. \end{aligned}$$

805 Therefore, we denote by $\mathcal{N}(\Delta; \mathcal{W}_{d,D,L,M,B}, \rho)$ a Δ -covering of $\mathcal{W}_{d,D,L,M,B}$ under metric
806 $\rho(\mathbf{W}_1, \mathbf{W}_2) = \|\mathbf{W}_1 - \mathbf{W}_2\|$ for some $\Delta > 0$. Then

$$\begin{aligned} & \sup_{\mathbf{w} \in \mathcal{W}_{d,D,L,M,B}} \left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{z}_i; \mathbf{w}) - \mathbb{E}[f(\mathbf{z}; \mathbf{w})] \right| \\ &\leq \sup_{\mathbf{w} \in \mathcal{N}(\Delta; \mathcal{W}_{d,D,L,M,B}, \rho)} \left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{z}_i; \mathbf{w}) - \mathbb{E}[f(\mathbf{z}; \mathbf{w})] \right| + \frac{L(B^2 + 1)^L \lambda_t}{(\sigma_t^2 + \tau^2 \lambda_t^2)} \cdot \Delta \cdot (L_n + \mathbb{E}[L_n]), \end{aligned}$$

807 where

$$L_n = \frac{1}{nd} \sum_{i=1}^n \left(\sqrt{m} \Pi \|\mathbf{A}\|_{\text{op}} + \|\boldsymbol{\varepsilon}_i\|_2 + \sigma_t^{-1} \tau^2 \lambda_t \|\mathbf{g}_i\|_2 \right) \cdot \left(\lambda_t \|\mathbf{A}\boldsymbol{\theta}_i\|_2 + \lambda_t \|\boldsymbol{\varepsilon}_i\|_2 + \sigma_t \|\mathbf{g}_i\|_2 \right).$$

808 Since $(\boldsymbol{\theta}_i, \boldsymbol{\varepsilon}_i, \mathbf{g}_i) \sim \pi_0^m \otimes \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I}_d) \otimes \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$, $\sigma_t^2 \leq 1$ and $\lambda_t^2 \leq 1$, we have $\mathbb{E}[L_n] \leq \bar{L}$ and
809 $L_n - \mathbb{E}[L_n]$ is $\text{SE}(\bar{L}/\sqrt{n}, \bar{L})$, for $\bar{L} = (m/d) \Pi^2 \|\mathbf{A}\|_{\text{op}}^2 + \sigma_t^{-2} (\tau^4 + 1)$. By Bernstein's inequality,
810 we conclude that with probability at least $1 - \eta/2$, we have

$$\begin{aligned} L_n + \mathbb{E}[L_n] &\leq C \cdot \bar{L} (1 + \sqrt{\log(2/\eta)/n} + \log(2/\eta)/n) \leq C \cdot \bar{L} (1 + \log(2/\eta)) \\ &= C \cdot ((m/d) \Pi^2 \|\mathbf{A}\|_{\text{op}}^2 + \sigma_t^{-2} (\tau^4 + 1)) \cdot (1 + \log(2/\eta)). \end{aligned}$$

811 for some numerical constant C .

812 Furthermore, note that we have

$$\begin{aligned} & f(\mathbf{z}; \mathbf{w}) \\ &= d^{-1} \|\sigma_t^{-1} \mathbf{g} + \bar{\mathbf{P}}_t(\text{ResN}_{\mathbf{W}}(\lambda_t \mathbf{x} + \sigma_t \mathbf{g}))\|_2^2 - d^{-1} \|(\sigma_t^{-1} - \sigma_t(\sigma_t^2 + \tau^2 \lambda_t^2)^{-1}) \mathbf{g} - \lambda_t(\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \boldsymbol{\varepsilon}\|_2^2 \\ &= d^{-1} \|\sigma_t^{-1} \mathbf{g} + \text{proj}_{\sqrt{m} \|\mathbf{A}\|_{\text{op}} \Pi \cdot \lambda_t (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1}}(\text{ResN}_{\mathbf{W}}(\mathbf{z}) + (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \mathbf{z}) - (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \mathbf{z} + \sigma_t^{-1} \mathbf{g}\|_2^2 \\ &\quad - d^{-1} \|(\sigma_t^{-1} - \sigma_t(\sigma_t^2 + \tau^2 \lambda_t^2)^{-1}) \mathbf{g} - \lambda_t(\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \boldsymbol{\varepsilon}\|_2^2 \\ &= d^{-1} \|\text{proj}_{\sqrt{m} \|\mathbf{A}\|_{\text{op}} \Pi \cdot \lambda_t (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1}}(\text{ResN}_{\mathbf{W}}(\mathbf{z}) + (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \mathbf{z}) - \lambda_t(\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \mathbf{A}\boldsymbol{\theta}\|_2^2 \\ &\quad + 2d^{-1} \langle (\sigma_t^{-1} - \sigma_t(\sigma_t^2 + \tau^2 \lambda_t^2)^{-1}) \mathbf{g} - \lambda_t(\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \boldsymbol{\varepsilon}, \\ &\quad \text{proj}_{\sqrt{m} \|\mathbf{A}\|_{\text{op}} \Pi \cdot \lambda_t (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1}}(\text{ResN}_{\mathbf{W}}(\mathbf{z}) + (\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \mathbf{z}) - \lambda_t(\sigma_t^2 + \tau^2 \lambda_t^2)^{-1} \mathbf{A}\boldsymbol{\theta} \rangle. \end{aligned}$$

813 As a consequence, $f(z; w) - \mathbb{E}_z[f(z, w)]$ is sub-Gaussian with variance proxy

$$C^2 \cdot \left(\frac{m \|\mathbf{A}\|_{\text{op}}^2 \Pi^2 \lambda_t^2}{d(\sigma_t^2 + \tau^2 \lambda_t^2)^2} + \frac{\tau^2 \lambda_t^2}{\sigma_t^2(\sigma_t^2 + \lambda_t^2 \tau^2)} \right)^2$$

814 for some other numerical constant C . Therefore, with probability at least $1 - \eta/2$, by sub-Gaussian
815 tail bound and by the bound $\log |\mathcal{N}(\Delta; \mathcal{W}_{d,D,L,M,B,\rho})| \leq [4dD + 2LDM] \cdot \log(1 + 2B/\Delta)$, we
816 have

$$\begin{aligned} & \sup_{\mathbf{w} \in \mathcal{N}(\Delta; \mathcal{W}_{d,D,L,M,B,\rho})} \left| \frac{1}{n} \sum_{i=1}^n f(z_i; w_i) - \mathbb{E}[f(z; w)] \right| \\ & \lesssim \left(\frac{m \|\mathbf{A}\|_{\text{op}}^2 \Pi^2 \lambda_t^2}{d(\sigma_t^2 + \tau^2 \lambda_t^2)^2} + \frac{\tau^2 \lambda_t^2}{\sigma_t^2(\sigma_t^2 + \lambda_t^2 \tau^2)} \right) \cdot \sqrt{\frac{[4dD + 2LDM] \cdot \log(1 + 2B/\Delta) + \log(2/\eta)}{n}}. \end{aligned}$$

817 Setting

$$\Delta = \left(\frac{m \|\mathbf{A}\|_{\text{op}}^2 \Pi^2 \lambda_t^2}{d(\sigma_t^2 + \tau^2 \lambda_t^2)^2} + \frac{\tau^2 \lambda_t^2}{\sigma_t^2(\sigma_t^2 + \lambda_t^2 \tau^2)} \right) \cdot \frac{(\sigma_t^2 + \tau^2 \lambda_t^2)}{nL(B^2 + 1)^L \lambda_t \cdot (md^{-1} \Pi^2 \|\mathbf{A}\|_{\text{op}}^2 + \sigma_t^{-2}(\tau^4 + 1))},$$

818 we conclude that with probability at least $1 - \eta$, when $n \geq \log(2/\eta)$, we have

$$\begin{aligned} & \sup_{\mathbf{w} \in \mathcal{N}(\mathcal{W}_{d,D,L,M,B,\rho,\Delta})} \left| \frac{1}{n} \sum_{i=1}^n f(z_i; w_i) - \mathbb{E}[f(z; w)] \right| \\ & \lesssim n^{-1} \cdot \left(\frac{m \|\mathbf{A}\|_{\text{op}}^2 \Pi^2 \lambda_t^2}{d(\sigma_t^2 + \tau^2 \lambda_t^2)^2} + \frac{\tau^2 \lambda_t^2}{\sigma_t^2(\sigma_t^2 + \lambda_t^2 \tau^2)} \right) \cdot (\log(2/\eta) + 1) + \left\{ \frac{m \|\mathbf{A}\|_{\text{op}}^2 \Pi^2 \lambda_t^2}{d(\sigma_t^2 + \tau^2 \lambda_t^2)^2} + \frac{\tau^2 \lambda_t^2}{\sigma_t^2(\sigma_t^2 + \lambda_t^2 \tau^2)} \right\} \\ & \quad \times \sqrt{\frac{(dD + LDM) \cdot [T + L \log(LB) + \log(nmT(\tau + 1)(\|\mathbf{A}\|_{\text{op}} \Pi + 1)\tau^{-1})] + \log(2/\eta)}{n}} \\ & \lesssim \left(\lambda_t^2 \|\mathbf{A}\|_{\text{op}}^2 \Pi^2 (\tau^{-4} + 1) \frac{m}{d} + \frac{\lambda_t^2}{\sigma_t^2} (1 + \tau^2) \right) \\ & \quad \times \sqrt{\frac{(dD + LDM) \cdot [T + L \log(LB) + \log(nmT(\tau + 1)(\|\mathbf{A}\|_{\text{op}} \Pi + 1)\tau^{-1})] + \log(2/\eta)}{n}}, \end{aligned}$$

819 where the inequalities above uses the definition that $\lambda_t = e^{-t}$, $\sigma_t^2 = 1 - e^{-2t}$ and $t \leq T$. This
820 concludes the proof of Proposition 7. \square

821 D.4 Uniform approximation of the denoiser

822 The lemma below tells us that denoiser functions can be uniformly approximated with a linear
823 combination of $\text{ReLU}(\cdot)$ with changing intercepts. Furthermore, such approximation can achieve
824 arbitrary precision.

825 **Lemma 5.** Assume π_0 is a probability distribution over \mathbb{R} that has bounded support, and $\gamma > 0$ is a
826 fixed constant. Define $F(\lambda) := \mathbb{E}_{(\beta, z) \sim \pi_0 \otimes \mathcal{N}(0,1)} [\beta \mid \beta + \gamma^{-1/2} z = \lambda \gamma^{-1}]$. Let $\Pi_{\min} := \inf_{\lambda} F(\lambda)$,
827 $\Pi_{\max} := \sup_{\lambda} F(\lambda)$, $\Pi := \max\{|\Pi_{\max}|, |\Pi_{\min}|\}$, and $\Delta := \Pi_{\max} - \Pi_{\min}$. One can verify that
828 $F(\cdot)$ is Π^2 -Lipschitz continuous and non-decreasing. For any $\zeta > 0$, we define

$$w_{\zeta} := \inf \left\{ w : \text{for all } \lambda_1 > \lambda_2 \geq w \text{ or } \lambda_1 < \lambda_2 \leq -w \text{ we have } |F(\lambda_1) - F(\lambda_2)| < \Delta / \lceil \Delta \zeta^{-1} \rceil \right\}. \quad (32)$$

829 Then there exists $\{a_j\}_{j \in \{0\} \cup \lceil \Delta \zeta^{-1} \rceil - 1}$ and $\{w_j\}_{j \in \lceil \Delta \zeta^{-1} \rceil - 1}$, such that

$$\sup_{\lambda \in \mathbb{R}} |F(\lambda) - f(\lambda)| \leq \zeta, \quad \text{where } f(\lambda) = \sum_{j=1}^{\lceil \Delta \zeta^{-1} \rceil - 1} a_j \text{ReLU}(\lambda - w_j) + a_0. \quad (33)$$

830 Furthermore, we have $\sup_{j \in \lceil \Delta \zeta^{-1} \rceil - 1} |w_j| \leq w_{\zeta}$, $|a_0| \leq \Pi$, and $|a_j| \leq 2\Pi^2$ for all $j \in \lceil \Delta \zeta^{-1} \rceil -$
831 1].

832 *Proof of Lemma 5.* When π_0 is a Dirac measure, we simply take $a_0 = \mathbb{E}[\beta]$. In other cases, one can
 833 verify that $F(\cdot)$ is strictly increasing, hence $\Pi_{\max} > \Pi_{\min}$. Then for any $\alpha \in (\Pi_{\min}, \Pi_{\max})$, there
 834 exists a unique $\mu_\alpha \in \mathbb{R}$, such that $F(\mu_\alpha) = \alpha$.

835 Let $a_0 = \Pi_{\min} + \Delta \lceil \Delta \zeta^{-1} \rceil^{-1}$. For $j \in [\lceil \Delta \zeta^{-1} \rceil - 1]$, we let

$$w_j = \mu_{-\Pi_{\min} + j\Delta / \lceil \Delta \zeta^{-1} \rceil}, \quad a_j = \frac{\Delta}{\lceil \Delta \zeta^{-1} \rceil (w_{j+1} - w_j)} - \frac{\Delta}{\lceil \Delta \zeta^{-1} \rceil (w_j - w_{j-1})}.$$

836 In the above equations, we make the convention that $w_0 = w_{\lceil \Delta \zeta^{-1} \rceil} = \infty$. With
 837 $\{a_j\}_{j \in \{0\} \cup [\lceil \Delta \zeta^{-1} \rceil - 1]}$ and $\{w_j\}_{j \in [\lceil \Delta \zeta^{-1} \rceil - 1]}$ defined as above, one can verify that Eq. (33) is true.
 838 Furthermore, since $\|F'\|_\infty \leq \Pi^2$, we have $|\Delta / \lceil \Delta \zeta^{-1} \rceil (w_{j+1} - w_j)| \leq \Pi^2$ for all possible j . This
 839 gives $|a_j| \leq 2\Pi^2$ for every j .

840 □

841 **Remark 1.** When $\pi_0 = \text{Unif}(\{\pm 1\})$, one can check that for any $\gamma > 0$, we have $F(x) = \tanh(x)$.
 842 In this case, one can verify that $|w_\zeta| \leq \log \lceil \zeta^{-1} \rceil$. In addition, we can further guarantee that
 843 $\sum_{j \in [\lceil \Delta \zeta^{-1} \rceil - 1]} |a_j| \leq 2$.

844 D.5 Approximation error of fixed point iteration

845 **Lemma 6.** Assume that $\mathbf{h} \in \mathbb{R}^d$, $\mathbf{U} \in \mathbb{R}^{d \times d}$ with $\|\mathbf{U}\|_{\text{op}} \leq A < \Pi^{-2}$ for some $\Pi > 0$. Further
 846 assume that $f_* : \mathbb{R} \mapsto \mathbb{R}$ is Π^2 -Lipschitz continuous and $f : \mathbb{R} \mapsto \mathbb{R}$ is a function satisfying

$$\sup_{u \in \mathbb{R}} |f(u) - f_*(u)| \leq \zeta. \quad (34)$$

847 Let $\hat{\mathbf{m}} \in \mathbb{R}^d$ satisfying $\|\hat{\mathbf{m}}\|_2 \leq \Pi\sqrt{d}$ be the unique fixed point of

$$\hat{\mathbf{m}} = f_*(\mathbf{U}\hat{\mathbf{m}} + \mathbf{h}). \quad (35)$$

848 Let $\tilde{\mathbf{m}}^0 = \mathbf{0}$ and

$$\tilde{\mathbf{m}}^k = f(\mathbf{U}\tilde{\mathbf{m}}^{k-1} + \mathbf{h}). \quad (36)$$

849 Then we have

$$\frac{1}{\sqrt{d}} \|\tilde{\mathbf{m}}^k - \hat{\mathbf{m}}\|_2 \leq \Pi \cdot (\Pi^2 A)^k + \frac{\zeta}{1 - \Pi^2 A}. \quad (37)$$

850 *Proof of Lemma 6.* By Eq. (34) and (36), we have

$$\tilde{\mathbf{m}}^k = f_*(\mathbf{U}\tilde{\mathbf{m}}^{k-1} + \mathbf{h}) + \zeta^k,$$

851 where $\|\zeta^k\|_2 \leq \sqrt{d}\zeta$. Comparing with Eq. (35), we get

$$\|\tilde{\mathbf{m}}^k - \hat{\mathbf{m}}\|_2 \leq \Pi^2 \|\mathbf{U}\|_{\text{op}} \|\tilde{\mathbf{m}}^{k-1} - \hat{\mathbf{m}}\|_2 + \|\zeta^k\|_2 \leq \Pi^2 A \cdot \|\tilde{\mathbf{m}}^{k-1} - \hat{\mathbf{m}}\|_2 + \sqrt{d}\zeta.$$

852 By the fact that $\|\tilde{\mathbf{m}}^0 - \hat{\mathbf{m}}\|_2 = \|\hat{\mathbf{m}}\|_2 \leq \Pi\sqrt{d}$, this gives Eq. (37), which concludes the proof of the
 853 lemma. □

854 D.6 Properties of two-phase time discretization scheme

855 The lemma below provides a bound related to the two-phase time discretization scheme that appears
 856 to be useful when deriving the sampling error bound.

857 **Lemma 7.** Consider the two-phase discretization scheme $(\kappa, N_0, N, T, \delta, \{t_k\}_{0 \leq k \leq N})$ and recall
 858 that $\gamma_k = t_{k+1} - t_k$ (Definition 1). Recall the definition $\lambda_t = e^{-t}$ and $\sigma_t^2 = 1 - e^{-2t}$. Then we have

$$\sum_{0 \leq k \leq N-1} \gamma_k \cdot \lambda_{T-t_k}^2 \sigma_{T-t_k}^{-4} \lesssim 1 + \delta^{-1}. \quad (38)$$

859 *Proof of Lemma 7.* Simple algebra yields

$$\sigma_t^{-2} = 1/[1 - e^{-2t}] \leq 10 \cdot [1 \vee (1/t)].$$

860 Note that $T - t_k \leq 1$ for all $k \geq N_0$ and $T - t_k \geq 1$ for all $k \leq N_0 - 1$ (c.f. Definition 1 for N_0).
861 Then the summation in the first phase has bound (we use the fact that $\kappa < 1$)

$$\begin{aligned} & \sum_{0 \leq k \leq N_0 - 1} \gamma_k \cdot \lambda_{T-t_k}^2 \sigma_{T-t_k}^{-4} \leq 100\kappa \sum_{0 \leq k \leq M-1} e^{-2(T-t_k)} \leq 100\kappa e^{-2} \sum_{k \geq 0} e^{-2k\kappa} \\ & \leq 100\kappa e^{-2} \frac{1}{1 - e^{-2\kappa}} \leq 100. \end{aligned}$$

862 Furthermore, the summation in the second phase yields (recall from Definition 1 that for $k \geq N_0$, we
863 have $T - t_{N_0+k} = (1 + \kappa)^{-k}$, $\gamma_{N_0+k} = \kappa/(1 + \kappa)^{k+1}$, and $\delta = (1 + \kappa)^{N_0-N}$)

$$\begin{aligned} & \sum_{N_0 \leq k \leq N-1} \gamma_k \lambda_{T-t_k}^2 \sigma_{T-t_k}^{-4} \leq 100 \sum_{N_0 \leq k \leq N-1} \gamma_k / (T - t_k)^2 \\ & = 100 \sum_{0 \leq k \leq N - N_0 - 1} [\kappa / (1 + \kappa)^{k+1}] \cdot (1 + \kappa)^{2k} = 100 \frac{\kappa}{\delta} \sum_{0 \leq k \leq N - N_0 - 1} (1 + \kappa)^{-2-k} \\ & \leq 100 \frac{\kappa}{\delta} \sum_{k=1}^{\infty} (1 + \kappa)^{-k} = 100/\delta. \end{aligned}$$

864 Combining the two inequalities above proves Eq. (38) and concludes the proof. \square

865 **E Proofs for Section 3: Ising models**

866 **E.1 Proof of Theorem 1**

867 **Approximate the minimizer of the free energy via an iterative algorithm**

868 We first show that we can approximate the minimizer of $\mathcal{F}_t^{\text{VI}}$ using a simple iterative algorithm.
869 Calculating the Hessian of $\mathcal{F}_t^{\text{VI}}$, we obtain

$$\nabla_{\mathbf{m}}^2 \mathcal{F}_t^{\text{VI}}(\mathbf{m}; \mathbf{z}) = \text{diag}\{(1 - m_i^2)_{i \in [d]}\} - \mathbf{A} + \mathbf{K} \succeq (1 - A) \cdot \mathbf{I}_d \succ 0, \quad \forall \mathbf{m} \in [-1, 1]^d,$$

870 where the inequalities are due to the fact that $\text{diag}\{(1 - m_i^2)_{i \in [d]}\} \succeq \mathbf{I}_d$ and the assumption that
871 $\|\mathbf{K} - \mathbf{A}\|_{\text{op}} \leq A < 1$. Therefore, $\mathcal{F}_t^{\text{VI}}(\cdot, \mathbf{z})$ is strongly convex in its first coordinate for all $\mathbf{z} \in \mathbb{R}^d$,
872 hence the critical equation

$$\nabla_{\mathbf{m}} \mathcal{F}_t^{\text{VI}}(\mathbf{m}; \mathbf{z}) = \tanh^{-1}(\mathbf{m}) - \mathbf{A}\mathbf{m} - \lambda_t \sigma_t^{-2} \mathbf{z} + \mathbf{K}\mathbf{m} = \mathbf{0},$$

873 can have at most one solution on $[-1, 1]^d$. Furthermore, $\nabla_{\mathbf{m}} \mathcal{F}_t^{\text{VI}}(\mathbf{m}; \mathbf{z}) = \mathbf{0}$ is equivalent to the
874 fixed point equation

$$\mathbf{m} = \tanh((\mathbf{A} - \mathbf{K})\mathbf{m} + \lambda_t \sigma_t^{-2} \mathbf{z}),$$

875 and $T(\mathbf{m}) = \tanh((\mathbf{A} - \mathbf{K})\mathbf{m} + \lambda_t \sigma_t^{-2} \mathbf{z})$ is a continuous mapping from $[-1, 1]^d$ to itself. Therefore,
876 there exists a solution of $\mathbf{m} = T(\mathbf{m})$ by Brouwer's fixed-point theorem. This implies that the above
877 fixed point equation has a unique solution $\hat{\mathbf{m}}_t(\mathbf{z}) \in [-1, 1]^d$.

878 Take $f : \mathbb{R} \rightarrow \mathbb{R}$ to be the function as derived by Lemma 5 achieving ζ -uniform approximation to
879 $\tanh(\cdot)$. We write $f(x) = \sum_{j=1}^{\lceil 2\zeta^{-1} \rceil - 1} a_j \text{ReLU}(x - w_j) + a_0$. Define iterative algorithm $\{\tilde{\mathbf{m}}^\ell\}_{\ell \geq 0}$
880 by

$$\tilde{\mathbf{m}}^0 = \mathbf{0}, \quad \tilde{\mathbf{m}}^\ell(\mathbf{z}) = \tilde{\mathbf{m}}^\ell = f((\mathbf{A} - \mathbf{K})\tilde{\mathbf{m}}^{\ell-1} + \lambda_t \sigma_t^{-2} \mathbf{z}). \quad (39)$$

881 Then by Lemma 6 with $\Pi = 1$, we obtain that

$$\|\tilde{\mathbf{m}}^\ell(\mathbf{z}) - \hat{\mathbf{m}}_t(\mathbf{z})\|_2 / \sqrt{d} \leq A^\ell + \zeta \cdot (1 - A)^{-1}. \quad (40)$$

882 **Represent the iterative algorithm as a ResNet**

883 Next, we show that $\tilde{\mathbf{m}}^\ell(\mathbf{z})$ defined as above takes the form of a ResNet.

884 **Lemma 8.** For all $\ell \in \mathbb{N}_+$ and $\delta \leq t \leq T$, there exists $\mathbf{W} \in \mathcal{W}_{d,D,\ell,M,B}$ with

$$D = 3d, \quad M = (\lceil 2\zeta^{-1} \rceil + 3)d,$$

$$B = (\lceil 2\zeta^{-1} \rceil - 1)(4 + \log \lceil \zeta^{-1} \rceil) + 8 + (1 - e^{-2\delta})^{-1} + \sqrt{d},$$

885 such that $(\lambda_t \tilde{\mathbf{m}}^\ell(\mathbf{z}) - \mathbf{z})/\sigma_t^2 = \text{ResN}_{\mathbf{W}}(\mathbf{z})$, where $\tilde{\mathbf{m}}^\ell$ is as defined in Eq. (39).

886 *Proof of Lemma 8.* Recall the definition of f as an approximation of \tanh as in Lemma 5. Recall
887 that a ResNet takes the form (ResNet). We shall choose the weight matrices appropriately such that
888 $\mathbf{u}^{(\ell)} = [\tilde{\mathbf{m}}^\ell; \sigma_t^{-2}\mathbf{z}; \mathbf{1}_d]^\top \in \mathbb{R}^{3d}$. In particular, for $\ell = 0$, we set

$$\mathbf{W}_{\text{in}} = \begin{bmatrix} \mathbf{0}_{d \times d} & \sigma_t^{-2}\mathbf{I}_d & \mathbf{0}_{d \times d} \\ \mathbf{0}_{1 \times d} & \mathbf{0}_{1 \times d} & \mathbf{1}_{1 \times d} \end{bmatrix}^\top \in \mathbb{R}^{3d \times (d+1)}.$$

889 For $\ell \geq 1$, we set

$$\mathbf{W}_1^{(\ell)} = \begin{bmatrix} a_i \mathbf{I}_d & \cdots & a_{\lceil 2\zeta^{-1} \rceil - 1} \mathbf{I}_d & -\mathbf{I}_d & \mathbf{I}_d & a_0 \mathbf{I}_d & -a_0 \mathbf{I}_d \\ \mathbf{0}_{d \times d} & \cdots & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} \\ \mathbf{0}_{d \times d} & \cdots & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} \end{bmatrix} \in \mathbb{R}^{3d \times (\lceil 2\zeta^{-1} \rceil + 3)d},$$

$$\mathbf{W}_2^{(\ell)} = \begin{bmatrix} \mathbf{A} - \mathbf{K} & \cdots & \mathbf{A} - \mathbf{K} & \mathbf{I}_d & -\mathbf{I}_d & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} \\ \lambda_t \mathbf{I}_d & \cdots & \lambda_t \mathbf{I}_d & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} \\ -w_1 \mathbf{I}_d & \cdots & -w_{\lceil 2\zeta^{-1} \rceil - 1} \mathbf{I}_d & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{I}_d & -\mathbf{I}_d \end{bmatrix}^\top \in \mathbb{R}^{(\lceil 2\zeta^{-1} \rceil + 3)d \times 3d}.$$

890 Finally, we take $\mathbf{W}_{\text{out}} = [\lambda_t \sigma_t^{-2} \mathbf{I}_d, -\mathbf{I}_d, \mathbf{0}_{d \times d}] \in \mathbb{R}^{d \times 3d}$.

891 By Lemma 5 and Remark 1, we have $\sum_{j=1}^{\lceil 2\zeta^{-1} \rceil - 1} |a_j| \leq 2$, $|a_0| \leq 1$, and $|w_j| \leq \log \lceil \zeta^{-1} \rceil$. Therefore,
892 $\|\mathbf{W}_{\text{in}}\|_{\text{op}} \leq \sqrt{d} + \sigma_t^{-2}$, $\|\mathbf{W}_{\text{out}}\|_{\text{op}} \leq 1 + \lambda_t \sigma_t^{-2}$, $\|\mathbf{W}_1^{(\ell)}\|_{\text{op}} \leq 2\lceil 2\zeta^{-1} \rceil + 2$ and $\|\mathbf{W}_2^{(\ell)}\|_{\text{op}} \leq$
893 $(\lceil 2\zeta^{-1} \rceil - 1)(2 + \log \lceil \zeta^{-1} \rceil) + 4$. Hence, $\|\mathbf{W}\| \leq (\lceil 2\zeta^{-1} \rceil - 1)(4 + \log \lceil \zeta^{-1} \rceil) + 8 + \sigma_t^{-2} + \sqrt{d}$.
894 Note that for $\delta \leq t \leq T$, it holds that $\sigma_t^{-2} \leq (1 - e^{-2\delta})^{-1}$. Therefore, we have

$$\|\mathbf{W}\| \leq B = (\lceil 2\zeta^{-1} \rceil - 1)(3 + \log \lceil \zeta^{-1} \rceil) + 8 + (1 - e^{-2\delta})^{-1} + \sqrt{d}.$$

895 This completes the proof of Lemma 8. \square

896 **Proof of Theorem 1**

897 Recall that we have $\hat{\mathbf{s}}_t(\mathbf{z}) = \text{P}_t[\text{ResN}_{\widehat{\mathbf{W}}}](\mathbf{z})$, where $\widehat{\mathbf{W}} = \text{argmin}_{\mathbf{W} \in \mathcal{W}} \hat{\mathbb{E}}[\|\text{P}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}) +$
898 $\sigma_t^{-1} \mathbf{g}\|_2^2]$ for $\mathcal{W} = \mathcal{W}_{d,D,L,M,B}$. Here, $\hat{\mathbb{E}}$ denotes averaging over the empirical data distribution. By
899 standard error decomposition analysis in empirical risk minimization theory, we have:

$$\mathbb{E}[\|\text{P}_t[\text{ResN}_{\widehat{\mathbf{W}}}](\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d \leq \inf_{\mathbf{W} \in \mathcal{W}} \mathbb{E}[\|\text{P}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d$$

$$+ 2 \sup_{\mathbf{W} \in \mathcal{W}} \left| \hat{\mathbb{E}}[\|\text{P}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d - \mathbb{E}[\|\text{P}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d \right|.$$

900 Furthermore, a standard identity in diffusion model theory shows:

$$\mathbb{E}[\|\hat{\mathbf{s}}_t(\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2]/d = \mathbb{E}[\|\hat{\mathbf{s}}_t(\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d + C, \quad C = \mathbb{E}[\|\mathbf{s}_t(\mathbf{z})\|_2^2]/d - \mathbb{E}[\|\sigma_t^{-1} \mathbf{g}\|_2^2]/d.$$

901 Combining the above yields:

$$\mathbb{E}[\|\hat{\mathbf{s}}_t(\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2]/d \leq \bar{\varepsilon}_{\text{app}}^2 + \bar{\varepsilon}_{\text{gen}}^2, \quad (41)$$

902 where $\bar{\varepsilon}_{\text{app}}^2$ is the approximation error and $\bar{\varepsilon}_{\text{gen}}^2$ is the generalization error,

$$\bar{\varepsilon}_{\text{app}}^2 = \inf_{\mathbf{W} \in \mathcal{W}} \mathbb{E}[\|\text{P}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2]/d,$$

$$\bar{\varepsilon}_{\text{gen}}^2 = 2 \sup_{\mathbf{W} \in \mathcal{W}} \left| \hat{\mathbb{E}}[\|\text{P}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d - \mathbb{E}[\|\text{P}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d \right|.$$

903 By Proposition 6 and take $D = 3d$ and $m = 0$, with probability at least $1 - \eta$, simultaneously for
 904 any $t \in \{T - t_k\}_{0 \leq k \leq N-1}$, we have

$$\bar{\varepsilon}_{\text{gen}}^2 \lesssim \frac{\lambda_t^2}{\sigma_t^4} \cdot \sqrt{\frac{[d^2 + LdM] \cdot [L \cdot \log(LB) + \log(\lambda_t^{-1})] + \log(N/\eta)}{n}}. \quad (42)$$

905 To bound $\bar{\varepsilon}_{\text{app}}^2$, by the identity that $\mathbf{s}_t(\mathbf{z}) = (\lambda_t \mathbf{m}_t(\mathbf{z}) - \mathbf{z})/\sigma_t^2$ and $\text{P}_t \text{ResN}_{\mathbf{W}}(\mathbf{z}) =$
 906 $\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\text{ResN}_{\mathbf{W}}(\mathbf{z}) + \sigma_t^{-2} \mathbf{z}) - \sigma_t^{-2} \mathbf{z}$, recalling $\tilde{\mathbf{m}}^L(\mathbf{z})$ as defined in Eq. (39), and by Lemma 8,
 907 we have

$$\begin{aligned} \bar{\varepsilon}_{\text{app}}^2 &= \inf_{\mathbf{W} \in \mathcal{W}} \mathbb{E}[\|\text{P}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2]/d \\ &= \inf_{\mathbf{W} \in \mathcal{W}} \mathbb{E}[\|\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\text{ResN}_{\mathbf{W}}(\mathbf{z}) + \sigma_t^{-2} \mathbf{z}) - \lambda_t \sigma_t^{-2} \mathbf{m}_t(\mathbf{z})\|_2^2]/d \\ &\leq \mathbb{E}[\|\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \tilde{\mathbf{m}}^L(\mathbf{z})) - \lambda_t \sigma_t^{-2} \mathbf{m}_t(\mathbf{z})\|_2^2]/d \\ &\lesssim \mathbb{E}[\|\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \tilde{\mathbf{m}}^L(\mathbf{z})) - \text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \hat{\mathbf{m}}(\mathbf{z}))\|_2^2]/d \\ &\quad + \mathbb{E}[\|\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \hat{\mathbf{m}}(\mathbf{z})) - \lambda_t \sigma_t^{-2} \mathbf{m}_t(\mathbf{z})\|_2^2]/d \end{aligned} \quad (43)$$

908 where the last inequality uses the triangle inequality. By Eq. (40) and the 1-Lipschitzness of
 909 $\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}$, the first quantity in the right-hand side is controlled by

$$\begin{aligned} &\mathbb{E}[\|\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \tilde{\mathbf{m}}^L(\mathbf{z})) - \text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \hat{\mathbf{m}}(\mathbf{z}))\|_2^2]/d \\ &\lesssim \frac{\lambda_t^2}{\sigma_t^4} \cdot (A^{2L} + \zeta^2(1-A)^{-2}) \lesssim \frac{\lambda_t^2}{\sigma_t^4} \cdot \left(A^{2L} + \frac{d^2}{(1-A)^2 M^2} \right), \end{aligned} \quad (44)$$

910 where the last inequality is by the fact that we can choose ζ such that $M = d \cdot ([2\zeta^{-1}] + 3)$, which
 911 gives $\zeta \leq 6d/M$. Furthermore, by Assumption 1 and by $\|\hat{\mathbf{m}}(\mathbf{z})\|_2 \leq \sqrt{d}$, the second quantity in the
 912 right-hand side is controlled by

$$\mathbb{E}[\|\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \hat{\mathbf{m}}(\mathbf{z})) - \lambda_t \sigma_t^{-2} \mathbf{m}_t(\mathbf{z})\|_2^2]/d \lesssim \frac{\lambda_t^2}{\sigma_t^4} \cdot \varepsilon_{\text{VI},t}^2(\mathbf{A}). \quad (45)$$

913 Combining Eq. (41), (42), (43), (44), (45) completes the proof of Theorem 1.

914 E.2 Proof of Corollary 1

915 Corollary 1 is a direct consequence of Theorem 1, Theorem 5, and Lemma 7.

916 E.3 Proofs for Section A.2

917 E.3.1 Proof of Lemma 1

918 Lemma 1 is a direct consequence of Lemma 9 below. Given Lemma 9, Lemma 1 holds by observing
 919 that when $\|\mathbf{A}\|_{\text{op}} < 1/2$, we have $(1 - \|\mathbf{A}\|_{\text{op}})^{-2} \leq 4$.

920 **Lemma 9.** Let $\mathbf{h} \in \mathbb{R}^d$, $\mathbf{A} \in \mathbb{R}^{d \times d}$ be symmetric with $\|\mathbf{A}\|_{\text{op}} < 1/2$. Consider the Ising model
 921 $\mu(\sigma) \propto \exp\{\langle \sigma, \mathbf{A}\sigma \rangle/2 + \langle \sigma, \mathbf{h} \rangle\}$ and denote $\mathbf{m} = \mathbb{E}_{\sigma \sim \mu}[\sigma]$. Let $\hat{\mathbf{m}}$ be the unique minimizer of
 922 the naive VB free energy

$$\hat{\mathbf{m}} = \operatorname{argmin}_{\mathbf{m} \in [-1, 1]^d} \left\{ \sum_{i=1}^d -\mathfrak{h}_{\text{bin}}(m_i) - \langle \mathbf{m}, \mathbf{A}\mathbf{m} \rangle/2 - \langle \mathbf{m}, \mathbf{h} \rangle \right\}.$$

923 Then we have

$$\frac{1}{d} \|\mathbf{m} - \hat{\mathbf{m}}\|_2^2 \leq \frac{1}{(1 - 2\|\mathbf{A}\|_{\text{op}})(1 - \|\mathbf{A}\|_{\text{op}})^2} \frac{\|\mathbf{A}\|_F^2}{d}.$$

924 *Proof of Lemma 9.* Denote $\ell_i(\sigma) = \sum_{j \neq i} A_{ij} \sigma_j + h_i$. Simple calculations yields $\mathbb{E}_{\mu}[\sigma_i | \{\sigma_j\}_{j \neq i}] =$
 925 $\tanh(\ell_i(\sigma))$, which implies that

$$\mathbb{E}_{\mu}[\sigma_i] = \mathbb{E}_{\mu}[\tanh(\ell_i(\sigma))].$$

926 By the fact that $\sup_{x \in \mathbb{R}} |(d^2/dx^2) \tanh(x)| \leq 1$ and by Taylor's expansion, we have

$$|\mathbb{E}_\mu[\tanh(\ell_i(\boldsymbol{\sigma}))] - \tanh(\mathbb{E}_\mu[\ell_i(\boldsymbol{\sigma})])|^2 \leq \text{Var}_\mu(\ell_i(\boldsymbol{\sigma})).$$

927 By Theorem 1 of Eldan et al. [2022], the Ising model satisfies a Poincare's Inequality with Poincare's
 928 coefficient to be $1/(1 - 2\|\mathbf{A}\|_{\text{op}})$ (we need to translate the Ising model to their setting, which leads
 929 to an additional 2 coefficient in front of $\|\mathbf{A}\|_{\text{op}}$). Therefore, the Poincare's inequality implies that

$$\text{Var}_\mu(\ell_i(\boldsymbol{\sigma})) \leq \frac{1}{1 - 2\|\mathbf{A}\|_{\text{op}}} \sum_{j \neq i} A_{ij}^2.$$

930 Combining the equations above, we get

$$\frac{1}{d} \left\| \mathbf{m} - \tanh(\mathbf{A}\mathbf{m} + \mathbf{h}) \right\|_2^2 = \frac{1}{d} \sum_{i=1}^d (\mathbb{E}_\mu[\sigma_i] - \tanh(\mathbb{E}_\mu[\ell_i(\boldsymbol{\sigma})]))^2 \leq \frac{1}{1 - 2\|\mathbf{A}\|_{\text{op}}} \frac{\|\mathbf{A}\|_F^2}{d} \equiv \varepsilon^2.$$

931 Furthermore, notice that $\hat{\mathbf{m}}$ is the unique minimizer of the naive VB free energy implies that
 932 $\hat{\mathbf{m}} = \tanh(\mathbf{A}\hat{\mathbf{m}} + \mathbf{h})$. Therefore, by the equation above, we get

$$\begin{aligned} \varepsilon &\geq \frac{1}{\sqrt{d}} \left\| (\mathbf{m} - \hat{\mathbf{m}}) - (\tanh(\mathbf{A}\mathbf{m} + \mathbf{h}) - \tanh(\mathbf{A}\hat{\mathbf{m}} + \mathbf{h})) \right\|_2 \\ &\geq \frac{1}{\sqrt{d}} \left(\|\mathbf{m} - \hat{\mathbf{m}}\|_2 - \|\tanh(\mathbf{A}\mathbf{m} + \mathbf{h}) - \tanh(\mathbf{A}\hat{\mathbf{m}} + \mathbf{h})\|_2 \right) \\ &\geq (1 - \|\mathbf{A}\|_{\text{op}}) \cdot \frac{1}{\sqrt{d}} \|\mathbf{m} - \hat{\mathbf{m}}\|_2. \end{aligned}$$

933 Combining the equations above concludes the proof of Lemma 9. □

934 E.3.2 Proof of Lemma 2

935 Lemma 2 is a direct consequence of the lemma below.

936 **Lemma 10** (Lemma 4.10 and Proposition 4.2 of El Alaoui et al. [2022]). *Let $\mathbf{J} \sim \text{GOE}(d)$ and
 937 $\beta < 1/2$. Let $\mathbf{x} \sim \mu(\mathbf{x}) \propto \exp\{\beta\langle \mathbf{x}, \mathbf{J}\mathbf{x} \rangle/2\}$ on $\{\pm 1\}^d$ and $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ independently. Let
 938 $\mathbf{z} = \lambda\mathbf{x} + \sigma\mathbf{g}$. Consider the posterior measure*

$$\mu(\mathbf{x}|\mathbf{z}) \propto \exp\{\beta\langle \mathbf{x}, \mathbf{J}\mathbf{x} \rangle/2 + (\lambda/\sigma^2)\langle \mathbf{x}, \mathbf{z} \rangle\},$$

939 and define $\mathbf{m}(\mathbf{z}) = \sum_{\mathbf{x} \in \{\pm 1\}^d} \mathbf{x} \mu(\mathbf{x}|\mathbf{z})$. Furthermore, consider the TAP free energy

$$\mathcal{F}_{\text{TAP}}(\mathbf{m}; \mathbf{z}, q) = \sum_{i=1}^d -\text{h}_{\text{bin}}(m_i) - \frac{\beta}{2} \langle \mathbf{m}, \mathbf{J}\mathbf{m} \rangle - \frac{\lambda}{\sigma^2} \langle \mathbf{z}, \mathbf{m} \rangle + \frac{\beta^2(1-q)}{2} \|\mathbf{m}\|_2^2,$$

940 take $q_\star = q_\star(\beta, \lambda, \sigma)$ to be the unique solution of

$$q_\star = \mathbb{E}_{G \sim \mathcal{N}(0,1)} \left[\tanh^2(\beta^2 q_\star + (\lambda^2/\sigma^2) + \sqrt{\beta^2 q_\star + (\lambda^2/\sigma^2)} G) \right],$$

941 and define $\hat{\mathbf{m}}(\mathbf{z}) = \text{argmin}_{\mathbf{m} \in [-1,1]^d} \mathcal{F}_{\text{TAP}}(\mathbf{m}; \mathbf{z}, q)$ to be the unique minimizer. Then we have

$$\|\mathbf{m}(\mathbf{z}) - \hat{\mathbf{m}}(\mathbf{z})\|_2^2/d \xrightarrow{p} 0.$$

942 **Remark 2.** We discuss the several seeming differences between Lemma 10 and [El Alaoui et al.,
 943 2022, Lemma 4.10].

944 • The parameter λ^2/σ^2 in Lemma 10 maps to the parameter t in [El Alaoui et al., 2022,
 945 Lemma 4.10]. The variable $(\lambda/\sigma^2)\mathbf{z} = (\lambda^2/\sigma^2)\mathbf{x} + (\lambda/\sigma)\mathbf{g}$ in Lemma 10 maps to the
 946 variable $\mathbf{y} \stackrel{d}{=} t\mathbf{x} + \sqrt{t} \cdot \mathbf{g}$ in [El Alaoui et al., 2022, Lemma 4.10].

947 • Lemma 10 takes $\hat{\mathbf{m}}(\mathbf{z})$ to be the unique global minimizer of $\mathcal{F}_{\text{TAP}}(\mathbf{m}; \mathbf{z}, q_\star)$, whereas
 948 [El Alaoui et al., 2022, Lemma 4.10] takes $\hat{\mathbf{m}}(\mathbf{z})$ to be a particular local minimizer of
 949 $\mathcal{F}_{\text{TAP}}(\mathbf{m}; \mathbf{z}, q_\star)$. However, when $\beta < 1/2$, it can be shown that \mathcal{F}_{TAP} is strongly convex
 950 with high probability, and hence the local minimizer is the global minimizer with high
 951 probability.

952 • [El Alaoui et al., 2022, Lemma 4.10] is proven under a different joint distribution of (\mathbf{J}, \mathbf{z}) .
 953 However, [El Alaoui et al., 2022, Proposition 4.2] shows that the distribution for [El Alaoui
 954 et al., 2022, Lemma 4.10] is contiguous to the distribution for Lemma 10, and hence the
 955 high probability event under the sampling distribution of [El Alaoui et al., 2022, Lemma
 956 4.10] can be translated to the corresponding high probability event under the sampling
 957 distribution of Lemma 10.

958 To prove Lemma 2, we take $c_t = \beta^2(1 - q_t)$ where q_t is the unique solution of

$$q_t = \mathbb{E}_{G \sim \mathcal{N}(0,1)} [\tanh^2(\beta^2 q_t + (\lambda_t^2/\sigma_t^2)) + \sqrt{\beta^2 q_t + (\lambda_t^2/\sigma_t^2)} G].$$

959 Hence by Lemma 10, for any $\beta < 1/2$, we have

$$\mathbb{E}_{\mathbf{z} \sim \mu_t} [\|\hat{\mathbf{m}}_t(\mathbf{z}) - \mathbf{m}_t(\mathbf{z})\|_2^2] / d \xrightarrow{P} 0, \quad d \rightarrow \infty.$$

960 Furthermore, note that $c_t \leq \beta^2$ and $\|\beta \mathbf{J}\|_{\text{op}} \leq 2\beta + \varepsilon$ with high probability for arbitrarily small ε .
 961 This ensures that $\|\beta \mathbf{J} - c_t \mathbf{I}_d\|_{\text{op}} \leq \|\beta \mathbf{J}\|_{\text{op}} + \beta^2 < 1$ when $\beta \leq 1/4$. This proves Lemma 2.

962 F Proofs for Section B: Generalization to other models

963 F.1 Proof of Theorem 2

964 Approximate the minimizer of the free energy via iterative algorithms

965 Once again, we first prove that we can approximately minimize the free energy by implementing a
 966 simple iterative algorithm. Recall that

$$\begin{aligned} \hat{\omega}_t(\mathbf{z}) &= \operatorname{argmin}_{\omega \in [-1,1]^{d+m}} \mathcal{F}_t^{\text{marginal}}(\omega; \mathbf{z}), \\ \mathcal{F}_t^{\text{marginal}}(\omega; \mathbf{z}) &:= \left\{ \sum_{i=1}^{d+m} -\text{h}_{\text{bin}}(\omega_i) - \frac{1}{2} \langle \omega, \mathbf{A} \omega \rangle - \frac{\lambda_t}{\sigma_t^2} \langle \mathbf{z}, \omega_{1:d} \rangle + \frac{1}{2} \langle \omega, \mathbf{K} \omega \rangle \right\}. \end{aligned}$$

967 Taking the gradient and the Hessian of $\mathcal{F}_t^{\text{marginal}}(\omega; \mathbf{z})$, we obtain

$$\begin{aligned} \nabla_{\omega} \mathcal{F}_t^{\text{marginal}}(\omega; \mathbf{z}) &= \tanh^{-1}(\omega) + (\mathbf{K} - \mathbf{A})\omega - \frac{\lambda_t}{\sigma_t^2} [\mathbf{z}; \mathbf{0}_m]^{\top}, \\ \nabla_{\omega}^2 \mathcal{F}_t^{\text{marginal}}(\omega; \mathbf{z}) &= \operatorname{diag}\{(1 - \omega_i^2)^{-1}\}_{i \in [d+m]} + \mathbf{K} - \mathbf{A}. \end{aligned}$$

968 Since $\|\mathbf{K} - \mathbf{A}\|_{\text{op}} \leq A < 1$, we can then conclude that $\nabla_{\omega}^2 \mathcal{F}_t^{\text{marginal}}(\omega; \mathbf{z}) \succeq (1 - A)\mathbf{I}_{d+m}$ for
 969 all $\mathbf{z} \in \mathbb{R}^d$, hence $\mathcal{F}_t^{\text{marginal}}(\cdot; \mathbf{z})$ is strongly-convex for all $\mathbf{z} \in \mathbb{R}^d$. This further implies that the
 970 fixed-point equation below

$$\omega = \tanh \left((\mathbf{A} - \mathbf{K})\omega + \frac{\lambda_t}{\sigma_t^2} [\mathbf{z}; \mathbf{0}_m]^{\top} \right)$$

971 has a unique solution. By Lemma 6, we obtain that if we run the iteration

$$\tilde{\omega}^0(\mathbf{z}) = \mathbf{0}, \quad \tilde{\omega}^k(\mathbf{z}) = f((\mathbf{A} - \mathbf{K})\tilde{\omega}^{k-1}(\mathbf{z}) + \lambda_t \sigma_t^{-2} [\mathbf{z}; \mathbf{0}_m]^{\top}), \quad (46)$$

972 where $\|f(\cdot) - \tanh(\cdot)\|_{\infty} \leq \zeta$, then

$$\frac{1}{\sqrt{d+m}} \|\tilde{\omega}^k(\mathbf{z}) - \hat{\omega}_t(\mathbf{z})\|_2 \leq A^k + \zeta(1 - A)^{-1}. \quad (47)$$

973 In particular, we require that $f(\cdot)$ is the function that we construct in Lemma 5.

974 Represent the iterative algorithm as a ResNet

975 Recall that $\hat{\mathbf{m}}_t(\mathbf{z}) = [\hat{\omega}_t(\mathbf{z})]_{1:d}$. We define $\tilde{\mathbf{m}}^{\ell}(\mathbf{z}) := [\tilde{\omega}^{\ell}(\mathbf{z})]_{1:d}$. In what follows, we show that
 976 $(\lambda_t \tilde{\mathbf{m}}^{\ell}(\mathbf{z}) - \mathbf{z})/\sigma_t^2$ can be expressed as a ResNet that takes input \mathbf{z} .

977 **Lemma 11.** For all $\ell \in \mathbb{N}_+$ and $\delta \leq t \leq T$, there exists $\mathbf{W} \in \mathcal{W}_{d,D,\ell,M,B}$ with

$$\begin{aligned} D &= 3(d+m), \quad M = (\lceil 2\zeta^{-1} \rceil + 1)(d+m), \\ B &= (\lceil 2\zeta^{-1} \rceil - 1) \cdot (\log \lceil \zeta^{-1} \rceil + 4) + 8 + \sqrt{d+m} + (1 - e^{-2\delta})^{-1}, \end{aligned}$$

978 such that $(\lambda_t \tilde{\mathbf{m}}^\ell(\mathbf{z}) - \mathbf{z})/\sigma_t^2 = \text{ResN}_{\mathbf{W}}(\mathbf{z})$, where $\tilde{\mathbf{m}}^\ell$ is as defined in Eq. (46).

979 *Proof of Lemma 11.* Recall the definition of f as an approximation of \tanh as in Lemma 5. The
 980 proof of this lemma is similar to that of Lemma 8. To be specific, we will select the weight
 981 matrices $\{\mathbf{W}_1^{(\ell)}, \mathbf{W}_2^{(\ell)}, \mathbf{W}_{\text{in}}, \mathbf{W}_{\text{out}}\}$ appropriately such that $\mathbf{u}^{(\ell)} = [\tilde{\omega}^\ell(\mathbf{z}); \sigma_t^{-2}[\mathbf{z}; \mathbf{0}_m]^\top; \mathbf{1}_{d+m}]^\top \in$
 982 $\mathbb{R}^{3(d+m)}$. When $\ell = 0$, this can be achieved by setting

$$\mathbf{W}_{\text{in}} = \begin{bmatrix} \mathbf{0}_{d \times (d+m)} & \sigma_t^{-2}[\mathbf{I}_d, \mathbf{0}_{d \times m}] & \mathbf{0}_{d \times (d+m)} \\ \mathbf{0}_{1 \times (d+m)} & \mathbf{0}_{1 \times (d+m)} & \mathbf{1}_{1 \times (d+m)} \end{bmatrix} \in \mathbb{R}^{(d+1) \times 3(d+m)}.$$

983 Also, recall that $f(x) = \sum_{j=1}^{\lceil 2\zeta^{-1} \rceil - 1} \text{ReLU}(x - w_j) + a_0$. Therefore, for $\ell \in \mathbb{N}_+$, we simply set

$$\begin{aligned} \mathbf{W}_1^{(\ell)} &= \begin{bmatrix} a_i \mathbf{I}_d & \cdots & a_{\lceil 2\zeta^{-1} \rceil - 1} \mathbf{I}_{d+m} & -\mathbf{I}_{d+m} & \mathbf{I}_{d+m} & a_0 \mathbf{I}_{d+m} & -a_0 \mathbf{I}_{d+m} \\ \mathbf{0}_{(d+m) \times (d+m)} & \cdots & \mathbf{0}_{(d+m) \times (d+m)} & \mathbf{0}_{(d+m) \times (d+m)} & \mathbf{0}_{(d+m) \times (d+m)} & \mathbf{0}_{(d+m) \times (d+m)} & \mathbf{0}_{(d+m) \times (d+m)} \\ \mathbf{0}_{(d+m) \times (d+m)} & \cdots & \mathbf{0}_{(d+m) \times (d+m)} & \mathbf{0}_{(d+m) \times (d+m)} & \mathbf{0}_{(d+m) \times (d+m)} & \mathbf{0}_{(d+m) \times (d+m)} & \mathbf{0}_{(d+m) \times (d+m)} \end{bmatrix} \\ &\in \mathbb{R}^{3(d+m) \times (\lceil 2\zeta^{-1} \rceil + 3)(d+m)}, \\ \mathbf{W}_2^{(\ell)} &= \begin{bmatrix} \mathbf{A} - \mathbf{K} & \cdots & \mathbf{A} - \mathbf{K} & \mathbf{I}_{d+m} & -\mathbf{I}_{d+m} & \mathbf{0}_{(d+m) \times (d+m)} & \mathbf{0}_{(d+m) \times (d+m)} \\ \lambda_t \mathbf{I}_{d+m} & \cdots & \lambda_t \mathbf{I}_{d+m} & \mathbf{0}_{(d+m) \times (d+m)} & \mathbf{0}_{(d+m) \times (d+m)} & \mathbf{0}_{(d+m) \times (d+m)} & \mathbf{0}_{(d+m) \times (d+m)} \\ -w_1 \mathbf{I}_{d+m} & \cdots & -w_{\lceil 2\zeta^{-1} \rceil - 1} \mathbf{I}_{d+m} & \mathbf{0}_{(d+m) \times (d+m)} & \mathbf{0}_{(d+m) \times (d+m)} & \mathbf{I}_{d+m} & -\mathbf{I}_{d+m} \end{bmatrix}^\top \\ &\in \mathbb{R}^{(\lceil 2\zeta^{-1} \rceil + 3)(d+m) \times 3(d+m)}. \end{aligned}$$

984 Finally, we take $\mathbf{W}_{\text{out}} = [\lambda_t \sigma_t^{-2} \mathbf{I}_d, \mathbf{0}_{d \times m}, -\mathbf{I}_d, \mathbf{0}_{d \times (d+2m)}] \in \mathbb{R}^{d \times 3(d+m)}$.

985 Next, we upper bound the norm of the residual network. By Lemma 5 and Remark 1, we have
 986 $\sum_{j=1}^{\lceil 2\zeta^{-1} \rceil - 1} |a_j| \leq 2$, $|a_0| \leq 1$, $|w_j| \leq \log \lceil \zeta^{-1} \rceil$. Therefore,

$$\begin{aligned} \|\mathbf{W}_{\text{in}}\|_{\text{op}} &\leq \sqrt{d+m} + \sigma_t^{-2}, \quad \|\mathbf{W}_{\text{out}}\|_{\text{op}} \leq 1 + \lambda_t \sigma_t^{-2}, \\ \|\mathbf{W}_1^{(\ell)}\|_{\text{op}} &\leq 2\lceil 2\zeta^{-1} \rceil + 2, \quad \|\mathbf{W}_2^{(\ell)}\|_{\text{op}} \leq (\lceil 2\zeta^{-1} \rceil - 1) \cdot (\log \lceil \zeta^{-1} \rceil + 2) + 4. \end{aligned}$$

987 This implies that

$$\|\mathbf{W}\| \leq B = (\lceil 2\zeta^{-1} \rceil - 1) \cdot (\log \lceil \zeta^{-1} \rceil + 4) + 8 + \sqrt{d+m} + (1 - e^{-2\delta})^{-1}.$$

988 This completes the proof of Lemma 11. \square

989 Proof of Theorem 2

990 Similar to the proof of Theorem 1, we obtain

$$\mathbb{E}[\|\hat{\mathbf{s}}_t(\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2]/d \leq \bar{\varepsilon}_{\text{app}}^2 + \bar{\varepsilon}_{\text{gen}}^2, \quad (48)$$

991 where $\bar{\varepsilon}_{\text{app}}^2$ is the approximation error and $\bar{\varepsilon}_{\text{gen}}^2$ is the generalization error,

$$\begin{aligned} \bar{\varepsilon}_{\text{app}}^2 &= \inf_{\mathbf{W} \in \mathcal{W}} \mathbb{E}[\|\mathbb{P}_t \text{ResN}_{\mathbf{W}}(\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2]/d, \\ \bar{\varepsilon}_{\text{gen}}^2 &= 2 \sup_{\mathbf{W} \in \mathcal{W}} \left| \hat{\mathbb{E}}[\|\mathbb{P}_t \text{ResN}_{\mathbf{W}}(\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d - \mathbb{E}[\|\mathbb{P}_t \text{ResN}_{\mathbf{W}}(\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2]/d \right|. \end{aligned}$$

992 By Proposition 6 and take $D = 3(d+m)$, with probability at least $1 - \eta$, simultaneously for any
 993 $t \in \{T - t_k\}_{0 \leq k \leq N-1}$, we have

$$\bar{\varepsilon}_{\text{gen}}^2 \lesssim \frac{\lambda_t^2}{\sigma_t^4} \cdot \sqrt{\frac{[ML + d](d+m) \cdot [L \cdot \log(LB) + \log(\lambda_t^{-1})] + \log(N/\eta)}{n}}. \quad (49)$$

994 To bound $\varepsilon_{\text{app}}^2$, by the identity that $\mathbf{s}_t(\mathbf{z}) = (\lambda_t \mathbf{m}_t(\mathbf{z}) - \mathbf{z})/\sigma_t^2$ and $\text{P}_t \text{ResN}_{\mathbf{W}}(\mathbf{z}) =$
 995 $\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\text{ResN}_{\mathbf{W}}(\mathbf{z}) + \sigma_t^{-2} \mathbf{z}) - \sigma_t^{-2} \mathbf{z}$, recalling $\tilde{\mathbf{m}}^L(\mathbf{z}) = \tilde{\omega}_{1:d}^L(\mathbf{z})$ as defined in Eq. (46),
 996 and by Lemma 11, we have

$$\begin{aligned} \varepsilon_{\text{app}}^2 &= \inf_{\mathbf{W} \in \mathcal{W}} \mathbb{E}[\|\text{P}_t \text{ResN}_{\mathbf{W}}(\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2]/d \\ &\lesssim \mathbb{E}[\|\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \tilde{\mathbf{m}}^L(\mathbf{z})) - \text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \hat{\mathbf{m}}(\mathbf{z}))\|_2^2]/d \\ &\quad + \mathbb{E}[\|\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \hat{\mathbf{m}}(\mathbf{z})) - \lambda_t \sigma_t^{-2} \mathbf{m}_t(\mathbf{z})\|_2^2]/d. \end{aligned} \quad (50)$$

997 By Eq. (47), the 1-Lipschitzness of $\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}$, and the definition that $\hat{\mathbf{m}}_t(\mathbf{z}) = [\hat{\omega}_t(\mathbf{z})]_{1:d}$ and
 998 $\tilde{\mathbf{m}}^\ell(\mathbf{z}) = [\tilde{\omega}^{(\ell)}(\mathbf{z})]_{1:d}$, the first quantity on the right-hand side is controlled by

$$\begin{aligned} &\mathbb{E}[\|\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \tilde{\mathbf{m}}^L(\mathbf{z})) - \text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \hat{\mathbf{m}}(\mathbf{z}))\|_2^2]/d \\ &\lesssim \frac{d+m}{d} \cdot \frac{\lambda_t^2}{\sigma_t^4} \cdot (A^{2L} + \zeta^2(1-A)^{-2}) \lesssim \frac{d+m}{d} \cdot \frac{\lambda_t^2}{\sigma_t^4} \cdot \left(A^{2L} + \frac{(d+m)^2}{(1-A)^2 M^2} \right), \end{aligned} \quad (51)$$

999 where the last inequality is by the fact that we can choose ζ such that $\zeta \leq 6(d+m)/M$. Furthermore,
 1000 by Assumption 2 and by $\|\hat{\mathbf{m}}(\mathbf{z})\|_2 \leq \sqrt{d}$, the second quantity in the right-hand side is controlled by

$$\mathbb{E}[\|\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \hat{\mathbf{m}}(\mathbf{z})) - \lambda_t \sigma_t^{-2} \mathbf{m}_t(\mathbf{z})\|_2^2]/d \lesssim \frac{\lambda_t^2}{\sigma_t^4} \cdot \varepsilon_{\text{V1},t}^2(\mathbf{A}). \quad (52)$$

1001 Combining Eq. (48), (49), (50), (51), (52) completes the proof of the score estimation result in
 1002 Theorem 2. The KL divergence bound is a direct consequence of score estimation error, Theorem 5,
 1003 and Lemma 7. This concludes the proof.

1004 F2 Proof of Theorem 3

1005 Approximate the minimizer of the free energy via iterative algorithm

1006 We define

$$\mathcal{F}_t^{\text{cond}}(\mathbf{m}; \mathbf{z}, \boldsymbol{\theta}) := \sum_{i=1}^d -\text{h}_{\text{bin}}(m_i) - \frac{1}{2} \langle \mathbf{m}, \mathbf{A}_{11} \mathbf{m} \rangle - \langle \mathbf{m}, \mathbf{A}_{12} \boldsymbol{\theta} \rangle - \frac{\lambda_t}{\sigma_t^2} \langle \mathbf{z}, \mathbf{m} \rangle + \frac{1}{2} \langle \mathbf{m}, \mathbf{K} \mathbf{m} \rangle.$$

1007 Taking the gradient and the Hessian of $\mathcal{F}_t^{\text{cond}}$, we obtain

$$\begin{aligned} \nabla_{\mathbf{m}} \mathcal{F}_t^{\text{cond}}(\mathbf{m}; \mathbf{z}, \boldsymbol{\theta}) &= \tanh^{-1}(\mathbf{m}) + (\mathbf{K} - \mathbf{A}_{11}) \mathbf{m} - \mathbf{A}_{12} \boldsymbol{\theta} - \frac{\lambda_t}{\sigma_t^2} \mathbf{z}, \\ \nabla_{\mathbf{m}}^2 \mathcal{F}_t^{\text{cond}}(\mathbf{m}; \mathbf{z}, \boldsymbol{\theta}) &= \text{diag}\{((1 - m_i^2)^{-1})_{i \in [d]}\} + \mathbf{K} - \mathbf{A}_{11}. \end{aligned}$$

1008 When $\|\mathbf{K} - \mathbf{A}_{11}\|_{\text{op}} \leq A < 1$, we always have $\nabla_{\mathbf{m}}^2 \mathcal{F}_t^{\text{cond}}(\mathbf{m}; \mathbf{z}, \boldsymbol{\theta}) \succeq (1-A)\mathbf{I} \succ 0$. That is to say,
 1009 $\mathcal{F}_t^{\text{cond}}(\cdot; \mathbf{z}, \boldsymbol{\theta})$ is strongly convex, hence

$$\mathbf{m} = \tanh \left((\mathbf{A}_{11} - \mathbf{K}) \mathbf{m} + \mathbf{A}_{12} \boldsymbol{\theta} + \frac{\lambda_t}{\sigma_t^2} \mathbf{z} \right)$$

1010 has a unique solution. We then can apply Lemma 6, and conclude that if we run iteration

$$\tilde{\mathbf{m}}^0(\mathbf{z}; \boldsymbol{\theta}) = \mathbf{0}, \quad \tilde{\mathbf{m}}^\ell(\mathbf{z}; \boldsymbol{\theta}) = f((\mathbf{A}_{11} - \mathbf{K}) \tilde{\mathbf{m}}^{\ell-1}(\mathbf{z}; \boldsymbol{\theta}) + \mathbf{A}_{12} \boldsymbol{\theta} + \lambda_t \sigma_t^{-2} \mathbf{z}) \quad (53)$$

1011 for some $\|f - \tanh\|_\infty \leq \zeta$, it then holds that

$$\frac{1}{\sqrt{d}} \|\tilde{\mathbf{m}}^\ell(\mathbf{z}; \boldsymbol{\theta}) - \hat{\mathbf{m}}_t(\mathbf{z}; \boldsymbol{\theta})\|_2 \leq A^\ell + \zeta(1-A)^{-1}. \quad (54)$$

1012 As usual, we require $f(\cdot)$ satisfies all other conditions from Lemma 5.

1013 **Represent the iterative algorithm as a ResNet**

1014 Next, we show that $(\lambda_t \tilde{\mathbf{m}}^\ell(\mathbf{z}; \boldsymbol{\theta}) - \mathbf{z})/\sigma_t^2$ can be expressed as a ResNet as in (ResNet-Conditional)
 1015 that has input $(\mathbf{z}, \boldsymbol{\theta})$.

1016 **Lemma 12.** For all $\ell \in \mathbb{N}_+$ and $\delta \leq t \leq T$, there exists $\mathbf{W} \in \mathcal{W}_{d,m,D,\ell,M,B}$ with

$$D = 4d, \quad M = (\lceil 2\zeta^{-1} \rceil + 3)d,$$

$$B = (\lceil 2\zeta^{-1} \rceil - 1) \cdot (\log \lceil \zeta^{-1} \rceil + 4 + \|\mathbf{A}_{12}\|_{\text{op}}) + 8 + (1 - e^{-2\delta})^{-1} + \|\mathbf{A}_{12}\|_{\text{op}} + \sqrt{d},$$

1017 such that $(\lambda_t \tilde{\mathbf{m}}^\ell(\mathbf{z}; \boldsymbol{\theta}) - \mathbf{z})/\sigma_t^2 = \text{ResN}_{\mathbf{W}}(\mathbf{z}, \boldsymbol{\theta})$, where $\tilde{\mathbf{m}}^\ell$ is as defined in Eq. (53).

1018 *Proof of Lemma 12.* Recall the definition of f as an approximation of tanh as in Lemma 5. We shall
 1019 choose the weight matrices such that $\mathbf{u}^{(\ell)} = [\tilde{\mathbf{m}}^\ell(\mathbf{z}; \boldsymbol{\theta}); \sigma_t^{-2}\mathbf{z}; \mathbf{A}_{12}\boldsymbol{\theta}; \mathbf{1}_d] \in \mathbb{R}^{4d}$. For $\ell = 0$, we
 1020 simply set

$$\mathbf{W}_{\text{in}} = \begin{bmatrix} \mathbf{0}_{d \times d} & \mathbf{0}_{d \times m} & \mathbf{0}_{d \times 1} \\ \sigma_t^{-2}\mathbf{I}_d & \mathbf{0}_{d \times m} & \mathbf{0}_{d \times 1} \\ \mathbf{0}_{d \times d} & \mathbf{A}_{12} & \mathbf{0}_{d \times 1} \\ \mathbf{0}_{d \times d} & \mathbf{0}_{d \times m} & \mathbf{1}_{d \times 1} \end{bmatrix} \in \mathbb{R}^{4d \times (d+m+1)}.$$

1021 For $\ell \geq 1$, we let

$$\mathbf{W}_1^{(\ell)} = \begin{bmatrix} a_i \mathbf{I}_d & \cdots & a_{\lceil 2\zeta^{-1} \rceil - 1} \mathbf{I}_d & -\mathbf{I}_d & \mathbf{I}_d & a_0 \mathbf{I}_d & -a_0 \mathbf{I}_d \\ \mathbf{0}_{d \times d} & \cdots & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} \\ \mathbf{0}_{d \times d} & \cdots & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} \\ \mathbf{0}_{d \times d} & \cdots & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} \end{bmatrix} \in \mathbb{R}^{4d \times (\lceil 2\zeta^{-1} \rceil + 3)d},$$

$$\mathbf{W}_2^{(\ell)} = \begin{bmatrix} \mathbf{A}_{11} - \mathbf{K} & \cdots & \mathbf{A}_{11} - \mathbf{K} & \mathbf{I}_d & -\mathbf{I}_d & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} \\ \lambda_t \mathbf{I}_d & \cdots & \lambda_t \mathbf{I}_d & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} \\ \mathbf{A}_{12} & \cdots & \mathbf{A}_{12} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} \\ -w_1 \mathbf{I}_d & \cdots & -w_{\lceil 2\zeta^{-1} \rceil - 1} \mathbf{I}_d & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times d} & \mathbf{I}_d & -\mathbf{I}_d \end{bmatrix}^\top \in \mathbb{R}^{(\lceil 2\zeta^{-1} \rceil + 3)d \times 4d}.$$

1022 Finally, we let $\mathbf{W}_{\text{out}} = [\lambda_t \sigma_t^{-2} \mathbf{I}_d, -\mathbf{I}_d, \mathbf{0}_{d \times d}, \mathbf{0}_{d \times d}] \in \mathbb{R}^{d \times 4d}$. By Lemma 5 and Remark 1, we have
 1023 $\sum_{j=1}^{\lceil 2\zeta^{-1} \rceil - 1} |a_j| \leq 2$, $|a_0| \leq 1$, $|w_j| \leq \log \lceil \zeta^{-1} \rceil$. Therefore,

$$\|\mathbf{W}_{\text{out}}\|_{\text{op}} \leq \lambda_t \sigma_t^{-2} + 1, \quad \|\mathbf{W}_{\text{in}}\|_{\text{op}} \leq \sqrt{d} + \sigma_t^{-2} + \|\mathbf{A}_{12}\|_{\text{op}},$$

$$\|\mathbf{W}_1^{(\ell)}\|_{\text{op}} \leq 2\lceil 2\zeta^{-1} \rceil + 2, \quad \|\mathbf{W}_2^{(\ell)}\|_{\text{op}} \leq (\lceil 2\zeta^{-1} \rceil - 1) \cdot (\log \lceil \zeta^{-1} \rceil + 2 + \|\mathbf{A}_{12}\|_{\text{op}}) + 4.$$

1024 As a result, we conclude that

$$\|\mathbf{W}\| \leq B = (\lceil 2\zeta^{-1} \rceil - 1) \cdot (\log \lceil \zeta^{-1} \rceil + 4 + \|\mathbf{A}_{12}\|_{\text{op}}) + 8 + (1 - e^{-2\delta})^{-1} + \|\mathbf{A}_{12}\|_{\text{op}} + \sqrt{d}.$$

1025 We have completed the proof of Lemma 12. \square

1026 **Proof of Theorem 3**

1027 Similar to the proof of Theorem 1, we obtain

$$\mathbb{E}_{\boldsymbol{\theta}, \mathbf{z}} [\|\hat{\mathbf{s}}_t(\mathbf{z}; \boldsymbol{\theta}) - \mathbf{s}_t(\mathbf{z}; \boldsymbol{\theta})\|_2^2] / d \leq \bar{\varepsilon}_{\text{app}}^2 + \bar{\varepsilon}_{\text{gen}}^2, \quad (55)$$

1028 where $\bar{\varepsilon}_{\text{app}}^2$ is the approximation error and $\bar{\varepsilon}_{\text{gen}}^2$ is the generalization error,

$$\bar{\varepsilon}_{\text{app}}^2 = \inf_{\mathbf{W} \in \mathcal{W}} \mathbb{E}_{\boldsymbol{\theta}, \mathbf{z}} [\|\mathbb{P}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}, \boldsymbol{\theta}) - \mathbf{s}_t(\mathbf{z}; \boldsymbol{\theta})\|_2^2] / d,$$

$$\bar{\varepsilon}_{\text{gen}}^2 = 2 \sup_{\mathbf{W} \in \mathcal{W}} \left| \hat{\mathbb{E}} [\|\mathbb{P}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}, \boldsymbol{\theta}) + \sigma_t^{-1} \mathbf{g}\|_2^2] / d - \mathbb{E}_{\boldsymbol{\theta}, \mathbf{z}} [\|\mathbb{P}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}, \boldsymbol{\theta}) + \sigma_t^{-1} \mathbf{g}\|_2^2] / d \right|.$$

1029 By Proposition 6 and take $D = 4d$, with probability at least $1 - \eta$, simultaneously for any $t \in$
 1030 $\{T - t_k\}_{0 \leq k \leq N-1}$, we have

$$\bar{\varepsilon}_{\text{gen}}^2 \lesssim \frac{\lambda_t^2}{\sigma_t^4} \cdot \sqrt{\frac{(MdL + d(d+m)) \cdot [L \cdot \log(LBd^{-1}(m+d)) + \log(\lambda_t^{-1})] + \log(N/\eta)}{n}}. \quad (56)$$

1031 To bound $\bar{\varepsilon}_{\text{app}}^2$, by the identity that $\mathbf{s}_t(\mathbf{z}; \boldsymbol{\theta}) = (\lambda_t \mathbf{m}_t(\mathbf{z}; \boldsymbol{\theta}) - \mathbf{z})/\sigma_t^2$ and $\text{P}_t \text{ResN}_{\mathbf{W}}(\mathbf{z}, \boldsymbol{\theta}) =$
 1032 $\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\text{ResN}_{\mathbf{W}}(\mathbf{z}, \boldsymbol{\theta}) + \sigma_t^{-2} \mathbf{z}) - \sigma_t^{-2} \mathbf{z}$, recalling $\tilde{\mathbf{m}}^L(\mathbf{z})$ as defined in Eq. (53), and by
 1033 Lemma 12, we have

$$\begin{aligned} \bar{\varepsilon}_{\text{app}}^2 &= \inf_{\mathbf{W} \in \mathcal{W}} \mathbb{E}_{\boldsymbol{\theta}, \mathbf{z}} [\|\text{P}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}, \boldsymbol{\theta}) - \mathbf{s}_t(\mathbf{z}; \boldsymbol{\theta})\|_2^2] / d \\ &\lesssim \mathbb{E}_{\boldsymbol{\theta}, \mathbf{z}} [\|\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \tilde{\mathbf{m}}^L(\mathbf{z}; \boldsymbol{\theta})) - \text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \hat{\mathbf{m}}(\mathbf{z}; \boldsymbol{\theta}))\|_2^2] / d \\ &\quad + \mathbb{E}_{\boldsymbol{\theta}, \mathbf{z}} [\|\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \hat{\mathbf{m}}(\mathbf{z}; \boldsymbol{\theta})) - \lambda_t \sigma_t^{-2} \mathbf{m}_t(\mathbf{z}; \boldsymbol{\theta})\|_2^2] / d. \end{aligned} \quad (57)$$

1034 By Eq. (54) and the 1-Lipschitzness of $\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}$, the first quantity on the right-hand side is
 1035 controlled by

$$\begin{aligned} &\mathbb{E}_{\boldsymbol{\theta}, \mathbf{z}} [\|\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \tilde{\mathbf{m}}^L(\mathbf{z}; \boldsymbol{\theta})) - \text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \hat{\mathbf{m}}(\mathbf{z}; \boldsymbol{\theta}))\|_2^2] / d \\ &\lesssim \frac{\lambda_t^2}{\sigma_t^4} \cdot (A^{2L} + \zeta^2(1-A)^{-2}) \lesssim \frac{\lambda_t^2}{\sigma_t^4} \cdot \left(A^{2L} + \frac{d^2}{(1-A)^2 M^2} \right), \end{aligned} \quad (58)$$

1036 where the last inequality is by the fact that we can choose ζ such that $\zeta \leq 6d/M$. Furthermore, by
 1037 Assumption 3 and by $\|\hat{\mathbf{m}}(\mathbf{z}; \boldsymbol{\theta})\|_2 \leq \sqrt{d}$, the second quantity in the right-hand side is controlled by

$$\mathbb{E}_{\boldsymbol{\theta}, \mathbf{z}} [\|\text{proj}_{\lambda_t \sigma_t^{-2} \sqrt{d}}(\lambda_t \sigma_t^{-2} \hat{\mathbf{m}}(\mathbf{z}; \boldsymbol{\theta})) - \lambda_t \sigma_t^{-2} \mathbf{m}_t(\mathbf{z}; \boldsymbol{\theta})\|_2^2] / d \lesssim \frac{\lambda_t^2}{\sigma_t^4} \cdot \varepsilon_{\text{VI}, t}^2(\mathbf{A}). \quad (59)$$

1038 Combining Eq. (55), (56), (57), (58), (59) completes the proof of the score estimation result in
 1039 Theorem 3. The KL divergence bound is a direct consequence of score estimation error, Theorem 5,
 1040 and Lemma 7. This concludes the proof.

1041 To prove the second result of the bound of the expected KL divergence, we simply notice that by
 1042 Theorem 5, conditioning on every $\boldsymbol{\theta}$ we have

$$\frac{1}{d} \text{KL}(\mu_{\delta}(\cdot | \boldsymbol{\theta}), \hat{\mu}(\cdot | \boldsymbol{\theta})) \lesssim \varepsilon^2 + \kappa^2 N + \kappa T + e^{-2T},$$

1043 where

$$\varepsilon^2 = \frac{1}{d} \sum_{k=0}^{N-1} \gamma_k \mathbb{E} [\|\hat{\mathbf{s}}_{T-t_k}(\mathbf{z}; \boldsymbol{\theta}) - \mathbf{s}_{T-t_k}(\mathbf{z}; \boldsymbol{\theta})\|_2^2 | \boldsymbol{\theta}].$$

1044 The proof is complete of Theorem 3 by simply integrating over $\boldsymbol{\theta}$.

1045 F.3 Proof of Theorem 4

1046 Relationship of the score function \mathbf{s}_t to the denoiser \mathbf{e}_t

1047 We first compute the score function $\mathbf{s}_t(\mathbf{z}) = \nabla_{\mathbf{z}} \log \mu_t(\mathbf{z})$, for $\mathbf{x} = \mathbf{A}\boldsymbol{\theta} + \boldsymbol{\varepsilon}$ and $\mathbf{z} = \lambda_t \mathbf{x} + \sigma_t \mathbf{g}$,
 1048 where $\mathbf{g} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ is independent of $(\boldsymbol{\theta}, \boldsymbol{\varepsilon}) \sim \pi_0^m \otimes \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I}_d)$. Note that

$$\begin{aligned} \mathbb{E}[\mathbf{x} | \mathbf{z}] &= \mathbb{E}[\mathbf{A}\boldsymbol{\theta} + \boldsymbol{\varepsilon} | \lambda_t \mathbf{A}\boldsymbol{\theta} + \lambda_t \boldsymbol{\varepsilon} + \sigma_t \mathbf{g}] = \mathbf{A} \mathbb{E}[\boldsymbol{\theta} | \mathbf{z}] + \mathbb{E}[\boldsymbol{\varepsilon} | \lambda_t \mathbf{A}\boldsymbol{\theta} + \lambda_t \boldsymbol{\varepsilon} + \sigma_t \mathbf{g}] \\ &= \mathbf{A} \mathbb{E}[\boldsymbol{\theta} | \mathbf{z}] + \frac{\lambda_t \tau^2}{\lambda_t^2 \tau^2 + \sigma_t^2} \mathbb{E}[\lambda_t \boldsymbol{\varepsilon} + \sigma_t \mathbf{g} | \lambda_t \mathbf{A}\boldsymbol{\theta} + \lambda_t \boldsymbol{\varepsilon} + \sigma_t \mathbf{g}] \\ &= \mathbf{A} \mathbb{E}[\boldsymbol{\theta} | \mathbf{z}] + \frac{\lambda_t \tau^2}{\lambda_t^2 \tau^2 + \sigma_t^2} \cdot (\mathbf{z} - \lambda_t \mathbf{A} \mathbb{E}[\boldsymbol{\theta} | \mathbf{z}]) = \frac{\sigma_t^2}{\lambda_t^2 \tau^2 + \sigma_t^2} \mathbf{A} \mathbb{E}[\boldsymbol{\theta} | \mathbf{z}] + \frac{\lambda_t \tau^2}{\lambda_t^2 \tau^2 + \sigma_t^2} \mathbf{z}. \end{aligned}$$

1049 By (Denoiser), we obtain

$$\mathbf{s}_t(\mathbf{z}) = \frac{\lambda_t}{\sigma_t^2} \mathbb{E}[\mathbf{x} | \mathbf{z}] - \frac{1}{\sigma_t^2} \mathbf{z} = -\frac{1}{\tau^2 \lambda_t^2 + \sigma_t^2} \mathbf{z} + \frac{\lambda_t}{\tau^2 \lambda_t^2 + \sigma_t^2} \mathbf{A} \cdot \mathbb{E}[\boldsymbol{\theta} | \mathbf{z}].$$

1050 We notice the equality in distribution $\mathbf{z}/\lambda_t \stackrel{d}{=} \mathbf{z}_* = \mathbf{A}\boldsymbol{\theta} + \bar{\boldsymbol{\varepsilon}}$ where $(\boldsymbol{\theta}, \bar{\boldsymbol{\varepsilon}}) \sim \pi_0^m \otimes \mathcal{N}(\mathbf{0}, \bar{\tau}_t^2 \mathbf{I}_d)$ (this
 1051 \mathbf{z}_* is as defined in Assumption 4). This implies

$$\mathbf{s}_t(\mathbf{z}) = -\frac{1}{\tau^2 \lambda_t^2 + \sigma_t^2} \mathbf{z} + \frac{\lambda_t}{\tau^2 \lambda_t^2 + \sigma_t^2} \mathbf{A} \cdot \mathbf{e}_t(\mathbf{z}/\lambda_t), \quad (60)$$

1052 where \mathbf{e}_t is as defined in Eq. (15).

1053 **Existence of a unique minimizer of the VI free energy**

1054 We analyze the VI free energy. We define

$$\mathcal{F}_t^{\text{sparse}}(\mathbf{e}; \mathbf{z}_*) := \sum_{i=1}^m \max_{\lambda} \left[\lambda e_i - \log \mathbb{E}_{\beta \sim \pi_0} [e^{\lambda \beta - \beta^2 \nu_t / 2}] \right] + \frac{1}{2\bar{\tau}_t^2} \|\mathbf{z}_* - \mathbf{A}\mathbf{e}\|_2^2 - \frac{1}{2} \langle \mathbf{e}, \mathbf{K}_t \mathbf{e} \rangle.$$

1055 Let $G_t(\lambda) = \log \mathbb{E}_{\beta \sim \pi_0} [e^{\lambda \beta - \beta^2 \nu_t / 2}]$, and $\lambda_i = \arg \max_{\lambda} [\lambda e_i - G_t(\lambda)]$, then $e_i = G'_t(\lambda_i)$. There-
1056 fore,

$$\begin{aligned} \frac{d}{de_i} [\lambda_i e_i - G_t(\lambda_i)] &= \lambda_i + \frac{e_i}{G'_t(\lambda_i)} - \frac{G'_t(\lambda_i)}{G''_t(\lambda_i)} = \lambda_i, \\ \frac{d^2}{d^2 e_i} [\lambda_i e_i - G_t(\lambda_i)] &= \frac{1}{G''_t(\lambda_i)}. \end{aligned}$$

1057 Hence, we have

$$\begin{aligned} \nabla_{\mathbf{e}} \mathcal{F}_t^{\text{sparse}}(\mathbf{e}; \mathbf{z}_*) &= (G'_t)^{-1}(\mathbf{e}) - \frac{1}{\bar{\tau}_t^2} \mathbf{A}^\top \mathbf{z}_* + \frac{1}{\bar{\tau}_t^2} \mathbf{A}^\top \mathbf{A} \mathbf{e} - \mathbf{K}_t \mathbf{e}, \\ \nabla_{\mathbf{e}}^2 \mathcal{F}_t^{\text{sparse}}(\mathbf{e}; \mathbf{z}_*) &= \text{diag}\{(G''_t(\lambda_i)^{-1})_{i \in [m]}\} + \frac{1}{\bar{\tau}_t^2} \mathbf{A}^\top \mathbf{A} - \mathbf{K}_t. \end{aligned}$$

1058 Note that $G''_t(\lambda_i) = \text{Var}_{(\beta, z) \sim \pi_0 \otimes \mathcal{N}(0, 1)} [\beta \mid \beta + \nu_t^{-1/2} z = \lambda \nu_t^{-1}] \leq \Pi^2$. In addition, note that
1059 $|G'_t(\lambda)| = |\mathbb{E}[\beta \mid \beta + \nu_t^{-1/2} z = \lambda \nu_t^{-1}]| \leq \Pi$ for all λ . By assumption, $\|\bar{\tau}_t^{-2} \mathbf{A}^\top \mathbf{A} - \mathbf{K}_t\|_{\text{op}} < \Pi^{-2}$,
1060 hence $\nabla_{\mathbf{e}}^2 \mathcal{F}_t^{\text{sparse}}(\mathbf{e}; \mathbf{z}_*)$ is positive-definite and $\mathcal{F}_t^{\text{sparse}}(\cdot; \mathbf{z}_*)$ as a function of \mathbf{e} is strongly convex.
1061 That is to say, the equation

$$\mathbf{e} = G'_t \left((-\bar{\tau}_t^{-2} \mathbf{A}^\top \mathbf{A} + \mathbf{K}_t) \mathbf{e} + \bar{\tau}_t^{-2} \mathbf{A}^\top \mathbf{z}_* \right)$$

1062 has a unique fixed point $\hat{\mathbf{e}}_t(\mathbf{z}_*)$.

1063 **Approximate the minimizer of the free energy via iterative algorithm**

1064 We denote by $f_t(\cdot)$ the function obtained from Lemma 5 that achieves ζ -uniform approximation to
1065 $G'_t(\cdot)$. By Lemma 6, we conclude that if we implement the following iteration

$$\tilde{\mathbf{e}}^0(\mathbf{z}_*) = \mathbf{0}, \quad \tilde{\mathbf{e}}^{\ell+1}(\mathbf{z}_*) = f_t \left((-\bar{\tau}_t^{-2} \mathbf{A}^\top \mathbf{A} + \mathbf{K}_t) \tilde{\mathbf{e}}^\ell(\mathbf{z}_*) + \bar{\tau}_t^{-2} \mathbf{A}^\top \mathbf{z}_* \right), \quad (61)$$

1066 then for all $\ell \in \mathbb{N}_+$, we have

$$\frac{1}{\sqrt{m}} \|\tilde{\mathbf{e}}^\ell(\mathbf{z}_*) - \hat{\mathbf{e}}_t(\mathbf{z}_*)\|_2 \leq \Pi \cdot (\Pi^2 A)^\ell + \frac{\zeta}{1 - \Pi^2 A}. \quad (62)$$

1067 **Represent the iterative algorithm as a ResNet**

1068 We then show that $\mathbf{s}_t(\mathbf{z}) = (\lambda_t \mathbf{A} \mathbf{e}_t(\mathbf{z}/\lambda_t) - \mathbf{z}) / (\tau^2 \lambda_t^2 + \sigma_t^2)$ (c.f. Eq. (60)) can be expressed as a
1069 ResNet that takes input \mathbf{z} .

1070 **Lemma 13.** For all $t \in \{T - t_k\}_{0 \leq k \leq N-1}$ and $\ell \in \mathbb{N}_+$, there exists $\mathbf{W} \in \mathcal{W}_{d, D, \ell, M, B}$, with

$$D = 3m + d, \quad M = (\lceil 2\Pi\zeta^{-1} \rceil + 3)m,$$

$$B = (\lceil 2\Pi\zeta^{-1} \rceil - 1) \cdot (A + 1 + 2\Pi^2 + w_\zeta) + 2\Pi + 6 + (\|\mathbf{A}\|_{\text{op}} + 1) / (1 - e^{-2\delta}) + \bar{\tau}_t^{-2} \lambda_t^{-1} \|\mathbf{A}\|_{\text{op}} + \sqrt{m},$$

1071 such that $(\lambda_t \mathbf{A} \tilde{\mathbf{e}}^\ell(\mathbf{z}/\lambda_t) - \mathbf{z}) / (\tau^2 \lambda_t^2 + \sigma_t^2) = \text{ResN}_{\mathbf{W}}(\mathbf{z})$. Here, $\tilde{\mathbf{e}}^\ell$ is as defined in Eq. (61), and
1072 w_ζ is given by

$$w_\zeta = \sup_{t \in \{T - t_k\}_{0 \leq k \leq N-1}} \inf \left\{ w : \text{for all } \lambda_1 > \lambda_2 \geq w \text{ or } \lambda_1 < \lambda_2 \leq -w \text{ we have } |G'_t(\lambda_1) - G'_t(\lambda_2)| < \zeta \right\}.$$

1073 *Proof of Lemma 13.* Recall that the ResNet is defined as (ResNet). Recall the definition of f_t as an
1074 approximation of G'_t as in Lemma 5. We shall choose the weight matrices appropriately, such that

1075 $\mathbf{u}^{(\ell)} = [\tilde{\mathbf{e}}^\ell(\mathbf{z}/\lambda_t); \bar{\tau}_t^{-2} \mathbf{A}^\top \mathbf{z}/\lambda_t; \mathbf{1}_m; \mathbf{z}] \in \mathbb{R}^{3m+d}$. For $\ell = 0$, we set

$$\mathbf{W}_{\text{in}} = \begin{bmatrix} \mathbf{0}_{d \times m} & \bar{\tau}_t^{-2} \lambda_t^{-1} \mathbf{A} & \mathbf{0}_{d \times m} & \mathbf{I}_d \\ \mathbf{0}_{1 \times m} & \mathbf{0}_{1 \times m} & \mathbf{1}_{1 \times m} & \mathbf{0}_{1 \times d} \end{bmatrix}^\top \in \mathbb{R}^{(3m+d) \times (d+1)}.$$

1076 For $\ell \geq 1$, we set

$$\mathbf{W}_1^{(\ell)} = \begin{bmatrix} a_i \mathbf{I}_m & \cdots & a_{\lceil 2\Pi\zeta^{-1} \rceil - 1} \mathbf{I}_m & -\mathbf{I}_m & \mathbf{I}_m & a_0 \mathbf{I}_m & -a_0 \mathbf{I}_m \\ \mathbf{0}_{m \times m} & \cdots & \mathbf{0}_{m \times m} & \mathbf{0}_{m \times m} & \mathbf{0}_{m \times m} & \mathbf{0}_{m \times m} & \mathbf{0}_{m \times m} \\ \mathbf{0}_{m \times m} & \cdots & \mathbf{0}_{m \times m} & \mathbf{0}_{m \times m} & \mathbf{0}_{m \times m} & \mathbf{0}_{m \times m} & \mathbf{0}_{m \times m} \\ \mathbf{0}_{d \times m} & \cdots & \mathbf{0}_{d \times m} & \mathbf{0}_{d \times m} & \mathbf{0}_{d \times m} & \mathbf{0}_{d \times m} & \mathbf{0}_{d \times m} \end{bmatrix} \in \mathbb{R}^{(3m+d) \times (\lceil 2\Pi\zeta^{-1} \rceil + 3)m},$$

$$\mathbf{W}_2^{(\ell)} = \begin{bmatrix} -\bar{\tau}_t^{-2} \mathbf{A}^\top \mathbf{A} + \mathbf{K}_t & \cdots & -\bar{\tau}_t^{-2} \mathbf{A}^\top \mathbf{A} + \mathbf{K}_t & \mathbf{I}_m & -\mathbf{I}_m & \mathbf{0}_{m \times m} & \mathbf{0}_{m \times m} \\ \mathbf{I}_m & \cdots & \mathbf{I}_m & \mathbf{0}_{m \times m} & \mathbf{0}_{m \times m} & \mathbf{0}_{m \times m} & \mathbf{0}_{m \times m} \\ -w_1 \mathbf{I}_m & \cdots & -w_{\lceil 2\Pi\zeta^{-1} \rceil - 1} \mathbf{I}_m & \mathbf{0}_{m \times m} & \mathbf{0}_{m \times m} & \mathbf{I}_m & -\mathbf{I}_m \\ \mathbf{0}_{d \times m} & \cdots & \mathbf{0}_{d \times m} & \mathbf{0}_{d \times m} & \mathbf{0}_{d \times m} & \mathbf{0}_{d \times m} & \mathbf{0}_{d \times m} \end{bmatrix}^\top \in \mathbb{R}^{(\lceil 2\Pi\zeta^{-1} \rceil + 3)m \times (3m+d)}.$$

1077 For the output layer, we let $\mathbf{W}_{\text{out}} = [\lambda_t \mathbf{A} / (\sigma_t^2 + \tau^2 \lambda_t^2), \mathbf{0}_{d \times m}, \mathbf{0}_{d \times m}, -(\tau^2 \lambda_t^2 + \sigma_t^2)^{-1} \mathbf{I}_d] \in$
 1078 $\mathbb{R}^{d \times (3m+d)}$.

1079 The following upper bounds are straightforward:

$$\|\mathbf{W}_{\text{in}}\|_{\text{op}} \leq \bar{\tau}_t^{-2} \lambda_t^{-1} \|\mathbf{A}\|_{\text{op}} + \sqrt{m} + 1, \quad \|\mathbf{W}_{\text{out}}\|_{\text{op}} \leq (\|\mathbf{A}\|_{\text{op}} + 1) / (1 - e^{-2\delta}),$$

$$\|\mathbf{W}_1^{(\ell)}\|_{\text{op}} \leq 2\Pi^2 (\lceil 2\Pi\zeta^{-1} \rceil - 1) + 2\Pi + 2, \quad \|\mathbf{W}_2^{(\ell)}\|_{\text{op}} \leq (\lceil 2\Pi\zeta^{-1} \rceil - 1) \cdot (A + 1 + w_\zeta) + 4.$$

In summary, we have

$$\|\mathbf{W}\| \leq (\lceil 2\Pi\zeta^{-1} \rceil - 1) \cdot (A + 1 + 2\Pi^2 + w_\zeta) + 2\Pi + 6 + (\|\mathbf{A}\|_{\text{op}} + 1) / (1 - e^{-2\delta}) + \bar{\tau}_t^{-2} \lambda_t^{-1} \|\mathbf{A}\|_{\text{op}} + \sqrt{m}.$$

1080 This concludes the proof of Lemma 13. \square

1081 Proof of Theorem 4

1082 Similar to the proof of Theorem 1, we obtain

$$\mathbb{E}[\|\hat{\mathbf{s}}_t(\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2] / d \leq \bar{\varepsilon}_{\text{app}}^2 + \bar{\varepsilon}_{\text{gen}}^2, \quad (63)$$

1083 where $\bar{\varepsilon}_{\text{app}}^2$ is the approximation error and $\bar{\varepsilon}_{\text{gen}}^2$ is the generalization error:

$$\bar{\varepsilon}_{\text{app}}^2 = \inf_{\mathbf{W} \in \mathcal{W}} \mathbb{E}[\|\bar{\mathbf{P}}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2] / d,$$

$$\bar{\varepsilon}_{\text{gen}}^2 = 2 \sup_{\mathbf{W} \in \mathcal{W}} \left| \mathbb{E}[\|\hat{\mathbf{P}}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2] / d - \mathbb{E}[\|\mathbf{P}_t \text{ResN}_{\mathbf{W}}(\mathbf{z}) + \sigma_t^{-1} \mathbf{g}\|_2^2] / d \right|.$$

1084 Applying Proposition 7 and taking $D = 3m + d$, we conclude that with probability at least $1 - \eta$,
 1085 simultaneously for any $t \in \{T - t_k\}_{0 \leq k \leq N-1}$, when $n \geq \log(2/\eta)$, we have

$$\bar{\varepsilon}_{\text{gen}}^2 \lesssim \left(\lambda_t^2 \|\mathbf{A}\|_{\text{op}}^2 \Pi^2 (\tau^{-4} + 1) \frac{m}{d} + \frac{\lambda_t^2}{\sigma_t^2} (1 + \tau^2) \right)$$

$$\times \sqrt{\frac{(dD + LDM) \cdot [T + L \log(LB) + \log(nmT(\tau + 1)(\|\mathbf{A}\|_{\text{op}} \Pi + 1)\tau^{-1})] + \log(2N/\eta)}{n}}.$$
(64)

1086 where we choose

$$B = M/m \cdot (A + 1 + 2\Pi^2 + w_\star) + 2\Pi + 6 + (\|\mathbf{A}\|_{\text{op}} + 1) / (1 - e^{-2\delta}) + \tau^{-2} \|\mathbf{A}\|_{\text{op}} + \sqrt{m},$$

$$w_\star = \sup_{t \in \{T - t_k\}_{0 \leq k \leq N-1}} \inf \left\{ w : \text{for all } \lambda_1 > \lambda_2 \geq w \text{ or } \lambda_1 < \lambda_2 \leq -w, |G'_t(\lambda_1) - G'_t(\lambda_2)| < M / (6m\Pi) \right\}.$$
(65)

1087 We next upper bound $\bar{\varepsilon}_{\text{app}}^2$. Recall Eq. (60) and $\bar{\tau}_t^2 = \tau^2 + \sigma_t^2 / \lambda_t^2$, we have $\mathbf{s}_t(\mathbf{z}) =$
 1088 $-\lambda_t^{-2} \bar{\tau}_t^{-2} \mathbf{z} + \lambda_t^{-1} \bar{\tau}_t^{-2} \mathbf{A} \bar{\mathbf{e}}_t(\mathbf{z}_*)$ (recall that $\mathbf{z}_* = \mathbf{z} / \lambda_t$) and recall $\bar{\mathbf{P}}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}) =$
 1089 $\text{proj}_{\sqrt{m} \|\mathbf{A}\|_{\text{op}} \Pi \lambda_t^{-1} \bar{\tau}_t^{-2}}(\text{ResN}_{\mathbf{W}}(\mathbf{z}) + \lambda_t^{-2} \bar{\tau}_t^{-2} \mathbf{z}) - \lambda_t^{-2} \bar{\tau}_t^{-2} \mathbf{z}$. According to Lemma 13, recalling $\tilde{\varepsilon}^L$

1090 as defined in Eq. (61), we have

$$\begin{aligned}
\bar{\varepsilon}_{\text{app}}^2 &= \inf_{\mathbf{W} \in \mathcal{W}} \mathbb{E}[\|\bar{\mathbf{P}}_t[\text{ResN}_{\mathbf{W}}](\mathbf{z}) - \mathbf{s}_t(\mathbf{z})\|_2^2]/d \\
&= \inf_{\mathbf{W} \in \mathcal{W}} \mathbb{E}[\|\text{proj}_{\sqrt{m}\|\mathbf{A}\|_{\text{op}}\lambda_t^{-1}\bar{\tau}_t^{-2}}(\text{ResN}_{\mathbf{W}}(\mathbf{z}) + \lambda_t^{-2}\bar{\tau}_t^{-2}\mathbf{z}) - \lambda_t^{-1}\bar{\tau}_t^{-2}\mathbf{A}\bar{\mathbf{e}}_t(\mathbf{z}_*)\|_2^2]/d \\
&\leq \mathbb{E}[\|\text{proj}_{\sqrt{m}\|\mathbf{A}\|_{\text{op}}\lambda_t^{-1}\bar{\tau}_t^{-2}}(\lambda_t^{-1}\bar{\tau}_t^{-2}\mathbf{A}\tilde{\mathbf{e}}^L(\mathbf{z}_*)) - \lambda_t^{-1}\bar{\tau}_t^{-2}\mathbf{A}\bar{\mathbf{e}}_t(\mathbf{z}_*)\|_2^2]/d \\
&\lesssim \mathbb{E}[\|\text{proj}_{\sqrt{m}\|\mathbf{A}\|_{\text{op}}\lambda_t^{-1}\bar{\tau}_t^{-2}}(\lambda_t^{-1}\bar{\tau}_t^{-2}\mathbf{A}\tilde{\mathbf{e}}^L(\mathbf{z}_*)) - \text{proj}_{\sqrt{m}\|\mathbf{A}\|_{\text{op}}\lambda_t^{-1}\bar{\tau}_t^{-2}}(\lambda_t^{-1}\bar{\tau}_t^{-2}\mathbf{A}\hat{\mathbf{e}}(\mathbf{z}_*))\|_2^2]/d \\
&\quad + \mathbb{E}[\|\text{proj}_{\sqrt{m}\|\mathbf{A}\|_{\text{op}}\lambda_t^{-1}\bar{\tau}_t^{-2}}(\lambda_t^{-1}\bar{\tau}_t^{-2}\mathbf{A}\hat{\mathbf{e}}(\mathbf{z}_*)) - \lambda_t^{-1}\bar{\tau}_t^{-2}\mathbf{A}\bar{\mathbf{e}}_t(\mathbf{z}_*)\|_2^2]/d
\end{aligned} \tag{66}$$

1091 where the last inequality is by the triangle inequality. By Eq. (62) and the 1-Lipschitzness of $\text{proj}(\cdot)$,
1092 we obtain that the first term in the right-hand side above is upper bounded by

$$\begin{aligned}
&\mathbb{E}[\|\text{proj}_{\sqrt{m}\|\mathbf{A}\|_{\text{op}}\lambda_t^{-1}\bar{\tau}_t^{-2}}(\lambda_t^{-1}\bar{\tau}_t^{-2}\mathbf{A}\tilde{\mathbf{e}}^L(\mathbf{z}_*)) - \text{proj}_{\sqrt{m}\|\mathbf{A}\|_{\text{op}}\lambda_t^{-1}\bar{\tau}_t^{-2}}(\lambda_t^{-1}\bar{\tau}_t^{-2}\mathbf{A}\hat{\mathbf{e}}(\mathbf{z}_*))\|_2^2]/d \\
&\lesssim \frac{m\|\mathbf{A}\|_{\text{op}}^2}{d\lambda_t^2\bar{\tau}_t^4} \cdot (\Pi^2 \cdot (\Pi^2 A)^{2L} + \zeta^2(1 - \Pi^2 A)^{-2}) \lesssim \frac{m\|\mathbf{A}\|_{\text{op}}^2}{d\lambda_t^2\bar{\tau}_t^4} \cdot \left(\Pi^2 \cdot (\Pi^2 A)^{2L} + \frac{m^2\Pi^2}{(1 - \Pi^2 A)^2 M^2} \right).
\end{aligned} \tag{67}$$

1093 In the above display, the last inequality is by the fact that we can choose ζ such that $M = m \cdot$
1094 $(\lceil 2\Pi\zeta^{-1} \rceil + 3)$, which implies that $2m\Pi/M \leq \zeta \leq 6m\Pi/M$. Furthermore, by Assumption 4 and by
1095 the fact that $\|\hat{\mathbf{e}}(\mathbf{z}_*)\|_2 \leq \sqrt{m}\Pi$, we obtain that the second quantity in the right-hand side of Eq. (66)
1096 is controlled by

$$\mathbb{E}[\|\text{proj}_{\sqrt{m}\|\mathbf{A}\|_{\text{op}}\lambda_t^{-1}\bar{\tau}_t^{-2}}(\lambda_t^{-1}\bar{\tau}_t^{-2}\mathbf{A}\hat{\mathbf{e}}(\mathbf{z}_*)) - \lambda_t^{-1}\bar{\tau}_t^{-2}\mathbf{A}\bar{\mathbf{e}}_t(\mathbf{z}_*)\|_2^2]/d \lesssim \frac{m\|\mathbf{A}\|_{\text{op}}^2}{d\lambda_t^2\bar{\tau}_t^4} \cdot \varepsilon_{\text{VI},t}^2(\mathbf{A}). \tag{68}$$

1097 Finally, we combine Eq. (63), (64), (66), (67), (68). This completes the proof of Theorem 4.

1098 F.4 Proof of Lemma 3

1099 Consider the sparse coding problem $\mathbf{z}_* = \mathbf{A}\boldsymbol{\theta} + \bar{\boldsymbol{\varepsilon}} \in \mathbb{R}^d$ with dictionary $\mathbf{A} \in \mathbb{R}^{d \times m}$, sparse
1100 representation $\boldsymbol{\theta} \in \mathbb{R}^m$, and noise $\bar{\boldsymbol{\varepsilon}} \in \mathbb{R}^d$. Assume that the model satisfies the following assumption.

1101 **Assumption 7** (Simplified version of Assumption 1 - 4 of Li et al. [2023b]). Assume that $\mathbf{A} =$
1102 $\mathbf{Q}\mathbf{D}\mathbf{O}^\top$ is the singular value decomposition of \mathbf{A} , where $\mathbf{Q} \in \mathbb{R}^{d \times d}$ and $\mathbf{O} \in \mathbb{R}^{m \times m}$ are orthogonal
1103 and $\mathbf{D} \in \mathbb{R}^{d \times m}$ is diagonal with diagonal elements $\{d_i\}_{i \in [\min\{d, m\}]}$. Assume that \mathbf{Q}, \mathbf{D} are
1104 deterministic, $\mathbf{O}, \boldsymbol{\theta}, \boldsymbol{\varepsilon}$ are mutually independent, and $\mathbf{O} \sim \text{Haar}(\text{SO}(m))$ is uniformly distributed
1105 on the special orthogonal group. As $d, m \rightarrow \infty$, we assume $\mu_{\mathbf{D}} \xrightarrow{W} \mathbf{D}$ where $\mu_{\mathbf{D}}$ is the empirical
1106 distribution of coordinates of \mathbf{D} , \mathbf{D} is a random variable with $\text{supp}\{\mathbf{D}^2\} \subseteq [d_-, d_+]$ and $0 <$
1107 $d_- < d_+ < \infty$, and \xrightarrow{W} denotes Wasserstein- p convergence. Furthermore, $\min_i \{d_i^2\} \rightarrow d_-$ and
1108 $\max_i \{d_i^2\} \rightarrow d_+$. We further assume $\theta_i \sim_{\text{iid}} \pi_0$ with $\mathbb{E}_{\pi_0}[\theta] = 0$, $\mathbb{E}_{\pi_0}[\theta^2] > 0$, and π_0 is compactly
1109 supported. Finally, we have $\bar{\varepsilon}_i \sim_{\text{iid}} \mathcal{N}(0, \bar{\tau}^2)$.

1110 Denote the posterior mean of $\boldsymbol{\theta}$ given $(\mathbf{A}, \mathbf{z}_*)$ by $\mathbf{e}(\mathbf{z}_*) = \mathbb{E}[\boldsymbol{\theta}|\mathbf{z}_*]$. Theorem 1.11 of Li et al. [2023b]
1111 proves the following.

1112 **Lemma 14** (Theorem 1.11 of Li et al. [2023b]). Let Assumption 7 hold. There exists $\bar{\tau}_0^2$ that depends
1113 on $(\alpha, \pi_0, \mathbf{D})$, such that the following happens. For any $\bar{\tau}^2 \geq \bar{\tau}_0^2$, there exists $\nu_* = (\alpha, \pi_0, \mathbf{D}, \bar{\tau}^2)$ that
1114 depends on $(\alpha, \pi_0, \mathbf{D}, \bar{\tau}^2)$ such that, taking $G(\lambda) = \log \mathbb{E}_{\beta \sim \pi_0}[e^{\lambda\beta - \beta^2\nu_*/2}]$ we have almost surely

$$\lim_{d, m \rightarrow \infty} \mathbb{E}_{\mathbf{z}_*} \left[\left\| \mathbf{e}(\mathbf{z}_*) - G'(-\bar{\tau}^{-2}((\mathbf{A}^\top \mathbf{A} - \nu_* \mathbf{I}_m)\mathbf{e}(\mathbf{z}_*) - \mathbf{A}^\top \mathbf{z}_*)) \right\|_2^2 \right] \mathbf{A} = 0.$$

1115 Furthermore, for any fixed $(\pi_0, \alpha, \mathbf{D})$, we have $\sup_{\bar{\tau}^2 \geq \bar{\tau}_0^2} \nu_*(\bar{\tau}^2) < \infty$.

1116 We remark that Theorem 1.11 of Li et al. [2023b] assumes the fixed noise level $\bar{\tau}^2 = 1$. However, a
1117 simple rescaling argument could extend the result to general $\bar{\tau}^2$.

1118 Given Lemma 14, we are now ready to prove Lemma 3. Taking $\bar{\tau}^2 = \bar{\tau}_t^2 = \tau^2 + \sigma_t^2/\lambda_t^2$, $\nu_t = \nu_*(\bar{\tau}_t^2)$,
1119 $G_t = G$, and $\mathbf{K}_t = \bar{\tau}_t^{-2}\nu_*(\bar{\tau}_t^2)$, we note that the minimizer of the VI free energy $\hat{\mathbf{e}}_t(\mathbf{z}_*) \in [-\Pi, \Pi]^m$
1120 should satisfy

$$\hat{\mathbf{e}}_t(\mathbf{z}_*) = G'_t(-\bar{\tau}_t^{-2}((\mathbf{A}^\top \mathbf{A} - \nu_t \mathbf{I}_m)\hat{\mathbf{e}}_t(\mathbf{z}_*) - \mathbf{A}^\top \mathbf{z}_*)).$$

1121 For the posterior mean $\mathbf{e}_t(\mathbf{z}_*) \in [-\Pi, \Pi]^m$, we have

$$\begin{aligned}
& \left\| \mathbf{e}_t(\mathbf{z}_*) - G'_t \left(-\bar{\tau}_t^{-2} ((\mathbf{A}^\top \mathbf{A} - \nu_t \mathbf{I}_m) \mathbf{e}_t(\mathbf{z}_*) - \mathbf{A}^\top \mathbf{z}_*) \right) \right\|_2 \\
& \geq \left\| \mathbf{e}_t(\mathbf{z}_*) - \hat{\mathbf{e}}_t(\mathbf{z}_*) \right\|_2 \\
& \quad - \left\| G'_t \left(-\bar{\tau}_t^{-2} ((\mathbf{A}^\top \mathbf{A} - \nu_t \mathbf{I}_m) \mathbf{e}_t(\mathbf{z}_*) - \mathbf{A}^\top \mathbf{z}_*) \right) - G'_t \left(-\bar{\tau}_t^{-2} ((\mathbf{A}^\top \mathbf{A} - \nu_t \mathbf{I}_m) \hat{\mathbf{e}}_t(\mathbf{z}_*) - \mathbf{A}^\top \mathbf{z}_*) \right) \right\|_2 \\
& \geq \left(1 - \Pi^2 \bar{\tau}_t^{-2} \|\mathbf{A}^\top \mathbf{A} - \nu_t \mathbf{I}_m\|_{\text{op}} \right) \|\mathbf{e}_t(\mathbf{z}_*) - \hat{\mathbf{e}}_t(\mathbf{z}_*)\|_2,
\end{aligned}$$

1122 where the last inequality used the fact that G'_t is Π^2 -Lipschitz. Notice that by Lemma 14,

1123 $\sup_{\bar{\tau}^2 \geq \bar{\tau}_0^2} \nu_*(\bar{\tau}^2) = \nu < \infty$, and $\|\mathbf{A}^\top \mathbf{A}\|_{\text{op}} = \max_i d_i^2$ bounded almost surely by some $D < \infty$ per

1124 Assumption 7. Therefore, when $\tau_0^2 \geq 2\Pi^2(D + \nu)$, we have $1 - \Pi^2 \bar{\tau}_t^{-2} \|\mathbf{A}^\top \mathbf{A} - \nu_* \mathbf{I}_m\|_{\text{op}} \geq 1/2$

1125 for any $\tau^2 \geq \tau_0^2$ and any t . This gives

$$\left\| \mathbf{e}_t(\mathbf{z}_*) - G'_t \left(-\bar{\tau}_t^{-2} ((\mathbf{A}^\top \mathbf{A} - \nu_t \mathbf{I}_m) \mathbf{e}_t(\mathbf{z}_*) - \mathbf{A}^\top \mathbf{z}_*) \right) \right\|_2 \geq \|\mathbf{e}_t(\mathbf{z}_*) - \hat{\mathbf{e}}_t(\mathbf{z}_*)\|_2 / 2.$$

1126 Furthermore, by Lemma 14, the posterior mean $\mathbf{e}_t(\mathbf{z}_*)$ satisfies

$$\lim_{d, m \rightarrow \infty} \mathbb{E}_{\mathbf{z}_*} \left[\left\| \mathbf{e}_t(\mathbf{z}_*) - G'_t \left(-\bar{\tau}_t^{-2} ((\mathbf{A}^\top \mathbf{A} - \nu_t \mathbf{I}_m) \mathbf{e}_t(\mathbf{z}_*) - \mathbf{A}^\top \mathbf{z}_*) \right) \right\|_2^2 \middle| \mathbf{A} \right] = 0.$$

1127 This implies that

$$\lim_{d, m \rightarrow \infty} \mathbb{E}_{\mathbf{z}_*} [\|\mathbf{e}_t(\mathbf{z}_*) - \hat{\mathbf{e}}_t(\mathbf{z}_*)\|_2^2 | \mathbf{A}] = 0,$$

1128 which concludes the proof of Lemma 3.