COMPUTATIONAL BOTTLENECKS FOR DENOISING DIFFUSIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Denoising diffusions sample from a probability distribution μ in \mathbb{R}^d by constructing a stochastic process $(\hat{x}_t : t \geq 0)$ in \mathbb{R}^d such that \hat{x}_0 is easy to sample, but the distribution of \hat{x}_T at large T approximates μ . The drift $m : \mathbb{R}^d \times \mathbb{R} \to \mathbb{R}^d$ of this diffusion process is learned by minimizing a score-matching objective.

Is every probability distribution μ , for which sampling is tractable, also amenable to sampling via diffusions? We address this question by studying its relation to information-computation gaps in statistical estimation. Earlier work in this area constructs broad families of distributions μ for which sampling is easy, but approximating the drift m(y,t) is conjectured to be intractable, and provides rigorous evidence for intractability.

We prove that this implies a failure of sampling via diffusions. First, there exist drifts whose score matching objective is superpolynomially close to the optimum value among polynomial time drifts and yet produce samples with distribution that is very far from the target μ . Second, any polynomial-time drift that is also Lipschitz continuous results in equally incorrect sampling.

We instantiate our results on the toy problem of sampling a sparse low-rank matrix, and further demonstrate empirically the failure of diffusion-based sampling. Our work implies that caution should be used in adopting diffusion sampling when other approaches are available.

1 Introduction

1.1 BACKGROUND

Diffusion sampling (DS) (Song & Ermon, 2019; Ho et al., 2020) has emerged as a central paradigm in generative artificial intelligence (AI). Given a target distribution μ on \mathbb{R}^d , we want to sample $x \sim \mu$. Diffusions achieve this goal by generating trajectories of a stochastic process (\hat{x}_t) whose state \hat{x}_T at large T is approximately distributed according to μ . This suggests a natural question:

Q: Are there distributions μ for which sampling via diffusions fails even if sampling from μ is easy?

In order to explain how DS might fail, it is useful to recall the setup and introduce some notations¹. The basic DS approach implements an approximation of the following stochastic differential equation (SDE), with initialization $y_0 = 0$:

$$dy_t = m(y_t; t)dt + dB_t, \tag{1}$$

$$m(y,t) := \mathbb{E}\{x|tx + \sqrt{t}g = y\},\qquad(2)$$

where $(B_t)_{t\geq 0}$ is Brownian motion (BM) and in Eq. (2) $x \sim \mu$ is independent of $g \sim N(0, I_d)$.

It is not hard to show that, if y_t is generated according to the above SDE, then there exists $x \sim \mu$ and an independent standard BM $(W_t)_{t>0}$ (different from $(B_t)_{t>0}$) such that

$$\mathbf{y}_t = t\,\mathbf{x} + \mathbf{W}_t\,. \tag{3}$$

¹We follow the formulation of Montanari (2023), which does not require time reversal.

Therefore, running the diffusion (1) until some large time T, and returning y_T/T or $m(y_T, T)$ yields a sample approximately distributed according to μ .

In practice, the function m is generally not accessible (cf. discussion below (6)), and is replaced by an approximation $\hat{m}(y,t)$. We can implement an Euler discretization of the SDE (1):

$$\hat{\mathbf{y}}_{t+\Delta} = \hat{\mathbf{y}}_t + \hat{\mathbf{m}}(\hat{\mathbf{y}}_t, t)\Delta + \sqrt{\Delta}\,\hat{\mathbf{z}}_t\,,\tag{4}$$

with Δ a small stepsize, and $(\hat{z}_t)_{t \in \mathbb{N}\Delta} \sim_{iid} \mathbb{N}(\mathbf{0}, \mathbf{I}_d)$. After iterating (4) up to a large time T, we output $\hat{x}_T = \hat{m}(\hat{y}_T, T)$. We refer to \hat{x}_T as a diffusion sample.

Diffusions reduce the problem of sampling from a distribution μ to that of approximating the conditional expectation m (Eq. (2)) by \hat{m} . The mapping $y \mapsto m(y,t)$ is the Bayes-optimal estimator of x in Gaussian noise:

$$\boldsymbol{m}(\cdot,t) = \underset{\boldsymbol{\varphi}:\mathbb{R}^n \to \mathbb{R}^n}{\operatorname{arg\,min}} \mathbb{E}\{\|\boldsymbol{\varphi}(\boldsymbol{y}_t) - \boldsymbol{x}\|^2\}.$$
 (5)

In words, we are given a Gaussian observation $y_t \sim N(tx, tI_d)$ (for a single t) and want to estimate x as to minimize mean square error (MSE). This is also known as the 'score-matching objective'.

The minimization in Eq. (5) has to be modified for two reasons: *First*, in general we do not know the distribution of x over which the expectation in (5) is taken; we only have a sample $(x_i)_{i < N} \sim_{iid} \mu$. We thus replace the MSE by its sample version:

minimize
$$\frac{1}{N} \sum_{i=1}^{N} \| \boldsymbol{\varphi}(\boldsymbol{y}_{i,t}) - \boldsymbol{x}_i \|^2$$
, subj. to $\boldsymbol{\varphi} \in \mathcal{N}$, (6)

where $y_{i,t} = tx_i + \sqrt{t}g_i$ for $(g_i)_{i \leq N} \sim_{iid} N(0, I_d)$. The minimization in (5) must be restricted to a function class \mathscr{N} (e.g. neural nets). A (near)-optimal solution to (6) will be \hat{m} .

Second, to efficiently implement the generative process (4), \hat{m} should be computable in polynomial time. For this reason, \mathcal{N} must be a set of such functions. This is a purely computational constraint, and is present even if we have access to μ (i.e., for $N=\infty$).

Most of the literature on diffusion sampling studies how samples quality deteriorates because of finite sample size N or non-vanishing step size Δ . Here we focus on a more fundamental limitation that arises because \hat{m} must be computable in polynomial time (the second remark above).

A key remark here is that the ideal drift m(y,t) is the Bayes-optimal denoiser, see (5). Namely it is the optimal function to estimate x with prior distribution μ from noisy observations $y_t \sim N(tx,tI_d)$: t can be interpreted as the signal-to-noise ratio (SNR) of this denoising problem. We will say that an *information-computation gap* arises for this problem (at SNR t) there exists a constant gap(t) > 0 such that, for all polynomial-time algorithms \hat{m} , if d is large enough

$$\mathbb{E}\left\{\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t) - \boldsymbol{x}\|^2\right\} \ge \inf_{\boldsymbol{\varphi}} \mathbb{E}\left\{\|\boldsymbol{\varphi}(\boldsymbol{y}_t) - \boldsymbol{x}\|^2\right\} + \mathsf{gap}(t). \tag{7}$$

Recent literature provides many instances of statistical estimation problems for which an information-computation gap is shown to exist (Brennan et al., 2018; Bandeira et al., 2022; Celentano & Montanari, 2022; Schramm & Wein, 2022) conditional on certain widely accepted conjectures. We stress that the conditional/conjectural nature of these results is, so far, unavoidable, a situation analogous to classical complexity theory that relies on $P \neq NP$. Several of the problems for which a gap arises take the form of estimating $x \sim \mu$ from observations $y_t = tx + \sqrt{t} g$.

Koehler & Vuong (2024) already pointed out informally that denoising problems presenting an information-computation gap can result into a failure of DS. As a concrete example, they suggested the spiked Wigner model (c.f. next section). While this informal remark is natural, making it mathematically precise is far from obvious. In fact –strictly speaking– the remark is **false**. If sampling from μ is easy, then the drift m(y,t) can be constructed to return (for all $t \geq t_0$) a fixed random sample $x \sim \mu$. Then the diffusion will sample correctly. However such m will be very far from an optimal denoiser. (See Proposition 2.1 for formal version of this counter-example.)

We also note that several earlier papers provided examples of probability distributions μ from physics and Bayesian statistics for which Gibbs sampling is expected to succeed, but DS appears

to fail (Montanari et al., 2007; Ricci-Tersenghi & Semerjian, 2009; Ghio et al., 2024; Huang et al., 2024). None of these papers presented a formal claim either.

The present paper fills this gap in the literature. We prove two general results that hold for any distribution μ that presents a certain version of information-computation gap (see formal statements below). *First*, we prove that there exists drifts that are approximate optimizers of the score matching objective (6) among polynomial time algorithms (up to an sub-polynomially small error) and yet lead to completely incorrect sampling. *Second*, we show that *every* polynomial-time computable drift that is a near optimum of score matching and is also Lipschitz continuous leads to incorrect sampling. *Finally*, we ilustrate the applicability of our theorems by studying a toy example, namely sampling a sparse low-rank matrix.

We emphasize that this failure of DS is of computational of nature and purely related to the requirement to approximate the Bayes optimal denoiser m(y, t) by a polytime computable function.

1.2 Summary of results

Recall that the Wasserstein-1 distance between two measures μ_1, μ_2 on \mathbb{R}^d is defined as

$$W_1(\mu_1, \mu_2) := \inf_{\gamma \in \mathcal{C}(\mu_1, \mu_2)} \int \|\boldsymbol{x}_1 - \boldsymbol{x}_2\|_2 \, \gamma(\mathrm{d}\boldsymbol{x}_1, \mathrm{d}\boldsymbol{x}_2) \,,$$

with the infimum taken over couplings on μ_1 and μ_2 . Given random vectors $\boldsymbol{X}_1, \boldsymbol{X}_2$ we denote by $W_1(\boldsymbol{X}_1, \boldsymbol{X}_2)$ the W_1 -distance of their distributions. We prove lower bounds on the W_1 to show incorrect sampling. Since we only consider distributions μ supported on vectors with bounded norm, a lower bound on W_1 implies lower bounds on TV distance and KL divergence. Hence our impossibility results are stated in a strong form.

As a running example/application, we will let μ to be the following distribution over $n \times n$ sparse low-rank matrices. Let $B_{n,k} := \{ \boldsymbol{u} \in \{0, \pm 1/\sqrt{k}\}^n | \|\boldsymbol{u}\|_0 = k \}$ be the set of $0/\pm (1/\sqrt{k})$ unit vectors with k nonzero entries ($\|\boldsymbol{u}\|_0$ denotes the number of nonzeros in \boldsymbol{u}). We define the target distribution $\mu = \mu_{n,k}$ to be the distribution of $\boldsymbol{x} = \boldsymbol{u}\boldsymbol{u}^\mathsf{T}$ when $\boldsymbol{u} \sim \mathrm{Unif}(B_{n,k})$. Note that $\boldsymbol{x} \in \mathbb{R}^{n \times n}$ is a matrix that we identify with a vector in \mathbb{R}^d for $d = n^2$. Sampling from μ is trivial: just sample a vector with entries in $\{0, 1/\sqrt{k}, -1/\sqrt{k}\}$ and exactly k non-zero entries, and let $\boldsymbol{x} = \boldsymbol{u}\boldsymbol{u}^\mathsf{T}$. However, rigorous evidence supports the claim that —for certain scalings of k, t with n—polynomial-time algorithms cannot approach the Bayes-optimal error (Butucea et al., 2015; Ma & Wu, 2015; Cai et al., 2017; Brennan et al., 2018; Schramm & Wein, 2022).

We will prove two sets of main results that hold for distributions μ such that the denoising problem presents an information-computation gap:

- **1.** (Theorem 1, Corollaries 3.2, C.1) Near optimizers of score-matching can sample incorrectly. We prove that there exists $\hat{m}: \mathbb{R}^{n \times n} \times \mathbb{R} \to \mathbb{R}^{n \times n}$ such that:
- M1. $\hat{\boldsymbol{m}}(\cdot)$ can be evaluated in polynomial time.
- M2. The estimation error achieved by \hat{m} (namely, $\mathbb{E}\{\|\hat{m}(y_t,t)-x\|^2\}$) is close to the optimal estimation error achieved by polynomial-time algorithms. Hence $\hat{m}(\cdot,t)$ will be a near minimizer of the score-matching objective (5) (integrated over t).
- M3. Samples \hat{x}_T generated by the discretized diffusions (4) with drift $\hat{m}(\cdot, t)$ at some large time T have distribution that is very far from the target μ ('as far as it can be' in W_1 distance.)
- **2.** (Theorem 3, Corollary 5.1) All (sufficiently) Lipschitz score-matching optimizers sample incorrectly. More precisely, we prove that any denoiser that near optimizes the score matching among polytime algorithms, acts optimally on pure noise data, and is C/t-Lipschitz for $t > t_1$ (for any constant C a suitable t_1), samples incorrectly.

Additionally, (Theorems 2, 5), we prove a reduction from estimation to DS. Namely, if accurate, polytime DS is possible, then near Bayes optimal estimation of x from $y_t = tx + \sqrt{t}g$ must also be possible in polynomial time for all t. The contrapositive of this statement implies that if an information-computation gap exists, then (near)-correct DS is impossible in polynomial time.

Roadmap. The rest of the paper is as follows. In Section 2 we motivate our setting and assumptions, and discuss some limitations of our results. In Section 3 we state formally our results (for technical

reasons we state two separate results depending on the growth of k with n.) Section 4 presents the general reduction from estimation to diffusion sampling. Section 5 proves that all Lipschitz score matching optimizers fail. Section 6 provides a numerical experiment of a neural network \hat{m} that outperforms (conjectured) asymptotically optimal denoisers for finite n, yet still samples poorly.

Notation. Throughout, $a_n \ll b_n$ means $a_n/b_n \to 0$. We refer to Appendix A for notations.

2 DISCUSSION

Setting. Our results indicate that a standard application of denoising diffusions methodology will fail to sample from μ when the associated denoising problem presents an information-computation gap. The example $\mu_{n,k}$ of sparse low-rank matrices shows that DS can fail in cases in which sampling from μ is trivial.

Our example also shows that the latent structure of the distribution can be exploited to construct a better algorithm. Namely, one can use diffusions to sample $u \sim \mathrm{Unif}(B_{n,k})$ (the posterior expectation m(y,t) is polytime-computable) and then generate $x = uu^{\mathsf{T}}$. On the other hand, identifying such latent structures from data can be hard in general, both statistically and computationally.

Limitations. We prove that there exists drifts $\hat{m}(\cdot,t)$ that lead to poor sampling, despite being nearly optimal (among poly-time algorithms) in terms of the score matching objective (5). In particular, these bad drifts will be near optimal solutions of the problem of (6), as long as \mathscr{N} only contains polytime methods and is rich enough to approximate them. We further exclude the existence of Lipschitz drifts $\hat{m}(\cdot,t)$ that also satisfy conditions M1 and M2 but yield good generative sampling.

In principle there could still be non-Lipschitz polytime drifts that are near score matching optimizers and sample well. However if such drifts exist, our results suggest that minimizing the score matching objective is not the right approach to find them (since the difference in value with bad drifts will be superpolynomially small).

Correct samplers violating M2. If we drop condition M2, i.e. we accept drifts that are bad for the score-matching objective, then it is possible to construct drifts that can be evaluated in polynomial time and yield good sampling. This is stated formally below and proven in Appendix I.

Proposition 2.1. Suppose that a discretized SDE $(\hat{y}_{\ell\Delta})_{\ell\geq 0}$ per (4) is generated, with step size $\Delta>0$ and noise stream $\hat{z}_t\stackrel{i.i.d.}{\sim} \mathsf{N}(0,\mathbf{I}_{n\times n})$. Then for every n,k, there exists a function $\hat{\boldsymbol{m}}(\boldsymbol{y},t)=\hat{\boldsymbol{m}}(\boldsymbol{y},t;\hat{\boldsymbol{z}}_1)$ parametrized by $\hat{\boldsymbol{z}}_1$ (with no additional randomness) such that: $(i) \mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t)-\boldsymbol{x}\|^2]=2(1-o(1))$ uniformly for every $t\geq 0$ (sub-optimal score-matching); $(ii) \ W_1(\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_{\ell\Delta},\ell\Delta),\boldsymbol{x})=0$ for all $\ell\geq 0$ $(\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_t,t)$ is an approximate sample of \boldsymbol{x}); $(iii) \lim_{\ell\to\infty} W_1(\hat{\boldsymbol{y}}_{\ell\Delta}/(\ell\Delta),\boldsymbol{x})=0$ $(\hat{\boldsymbol{y}}_t/t$ is an approximate sample of \boldsymbol{x} at large time).

The drift constructed in this proposition has very poor value of the score-matching objective.

Further related work. A number of groups proved positive results on diffusion sampling. Alaoui et al. (2022); Chen et al. (2023b); Montanari & Wu (2023); Lee et al. (2023); Benton et al. (2023) provide reductions from diffusion sampling to score estimation. Chen et al. (2023a); Shah et al. (2023); Mei & Wu (2025); Li et al. (2024) give end-to-end guarantees for classes of distributions μ .

The computational bottleneck that we study here has been observed before in the context of certain Gibbs measures and Bayes posterior distributions Ghio et al. (2024); Alaoui et al. (2023); Huang et al. (2024), and random constraint satisfaction problems Montanari et al. (2007); Ricci-Tersenghi & Semerjian (2009) (the later papers use sequential sampling rather than diffusion sampling).

Our work provides an approach to rigorize the latter line of work.

3 NEAR-OPTIMAL POLYTIME DRIFTS WITH INCORRECT DIFFUSION SAMPLING

Given an arbitrary polytime computable drift \hat{m}_0 , we will construct a different polytime drift \hat{m} , with nearly equal score matching objective and yet incorrect sampling. In Subsection 3.1, we state

our assumptions and general result. In Subsection 3.2, we apply the general theorem to the example of sampling sparse low-rank matrices. We also indicate several other similar examples.

In what follows (x, y_t) will always be distributed according to the ideal diffusion process of (3), which also satisfies (2). In particular $x \sim \mu$, $y_t = tx + W_t$, for $(W_t)_{t\geq 0}$ a BM. On the other hand, (\hat{y}_t) will denote the process generated with the implemented procedure (4).

3.1 General result

Throughout, we will consider distributions μ that are supported on $\mathsf{B}^d(1) := \{ \boldsymbol{x} \in \mathbb{R}^d : \|\boldsymbol{x}\|_2 \le 1 \}$. We will state our assumptions and results having in mind the case of measures that are roughly centered: $\mathbb{E}_{\boldsymbol{x} \sim \mu}[\boldsymbol{x}] = \int \boldsymbol{x} \, \mu(\mathrm{d}\boldsymbol{x}) \approx \mathbf{0}$, although this condition is not formally needed.

Our first main assumption is that any polynomial-time algorithm to estimate x from $y_t \sim N(tx,tI_d)$ fails when t is below a certain threshold $t_{\rm alg}$. When $t/t_{\rm alg} < 1$, we expect that polytime algorithms will not perform better (in score-matching, c.f. (5)) than the best constant estimator of x, namely $\mathbb{E}_{x\sim\mu}[x]$. In the case $\|\mathbb{E}_{x\sim\mu}[x]\|\approx 0$, it follows that polytime algorithms \hat{m}_0 with good score-matching will have small norm $\|\hat{m}_0(y_t,t)\|$. This small-norm property is captured by our assumption. More details are discussed at the beginning of Subsection 3.2.1, and Proposition B.1.

Assumption 1 (Small norm below threshold). Let $y_t = tx + W_t$, for $(x, (W_t)_{t \ge 0}) \sim \mu \otimes BM$. Then, there exists a function $\eta_1 : \mathbb{N} \to \mathbb{R}$ (which we refer to as 'rate') such that $\eta_1(d) = o_d(1)$ and, for any $\varepsilon, \gamma > 0$,

$$\int_0^{(1-\gamma)t_{\text{alg}}} \mathbb{P}(\|\hat{\boldsymbol{m}}_0(\boldsymbol{y}_t,t)\| \geq \varepsilon) dt = O(\eta_1(d)).$$

Our second assumption is that polytime *detection* is reliable for t above t_{alg} . By detection, we consider the following hypothesis testing problem. Given $y \in \mathbb{R}^d$, we test if y is distributed as $tx + \sqrt{t} N(\mathbf{0}, I_d)$ or as $N(a, tI_d)$ for ||a|| small, where a might depend on the Gaussian noise.

Assumption 2 (Hypothesis testing succeeds above threshold). For $c \in (0,1)$, define $\mathcal{A}_d(c) = \{ \boldsymbol{a} \in \mathbb{R}^d : \|\boldsymbol{a}\| \leq c \, t_{\text{alg}} \}$. We assume there exists $\delta, \eta_2 : \mathbb{N} \to \mathbb{R}$ (which we refer to as rates), and a polytime binary test function $\phi : \mathbb{R}^d \times \mathbb{R}_{>0} \to \{0,1\}$ such that:

- 1. (Sharp detection threshold) $\delta(d) = o_d(1)$.
- 2. For the process $(y_t = tx + W_t)$, ϕ rejects with high probability:

$$\int_{t_{\text{alo}}(1+\delta)}^{\infty} \mathbb{P}(\phi(\boldsymbol{y}_t, t) = 0) dt = O(\eta_2(d)).$$

3. Uniformly over the set $A_d(c)$, ϕ fails to reject with high probability. Namely:

$$\mathbb{P}\big(\exists t \geq t_{\mathrm{alg}}(1+\delta) \text{ such that } \sup_{\boldsymbol{a} \in \mathcal{A}_{d}(c)} \phi(\boldsymbol{a} + \boldsymbol{W}_{t}) = 1\big) = o(1) \,.$$

Remark. Since we try to state our theorem in the strongest form, Assumptions 1 and 2 do not take the same form as the information-computation gap (7). Nevertheless, it can be proven that (for a broad class of problems) these assumptions cannot hold unless an information-computation gap is present. We leave this point for future work.

We state our first main result. It stipulates that we can construct a polytime algorithm which has 'essentially' the same score-matching objective as \hat{m}_0 yet yields bad samples.

Theorem 1. Let μ be a probability measure supported on $\mathsf{B}^d(1)$ such that $\liminf_{d\to\infty}\int\|\boldsymbol{x}\|\,\mu(\mathrm{d}\boldsymbol{x})=\alpha>0$. Assume that there exist $t_{\mathrm{alg}}=t_{\mathrm{alg}}(d)>0$, a drift $\hat{\boldsymbol{m}}_0:\mathbb{R}^d\times\mathbb{R}\to\mathbb{R}$, and functions $\eta_1(d),\delta(d),\varepsilon_2(d)=o_d(1)$, such that following conditions hold: $(i)\sup_{\boldsymbol{y},t}\|\hat{\boldsymbol{m}}_0(\boldsymbol{y},t)\|\leq 1$. (ii) Assumption 1 holds with rate $\eta_1(d)$. (iii) Assumption 2 holds with rates $\delta(d),\eta_2(d)$.

Then there exists a modified drift \hat{m} such that

M1. $\hat{\boldsymbol{m}}(\cdot)$ can be evaluated in polynomial time.

M2. If $y_t = tx + B_t$ is the true diffusion (equivalently given by (1)), then

$$\int_0^\infty \mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) - \hat{\boldsymbol{m}}_0(\boldsymbol{y}_t,t)\|^2] dt = O(\eta_1(d) \vee \eta_2(d)).$$

M3. For any step size $\Delta = \Delta_n > 0$, we have incorrect sampling:

$$\inf_{t \in \mathbb{N} \cdot \Delta, t \ge (1+\delta)t_{\text{alg}}} W_1(\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_t, t), \boldsymbol{x}) \ge \alpha - o_d(1).$$
(8)

The proof is presented in Appendix F. The main idea is to let $\hat{\boldsymbol{m}}(\boldsymbol{y},t)$ be $\hat{\boldsymbol{m}}_0(\boldsymbol{y},t)\mathbf{1}_{\parallel\hat{\boldsymbol{m}}_0(\boldsymbol{y},t)\parallel\leq\varepsilon}$ for $t\leq (1-\gamma)t_{\rm alg}$, and $\hat{\boldsymbol{m}}_0(\boldsymbol{y},t)\phi(\boldsymbol{y},t)$ for $t\geq (1+\delta)t_{\rm alg}$, with small constants (γ,ε) and test ϕ .

Remark. It makes sense to assume that $\|\hat{m}_0(\cdot,\cdot)\| \le 1$. Since $\operatorname{supp}(\mu) \subseteq \mathsf{B}^d(1)$ and the latter is a convex set, projecting any \hat{m}_0 onto this set yields a smaller MSE.

3.2 Example: Sampling low-rank matrices

We state two separate results for the probability distribution $\mu=\mu_{n,k}$ described in the introduction, depending on the scaling of k with n: in Section 3.2.1 we assume $\sqrt{n} \ll k \ll n$; while in Appendix C we assume $k \ll \sqrt{n}$. Indeed, the nature of the problem changes at the threshold $k \asymp \sqrt{n}$.

A crucial role will be played by the following threshold

$$t_{\text{alg}}(n,k) := \begin{cases} k^2 \log\left(\frac{n}{k^2}\right) & \text{if } k \ll \sqrt{n} \\ \frac{n}{2} & \text{if } \sqrt{n} \ll k \ll n \end{cases}$$
 (9)

It is expected that for $t \leq (1 - \delta)t_{\text{alg}}(n, k)$ and δ any fixed constant, no polytime algorithm can estimate \boldsymbol{x} significantly better than the estimator $\hat{\boldsymbol{m}}_{\text{null}} = \mathbb{E}[\boldsymbol{x}] \approx \boldsymbol{0}$ for $k \ll n$ (see Conjecture 3.1).

Since $\|x\|_F = 1$ for $x \sim \mu$, the Bayes denoiser m(y,t) = m(y) does not depend on t (this can be seen by Bayes rule). From now on, we refer to $\|x\| = \|x\|_F$ as the Frobenius norm.

3.2.1 Moderately sparse regime: $\sqrt{n} \ll k \ll n$

Assumption 1 states that, for $y_t = tx + W_t$, the estimated drift $\hat{m}_0(y_t, t)$ should have small norm with high probability. This condition holds under the well-accepted Conjecture 3.1 below on information-computation gaps. In fact, a simple consequence of this conjecture is that any polytime \hat{m} matching this error must satisfy $\mathbb{E}\{\|\hat{m}(y_t,t)\|^2\} = o_n(1)$ (see Proposition B.1).

Conjecture 3.1. For $\sqrt{n} \ll k \ll n$, there exists $\underline{k}_n \ll n$ such that the following holds for any $k = k_n$, with $\underline{k}_n \leq k_n \ll n$. Let $\{\hat{m}_n\}_{n\geq 1}$, $\hat{m}_n : \mathbb{R}^{n\times n} \times \mathbb{R} \to \mathbb{R}^{n\times n}$ be any sequence of polytime algorithms (polynomial time in n). Then for any $\delta > 0$, we have

$$\inf_{t \le (1-\delta)t_{\text{alg}}} \mathbb{E}\{\|\hat{\boldsymbol{m}}_n(\boldsymbol{y}_t, t) - \boldsymbol{x}\|^2\} \ge 1 - o_n(1).$$
(10)

We refer to Ma & Wu (2015); Cai et al. (2017); Hopkins et al. (2017); Brennan et al. (2018); Schramm & Wein (2022); Kunisky et al. (2019) for evidence towards this conjecture. Next, we provide the following implication of Theorem 1, whose proof is in Appendix G.

Corollary 3.2. Assume $\sqrt{n} \ll k \ll n$, so that $t_{\text{alg}}(n,k) := n/2$ per (9). Let $\hat{\boldsymbol{m}}_0$ be an arbitrary poly-time algorithm such that $\sup_{\boldsymbol{y},t} \|\hat{\boldsymbol{m}}_0(\boldsymbol{y},t)\|_F \leq 1$ and Assumption 1 holds with rate η_1 such that $\eta_1 \ll n^{-D} \forall D > 0$. Then there exists an estimator $\hat{\boldsymbol{m}}$ such that:

M1. $\hat{\boldsymbol{m}}(\cdot)$ can be evaluated in polynomial time.

M2. If $y_t = tx + B_t$ is the true diffusion (equivalently given by (1)), then, for every D > 0,

$$\int_0^\infty \mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) - \hat{\boldsymbol{m}}_0(\boldsymbol{y}_t,t)\|^2] dt = O(n^{-D}).$$

M3. There exists $\delta = o_n(1)$ such that, for any step size $\Delta = \Delta_n > 0$, we have incorrect sampling:

$$\inf_{t \in \mathbb{N} \cdot \Delta, t \ge (1+\delta)t_{\text{alg}}} W_1(\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_t, t), \boldsymbol{x}) \ge 1 - o_n(1). \tag{11}$$

To connect the last corollary with the introduction, we recall two facts from the literature on submatrix estimation: (i) The Bayes estimator $m(y_t)$ achieves small MSE in a large interval above $t_{\rm alg}$ (Proposition 3.3); (ii) No polytime estimator is expected to perform better than the null estimator below $t_{\rm alg}$ (Conjecture 3.1). Regarding (i), we state a characterization of the Bayes optimal error. The proof is analogous to the main result in Butucea et al. (2015), which considers the case of asymmetric matrices. (For $k \le n^a$, a < 5/6, see also Barbier et al. (2020).)

Proposition 3.3 (Modification of Butucea et al. (2015)). Let m(y) be the posterior mean estimator in Eq. (2). Assume $1 \ll k \ll n$, and define $t_{\text{Bayes}}(n,k) := 2k \log(n/k)$. Then, for any $\delta > 0$, we have $\inf_{t \leq (1-\delta)t_{\text{Baves}}} \mathbb{E}\{\|m(y_t) - x\|^2\} = 1 - o_n(1)$, $\sup_{t \geq (1+\delta)t_{\text{Baves}}} \mathbb{E}\{\|m(y_t) - x\|^2\} = o_n(1)$.

In other words, for $2k \log(n/k) \ll t \ll n$, the optimal estimator can estimate the signal x accurately, but we expect that no polytime algorithm can achieve the same.

3.2.2 Very sparse regime: $k \ll \sqrt{n}$, and other examples

In the very sparse regime $k \ll \sqrt{n}$, we prove a result similar to Corollary 3.2 (Corollary C.1).

Other examples. We mention a few examples where it is relatively straightforward to apply Theorem 1, following the blueprint in Corollary 3.2. (i) Sampling low rank tensors, e.g. $x = u^{\otimes q} \in (\mathbb{R}^n)^{\otimes q}, q \geq 3$ when $u \sim \mathrm{Unif}(\{+1/\sqrt{n}, -1/\sqrt{n}\}^n)$ or u is uniform on the unit sphere; the corresponding denoising problem is known as tensor PCA (Montanari & Richard, 2014) (in this case $d = n^q$). (ii) Sampling elements of random linear subspaces of $\{0,1\}^d$: $x = Gu \mod 2$, where $G \in \{0,1\}^{d \times \ell}$ is a fixed (known) uniformly random matrix and $u \sim \mathrm{Unif}(\{0,1\}^\ell), \ell = rn$ for $r \in (0,1)$ a constant; the corresponding denoising problem amounts to decoding random linear codes (Richardson & Urbanke, 2008; Ghazi & Lee, 2017) (this example fits our framework after centering). We give two classes of examples for which applying Theorem 1 requires additional technical work (defer to future publications): (iii) Sampling from Bayesian posteriors, e.g. posterior of a low-rank plus noise estimation problem that presents an information-computation gap (Lelarge & Miolane, 2017; Montanari & Wu, 2023; Ghio et al., 2024); (iv) Sampling solutions of random constraint satisfaction problems (Montanari et al., 2007; Ghio et al., 2024).

4 REDUCTION OF ESTIMATION TO DIFFUSION-BASED SAMPLING

To complement previous results, we prove a general reduction: if diffusion sampling can be performed in polynomial time with sufficient accuracy, then we can perform also denoising. The contrapositive of this statement aligns with results in previous sections.

To avoid unessential complications, in this section we assume μ to be supported on the unit sphere $\mathbb{S}^{d-1}=\{\boldsymbol{x}:\|\boldsymbol{x}\|=1\}$. We denote by $P_{\boldsymbol{y}}^T$ the law of $(\boldsymbol{y}_t)_{0\leq t\leq T}$ where \boldsymbol{y}_t given by Eq. (1) and by $P_{\hat{\boldsymbol{y}}}^{T,\Delta}$ the law of $(\hat{\boldsymbol{y}}_t)_{0\leq t\leq T}$, which is the discretized diffusion trajectory defined in (4) (interpolated linearly outside $\mathbb{N}\cdot\Delta$).

It is further useful to define $\overline{P}_{\hat{y}}^{T,\Delta}$ to be the law of the SDE interpolating that of (4):

$$d\hat{\mathbf{y}}_t = \hat{\mathbf{m}}(\hat{\mathbf{y}}_{|t|_{\Delta}}, \lfloor t \rfloor_{\Delta}) dt + d\mathbf{B}_t, \qquad (12)$$

where $\lfloor t \rfloor_{\Delta} := \max\{s \in \mathbb{N} \cdot \Delta : s \leq t\}.$

Theorem 2. Assume that $\hat{\boldsymbol{m}}(\cdot,\cdot)$ has complexity χ and that for any $T \leq \theta d$, $D_{\mathrm{KL}}(\overline{P}_{\hat{\boldsymbol{y}}}^{T,\Delta} \| P_{\boldsymbol{y}}^T) \leq \varepsilon$

Then for any $\sigma>0$ there exists an algorithm a randomized algorithm \hat{m}_+ with complexity $(N\chi\cdot T/\Delta)$ that approximates the posterior expectation:

$$\mathbb{E}\{\|\hat{\boldsymbol{m}}_{+}(\boldsymbol{y}) - \boldsymbol{m}(\boldsymbol{y})\|^{2}\} \leq 2\overline{\varepsilon} + 2N^{-1}.$$
(13)

Here $\bar{\varepsilon} := \sqrt{2\varepsilon} + \varepsilon_0(\theta)$ and $\varepsilon_0(\theta) := \mathbb{E}\|P_{\boldsymbol{x}|\boldsymbol{y}} - \mathsf{N}(\boldsymbol{0}, (\theta d)^{-1}\boldsymbol{I}_d) * P_{\boldsymbol{x}|\boldsymbol{y}}\|_{\mathsf{TV}}$ is the expected TV distance between $P_{\boldsymbol{x}|\boldsymbol{y}}$ and the convolution of $P_{\boldsymbol{x}|\boldsymbol{y}}$.

The proof of this result is presented in Appendix S, along with a modification.

5 ALL LIPSCHITZ POLYTIME ALGORITHMS FAIL

In Section 3 (Theorem 1 and Corollary 3.2) we proved that there exist near-optimizers of the score matching objective that perform poorly. However, we did not rule out the possibility that the optimal (in the sense of score-matching) polytime drift \hat{m} will perform well. We next show that this is not the case, under an additional assumption, namely that the drift $\hat{m}(\cdot;t)$ is Lipschitz continuous for $t \geq (1+\delta)t_{\text{alg}}$. Proof is given in Appendix V. (We assume the Lipschitz constant to be C/t, because the input of the denoiser is $y_t = tx + W_t$, and hence the two t-dependent factors cancel.)

Theorem 3. Let μ be supported on $\mathsf{B}^d(1) = \{x : \|x\| \le 1\}$, $\int x \, \mu(\mathrm{d}x) = \mathbf{0}$, and $\liminf_{d\to\infty} \int \|x\| \mu(\mathrm{d}x) = \alpha > 0$. Let $\hat{m} : \mathbb{R}^d \times \mathbb{R}_{\ge 0} \to \mathbb{R}^d$ be a polytime denoiser such that $\sup_{u,t} \|\hat{m}(y,t)\| \le 1$ (below W_t is a standard BM):

1. $\hat{\boldsymbol{m}}$ is nearly optimal, namely for $\boldsymbol{y}_t = t\boldsymbol{x} + \boldsymbol{W}_t$, and every $\gamma > 0$

$$\sup_{t \le (1-\gamma)t_{\text{alg}}} \left| \mathbb{E} \left\{ \| \hat{\boldsymbol{m}}(\boldsymbol{y}_t, t) - \boldsymbol{x} \|^2 \right\} - \mathbb{E}[\| \boldsymbol{x} \|^2] \right| = o(t_{\text{alg}}^{-1}),$$
(14)

$$\sup_{t \ge (1+\gamma)t_{\text{alg}}} \mathbb{E}\{\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t, t) - \boldsymbol{x}\|^2\} = o(1),$$
(15)

and that for every $t \geq 0, c \in [0, 1]$, $\mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t, t) - \boldsymbol{x}\|^2] \leq \mathbb{E}[\|c\hat{\boldsymbol{m}}(\boldsymbol{y}_t, t) - \boldsymbol{x}\|^2]$.

- 2. ($\hat{\boldsymbol{m}}$ is small on pure noise.) For some $\delta=o(1)$, $\int_{(1+\delta)t_{\mathrm{alg}}}^{\infty}\mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{W}_t,t)\|^2]dt=o(1)$
- 3. $\hat{\boldsymbol{m}}(\cdot,t)$ is C/t-Lipschitz for some constant C and all $t \geq (1+\delta)t_{\text{alg}}$.

Then, for every constant $C_0 > 0$ and step size $\Delta > 0$:

$$\inf_{t \in \mathbb{N} \Delta \cap [0, C_0 t_{\text{alg}}]} W_1(\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_t, t), \boldsymbol{x}) \ge \alpha - o(1)$$

We apply the above theorem to our running example of sampling sparse low-rank matrices. In order to make sure that condition 2 in the theorem is verified, we introduce a variant $\overline{\mu}_{n,k}$ of $\mu_{n,k}$ (all conclusions stated for μ , e.g., Theorem 1, Corollaries 3.2, C.1 hold for $\overline{\mu}_{n,k}$ as well.) Letting $\mu_{n,k}^0$ be the centered version of $\mu_{n,k}$; we define $\overline{\mu}_{n,k} = \frac{1}{2}\,\delta_0 + \frac{1}{2}\,\mu_{n,k}^0$. In words, with probability 1/2 we let x=0 and with probability 1/2 we draw $x=\tilde{x}-\mathbb{E}[\tilde{x}], \tilde{x}\sim \mu_{n,k}$, a sparse rank-one matrix, as in previous sections. As mentioned, this mixture distribution $\overline{\mu}_{n,k}$ is mainly to satisfy condition 2 of Theorem 3. Indeed, we have the following decomposition

$$\mathbb{E}_{\boldsymbol{x} \sim \overline{\mu}_{n,k}}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) - \boldsymbol{x}\|^2] = \frac{1}{2}\mathbb{E}_{\boldsymbol{x} \sim \mu_{n,k}^0}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) - \boldsymbol{x}\|^2] + \frac{1}{2}\mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{W}_t,t)\|^2],$$

which shows that, to get $\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) \approx \boldsymbol{x}$ under the mixture distribution, we also need $\hat{\boldsymbol{m}}(\boldsymbol{W}_t,t) \approx \boldsymbol{0}$. More concretely, we can get explicit rates on $\mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{W}_t,t)\|^2]$ for t above $t_{\rm alg}$ by enforcing that $\hat{\boldsymbol{m}}$ cannot be improved by multiplying by certain hypothesis tests. The full result is as follows.

Corollary 5.1. Assume \underline{k}_n exists as in Conjecture 3.1. Let $k = k_n$ be such that $\underline{k}_n \vee \sqrt{n} \leq k_n \ll n$ (moderately sparse regime). Let \hat{m}_n be a polytime denoiser such that for every fixed constant $\gamma > 0$:

1. $\hat{\boldsymbol{m}}_n$ is nearly optimal, namely (for $\boldsymbol{y}_t = t\boldsymbol{x} + \boldsymbol{W}_t$, \boldsymbol{W}_t standard BM)

$$\sup_{t \le (1-\gamma)t_{\text{alg}}} \left| \mathbb{E} \left\{ \| \hat{\boldsymbol{m}}(\boldsymbol{y}_t, t) - \boldsymbol{x} \|^2 \right\} - \mathbb{E}[\| \boldsymbol{x} \|^2] \right| = o_n(n^{-1}),$$
 (16)

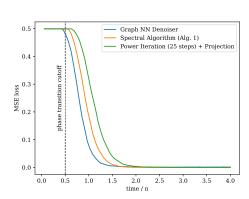
$$\sup_{t \ge (1+\gamma)t_{\text{alg}}} \mathbb{E}\{\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t, t) - \boldsymbol{x}\|^2\} = o_n(1),$$
(17)

and further, for any $t \geq 0$ the MSE of $\hat{\mathbf{m}}_n$ smaller or equal than the MSE of $c(\lambda_1(\mathbf{y}_t))\hat{\mathbf{m}}_n(\mathbf{y}_t,t)$ for any polytime function $c(\cdot)$ of the maximum eigenvalue of $(\mathbf{y}_t + \mathbf{y}_t^{\mathsf{T}})/\sqrt{2}$, and than the MSE of $\mathbf{P}_{\mathsf{B}}\hat{\mathbf{m}}_n$, for \mathbf{P}_{B} the projection onto the unit ball.

2. $\hat{m}_n(\cdot,t): \mathbb{R}^{n\times n} \to \mathbb{R}^{n\times n}$ is C/t-Lipschitz for some constant C and all $t \geq (1+\delta)t_{\text{alg.}}$

Then, for every constant $C_0 > 0$, and step size $\Delta > 0$

$$\inf_{t \in \mathbb{N} \cdot \Delta \cap [0, C_0 n]} W_1(\hat{\boldsymbol{n}}(\hat{\boldsymbol{y}}_t, t), \boldsymbol{x}) \ge \frac{1}{2} - o_n(1).$$
 (18)



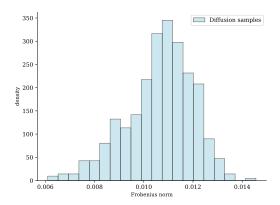


Figure 1: Generating sparse rank-one matrices $x \sim \tilde{\mu}_{n,k}$ using denoising diffusions, for n=350, k=20. Left: MSE of various denoisers (vertical line corresponds to the algorithmic threshold $t_{\rm alg}$.) Right: Frobenius norms of generated samples.

The proof is given in Appendix W. We note that the error of polytime denoisers in (16) (and sampling error of Eq. 18) is 1/2 instead of 1 because the best constant denoiser achieves error 1/2.

Corollary 5.1 does not rule out the possibility that there exists a near-optimizer of score matching that violates the Lipschitz condition and samples well. However, for $t \geq (1+\delta)t_{\rm alg}$ accurate estimation is possible with Lipschitz algorithms, and indeed many natural methods are in this class (e.g. neural nets with bounded number of layers and suitable operator norm bounds on the weights.)

6 NUMERICAL ILLUSTRATION

The theory developed in the previous section yields a concrete prediction of the failure mode of DS when applied to the distribution $\tilde{\mu}_{n,k} = (1/2)\delta_0 + (1/2)\mu_{n,k}$ (with $\mu_{n,k}$ the law of $x = uu^T$, $u \sim \text{Unif}(B_{n,k})$). Namely (for large n, and $\sqrt{n} \ll k \ll n$):

- 1. Given sufficient model complexity and training samples, we expect the learnt denoiser $\hat{m}_n(t,\cdot)$ to achieve MSE close to 1/2 for $t < (1-\delta)t_{\rm alg}$, and close to 0 for $t > (1+\delta)t_{\rm alg}$.
- 2. We expect DS based on such a denoiser to generate samples concentrated around 0.

We tested these predictions in a numerical experiment. We considered three polytime denoisers:

- (a) The spectral-plus-projection denoiser of Algorithm 2;
- (b) A modification of the latter whereby the eigenvector calculation is replaced by 25 iterations of power method;
- (c) A learned graph neural network (GNN) (Scarselli et al., 2008; Kipf & Welling, 2016).

We carry out experiments with denoiser (b) because ℓ iterations of power method can be approximated by an ℓ -layers GNN. Hence, method (b) provides a baseline for GNN denoisers.

Figure 1, left frame, reports the MSE achieved by the three denoisers (a), (b), (c) as a function of t/n, for n=350, k=20. As GNNs are permutation-equivariant, we are training on $\approx 3\%$ of all possible outcomes, for n=350 and k=20. We observe that the GNN denoiser outperforms both the spectral algorithm and its approximation via power iteration. However, none of the three approaches can overcome the barrier at $t_{\rm alg}=n/2$, while they perform reasonably well above that threshold. This confirms the prediction at point 1 above.

On the right, we plot the histogram of Frobenius norms of samples generated with the GNN denoiser. These values are close to 0, which confirms the prediction at point 2 above. By using $\|\cdot\|_F$ as a 1-Lipschitz test function, we obtain that the Wasserstein distance between diffusion samples and the target distribution is at least 0.48 (the asymptotic prediction from theory is 0.50).

REFERENCES

- Ahmed El Alaoui, Andrea Montanari, and Mark Sellke. Sampling from the sherrington-kirkpatrick gibbs measure via algorithmic stochastic localization. In 2022 IEEE 63rd Annual Symposium on Foundations of Computer Science (FOCS), pp. 323–334. IEEE, 2022.
- Ahmed El Alaoui, Andrea Montanari, and Mark Sellke. Sampling from mean-field gibbs measures via diffusion processes. *arXiv:2310.08912*, 2023.
- Jinho Baik, Gérard Ben Arous, and Sandrine Péché. Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices. *Ann. Probab.*, 33(1):1643–1697, 2005.
- Afonso S. Bandeira and Ramon van Handel. Sharp nonasymptotic bounds on the norm of random matrices with independent entries. *The Annals of Probability*, 44(4), July 2016. ISSN 0091-1798. doi: 10.1214/15-aop1025. URL http://dx.doi.org/10.1214/15-AOP1025.
- Afonso S Bandeira, Ahmed El Alaoui, Samuel Hopkins, Tselil Schramm, Alexander S Wein, and Ilias Zadik. The franz-parisi criterion and computational trade-offs in high dimensional statistics. *Advances in Neural Information Processing Systems*, 35:33831–33844, 2022.
- Jean Barbier, Nicolas Macris, and Cynthia Rush. All-or-nothing statistical and computational phase transitions in sparse spiked matrix estimation. *Advances in Neural Information Processing Systems*, 33:14915–14926, 2020.
- Joe Benton, Valentin De Bortoli, Arnaud Doucet, and George Deligiannidis. Linear convergence bounds for diffusion models via stochastic localization. *arXiv:2308.03686*, 2023.
- Matthew Brennan, Guy Bresler, and Wasim Huleihel. Reducibility and computational lower bounds for problems with planted sparse structure. In *Conference On Learning Theory*, pp. 48–166. PMLR, 2018.
- Cristina Butucea, Yuri I Ingster, and Irina A Suslina. Sharp variable selection of a sparse submatrix in a high-dimensional noisy matrix. *ESAIM: Probability and Statistics*, 19:115–134, 2015.
- T Tony Cai, Tengyuan Liang, and Alexander Rakhlin. Computational and statistical boundaries for submatrix localization in a large noisy matrix. *The Annals of Statistics*, 45(4):1403–1430, 2017.
- Michael Celentano and Andrea Montanari. Fundamental barriers to high-dimensional regression with convex penalties. *The Annals of Statistics*, 50(1):170–196, 2022.
- 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*, pp. 4735–4763. PMLR, 2023a.
- 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. In *International Conference on Learning Representations*, 2023b.
- Yash Deshpande and Andrea Montanari. Sparse pca via covariance thresholding. *Journal of Machine Learning Research*, 17(141):1–41, 2016.
- Badih Ghazi and Euiwoong Lee. Lp/sdp hierarchy lower bounds for decoding random ldpc codes. *IEEE Transactions on Information Theory*, 64(6):4423–4437, 2017.
- Davide Ghio, Yatin Dandi, Florent Krzakala, and Lenka Zdeborová. Sampling with flows, diffusion, and autoregressive neural networks from a spin-glass perspective. *Proceedings of the National Academy of Sciences*, 121(27):e2311810121, 2024.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
 - Samuel B Hopkins, Pravesh K Kothari, Aaron Potechin, Prasad Raghavendra, Tselil Schramm, and David Steurer. The power of sum-of-squares for detecting hidden structures. In 2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS), pp. 720–731. IEEE, 2017.

- Brice Huang, Andrea Montanari, and Huy Tuan Pham. Sampling from spherical spin glasses in total variation via algorithmic stochastic localization. *arXiv:2404.15651*, 2024.
 - Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
 - Frederic Koehler and Thuy-Duong Vuong. Sampling multimodal distributions with the vanilla score: Benefits of data-based initialization. In *The Twelfth International Conference on Learning Representations*, 2024.
 - Dmitriy Kunisky, Alexander S Wein, and Afonso S Bandeira. Notes on computational hardness of hypothesis testing: Predictions using the low-degree likelihood ratio. In *ISAAC Congress* (*International Society for Analysis, its Applications and Computation*), pp. 1–50. Springer, 2019.
 - Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence of score-based generative modeling for general data distributions. In *International Conference on Algorithmic Learning Theory*, pp. 946–985. PMLR, 2023.
 - Marc Lelarge and Léo Miolane. Fundamental limits of symmetric low-rank matrix estimation, 2017. URL https://arxiv.org/abs/1611.03888.
 - Gen Li, Zhihan Huang, and Yuting Wei. Towards a mathematical theory for consistency training in diffusion models. *arXiv:2402.07802*, 2024.
 - Zongming Ma and Yihong Wu. Computational barriers in minimax submatrix detection. *The Annals of Statistics*, pp. 1089–1116, 2015.
 - Song Mei and Yuchen Wu. Deep networks as denoising algorithms: Sample-efficient learning of diffusion models in high-dimensional graphical models. *IEEE Transactions on Information Theory*, 2025.
 - Andrea Montanari. Sampling, diffusions, and stochastic localization. arXiv, 2023.
 - Andrea Montanari and Emile Richard. A statistical model for tensor pca. Advances in neural information processing systems, 27, 2014.
 - Andrea Montanari and Yuchen Wu. Posterior Sampling in High Dimension via Diffusion Processes. *arXiv:2304.11449*, 2023.
 - Andrea Montanari, Federico Ricci-Tersenghi, and Guilhem Semerjian. Solving constraint satisfaction problems through belief propagation-guided decimation. *arXiv:0709.1667*, 2007.
 - Hoi Nguyen, Terence Tao, and Van Vu. Random matrices: tail bounds for gaps between eigenvalues. *Probability Theory and Related Fields*, 167, 04 2017. doi: 10.1007/s00440-016-0693-5.
 - Minyu Peng. Eigenvalues of deformed random matrices. arXiv:1205.0572, 2012.
 - Federico Ricci-Tersenghi and Guilhem Semerjian. On the cavity method for decimated random constraint satisfaction problems and the analysis of belief propagation guided decimation algorithms. *Journal of Statistical Mechanics: Theory and Experiment*, 2009(09):P09001, 2009.
 - Tom Richardson and Ruediger Urbanke. Modern coding theory. Cambridge University Press, 2008.
 - Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE transactions on neural networks*, 20(1):61–80, 2008.
 - Tselil Schramm and Alexander S Wein. Computational barriers to estimation from low-degree polynomials. *The Annals of Statistics*, 50(3):1833–1858, 2022.
 - Kulin Shah, Sitan Chen, and Adam Klivans. Learning mixtures of gaussians using the ddpm objective. *Advances in Neural Information Processing Systems*, 36:19636–19649, 2023.
 - Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.

A NOTATIONS

 Throughout the paper it will be understood that we are considering sequences of problems indexed by n, where $x \in \mathbb{R}^{n \times n}$ and the sparsity index $k = k_n$ diverges as well. We write $f(n) \ll g(n)$ or f(n) = o(g(n)) if $f(n)/g(n) \to 0$ and $f(n) \lesssim g(n)$ or f(n) = O(g(n)) if $f(n)/g(n) \le C$ for a constant C. Finally $f(n) = \Theta(g(n))$ or $f(n) \approx g(n)$ if $1/C \le f(n)/g(n) \le C$.

We write $W \sim \mathsf{GOE}(n)$ if $W = W^\mathsf{T}$ is a random symmetric matrix with $(W_{ij})_{i \leq j \leq n}$ independent entries $W_{ii} \sim \mathsf{N}(0,2)$, and $W_{ij} \sim \mathsf{N}(0,1)$ for i < j. We say that $(W_t : t \geq 0)$ is a $\mathsf{GOE}(n)$ process if $W_t \in \mathbb{R}^{n \times n}$ is a symmetric matrix with entries above and on the diagonal $(W_t(i,j) : i < j \leq n; W_t(i,i)/\sqrt{2} : i \leq n; t \geq 0)$ forming a collection of n(n+1)/2 independent BMs

We use C, C_i, c_i, \ldots to denote absolute constants, whose value can change from line to line.

B A SIMPLE CONSEQUENCE OF CONJECTURE 3.1

We state and prove the following proposition.

Proposition B.1. Suppose that Conjecture 3.1 holds for a distribution μ with $\mathbb{E}_{\boldsymbol{x} \sim \mu}[\|\boldsymbol{x}\|^2] = 1$. Then for any sequence of times $t = t_n \leq (1 - \delta)t_{\text{alg}}$,

$$\mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) - \boldsymbol{x}\|^2] = 1 - o(1) \Rightarrow \mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t)\|^2] = o(1)$$

In words, if \hat{m} is (near)-optimal in score matching for $t \leq (1 - \delta)t_{alg}$, then $\|\hat{m}(y_t, t)\|$ is small.

Before giving the proof, we remark that the full Conjecture 3.1 is not needed. It suffices for \hat{m} to have a weaker property; namely, that for any fixed constants $c \in [-1, 1]$ and $\delta \in (0, 1)$,

$$\inf_{t \leq (1-\delta)t_{\text{alg}}} \mathbb{E}[\|c\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) - \boldsymbol{x}\|^2] \geq 1 - o(1)$$

Proof. Fix $c \in [-1, 1]$ to be a constant chosen later. From the property of \hat{m} , we get from Cauchy-Schwarz that

$$\frac{1}{2}\mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t)\|^2] - \mathbb{E}[\|\boldsymbol{x}\|^2] \le 1 - o(1) \Rightarrow \mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t)\|^2] \le 4 - o(1)$$

We use Conjecture 3.1 for the sequence of estimators $c\hat{m}$, which states that uniformly over $t \leq (1-\delta)t_{\text{ale}}$:

$$\mathbb{E}[\|c\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t)-\boldsymbol{x}\|^2] \ge 1 - o(1) \Rightarrow c^2 \mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t)\|^2] - 2c \mathbb{E}[\langle \hat{\boldsymbol{m}}(\boldsymbol{y}_t,t), \boldsymbol{x} \rangle] \ge -o(1)$$

Suppose for sake of contradiction, that $\limsup_{n\to\infty} |\mathbb{E}[\langle \hat{\boldsymbol{m}}(\boldsymbol{y}_t,t),\boldsymbol{x}\rangle]| \geq \beta > 0$. Without loss of generality, we consider the subsequence (n_k) such that $\mathbb{E}[\langle \hat{\boldsymbol{m}}(\boldsymbol{y}_t,t),\boldsymbol{x}\rangle] \geq \beta/2$. Along this subsequence, we have

$$4c^2 - c\beta \ge -o(1)$$

for all $c \in [-1, 1]$. However, we know that this is not true for c > 0 small enough; specifically, take $c < \beta/8$ so that we have $-c\beta/2 \ge -o(1)$, contradiction. Hence $\mathbb{E}[\langle \hat{\boldsymbol{m}}(\boldsymbol{y}_t, t), \boldsymbol{x} \rangle] = o(1)$. From the property of $\hat{\boldsymbol{m}}$, we obtain the conclusion.

C APPLYING THEOREM 1 TO VERY SPARSE MATRICES

As mentioned in Section D.3, we state and prove an analogous version of Corollary 3.2 in the very sparse case. One different aspect from the moderate case is that k can be smaller asymptotically: in particular, k can be sub-polynomial in n. Therefore, we first give a modification of Assumption 1.

Assumption 3. Consider $\underline{k}_n \ll k \ll n$ for \underline{k}_n in Conjecture 3.1. Let $\mathbf{y}_t = t\mathbf{x} + \mathbf{W}_t$ for (\mathbf{W}_t) sBM independent of \mathbf{x} . Then a near-optimal estimator $\hat{\mathbf{m}}_0(\mathbf{y},t)$ in score-matching satisfies: for every pair $(\gamma, \varepsilon) \in (0,1)$,

$$\int_{0}^{(1-\gamma)t_{\text{alg}}} \mathbb{P}(\|\hat{\boldsymbol{m}}_{0}(\boldsymbol{y}_{t},t)\| \geq \varepsilon) \, \mathrm{d}t = O(k^{-D})$$

for every fixed D > 0.

Corollary C.1. Assume $(\log n)^2 \ll k \ll n$, so that $t_{\rm alg}(n,k) := k^2 \log(n/k^2)$. Let $\hat{\boldsymbol{m}}_0$ be an arbitrary poly-time algorithm such that $\sup_{\boldsymbol{y},t} \|\hat{\boldsymbol{m}}_0(\boldsymbol{y},t)\|_F \leq 1$ and Assumption 3 holds. Then there exists an estimator $\hat{\boldsymbol{m}}$ such that

M1. $\hat{m}(\cdot)$ can be evaluated in polynomial time.

M2. If $y_t = tx + B_t$ is the true diffusion (equivalently given by (1)), then, for every D > 0,

$$\int_0^\infty \mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) - \hat{\boldsymbol{m}}_0(\boldsymbol{y}_t,t)\|^2] dt = O(k^{-D}).$$

M3. There exists $\delta = o_n(1)$ such that, for any step size $\Delta = \Delta_n > 0$, we have incorrect sampling:

$$\inf_{t \in \mathbb{N} \cdot \Delta, t \ge (1+\delta)t_{\text{alg}}} W_1(\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_t, t), \boldsymbol{x}) \ge 1 - o_n(1). \tag{19}$$

Proof. By the blueprint Theorem 1, we find (a sequence of) hypothesis tests $\phi(\boldsymbol{y},t)$ indexed by t such that Assumption 2 holds. We choose a rate $\delta_n = o_n(1)$ slow enough, and ε_n be the resulting sequence, such that Proposition H.1 holds. We now describe $\phi(\boldsymbol{y},t)$, based on Algorithm 1, from time $t = (1+\delta)t_{\rm alg} = (1+\delta)k^2\log(n/k^2)$ upto t = n:

- Let $s = \sqrt{(1+\varepsilon_n)\log(n/k^2)}$. Compute $\mathbf{y}_+ = \mathbf{y} + \sqrt{\varepsilon_n t}\mathbf{g}$ and $\mathbf{y}_- = \mathbf{y} \sqrt{t/\varepsilon_n}\mathbf{g}$, with $\mathbf{g} \sim \mathsf{N}(\mathbf{0}, \mathbf{I})$. Then, compute $\mathbf{A}_+ = (\mathbf{y}_+ + \mathbf{y}_+^\mathsf{T})/(2\sqrt{t})$, and $\mathbf{A}_- = (\mathbf{y}_- + \mathbf{y}_-)/(2\sqrt{t})$.
- Let v be the leading eigenvector of $\eta_s(\mathbf{A}_+)$. Then, let $\hat{v} = \mathbf{A}_-v$. Let \hat{S} be the set of k indices of \hat{v} with largest magnitude, and compute w such that $w_i = (1/\sqrt{k})\operatorname{sign}(\hat{v}_i)\mathbf{1}_{i\in\hat{S}}$.
- Finally, reject iff $\langle w, yw \rangle \ge \beta t$, for some $1 > \beta > c$.

From Proposition H.1, we know that

$$\sup_{t \ge (1+\delta)t_{\text{alg}}} \mathbb{P}(\boldsymbol{w}(\boldsymbol{y}_t,t) \ne \boldsymbol{x}) \ll n^{-D}$$

for every D>0. On the event that ${\boldsymbol w}({\boldsymbol y}_t,t)={\boldsymbol x},$ we get that

$$\langle \boldsymbol{w}(\boldsymbol{y}_t,t),\boldsymbol{y}_t\boldsymbol{w}(\boldsymbol{y}_t,t)\rangle = t + \langle \boldsymbol{w}(\boldsymbol{y}_t,t),\boldsymbol{W}_t\boldsymbol{w}(\boldsymbol{y}_t,t)\rangle \geq t - \sup_{\boldsymbol{v}:\|\boldsymbol{v}\|_0 = k, v_i \in \{0,\pm 1/\sqrt{k}\}} \langle \boldsymbol{v},\boldsymbol{W}_t\boldsymbol{v}\rangle$$

for (W_t) standard Brownian motion. From Lemma G.1, we get that with error probability at most $\binom{n}{k}^{-D}$ for some D, we get that

$$\langle \boldsymbol{w}(\boldsymbol{y}_t,t), \boldsymbol{y}_t \boldsymbol{w}(\boldsymbol{y}_t,t) \rangle = t + \langle \boldsymbol{w}(\boldsymbol{y}_t,t), \boldsymbol{W}_t \boldsymbol{w}(\boldsymbol{y}_t,t) \rangle \ge t - C \sqrt{t \log \binom{n}{k}} = t(1 - o(1))$$

Therefore, we obtain that for $\beta < 1$,

$$\sup_{t \ge (1+\delta)t_{\text{alg}}} \mathbb{P}(\phi(\boldsymbol{y}_t, t) = 0) \ll n^{-D}$$

for any D>0. After time t=n, we use the same tests ϕ_1,ϕ_2 as documented in the proof of Corollary 3.2, as $t=n>(1+\delta)(n/2)$, where n/2 is the algorithmic threshold of the moderately sparse case. The reason we can do this is that the spectral method, as in Algorithm 2, works even when $k\ll \sqrt{n}$ (although the threshold for this algorithm is asymptotically worse than $k^2\log(n/k^2)$). Furthermore, the size of the perturbation $a\in \mathcal{A}_d(c)$ is at most $\|a\|\leq ct_{\rm alg}=ck^2\log(n/k^2)\ll c(n/2)$.

Consequently, the first condition of Assumption 2 holds with rate n^{-D} for every D > 0. To deal with the second condition, note simply that w is a k-sparse vector. A close inspection of the proof of Corollary 3.2 shows that it does not really matter how w is computed; the main idea is simply that for all $a \in \mathcal{A}_d(c)$,

$$\langle \boldsymbol{w}, (\boldsymbol{a} + \boldsymbol{B}_t) \boldsymbol{w} \rangle \le \|\boldsymbol{a}\| + \langle \boldsymbol{w}, \boldsymbol{B}_t \boldsymbol{w} \rangle \le c t_{\text{alg}} + \langle \boldsymbol{w}, \boldsymbol{B}_t \boldsymbol{w} \rangle \le c t_{\text{alg}} + C \sqrt{t \log \binom{n}{k}}$$

for each t. Of course, we have to bound this simultaneously for all t, and this is done in the proof of Corollary 3.2; c.f. Appendix G.

CONCRETE EXAMPLES: DENOISERS FOR SPARSE LOW-RANK MATRICES

D.1 ALGORITHMS

703 704

705 706

708

709 710

711 712

713

714

715 716 717

718

719

720

721

727

728

729

730

731 732

733

734 735

736

738

739

740

741

742

743

744 745

746 747

748

749 750 751

752

753

754

755

In this section, we provide the detailed pseudocode for Algorithms 1 and 2. In Algorithm 1 we use the following soft-thresholding function, with a parameter s:

$$\eta_s(y) = \text{sign}(y) \max(|y| - t, 0) = \text{sign}(y)(|y| - t)_+$$

Algorithm 1 Submatrix Estimation Algorithm (very sparse regime)

- 1: **Input:** Data y_t ; time t; parameters s, ε
- 2: Output: Estimate of \boldsymbol{x} : $\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t)$
- 3: Let $g_t \sim \mathsf{N}(0, t I_{n \times n})$ and compute $y_{t,+} := y_t + \sqrt{\varepsilon} g_t, y_{t,-} := y_t g_t/\sqrt{\varepsilon}$
- 4: Symmetrize: $A_{t,+}=(y_{t,+}+y_{t,+}^{\mathsf{T}})/(2\sqrt{t}), A_{t,-}=(y_{t,-}+y_{t,-}^{\mathsf{T}})/(2\sqrt{t})$ 5: Compute top eigenvector of $\eta_s(A_{t,+})$, denoted if by v_t
- 6: If $t \ge t_{\text{alg}} \lor 1$ and $\lambda_1\left(\eta_s\left(\boldsymbol{A}_{t,+}\right)\right) > k + \frac{\sqrt{t}}{s}$, continue; otherwise return $\hat{\boldsymbol{m}}(\boldsymbol{y},t) := \boldsymbol{0}$
- 7: Compute the vector $\hat{\boldsymbol{v}}_t := \boldsymbol{A}_{t,-} \boldsymbol{v}_t$
 - 8: Let \hat{S} be the set of k indices i of largest values of $|\hat{v}_{t,i}|$, and compute vector w such that $w_i = \operatorname{sign}(\hat{v}_{t,i}) \mathbf{1}_{i \in \hat{S}}$
 - 9: **return** $\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) := \mathbf{1}_{\hat{\boldsymbol{S}}} \mathbf{1}_{\hat{\boldsymbol{S}}}^{\mathsf{T}}/k$

Algorithm 2 Submatrix Estimation Algorithm (moderately sparse regime)

- 1: **Input:** Data y_t ; time t; parameter ε
- 2: Output: Estimate of x: $\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t)$
- 3: If $t \ge t_{\text{alg}}$, continue; otherwise return $\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) = \boldsymbol{0}$
- 4: Symmetrize: $\boldsymbol{A}_t = (\boldsymbol{y}_t + \boldsymbol{y}_t^\mathsf{T})/(2\sqrt{t})$
- 5: If $t \geq n^2$ and $\lambda_1(\boldsymbol{A}_t) \leq \sqrt{t}/2$, return 0; otherwise continue 6: Compute top eigenvector of \boldsymbol{A}_t , denoted if by \boldsymbol{v}_t
- 7: Compute \hat{S} by $\hat{S}:=\left\{i\in[n]:\ |v_{t,i}|\geq \frac{\varepsilon}{\sqrt{k}}\right\}$
- 8: Compute vector \boldsymbol{w} such that $w_i = \operatorname{sign}(v_{t,i}) \mathbf{1}_{i \in \hat{S}}$
- 9: **return** $\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) := \boldsymbol{w}\boldsymbol{w}^{\mathsf{T}}/|\hat{S}| \text{ if } |\hat{S}| \geq k/2; \text{ otherwise return } \boldsymbol{0}$

Moderately sparse regime: $\sqrt{n} \ll k \ll n$

Since Theorem 3.2 is somewhat abstract, we complement it with an explicit example of \hat{m} : namely, it is a modification of a standard spectral estimator. While achieving near optimal estimation error (among polytime algorithms), \hat{m} fails to generate samples from the correct distribution.

Proposition D.1. Assume $\sqrt{n} \ll k \ll n$, so that $t_{alc}(n,k) := n/2$ per (9). Then the estimator \hat{m} defined in Algorithm 2 satisfies the following:

- M1. $\hat{\boldsymbol{m}}(\cdot)$ can be evaluated in polynomial time.
- **M2.** For any $\delta > 0$, there exists $c = c(\delta)$, $C = C(\delta)$ such that

$$\inf_{t \leq (1-\delta)t_{\mathrm{alg}}} \mathbb{E} \big\{ \| \hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) - \boldsymbol{x} \|^2 \big\} = 1 - o_n(1) \,, \quad \sup_{t \geq (1+\delta)t_{\mathrm{alg}}} \mathbb{E} \big\{ \| \hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) - \boldsymbol{x} \|^2 \big\} \leq C \, e^{-cn/k} \,.$$

M3. For any $\Delta > 0$, we have incorrect sampling: $\inf_{t \in \mathbb{N} \cdot \Delta} W_1(\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_t, t), \boldsymbol{x}) = 1 - o_n(1)$.

Therefore, we enforce that $\hat{m} \equiv 0$ for $t < t_{\rm alg}$. Recall that this implies Assumption 1 holds trivially. Regarding the specific design of \hat{m} , Algorithm 2 uses a thresholded spectral approach. We compute the leading eigenvector of (the symmetrized version of) y_t , call it $v_t \in \mathbb{R}^n$. We then estimate the support S of the latent rank-one matrix x using the entries of v_t with largest magnitude. D.3 Very sparse regime: $k \ll \sqrt{n}$

 We have an analogous result for the very sparse regime, where the sparsity level $k \ll \sqrt{n}$.

Proposition D.2. Assume $(\log n)^{5/2} \lesssim k \ll \sqrt{n}$, and note that here $t_{\text{alg}}(n,k) = k^2 \log(n/k^2)$, per (9). Then the randomized estimator $\hat{\mathbf{m}} : \mathbb{R}^{n \times n} \times \mathbb{R} \to \mathbb{R}^{n \times n}$ of Algorithm 1 satisfies the following:

M1. $\hat{\boldsymbol{m}}(\cdot)$ can be evaluated in polynomial time.

M2. For any $\delta > 0$ and D > 0:

$$\inf_{t \leq (1-\delta)t_{\rm alg}} \mathbb{E} \big\{ \| \hat{\boldsymbol{m}}(\boldsymbol{y}_t, t) - \boldsymbol{x} \|^2 \big\} = 1 - o_n(1) \,, \quad \sup_{t \geq (1+\delta)t_{\rm alg}} \mathbb{E} \big\{ \| \hat{\boldsymbol{m}}(\boldsymbol{y}_t, t) - \boldsymbol{x} \|^2 \big\} \ll n^{-D} \,.$$

M3. For any $\Delta > 0$, we have incorrect sampling: $\inf_{t \in \mathbb{N} \cdot \Delta} W_1(\hat{\boldsymbol{n}}(\hat{\boldsymbol{y}}_t, t), \boldsymbol{x}) = 1 - o_n(1)$.

The pseudocode for the estimator $\hat{\boldsymbol{m}}(\cdot)$ that is constructed in the above is given as Algorithm 1. This is based on a standard approach in the literature Deshpande & Montanari (2016); Cai et al. (2017), with some modifications to allow for its analysis in the diffusion setting. The main steps are as follows: (1) Perform Gaussian data splitting of \boldsymbol{y}_t into $\boldsymbol{y}_{t,+}, \boldsymbol{y}_{t,-}$, see Line 3 of Algorithm 1, with most of the information preserved in $\boldsymbol{y}_{t,+}$. (2) Use entrywise soft thresholding $\eta_s(x) = (|x|-s)_+ \operatorname{sign}(x)$ to reduce the noise in the symmetrized version of $\boldsymbol{y}_{t,+}$. (3) Compute a first estimate of the latent vector $\mathbf{1}_S$ by the principal eigenvector of the above matrix. (4) Refine this estimate using the remaining information $\boldsymbol{y}_{t,-}$.

We point out that Proposition 3.3 remains true in the regime $\sqrt{n} \ll k \ll n$, and hence we observe a gap between $t_{\rm alg}(n,k)$ and $t_{\rm Bayes}(n,k)$ in this regime as well.

E PROOF OF PROPOSITION D.1

E.1 PROPERTIES OF THE ESTIMATOR $\hat{\boldsymbol{m}}(\cdot)$

Proposition E.1. Assume $\sqrt{n} \ll k \ll n$, and note that in this case $t_{\text{alg}}(n,k) = n/2$. Let $\hat{\boldsymbol{m}}(\cdot)$ be the estimator of Algorithm 2 with input parameter ε . For every $\delta > 0$, there exists $\varepsilon > 0$ such that

$$\sup_{t \ge (1+\delta)t_{\text{alg}}} \mathbb{E}\left[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t, t) - \boldsymbol{x}\|^2\right] \le C e^{-n\varepsilon^2/64k}. \tag{20}$$

The proof of this proposition is standard, and will be presented in Appendix N. We note that the rate in Equation (20) gets slower the closer k is to n; it is super-polynomial if $n \gg k \log n$.

By definition, when $\sqrt{n} \ll k \ll n$ and $t < t_{\text{alg}}(n,k)$, the Algorithm 2 returns $\hat{\boldsymbol{m}}(\boldsymbol{y},t) = \boldsymbol{0}$, so we automatically have the following result.

Proposition E.2. For any fixed $\delta > 0$, and $t \leq (1 - \delta)t_{alg}$, we have $\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t, t) - \boldsymbol{x}\| = 1$.

E.2 AUXILIARY LEMMAS

The following lemmas are needed for the analysis of the generated diffusion. Their proofs are deferred to Appendices O, P, Q, R.

Lemma E.3. Let W_t be a GOE process. Then for each time $t_0 \ge 0$,

$$\mathbb{P}\left(\max_{0 \le t \le t_0} \|\boldsymbol{W}_t\|_{\text{op}} \ge 16\sqrt{t_0 n}\right) \le 2\exp\left(-32n\right).$$

Lemma E.4. Let W_t be a GOE process, and let v_t be any eigenvector of W_t for every $t \ge 0$. Define the set

$$A(\boldsymbol{v}_t; C) = \left\{ i : 1 \le i \le n, |v_{ti}| \ge \frac{C\sqrt{\log(n/k)}}{\sqrt{n}} \right\}$$

Then for any C > 4, we have

$$\mathbb{P}\left(|A(\boldsymbol{v}_t;C)| \geq \max\{\sqrt{k},k^2/n\}\right) = O\left(\exp(-(1/3)n^{1/4})\right)$$

As a consequence, using this eigenvector, \hat{m} will evaluate to 0 per line 8 of Algorithm 2.

Lemma E.5. Let W_t be a GOE process, and for each t, let v_t be a top eigenvector of W_t . Then for any times $t_0 \le t_1$, with probability at least $1 - 2 \exp(-32n)$,

$$\sup_{t_0 \le t \le t_1} |\langle \boldsymbol{v}_t, \boldsymbol{W}_{t_0} \boldsymbol{v}_t \rangle - \lambda_1(\boldsymbol{W}_{t_0})| \le 32\sqrt{n(t_1 - t_0)}.$$

Lemma E.6 (Concentration for deformed GOE model). Consider the model $\mathbf{Y} = \theta v v^{\mathsf{T}} + \mathbf{W}$ for $\mathbf{W} \sim \mathsf{GOE}(n)/\sqrt{n}$ and $\theta > 1$ a constant, v a unit vector. Let $v_1(\mathbf{Y})$ be the top eigenvector of \mathbf{Y} . Define $(x^*, u^*) = (\theta + 1/\theta, 1 - 1/\theta^2)$. For any closed set F such that $d((x^*, u^*), F) > 0$, there exists a constant c > 0 such that

$$\mathbb{P}\left((\lambda_1(\boldsymbol{Y}), \langle \boldsymbol{v}_1(\boldsymbol{Y}), \boldsymbol{v} \rangle^2) \in F\right) \le \exp(-cn)$$

for all n large enough.

We only use Lemma E.6 for the alignment $\langle v_1(Y), v \rangle^2$.

E.3 ANALYSIS OF THE DIFFUSION PROCESS: PROOF OF PROPOSITION D.1

We will prove Theorem D.1 for $1 \gg \varepsilon \ge C \sqrt{\log(n/k)/(n/k)}$ for some sufficiently large constant C.

Suppose that we generate the following diffusion, with $(z_t)_{t\geq 0}$ a standard n^2 -dimensional BM, and $\hat{y}_0 = 0$:

$$\hat{\boldsymbol{y}}_{\ell\Delta} = \hat{\boldsymbol{y}}_{(\ell-1)\Delta} + \Delta \cdot \hat{\boldsymbol{m}} \left(\hat{\boldsymbol{y}}_{(\ell-1)\Delta}, (\ell-1)\Delta \right) + \left(\boldsymbol{z}_{\ell\Delta} - \boldsymbol{z}_{(\ell-1)\Delta} \right).$$

We will prove that the generated diffusion never passes the termination conditions (c.f. Algorithm 2, lines 3, 5, 8).

E.3.1 Analysis up to an intermediate time

Define $t_{\text{between}} = n^2$. Following the same strategy with Section H.2, we will first show that $\hat{m} = 0$ up to t_{between} with high probability by analyzing only the noise process (in short, if $\hat{m} = 0$ always, our generated diffusion coincides with the noise process). Our strategy is of the same nature as that of Section H.2. Indeed, we will attempt to prove that $\hat{m} = 0$ simultaneously for all t, with high probability. In this phase $(0 \le t \le t_{\text{between}})$, we will show that $|\hat{S}| < k/2$ (c.f. definition in Algorithm 2, lines 7, 8) for $0 \le t \le t_{\text{between}}$, with high probability (v_t is the top eigenvector of A_t , c.f. Algorithm 2). Note that line 5 of Algorithm 2 is not relevant in this phase. We first show this for a sequence of time points $\{t_\ell\}_{\ell \ge 1}$, then control the in-between fluctuations. We can set t_1 to be any value in [0, n/2), as the algorithm returns 0 if $t < t_{\text{alg}} = n/2$ anyway. We denote the GOE process

$$\boldsymbol{W}_t = \frac{\boldsymbol{B}_t + \boldsymbol{B}_t^\mathsf{T}}{2} = \sqrt{t} \boldsymbol{A}_t.$$

It is clear that the eigenvectors of W_t and A_t coincide.

We choose the following time points:

$$t_{\ell} = \frac{n}{2} - 1 + \frac{\ell}{n^4}.$$

To exceed $t_{\text{between}} = n^2$, we will need n^6 values of ℓ . By union bound from Lemma E.4 (recall also the definition of the set A(v; C) from this Lemma),

$$\mathbb{P}\left(\exists 1 \le \ell \le n^6 : |A(\mathbf{v}_{t_{\ell}}; C)| \ge \max\{\sqrt{k}, k^2/n\}\right) \le O\left(\exp(-(1/3)n^{1/4} + 6\log n)\right)$$
 (21)

Next, we will control the in-between fluctuations; specifically, we would like to show that $\max_{t_\ell \leq t \leq t_{\ell+1}} |A(\boldsymbol{v}_t;C)| \leq C_0 \max\{\sqrt{k},k^2/n\}$ simultaneously for many values of ℓ (with high probability), for some constant $C_0 > 0$. Our approach is as follows.

(i) Let v_t be a top eigenvector of W_t . If t is close to t_ℓ , then v_t is an approximate solution to the equation (in v):

$$\boldsymbol{v}^\mathsf{T} \boldsymbol{W}_{t_\ell} \boldsymbol{v} = \lambda_1(\boldsymbol{W}_{t_\ell})$$

- (ii) \boldsymbol{v}_t can be written in the coordinate system of the orthonormal eigenvectors $\boldsymbol{U}_{t_\ell} = [\boldsymbol{u}_1|\cdots|\boldsymbol{u}_n]$ of \boldsymbol{W}_{t_ℓ} , corresponding to decreasing eigenvalues $\lambda_1(\boldsymbol{W}_{t_\ell}) \geq \cdots \geq \lambda_n(\boldsymbol{W}_{t_\ell})$. Namely, $\boldsymbol{v}_t = \boldsymbol{U}_{t_\ell} \boldsymbol{U}_{t_\ell}^\mathsf{T} \boldsymbol{v}_t = \boldsymbol{U}_{t_\ell} \boldsymbol{w}$ with $\|\boldsymbol{w}\| = 1$.
- (iii) Let m be a (sufficiently large) constant integer with $1 \le m \le n$. The first m components of w take up 1 o(1) in L_2 -norm by (i) and (ii) with overwhelming probability, from which we can simply use triangle inequality to upper bound $|A(v_t; C)|$ according to Lemma E.4 for u_1, \dots, u_m , which incurs only a constant factor of error probability, by union bound.

Define $p_n = P\left(|\lambda_1(\boldsymbol{W}_1) - \lambda_7(\boldsymbol{W}_1)| \le n^{-C'-1/2}\right)$ for any C' > 0 (here we take m = 7). We use the following result (we have accounted for the scaling).

Lemma E.7 (Corollary 2.5, Nguyen et al. (2017)). Let $W_1 \sim (1/\sqrt{2})\mathsf{GOE}(n)$. For any fixed $l \geq 1, C' > 0$, there exists a constant $c_0 = c_0(l, C')$ such that

$$\mathbb{P}\left(\lambda_1(\mathbf{W}_1) - \lambda_{1+l}(\mathbf{W}_1) \le \frac{1}{2}n^{-C'-1/2}\right) \le c_0 n^{-C' \cdot \frac{l^2+2l}{3}}.$$

We materialize our approach above. We can write, with $v_t = U_{t_\ell} w$:

$$oldsymbol{v}_t^\mathsf{T} oldsymbol{W}_{t_\ell} oldsymbol{v}_t = oldsymbol{w}^\mathsf{T} oldsymbol{D}_{t_\ell} oldsymbol{w} = \sum_{i=1}^n (D_{t_\ell})_{ii} w_i^2.$$

We then obtain that p

$$\boldsymbol{v}_t^\mathsf{T} \boldsymbol{W}_{t_\ell} \boldsymbol{v}_t - \lambda_1(\boldsymbol{W}_{t_\ell}) \leq \sum_{i=8}^n (\lambda_i(\boldsymbol{W}_{t_\ell}) - \lambda_1(\boldsymbol{W}_{t_\ell})) w_i^2 < -\frac{1}{2} \sqrt{t_\ell} n^{-3/2} \sum_{i=8}^n w_i^2$$

with probability at least $1 - p_n \ge 1 - c_0 n^{-8}$, from Lemma E.7 and $W_{t_\ell} \sim \sqrt{t_\ell} W_1$. Now from Lemma E.5, we know that with probability at least $1 - 2 \exp(-32n)$,

$$\boldsymbol{v}_t^\mathsf{T} \boldsymbol{W}_{t_\ell} \boldsymbol{v}_t - \lambda_1(\boldsymbol{W}_{t_\ell}) \ge -32\sqrt{n(t_{\ell+1} - t_\ell)}.$$

With probability at least $1 - c_0 n^{-8} - 2 \exp(-32n)$, both of these statements are true, uniformly over $t_{\ell} \le t \le t_{\ell+1}$, leading to

$$64\sqrt{\frac{t_{\ell+1}-t_{\ell}}{t_{\ell}}} \ge n^{-3/2} \sum_{i=0}^{n} w_i^2.$$

A simple bit of algebra shows that

$$\sqrt{\frac{t_{\ell+1} - t_{\ell}}{t_{\ell}}} \le 2n^{-5/2} \Rightarrow \sum_{i=8}^{n} w_i^2 \le 128n^{-1}.$$

Consider the first 7 eigenvectors $\{m{u}_{t_\ell,i}\}_{i=1}^7$ of $m{W}_{t_\ell}$. Let

$$A = \bigcup_{i=1}^{7} A(\boldsymbol{u}_{t_{\ell},i}; C).$$

From Lemma E.4 and a union bound, that $|A| \leq 7 \max\{\sqrt{k}, k^2/n\}$ with probability at least $1 - O(\exp(-(1/3)n^{1/4}))$. For every $j \in A^c$, we have

$$|v_{tj}| \le \sum_{i=1}^{n} |w_i| \cdot |u_{t_{\ell},i,j}| < \sum_{i=1}^{7} |w_i| \cdot \frac{C\sqrt{\log(n/k)}}{\sqrt{n}} + \sum_{i=8}^{n} |w_i| \le \frac{C'\sqrt{\log(n/k)}}{\sqrt{n}} + \frac{128}{\sqrt{n}} < C''\sqrt{\frac{\log(n/k)}{n}}$$

for a large enough constant C'' > 0. This means that with probability at least $1 - O(n^{-8})$,

$$\sup_{t_{\ell} \le t \le t_{\ell+1}} |A(v_t; C'')| \le 7 \max\{\sqrt{k}, k^2/n\} < k/2$$

From Equation (21) and a union bound over $\ell \geq 1$, we know that with high probability,

$$\sup_{n/2-1 \le t \le n^2} |A(v_t; C)| \le 7 \max\{\sqrt{k}, k^2/n\}$$

for some absolute constant C>0, meaning that $\hat{m{m}}=0$ up to $t_{ ext{between}}$, as long as

$$\varepsilon > \frac{C\sqrt{\log(n/k)}}{\sqrt{n/k}}$$

E.3.2 Analysis to the infinite horizon

We will prove that simultaneously for all $t \ge n^2$, Algorithm 2 always terminates at line 5, or that $\lambda_1(\mathbf{W}_t) \le t/2$. Similar to Subsection E.3.1, we choose the following sequence of time points for all $\ell \ge 1$:

$$t_{\ell}^{(2)} = n^2 + \ell - 1$$

By standard Gaussian concentration and the Bai-Yin theorem, we get, for instance, the following tail bound (constants are loose) for all $x \ge 0$:

$$\mathbb{P}\left(\lambda_1\left(\frac{\boldsymbol{W}_t}{\sqrt{t}}\right) \ge 4\sqrt{n} + x\right) \le 2\exp\left(-\frac{x^2}{2}\right)$$

Set $x = \frac{\sqrt{t}}{8}$. Since $t \ge n^2$, we have $x \ge n/8$, and so $x + 4\sqrt{n} \le 2x$ for n large enough. Consequently,

$$\mathbb{P}\left(\lambda_1\left(\frac{\boldsymbol{W}_t}{\sqrt{t}}\right) \geq \frac{\sqrt{t}}{4}\right) \leq 2\exp\left(-c_1t\right)$$

for some universal constant $c_1 > 0$. A union bound for the chosen points gives:

$$\mathbb{P}\left(\exists \ell \ge 1 : \lambda_1\left(\boldsymbol{W}_{t_{\ell}^{(2)}}\right) \ge t_{\ell}^{(2)}/4\right) \le 2\sum_{\ell=1}^{\infty} \exp\left(-c_1 t_{\ell}^{(2)}\right) = 2\exp(-c_1 n^2) \sum_{\ell=1}^{\infty} \exp(-c_1 (\ell-1)) \lesssim \exp(-c_1 n^2)$$
(22)

Next we control the in-between fluctuations. From a simple modification of Lemma E.3, we have

$$\mathbb{P}\left(\sup_{t_{\ell}^{(2)} \leq t \leq t_{\ell+1}^{(2)}} \left| \lambda_1 \left(\boldsymbol{W}_{t_{\ell}^{(2)}} \right) - \lambda_1 \left(\boldsymbol{W}_{t} \right) \right| \geq 16\sqrt{\left(t_{\ell+1}^{(2)} - t_{\ell}^{(2)}\right) \cdot t_{\ell}^{(2)} n} \right) \leq 2 \exp\left(-32nt_{\ell}^{(2)}\right)$$

so that by union bound

$$\mathbb{P}\left(\exists \ell \geq 1: \sup_{t_{\ell}^{(2)} \leq t \leq t_{\ell+1}^{(2)}} \left| \lambda_1\left(\boldsymbol{W}_{t_{\ell}^{(2)}}\right) - \lambda_1\left(\boldsymbol{W}_{t}\right) \right| \geq 16\sqrt{\left(t_{\ell+1}^{(2)} - t_{\ell}^{(2)}\right) \cdot t_{\ell}^{(2)} n} \right) \lesssim \exp(-32n^3) \tag{23}$$

Consider the intersection of events described in Equations (22) and (23):

$$A = \left\{ \exists \ell \ge 1 : \lambda_1 \left(\mathbf{W}_{t_{\ell}^{(2)}} \right) \ge t_{\ell}^{(2)} / 4 \right\} \cup \left\{ \exists \ell \ge 1 : \sup_{t_{\ell}^{(2)} \le t \le t_{\ell+1}^{(2)}} \left| \lambda_1 \left(\mathbf{W}_{t_{\ell}^{(2)}} \right) - \lambda_1 \left(\mathbf{W}_{t} \right) \right| \ge 16 \sqrt{(t_{\ell+1}^{(2)} - t_{\ell}^{(2)}) \cdot t_{\ell}^{(2)} n} \right\}$$

For each $t \ge n^2$, let t_ℓ be largest such that $t_\ell \le t < t_{\ell+1}$. On A, we have

$$\lambda_1(\boldsymbol{W}_t) \le \frac{t_\ell^{(2)}}{4} + 16\sqrt{(t_{\ell+1}^{(2)} - t_\ell^{(2)}) \cdot t_\ell^{(2)} n} = \frac{t_\ell^{(2)}}{4} + 16\sqrt{t_\ell^{(2)} n} \le \frac{t_\ell^{(2)}}{2} \le \frac{t}{2}$$

for n large enough, since $t_{\ell}^{(2)} \ge n^2 \gg n$. Hence the algorithm always returns $\mathbf{0}$ with high probability, and we are done.

F Proof of Theorem 1

We define \hat{m} as follows:

$$\hat{\boldsymbol{m}}(\boldsymbol{y},t) = \begin{cases} \hat{\boldsymbol{m}}_0(\boldsymbol{y},t) \mathbf{1}_{\parallel \hat{\boldsymbol{m}}_0(\boldsymbol{y},t) \parallel \leq \varepsilon} & \text{if } t \leq (1-\gamma)t_{\text{alg}}, \\ \hat{\boldsymbol{m}}_0(\boldsymbol{y},t) & \text{if } (1-\gamma)t_{\text{alg}} < t < (1+\delta)t_{\text{alg}}, \\ \hat{\boldsymbol{m}}_0(\boldsymbol{y},t)\phi(\boldsymbol{y},t) & \text{if } t \geq (1+\delta)t_{\text{alg}}, \end{cases}$$

First, we check that Condition M2 holds.

$$\int_0^\infty \mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) - \hat{\boldsymbol{m}}_0(\boldsymbol{y}_t,t)\|^2]dt$$

$$= \int_0^{(1-\gamma)t_{\text{alg}}} \mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) - \hat{\boldsymbol{m}}_0(\boldsymbol{y}_t,t)\|^2]dt + \int_{(1+\delta)t_{\text{alg}}}^{\infty} \mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) - \hat{\boldsymbol{m}}_0(\boldsymbol{y}_t,t)\|^2]dt$$

$$\leq \int_0^{(1-\gamma)t_{\text{alg}}} \mathbb{P}(\|\hat{\boldsymbol{m}}_0(\boldsymbol{y}_t,t)\| > \varepsilon)dt + \int_{(1+\delta)t_{\text{alg}}}^{\infty} \mathbb{P}(\phi(\boldsymbol{y},t) = 0) dt$$

$$= O(\eta_1(d) + \eta_2(d))$$

where in (a) we use the fact that ϕ is binary and Conditions i), ii) and iii.1). Secondly, we check that Condition M3 holds. To reduce notational clutter, we assume that $t_1/\Delta = \ell_0$ is an integer. Then, we can write

$$\hat{\boldsymbol{y}}_{t_1} = \boldsymbol{B}_{t_1} + \Delta \sum_{i=1}^{\ell_0} \hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_{i\Delta}, i\Delta)$$
 (24)

where the drift accumulation term is bounded by, from Condition i):

$$\left\| \Delta \sum_{i=1}^{\ell_0} \hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_{i\Delta}, i\Delta) \right\| \leq \varepsilon (1-\gamma) t_{\text{alg}} + (\gamma + \delta) t_{\text{alg}} \leq c \cdot t_{\text{alg}}$$

for every $c \in (0,1)$ by taking ε, γ small enough, as $\delta = o_d(1)$. Suppose that from Condition iii.2), the event $\{\phi(\boldsymbol{a}+\boldsymbol{B}_t)=0 \text{ for all } t\geq t_1, \boldsymbol{a}\in\mathcal{A}_d(c)\}$ holds. Then it is clear that $\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_{t_1},t_1)=\phi(\hat{\boldsymbol{y}}_{t_1},t_1)=0$ by definition of $\hat{\boldsymbol{m}}$. Suppose the inductive hypothesis that $\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_{t_1+i\Delta},t_1+i\Delta)=0$ for all $0\leq i\leq k$. Then from Eq. (24), we get that $\hat{\boldsymbol{y}}_{t_1+(k+1)\Delta}-\boldsymbol{B}_{t_1+(k+1)\Delta}\in\mathcal{A}_d(c)$ and so $\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_{t_1+(k+1)\Delta},t_1+(k+1)\Delta)=0$. Consequently, $\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_t,t)=0$ for all $t=\ell\Delta,t\geq t_1$. By Condition iii.2), the preceding event holds with high probability, so that for each $t=\ell\Delta,t\geq t_1$, $\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_t,t)=0$ with high probability. By boundedness of $\hat{\boldsymbol{m}}$ and definition of the Wasserstein-1 distance used on the function $f(\cdot)=\|\cdot\|$, we obtain that $W_1(\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_t,t),\boldsymbol{x})\geq \alpha-o(1)$. The proof ends here.

G Proof of Corollary 3.2

G.1 AUXILIARY LEMMAS

We will use the following lemmas, whose proofs are deferred to Appendix K.

Lemma G.1. Let $W \sim \mathsf{GOE}(n, 1/2)$, and $C > \sqrt{2}$ some positive constant. Then we have

$$\mathbb{P}\left(\max_{\boldsymbol{v}\in\Omega_{n,k}}|\langle\boldsymbol{v},\boldsymbol{W}\boldsymbol{v}\rangle|\geq C\sqrt{\log\binom{n}{k}}\right)\leq 2\binom{n}{k}^{-C^2/2+2}$$

We state the following non-asymptotic result from Peng (2012).

Lemma G.2 (Theorem 3.1, Peng (2012).). Let $y = \theta u u^T + \mathsf{GOE}(n, 1/n)$, and denote by $\lambda_1(y)$ the top eigenvalue of y. Letting $\xi(\theta) := \theta + \theta^{-1}$, the following holds for every $x \ge 0$ and $\theta > 1$:

$$\mathbb{P}\left(\lambda_1(\boldsymbol{y}) \le \xi(\theta) - x - \frac{2}{n}\right) \le \exp\left(-\frac{(n-1)(\theta-1)^4}{16\theta^2}\right) + 8\exp\left(-\frac{1}{4} \cdot \frac{(n-1)(\theta-1)^5 x^2}{(\theta+1)^3}\right).$$

In Appendix K, we will use the last lemma to prove the bound below.

Lemma G.3 (Alignment bound). Let $y = \theta u u^{\mathsf{T}} + \mathsf{GOE}(n, 1/n)$, where $\theta = \sqrt{1 + \delta}$. Let v_1 be the top eigenvector of y. Then, there exist constants C, c > 0 such that the following holds for any $\delta = \delta_n$ with $n^{-c} \ll \delta \ll 1$ for some c > 0 small enough:

$$\mathbb{P}\left(|\langle \boldsymbol{v}_1, \boldsymbol{u}\rangle| \le c\delta^2\right) \le C e^{-cn^{1/3}}.$$
 (25)

We also use the following lemma, which is implied in the proof of Lemma E.4.

Lemma G.4. Let $g \sim N(0, I_n)$ and define the set

$$\mathcal{L}(\boldsymbol{g}; C) = \left\{ i : 1 \le i \le n, |g_i| \ge C\sqrt{\log(n/k)} \right\}.$$

Assume $\sqrt{n} \ll k \ll n$. Then, for any C large enough, there exists C_* such that

$$\mathbb{P}\left(|\mathcal{L}(\boldsymbol{g};C)| \ge \left(\sqrt{k} \lor \frac{k^2}{n}\right)\right) \le C_* e^{-n^{1/4}}.$$

G.2 Proof of Corollary 3.2

 Let $\delta = o_n(1)$ be a parameter to be chosen later and recall that $t_{\rm alg} = n/2$. We define $\hat{\boldsymbol{m}}$ as follows:

$$\hat{\boldsymbol{m}}(\boldsymbol{y},t) = \begin{cases} \hat{\boldsymbol{m}}_0(\boldsymbol{y},t) \mathbf{1}_{\parallel \hat{\boldsymbol{m}}_0(\boldsymbol{y},t) \parallel \leq \varepsilon} & \text{if } t \leq (1-\gamma)t_{\text{alg}}, \\ \hat{\boldsymbol{m}}_0(\boldsymbol{y},t) & \text{if } (1-\gamma)t_{\text{alg}} < t < (1+\delta)t_{\text{alg}}, \\ \hat{\boldsymbol{m}}_0(\boldsymbol{y},t)\phi_1(\boldsymbol{y},t) & \text{if } (1+\delta)t_{\text{alg}} \leq t < n^4, \\ \hat{\boldsymbol{m}}_0(\boldsymbol{y},t)\phi_2(\boldsymbol{y},t) & \text{if } t \geq n^4, \end{cases}$$

where $\phi_1, \phi_2 : \mathbb{R}^{n \times n} \times \mathbb{R}_{\geq 0} \to \{0, 1\}$ are defined below and γ is given in Assumption 1. It will be clear from the constructions below that ϕ_1, ϕ_2 can be evaluated in polynomial time. To establish the relationship between Corollary 3.2 and Theorem 1, we make a few remarks:

- Assumption 1 is used in (27);
- By defining the hypothesis test ϕ such that

$$\phi(\boldsymbol{y}, t) = \begin{cases} \phi_1(\boldsymbol{y}, t) & \text{if } (1 + \delta)t_{\text{alg}} \le t < n^4 \\ \phi_2(\boldsymbol{y}, t) & \text{if } t \ge n^4 \end{cases}$$

we recover the desired properties of Assumption 2, with $\eta_2(n) = O(n^{-D})$ for any D > 0.

In order to prove claim M2, we need to bound $J_n(0,\infty)$, whereby, for $t_a \leq t_b$,

$$J_n(t_a, t_b) := \int_{t_a}^{t_b} \mathbb{E}\left[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t, t) - \hat{\boldsymbol{m}}_0(\boldsymbol{y}_t, t)\|^2\right] dt$$

Setting $t_1 = (1 + \delta)t_{\text{alg}}$ and $t_2 = n^4$, we write

$$J_n(0,\infty) = J_n(0,t_1) + J_n(t_1,t_2) + J_n(t_2,\infty),$$
(26)

and will bound each of the three terms separately.

Bounding $J_n(0, t_1)$. By Assumption 1, we know that

$$J_n(0, t_1) \le \int_0^{(1-\gamma)t_0} \mathbb{P}(\|\hat{\boldsymbol{m}}_0(\boldsymbol{y}_t, t)\| > \varepsilon) \, \mathrm{d}t = O(n^{-D}),$$
 (27)

for every D > 0.

Bounding $J_n(t_1, t_2)$. For a matrix $\mathbf{y} \in \mathbb{R}^{n \times n}$ and time point t, we define $\phi_1(\mathbf{y}, t)$ according to the following procedure:

- 1. Compute the symmetrized matrix $\mathbf{A} = (1/2)(\mathbf{y} + \mathbf{y}^{\mathsf{T}})$;
- 2. Compute its top eigenvector v (choose at random if this is not unique).
- 3. Let $S \subseteq [n]$ be the k positions in v with the largest magnitude, and define $\hat{v} \in \mathbb{R}^n$ with $\hat{v}_i = (1/\sqrt{k})\operatorname{sign}(v_i) \mathbf{1}_{i \in S}$;
- 4. Compute the test statistic $s := \langle \hat{v}, A\hat{v} \rangle$, and return 1 if $s \geq \beta t$; 0 otherwise.

Here $\beta \in (0,1)$ is a fixed constant to be chosen later.

We will show that $\phi_1(y_t,t)=1$ with overwhelming probability for the true model $y_t=tx+B_t$. Define

$$oldsymbol{A}_t = rac{oldsymbol{y}_t + oldsymbol{y}_t^\mathsf{T}}{2} = toldsymbol{u}oldsymbol{u}^\mathsf{T} + oldsymbol{W}_t$$

where we recall that $u \sim \operatorname{Unif}(\Omega_{n,k})$ and W_t is a $\operatorname{GOE}(n)$ process. Let v_t, \hat{v}_t be the top eigenvector and thresholded vector of A_t , respectively.

We have

$$s_t := \langle \hat{\boldsymbol{v}}_t, \boldsymbol{A}_t \hat{\boldsymbol{v}}_t \rangle = t \cdot \langle \boldsymbol{u}, \hat{\boldsymbol{v}}_t \rangle^2 + \langle \hat{\boldsymbol{v}}_t, \boldsymbol{W}_t \hat{\boldsymbol{v}}_t \rangle. \tag{28}$$

Using Lemma G.1, we know that $|\langle \hat{\boldsymbol{v}}_t, \boldsymbol{W} \hat{\boldsymbol{v}}_t \rangle| \leq 4\sqrt{k \log(n/k)} \cdot \sqrt{t}$ with probability at least $1 - \binom{n}{k}^{-6}$, say. Further

$$v_t = \langle v_t, u \rangle u + \sqrt{1 - \langle v_t, u \rangle^2} w$$

where w is a uniformly random unit vector orthogonal to u, independent of $\langle v_t, u \rangle$. Alternatively, there exists $g \sim N(0, I_n)$, independent of $\langle v_t, u \rangle$ so that:

$$oldsymbol{w} = rac{(oldsymbol{I}_n - oldsymbol{u} oldsymbol{u}^\mathsf{T}) oldsymbol{g}}{\|(oldsymbol{I}_n - oldsymbol{u} oldsymbol{u}^\mathsf{T}) oldsymbol{g}\|}$$

For every $1 \le i \le n$, we have

$$|((\boldsymbol{I}_n - \boldsymbol{u}\boldsymbol{u}^\mathsf{T})\boldsymbol{g})_i| \leq |g_i| + \frac{|\langle \boldsymbol{u}, \boldsymbol{g} \rangle|}{\sqrt{k}}.$$

Since by assumption $k \log(n/k) \ll n$, with probability at least $1 - C_1 \exp(-k/2)$, $|\langle \boldsymbol{u}, \boldsymbol{g} \rangle| \leq \sqrt{k \log(n/k)}$ and $||(\boldsymbol{I}_n - \boldsymbol{u}\boldsymbol{u}^{\mathsf{T}})\boldsymbol{g}|| \geq \sqrt{n}/2$, so that

$$i \in \mathcal{L}(\boldsymbol{g}; C) \implies |w_i| \le (2C+2) \cdot \sqrt{\frac{1}{n} \log(n/k)}$$
.

Therefore, by using Lemma G.4 and $k \log(n/k) \gg n^{1/2}$, we get that with probability at least $1 - C_* \exp(-n^{1/4})$, $|\mathcal{A}(\boldsymbol{w}; C)| \leq \max\{\sqrt{k}, k^2/n\}$ for some constant C > 0, where

$$\mathcal{A}(\boldsymbol{w};C) := \left\{ i : 1 \le i \le n, |w_i| \ge C \sqrt{\frac{1}{n} \log(n/k)} \right\}.$$

By Lemma G.3, we know that with probability at least $1 - C_* \exp(-cn^{1/3})$ for some $C_*, c > 0$, $|\langle \boldsymbol{v}_{t_0(1+\delta)}, \boldsymbol{u} \rangle| \ge c\delta^2$. Since $|\langle \boldsymbol{v}_t, \boldsymbol{u} \rangle|$ is stochastically increasing with t (Fact N.2), we actually have that for any $t \ge (1+\delta)t_{\text{alg}}$, with the same probability $|\langle \boldsymbol{v}_t, \boldsymbol{u} \rangle| \ge c\delta^2$.

On the event $\mathcal{E}_{\delta} := \{ |\langle \boldsymbol{v}_t, \boldsymbol{u} \rangle| \geq c\delta^2 \}$, further suppose without loss of generality that $\langle \boldsymbol{v}_t, \boldsymbol{u} \rangle > 0$. As a consequence of the result above, if $i \in \text{supp}(\boldsymbol{u})$ and $i \notin \mathcal{A}(\boldsymbol{w}; C)$, then

$$u_i > 0, i \notin \mathcal{A}(\boldsymbol{w}; C) \implies v_{t,i} \ge \frac{c\delta^2}{\sqrt{k}} - C\sqrt{\frac{1}{n}\log(n/k)},$$

$$u_i < 0, i \notin \mathcal{A}(\boldsymbol{w}; C) \implies v_{t,i} \le -\frac{c\delta^2}{\sqrt{k}} + C\sqrt{\frac{1}{n}\log(n/k)}$$

Similarly, if $i \notin \text{supp}(\boldsymbol{u})$, then

$$u_i = 0, i \notin \mathcal{A}(\boldsymbol{w}; C) \Rightarrow |v_{t,i}| \leq C \sqrt{\frac{1}{n} \log(n/k)}.$$

We next choose $\delta = \delta_n$ such that $\sqrt{(k/n)\log(n/k)} \ll \delta_n \ll 1$. Hence, we conclude that

$$i \in \operatorname{supp}(\boldsymbol{u}) \setminus \mathcal{A}(\boldsymbol{w}; C) \Rightarrow \hat{v}_{t_i} = \operatorname{sign}(\langle \boldsymbol{v}_t, \boldsymbol{u} \rangle) \cdot u_i,$$
 (29)

whence

$$|\langle \boldsymbol{u}, \hat{\boldsymbol{v}}_t \rangle| \ge 1 - \frac{2}{k} |\mathcal{A}(\boldsymbol{w}; C)| = 1 - o_n(1), \qquad (30)$$

with probability at least $1 - C_* \exp(-n^{1/4}/3)$. On this event, by Eq. (28) we obtain that

$$s_t > t \cdot (1 - o_n(1))^2 - 4 \cdot \sqrt{k \log(n/k)} \cdot \sqrt{t}$$

For $t \ge (1+\delta)t_{\text{alg}} = (1+\delta)n/2$, we this implies that, for any fixed $\beta \in (0,1)$, with probability at least $1 - C_* \exp(-n^{1/4}/3)$ (possibly adjusting the constant C_*).

Recalling that $\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) = \hat{\boldsymbol{m}}_0(\boldsymbol{y}_t,t)\phi_1(\boldsymbol{y}_t,t)$ for $t \in [t_1,t_2]$, and $\|\hat{\boldsymbol{m}}_0(\boldsymbol{y},t)\| \leq 1$, we have, for $t_1 = (1+\delta)t_{\text{alg}}$, $t_2 = n^4$ as defined above,

$$J_n(t_1, t_2) = \int_{t_1}^{t_2} \mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t, t) - \hat{\boldsymbol{m}}_0(\boldsymbol{y}_t, t)\|^2] dt \le \int_{t_1}^{t_2} \mathbb{P}(\phi_1(\boldsymbol{y}_t, t) = 0) \le C_* n^4 e^{-cn^{1/3}}.$$
(31)

Bounding $J_n(t_2, \infty)$. When $t \ge t_2$, we use a simple eigenvalue test. For a matrix y, and time point t, we define $\phi_2(y, t)$ according to the following procedure:

- 1. Compute symmetrized matrix $\mathbf{A} = (1/2)(\mathbf{y} + \mathbf{y}^{\mathsf{T}})$.
- 2. Compute top eigenvalue $\lambda_1(A)$.
- 3. Return 1 if $\lambda_1(\mathbf{A}) > t/2$, and 0 otherwise.

Under the true model $y_t = tx - B_t$, we have:

$$\lambda_1(\boldsymbol{A}_t) \ge t - \|\boldsymbol{W}_t\|_{\text{op}} \ge t - t^{2/3}$$

with error probability given by

$$\mathbb{P}\left(\|\boldsymbol{W}_t\|_{\text{op}} \ge t^{2/3}\right) \le C \exp\left(-ct^{1/3}\right),\,$$

for constants C, c > 0. Thus, we get that

$$J_n(t_2, \infty) \le 2 \int_{t_2}^{\infty} \mathbb{P}(\|\boldsymbol{W}_t\|_{\text{op}} \ge t^{2/3}) \, \mathrm{d}t \le C_* \int_{t_2}^{\infty} e^{-ct^{1/3}} \, \mathrm{d}t$$

$$\le C_{**} e^{-ct_2^{1/3}} \le C_{**} e^{-cn^{4/3}} \,. \tag{32}$$

Claim M2 of Theorem 3.2 follows from Eqs. (27), (31), (32).

We finally consider claim M3. Equation (4) yields, for every $\ell \geq 1$:

$$\hat{\mathbf{y}}_{\ell\Delta} = \Delta \sum_{i=1}^{\ell} \hat{\mathbf{m}}(\hat{\mathbf{y}}_{(i-1)\Delta}, (i-1)\Delta) + \sqrt{\Delta} \sum_{i=1}^{\ell} \mathbf{g}_{i\Delta}.$$
 (33)

We define

$$\bar{\boldsymbol{m}}_{t_1} := \Delta \sum_{i=1}^{\ell_1} \hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_{(i-1)\Delta}, (i-1)\Delta), \quad \ell_1 := \lfloor t_{\text{alg}}(1+\delta)/\Delta \rfloor. \tag{34}$$

We further define the following auxiliary process for $t \ge \ell_1 \Delta = \lfloor t_1/\Delta \rfloor \Delta$:

$$\tilde{\mathbf{y}}_t = \bar{\mathbf{m}}_{t_1} + \mathbf{B}_t \,, \tag{35}$$

where $(\boldsymbol{B}_t:\,t\geq 0)$ is a BM such that $\boldsymbol{B}_{j\Delta}=\sqrt{\Delta}\sum_{i=1}^k\boldsymbol{g}_{i\Delta}$. In particular, $\tilde{\boldsymbol{y}}_{\ell_1\Delta}=\hat{\boldsymbol{y}}_{\ell_1\Delta}$.

From triangle inequality, using Assumption 1, the condition $\|\hat{m}_0(y,t)\|_F \leq 1$, and that by construction $\|\hat{m}(y,t)\|_F \leq \varepsilon$ for $t \leq (1-\gamma)t_{\text{alg}}$, we get

$$\|\bar{\boldsymbol{m}}\|_F \leq \varepsilon (1-\gamma)t_{\text{alg}} + (\gamma+\delta)t_{\text{alg}}$$
.

We claim that (with high probability) $\phi_1(\tilde{y}_t,t)=0$ simultaneously for all $t\in[t_1,t_2]$. As a consequence, recalling the definition of \hat{m} , we obtain that $\tilde{y}_t=\hat{y}_t$ for all $t\in[t_1,t_2]$. In order to prove the claim, define, for $\ell > 1$:

$$\mathcal{T}_n := \left\{ t_\ell^+ = t_1 + \frac{\ell - 1}{n} : \ell \in \mathbb{N} \right\} \cap [t_1, t_2].$$

For every $t \in [t_1, t_2]$, consider the thresholded vector \hat{v}_t in the definition of the test ϕ_1 . We know that

$$\langle \hat{\boldsymbol{v}}_t, \tilde{\boldsymbol{y}}_t \hat{\boldsymbol{v}}_t \rangle = \langle \hat{\boldsymbol{v}}_t, \bar{\boldsymbol{m}} \hat{\boldsymbol{v}}_t \rangle + \langle \hat{\boldsymbol{v}}_t, \boldsymbol{B}_t \hat{\boldsymbol{v}}_t \rangle \le \varepsilon (1 - \gamma) t_{\text{alg}} + (\gamma + \delta) t_{\text{alg}} + \max_{\boldsymbol{v} \in \Omega_{n,k}} \langle \boldsymbol{v}, \boldsymbol{B}_t \boldsymbol{v} \rangle. \tag{36}$$

Using Lemma G.1, we get that

$$\mathbb{P}\left(\max_{\boldsymbol{v}\in\Omega_{n,k}}|\langle \boldsymbol{v},\boldsymbol{B}_t\boldsymbol{v}\rangle|\geq C\sqrt{\log\binom{n}{k}}\sqrt{t}\right)\leq 2\binom{n}{k}^{-C^2/2+2}.$$

Note that $|\mathcal{T}_n| \leq n^5$, whence

$$\mathbb{P}\left(\exists t \in \mathcal{T}_n : \max_{\boldsymbol{v} \in \Omega_{n,k}} |\langle \boldsymbol{v}, \boldsymbol{W}_{t_{\ell}} \boldsymbol{v} \rangle| \ge C \sqrt{\log \binom{n}{k}} \sqrt{t_{\ell}}\right) \le 2 \binom{n}{k}^{-C^2/2+2} n^5.$$

For $t \in [t_{\ell}, t_{\ell+1}]$, we have

$$\max_{t_\ell \leq t \leq t_{\ell+1}} \left\{ \max_{\boldsymbol{v} \in \Omega_{n,k}} |\langle \boldsymbol{v}, \boldsymbol{W}_t \boldsymbol{v} \rangle| - \max_{\boldsymbol{v} \in \Omega_{n,k}} |\langle \boldsymbol{v}, \boldsymbol{W}_{t_\ell} \boldsymbol{v} \rangle| \right\} \leq \max_{t_\ell \leq t \leq t_{\ell+1}} \|\boldsymbol{W}_t - \boldsymbol{W}_{t_\ell}\|_{\text{op}} \,.$$

Using Lemma E.3, we get that $\max_{t_{\ell} \le t \le t_{\ell+1}} \| \boldsymbol{W}_t - \boldsymbol{W}_{t_{\ell}} \|_{\text{op}} \le 16 \sqrt{(t_{\ell+1} - t_{\ell})n} = 16$ with probability at least $1 - 2 \exp(-32n)$. Taking a union bound over \mathcal{T}_n , we get that

$$\mathbb{P}\left(\exists t \in [t_1, t_2] : \max_{\boldsymbol{v} \in \Omega_{n,k}} |\langle \boldsymbol{v}, \boldsymbol{W}_t \boldsymbol{v} \rangle| \ge C \sqrt{\log \binom{n}{k}} \sqrt{t}\right) \le 2 \binom{n}{k}^{-C^2/2+2} n^5 + 2n^5 e^{-32n},$$

for possibly a different constant C > 0.

Using Eq. (36) and the last estimate, we obtain that

$$\mathbb{P}\left(\exists t \in [t_1, t_2] : \langle \hat{\boldsymbol{v}}_t, \tilde{\boldsymbol{y}}_t \hat{\boldsymbol{v}}_t \rangle \ge b_n t_{\text{alg}} + C \sqrt{\log \binom{n}{k}} \sqrt{t}\right) \le 2 \binom{n}{k}^{-C^2/2 + 2} n^5 + 2n^5 \exp(-32n),$$
(37)

where $b_n := \varepsilon(1-\gamma) + (\gamma+\delta_n)$. Since $\delta_n = o_n(1)$, we have $b_n \to (1-\gamma)\varepsilon + \gamma$. Further, by choosing γ, ε small enough, we get that $\limsup_n b_n < c/2$, say. Hence, we get that for $t \ge t_1$, and all n large enough

$$b_n t_{\rm alg} + C \sqrt{\log \binom{n}{k}} \sqrt{t} < ct/2$$

for the constant c of Assumption 2. Therefore Eq. (37) implies that $\phi_1(\tilde{y}_t, t) = 0$ simultaneously for all $t \in [t_1, t_2]$, with high probability, by choosing $\beta > c$. We conclude that, with high probability $\hat{y}_t = \tilde{y}_t$ throughout $t \in [t_1, t_2]$.

Finally, we extend the analysis to $t \in [t_2, \infty)$ by proving that, with high probability, $\phi_2(\tilde{y}_t, t) = 0$ and hence $\hat{y}_t = \tilde{y}_t$ for all $t \in [t_2, \infty)$. We use (for $A_t = (\tilde{y}_t + \tilde{y}_t^\mathsf{T})/2$, $W_t = (B_t + B_t^\mathsf{T})/2$)

$$\lambda_1(\boldsymbol{A}_t) \le \varepsilon (1 - \gamma) t_{\text{alg}} + (\delta + \gamma) t_{\text{alg}} + \lambda_1(\boldsymbol{W}_t)$$
(38)

Following exactly the argument as in the proof of Proposition D.1 (in particular, Subsection E.3.2), we get that $\lambda_1(\boldsymbol{W}_t) \leq t/3$ simultaneously for all $t \geq t_2$, with high probability. On this event, $\lambda_1(\boldsymbol{A}_t) < t/2$ with high probability (because $t/6 \gg (1-\gamma)t_{\rm alg}(1-\varepsilon) + (\gamma+\delta)t_{\rm alg}$). Hence, with high probability $\phi_2(\tilde{\boldsymbol{y}}_t,t)=0$ and hence $\hat{\boldsymbol{y}}_t=\tilde{\boldsymbol{y}}_t$ for all $t\in[t_2,\infty)$.

We thus proved that, with high probability, $\hat{\boldsymbol{y}}_t = \tilde{\boldsymbol{y}}_t$ for all $t \geq t_1 = (1+\delta)t_{\text{alg}}$, whence $\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_t,t) = 0$ as well (because $\phi_1(\tilde{\boldsymbol{y}}_t,t) = 0$ for $t \in [t_1,t_2]$ and $\phi_2(\tilde{\boldsymbol{y}}_t,t) = 0$ for $t \in [t_2,\infty)$). Claim M3 thus follows.

H Proof of Corollary D.2

H.1 PROPERTIES OF THE ESTIMATOR $\hat{\boldsymbol{m}}(\cdot)$

Proposition H.1. Assume $(\log n)^2 \ll k \ll \sqrt{n}$ and let $\hat{\boldsymbol{m}}(\cdot)$ be the estimator defined in Algorithm 1. Recall that in this case, $t_{\text{alg}} = k^2 \log(n/k^2)$. Then for any $\delta > 0$ there exists $\varepsilon > 0$ such that, letting $s = \sqrt{(1+\varepsilon)\log(n/k^2)}$, we have

$$\sup_{t \ge (1+\delta)t_{\text{alg}}} \mathbb{P}(\hat{\boldsymbol{m}}(\boldsymbol{y}_t, t) \neq \boldsymbol{x}) = O(n^{-D})$$
(39)

for any fixed D > 0.

 The proof of this proposition is a modification of the one in Cai et al. (2017), and will be presented in Appendix J. Note that Proposition H.1 directly implies the first inequality in Condition M2v of Theorem D.2, as $\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t)-\boldsymbol{x}\| \leq 2$.

By definition, when $t < t_{\text{alg}}$, the algorithm will return $\hat{m} = 0$, so we automatically have the following, which implies the second inequality in Condition M2v.

Proposition H.2. For any fixed $\delta > 0$, and $t \leq (1 - \delta)t_{alg}$, we have $\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t, t) - \boldsymbol{x}\| = 1$.

In the proof of Theorem D.2, we will also make use of the following estimates, whose proof is deferred to the Appendix M.

Lemma H.3. Let $(w_t : t \ge 0)$ be a process defined as

$$oldsymbol{w}_t = rac{1}{2} \left\{ (oldsymbol{B}_t + \sqrt{arepsilon} oldsymbol{g}_t) + (oldsymbol{B}_t + \sqrt{arepsilon} oldsymbol{g}_t)^\mathsf{T}
ight\}$$

where B, g are independent n^2 -dimensional BMs, and $0 \le \varepsilon < 1$. Then, for any $0 \le t_0 \le t_1$, and $t \ge 0$, $s \ge 1$,

$$\mathbb{P}\Big(\max_{t_0 \le t \le t_1} \|\boldsymbol{w}_t - \boldsymbol{w}_{t_0}\|_F \ge 4\sqrt{(t_1 - t_0)s} \cdot n\Big) \le 2e^{-n^2s/4},$$
(40)

$$\mathbb{P}(\|\boldsymbol{w}_t\|_F \ge 4\sqrt{ts} \cdot n) \le 2e^{-n^2s/4}. \tag{41}$$

H.2 Analysis of the diffusion process: Proof of Corollary D.2

We are left to prove that Condition M3v of Corollary D.2 holds.

For that purpose, we make the following choices about Algorithm 1:

- (C1) We select the constants in the algorithm to be $\varepsilon_n = o_n(1)$ and $s_n = \sqrt{(1 + \varepsilon_n) \log(n/k^2)}$. We will use the shorthands $s = s_n$ and $\varepsilon = \varepsilon_n$, unless there is ambiguity.
- (C2) The process $(g_t)_{t\geq 0}$ used in Algorithm 1 follows a n^2 -dimensional BM.

Note that Propositions H.1, H.2 hold under these choices, and in particular $g_t \sim \mathsf{N}(0, tI_{n \times n})$ at all times. Also the sequence of random vectors $g_{\ell\Delta}$, $\ell \in \mathbb{N}$ can be generated via $g_{\ell\Delta} = g_{(\ell-1)\Delta} + \sqrt{\Delta}\hat{g}_{\ell}$, for some i.i.d. standard normal vectors $\{\hat{g}_{\ell}\}_{\ell>0}$.

Letting $(z_t)_{t\geq 0}$ a standard BM in $\mathbb{R}^{n\times n}$, and $\hat{y}_0=0$ we can rewrite the approximate diffusion (4) as follows (for $t\in\mathbb{N}\cdot\Delta$)

$$\hat{\mathbf{y}}_{t+\Delta} = \hat{\mathbf{y}}_t + \Delta \cdot \hat{\mathbf{m}} \left(\hat{\mathbf{y}}_t, t \right) + \left(\mathbf{z}_{t+\Delta} - \mathbf{z}_t \right). \tag{42}$$

We further define

$$\boldsymbol{w}_{t} = \frac{1}{2} \left\{ (\boldsymbol{z}_{t} + \sqrt{\varepsilon} \boldsymbol{g}_{t}) + (\boldsymbol{z}_{t} + \sqrt{\varepsilon} \boldsymbol{g}_{t})^{\mathsf{T}} \right\}. \tag{43}$$

It is easy to see that $(c(\varepsilon)\boldsymbol{w}_t:t\geq 0)$ is a $\mathsf{GOE}(n)$ process for $c(\varepsilon):=((1+\varepsilon)/2)^{-1/2}$. The key technical estimate in the proof of Theorem D.2, Condition M3v is stated in the next proposition.

Proposition H.4. Let $(w_t: t \ge 0)$ be defined as per Eq. (43), and assume $\varepsilon_n = o_n(1)$ and $s_n = \sqrt{(1+\varepsilon_n)\log(n/k^2)}$. Further, assume $k \ge C(\log n)^{5/2}$ for some sufficiently large absolute constant C > 0. Then

$$\lim_{n \to \infty} \mathbb{P}\left\{ \|\eta_s(\boldsymbol{w}_t/\sqrt{t})\|_{\text{op}} \le k + \sqrt{t}/s \ \forall t \ge 1 \right\} = 1.$$
 (44)

Before proving this proposition, let us show that it implies Condition M3v of Theorem D.2. Indeed we claim that, with high probability, for all $\ell \in \mathbb{N}$, $\hat{m}(\hat{y}_{\ell\Delta}, \ell\Delta) = 0$ and $\hat{y}_{\ell\Delta} = z_{\ell\Delta}$. This is proven by induction over ℓ . Indeed, if it holds up to a certain $\ell - 1 \in \mathbb{N}$, then we have $\hat{y}_{\ell\Delta} = z_{\ell\Delta}$ by Eq. (42) whence it follows that $A_{t,+} = w_t/\sqrt{t}$, for $t = \ell\Delta$ (c.f. Algorithm 1, line 4) and therefore $\hat{m}(\hat{y}_t, t) = 0$ by Proposition H.4 (because the condition in Algorithm 1, line 6, is never passed).

We therefore have

$$\inf_{\ell \ge 0} W_1(\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_{\ell\Delta}, \ell\Delta), \boldsymbol{x}) \ge \inf_{\ell \ge 0} \mathbb{P}(\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_{\ell\Delta}, \ell\Delta) = \boldsymbol{0}) = 1 - o_n(1). \tag{45}$$

This concludes the proof of Theorem D.2. We next turn to proving Proposition H.4.

Proof of Proposition H.4. We follow a strategy analogous to the proof of the Law of Iterated Logarithm. We choose a sparse sequence of time points $\{t_\ell\}_{\ell=1}^{\infty}$, and (i) establish the statement jointly for these time points, and (ii) control deviations in between. In particular, we consider

$$t_{\ell} = \left(1 + \frac{\ell - 1}{n^3}\right)^2$$

for all $\ell \geq 1$.

We first show that simultaneously for all $\ell \geq 1$, we have $\max_{i,j} |(\boldsymbol{w}_{t_\ell})_{ij}/\sqrt{t_\ell}| \leq 8t_\ell^{1/4}\sqrt{\log n}$. We have, by sub-gaussianity of $(\boldsymbol{w}_{t_\ell})_{ij}$ and a union bound (here we account also for the case where i=j, in which there is an inflated variance), along with $\varepsilon = o_n(1)$: using the bound $2xy \geq x+y$ when $x,y \geq 1$ for $x=t_\ell^{1/2},y=\log n$. Taking a union bound once again over ℓ , we have

$$\mathbb{P}\left(\exists \ell \geq 1, 1 \leq i, j \leq n : |(\boldsymbol{w}_{t_{\ell}})_{ij}| \geq 8t_{\ell}^{3/4} \sqrt{\log n}\right) \leq n^{-6} \cdot \sum_{\ell=0}^{\infty} \exp\left(-8 \cdot \left(1 + \frac{\ell}{n^3}\right)\right)$$

We have, as the summands form a decreasing function of ℓ integer:

$$\sum_{\ell=0}^{\infty} \exp\left(-8 \cdot \left(1 + \frac{\ell}{n^3}\right)\right) \le C + \int_0^{\infty} \exp\left(-\frac{8x}{n^3}\right) dx \le Cn^3. \tag{46}$$

We thus obtain that

$$\mathbb{P}\left(\exists \ell \ge 1, 1 \le i, j \le n : |(\boldsymbol{w}_{t_{\ell}})_{ij}| \ge 8t_{\ell}^{3/4} \sqrt{\log n}\right) = O(n^{-3}). \tag{47}$$

The point of this calculation is that simultaneously for all $\ell \geq 1$, we can truncate the entries of $\eta_s(\boldsymbol{w}_{t_\ell}/\sqrt{t_\ell})$ by $8t_\ell^{1/4}\sqrt{\log n}$ without worry.

Namely, for each $\ell \geq 1$, $\vartheta_{\ell} = 8t_{\ell}^{1/4} \sqrt{\log n}$ we define $\tilde{w}_{t_{\ell}} \in \mathbb{R}^{n \times n}$ by

$$(\tilde{\boldsymbol{w}}_{t_{\ell}})_{ij} := \begin{cases} \eta_{s}(\boldsymbol{w}_{t_{\ell}}/\sqrt{t_{\ell}}) & \text{if } |\eta_{s}(\boldsymbol{w}_{t_{\ell}}/\sqrt{t_{\ell}})| \leq \vartheta_{\ell}, \\ \vartheta_{\ell} & \text{if } \eta_{s}(\boldsymbol{w}_{t_{\ell}}/\sqrt{t_{\ell}}) > \vartheta_{\ell}, \\ -\vartheta_{\ell} & \text{if } \eta_{s}(\boldsymbol{w}_{t_{\ell}}/\sqrt{t_{\ell}}) < -\vartheta_{\ell}. \end{cases}$$
(48)

By Eq. (47), we have

$$\mathbb{P}\left(\exists \ell \ge 1 : \eta_s(\boldsymbol{w}_{t_\ell}/\sqrt{t_\ell}) \ne \tilde{\boldsymbol{w}}_{t_\ell}\right) = O(n^{-3}). \tag{49}$$

We have from Bandeira & van Handel (2016), for every $x \ge 0$:

$$\mathbb{P}\left(\|\tilde{\boldsymbol{w}}_{t_{\ell}}\|_{\text{op}} \ge 4\sigma + x\right) \le n \exp\left(-\frac{cx^2}{\sigma_{\star}^2}\right),\tag{50}$$

for some absolute constant c > 0, where

$$\sigma^2 := \max_{i \le n} \sum_{j=1}^n \mathbb{E}\left[(\tilde{\boldsymbol{w}}_{t_\ell})_{ij}^2 \right] \le \sum_{j=1}^n \mathbb{E}\left[\eta_s \left(\frac{\boldsymbol{w}_{t_\ell}}{\sqrt{t_\ell}} \right)_{ij}^2 \right], \tag{51}$$

$$\sigma_{\star} := \max_{i,j \le n} |(\tilde{\boldsymbol{w}}_{t_{\ell}})_{ij}| \le 8t_{\ell}^{1/4} \sqrt{\log n}.$$
 (52)

It can be seen from an immediate Gaussian calculation that, for $i \neq j$ and $Z \sim N(0,1)$:

$$\mathbb{E}\left[\eta_s \left(\frac{\boldsymbol{w}_{t_\ell}}{\sqrt{t_\ell}}\right)_{ij}^2\right] = \int_0^\infty 4z \cdot \mathbb{P}\left(\sqrt{\frac{1+\varepsilon}{2}}Z \ge z + s\right) dz$$

$$\stackrel{(a)}{\leq} \int_0^\infty 4z \cdot \frac{1}{z+s} \cdot \exp\left(-\frac{(z+s)^2}{1+\varepsilon}\right) dz$$

$$\stackrel{(b)}{\ll} \frac{1}{s} \exp\left(-\frac{s^2}{1+\varepsilon}\right) \le \exp\left(-\frac{s^2}{1+\varepsilon}\right)$$

Here in (a) we employ the Mill's ratio bound, and (b) follows from $z + s \ge s$ and $s \to \infty$.

Proceeding analogously for the diagonal entries of $\eta_s(\boldsymbol{w}_{t_\ell}/\sqrt{t_\ell})$, we obtain that $\sigma \ll \sqrt{n} \exp(-s^2/(2(1+\varepsilon))) = k$ by definition of s.

We set $x = k/3 + \sqrt{t_\ell}/(3s)$. Since $x \gg \sigma$, we have $4\sigma + x \le (3/2)x$ if n, k are sufficiently large. Using Eq. (50) we obtain that, for some universal constants c, c', c'' > 0:

$$\mathbb{P}\left(\|\tilde{\boldsymbol{w}}_{t_{\ell}}\|_{\operatorname{op}} \geq \frac{k}{2} + \frac{\sqrt{t_{\ell}}}{2s}\right) \leq n \exp\left(-\frac{c\left(\frac{k}{3} + \frac{\sqrt{t_{\ell}}}{3s}\right)^{2}}{64t_{\ell}^{1/2}\log n}\right) \stackrel{(a)}{\leq} n \exp\left(-\frac{c'k}{s\log n} - \frac{c''t_{\ell}^{1/2}}{s^{2}\log n}\right).$$

In step (a), we simply expand the squared term in the numerator and drop the quadratic term in k. Now, taking a union bound over $\ell \geq 1$, we get that (similar to Eq. (46))

$$\mathbb{P}\left(\exists \ell \geq 1 : \|\tilde{\boldsymbol{w}}_{t_{\ell}}\|_{\text{op}} \geq \frac{k}{2} + \frac{\sqrt{t_{\ell}}}{2s}\right) \leq n \exp\left(-\frac{c'k}{s\log n}\right) \sum_{\ell=1}^{\infty} \exp\left(-\frac{c''t_{\ell}^{1/2}}{s^2\log n}\right) \\
\leq n \exp\left(-\frac{c'k}{s\log n}\right) \left(O(1) + \int_{0}^{\infty} \exp\left(-\frac{c''x}{s^2n^3\log n}\right) dx\right) \\
= O\left(\exp\left(-\frac{c'k}{s\log n}\right) \cdot s^2n^4\log n\right) \\
= o_n(1),$$

where the last estimate holds if $k \ge C(\log n)^{5/2}$ for some sufficiently large C > 0. In conclusion, using the last display and Eq. (49) we have shown that

$$\mathbb{P}\left(\exists \ell \ge 1 : \left\| \eta_s \left(\frac{\boldsymbol{w}_{t_\ell}}{\sqrt{t_\ell}} \right) \right\|_{\text{op}} \ge \frac{k}{2} + \frac{\sqrt{t_\ell}}{2s} \right) = o_n(1).$$
 (53)

Now we control the in-between fluctuations. Noting that $\eta_s(\cdot)$ is a 1-Lipschitz function, we have the following crude bound:

$$\begin{aligned} \max_{t_{\ell} \leq t \leq t_{\ell+1}} \left\| \eta_s \left(\frac{\boldsymbol{w}_t}{\sqrt{t}} \right) - \eta_s \left(\frac{\boldsymbol{w}_{t_{\ell}}}{\sqrt{t_{\ell}}} \right) \right\|_{\text{op}} &\leq \max_{t_{\ell} \leq t \leq t_{\ell+1}} \left\| \frac{\boldsymbol{w}_t}{\sqrt{t}} - \frac{\boldsymbol{w}_{t_{\ell}}}{\sqrt{t_{\ell}}} \right\|_F \\ &\leq \frac{\max_{t_{\ell} \leq t \leq t_{\ell+1}} \|\boldsymbol{w}_t - \boldsymbol{w}_{t_{\ell}}\|_F}{\sqrt{t_{\ell}}} + \|\boldsymbol{w}_{t_{\ell}}\|_F \left(\frac{1}{\sqrt{t_{\ell}}} - \frac{1}{\sqrt{t_{\ell+1}}} \right) . \end{aligned}$$

From Lemma H.3, we obtain that

$$\mathbb{P}\left(\max_{t_{\ell} \leq t \leq t_{\ell+1}} \left\| \eta_s(\boldsymbol{w}_t/\sqrt{t}) - \eta_s(\boldsymbol{w}_{t_{\ell}}/\sqrt{t_{\ell}}) \right\|_{\text{op}} \geq 4n \cdot \sqrt{t_{\ell+1} - t_{\ell}} + 4n \cdot \sqrt{t_{\ell}} \cdot \left(1 - \sqrt{\frac{t_{\ell}}{t_{\ell+1}}}\right)\right) \leq 4e^{-n^2 t_{\ell}/4}.$$

By definition of t_{ℓ} , simple algebra reveals that (we also use the fact that $n^{-1/2} \ll s^{-1}$):

$$4n \cdot \sqrt{t_{\ell+1} - t_{\ell}} + 4n \cdot \sqrt{t_{\ell}} \cdot \left(1 - \sqrt{\frac{t_{\ell}}{t_{\ell+1}}}\right) \le \frac{\sqrt{t_{\ell}}}{2s}$$
.

By union bound over $\ell \geq 1$,

$$\mathbb{P}\left(\exists \ell \geq 1 : \max_{t_{\ell} \leq t \leq t_{\ell+1}} \left\| \eta_s(\boldsymbol{w}_t/\sqrt{t}) - \eta_s(\boldsymbol{w}_{t_{\ell}}/\sqrt{t_{\ell}}) \right\|_{\text{op}} \geq \frac{k}{2} + \frac{\sqrt{t_{\ell}}}{2s} \right) \\
\leq 4 \sum_{\ell=1}^{\infty} \exp\left(-\frac{n^2}{8} - \frac{t_{\ell}}{8}\right) \\
= 4 \exp(-n^2/8) \sum_{\ell=1}^{\infty} \exp\left(-\frac{1}{8} \left(1 + \frac{\ell - 1}{n^3}\right)^2\right)$$

$$\leq 4\exp(-n^2/8)\left(O(1) + \int_0^\infty \exp\left(-\frac{x^2}{8n^6}\right)dx\right) = O(\exp(-n^2/8)n^3) = o(1).$$

Using this estimate together with Eq. (53), we conclude that with high probability the following holds simultaneously for all $t \ge 1$. Letting ℓ be largest such that $t_{\ell} \le t$:

$$\left\| \eta_s(\boldsymbol{w}_t/\sqrt{t}) \right\|_{\text{op}} \leq \left\| \eta_s(\boldsymbol{w}_{t_\ell}/\sqrt{t_\ell}) \right\|_{\text{op}} + \left\| \eta_s(\boldsymbol{w}_t/\sqrt{t}) - \eta_s(\boldsymbol{w}_{t_\ell}/\sqrt{t_\ell}) \right\|_{\text{op}} \leq k + \frac{\sqrt{t_\ell}}{s} \leq k + \frac{\sqrt{t}}{s},$$

and this finishes the proof.

We remark that Assumption (C2) in the proof above is technically not needed, meaning that the additional noise stream g_t can in fact be discarded entirely: an appropriate thresholding of v_t , the top eigenvector of $\eta_s(A_{t,+})$, as in Algorithm 2, will also suffice to satisfy all conditions of Theorem D.2, although x will not be recovered exactly; some o(k) positions outside the support of x will also be chosen, at most. The reason for this is that the alignment $|\langle v_t, u \rangle| = 1 - o_n(1)$ already, from a close inspection of our proof of Proposition H.1. Regarding the proof of Proposition H.4 above, one can easily realize that even if $\varepsilon = 0$, it will go through without any modification. We choose to keep our formulation of Algorithm 1 as faithful to the original work of Cai et al. (2017) as possible to discuss a variety of approaches, and leave this to the interested reader.

I Proof of Proposition 2.1

We take the first row of \hat{z}_1 , and let $A = \{z_{11}, \cdots, z_{1n}\}$. Let $r_j = \operatorname{rank}(z_{1j})$ denote the rank of z_{1j} with respect to the elements of A. Then since $z_{1j} \sim \mathsf{N}(0,1)$ across j, the collection of the first k ranks $A_k = \{r_1, \cdots, r_k\}$ constitutes a sample without replacement from [n]. Construct v a binary vector such that $v_i = 1$ if and only if $i \in A_k$, and let u be a randomized-sign vector version of $(1/\sqrt{k})v$. Let

$$\hat{\boldsymbol{m}}(\boldsymbol{y},t;\boldsymbol{g}_1) = \boldsymbol{u}\boldsymbol{u}^\mathsf{T} = \boldsymbol{x}' \tag{54}$$

then it is clear that $\hat{m}(y,t;\hat{z}_1) \sim x$ and is independent of x (as it is a function of only \hat{z}_1). The identity from (i) follows accordingly. To see that this error is clearly sub-optimal compared to polynomial time algorithms, observe that $\hat{m}=0$ is a polynomial time drift, which achieves error 1 at every t.

Point (ii) also follows immediately. Indeed, for every $\ell > 0$,

$$W_1(\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_{\ell\Delta}, \ell\Delta), \boldsymbol{x}) = W_1(\boldsymbol{x}', \boldsymbol{x}) = 0$$

Lastly, regarding point (iii), note that since x' is not dependent on t, we have, for every $\ell \geq 1$,

$$\hat{m{y}}_{\ell\Delta} = \hat{m{y}}_{(\ell-1)\Delta} + \Deltam{x}' + \sqrt{\Delta}\hat{m{z}}_{\ell\Delta}$$

Simple induction gives

$$\hat{m{y}}_{\ell\Delta} = (\ell\Delta)m{x}' + \sqrt{\Delta}\sum_{j=1}^\ell \hat{m{z}}_{j\Delta}$$

We take the coupling of $(\hat{y}_{\ell\Delta}/(\ell\Delta), x)$ such that x' = x. Then by definition of the Wasserstein-2 metric,

$$W_2(\hat{m{y}}_{\ell\Delta}/(\ell\Delta), m{x})^2 \leq \mathbb{E}\left[\left\|rac{\sqrt{\Delta}\sum_{j=1}^\ell \hat{m{z}}_{j\Delta}}{\ell\Delta}
ight\|^2
ight] = rac{n^2}{\ell\Delta}$$

It is clear that as $\ell \to \infty$, this quantity converges to 0. Hence we are done.

J Proof of Proposition H.1

We conduct our analysis conditional on u (recall that $x = uu^{\mathsf{T}}$), and S be the support of u. Let $A_0 = \sqrt{t}uu^{\mathsf{T}}$ and notice that

$$A_{t,+} = A_0 + \sigma_+ Z, \quad A_{t,-} = A_0 + \sigma_- W,$$
 (55)

where $\sigma_+^2 := (1+\varepsilon)/2$, $\sigma_-^2 := (1+\varepsilon)/(2\varepsilon)$, and $\boldsymbol{Z}, \boldsymbol{W} \sim \mathsf{GOE}(n)$ are independent random matrices.

We have

$$\eta_s(\mathbf{A}_{t,+}) = \mathbf{A}_0 + \eta_s(\sigma_+ \mathbf{Z}) + \mathbb{E}[\mathbf{B}] + (\mathbf{B} - \mathbb{E}[\mathbf{B}]), \tag{56}$$

where

$$B_{ij} = \eta_s \left(\sqrt{t} \cdot u_i u_j + \sigma_+ Z_{ij} \right) - \sqrt{t} \cdot u_i u_j - \eta_s (\sigma_+ Z_{ij}). \tag{57}$$

Our first order of business is to analyze $\mathbb{E}[B]$. If $i \notin S$ or $j \notin S$, we have $\mathbb{E}[B_{ij}] = 0$. On the other hand, if $i, j \in S$, then

(i) Case 1: $u_i u_j = 1/k$

In this case, we have $\mathbb{E}[B_{ij}] = -b_0 + b_1 \mathbf{1}_{i=j}$ where (below $G \sim \mathsf{N}(0,1)$)

$$b_0 := -\mathbb{E}\Big\{\eta_s\Big(\frac{\sqrt{t}}{k} + \sigma_+ G\Big) - \frac{\sqrt{t}}{k}\Big\}, \quad b_1 := \mathbb{E}\Big\{\eta_s\Big(\frac{\sqrt{t}}{k} + \sqrt{2}\sigma_+ G\Big) - \eta_s\Big(\frac{\sqrt{t}}{k} + \sigma_+ G\Big)\Big\}. \tag{58}$$

Recalling that σ_+ is bounded and bounded away from 0 (without loss of generality we can assume $\varepsilon < 1/2$) and $s, \sqrt{t}/k - s$ grows with n, k, so that $\eta_s(\sqrt{t}/k + Z_{ij}) = \sqrt{t}/k + Z_{ij} - s$ with high probability; hence $|B_{ij} + s| = o_P(s)$ (as $Z_{ij} = o(s)$ with high probability). Noting that $|B_{ij}| \le 2s$, we get $b_0 = s(1+o(1))$ and $b_1 = o(s)$ (distribution on diagonal is different).

- (ii) Case 2: $u_i u_j = -1/k \Rightarrow i \neq j$ By a similar reasoning, we have $\mathbb{E}[B_{ij}] = b_0 = s(1 + o(1))$.
- We can thus rewrite

$$\eta_s(\mathbf{A}_{t,+}) = (\sqrt{t} - kb_0) \cdot \mathbf{u}\mathbf{u}^\mathsf{T} + b_1 \cdot \mathbf{P}_S + \eta_s(\sigma_+ \mathbf{Z}) + (\mathbf{B} - \mathbb{E}[\mathbf{B}]), \tag{59}$$

where $(\mathbf{P}_S)_{ij} = 1$ if $i = j \in S$ and = 0 otherwise.

Next, we analyze the operator norm of $\eta_s(\sigma_+ \mathbf{Z})$. Let $\tilde{\mathbf{Z}} = (m\mathbf{Z}_{ij})_{i,j \le n}$ be defined as

$$mZ_{ij} = \eta_s(\sigma_+ Z_{ij}) \mathbf{1}(|\eta_s(\sigma_+ Z_{ij})| \le C \log n).$$

$$(60)$$

for some constant C>0; we have $\max_{ij\leq n}|Z_{ij}|\leq C\log n$ with error probability at most $\exp(-c(\log n)^2)\ll n^{-D}$ for any fixed D>0. We have $\tilde{\boldsymbol{Z}}=\eta_s(\boldsymbol{Z})$. By Bandeira & van Handel (2016), there exists an absolute constant c>0 such that for every u>0:

$$\mathbb{P}\left(\|\tilde{\boldsymbol{Z}}\|_{\text{op}} \ge 4\sigma + u\right) \le n \exp\left(-\frac{cu^2}{L^2}\right),\tag{61}$$

where

$$\sigma^2 = \max_{i \le n} \sum_{j=1}^n \mathbb{E}[\eta_s(Z_{ij})^2], \qquad (62)$$

$$L = \max_{i,j \le n} \|\boldsymbol{m}\boldsymbol{Z}_{ij}\|_{\infty} \le C \log n.$$
 (63)

An immediate Gaussian calculation yields, for $i \neq j$:

$$\mathbb{E}[\eta_s(\sigma_+ Z_{ij})^2] = \int_0^\infty 4z \cdot \mathbb{P}(\sigma_+ Z_{ij} \ge z + s) \mathrm{d}z \le C_1 e^{-s^2/(1+\varepsilon)}. \tag{64}$$

for some constant $C_1 > 0$.

Proceeding analogously for $\eta_s(\sigma_+ Z_{ii})$ and substituting in Eq. (61), we get $\sigma^2 \leq 2C_1 n \exp\{-s^2/(1+\varepsilon)\}$. Applying Eq. (61) there exists an absolute constant C, C' > 0 such that, by taking u = C'k, with probability at least $1 - \exp(-ck^2/(\log n)^2)$,

$$\|\eta_s(\sigma_+ \mathbf{Z})\|_{\text{op}} \le C\left(\sqrt{n}\exp\{-s^2/2(1+\varepsilon)\} \vee \sqrt{\log n}\right)$$
(65)

Note that the error probability is at most $\exp(-ck)$, because we already know that $k \gg (\log n)^2$.

Lastly, consider $B - \mathbb{E}[B]$. By Eq. (57) we know that the entries of this matrix are independent with mean 0 and bounded by 2s, hence subgaussian. Further only a $k \times k$ submatrix is nonzero, so that

$$\|\boldsymbol{B} - \mathbb{E}[\boldsymbol{B}]\|_{\text{op}} \le C_1 \sqrt{k} s \,, \tag{67}$$

with high probability (for instance, the operator norm tail bound above can be applied once more, which gives an error probability of at most $C \exp(-ck)$ for some absolute constant C, c > 0).

Summarizing, we proved that

$$\eta_s(\mathbf{A}_{t,+}) = (\sqrt{t} - kb_0) \cdot \mathbf{u}\mathbf{u}^\mathsf{T} + \mathbf{\Delta}, \tag{68}$$

$$\|\Delta\|_{\text{op}} \le C(k + \sqrt{k}s) \le C'k\,,\tag{69}$$

where in the last step we used $k \gg (\log n)^2$

Recall that v_t denotes the top eigenvector of $\eta_s(A_{t,+})$. By Davis-Kahan,

$$\min_{a \in \{+1, -1\}} \| \boldsymbol{v}_t - a \boldsymbol{u} \| \le \frac{Ck}{\sqrt{t} - kb_0}$$
(70)

$$\stackrel{(a)}{\leq} \frac{Ck}{\sqrt{t - (1 + \varepsilon)ks}} \tag{71}$$

$$\stackrel{(b)}{\leq} \varepsilon, \tag{72}$$

where in (a) we used the fact that $b_0 = s + o(s)$ and in (b) the fact that $t \ge (1 + \delta)k^2 \log(n/k^2)$, whereby we can assume $\delta \ge C\varepsilon$ for C a sufficiently large absolute constant. Recalling the definition of the score \hat{v}_t :

$$\hat{\boldsymbol{v}}_t = \boldsymbol{A}_{t,-} \boldsymbol{v}_t = \sqrt{t} \boldsymbol{u} \langle \boldsymbol{u}, \boldsymbol{v}_t \rangle + \sigma_- \boldsymbol{W} \boldsymbol{v}_t$$

where we know that $G = Wv_t \sim N(0, I_n)$ by independence of W and v_t . Assuming to be definite that the sign of the eigenvector is chosen so that the last bound holds with a = +1, we get that $\langle u, v_t \rangle \geq 1 - \varepsilon^2$. We get that for every $j \in S$:

$$u_j > 0 \Rightarrow \hat{v}_{t,j} \ge (1 - \varepsilon^2) \sqrt{\frac{t}{k}} - \sigma_- |G_j| \ge (1 - \varepsilon^2) \sqrt{\frac{t}{k}} - \frac{C}{\varepsilon} \log n$$
 (73)

$$u_j < 0 \Rightarrow \hat{v}_{t,j} \le -(1 - \varepsilon^2)\sqrt{\frac{t}{k}} + \sigma_-|G_j| \le -(1 - \varepsilon^2)\sqrt{\frac{t}{k}} + \frac{C}{\varepsilon}\log n$$
 (74)

where we use a union bound to get $|G_j| \le C \log n$ for all $j \le n$, with probability at least $1 - \exp(-c(\log n)^2)$. Similarly, for all $i \notin S$,

$$|\hat{v}_{t,j}| \le \sigma_- |G_j| \le \frac{C}{\sqrt{\varepsilon}} \log n$$

These calculations reveal that: (i) the entries with the largest magnitudes are the elements of S, and (ii) if u_i and $\hat{v}_{t,i}$ share the same sign for all $i \in S$. On this event, $\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) - \boldsymbol{x}\| = 0$.

Lastly, we claim that the top eigenvalue of $\eta_s(A_{t,+})$ is larger than $k + \sqrt{t}/s$. From triangle inequality applied to Eq. (59), we have

$$\lambda_1(\eta_s(\mathbf{A}_{t,+})) \ge (\sqrt{t} - kb_0) - b_1 - \|\eta_s(\sigma_+ \mathbf{Z})\|_{\text{op}} - \|\mathbf{B} - \mathbb{E}[\mathbf{B}]\|_{\text{op}}$$
 (75)

$$\stackrel{(a)}{\geq} (\sqrt{t} - (1+\varepsilon)ks) - Ck - C\sqrt{k}s \tag{76}$$

$$\stackrel{(b)}{\geq} \sqrt{t} - (1 + C_0 \varepsilon) ks \tag{77}$$

$$\stackrel{(c)}{>} \varepsilon ks$$
. (78)

Here (a) follows from Eqs. (66) and (67), (b) because $k \gg 1$ and $s \gg 1$ and (c) follows by taking $\delta \geq C\varepsilon$ for C a sufficiently large absolute constant. The claim follows because $ks \gg k$ and $ks \gg \sqrt{t}$, and that $\exp(-c(\log n)^2)$ is a super-polynomially small rate.

K AUXILIARY LEMMAS FOR SECTION G

K.1 Proof of Lemma G.1

We let $B \sim N(\mathbf{0}, I_{n^2})$ so that $W = (B + B^T)/2$. For $v \in \Omega_{n,k}$ we have $\langle v, Wv \rangle = \langle v, Bv \rangle \sim N(0,1)$. We thus have, by Gaussian tail bounds and a triangle inequality:

$$\mathbb{P}\left(|\langle \boldsymbol{v}, \boldsymbol{W} \boldsymbol{v} \rangle| \geq C \sqrt{\log \binom{n}{k}}\right) \leq 2 \exp\left(-\frac{C^2}{2} \log \binom{n}{k}\right) \;.$$

Taking the union bound over $v \in \Omega_{n,k}$ gives the desired statement, since the cardinality of this set is $\binom{n}{k} 2^k$.

K.2 Proof of Lemma G.3

Using Lemma G.2, we can take $x = n^{-1/4}$, say, and $\theta = \sqrt{1+\delta}$ for $\delta \ge n^{-c_0}$ for some small enough $c_0 > 0$, to get that

$$\mathbb{P}\left(\lambda_1(\boldsymbol{y}) \le \theta + 1/\theta - n^{-1/4} - 2/n\right) \le C \exp(-cn^{1/3}),$$

for some absolute constants C, c > 0.

We have the following identity, letting $W \sim \mathsf{GOE}(n, 1/n)$:

$$\langle oldsymbol{u}, oldsymbol{v}_1
angle^2 = rac{1}{ heta^2 \langle oldsymbol{u}, (\lambda_1(oldsymbol{y}) oldsymbol{I} - oldsymbol{W})^{-2} oldsymbol{u}
angle} \geq rac{1}{ heta^2 \cdot \|(\lambda_1(oldsymbol{y}) oldsymbol{I} - oldsymbol{W})^{-2}\|_{ ext{op}}} \, .$$

By standard Gaussian concentration, we know that, for any $\Delta>0$

$$\mathbb{P}(\|\boldsymbol{W}\|_{\text{op}} \ge 2 + \Delta) \le C \exp(-cn\Delta^2).$$

In this inequality, we take

$$\Delta = \frac{1}{4} (\theta + 1/\theta - n^{-1/4} - 2/n - 2).$$

Note that with $\theta = \sqrt{1+\delta}$ and $\delta = o_n(1)$, we know that $\theta + 1/\theta - 2 = \Theta(\delta^2)$, so that $\Delta = \Theta(\delta^2)$ if $\delta \ge n^{-c_0}$ with $c_0 \le 1/8$. Hence, by a union bound on the two concentration inequalities,

$$\mathbb{P}\left(\lambda_{\min}(\lambda_1(\boldsymbol{y})\boldsymbol{I} - \boldsymbol{W}\right) \le 2\Delta\right) \le C \exp(-cn^{1/3})$$

and on the complement of this event, we know that

$$\langle \boldsymbol{v}_1, \boldsymbol{u} \rangle^2 \ge \frac{4\Delta^2}{\theta^2} = \Theta(\delta^4)$$

since $\theta = \Omega(1)$, and so $|\langle \boldsymbol{v}_1, \boldsymbol{u} \rangle| = \Omega(\delta^2)$.

L Proof of Proposition E.1

In our proof, we will use the following elementary facts.

Fact L.1. For any deterministic unit vector \mathbf{u} , a unit vector \mathbf{v} is uniformly random on the orthogonal subspace to \mathbf{u} if and only if $\langle \mathbf{v}, \mathbf{u} \rangle = 0$ and $\mathbf{v} \stackrel{\text{d}}{=} \mathbf{Q} \mathbf{v}$ for every orthogonal matrix \mathbf{Q} such that $\mathbf{Q} \mathbf{u} = \mathbf{u}$.

Fact L.2. Let A be a symmetric matrix, and u a unit vector. Denote $B_{\alpha} = \alpha u u^{\mathsf{T}} + A$, and let $v(\alpha)$ be a top eigenvector of B_{α} . Then $f(\alpha) = |\langle v(\alpha), u \rangle|$ is an increasing function of $\alpha > 0$.

Let $\boldsymbol{u} \sim \text{Unif}(\Omega_{n,k})$. Recall that $\boldsymbol{A}_t = \sqrt{t}\boldsymbol{u}\boldsymbol{u}^\mathsf{T} + \sqrt{t}\boldsymbol{W}$ where $\boldsymbol{W} \sim \mathsf{GOE}(n,1/2)$.

We conduct our analysis conditional on \boldsymbol{u} . Let \boldsymbol{v}_t be a top eigenvector of \boldsymbol{A}_t . For $t=(1+\delta)n/2$, $|\langle \boldsymbol{v}_t, \boldsymbol{u} \rangle| \stackrel{a.s.}{\to} \sqrt{\delta/(1+\delta)}$, so that with high probability $|\langle \boldsymbol{v}_t, \boldsymbol{u} \rangle| \geq \sqrt{\delta}/(2\sqrt{1+\delta})$. If $t \geq (1+\delta)n/2$, we can use Fact N.2 to obtain the same result. By choosing ε such that $2\varepsilon < 1$

 $\sqrt{\delta/(1+\delta)}$, we know from standard concentration of the alignment (Lemma E.6) that $|\langle \boldsymbol{v}_t, \boldsymbol{u} \rangle| \geq 2\varepsilon$ with probability at least $1 - \exp(-cn)$ for some c > 0 possibly dependent on (ε, δ) .

By rotational invariance of W, $w_t := \frac{v_t - \langle v_t, u \rangle u}{\|v_t - \langle v_t, u \rangle u\|}$ is uniformly random on the orthogonal subspace to u. hence, there exists $g \sim \mathsf{N}(\mathbf{0}, I_n)$, such that

$$oldsymbol{w}_t \sim rac{(oldsymbol{I}_n - oldsymbol{u} oldsymbol{u}^{\mathsf{T}}) oldsymbol{g}}{\|(oldsymbol{I}_n - oldsymbol{u} oldsymbol{u}^{\mathsf{T}}) oldsymbol{g}\|}$$

Since $\|(\boldsymbol{I}_n - \boldsymbol{u}\boldsymbol{u}^\mathsf{T})\boldsymbol{g}\| \sim \|\boldsymbol{g}'\|$ for some $\boldsymbol{g}' \sim \mathsf{N}(\boldsymbol{0}, \boldsymbol{I}_{n-1})$, we have, for some constant c > 0,

$$\mathbb{P}\left(\|(\boldsymbol{I}_n - \boldsymbol{u}\boldsymbol{u}^\mathsf{T})\boldsymbol{g}\| \le \frac{\sqrt{n}}{2}\right) \le \exp\left(-cn\right).$$

Further, for every $1 \le i \le n$, we know that

$$|((\boldsymbol{I}_n - \boldsymbol{u}\boldsymbol{u}^\mathsf{T})\boldsymbol{g})_i| \leq |g_i| + \|\boldsymbol{u}\|_{\infty} \cdot |\langle \boldsymbol{u}, \boldsymbol{g} \rangle| \leq |g_i| + \frac{|\langle \boldsymbol{u}, \boldsymbol{g} \rangle|}{\sqrt{k}}.$$

We next show that, with the claimed probability, only a few entries of w_t can have large magnitude. As a result, less than ℓ entries of u can be estimated incorrectly (with $\ell = 1$ if $k \ll n/\log n$).

Define $\ell = \lceil n \exp(-a_n \cdot n/k) \rceil \ge 1$ (with a_n a sequence to be chosen later) and $g_{(\ell)}^{abs}$ as the ℓ -th largest value among the $|g_i|$'s. We have

$$\mathbb{P}\left(|\langle \boldsymbol{u}, \boldsymbol{g} \rangle| \ge \sqrt{na_n}\right) \le \exp\left(-\frac{na_n}{2}\right).$$

Furthermore, from a union bound, we get that

$$\begin{split} \mathbb{P}\left(g_{(\ell)}^{\text{abs}} \geq \frac{2\sqrt{na_n}}{\sqrt{k}}\right) &\leq \binom{n}{\ell} \cdot \exp\left(-\frac{2n\ell \cdot a_n}{k}\right) \\ &\leq \left(\frac{en}{\ell}\right)^{\ell} \exp\left(-\frac{2n\ell \cdot a_n}{k}\right) \\ &= \exp\left(-\frac{2n\ell \cdot a_n}{k} + \ell \log \frac{n}{\ell} + \ell\right) \,. \end{split}$$

By definition, we know that $\ell \ge \max\{n \exp(-a_n \cdot n/k), 1\}$, so that

$$\frac{2na_n}{k} - \log \frac{n}{\ell} - 1 \ge \frac{2na_n}{k} - \min \left\{ \log n, a_n \cdot n/k \right\} - 1 \ge \frac{na_n}{2k}$$

as long as $na_n \gg k$. This means that

$$\mathbb{P}\left(g^{\mathrm{abs}}_{(\ell)} \geq \frac{2\sqrt{na_n}}{\sqrt{k}}\right) \leq \exp\left(-\frac{na_n}{2k}\right)\,.$$

Define the set

$$\mathcal{A}_n(t) := \left\{ i \le n : |w_{ti}| \ge 6\sqrt{\frac{a_n}{k}} \right\}.$$

By the bounds above we have

$$\mathbb{P}(|\mathcal{A}_n(t)| \le \ell - 1) \ge 1 - e^{-cn} - e^{-na_n/2k} - e^{-na_n/2}.$$

Suppose that $|\langle v_t, u \rangle| \ge 2\varepsilon$ also holds, and suppose without loss of generality that $\langle v_t, u \rangle \ge 2\varepsilon$. Then, we have (as long as $6\sqrt{a_n} \le (9/10)\varepsilon$)

$$i \in S, i \notin \mathcal{A}_n(t) \, u_i > 0 \Rightarrow v_{ti} \ge \frac{\langle \boldsymbol{v}_t, \boldsymbol{u} \rangle}{\sqrt{k}} - \frac{6\sqrt{a_n}}{\sqrt{k}} > \frac{\varepsilon}{\sqrt{k}} \Rightarrow i \in \hat{S}, \operatorname{sign}(v_{ti}) > 0,,$$
$$i \in S, i \notin \mathcal{A}_n(t), u_i < 0 \Rightarrow v_{ti} \le -\frac{\langle \boldsymbol{v}_t, \boldsymbol{u} \rangle}{\sqrt{k}} + \frac{6\sqrt{a_n}}{\sqrt{k}} < -\frac{\varepsilon}{\sqrt{k}} \Rightarrow i \in \hat{S}, \operatorname{sign}(v_{ti}) < 0.$$

Analogously, for $i \notin S$, $i \notin A_n(t)$, we have

$$|v_{ti}| \le \frac{6\sqrt{a_n}}{\sqrt{k}} < \frac{\varepsilon}{\sqrt{k}}.$$

and we obtain that at most $\ell-1$ positions could be mis-identified.

Next, we show that the termination condition (Line 5, Algorithm 2) does not trigger for each $t \ge n^2$ (with high probability). We write

$$oldsymbol{A}_t = rac{oldsymbol{y}_t + oldsymbol{y}_t^\mathsf{T}}{2\sqrt{t}} = \sqrt{t}oldsymbol{u}oldsymbol{u}^\mathsf{T} + \left(rac{oldsymbol{B}_t + oldsymbol{B}_t^\mathsf{T}}{2\sqrt{t}}
ight)$$

From Weyl's inequality:

$$\lambda_1(\boldsymbol{A}_t) \geq \sqrt{t} - \left\| \frac{\boldsymbol{B}_t + \boldsymbol{B}_t^\mathsf{T}}{2\sqrt{t}} \right\|_{\mathrm{op}}$$

From standard operator norm results for GOE matrices (as $(\boldsymbol{B}_t + \boldsymbol{B}_t^\mathsf{T})/\sqrt{2t} \sim \mathsf{GOE}(n)$), we know that $\|(\boldsymbol{B}_t + \boldsymbol{B}_t^\mathsf{T})/(2\sqrt{t})\|_{\mathsf{op}} \leq 2\sqrt{n}$ with probability at least $1 - \exp(-cn)$, for some c > 0. Hence $\lambda_1(\boldsymbol{A}_t) \geq \sqrt{t} - 2\sqrt{n} > \sqrt{t}/2$ as $t \geq n^2 \gg n$.

We obtain that

$$\mathbb{E}\left[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t)-\boldsymbol{x}\|^2\right] = O\left(\frac{\ell-1}{k} + \exp\left(-\frac{na_n}{k}\right)\right) = O\left(\exp(-n\varepsilon^2/64k)\right)$$

where we picked $a_n > \varepsilon^2/64$ satisfying the bounds outlined above, namely (i) $6\sqrt{a_n} \le 0.9\varepsilon$, and (ii) $na_n \gg k$. Notice that $na_n/k \gg \log(na_n/k)$ if $na_n \gg k$, and $\ell - 1 \le n \cdot \exp(-a_n \cdot n/k)$.

M PROOF OF LEMMA H.3

Proof. By the Markov Property, we know that $\max_{t_{\ell} \leq t \leq t_{\ell+1}} \| \boldsymbol{w}_t - \boldsymbol{w}_{t_{\ell}} \|_F \stackrel{d}{=} \max_{0 < t < t_{\ell+1} - t_{\ell}} \| \boldsymbol{w}_t \|_F$. By Gaussian concentration, we have

$$\mathbb{P}\left(\max_{0 \le t \le t_{\ell+1} - t_{\ell}} \|\boldsymbol{w}_t\|_F - \mathbb{E}\left[\max_{0 \le t \le t_{\ell+1} - t_{\ell}} \|\boldsymbol{w}_t\|_F\right] \ge x\right) \le 2\exp\left(-\frac{x^2}{4(t_{\ell+1} - t_{\ell})}\right)$$

This can be proven, e.g. by discretizing the interval $[0, t_{\ell+1} - t_{\ell}]$ into r equal-length intervals and employing standard Gaussian concentration on vectors (then pushing $r \to \infty$). As the argument is standard, we omit the proof for brevity.

Now we bound $\mathbb{E}\left[\max_{0 \leq t \leq t_{\ell+1} - t_{\ell}} \|\boldsymbol{w}_t\|_F\right]$. We know that $\|\boldsymbol{w}_t\|_F$ is a non-negative submartingale, so that from Doob's inequality:

$$\mathbb{E}\left[\max_{0 \le t \le t_{\ell+1} - t_{\ell}} \|\boldsymbol{w}_{t}\|_{F}^{2}\right] \le 4\mathbb{E}[\|\boldsymbol{w}_{t_{\ell+1} - t_{\ell}}\|_{F}^{2}] \le 9(t_{\ell+1} - t_{\ell})n^{2}$$

so that from Cauchy-Schwarz, $\mathbb{E}\left[\max_{0 \leq t \leq t_{\ell+1} - t_{\ell}} \|\boldsymbol{w}_t\|_F\right] \leq 3\sqrt{t_{\ell+1} - t_{\ell}}n$. Hence as $t_{\ell} \geq 1$,

$$\mathbb{P}\left(\max_{t_{\ell} \leq t \leq t_{\ell+1}} \|\boldsymbol{w}_{t} - \boldsymbol{w}_{t_{\ell}}\|_{F} \geq 4\sqrt{(t_{\ell+1} - t_{\ell})t_{\ell}} \cdot n\right) \leq 2\exp\left(-\frac{n^{2}t_{\ell}}{4}\right)$$

The second tail bound follows immediately (at least, the proof would be analogous to the preceding display). \Box

N Proof of Proposition E.1

In our proof, we will use the following facts; since they are elementary, we omit the proof.

Fact N.1. For any deterministic unit vector \mathbf{u} , a unit vector \mathbf{v} is uniformly random on the orthogonal subspace to \mathbf{u} if and only if $\langle \mathbf{v}, \mathbf{u} \rangle = 0$ and $\mathbf{v} \sim \mathbf{Q} \mathbf{v}$ for every orthogonal matrix \mathbf{Q} such that $\mathbf{Q} \mathbf{u} = \mathbf{u}$.

Fact N.2. Let A be a symmetric matrix, and u a unit vector. Denote $B_{\alpha} = \alpha u u^{\mathsf{T}} + A$, and let $v(\alpha)$ be a top eigenvector of B_{α} . Then $f(\alpha) = |\langle v(\alpha), u \rangle|$ is an increasing function of $\alpha > 0$.

We suppose that $t \ge (1+\delta)^2 n/2$ instead of $(1+\delta)n/2$, for notational convenience. Let u be a unit vector of random signs generated from a uniformly random set $S \subset [n]$ of size k. Since scaling does not change the eigenvectors, we instead consider the matrix

$$ilde{m{Y}}_t = rac{m{A}_t}{\sqrt{n}} = \sqrt{rac{t}{n}}m{u}m{u}^\mathsf{T} + \left(rac{m{B}_t + m{B}_t^\mathsf{T}}{2\sqrt{tn}}
ight) = \sqrt{rac{t}{n}}m{u}m{u}^\mathsf{T} + m{W}_n$$

where it is clear that $\sqrt{2}W_n \sim \mathsf{GOE}(n)$.

We conduct our analysis conditional on \boldsymbol{u} . Let \boldsymbol{v}_t be a top eigenvector of $\tilde{\boldsymbol{Y}}_t$. We know that when $t=(1+\delta)n/2$, $|\langle \boldsymbol{v}_t, \boldsymbol{u}\rangle| \overset{a.s.}{\to} \sqrt{\frac{\delta}{1+\delta}}$, so that with high probability $|\langle \boldsymbol{v}_t, \boldsymbol{u}\rangle| \geq \frac{\sqrt{\delta}}{2\sqrt{1+\delta}}$. If $t\geq (1+\delta)n/2$, we can use Fact N.2 to obtain the same result. By choosing ε such that $2\varepsilon < \sqrt{\delta/(1+\delta)}$, we know from standard concentration of the alignment (Lemma E.6) that $|\langle \boldsymbol{v}_t, \boldsymbol{u}\rangle| \geq 2\varepsilon$ with probability at least $1-\exp(-cn)$ for some c>0 possibly dependent on (ε,δ) .

By rotational invariance of W_n , we know that $v_t \sim Qv_t$ for any orthogonal matrix Q such that Qu = u. We obtain that $v_t - \langle v_t, u \rangle u$ also has this property, so by Fact N.1, we get that $w_t = \frac{v_t - \langle v_t, u \rangle u}{\|v_t - \langle v_t, u \rangle u\|}$ is uniformly random on the orthogonal subspace to u. We can write, with $q \sim \mathsf{N}(0, I_n)$:

$$oldsymbol{w}_t \sim rac{(oldsymbol{I}_n - oldsymbol{u} oldsymbol{u}^{\mathsf{T}}) oldsymbol{g}}{\|(oldsymbol{I}_n - oldsymbol{u} oldsymbol{u}^{\mathsf{T}}) oldsymbol{g}\|}$$

We first deal with the denominator. From triangle inequality, we know that $\|(\boldsymbol{I}_n - \boldsymbol{u}\boldsymbol{u}^\mathsf{T})\boldsymbol{g}\| \ge \|\boldsymbol{g}\| - \|\boldsymbol{u}\boldsymbol{u}^\mathsf{T}\boldsymbol{g}\| = \|\boldsymbol{g}\| - |\langle \boldsymbol{u}, \boldsymbol{g} \rangle|$. Since $\langle \boldsymbol{u}, \boldsymbol{g} \rangle \sim \mathsf{N}(0, 1)$, we have from standard sub-exponential concentration on $\|\boldsymbol{g}\|$ that

$$\mathbb{P}\left(\|(\boldsymbol{I}_n - \boldsymbol{u}\boldsymbol{u}^\mathsf{T})\boldsymbol{g}\| \le \frac{\sqrt{n}}{2}\right) \le \exp\left(-cn\right)$$

for some constant c > 0. For every $1 \le i \le n$, we know that

$$|((\boldsymbol{I}_n - \boldsymbol{u}\boldsymbol{u}^\mathsf{T})\boldsymbol{g})_i| \leq |g_i| + \|\boldsymbol{u}\|_{\infty} \cdot |\langle \boldsymbol{u}, \boldsymbol{g} \rangle| \leq |g_i| + \frac{|\langle \boldsymbol{u}, \boldsymbol{g} \rangle|}{\sqrt{k}}$$

Define $\ell = \lceil n \exp(-a_n \cdot n/k) \rceil \ge 1$ and $g_{(\ell)}^{abs}$ as the ℓ -th largest value among the $|g_i|$'s. We have

$$\mathbb{P}\left(\left|\left\langle \boldsymbol{u},\boldsymbol{g}\right\rangle\right| \geq \sqrt{na_n}\right) \leq \exp\left(-\frac{na_n}{2}\right)$$

Furthermore, from a union bound, we get that

$$\mathbb{P}\left(g_{(\ell)}^{\mathsf{abs}} \geq \frac{2\sqrt{na_n}}{\sqrt{k}}\right) \leq \binom{n}{\ell} \cdot \exp\left(-\frac{2n\ell \cdot a_n}{k}\right) \leq \left(\frac{en}{\ell}\right)^\ell \exp\left(-\frac{2n\ell \cdot a_n}{k}\right) = \exp\left(-\frac{2n\ell \cdot a_n}{k} + \ell \log \frac{n}{\ell} + \ell\right)$$

By definition, we know that $\ell \ge \max\{n \exp(-a_n \cdot n/k), 1\}$, so that

$$\frac{2na_n}{k} - \log \frac{n}{\ell} - 1 \ge \frac{2na_n}{k} - \min \left\{ \log n, a_n \cdot n/k \right\} - 1 \ge \frac{na_n}{k}$$

as long as $na_n \gg k$. With probability at least $1 - \exp(-cn) - \exp(-na_n/k) - \exp(-na_n/2)$, we have that at most $\ell - 1$ positions i in a "bad" set A have $|w_{ti}| \geq 6\sqrt{a_n/k}$. Suppose that $|\langle \boldsymbol{v}_t, \boldsymbol{u} \rangle| \geq 2\varepsilon$ also holds, and suppose without loss of generality that $\langle \boldsymbol{v}_t, \boldsymbol{u} \rangle \geq 2\varepsilon$. Then, we have

$$i \in S, i \notin A, u_i > 0 \Rightarrow v_{ti} \ge \frac{\langle \boldsymbol{v}_t, \boldsymbol{u} \rangle}{\sqrt{k}} - \frac{6\sqrt{a_n}}{\sqrt{k}} > \frac{\varepsilon}{\sqrt{k}} \Rightarrow i \in \hat{S}, \operatorname{sign}(v_{ti}) > 0$$

$$i \in S, i \notin A, u_i < 0 \Rightarrow v_{ti} \le -\frac{\langle v_t, u \rangle}{\sqrt{k}} + \frac{6\sqrt{a_n}}{\sqrt{k}} < -\frac{\varepsilon}{\sqrt{k}} \Rightarrow i \in \hat{S}, \operatorname{sign}(v_{ti}) < 0$$

as long as $6\sqrt{a_n} < 3\varepsilon/4$, say. Analogously, for $i \notin S, i \notin A$, we have

$$|v_{ti}| \le \frac{6\sqrt{a_n}}{\sqrt{k}} < \frac{\varepsilon}{\sqrt{k}}$$

and we obtain that at most $\ell-1$ positions could be mis-identified.

For completeness, we show that the termination condition (Line 5, Algorithm 2) does not trigger for each $t \ge n^2$ (with high probability). We write

$$oldsymbol{A}_t = rac{oldsymbol{y}_t + oldsymbol{y}_t^\mathsf{T}}{2\sqrt{t}} = \sqrt{t}oldsymbol{u}oldsymbol{u}^\mathsf{T} + \left(rac{oldsymbol{B}_t + oldsymbol{B}_t^\mathsf{T}}{2\sqrt{t}}
ight)$$

From Weyl's inequality:

$$\lambda_1(oldsymbol{A}_t) \geq \sqrt{t} - \left\| rac{oldsymbol{B}_t + oldsymbol{B}_t^\mathsf{T}}{2\sqrt{t}}
ight\|_{\mathrm{op}}$$

From standard operator norm results for GOE matrices (as $(B_t + B_t^{\mathsf{T}})/\sqrt{2t} \sim \mathsf{GOE}(n)$), we know that $\|(B_t + B_t^{\mathsf{T}})/(2\sqrt{t})\|_{\mathsf{op}} \leq 2\sqrt{n}$ with probability at least $1 - \exp(-cn)$, for some c > 0. Hence $\lambda_1(A_t) \geq \sqrt{t} - 2\sqrt{n} > \sqrt{t}/2$ as $t \geq n^2 \gg n$.

We obtain that

$$\mathbb{E}\left[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t)-\boldsymbol{x}\|^2\right] = O\left(\frac{\ell-1}{k} + \exp\left(-\frac{na_n}{k}\right)\right) = O\left(\exp(-a_n \cdot n/k)\right)$$

where we notice that $na_n/k \gg \log(na_n/k)$ if $na_n \gg k$.

O PROOF OF LEMMA E.3

By Gaussian concentration, we have

$$\mathbb{P}\left(\max_{0 \leq t \leq t_{\ell+1} - t_{\ell}} \|\boldsymbol{W}_t\|_{\text{op}} - \mathbb{E}\left[\max_{0 \leq t \leq t_{\ell+1} - t_{\ell}} \|\boldsymbol{W}_t\|_{\text{op}}\right] \geq x\right) \leq 2\exp\left(-\frac{x^2}{2(t_{\ell+1} - t_{\ell})}\right)$$

This can be proven, e.g. by discretizing the interval $[0, t_{\ell+1} - t_{\ell}]$ into r equal-length intervals and employing Gaussian concentration on vectors (then pushing $r \to \infty$). As the argument is standard, we omit the proof for brevity.

To evaluate $\mathbb{E}[\max_{0 \le t \le t_{\ell+1} - t_{\ell}} \| \boldsymbol{W}_t \|_{\text{op}}]$, we recognize that $\| \boldsymbol{W}_t \|_{\text{op}}$ is a submartingale, so that from Doob's inequality:

$$\mathbb{E}\left[\max_{0 \leq t \leq t_{\ell+1} - t_{\ell}} \|\boldsymbol{W}_{t}\|_{\text{op}}^{2}\right] \leq 4\mathbb{E}\left[\|\boldsymbol{W}_{t_{\ell+1} - t_{\ell}}\|_{\text{op}}^{2}\right]$$

Once again from Gaussian concentration,

$$\mathbb{P}\left(|\|\boldsymbol{W}_1\|_{\text{op}} - \mathbb{E}[\|\boldsymbol{W}_1\|_{\text{op}}]| \ge x\right) \le 2\exp\left(\frac{-x^2}{2}\right)$$

so that $\mathbb{P}\left(\|\boldsymbol{W}_{t_{\ell+1}-t_{\ell}}\|_{\text{op}} - \mathbb{E}[\|\boldsymbol{W}_{t_{\ell+1}-t_{\ell}}\|_{\text{op}}]| \geq x\right) \leq 2\exp\left(-x^2/(2(t_{\ell+1}-t_{\ell}))\right)$. Hence $\|\boldsymbol{W}_{t_{\ell+1}-t_{\ell}}\|_{\text{op}}$ is $(t_{\ell+1}-t_{\ell})$ -subgaussian, implying that $\operatorname{Var}(\|\boldsymbol{W}_{t_{\ell+1}-t_{\ell}}\|_{\text{op}}) \leq 6(t_{\ell+1}-t_{\ell})$. As $\mathbb{E}[\|\boldsymbol{W}_{t_{\ell+1}-t_{\ell}}\|_{\text{op}}]^2 \sim 4(t_{\ell+1}-t_{\ell})n$ (one can obtain this from the Bai-Yin Theorem along with sub-gaussianity, for instance), we get that $\mathbb{E}\left[\|\boldsymbol{W}_{t_{\ell+1}-t_{\ell}}\|^2\right] \leq 16(t_{\ell+1}-t_{\ell})n$ eventually as n gets large.

From Cauchy-Schwarz inequality, we get that

$$\mathbb{E}\left[\max_{0 \le t \le t_{\ell+1} - t_{\ell}} \|\boldsymbol{W}_t\|_{\text{op}}\right] \le 8\sqrt{t_{\ell+1} - t_{\ell}}\sqrt{n}$$

We conclude that

$$\mathbb{P}\left(\max_{0 \le t \le t_{\ell+1} - t_{\ell}} \|\boldsymbol{W}_t\|_{\text{op}} \ge 16\sqrt{(t_{\ell+1} - t_{\ell})n}\right) \le 2\exp\left(-32n\right)$$

P PROOF OF LEMMA E.4

First, by orthogonal invariance of W_t , we know that v_t is uniformly random over the unit sphere \mathbb{S}^{n-1} . We can write, using $g \sim \mathsf{N}(0, I_n)$, the following representation

$$oldsymbol{v}_t \sim rac{oldsymbol{g}}{\|oldsymbol{g}\|}$$

As in the statement of the Lemma, we define the following set, for $v \in \mathbb{R}^n$ and C > 0:

$$A(\boldsymbol{v};C) = \left\{ i : 1 \le i \le n, |v_i| \ge \frac{C\sqrt{\log(n/k)}}{\sqrt{n}} \right\}$$

As with the proof of Proposition E.1, we first deal with the denominator $\|g\|$: indeed, sub-exponential concentration gives us

$$\mathbb{P}\left(\sum_{j=1}^{n} g_j^2 \le \frac{n}{2}\right) \le 2\exp(-n/8) \tag{79}$$

This leads us to define another set

$$B(\boldsymbol{g}; C) = \left\{ i : 1 \le i \le n, |g_i| \ge C\sqrt{\log(n/k)} \right\}$$

Let $p_n = \mathbb{P}(|g_1| \geq C\sqrt{\log(n/k)})$, then we have $|B(g;C)| \sim \text{Bin}(n,p_n)$. From Gaussian tail bounds, we know that $p_n \leq (n/k)^{-C^2/2}$. We now use a Chernoff bound of the following form: for every $x \geq 4\mathbb{E}[X]$, where $X \sim \text{Bin}(n,p)$, then

$$\mathbb{P}\left(X \ge x\right) \le \exp\left(-x/3\right)$$

It is clear that $np_n \ll k^2/n \le \max\{k^2/n, \sqrt{k}\}$ when C > 2, so that we have

$$\mathbb{P}\left(|B(\boldsymbol{g};C)| \geq \max\{\sqrt{k},k^2/n\}\right) \leq \exp\left(-\frac{1}{3}\max\{\sqrt{k},k^2/n\}\right) \leq \exp\left(-\frac{1}{3}n^{1/4}\right)$$

Therefore, with each fixed t, by union bound with probability at least $1 - O(\exp(-\sqrt{n}))$, we have, for a possibly different C > 0, $|A(v_t; C)| \le \max\{\sqrt{k}, k^2/n\}$. Our proof ends here, as $\max\{\sqrt{k}, k^2/n\} \ll k/2$ for $\sqrt{n} \ll k \ll n$.

O Proof of Lemma E.5

We know that

$$\begin{aligned} \boldsymbol{v}_t^\mathsf{T} \boldsymbol{W}_{t\ell} \boldsymbol{v}_t = & \boldsymbol{v}_t^\mathsf{T} \boldsymbol{W}_t \boldsymbol{v}_t - \boldsymbol{v}_t^\mathsf{T} (\boldsymbol{W}_t - \boldsymbol{W}_{t\ell}) \boldsymbol{v}_t \\ = & \lambda_1 (\boldsymbol{W}_t) - \boldsymbol{v}_t^\mathsf{T} (\boldsymbol{W}_t - \boldsymbol{W}_{t\ell}) \boldsymbol{v}_t \\ = & \lambda_1 (\boldsymbol{W}_{t\ell}) - \boldsymbol{v}_t^\mathsf{T} (\boldsymbol{W}_t - \boldsymbol{W}_{t\ell}) \boldsymbol{v}_t + (\lambda_1 (\boldsymbol{W}_t) - \lambda_1 (\boldsymbol{W}_{t\ell})) \end{aligned}$$

from which we obtain from Weyl's inequality that

$$\sup_{t_{\ell} \le t \le t_{\ell+1}} \left| \boldsymbol{v}_t^\mathsf{T} \boldsymbol{W}_{t_{\ell}} \boldsymbol{v}_t - \lambda_1(\boldsymbol{W}_{t_{\ell}}) \right| \le 2 \sup_{t_{\ell} \le t \le t_{\ell+1}} \| \boldsymbol{W}_t - \boldsymbol{W}_{t_{\ell}} \|_{\mathsf{op}} \le 32 \sqrt{(t_{\ell+1} - t_{\ell})n}$$

with probability at least $1 - 2\exp(-32n)$.

R PROOF OF LEMMA E.6

By Weyl's inequality, $W \mapsto \lambda_1(Y)$ (with $Y = \theta v v^{\mathsf{T}} + W$) is a 1-Lipschitz function and therefore, by Borell inequality (and Baik et al. (2005)), letting $\lambda_*(\theta) := \theta + 1/\theta$, for any $\varepsilon > 0$,

$$\mathbb{P}(|\lambda_1(Y) - \lambda_*(\theta)| \ge \varepsilon) \le 2e^{-n\varepsilon^2/4}. \tag{80}$$

To prove concentration of $\langle v_1(Y), v \rangle^2$, note that simple linear algebra yields

$$\frac{1}{\langle \boldsymbol{v}_1(\boldsymbol{Y}), \boldsymbol{v} \rangle^2} = \langle \boldsymbol{v}, (\lambda_1(\boldsymbol{Y})\boldsymbol{I} - \boldsymbol{W})^{-2} \boldsymbol{v} \rangle =: F(\boldsymbol{W}). \tag{81}$$

It is therefore sufficient to prove that F(W) concentrates around a value that is bounded away from 0. Fix $\varepsilon_0 > 0$ such that $2 + 3\varepsilon_0 < \lambda_*(\theta)$ and define the event

$$\mathcal{E} := \left\{ \boldsymbol{W} : \|\boldsymbol{W}\|_{\text{op}} \le 2 + \varepsilon_0, \ |\lambda_1(\boldsymbol{Y}) - \lambda_*| \le \varepsilon_0 \right\}. \tag{82}$$

By the Bai-Yin law and Gaussian concentration (plus the above concentration of λ_1), $\mathbb{P}(\mathcal{E}) \geq 1 - 2e^{-c(\varepsilon_0)n}$ for some $c(\varepsilon_0) > 0$. Further, it is easy to check that $F(\boldsymbol{W})$ is Lipschitz on \mathcal{E} , whence the concentration of $\langle \boldsymbol{u}, \boldsymbol{v}_1(\boldsymbol{W}) \rangle^2$ follows by another application of Borell inequality.

S PROOFS OF REDUCTION RESULTS

S.1 Proof of Theorem 2

We state and prove a more detailed version of Theorem 2.

Theorem 4. Assume that $\hat{\boldsymbol{n}}(\cdot,\cdot)$ has complexity χ and that for any $T \leq \theta d$, $D_{\text{KL}}(\overline{P}_{\hat{\boldsymbol{y}}}^{T,\Delta} || P_{\boldsymbol{y}}^T) \leq \varepsilon$ (where $\overline{P}_{\hat{\boldsymbol{y}}}^{T,\Delta}$ is the continuous time process obtained by Brownian-linear interpolation of Eq. (4)).

Then for any $\sigma > 0$ there exists an algorithm with complexity $O(\chi \cdot T/\Delta)$, that takes as input $y = x + \sigma g$, $(x, g) \sim \mu \otimes \mathsf{N}(0, I)$, and outputs \hat{x} , such that

$$\mathbb{E}\|\mathbf{P}_{\boldsymbol{x}|\boldsymbol{y}} - \mathbf{P}_{\hat{\boldsymbol{x}}|\boldsymbol{y}}\|_{\mathsf{TV}} \le \sqrt{2\varepsilon} + \varepsilon_0(\theta) =: \overline{\varepsilon}, \tag{83}$$

where $\varepsilon_0(\theta) := \mathbb{E}\|P_{\boldsymbol{x}|\boldsymbol{y}} - \mathsf{N}(\boldsymbol{0}, (\theta d)^{-1}\boldsymbol{I}_d) * P_{\boldsymbol{x}|\boldsymbol{y}}\|_{\mathsf{TV}}$ is the expected TV distance between $P_{\boldsymbol{x}|\boldsymbol{y}}$ and the convolution of $P_{\boldsymbol{x}|\boldsymbol{y}}$ with a Gaussian with variance $1/(\theta d)$. As a consequence, there exists a randomized algorithm $\hat{\boldsymbol{m}}_+$ with complexity $(N\chi \cdot T/\Delta)$ that approximates the posterior expectation:

$$\mathbb{E}\left\{\|\hat{\boldsymbol{m}}_{+}(\boldsymbol{y}) - \boldsymbol{m}(\boldsymbol{y})\|^{2}\right\} \leq 2\overline{\varepsilon} + 2N^{-1}.$$
(84)

Proof. The algorithm consists in running the discretized diffusion (4) with initialization $\hat{y}_{t_0} = y/\sigma^2$ at $t = t_0 := 1/\sigma^2$. To avoid notational burden, we will assume $(T - t_0)/\Delta$ to be an integer. Let $\hat{y}_{t_0}^*$ be generated by the discretized diffusion with initialization at \hat{y}_0 at t = 0. Note that the distribution of \hat{y}_{t_0} is the same as the one of $t_0x + \sqrt{t}g$ and hence by Assumption (b), and Pinsker's inequality

$$\|P_{\hat{\boldsymbol{y}}_{t_0}} - P_{\hat{\boldsymbol{y}}_{t_0}^*}\|_{TV} \le \sqrt{\frac{1}{2}D_{KL}(P_{\hat{\boldsymbol{y}}_{t_0}^*}\|P_{\hat{\boldsymbol{y}}_{t_0}})} \le \sqrt{\frac{1}{2}D_{KL}(\overline{P}_{\hat{\boldsymbol{y}}}^{T,\Delta}\|P_{\hat{\boldsymbol{y}}}^T)} \le \sqrt{\frac{\varepsilon}{2}}.$$
 (85)

Hence $\hat{y}_{t_0}, \, \hat{y}_{t_0}^*$ can be coupled so that $\mathbb{P}(\hat{y}_{t_0} \neq \hat{y}_{t_0}^*) \leq \sqrt{\varepsilon/2}$.

We extend this to a coupling of $(\hat{\boldsymbol{y}}_t^*)_{t_0 \leq t \leq T}$ and $(\hat{\boldsymbol{y}}_t)_{t_0 \leq t \leq T}$ in the obvious way: we generate the two trajectories according to the discretized diffusion (4) with the same randomness $\hat{\boldsymbol{z}}_t$. Therefore $\mathbb{P}(\hat{\boldsymbol{y}}_T \neq \hat{\boldsymbol{y}}_T^*) \leq \sqrt{\varepsilon/2}$. Another application of the assumption $D_{\mathrm{KL}}(\overline{\mathbb{P}}_{\hat{\boldsymbol{y}}}^{T,\Delta} \| \mathbb{P}_{\boldsymbol{y}}^T) \leq \varepsilon$ and Pinsker's inequality yields $\mathbb{P}(\boldsymbol{y}_T \neq \hat{\boldsymbol{y}}_T^*) \leq \sqrt{\varepsilon/2}$, for $\boldsymbol{y}_T \stackrel{\mathrm{d}}{=} T\boldsymbol{x} + \sqrt{T}\boldsymbol{g}'$ with $(\boldsymbol{x}, \boldsymbol{g}') \sim \mu \otimes \mathsf{N}(\boldsymbol{0}, \boldsymbol{I})$. We conclude by triangle inequality $\mathbb{P}(\boldsymbol{y}_T \neq \hat{\boldsymbol{y}}_T) \leq 2\sqrt{\varepsilon/2}$, which coincides with the claim (83).

Finally, Eq. (84) follows by generating N i.i.d. copies $\hat{x}_1, \dots, \hat{x}_N$ using the above procedure, and letting $\hat{m}(y)$ be their empirical average.

S.2 Proof of Theorem 5

The next statement makes a weaker assumption on the accuracy of the diffusion sampler (transportation instead of KL distance), but in exchange assumes the approximate drift $\hat{\boldsymbol{m}}$ to be Lipschitz. We note that $\mathrm{Lip}(\boldsymbol{m}(\,\cdot\,,t)) = \sup_{\boldsymbol{y}} \|\mathrm{Cov}(\boldsymbol{x}|\boldsymbol{y}_t = \boldsymbol{y})\|_{\mathrm{op}}$, and the latter is of O(1/d) (for instance) if the coordinates of \boldsymbol{x} are weakly dependent under the posterior.

Theorem 5. Assume that $\hat{\boldsymbol{m}}(\cdot, \cdot)$ has computational complexity χ and satisfies the following: (a) For every $t \geq 1/\sigma^2$, $\boldsymbol{y} \mapsto \hat{\boldsymbol{m}}(\boldsymbol{y},t)$ is L/d-Lipschitz. (b) There is a stepsize Δ such that $W_1(\mathrm{P}^{T,\Delta}_{\hat{\boldsymbol{y}}},\mathrm{P}^T_{\boldsymbol{y}}) \leq \varepsilon$ for any $T \leq \theta d$.

Then for any $\sigma > 0$ there exists an algorithm with complexity $O(\chi \cdot T/\Delta)$, that takes as input $y = x + \sigma g$, $(x, g) \sim \mu \otimes \mathsf{N}(0, I)$, and outputs \hat{x} , such that

$$\mathbb{E}_{\boldsymbol{y}}W_1(P_{\boldsymbol{x}|\boldsymbol{y}}, P_{\hat{\boldsymbol{x}}|\boldsymbol{y}}) \le 2e^{\theta L}\varepsilon + \frac{1}{\sqrt{\theta}} =: \overline{\varepsilon}.$$
 (86)

As a consequence, Eq. (13) holds also in this case with the new definition of $\overline{\varepsilon}$.

The algorithm consists in running the discretized diffusion (4) with initialization $\hat{y}_{t_0} = y/\sigma^2$ at $t = t_0 := 1/\sigma^2$. To avoid notational burden, we will assume $(T - t_0)/\Delta$ to be an integer. Let $\hat{y}_{t_0}^*$ be generated by the discretized diffusion with initialization at \hat{y}_0 at t = 0. Note that the distribution of \hat{y}_{t_0} is the same as the one of $t_0x + \sqrt{t}g$ and hence by Assumption (b),

$$W_1(P_{\hat{y}_{t_0}}, P_{\hat{y}_{t_0}^*}) \le W_1(P_{T,\Delta}^{\hat{y}}, P_T^{y}) \le \varepsilon.$$
 (87)

In other words there exists a coupling of $\hat{y}_{t_0}^*$ and \hat{y}_{t_0} such that $\mathbb{E}\|\hat{y}_{t_0}^* - \hat{y}_{t_0}\|_2 \leq \varepsilon$.

We extend this to a coupling of $(\hat{y}_t^*)_{t_0 \le t \le T}$ and $(\hat{y}_t)_{t_0 \le t \le T}$ in the obvious way: we generate the two trajectories according to the discretized diffusion (4) with the same randomness \hat{z}_t . A simple recursive argument (using the Lipschitz property of \hat{m} , in Assumption (a)) then yields

$$\mathbb{E}\|\hat{\boldsymbol{y}}_{T}^{*} - \hat{\boldsymbol{y}}_{T}\|_{2} \le \left(1 + L\Delta/d\right)^{T/\Delta} \varepsilon \le e^{LT/d} \varepsilon. \tag{88}$$

(See for instance Montanari & Wu (2023) or Alaoui et al. (2023) for examples of this calculation.) Let now $\boldsymbol{y}_T \stackrel{\mathrm{d}}{=} T\boldsymbol{x} + \sqrt{T}\boldsymbol{g}'$ for $(\boldsymbol{x},\boldsymbol{g}') \sim \mu \otimes \mathsf{N}(\boldsymbol{0},\boldsymbol{I})$. Another application of Assumption (a) implies that this can be coupled to $\hat{\boldsymbol{y}}_T^*$ so that $\mathbb{E}\|\boldsymbol{y}_T - \hat{\boldsymbol{y}}_T^*\| \leq \varepsilon$, and therefore

$$\mathbb{E}\|\hat{\mathbf{y}}_T - \mathbf{y}_T\|_2 \le 2e^{LT/d}\varepsilon. \tag{89}$$

As output, we return $\hat{x} = \hat{y}_T/T$. Using $\mathbb{E}||y_T - x|| = \mathbb{E}||g||/\sqrt{T}$ and $T = \theta d$,

$$\mathbb{E}\|\boldsymbol{x} - \hat{\boldsymbol{x}}\| \le 2e^{\theta L}\varepsilon + \frac{1}{\sqrt{\theta}}.$$
 (90)

Since the coupling has been constructed conditionally on y, the claim (86) follows.

Finally, Eq. (13) follows by generating N i.i.d. copies $\hat{x}_1, \dots, \hat{x}_N$ using the above procedure, and letting $\hat{m}(y)$ be their empirical average.

T PROOF OF LEMMA G.1

We let $B \sim N(0, I_{n^2})$ so that $W \sim (B + B^{\intercal})/2$. We consider a non-random vector v fitting the description, and note that

$$|\langle \boldsymbol{v}, \boldsymbol{W} \boldsymbol{v} \rangle| \leq \frac{1}{2} \left\{ |\langle \boldsymbol{v}, \boldsymbol{B} \boldsymbol{v} \rangle| + |\langle \boldsymbol{v}, \boldsymbol{B}^\intercal \boldsymbol{v} \rangle| \right\}$$

We know that $\langle \boldsymbol{v}, \boldsymbol{B}\boldsymbol{v} \rangle, \langle \boldsymbol{v}, \boldsymbol{B}^{\intercal}\boldsymbol{v} \rangle \sim \mathsf{N}(0,1)$. We thus have, by Gaussian tail bounds and a triangle inequality:

$$\mathbb{P}\left(|\langle \boldsymbol{v}, \boldsymbol{W} \boldsymbol{v} \rangle| \geq C \sqrt{\log \binom{n}{k}}\right) \leq 2 \exp\left(-\frac{C^2}{2} \log \binom{n}{k}\right)$$

Union bounding over the set of all such vectors gives us the desired statement, as the cardinality of this set is $\binom{n}{k} 2^k$.

U PROOF OF LEMMA G.3

Proof of Lemma G.3: Using Lemma G.2, we can take $x = n^{-1/4}$, say, and $\theta = \sqrt{1+\delta}$ for $\delta = o_n(1)$, to get that

$$\mathbb{P}\left(\lambda_1(\boldsymbol{y}) \leq \theta + 1/\theta - n^{-1/4} - 2/n\right) \leq C \exp(-cn^{1/2})$$

for some absolute constants C, c > 0.

We have the following identity, letting $W \sim \mathsf{GOE}(n, 1/n)$:

$$\langle \boldsymbol{u}, \boldsymbol{v}_1 \rangle^2 = \frac{1}{\theta^2 \langle \boldsymbol{u}, (\lambda_1(\boldsymbol{y})\boldsymbol{I} - \boldsymbol{W})^{-2}\boldsymbol{u} \rangle} \geq \frac{1}{\theta^2 \cdot \|(\lambda_1(\boldsymbol{y})\boldsymbol{I} - \boldsymbol{W})^{-2}\|_{\text{op}}}$$

By standard Gaussian concentration, we know that

$$\mathbb{P}(\|\boldsymbol{W}\|_{\text{op}} \ge 2 + x) \le C \exp(-cnx^2)$$

We take

$$x = \frac{\theta + 1/\theta - n^{-1/4} - 2/n - 2}{4}$$

Note that with $\theta = \sqrt{1+\delta}$ and $\delta = o_n(1)$, we know that $\theta + 1/\theta - 2 = \Theta(\delta^2)$, so that $x = \Theta(\delta^2)$ if $\delta \gg n^{-1/8}$. Furthermore we have, by Theorem 1 above,

$$\mathbb{P}\left(\lambda_{\min}(\lambda_1(\boldsymbol{y})\boldsymbol{I} - \boldsymbol{W}\right) \le 2x\right) \le C \exp(-cn^{1/2})$$

and on the complement of this event, we know that

$$\langle \boldsymbol{v}_1, \boldsymbol{u} \rangle^2 \ge \frac{4x^2}{\theta^2} = \Theta(\delta^4)$$

since $\theta = \Omega(1)$, and so $|\langle \boldsymbol{v}_1, \boldsymbol{u} \rangle| = \Omega(\delta^2)$. Hence we are done.

V PROOF OF THEOREM 3

The optimality of \hat{m} with respect to scalings $c\hat{m}$ implies, by Pythagoras' theorem:

$$\mathbb{E}\{\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t)-\boldsymbol{x}\|^2\} = \mathbb{E}\{\|\boldsymbol{x}\|^2\} - \mathbb{E}\{\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t)\|^2\},$$

whence, using assumption (16), we obtain that

$$\sup_{t \le (1-\gamma)t_{\text{alg}}} \mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t, t)\|^2] = o(t_{\text{alg}}^{-1}). \tag{91}$$

Recall that (\hat{y}_t) is the generated diffusion, defined in Eq.(4). From Girsanov's formula on $[0, (1-\gamma)t_{alg}]$, we get that:

$$\mathsf{KL}\left((\boldsymbol{y}_t)_{t\in\mathbb{N}\Delta\cap[0,(1-\gamma)t_{\mathrm{alg}}]}\|(\hat{\boldsymbol{y}}_t)_{t\in\mathbb{N}\Delta\cap[0,(1-\gamma)t_{\mathrm{alg}}]}\right) = \frac{\Delta}{2}\sum_{t\in\mathbb{N}\Delta\cap[0,(1-\gamma)t_{\mathrm{alg}}]}\mathbb{E}[\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t)\|^2] = o(1)$$

From Eq. (97), we get from Markov's inequality that with high probability,

$$\frac{\Delta}{2} \sum_{t \in \mathbb{N} \Delta \cap [0, (1-\gamma)t_{\text{alg}}]} \|\hat{\boldsymbol{m}}(\boldsymbol{y}_t, t)\|^2 = o(1) \stackrel{(a)}{\Rightarrow} \frac{\Delta}{2} \sum_{t \in \mathbb{N} \Delta \cap [0, (1-\gamma)t_{\text{alg}}]} \|\hat{\boldsymbol{m}}(\boldsymbol{y}_t, t)\| = o(\sqrt{t_{\text{alg}}}),$$

where (a) follows by Cauchy-Schwarz. By Pinsker's inequality on Eq. (98), we obtain that the same event holds for (\hat{y}_t) with high probability:

$$\frac{\Delta}{2} \sum_{t \in \mathbb{N} \Delta \cap [0, (1-\gamma)t_{\text{alg}}]} \|\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_t, t)\| = o(\sqrt{t_{\text{alg}}}).$$

Fix a constant $\varepsilon_0 > 0$ to be chosen later. By taking the constant γ to be close enough to 1, we get that for $t_b := \min\{\ell\Delta : \ell\Delta \geq (1+\delta)t_{\text{alg}}\}$:

$$\hat{\boldsymbol{y}}_{t_b} = \boldsymbol{B}_{t_b} + \Delta \sum_{t \in \mathbb{N} \Delta \cap [0, t_b]} \hat{\boldsymbol{m}}_n(\hat{\boldsymbol{y}}_t, t) := \boldsymbol{m}_0 + \boldsymbol{B}_{t_b}$$

with $\mathbb{P}(\|\boldsymbol{m}_0\| \geq \varepsilon_0 t_{\text{alg}}) = o(1)$. Next we couple $(\hat{\boldsymbol{y}}_t : t \geq t_b)$ to $(\hat{\boldsymbol{y}}_t^0 : t \geq t_b)$ defined by letting $\hat{\boldsymbol{y}}_{t_b}^0 = \boldsymbol{B}_{t_b}$ and, for $t \in \mathbb{N}\Delta \cap [t_b, \infty)$,

$$\hat{\boldsymbol{y}}_{t+\Delta}^0 = \hat{\boldsymbol{y}}_t^0 + \hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_t^0, t)\Delta + \boldsymbol{B}_{t+\Delta} - \boldsymbol{B}_t$$
.

By the assumed Lipschitz property of \hat{m} and Gromwall's lemma:

$$\|\hat{\boldsymbol{y}}_{t} - \hat{\boldsymbol{y}}_{t}^{0}\| \leq \prod_{t' \in \mathbb{N}\Delta \cap [t_{b}, t]} \left(1 + \frac{C\Delta}{t'}\right) \cdot \|\boldsymbol{m}_{0}\|$$

$$\leq \left(\frac{t}{t_{\text{alg}}}\right)^{C'} \|\boldsymbol{m}_{0}\| \leq C' \varepsilon_{0} t_{\text{alg}}, \tag{92}$$

where the last inequality holds for some absolute constant C' and all $t \leq Ct_{alg}$, on the high probability event $\|\boldsymbol{m}_0\| \leq \varepsilon_0 t_{alg}$.

We are now in position to finish the proof of the theorem. We couple the process $(\hat{y}_t^0: t \ge t_b)$ defined above with $(B_t: t \ge t_b)$ to get

$$\mathsf{KL}(\boldsymbol{B}_{t+\Delta}\|\hat{\boldsymbol{y}}_{t+\Delta}^{0}) \leq \mathsf{KL}(\boldsymbol{B}_{t}\|\boldsymbol{y}_{t}^{0}) + C\,\mathbb{E}\big\{\|\hat{\boldsymbol{m}}(\boldsymbol{B}_{t},t)\|^{2}\big\} \cdot \Delta$$

Using $\mathsf{KL}(\boldsymbol{B}_{t_b} \| \hat{\boldsymbol{y}}_{t_b}^0) = 0$, summing the last inequality over $t \geq t_b$, and applying Pinsker's inequality we obtain, with C' a suitably large constant

$$\sup_{t\in\mathbb{N}\Delta\cap[t_{\mathrm{alg}}(1+\delta),\infty)}\mathsf{TV}(\hat{\boldsymbol{y}}_t^0,\boldsymbol{B}_t)=o(1)\,. \tag{93}$$

Putting together this bound and Eq. (99) (which holds with high probability) we obtain that

$$\mathbb{P}\left(\max_{t \in [t_{\text{alg}}(1+\delta), Ct_{\text{alg}}]} \|\hat{\boldsymbol{y}}_t - \boldsymbol{B}_t\| \ge C' \varepsilon_0 t_{\text{alg}}\right) = o_n(1). \tag{94}$$

We collapse $C'\varepsilon_0$ into ε_0 , as ε_0 is arbitrary. Using once more Lemma W.1 and the Lipschitz property of \hat{m} , we obtain that, for $\hat{x}_t = \hat{m}(\hat{y}_t, t)$,

$$\mathbb{P}\left(\max_{t \in [t_{\text{alg}}(1+\delta), Ct_{\text{alg}}]} \|\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_t, t)\| \ge \varepsilon_0\right) = o(1), \tag{95}$$

which implies that

$$\inf_{t \in [t_{\text{alg}}(1+\delta), Ct_{\text{alg}}]} W_1(\hat{\boldsymbol{m}}(\hat{\boldsymbol{y}}_t, t), \boldsymbol{x}) \ge \alpha - \varepsilon_0 + o(1)$$

By taking $\varepsilon_0 \downarrow 0$, we obtain the claim of the theorem.

W Proof of Corollary 5.1

In order to simplify some of the formulas below we center $\mu_{n,k}$. Namely, we redefine $\mu_{n,k}$ to be the distribution of $\mathbf{x} = \mathbf{u}\mathbf{u}^\mathsf{T} - \mathbb{E}[\mathbf{u}\mathbf{u}^\mathsf{T}]$ when $\mathbf{u} \sim \mathrm{Unif}(B_{n,k})$.

Throughout this proof, C denotes a generic constant which depends on the constants in the assumptions, and is allowed to change from line to line. We will write $\mathbb{E}_{n,k}$ for expectation under $\mu_{n,k}$ and \mathbb{E} for expectation under $\overline{\mu}_{n,k} = \frac{1}{2}\mu_{n,k} + \frac{1}{2}\delta_{\mathbf{0}}$. Further $\mathbf{y}_t = t\mathbf{x} + \mathbf{W}_t$, where the distribution of \mathbf{x} is either $\mu_{n,k}$ or $\overline{\mu}_{n,k}$ as indicated.

The optimality of \hat{m}_n with respect to scalings $c\hat{m}_n$ implies, by Pythagoras' theorem:

$$\mathbb{E}\left\{\|\hat{\boldsymbol{m}}_n(\boldsymbol{y}_t,t)-\boldsymbol{x}\|^2\right\} = \mathbb{E}\{\|\boldsymbol{x}\|^2\} - \mathbb{E}\left\{\|\hat{\boldsymbol{m}}_n(\boldsymbol{y}_t,t)\|^2\right\},\,$$

whence, using assumption (16), we obtain that

$$\sup_{t \le (1-\gamma)t_{\text{alg}}} \mathbb{E}[\|\hat{\boldsymbol{m}}_n(\boldsymbol{y}_t, t)\|^2] = o_n(n^{-1}).$$
(96)

By Law of Total Probability, we have

$$\mathbb{E}[\|\hat{\boldsymbol{m}}_n(\boldsymbol{y}_t,t)\|^2] = \frac{1}{2}\mathbb{E}_{\boldsymbol{x} \sim \mu_{n,k}}[\|\hat{\boldsymbol{m}}_n(\boldsymbol{y}_t,t)\|^2] + \frac{1}{2}\mathbb{E}[\|\hat{\boldsymbol{m}}_n(\boldsymbol{W}_t,t)\|^2],$$

from which we get

$$\sup_{t \le (1-\gamma)t_{\text{alg}}} \mathbb{E}[\|\hat{\mathbf{m}}_n(\mathbf{W}_t, t)\|^2] = o_n(n^{-1})$$
(97)

From Girsanov's formula on $[0,(1-\gamma)t_{\mathrm{alg}}],$ we get that

$$\mathsf{KL}\left((\boldsymbol{W}_{t})_{t\in\mathbb{N}\Delta\cap[0,(1-\gamma)t_{\mathrm{alg}}]}\|(\hat{\boldsymbol{y}}_{t})_{t\in\mathbb{N}\Delta\cap[0,(1-\gamma)t_{\mathrm{alg}}]}\right) = \frac{\Delta}{2}\sum_{t\in\mathbb{N}\Delta\cap[0,(1-\gamma)t_{\mathrm{alg}}]}\mathbb{E}[\|\hat{\boldsymbol{m}}_{n}(\boldsymbol{W}_{t},t)\|^{2}] = o_{n}(1)$$
(98)

due to the fact that $t_{\text{alg}} = n/2$. From Eq. (97), we get from Markov's inequality that with high probability,

$$\frac{\Delta}{2} \sum_{t \in \mathbb{N} \Delta \cap [0, (1-\gamma)t_{\text{alg}}]} \|\hat{\boldsymbol{m}}_n(\boldsymbol{W}_t, t)\|^2 = o_n(1) \stackrel{(a)}{\Rightarrow} \frac{\Delta}{2} \sum_{t \in \mathbb{N} \Delta \cap [0, (1-\gamma)t_{\text{alg}}]} \|\hat{\boldsymbol{m}}_n(\boldsymbol{W}_t, t)\| = o_n(\sqrt{n}),$$

where (a) follows by Cauchy-Schwarz. By Pinsker's inequality on Eq. (98), we obtain that the same event holds for (\hat{y}_t) with high probability:

$$\frac{\Delta}{2} \sum_{t \in \mathbb{N} \Delta \cap [0, (1-\gamma)t_{\text{alg}}]} \|\hat{\boldsymbol{m}}_n(\hat{\boldsymbol{y}}_t, t)\| = o_n(\sqrt{n}).$$

Fix a constant $\varepsilon_0 > 0$ to be chosen later. By taking the constant γ to be close enough to 1, we get that for $t_b := \min\{\ell\Delta : \ell\Delta \ge (1+\delta)t_{\rm alg}\}$:

$$\hat{\boldsymbol{y}}_{t_b} = \boldsymbol{B}_{t_b} + \Delta \sum_{t \in \mathbb{N} \Delta \cap [0, t_b]} \hat{\boldsymbol{m}}_n(\hat{\boldsymbol{y}}_t, t) := \boldsymbol{m}_0 + \boldsymbol{B}_{t_b}$$

with $\mathbb{P}(\|\boldsymbol{m}_0\| \geq \varepsilon_0 n) = o_n(1)$, and $(\hat{\boldsymbol{y}}_t)$ is the generated diffusion, defined in Eq.(4). Next we couple $(\hat{\boldsymbol{y}}_t: t \geq t_b)$ to $(\hat{\boldsymbol{y}}_t^0: t \geq t_b)$ defined by letting $\hat{\boldsymbol{y}}_{t_b}^0 = \boldsymbol{B}_{t_b}$ and, for $t \in \mathbb{N}\Delta \cap [t_b, \infty)$,

$$\hat{m{y}}_{t+\Delta}^0 = \hat{m{y}}_t^0 + \hat{m{m}}_n(\hat{m{y}}_t^0,t)\Delta + m{B}_{t+\Delta} - m{B}_t$$
 .

By the assumed Lipschitz property of \hat{m} and Gromwall's lemma:

$$\|\hat{\boldsymbol{y}}_{t} - \hat{\boldsymbol{y}}_{t}^{0}\| \leq \prod_{t' \in \mathbb{N}\Delta \cap [t_{b}, t]} \left(1 + \frac{C\Delta}{t'}\right) \cdot \|\boldsymbol{m}_{0}\|$$

$$\leq \left(\frac{t}{t_{\text{alg}}}\right)^{C'} \|\boldsymbol{m}_{0}\| \leq C' \varepsilon_{0} n, \qquad (99)$$

where the last inequality holds for some absolute constant C' and all $t \leq Cn = (2C)t_{alg}$, on the high probability event $\|\boldsymbol{m}_0\| \leq \varepsilon_0 n$.

In order to finish the proof, we state and prove a useful lemma. In a nutshell, \hat{m}_n resists improvements from eigenvalue hypothesis tests:

Lemma W.1. Under the assumptions of Theorem 5.1, assume that δ_n vanishes slowly enough. Then, for $t \geq (1+\delta)t_{alg}$,

$$\mathbb{E}\left\{\|\hat{\boldsymbol{m}}_{n}(\boldsymbol{B}_{t},t)\|^{2}\right\} \leq Ce^{-(\sqrt{t}-\sqrt{t_{\text{alg}}})^{4}/Cn}.$$
(100)

Proof. Let $\lambda_1(\boldsymbol{y}_t)$ be the maximum eigenvalue of $(\boldsymbol{y}_t + \boldsymbol{y}_T^\mathsf{T})/\sqrt{2}$, $\lambda_*(t) := \sqrt{2}(\sqrt{t} + \sqrt{t_{\text{alg}}})^2$ and $\phi(\boldsymbol{y}_t) := \mathbf{1}(\lambda_1(\boldsymbol{y}_t) > \lambda_*(t))$. Concentration results about spiked GOE matrices imply, for all $t \geq (1+\delta)t_{\text{alg}}$,

$$\mathbb{P}_{n,k}(\phi(\boldsymbol{y}_t) = 0) \le Ce^{-n(\sqrt{t} - \sqrt{t_{\text{alg}}})^4/C}, \quad \mathbb{P}(\phi(\boldsymbol{B}_t) = 0) \le C \exp\left\{-\frac{1}{Cn}(\sqrt{t} - \sqrt{t_{\text{alg}}})^4\right\}. \tag{101}$$

(To simplify notations, we omit the dependence of ϕ on t.)

By assumption, the MSE of $\hat{\boldsymbol{m}}_n(\boldsymbol{y}_t,t)$ is not larger than the one of $\hat{\boldsymbol{m}}_n(\boldsymbol{y}_t,t)\phi(\boldsymbol{y}_t)$. Letting $\overline{\phi}(\boldsymbol{y}_t):=1-\phi(\boldsymbol{y}_t)$:

$$\mathbb{E}\{\|\hat{\boldsymbol{m}}(\boldsymbol{y}_{t},t)-\boldsymbol{x}\|^{2}\} \leq \mathbb{E}\{\|\hat{\boldsymbol{m}}(\boldsymbol{y}_{t},t)\phi(\boldsymbol{y}_{t})-\boldsymbol{x}\|^{2}\}$$

$$= \mathbb{E}\{\|\hat{\boldsymbol{m}}(\boldsymbol{y}_{t},t)-\boldsymbol{x}\|^{2}\phi(\boldsymbol{y}_{t})\} + \mathbb{E}\{\|\boldsymbol{x}\|^{2}\overline{\phi}(\boldsymbol{y}_{t})\}$$

$$= \mathbb{E}\{\|\hat{\boldsymbol{m}}(\boldsymbol{y}_{t},t)-\boldsymbol{x}\|^{2}\phi(\boldsymbol{y}_{t})\} + \mathbb{P}_{n,k}(\phi(\boldsymbol{y}_{t})=0), \qquad (102)$$

whence

$$\mathbb{E}\left\{\|\hat{\boldsymbol{m}}(\boldsymbol{y}_t,t) - \boldsymbol{x}\|^2 \overline{\phi}(\boldsymbol{y}_t)\right\} \le \mathbb{P}_{n,k}(\phi(\boldsymbol{y}_t) = 0). \tag{103}$$

On the other hand

$$\mathbb{E}\{\|\hat{\boldsymbol{m}}(\boldsymbol{y}_{t},t)-\boldsymbol{x}\|^{2}\overline{\phi}(\boldsymbol{y}_{t})\} \geq \mathbb{E}\{\|\hat{\boldsymbol{m}}(\boldsymbol{y}_{t},t)\|^{2}\boldsymbol{1}_{\boldsymbol{x}=0}\overline{\phi}(\boldsymbol{y}_{t})\}$$

$$= \frac{1}{2}\mathbb{E}\{\|\hat{\boldsymbol{m}}(\boldsymbol{B}_{t},t)\|^{2}\overline{\phi}(\boldsymbol{W}_{t})\}$$

$$\geq \frac{1}{2}\mathbb{E}\{\|\hat{\boldsymbol{m}}(\boldsymbol{B}_{t},t)\|^{2}\} - \frac{1}{2}\mathbb{P}(\phi(\boldsymbol{W}_{t})=1). \tag{104}$$

Putting together Eqs. (101), (103), (104), we obtain (eventually adjusting the constant C)

$$\mathbb{E}\left\{\|\hat{\boldsymbol{m}}(\boldsymbol{B}_t,t)\|^2\right\} \leq \mathbb{P}\left(\phi(\boldsymbol{B}_t) = 1\right) + 2\,\mathbb{P}_{n,k}\left(\phi(\boldsymbol{y}_t) = 0\right)$$
$$\leq C\exp\left\{-\frac{1}{Cn}(\sqrt{t} - \sqrt{t_{\text{alg}}})^4\right\}.$$

We are now in position to finish the proof of the theorem. We couple the process $(\hat{y}_t^0: t \ge t_b)$ defined above with $(B_t: t \ge t_b)$ to get

$$\begin{split} \mathsf{KL}(\boldsymbol{B}_{t+\Delta} \| \hat{\boldsymbol{y}}_{t+\Delta}^0) & \leq \mathsf{KL}(\boldsymbol{B}_t \| \boldsymbol{y}_t^0) + C \, \mathbb{E} \big\{ \| \hat{\boldsymbol{m}}(\boldsymbol{B}_t, t) \|^2 \big\} \cdot \Delta \\ & \leq \mathsf{KL}(\boldsymbol{B}_t \| \hat{\boldsymbol{y}}_t^0) + C \Delta \exp \Big\{ - \frac{1}{Cn} (\sqrt{t} - \sqrt{t_{\text{alg}}})^4 \Big\} \,. \end{split}$$

Using $\mathsf{KL}(\boldsymbol{B}_{t_b} \| \hat{\boldsymbol{y}}_{t_b}^0) = 0$, summing the last inequality over $t \geq t_b$, and applying Pinsker's inequality we obtain for $\delta_n \geq C'(\log n/n)^{1/4}$ with C' a suitably large constant

$$\sup_{t \in \mathbb{N} \Delta \cap [t_{\text{alg}}(1+\delta), \infty)} \mathsf{TV}(\hat{\boldsymbol{y}}_t^0, \boldsymbol{B}_t) = o_n(1). \tag{105}$$

Putting together this bound and Eq. (99) (which holds with high probability) we obtain that

$$\mathbb{P}\left(\max_{t \in [t_{\text{alg}}(1+\delta), Cn]} \|\hat{\boldsymbol{y}}_t - \boldsymbol{B}_t\| \ge C' \varepsilon_0 n\right) = o_n(1). \tag{106}$$

We collapse $C'\varepsilon_0$ into ε_0 , as ε_0 is arbitrary. Using once more Lemma W.1 and the Lipschitz property of \hat{m}_n , we obtain that, for $\hat{x}_t = \hat{m}(\hat{y}_t, t)$,

$$\mathbb{P}\left(\max_{t \in [t_{\text{alg}}(1+\delta), C_n]} \|\hat{\boldsymbol{m}}_n(\hat{\boldsymbol{y}}_t, t)\| \ge \varepsilon_0\right) = o_n(1), \tag{107}$$

which implies that

$$\inf_{t \in [t_{\text{alg}}(1+\delta), Cn]} W_1(\hat{\boldsymbol{m}}_n(\hat{\boldsymbol{y}}_t, t), \boldsymbol{x}) \ge 1/2 - \varepsilon_0 + o_n(1)$$

By taking $\varepsilon_0 \downarrow 0$, we obtain the claim of the theorem.

X DETAILS OF NUMERICAL SIMULATIONS

Our GNN architecture uses node embeddings that are generated by 3 iterations of the power method and 10 message passing layers. Each message passing layer comprises of the 'message' and 'node-update' multi-layer perceptrons (MLPs), both of which are 2-layer neural networks with LeakyReLU

nonlinearity. We simply use the complete graph with self-loops for node embedding updates. We find that 'seeding' the node embeddings with iterations of power method is crucial for effective training.

During training of the denoiser, we sample time points t as follows: choose a time threshold t_{\star} , and sample so that times $t > t_{\star}$ are picked with total probability 0.95 (and times $t \leq t_{\star}$ are picked with total probability 0.05). Within each interval $(0, t_{\star}]$ and (t_{\star}, T) , times are chosen at random. This allows the neural network to initially prioritize learning in a low-noise regime. Several fine-tuning steps are taken, for which t_{\star} is gradually decreased to refine the network on lower SNR.

Empirically, training directly with 10 layers is difficult, due to its depth. We find that training initially with 7 layers, then subsequently introducing the later layers results in more stable training.

We train such a network using N=30000 samples x_i from the distribution $\overline{\mu}_{n,k}$, and evaluate their MSE on $N_{\text{test}}=15000$ samples.