

---

# A Solvable High-Dimensional Model of GAN

---

Chuang Wang<sup>1,2</sup>  
wangchuang@ia.ac.cn

Hong Hu<sup>2</sup>  
honghu@g.harvard.edu

Yue M. Lu<sup>2</sup>  
yuelu@seas.harvard.edu

1. State Key Laboratory of Pattern Recognition, Institute of Automation,  
Chinese Academy of Science, 95 Zhong Guan Cun Dong Lu, Beijing 100190, China
2. John A. Paulson School of Engineering and Applied Sciences, Harvard University  
33 Oxford Street, Cambridge, MA 02138, USA

## Abstract

We present a theoretical analysis of the training process for a single-layer GAN fed by high-dimensional input data. The training dynamics of the proposed model at both microscopic and macroscopic scales can be exactly analyzed in the high-dimensional limit. In particular, we prove that the macroscopic quantities measuring the quality of the training process converge to a deterministic process characterized by an ordinary differential equation (ODE), whereas the microscopic states containing all the detailed weights remain stochastic, whose dynamics can be described by a stochastic differential equation (SDE). This analysis provides a new perspective different from recent analyses in the limit of small learning rate, where the microscopic state is always considered deterministic, and the contribution of noise is ignored. From our analysis, we show that the level of the background noise is essential to the convergence of the training process: setting the noise level too strong leads to failure of feature recovery, whereas setting the noise too weak causes oscillation. Although this work focuses on a simple copy model of GAN, we believe the analysis methods and insights developed here would prove useful in the theoretical understanding of other variants of GANs with more advanced training algorithms.

## 1 Introduction

A generative adversarial network (GAN) [1] seeks to learn a high-dimensional probability distribution from samples. While there have been numerous advances on the application front [2–6], considerably less is known about the underlying theory and conditions that can explain or guarantee the successful trainings of GANs.

Recently, it has been a very active area of research to study either the equilibrium properties [7–9] or the training dynamics [10, 11]. Specifically, there is a line of works studying the dynamics of the gradient-based training algorithms *e.g.*, [11–16]. The basic idea is the following. The evolution of the learnable parameters in the training dynamics can be considered as a discrete-time process. With a proper time scaling, this discrete-time process converges to a deterministic continuous-time process as the learning rates tend to 0, which is characterized by an ordinary differential equation (ODE). By studying local stability of the ODE’s fixed points, [12] shows that oscillation in the training algorithm is due to the eigenvalues of the Jacobian of the gradient vector field with zero real part and large imaginary part. Due to this fact, various stabilization approaches are proposed, for example adding additional regularizers [13, 14], and using two timescale [15] training. Very recently, [16] argues that those stabilization techniques may encourage the algorithms to converge non-Nash stationary points. All above works consider a small-learning-rates limit, where the limiting process

is always deterministic. The stochasticity and the effect of the noise is essentially ignored, which may not reflect practical situations. Thus, a new analysis paradigm to study the dynamics with the consideration of the intrinsic stochasticity is needed.

In this paper, we present a *high-dimensional* and *exactly solvable* model of GAN. Its dynamics can be precisely characterized at both macroscopic and microscopic scales, where the former is deterministic and the latter remains stochastic. Interestingly, our theoretical analysis shows that injecting additional noise can stabilize the training. Specifically, our main technical contributions are twofold:

- We present an asymptotically exact analysis of the training process of the proposed GAN model. Our analysis is carried out on both the *macroscopic* and the *microscopic* levels. The macroscopic state measures the overall performance of the training process, whereas the microscopic state contains all the detailed weights information. In the high-dimensional limit ( $n \rightarrow \infty$ ), we show that the former converges to a deterministic process governed by an ordinary differential equation (ODE), whereas the latter stays stochastic described by a stochastic differential equation (SDE).
- We show that depending on the choice of the learning rates and the strength of noise, the training process can reach either a successful, a failed, an oscillating, or a mode-collapsing phase. By studying the stabilities of the fixed points of the limiting ODEs, we precisely characterize when each phase takes place. The analysis reveals a condition on the learning rates and the noise strength for successful training. We show that the level of the background noise is essential to the convergence of the training process: setting the noise level too strong (small signal-to-noise ratio) leads to failure of feature recovery, whereas setting the noise too weak (large signal-to-noise ratio) causes oscillation.

Our work builds upon a general analysis framework [17] for studying the scaling limits of high-dimensional exchangeable stochastic processes with applications to nonlinear regression problems. Similar techniques have also been used in the literature to study Monte Carlo methods [18], online perceptron learning [19, 20], online sparse PCA [21], subspace estimation [22], online ICA [23] and more recently, the supervised learning of two-layer neural networks [24], but to our best knowledge, this technique has not yet been used in analyzing GANs.

The rest of the paper is organized as follows. We present the proposed GAN model and the associated training algorithm in Section 2. Our main results are presented in Section 3, where we show that the macroscopic and microscopic dynamics of the training process converge to their respective limiting processes that are characterized by an ODE and SDE, respectively. In Section S-I, we analyze the stationary solutions of the limiting ODEs and precisely characterizes the long-term behaviors of the training process. We conclude in Section 5.

## 2 Formulations

In this section, we introduce the proposed GAN model and specify the associated training algorithm.

**Model for the real data.** In order to establish the theoretical analysis, we first impose a model for the probability distribution from which we draw our real data samples. We assume that the real data  $\mathbf{y}_k \in \mathbb{R}^n$ ,  $k = 0, 1, \dots$  are drawn according to the following generative model:

$$\mathbf{y}_k = \mathcal{G}(\mathbf{c}_k, \mathbf{a}_k; \mathbf{U}, \eta_T) \stackrel{\text{def}}{=} \mathbf{U}\mathbf{c}_k + \sqrt{\eta_T}\mathbf{a}_k, \quad (1)$$

where  $\mathbf{U} \in \mathbb{R}^{n \times d}$  is a deterministic unknown feature matrix with  $d$  features;  $\mathbf{c}_k \in \mathbb{R}^d$  is a random vector drawn from an unknown distribution  $\mathcal{P}_c$ ;  $\mathbf{a}_k$  is an  $n$ -dimensional random vector acting as the background noise; and  $\eta_T$  is a parameter to control the strength of noise. Without loss of generality<sup>1</sup>, we assume  $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_d$ , where  $\mathbf{I}_d$  is the  $d \times d$  identity matrix.

This generative model, referred to as the spiked covariance model [25] in the literature, is commonly used in the theoretical study of principal component analysis (PCA). We note that this model is not a trivial task for PCA even when  $d = 1$  if the variance of the noise  $\mathbf{a}_k$  is a non-zero constant. As

<sup>1</sup>If  $\mathbf{U}$  is not orthogonal, we can rewrite  $\mathbf{U}\mathbf{c}$  in (1) as  $(\mathbf{U}\mathbf{R})(\mathbf{R}^{-1}\mathbf{c})$ , where  $\mathbf{R}$  is a matrix that orthogonalizes and normalizes the columns of  $\mathbf{U}$ . We can then study an equivalent system where the new feature vector is  $\mathbf{R}^{-1}\mathbf{c}$ .

proved in [25], the best estimator can not perfectly recover the signal  $\mathbf{U}$  given an  $\mathcal{O}(n)$  number of samples  $\mathbf{y}_k$ . Thus, it is of sufficient interest to investigate whether a GAN can retrieve informative results for the principal components in the same scaling limit.

**The GAN model** The GAN we are going to analyze is defined as follows. We assume that the generator  $\mathcal{G}$  has the same linear structure as the real data model (1) given above:

$$\tilde{\mathbf{y}}_k = \mathcal{G}(\tilde{\mathbf{c}}_k, \tilde{\mathbf{a}}_k; \mathbf{V}, \eta_G) \quad (2)$$

but the parameters are different. Here,  $\tilde{\mathbf{y}}_k$  denotes a fake sample produced by the generator;  $\tilde{\mathbf{a}}_k$  is an  $n$ -dimensional random noise vector; the random variable  $\tilde{\mathbf{c}}_k$  is drawn from a fixed distribution  $\mathcal{P}_{\tilde{\mathbf{c}}}$ ;  $\eta_G$  is the noise strength; and the matrix  $\mathbf{V} \in \mathbb{R}^{n \times d}$  represents the parameters of the generator. (In an ideal case in which the generator learns the underlying true probability distribution perfectly, we have  $\mathbf{V} = \mathbf{U}$ .) Throughout the paper, we follow the notational convention that all the symbols that are decorated with a tilde (e.g.,  $\tilde{\mathbf{y}}_k, \tilde{\mathbf{c}}_k, \tilde{\mathbf{a}}_k$ ) denote quantities associated with the generator.

We define the discriminator  $\mathcal{D}$  of our GAN model as

$$\mathcal{D}(\mathbf{y}; \mathbf{w}) \stackrel{\text{def}}{=} \hat{D}(\mathbf{y}^\top \mathbf{w}).$$

Here,  $\mathbf{y}$  is an input vector, which can be either the real data  $\mathbf{y}_k$  from (1) or the fake one  $\tilde{\mathbf{y}}_k$  from (2);  $\hat{D} : \mathbb{R} \mapsto \mathbb{R}$  can be any function; and the vector  $\mathbf{w} \in \mathbb{R}^n$  represents the parameters associated with the discriminator. Later, we will show that the generator can learn multiple features even though the discriminator only has one feature vector  $\mathbf{w}$ . Discriminators with multiple features can also be analyzed in a similar way, but in this paper we consider the single-feature discriminator for simplicity.

**The training algorithm.** The proposed GAN model has two set of parameters  $\mathbf{V}$  and  $\mathbf{w}$  to be learned from the data. The training process is formulated as the following MinMax problem

$$\min_{\mathbf{V}} \max_{\mathbf{w}} \mathbb{E}_{\mathbf{y} \sim \mathcal{P}(\mathbf{y}; \mathbf{U})} \mathbb{E}_{\tilde{\mathbf{y}} \sim \tilde{\mathcal{P}}(\tilde{\mathbf{y}}; \mathbf{V})} \mathcal{L}(\mathbf{y}, \tilde{\mathbf{y}}; \mathbf{w}), \quad (3)$$

where the two probability distributions  $\mathcal{P}(\mathbf{y}; \mathbf{U})$  and  $\tilde{\mathcal{P}}(\tilde{\mathbf{y}}; \mathbf{V})$  represent the distributions of the real data  $\mathbf{y}$  and the fake data  $\tilde{\mathbf{y}}$  as specified by (1) and (2) respectively, and

$$\mathcal{L}(\mathbf{y}, \tilde{\mathbf{y}}; \mathbf{w}) \stackrel{\text{def}}{=} F(\hat{D}(\mathbf{y}^\top \mathbf{w})) - \tilde{F}(\hat{D}(\tilde{\mathbf{y}}^\top \mathbf{w})) - \frac{\lambda}{2} H(\mathbf{w}^\top \mathbf{w}) + \frac{\lambda}{2} \text{tr}(H(\mathbf{V}^\top \mathbf{V})) \quad (4)$$

with  $F(\cdot)$  and  $\tilde{F}(\cdot)$  being two functions that quantify the performance of the discriminator and  $\lambda > 0$  being a constant. The function  $H(\cdot)$  acts as a regularization term introduced to control the magnitude of the parameters  $\mathbf{w}$  and  $\mathbf{V}$ . It can be an arbitrary real-valued function, which is applied element-wisely if the input is a matrix.

We consider a standard training algorithm that uses the vanilla stochastic gradient descent/ascent (SGDA) to seek a solution of (3). To simplify the theoretical analysis, we consider an online (i.e., streaming) setting where each data sample  $\mathbf{y}_k$  is used only once. At step  $k$ , the model parameters  $\mathbf{w}_k$  and  $\mathbf{V}_k$  are updated using a new real sample  $\mathbf{y}_k$  and two fake samples  $\tilde{\mathbf{y}}_{2k}$  and  $\tilde{\mathbf{y}}_{2k+1}$ , according to

$$\begin{aligned} \mathbf{w}_{k+1} &= \mathbf{w}_k + \frac{\tau}{n} \nabla_{\mathbf{w}_k} \mathcal{L}(\mathbf{y}_k, \tilde{\mathbf{y}}_{2k}; \mathbf{w}_k) \\ \mathbf{V}_{k+1} &= \mathbf{V}_k - \frac{\tilde{\tau}}{n} \nabla_{\mathbf{V}_k} \mathcal{L}(\mathbf{y}_k, \mathcal{G}(\tilde{\mathbf{c}}_{2k+1}, \tilde{\mathbf{a}}_{2k+1}; \mathbf{V}_k; \eta_G); \mathbf{w}_k), \end{aligned} \quad (5)$$

where  $\tilde{\mathbf{c}}_{2k+1}, \tilde{\mathbf{a}}_{2k+1}$  are random variables that generates the fake sample  $\tilde{\mathbf{y}}_{2k+1}$  according to (2). The two parameters  $\tau$  and  $\tilde{\tau}$  in the above expressions control the learning rates of the discriminator and the generator, respectively. In (5), we only consider a single-step update for  $\mathbf{w}_k$ . This is a special case of Algorithm 1 in [1] with the batch-size  $m$  set to 1. We note that the analysis presented in this paper can be naturally extended to the mini-batch case where  $m$  is a finite number.

**Example 1.** We define  $F(\hat{D}(x)) = \tilde{F}(\hat{D}(x)) = x^2/2$ , and the regularizer function  $H(\mathbf{A}) = \log \cosh(\mathbf{A} - \mathbf{I})$ , where  $\mathbf{I}$  is the identity matrix with the same dimension of  $\mathbf{A}$ , and the function  $\log \cosh(\cdot)$  transforms the input matrix element-wisely. We use this specific regularizer to control the magnitude of the model parameters  $\mathbf{V}$  and  $\mathbf{w}$ . In practice, any convex function with its minimum reached at zero would be fine. Our choice  $\log \cosh(\mathbf{A} - \mathbf{I})$  here is just a convenient special case since its derivative  $H'(x) = \tanh(x)$  is smooth and bounded. Furthermore, we set the regularization parameter  $\lambda \rightarrow \infty$ , the original problem (3) becomes a constrained MinMax problem

$$\min_{\text{diag}(\mathbf{V}^\top \mathbf{V}) = \mathbf{I}_d} \max_{\|\mathbf{w}\|=1} \mathbb{E}_{\mathbf{y} \sim \mathcal{P}} \mathbb{E}_{\tilde{\mathbf{y}} \sim \tilde{\mathcal{P}}} \left[ (\mathbf{y}^\top \mathbf{w})^2 - (\tilde{\mathbf{y}}^\top \mathbf{w})^2 \right],$$

in which the diagonal operation  $\text{diag}(\mathbf{A})$  returns a matrix where the diagonal entries are the same as  $\mathbf{A}$  and the off-diagonal entries are all zero. The condition  $\text{diag}(\mathbf{V}^\top \mathbf{V}) = \mathbf{I}_d$  ensures that each column vector of  $\mathbf{V}$  is normalized.

### 3 Dynamics of the GAN

**Definition 1.** Let  $\mathbf{X}_k \stackrel{\text{def}}{=} [\mathbf{U}, \mathbf{V}_k, \mathbf{w}_k] \in \mathbb{R}^{n \times (2d+1)}$ . We call  $\mathbf{X}_k$  the *microscopic state* of the training process at iteration step  $k$ .

The microscopic state  $\mathbf{X}_k$  contains all the information about the training process. In fact, the sequence  $\{\mathbf{X}_k\}_{k=0,1,2,\dots}$  forms a Markov chain on  $\mathbb{R}^{n \times (2d+1)}$ . This can be easily verified from the update rule of  $\mathbf{X}_k$  as defined in (5), in which the real data  $\mathbf{y}_k$  and fake data  $\tilde{\mathbf{y}}_k$  are drawn according to (1) and (2) respectively. The Markov chain is driven by the initial state  $\mathbf{X}_0$  and the sequence of random variables  $\{(c_k, \mathbf{a}_k, \tilde{c}_{2k}, \tilde{\mathbf{a}}_{2k}, \tilde{c}_{2k+1}, \tilde{\mathbf{a}}_{2k+1})\}_{k=0,1,2,\dots}$ .

**Definition 2.** Let  $\mathbf{P}_k \stackrel{\text{def}}{=} \mathbf{U}^\top \mathbf{V}_k$ ,  $\mathbf{q}_k \stackrel{\text{def}}{=} \mathbf{U}^\top \mathbf{w}_k$ ,  $\mathbf{r}_k \stackrel{\text{def}}{=} \mathbf{V}_k^\top \mathbf{w}_k$ ,  $\mathbf{S}_k \stackrel{\text{def}}{=} \mathbf{V}_k^\top \mathbf{V}_k$ , and  $z_k \stackrel{\text{def}}{=} \mathbf{w}_k^\top \mathbf{w}_k$ . We call the tuple  $\{\mathbf{P}_k, \mathbf{q}_k, \mathbf{r}_k, \mathbf{S}_k, z_k\}$  the *macroscopic state* of the Markov chain  $\mathbf{X}_k$  at step  $k$ .

Those macroscopic quantities measure the cosine similarities among the feature vectors of the true model  $\mathbf{U}$ , the generator  $\mathbf{V}_k$  and the discriminator  $\mathbf{w}_k$ . For example, the cosine of the angle between the  $i$ th true feature (*i.e.*, the  $i$ th column of  $\mathbf{U}$ ) and the  $j$ th feature estimated in the generator (*i.e.*, the  $j$ th column of  $\mathbf{V}_k$ ) is  $[\mathbf{P}_k]_{i,j} / \sqrt{[\mathbf{S}_k]_{j,j}}$ , where  $[\mathbf{P}_k]_{i,j}$  is the inner product between the two feature vectors and  $\sqrt{[\mathbf{S}_k]_{j,j}}$  is the norm of the  $j$ th column of  $\mathbf{V}_k$ . (The columns of  $\mathbf{U}$  are unit vectors and need not be normalized here.) For simplicity, we introduce a compact notation for the macroscopic state:

$$\mathbf{M}_k \stackrel{\text{def}}{=} \mathbf{X}_k^\top \mathbf{X}_k = \begin{bmatrix} \mathbf{I} & \mathbf{P}_k & \mathbf{q}_k \\ \mathbf{P}_k^\top & \mathbf{S}_k & \mathbf{r}_k \\ \mathbf{q}_k^\top & \mathbf{r}_k^\top & z_k \end{bmatrix}. \quad (6)$$

In what follows, we investigate the dynamics of the training algorithm (5) at both the macroscopic and the microscopic levels. At the macroscopic level, by examining the cosines of the angles, we study how closely the model parameters  $\mathbf{V}_k$ ,  $\mathbf{w}_k$  associated with the generator and discriminator can align with the ground truth feature vectors, *i.e.*, the columns of  $\mathbf{U}$ . At the microscopic level, we study how the elements in the matrix  $\mathbf{V}_k$  and the vector  $\mathbf{w}_k$  evolve as a stochastic process. As our analysis will reveal, the mechanisms behind the two levels are different: the macroscopic dynamics is asymptotically deterministic whereas the microscopic dynamics stays stochastic even as  $n \rightarrow \infty$ .

#### 3.1 Macroscopic dynamics

We first study the asymptotic dynamics of the macroscopic state  $\mathbf{M}_k$ . Our theoretical analysis is carried out under the following assumptions.

- (A.1) The sequences of  $c_k \sim \mathcal{P}_c$  and  $\tilde{c}_k \sim \mathcal{P}_{\tilde{c}}$  for  $k = 0, 1, \dots$  are i.i.d. random variables with bounded moments of all orders, and  $\{c_k\}$  is independent of  $\{\tilde{c}_k\}$ .
- (A.2) The sequences  $\{\mathbf{a}_k\}$  and  $\{\tilde{\mathbf{a}}_k\}$  for  $k = 0, 1, \dots$  are both independent Gaussian vectors with zero mean and the covariance matrix  $\mathbf{I}_n$ . Moreover,  $\{\mathbf{a}_k\}$ ,  $\{\tilde{\mathbf{a}}_k\}$  are independent of  $\{c_k\}$  and  $\{\tilde{c}_k\}$ .
- (A.3) The first-order derivative of  $H(\cdot)$  and the derivatives up to fourth order of the functions  $F(\hat{D}(\cdot))$  and  $\tilde{F}(\hat{D}(\cdot))$  exist and they are also uniformly bounded.
- (A.4) Let  $[\mathbf{U}, \mathbf{V}_0, \mathbf{w}_0]$  be the initial microscopic state. For  $i = 1, 2, \dots, n$ , we have  $\mathbb{E}[\sum_{\ell=1}^d ([\mathbf{U}]_{i,\ell}^4 + [\mathbf{V}_0]_{i,\ell}^4 + [\mathbf{w}_0]_{i,\ell}^4)] \leq C/n^2$ , where  $C$  is a constant not depending on  $n$ .
- (A.5) The initial macroscopic state  $\mathbf{M}_0$  satisfies  $\mathbb{E}\|\mathbf{M}_0 - \mathbf{M}_0^*\| \leq C/\sqrt{n}$ , where  $\mathbf{M}_0^*$  is a deterministic matrix and  $C$  is a constant not depending on  $n$ .

We provide a few remarks on the above assumptions. In Assumption (A.1),  $\mathcal{P}_c$  and  $\mathcal{P}_{\tilde{c}}$  can be different. For example,  $c$  is Gaussian, and  $\tilde{c}$  is uniform on  $[-1, 1]^d$ . The assumption (A.2) can

be relaxed to non-Gaussian cases as long as all moments of  $\mathbf{a}_k$  and  $\tilde{\mathbf{a}}_k$  are bounded, but we use Gaussian assumption here to simplify the proof. The assumption (A.4) requires that the elements in the parameter matrix of real data  $\mathbf{U}$  and initial microscopic state  $\mathbf{X}_0$  are  $\mathcal{O}(1/\sqrt{n})$  numbers. Intuitively, this assumption ensures that  $\mathbf{U}$  and  $\mathbf{X}_0$  are generic matrices with  $\mathcal{O}(1)$  Frobenius norms (*i.e.*, not the matrices that most elements are zeros and only few elements are large numbers). The assumption (A.5) ensures that the initial macroscopic states converges to a deterministic value as the system size  $n$  goes to infinity. The following theorem proves that if the initial state is convergent, then the whole training process converges to a deterministic process as  $n \rightarrow \infty$ , which is characterized by an ODE.

**Theorem 1.** Fix  $T > 0$ . It holds under Assumptions (A.1)–(A.5) that

$$\max_{0 \leq k \leq nT} \mathbb{E} \|\mathbf{M}_k - \mathbf{M}(\frac{k}{n})\| \leq \frac{C(T)}{\sqrt{n}}, \quad (7)$$

where  $C(T)$  is a constant that depends on  $T$  but not on  $n$ , and  $\mathbf{M}(t) = \begin{bmatrix} \mathbf{I} & \mathbf{P}_t & \mathbf{q}_t \\ \mathbf{P}_t^\top & \mathbf{S}_t & \mathbf{r}_t \\ \mathbf{q}_t^\top & \mathbf{r}_t^\top & z_t \end{bmatrix} \in \mathbb{R}^{(2d+1) \times (2d+1)}$  is a deterministic function. Moreover,  $\mathbf{M}(t)$  is the unique solution of the following ODE:

$$\begin{aligned} \frac{d}{dt} \mathbf{P}_t &= \tilde{\tau}(\mathbf{q}_t \tilde{\mathbf{g}}_t^\top + \mathbf{P}_t \mathbf{L}_t) \\ \frac{d}{dt} \mathbf{q}_t &= \tau(\mathbf{g}_t - \mathbf{P}_t \tilde{\mathbf{g}}_t + \mathbf{q}_t h_t) \\ \frac{d}{dt} \mathbf{r}_t &= \tau(\mathbf{P}_t^\top \mathbf{g}_t - \mathbf{S}_t \tilde{\mathbf{g}}_t + \mathbf{r}_t h_t) + \tilde{\tau}(z_t \tilde{\mathbf{g}}_t + \mathbf{L}_t \mathbf{r}_t) \\ \frac{d}{dt} \mathbf{S}_t &= \tilde{\tau}(\mathbf{r}_t \tilde{\mathbf{g}}_t^\top + \tilde{\mathbf{g}}_t \mathbf{r}_t^\top + \mathbf{S}_t \mathbf{L}_t + \mathbf{L}_t \mathbf{S}_t) \\ \frac{d}{dt} z_t &= 2\tau(\mathbf{q}_t^\top \mathbf{g}_t - \mathbf{r}_t^\top \tilde{\mathbf{g}}_t + z_t h_t) + \tau^2 b_t \end{aligned} \quad (8)$$

with the initial condition  $\mathbf{M}(0) = \mathbf{M}_0^*$ , where

$$\begin{aligned} \mathbf{g}_t &= \langle \mathbf{c} f(\mathbf{c}^\top \mathbf{q}_t + e\sqrt{z_t \eta_T}) \rangle_{\mathbf{c}, e}, \quad \tilde{\mathbf{g}}_t = \langle \tilde{\mathbf{c}} \tilde{f}(\tilde{\mathbf{c}}^\top \mathbf{r}_t + e\sqrt{z_t \eta_G}) \rangle_{\tilde{\mathbf{c}}, e}, \quad \mathbf{L}_t = -\lambda \text{diag}(H'(\mathbf{S}_t)) \\ h_t &= \langle f'(\mathbf{c}^\top \mathbf{q}_t + e\sqrt{z_t \eta_T}) \rangle_{\mathbf{c}, e} - \langle \tilde{f}'(\tilde{\mathbf{c}}^\top \mathbf{r}_t + e\sqrt{z_t \eta_G}) \rangle_{\tilde{\mathbf{c}}, e} - \lambda H'(z_t), \\ b_t &= \eta_T \langle f^2(\mathbf{c}^\top \mathbf{q}_t + e\sqrt{z_t \eta_T}) \rangle_{\mathbf{c}, e} + \eta_G \langle \tilde{f}^2(\tilde{\mathbf{c}}^\top \mathbf{r}_t + e\sqrt{z_t \eta_G}) \rangle_{\tilde{\mathbf{c}}, e}. \end{aligned} \quad (9)$$

The two functions  $f, \tilde{f}$  stand for  $f(x) = \frac{d}{dx} F(\hat{D}(x))$  and  $\tilde{f}(x) = \frac{d}{dx} \tilde{F}(\hat{D}(x))$ , and  $f', \tilde{f}'$  and  $H'$  are derivatives of  $f, \tilde{f}$  and  $H$  respectively. The two constants  $\eta_T$  and  $\eta_G$  are the strength of the noise in the true data model and the generator, respectively. The brackets  $\langle \cdot \rangle_{\mathbf{c}, e}$  and  $\langle \cdot \rangle_{\tilde{\mathbf{c}}, e}$  denote the averages over the random variables  $\mathbf{c} \sim \mathcal{P}_{\mathbf{c}}$ ,  $\tilde{\mathbf{c}} \sim \mathcal{P}_{\tilde{\mathbf{c}}}$ , and  $e \sim \mathcal{N}(0, 1)$ , where  $\mathcal{P}_{\mathbf{c}}$  and  $\mathcal{P}_{\tilde{\mathbf{c}}}$  are the distributions involved in defining the generative model (1) and the generator (2).

This theorem implies that for each  $k = \lfloor tn \rfloor$  for some  $t \in [0, T]$ , the macroscopic state  $\mathbf{M}_k$  converges to a deterministic number  $\mathbf{M}(t)$ , and the convergence rate is  $\mathcal{O}(1/\sqrt{n})$ . The limiting ODE (8) for the macroscopic states involves  $\mathcal{O}(d^2)$  variables, where  $d$  is the number of internal features often assumed to be a finite number that is much less than  $n$ . This ODE is essentially different from the ODE derived in the small-learning-rate limit [11–16], in which the number of variables is  $\mathcal{O}(n)$ .

The complete proof can be found in the Supplementary Materials. We briefly sketch the proof here. First, we note that  $\mathbf{M}_k$  is a discrete-time stochastic process driven by the Markov chain  $\mathbf{X}_k$ . Then, we apply the martingale decomposition for  $\mathbf{M}_k$  and get

$$\mathbf{M}_{k+1} - \mathbf{M}_k = \frac{1}{n} \phi(\mathbf{M}_k) + (\mathbf{M}_{k+1} - \mathbb{E}_k \mathbf{M}_{k+1}) + [\mathbb{E}_k \mathbf{M}_{k+1} - \mathbf{M}_k - \frac{1}{n} \phi(\mathbf{M}_k)],$$

where the matrix-valued function  $\phi(\mathbf{M})$  represents the functions on the right hand sides of the ODE (8), and  $\mathbb{E}_k$  denotes the conditional expectation given the state of the Markov chain  $\mathbf{X}_k$ . Finally, we show the martingale  $\sum_{k'=0}^k (\mathbf{M}_{k'+1} - \mathbb{E}_{k'} \mathbf{M}_{k'})$  and the higher-order term  $\mathbb{E}_k \mathbf{M}_{k+1} - \mathbf{M}_k - \frac{1}{n} \phi(\mathbf{M}_k)$  have no contribution when  $n$  goes to infinity.

Due to the limitation of our current proof, the constant  $C(T)$  in (7) grows exponentially as  $T$  increases. This is not a problem for any finite  $T$ , but may cause some problem to study the long time behavior when  $T \rightarrow \infty$ . However, if we impose a sufficient large regularizer parameter  $\lambda$  to limit the norms of the microscopic weights  $\mathbf{V}_k$  and  $\mathbf{w}_k$ , then the macroscopic state  $\mathbf{M}_k$  is bounded

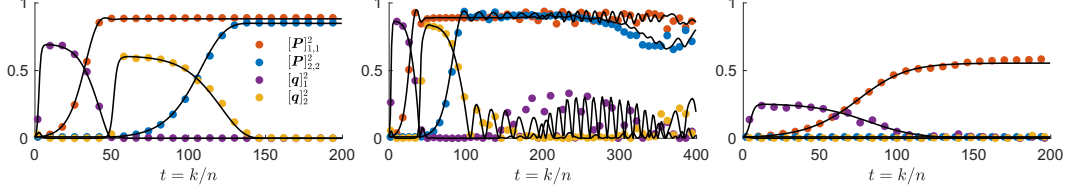


Figure 1: Macroscopic dynamics of the GAN with  $d = 2$  features:  $[P_k]_{i,j}$  is the cosine of the angle between  $i$ 'th column vector of the real feature matrix  $U_k$  and  $j$ 'th column vector of the generator's weight matrix  $V_k$ . Similarly,  $[q_k]_i$  is the cosine of angle between  $i$ 'th column vector of  $U_k$  and the discriminator's weight vector  $w_k$ . Colored dots are results from experiments, and the curves tracing these dots are our theoretical prediction by the ODE (8). From the left to right, the variance of background noise is  $\eta_T = \eta_G = 2, 1, 4$  respectively, and other parameters are the same. The left figure is an example of successful training, where two features (red and blue dots) are retrieved by the generator. The center figure shows an oscillating training. It happens when noise are weak. The right figure shows a mode collapsing state, in which only the first feature are estimated by the generator.

as  $[M_k]_{i,j}^2 \leq [M_k]_{i,i}[M_k]_{j,j}$ . In our experiments,  $\lambda > 1$  is sufficient. In this case, the constant  $C(T)$  is bounded not depending on  $T$ . In Example 1, when  $\lambda \rightarrow \infty$ ,  $[M_k]_{i,i} = 1$ , and therefore  $[M_k]_{i,j}^2 \leq 1$  and  $C(T) \leq (2d+1)^2$ , where the number of features  $d$  is considered a constant not growing with  $n$ . This justifies the fixed points analysis of the ODE as discussed in Section S-I, which reflects the long-time training behavior. A better proof strategy to get rid of this dependence of  $T$  is also possible, *e.g.*, [26].

**Numerical verification.** We verify the theoretical prediction given by the ODE (8) via numerical simulations under the settings stated in Example 1. The results are shown in Figure 1. The number of features is  $d = 2$ , and  $c_k$  and  $\tilde{c}_k$  are both Gaussian with zero mean and covariance  $\text{diag}([5, 3])$ . The dimension is  $n = 5,000$ , and the learning rates of the generator and discriminator are  $\tilde{\tau} = 0.04$  and  $\tau = 0.2$  respectively. After testing different noise strength  $\eta_T = \eta_G = 2, 1, 4$ , we have observed at least three nontrivial dynamical patterns: success, oscillating or mode collapsing. In all these experiments, our theoretical predictions match the actual trajectories of the macroscopic states pretty well.

Let us take a closer look at the successful case as shown in the left figure in Figure 1. The dynamics can be split into 4 stages. At the first stage, the discriminator learns the first feature of the true model. At this state,  $[q_t]_1$  quickly increases. At the second stage, the generator starts to learn the first feature and the discriminator is deceived. At this stage,  $[P_t]_{1,1}^2$  increases and  $[q_t]_1^2$  decreases. Once the discriminator completely forgets the first feature as  $[q_t]_1 \approx 0$ , the third state begins. The discriminator starts to learn the second feature as  $[q_t]_2^2$  increases. Then, at the last stage, the generator learns the second feature and the discriminator is fooled again. In this region,  $[P_t]_{2,2}^2$  increases and  $[q_t]_2^2$  decreases down to 0. Eventually, the generators learns both features and the discriminator is completely fooled. It ends up at a stationary state that  $q_t = 0$  and  $P_t$  is nearly an identity matrix. Interestingly, this experiment shows that the generator learn features sequentially given a single-feature discriminator. This may be a reason why in practice, the discriminator's structure can be much simpler than the generator's.

### 3.2 Microscopic dynamics

In this section, we study how the elements in  $X_k = [U, V_k, w_k]$  evolve during the training process. Instead of studying the trajectory of  $X_k$ , we study the evolution of the *empirical measure* of the microscopic states, which is defined as

$$\mu_k(\hat{u}, \hat{v}, \hat{w}) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \delta([\hat{u}^\top, \hat{v}^\top, \hat{w}] - \sqrt{n}[[U]_{i,:}, [V_k]_{i,:}, [w]_i])$$

where  $\delta(\cdot)$  is a Dirac measure on  $\mathbb{R}^{2d+1}$  and  $[U]_{i,:}, [V_k]_{i,:}$  are  $i$ th row of  $U$  and  $V_k$  respectively. The scaling factor  $\sqrt{n}$  in the Dirac measures is introduced because  $[U]_{i,\ell}, [V_k]_{i,\ell}$  and  $[w_k]_i$  are  $\mathcal{O}(1/\sqrt{n})$  quantities.

We next embed the discrete-time measure-valued stochastic process  $\mu_k$  into a continuous-time process by defining  $\mu_t^{(n)} \stackrel{\text{def}}{=} \mu_k(\hat{\mathbf{u}}, \hat{\mathbf{v}}, \hat{w})$  with  $k = \lfloor nt \rfloor$ . Following the general technical approach presented in [17], we can show that under the same assumptions as Theorem 1, given  $T > 0$ , the sequence of measure-valued process  $\{\{\mu_t^{(n)}\}_{t \in [0, T]}\}_n$  converges weakly to a deterministic process  $\{\mu_t\}_{t \in [0, T]}$ . In addition,  $\mu_t$  is the measure of the solution to the stochastic differential equation

$$\begin{aligned} d\hat{\mathbf{u}}_t &= 0 \\ d\hat{\mathbf{v}}_t &= \tilde{\tau}(\hat{w}_t \tilde{\mathbf{g}}_t + \mathbf{L}_t \hat{\mathbf{v}}_t) dt \\ d\hat{w}_t &= \tau(\hat{\mathbf{u}}_t^\top \mathbf{g}_t + \hat{\mathbf{v}}_t^\top \tilde{\mathbf{g}}_t + \hat{w}_t h_t) dt + \tau \sqrt{b_t} dB_t \end{aligned} \quad (10)$$

where  $(\hat{\mathbf{u}}_0, \hat{\mathbf{v}}_0, \hat{w}_0) \sim \mu_0$ ;  $B_t$  is the standard Brownian motion. The functions  $\mathbf{g}_t, \tilde{\mathbf{g}}_t, \mathbf{L}_t, h_t$  and  $b_t$  are defined in (9), in which the macroscopic quantities  $\mathbf{P}_t, \mathbf{S}_t, \mathbf{q}_t, z_t, \mathbf{r}_t$  are computed as follows

$$\mathbf{P}_t = \langle \mu_t, \hat{\mathbf{u}} \hat{\mathbf{v}}^\top \rangle, \quad \mathbf{S}_t = \langle \mu_t, \hat{\mathbf{v}} \hat{\mathbf{v}}^\top \rangle, \quad \mathbf{q}_t = \langle \mu_t, \hat{\mathbf{u}} \hat{w} \rangle, \quad z_t = \langle \mu_t, \hat{w}^2 \rangle, \quad \mathbf{r}_t = \langle \mu_t, \hat{\mathbf{v}} \hat{w} \rangle, \quad (11)$$

where  $\langle \mu_t, \cdot \rangle$  denotes the expectation with respect to the measure  $\mu_t$ .

The SDE (10) shows the intuitive meaning of the functions defined in (9):  $\mathbf{g}_t, \tilde{\mathbf{g}}_t, \mathbf{L}_t, h_t$  are drift coefficients of the SDE and  $b_t$  is the diffusion coefficient of the SDE. We also note that if one follows the analysis in the small-learning-rate limit [11–16], one will get an ODE for the microscopic states. Compared to our SDE formula, the diffusion term  $\tau \sqrt{b_t} dB_t$  is missing in those works, and therefore the effect of the noise can not be analyzed.

Moreover, the deterministic measure  $\mu_t$  is unique solution of the following PDE (given in its weak form): for any bounded smooth test function  $\varphi(\hat{\mathbf{u}}, \hat{\mathbf{v}}, \hat{w})$ ,

$$\begin{aligned} \frac{d}{dt} \langle \mu_t, \varphi(\hat{\mathbf{u}}, \hat{\mathbf{v}}, \hat{w}) \rangle &= \\ \tilde{\tau} \langle \mu_t, (\hat{w} \tilde{\mathbf{g}}_t^\top + \hat{\mathbf{v}}^\top \mathbf{L}_t) \nabla_{\hat{\mathbf{v}}} \varphi \rangle &+ \tau \langle \mu_t, (\hat{\mathbf{u}}^\top \mathbf{g}_t - \hat{\mathbf{v}}^\top \tilde{\mathbf{g}}_t + h_t \hat{w}) \frac{\partial}{\partial \hat{w}} \varphi \rangle + \frac{\tau^2}{2} b_t \langle \mu_t, \frac{\partial^2}{\partial \hat{w}^2} \varphi \rangle \end{aligned} \quad (12)$$

where  $\mathbf{q}_t, \mathbf{r}_t, \mathbf{S}_t$ , and  $z_t$  are defined in (11), and the functions  $\mathbf{g}_t, \tilde{\mathbf{g}}_t, b_t, h_t$  and  $\mathbf{L}_t$  are defined in (9). We refer readers to [17] for a general framework for rigorously establishing the above scaling limit.

The connection between the microscopic and macroscopic dynamics can also be derived from the weak formulation of the PDE. Let  $\varphi$  being each element of  $\hat{\mathbf{u}} \hat{\mathbf{v}}^\top, \hat{\mathbf{u}} \hat{w}, \hat{\mathbf{v}} \hat{w}, \hat{\mathbf{v}} \hat{\mathbf{v}}^\top, \hat{w}^2$ , and substituting those  $\varphi$  into the PDE (12), we can derive the ODE (8). In the setting of this paper, the macroscopic dynamics enjoys a closed ODE: We can predict the macroscopic states without solving the PDE nor SDE at microscopic scale. However, in a more general setting, e.g. when we add a regularizer other than the L2 type, the ODE itself may not be closed. In that case, one has to solve the PDE directly.

**Numerical verification.** We verify the predictions given by the PDE (12) by setting  $d = 1$  using a special choice of the  $(n \times 1)$ -dimensional target feature matrix  $\mathbf{U}$  whose elements are all  $1/\sqrt{n}$  with  $n = 10,000$ . We also set the initial condition  $\mu_0(\hat{\mathbf{v}}, \hat{w} | \hat{\mathbf{u}} = 1)$  to be a Gaussian distribution. (When  $d = 1$ , the macroscopic quantities  $\mathbf{P}_t, \mathbf{q}_t, \mathbf{r}_t, \mathbf{S}_t$  reduce to scalars, so we remove their boldface here.) In this case, the PDE (12) admits a particularly simple analytical solution: at any time  $t$ , the solution  $\mu_t(\hat{\mathbf{v}}, \hat{w} | \hat{\mathbf{u}} = 1)$  is a Gaussian distribution whose mean and covariance matrix are given by  $\mathbb{E}_{\mu_t(\hat{\mathbf{v}}, \hat{w} | \hat{\mathbf{u}} = 1)} \begin{bmatrix} \hat{\mathbf{v}} \\ \hat{w} \end{bmatrix} = \begin{bmatrix} P_t \\ q_t \end{bmatrix}, \mathbb{E}_{\mu_t(\hat{\mathbf{v}}, \hat{w} | \hat{\mathbf{u}} = 1)} \begin{bmatrix} \hat{\mathbf{v}} \\ \hat{w} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{v}} & \hat{w} \end{bmatrix} = \begin{bmatrix} S_t & r_t \\ r_t & z_t \end{bmatrix}$ . Figure 2 overlays the contours of the probability distribution  $\mu_t(\hat{\mathbf{v}}, \hat{w} | \hat{\mathbf{u}} = 1)$  at different times  $t$  over the point clouds of the actual experiment data  $(\sqrt{n}[\mathbf{w}_k]_i, \sqrt{n}[\mathbf{V}_k]_{i,1})$ . We can see that the theoretical prediction given by (12) has excellent agreement with simulation results.

## 4 Local Stability Analysis of the ODE for the Macroscopic States

In this section, we study how the parameters, such as the learning rates  $\tau$  and  $\tilde{\tau}$ , noise strength  $\eta_G$  and  $\eta_T$  affect the training algorithm. We will focus on the concrete model as described in Example 1 so that we can have analytical solutions.

In order to further reduce the degrees of freedom of the ODE (8), we let the regularization parameter  $\lambda \rightarrow \infty$ . In this case, the vector  $\mathbf{w}_k$  and all columns vectors of  $\mathbf{V}_k$  are always normalized. Thus

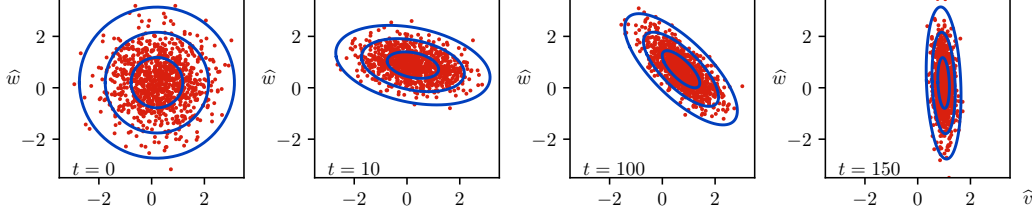


Figure 2: The evolution of the microscopic states at  $t = 0, 10, 100$ , and  $150$ . For each fixed  $t$ , the red points in the corresponding figure represent the values of  $(\hat{v}, \hat{w}) = (\sqrt{n}[\mathbf{V}_k]_{i,1}, \sqrt{n}[\mathbf{w}_k]_i)$  for  $i = 1, 2, \dots, n$ , where  $k = \lfloor nt \rfloor$ . The blue ellipses illustrate the contours corresponding to one, two, and three standard deviations of the 2-D Gaussian distribution predicted by the PDE (12).

$z_k = 1$  and  $[\mathbf{S}]_{i,i} = 1$ . The macroscopic state is then described by  $\mathbf{P}_k, \mathbf{q}_k, \mathbf{r}_k$  and off-diagonal terms of  $\mathbf{S}_k$ . Correspondingly, the ODE in Theorem 1 reduces to

$$\begin{cases} \frac{d}{dt} \mathbf{P}_t &= \tilde{\tau} (\mathbf{q}_t \mathbf{r}_t^\top \tilde{\Lambda} + \mathbf{P}_t \mathbf{L}_t) \\ \frac{d}{dt} \mathbf{q}_t &= \tau (\Lambda \mathbf{q}_t - \mathbf{P}_t \tilde{\Lambda} \mathbf{r}_t + h_t \mathbf{q}_t) \\ \frac{d}{dt} \mathbf{r}_t &= \tau (\mathbf{P}_t^\top \Lambda \mathbf{q}_t - \mathbf{S}_t \tilde{\Lambda} \mathbf{r}_t + h_t \mathbf{r}_t) + \tilde{\tau} (\tilde{\Lambda} + \mathbf{L}_t) \mathbf{r}_t \\ \frac{d}{dt} \mathbf{S}_t &= \tilde{\tau} (\mathbf{r}_t \mathbf{r}_t^\top \tilde{\Lambda} + \tilde{\Lambda} \mathbf{r}_t \mathbf{r}_t^\top + \mathbf{S}_t \mathbf{L}_t + \mathbf{L}_t \mathbf{S}_t) \end{cases} \quad (13)$$

where  $\Lambda$  and  $\tilde{\Lambda}$  are the covariance matrices of the distributions  $P_c$  and  $P_{\tilde{c}}$ , respectively; and

$$h_t = (1 - \frac{\tau \eta_G}{2}) \mathbf{r}_t^\top \tilde{\Lambda} \mathbf{r}_t - (1 + \frac{\tau \eta_T}{2}) \mathbf{q}_t^\top \Lambda \mathbf{q}_t - \tau \frac{\eta_G^2 + \eta_T^2}{2}, \quad \mathbf{L}_t = -\text{diag}(\mathbf{r}_t \mathbf{r}_t^\top \tilde{\Lambda}), \quad (14)$$

in which  $\eta_T$  and  $\eta_G$  are the variance of noise in the true data model and generator, respectively. The derivation from the ODE (8) to (13) is presented in the Supplementary Materials.

Next, we discuss under what conditions, the GAN can reach a desirable training state by studying local stability of a particular type of fixed points of the ODE (13). The perfect estimation of the generator corresponds to  $\mathbf{P}_t$  being an identity matrix (up to a permutation of rows and columns). A complete fail state relates to  $\mathbf{P} = \mathbf{0}$ . Furthermore, It is easy to verify that if  $\mathbf{q}_t = \mathbf{r}_t = \mathbf{0}$ , the ODE (13) will be stable for any  $\mathbf{P}_t = \mathbf{P}$ .

**Claim 1.** *The macroscopic states  $\mathbf{P}_t, \mathbf{q} = \mathbf{r} = \mathbf{0}$  for all valid  $\mathbf{P}_t$  are always the fixed points of the ODE (13). Furthermore, a sufficient condition that the perfect estimation state  $\mathbf{P}_t = \mathbf{I}, \mathbf{q} = \mathbf{r} = \mathbf{0}$  is locally stable and the failed state  $\mathbf{P}_t = \mathbf{0}, \mathbf{q} = \mathbf{r} = \mathbf{0}$  is unstable if*

$$\frac{1}{2} \max_{\ell} \{ \Lambda_{\ell} - \tilde{\Lambda}_{\ell} + \alpha \tilde{\Lambda}_{\ell} \} \leq \tau \bar{\eta}^2 < \min_{\ell} \Lambda_{\ell}, \quad (15)$$

where  $\alpha = \frac{\tilde{\tau}}{\tau}$ ,  $\bar{\eta}^2 = \frac{1}{2}(\eta_T^2 + \eta_G^2)$ , and  $\Lambda_{\ell} = [\Lambda]_{\ell,\ell}$ ,  $\tilde{\Lambda}_{\ell} = [\tilde{\Lambda}]_{\ell,\ell}$ .

The proof can be found in the Supplementary Materials. If the right inequality in (15) is violated, any feature  $\ell$  with the signal-to-noise ratio  $[\Lambda]_{\ell,\ell} < \tau \bar{\eta}^2$  is not learned by the generator resulting *mode collapsing*. The right figure in Figure 1 demonstrates this situations, where only one of the two features is recovered. If the left inequality in (15) is violated, the training processes can be trapped in an *oscillation phase*. This phenomenon is shown in the middle figure in Figure 1. This result indicates that proper background noise can help to avoid oscillation and stabilize the training process. In fact, the trick of injecting additional noise has been used in practice to train multi-layer GANs [27]. To our best knowledge, our paper is the first theoretical study on why noise can have such a positive effect via a dynamic perspective.

In experiments, the training is not ended at the perfect recovery point due to the presence of the noise but converges at another fixed point nearby. This is because the perfect state is marginally stable, as the Jacobian matrix always has zero eigenvalues. It indicates that there are other locally stable fixed points near  $\mathbf{P} = \mathbf{I}$ . In fact, all points in the hyper-rectangle region satisfying  $\mathbf{q} = \mathbf{r} = \mathbf{0}$  and  $|p_{\ell}^*| \leq |[\mathbf{P}]_{\ell,\ell}| \leq 1, \forall \ell = 1, 2, \dots, d$  are locally stable for some critical  $p_{\ell}^*$ . In the matched case when  $\Lambda_{\ell} = \tilde{\Lambda}_{\ell}$ , we have  $p_{\ell}^* = [(\Lambda_{\ell} - \tau \bar{\eta}^2)(\tilde{\Lambda}_{\ell} + \tau \bar{\eta}^2 - \alpha \tilde{\Lambda}_{\ell}) / (\Lambda_{\ell} \tilde{\Lambda}_{\ell})]^{1/2}$ ,  $\alpha = \frac{\tilde{\tau}}{\tau}$  and  $\bar{\eta}^2 =$



$\frac{1}{2}(\eta_T^2 + \eta_G^2)$ . Starting from a point near the origin, numerical solution of the ODE shows the training processes are ended up at the corner of this hyper-rectangle, *i.e.*,  $\mathbf{P}^* = \text{diag}(\{p_\ell^*, \ell = 1, 2, \dots, d\})$ . In the small-learning rate limit  $\tau \rightarrow 0$  and the learning rate ratio  $\alpha \rightarrow 0$ , we get the perfect recovery  $\mathbf{P}^* = \mathbf{I}$ . The limit  $\tau \rightarrow 0, \alpha \rightarrow 0$  was studied in the small-learning-rate analysis with the two-time scaling [15], and the result is consistent, but our analysis includes the situations with finite  $\tau$  and  $\alpha$ .

In addition, we provide a phase diagram analysis in a single-feature case  $d = 1$  in the Supplementary Materials. All possible fixed points in this case are enumerated and their local stability is analyzed. This helps us understand the successful recovery condition (15), which is the intersection of the informative phases that each feature can be recovered individually.

## 5 Conclusion

We present a simple high-dimensional model for GAN with an exactly analyzable training process. Using the tool of scaling limits of stochastic processes, we show that the macroscopic state associated with the training process converges to a deterministic process characterized as the unique solution of an ODE, whereas the microscopic state remains stochastic described by an SDE, whose time-varying probability measure is described by a limiting PDE.

Indeed, it is a common picture in statistical physics that the macroscopic states of large systems tend to converge to deterministic values due to self-averaging. These notions, especially the mean-field dynamics, have been applied to analyzing neural networks both in shallow [19, 20] and deep models [28]. However, this mean-field regime was not considered in previous analyses of GAN. For example, a series of recent works *e.g.*, [11–16] considers a different scaling regime where the learning rate goes to zero but the system dimension  $n$  stays fixed. In that regime, the microscopic dynamics are deterministic even with the presence of the microscopic noise. In contrast, we study the regime where the learning rate is fixed but the dimension  $n \rightarrow \infty$ . This setting allows us to quantify the effect of training noise in the learning dynamics.

In this paper, we only consider a linear generator with a latent variable  $\tilde{\mathbf{c}}$  drawn from a fixed distribution  $\mathcal{P}_{\tilde{\mathbf{c}}}$ , but our analysis can be extended to a more complex non-linear model with a learnable latent-variable distribution. Specifically, in order to compute derivatives w.r.t.  $\mathcal{P}_{\tilde{\mathbf{c}}}$ , the latent variable  $\tilde{\mathbf{c}} \sim \mathcal{P}_{\tilde{\mathbf{c}}}$  should be reparameterized by a deterministic function  $\tilde{\mathbf{c}} = f(\mathbf{z}; \boldsymbol{\theta})$ , where  $\boldsymbol{\theta}$  is a learnable parameter and  $\mathbf{z}$  is a random variable drawn from a simple and fixed distribution. For example, a Gaussian mixture with  $L$  equal-probability modes can be parameterized by  $\tilde{\mathbf{c}} = \sum_{\ell=1}^L (\boldsymbol{\mu}_\ell + \boldsymbol{\Sigma}_\ell \boldsymbol{\epsilon}_\ell) \beta_\ell$ , where  $\boldsymbol{\mu}_\ell$  and  $\boldsymbol{\Sigma}_\ell$  are two learnable parameters representing the mean and covariance of the  $\ell$ th mode respectively, and  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$ ;  $\beta_\ell$  is a random indicator variable where only one  $\beta_\ell$  for  $\ell = 1, 2, \dots, L$  is 1 and the others are 0. In practice,  $f(\mathbf{z}; \boldsymbol{\theta})$  is implemented by a multilayer neural network. Our analysis can be naturally extended to analyzing this model as long as the dimensions of  $\tilde{\mathbf{c}}$  and  $\boldsymbol{\theta}$  keep finite when the data dimension  $n$  goes to infinity. More challenging situations, where the dimension of  $\boldsymbol{\theta}$  is proportional to  $n$ , will be explored in future works.

Although our analysis is carried out in the asymptotic setting, numerical experiments show that our theoretical predictions can accurately capture the actual performance of the training algorithm at moderate dimensions. Our analysis also reveals several different phases of the training process that highly depend on the choice of the learning rates and noise strength. The analysis reveals a condition on the learning rates and the strength of noise to have successful training. Violating this condition results either oscillation or mode collapsing. Despite its simplicity, the proposed model of GAN provides a new perspective and some insights for the study of more realistic models and more involved training algorithms.

**Acknowledgments** This work was supported by the US Army Research Office under contract W911NF-16-1-0265 and by the US National Science Foundation under grants CCF-1319140, CCF-1718698, and CCF-1910410.

## References

- [1] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative Adversarial Nets,” in *Advances in Neural Information Processing System*, 2014, pp. 2672–2680.
- [2] M. Arjovsky, S. Chintala, and L. Bottou, “Wasserstein Generative Adversarial Networks,” *Proceedings of The 34th International Conference on Machine Learning*, pp. 1–32, 2017.
- [3] M. Lucic, K. Kurach, M. Michalski, S. Gelly, and O. Bousquet, “Are gans created equal? a large-scale study,” in *Advances in neural information processing systems*, 2018, pp. 698–707.
- [4] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi, “Photo-Realistic Single Image Super-Resolution Using a Generative Adversarial Network,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4681–4690.
- [5] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, “Image-to-Image Translation with Conditional Adversarial Networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1125–1134.
- [6] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee, “Generative Adversarial Text to Image Synthesis,” *33rd International Conference on Machine Learning*, pp. 1060–1069, 2016.
- [7] S. Arora, R. Ge, Y. Liang, and Y. Zhang, “Generalization and Equilibrium in Generative Adversarial Nets,” in *International Conference on Machine Learning*, 2017, pp. 224–232.
- [8] M. Arjovsky and L. Bottou, “Towards Principled Methods for Training Generative Adversarial Networks,” *arXiv preprint arXiv:1701.04862*, 2017.
- [9] S. Feizi, C. Suh, F. Xia, and D. Tse, “Understanding GANs: the LQG Setting,” *arXiv:1710.10793*, 2017.
- [10] J. Li, A. Madry, J. Peebles, and L. Schmidt, “Towards Understanding the Dynamics of Generative Adversarial Networks,” *arXiv preprint arXiv:1706.09884*, 2017.
- [11] L. Mescheder, A. Geiger, and S. Nowozin, “Which training methods for gans do actually converge?” in *International Conference on Machine Learning*, 2018, pp. 3478–3487.
- [12] L. Mescheder, S. Nowozin, and A. Geiger, “The numerics of GANs,” in *Advances in Neural Information Processing Systems*, 2017, pp. 1823–1833.
- [13] V. Nagarajan and J. Z. Kolter, “Gradient descent GAN optimization is locally stable,” in *Advances in Neural Information and Processing Systems*, 2017, pp. 5591–5600.
- [14] K. Roth, A. Lucchi, S. Nowozin, and T. Hofmann, “Stabilizing Training of Generative Adversarial Networks through Regularization,” in *Advances in Neural Information Processing Systems*, 2017, pp. 2015–2025.
- [15] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, “GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium,” in *Advances in Neural Information Processing Systems*, 2017, pp. 6629–6640.
- [16] E. V. Mazumdar, M. I. Jordan, and S. S. Sastry, “On finding local nash equilibria (and only local nash equilibria) in zero-sum games,” *arXiv preprint arXiv:1901.00838*, 2019.
- [17] C. Wang, J. Mattingly, and Y. M. Lu, “Scaling Limit: Exact and Tractable Analysis of Online Learning Algorithms with Applications to Regularized Regression and PCA,” *arXiv preprint arXiv:1712.04332*, 2017.
- [18] G. O. Roberts, A. Gelman, and W. R. Gilks, “Weak convergence and optimal scaling of random walk Metropolis algorithms,” *Annals of Applied Probability*, vol. 7, no. 1, pp. 110–120, 1997.
- [19] D. Saad and S. A. Solla, “Exact Solution for On-Line Learning in Multilayer Neural Networks,” *Phys. Rev. Lett.*, vol. 74, no. 21, pp. 4337–4340, 1995.
- [20] M. Biehl and H. Schwarze, “Learning by on-line gradient descent,” *Journal of Physics A*, vol. 28, no. 3, pp. 643–656, 1995.
- [21] C. Wang and Y. M. Lu, “Online Learning for Sparse PCA in High Dimensions: Exact Dynamics and Phase Transitions,” in *Information Theory Workshop (ITW), 2016 IEEE*, 2016, pp. 186–190.
- [22] C. Wang, Y. C. Eldar, and Y. M. Lu, “Subspace estimation from incomplete observations: A high-dimensional analysis,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 12, no. 6, pp. 1240–1252, Dec 2018.
- [23] C. Wang and Y. M. Lu, “The Scaling Limit of High-Dimensional Online Independent Component Analysis,” in *Advances in Neural Information Processing Systems*, 2017, pp. 6641–6650.
- [24] S. Mei, A. Montanari, and P.-M. Nguyen, “A Mean Field View of the Landscape of Two-Layers Neural Networks,” *arXiv preprint*, p. arXiv:1804.06561, 2018.
- [25] I. Johnstone and A. Lu, “On consistency and sparsity for principal components analysis in high dimensions,” *Journal of the American Statistical Association*, vol. 104, no. 486, pp. 682–693, 2009.
- [26] B. Jourdain, T. Lelièvre, and B. Miasojedow, “Optimal scaling for the transient phase of Metropolis Hastings algorithms: The longtime behavior,” *Bernoulli*, vol. 20, no. 4, pp. 1930–1978, 2014.
- [27] C. K. Sønderby, J. Caballero, L. Theis, W. Shi, and F. Huszár, “Amortised map inference for image super-resolution,” *International Conference on Learning Representations*, 2017.
- [28] P.-M. Nguyen, “Mean field limit of the learning dynamics of multilayer neural networks,” *arXiv preprint arXiv:1902.02880*, 2019.
- [29] P. Billingsley, *Convergence of probability measures*. John Wiley & Sons, 2013.
- [30] H. Kushner and G. G. Yin, *Stochastic approximation and recursive algorithms and applications*. Springer Science & Business Media, 2003, vol. 35.

# Supplementary Materials

These Supplementary Materials provide additional information, detailed derivations and proof of the results shown in the main text. Specifically, in Section S-I we provide a local stability analysis and draw the phase diagram in the case  $d = 1$  and  $d = 2$ . In Section S-II, we present a heuristic derivation of the stochastic differential equation (SDE) for the microscopic states. Next, in Section S-III, we show a derivation of the ODE for the macroscopic states from the weak formulation of the PDE. We then establish the full proof of the Theorem 1 in Section S-IV. Finally, we present the local stability analysis of the ODE's fixed points in Section S-V.

*Notation:* Throughout the paper, we use  $\mathbf{I}_d$  to denote the  $d \times d$  identity matrix. Depending on the context,  $\|\cdot\|$  denotes either the  $\ell_2$  norm of a vector or the spectral norm of a matrix. For any  $x \in \mathbb{R}$ , the floor operation  $\lfloor x \rfloor$  gives the largest integer that is smaller than or equal to  $x$ . We denote  $[v]_i$  the  $i$ th element of the vector  $v$  and denote  $[M]_{i,j}$  the element at  $i$ th row and  $j$ th column of the matrix  $M$ . Finally,  $C(T)$  denotes a constant that depends on the terminal time  $T$ , and  $C$  denotes a general constant that does not depend on  $T$  and  $n$ . Both  $C$  and  $C(T)$  can vary line to line.

## S-I Phase diagram for the case $d = 1$ and $d = 2$

In what follows, we provide a thorough study of all the fixed points of the ODE (13) when the number of feature  $d = 1$  and  $d = 2$ . In particular, three major phases are identified under different settings of the learning rates  $\tau$  and  $\tilde{\tau}$  with the fixed model parameters  $\eta_T, \eta_G, \Lambda$ , and  $\tilde{\Lambda}$ .

**Phase diagram for  $d = 1$ .** By analyzing the local stabilities of these fixed points as illustrated in Figure 3(a), we obtain the phase diagram as shown in Figure 3(b). For simplicity, we only present the result when  $\eta_T = \eta_G = 1$ , and  $\Lambda = \tilde{\Lambda}$ , which is denoted by  $\Lambda$  used in the remaining part of this section. Detailed derivations are presented in S-V.

Even in this simplest case, we find there are in total 5 types of fixed points, the locations of which are visualized in the 3-dimensional space  $(P, q, r)$  shown in Figure 3(a). Each type of the fixed points has an intuitive meaning in terms of the two-player game between  $\mathcal{G}$  and  $\mathcal{D}$ . We list the detailed information in Table 1, in which we define a function  $\beta(\tau) = \begin{cases} [1 + (\frac{\Lambda}{2} - \frac{\Lambda}{\tau})^{-1}]^{-1}, & \text{if } \tau \leq \frac{2\Lambda}{\Lambda+2} \\ +\infty, & \text{otherwise} \end{cases}$ .

*Noninformative phase:* We say that the ODE (13) is in a noninformative phase if either a type-1 or type-2 fixed point in Table 1 is stable. In this case,  $P = 0$ , which indicates that the generator's parameter vector  $\mathbf{V}$  has no correlation with the true feature vector  $\mathbf{U}$ . In Figure 3(b), the region labeled as noninfo-1 is the stable region for the type-1 fixed point, and noninfo-2 is the stable region for the type-2 fixed point. The two regions have no overlap. However, we note that in noninfo-1, the type-3 fixed points can also be stable, in which case the stationary point of the ODE is determined by the initial condition.

*Informative phase:* We say that the ODE (13) is in an informative phase if neither type-1 nor type-2 fixed point is stable, and if at least one fixed point of type-3 and type-5 is stable. In this case, it is guaranteed that  $P$  is nonzero, indicating that the generator can achieve non-vanishing correlation with the real feature vector. In addition, the stable regions for the type-3 and type-5 fixed points are disjoint. They are shown in Figure 3(b) as info-1 and info-2, respectively. The difference between the two region is that, in info-1,  $q$  is exactly 0 indicating that the discriminator is completely fooled, whereas in info-2,  $q$  is nonzero.

*Oscillating phase:* We say that the ODE (13) is in an oscillating phase if none of the fixed points in Table 1 is stable. In this phase, limiting cycles emerge and the system will oscillate on these cycles indefinitely. Moreover, we found two types of limiting cycles.

To further illustrate the phase transitions, we draw ODE trajectories and phase portraits in Figure 4 corresponding to different choices of the step sizes (from left to right,  $\tilde{\tau} = 0.03, 0.2, 0.4, 0.47$ ).

The two figures in the first column of Figure 4 show a case in the Info-1 phase. The bottom red dot in Figure 3(b) represents this configuration of the step sizes, where  $\tilde{\tau}/\tau$  is small. The top figure of Figure 4(a) shows the dynamics of  $P_t, q_t$  and  $r_t$ , and the bottom figure shows the phase portrait on

Table 1: List of the fixed points of the ODE (13) when  $d = 1$  and  $\Lambda = \tilde{\Lambda}$ .

Type	Location	Existence	Stable Region	Intuitive Interpretation
1	$P = q = 0,$ $r = 0$	always	$\tau > \Lambda^2, \frac{\tilde{\tau}}{\tau} < \frac{\tau + \Lambda}{\Lambda}$	Both $\mathcal{G}$ and $\mathcal{D}$ fail, and they are uncorrelated
2	$P = q = 0$ $r = \pm r^* \neq 0$	$\frac{\tilde{\tau}}{\tau} \geq \frac{\tau + \Lambda}{\Lambda}$ or $\frac{\tilde{\tau}}{\tau} \leq 1 - \frac{\tau}{2}$	$\max\{2, \frac{\tau + \Lambda}{\Lambda}\} \leq \frac{\tilde{\tau}}{\tau} \leq \frac{\tau}{\beta(\tau)}$	Both $\mathcal{G}$ and $\mathcal{D}$ fail, and they are correlated
3	$q = r = 0$ $ P  \in (0, 1]$	always	$ P  = 1$ is stable if $\frac{\tilde{\tau}}{\tau} \leq \min\{\frac{2\tau}{\Lambda}, \max\{\frac{\tau^2 \Lambda^{-1}}{ \tau - \Lambda }, 4\}\}$	$\mathcal{G}$ wins and $\mathcal{D}$ loses
4	$P = r = 0$ $q = \pm q^* \neq 0$	always	always unstable	$\mathcal{G}$ loses and $\mathcal{D}$ wins
5	None of $P, q$ or $r$ is zero	not always, at most 8 fixed points	can be computed numerically	Both $\mathcal{G}$ and $\mathcal{D}$ are informative

$P - q$  plane. Top figure of Figure 4.(a) shows an interesting phenomenon that dynamics are separated into two stages. At the first stage,  $q_t$  (red dots, cosine similarity between the true feature vector and discriminator's estimation) increases drastically from 0 to some value near 1, while  $P_t$  (blue dots, cosine similarity between the true feature vector and generator's estimation) almost doesn't change. Intuitively, at this stage, the discriminator learns the true model while the generator is unchanged. In the second stage, the generator start to fool the discriminator, where  $|P_t|$  increases and  $q_t$  decreases. In fact, these two-stage dynamics can be understood from the ODE (13): When  $\tau/\tau$  is small, the process can be decomposed into two processes in different time scales. In particular, the discriminator is associated with the faster dynamics as  $\tau \gg \tilde{\tau}$ , and the generator governs the slower dynamics. Figure 1 in the main text shows that this picture is still hold for multi-feature cases in the hierarchical dynamics.

The figures in the middle two columns of Figure 4 show the two types of limiting cycles that can emerge in the oscillating phase. The middle two red dots in Figure 3.(b) represents these configurations of the step sizes. The last column of Figure 4 shows another stable phase in Info-2. In this phase,  $\tau/\tau$  is relatively large. The two time-scale dynamics are mixed, and another type of stable fixed points emerges.

**Phase diagram for  $d = 2$ .** Figure 5 shows the phase diagram when  $d = 2$ . In particular, the two red lines between Info-1 and Noninfo-1 in Figure 5 are determined by the left inequality in (15). In Info-1, both feature vectors are recovered by the generator. The dynamics of this phase are shown in Figure 1.(a) in the main text. In the Half-info phase, only the feature vector with the larger signal-to-noise ratio is recovered. The dynamics of this phase are shown in Figure 1.(c) in the main text. The blue line between Info-1 and oscillating phases shows the boundary between oscillation state and stable state.

## S-II Heuristic derivations of the dynamics of the microscopic states

In this section, we derive the stochastic differential equations (10) in the main text for the microscopic states in a non-rigorous way. Specifically, we directly discard higher-order terms without any justification, in order to highlight the main ideas. In Section S-IV, we rigorously justify these steps by providing bounds on those terms.

Our starting point is the iterative algorithm (5) in the main text. Substituting the objective function  $\mathcal{L}$  defined in (4) into (5), we have

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \frac{\tau}{n} [\mathbf{y}_k f(\mathbf{y}_k^\top \mathbf{w}_k) - \tilde{\mathbf{y}}_{2k} \tilde{f}(\tilde{\mathbf{y}}_{2k}^\top \mathbf{w}_k) - \lambda \mathbf{w}_k H'(\mathbf{w}_k^\top \mathbf{w}_k)] \quad (\text{S-1})$$

$$\mathbf{V}_{k+1} = \mathbf{V}_k + \frac{\tilde{\tau}}{n} [\mathbf{w}_k \tilde{\mathbf{c}}_{2k+1}^\top \tilde{f}(\tilde{\mathbf{y}}_{2k+1}^\top \mathbf{w}_k) - \lambda \mathbf{V}_k \text{diag}(H'(\mathbf{V}_k^\top \mathbf{V}_k))], \quad (\text{S-2})$$

where  $\mathbf{y}_k$  and  $\tilde{\mathbf{y}}_k$  are true and fake samples generated according to (1) and (2) respectively. The two functions  $f, \tilde{f}$  stand for  $f(x) = \frac{d}{dx} F(\hat{D}(x))$  and  $\tilde{f}(x) = \frac{d}{dx} \tilde{F}(\hat{D}(x))$ . The function  $H'$  is derivative of  $H$ . If the input of  $H'(\cdot)$  is a matrix,  $H'$  applies to the input matrix element-wisely. The operation  $\text{diag}(\mathbf{A})$  is a diagonal matrix of  $\mathbf{A}$ , where the off-diagonal term are set to zero.

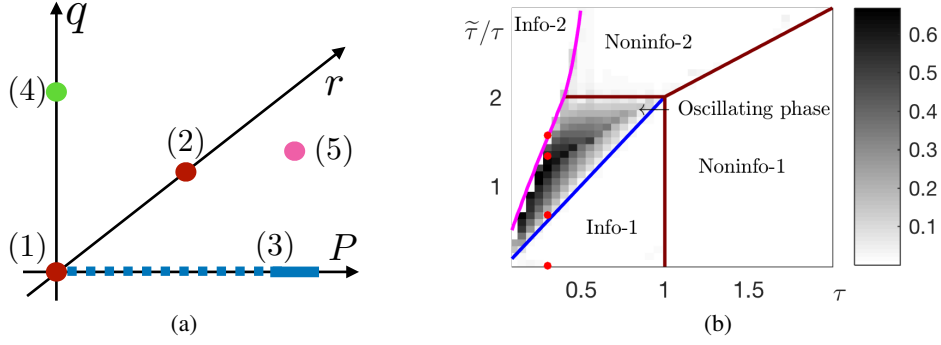


Figure 3: (a): The locations of the five types of fixed points of the ODE (13). Their properties are listed in Table 1. (b): The phase diagram for the stationary state of the ODE (13). The colored lines illustrate the theoretical prediction of the boundaries between the different phases. Simulations results for a single numerical experiment are also shown to illustrate the oscillating phase: Each grey square represents the value of  $\frac{1}{200} \int_{800}^{1000} [(P_t - \langle P_t \rangle)^2 + (q_t - \langle q_t \rangle)^2 + (r_t - \langle r_t \rangle)^2] dt$  where  $\langle P_t \rangle = \frac{1}{200} \int_{800}^{1000} P_t dt$ , and  $\langle q_t \rangle$  and  $\langle r_t \rangle$  are defined similarly. Note that the above quantity measures the variation (over time) of the training process as it approaches steady states. We see that the variation is indeed nonzero in the oscillating phase (see Figure 4), whereas the variation is close to zero in all other phases.

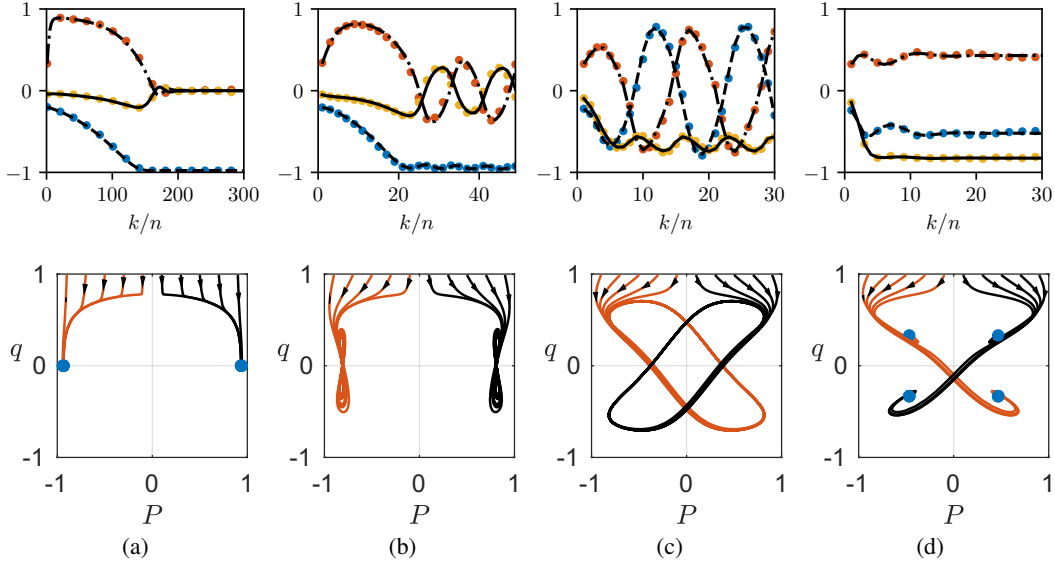


Figure 4: Macroscopic dynamics of Example 1 with  $d = 1$ . In the first row, the red, blue and yellow dots represent  $P_t$ ,  $q_t$ , and  $r_t$  respectively of the experimental results of a single trial. The black curves under the dots are theoretical predictions given by the ODE (13). We set a fix the discriminator's learning rate  $\tau = 0.3$  and vary the generator's learning rate  $\tilde{\tau} = 0.03, 0.2, 0.4, 0.47$  from left to right column. These parameter settings are marked by the four red dots in the phase diagram in Figure 3. The second row is the phase portraits of the trajectories shown in the first row onto the  $P$ - $q$  plane. Figure (a) shows a case in the phase of info-1, where a subset of type (3) fixed points are stable. Figure (b) and (c) are in the oscillating phase, and (d) is in info-2, where the fixed points of type-5 are stable. The blue dots in the figures show the stable fixed points.

We note that the elements of  $\mathbf{w}_k$  and  $\mathbf{V}_k$  are  $\mathcal{O}(\frac{1}{\sqrt{n}})$  number as the norm of  $\mathbf{w}_k$  and the norms of column vectors of  $\mathbf{V}_k$  are all  $\mathcal{O}(1)$  numbers. To investigate the dynamics of the microscopic state, it is convenient to rescale  $\mathbf{w}_k$  and  $\mathbf{V}_k$  by a factor of  $\sqrt{n}$ . We define  $\hat{\mathbf{u}}_i$  and  $\hat{\mathbf{v}}_{k,i}$  as the column view of

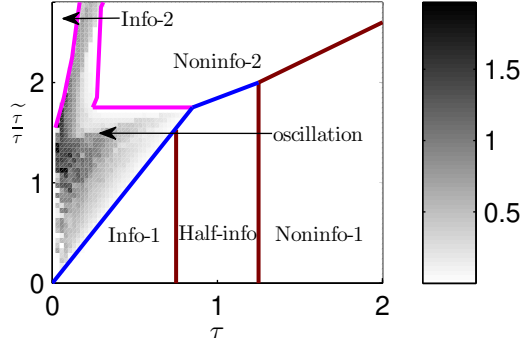


Figure 5: The phase diagram for the stationary states of the ODE (13) when  $d = 2$ . This phase diagram is generated by numerically computing the fixed points and eigenvalues of the Jacobian of the ODE (13).

the  $i$ 'th row of the matrices  $\sqrt{n}U$  and  $\sqrt{n}V_k$  respectively, and  $\hat{w}_{k,i} \stackrel{\text{def}}{=} \sqrt{n}[\mathbf{w}_{k+1}]_i$ . The update rule of  $((\hat{\mathbf{u}}_i, \hat{\mathbf{v}}_{k,i}, \hat{w}_{k,i})_{i=1,\dots,n})_{k=0,1,2,\dots}$  is

$$\hat{w}_{k+1,i} - \hat{w}_{k,i} = \frac{\tau}{n} \left[ (\hat{\mathbf{u}}_i^\top \mathbf{c}_k + \sqrt{n\eta_\Gamma} a_{k,i}) f_k - (\hat{\mathbf{v}}_{k,i}^\top \tilde{\mathbf{c}}_{2k} + \sqrt{n\eta_G} \tilde{a}_{2k,i}) \tilde{f}_{2k} - \lambda H'(z_k) \hat{w}_{k,i} \right], \quad (\text{S-3})$$

$$\hat{\mathbf{v}}_{k+1,i} - \hat{\mathbf{v}}_{k,i} = \frac{\tilde{\tau}}{n} \left[ \hat{w}_{k,i} \tilde{\mathbf{c}}_{2k+1} \tilde{f}_{2k+1} - \lambda \text{diag}(H'(\mathbf{S}_k)) \hat{\mathbf{v}}_{k,i} \right], \quad (\text{S-4})$$

where  $a_{k,i}, \tilde{a}_{k,i}$  are the  $i$ th elements of  $\mathbf{a}_k$  and  $\tilde{\mathbf{a}}_k$  respectively, and  $f_k$  and  $\tilde{f}_k$  are shorthands for

$$f_k = f(\mathbf{y}_k^\top \mathbf{w}_k / \sqrt{n}) = f\left(\mathbf{q}_k^\top \mathbf{c}_k + \sqrt{\frac{\eta_\Gamma}{n}} \sum_{j=1}^n a_{k,j} \hat{w}_{k,j}\right)$$

$$\tilde{f}_k = \tilde{f}(\mathbf{y}_k^\top \mathbf{w}_{\lfloor k/2 \rfloor} / \sqrt{n}) = \tilde{f}\left(\mathbf{r}_{\lfloor k/2 \rfloor}^\top \tilde{\mathbf{c}}_k + \sqrt{\frac{\eta_G}{n}} \sum_{j=1}^n \tilde{a}_{k,j} \hat{w}_{\lfloor k/2 \rfloor, j}\right),$$

respectively, and the empirical macroscopic quantities  $\mathbf{q}_k, \mathbf{r}_k, z_k$  and  $\mathbf{S}_k$  are defined as follows

$$\begin{aligned} \mathbf{q}_k &\stackrel{\text{def}}{=} \mathbf{U}^\top \mathbf{w}_k = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{u}}_i \hat{w}_i, & \mathbf{r}_k &\stackrel{\text{def}}{=} \mathbf{V}_k^\top \mathbf{w}_k = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{v}}_{k,i} \hat{w}_i, \\ z_k &\stackrel{\text{def}}{=} \mathbf{w}_k^\top \mathbf{w}_k = \frac{1}{n} \sum_{i=1}^n \hat{w}_{k,i}^2, & \mathbf{S}_k &\stackrel{\text{def}}{=} \mathbf{V}_k^\top \mathbf{V}_k = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{v}}_{k,i} \hat{\mathbf{v}}_{k,i}^\top, \\ \mathbf{P}_k &\stackrel{\text{def}}{=} \mathbf{U}^\top \mathbf{V}_k = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{u}}_i \hat{\mathbf{v}}_{k,i}^\top. \end{aligned} \quad (\text{S-5})$$

The matrix  $\mathbf{P}_k$  is not used in this section, but we put it here with the other macroscopic quantities for future reference.

Now we derive (10) from (S-3) and (S-4).

First, it is trivial to get the first equation of the SDE  $d\hat{\mathbf{u}}_t = 0$  in (10) in the main text, since  $\hat{\mathbf{u}}_i$  does not change over time.

Next, we derive the second equation in (10). Averaging over  $\tilde{\mathbf{c}}_{2k+1}$  and  $\tilde{\mathbf{a}}_{2k+1}$  on the both sides of (S-4), we get

$$\begin{aligned} &\langle \hat{\mathbf{v}}_{k+1,i} - \hat{\mathbf{v}}_{k,i} \rangle_{\tilde{\mathbf{c}}_{2k+1}, \tilde{\mathbf{a}}_{2k+1}} \\ &= \frac{\tilde{\tau}}{n} \left[ \left\langle \tilde{f}\left(\mathbf{r}_k^\top \tilde{\mathbf{c}} + \sqrt{\frac{\eta_G}{n}} \sum_{j=1}^n [\tilde{\mathbf{a}}]_j \hat{w}_{k,j}\right) \tilde{\mathbf{c}} \right\rangle_{\tilde{\mathbf{c}}, \tilde{\mathbf{a}}} \hat{w}_{k,i} - \lambda \text{diag}(H'(\mathbf{S}_k)) \hat{\mathbf{v}}_{k,i} \right]. \end{aligned}$$

The bracket  $\langle \cdot \rangle_{\tilde{\mathbf{c}}, \tilde{\mathbf{a}}}$  here denotes the average over  $\tilde{\mathbf{c}} \sim \mathcal{P}_{\tilde{\mathbf{c}}}$ , and standard Gaussian vector  $\tilde{\mathbf{a}}$ , where  $\tilde{\mathbf{c}}$  and  $\tilde{\mathbf{a}}$  are the random variables generating the fake sample in the generator as described in (2). Noting that  $\tilde{\mathbf{a}}$  is a Gaussian vector, the term  $\frac{1}{\sqrt{n}} \sum_{j=1}^n [\tilde{\mathbf{a}}]_j \hat{w}_{k,j}$  in the above equation is also a Gaussian

random variable, whose mean is zero and variance is  $z_k$ , which is defined in (S-5). Therefore, we have

$$\langle \hat{\mathbf{v}}_{k+1,i} - \hat{\mathbf{v}}_{k,i} \rangle_{\tilde{\mathbf{c}}_{2k+1}, \tilde{\mathbf{a}}_{2k+1}} = \frac{\tau}{n} [\tilde{\mathbf{g}}_k \hat{\mathbf{w}}_{k,i} + \mathbf{L}_k \hat{\mathbf{v}}_{k,i}], \quad (\text{S-6})$$

where

$$\tilde{\mathbf{g}}_k = \left\langle \tilde{f}(\mathbf{r}_k^\top \tilde{\mathbf{c}} + \sqrt{z_k \eta_G} e) \right\rangle_{\tilde{\mathbf{c}}, e} \quad (\text{S-7})$$

$$\mathbf{L}_k = -\lambda \text{diag}(H'(\mathbf{S}_k)), \quad (\text{S-8})$$

where  $\langle \cdot \rangle_{\tilde{\mathbf{c}}, e}$  denotes the average over  $\tilde{\mathbf{c}} \sim \mathcal{P}_{\tilde{\mathbf{c}}}$  and  $e \sim \mathcal{N}(0, 1)$ . In addition, from (S-4), we also know that the second moment

$$\left\langle (\hat{\mathbf{v}}_{k+1,i} - \hat{\mathbf{v}}_{k,i})^2 \right\rangle_{\tilde{\mathbf{c}}_{2k+1}, \tilde{\mathbf{a}}_{2k+1}} = \mathcal{O}(n^{-\frac{3}{2}}). \quad (\text{S-9})$$

The moments estimations (S-6) and (S-9) imply the second equation in (10) in the main text. Since the second moments growth smaller than  $\mathcal{O}(n^{-1})$ , the differential equation for  $\hat{\mathbf{v}}_t$  has no diffusion term.

Finally, we derive the last equation in (10) in the main text from the update rule of  $\hat{\mathbf{w}}_k$  (S-3). We observe that both the terms inside the function  $f$  and outside of  $f$  in (S-3) depend on  $a_{k,i}$ . Using Taylor's expansion, we linearize the contribution of  $a_{k,i}$  to the function  $f$ :

$$\begin{aligned} f_k &= f\left(\mathbf{q}_k^\top \mathbf{c}_k + \sqrt{\frac{\eta_I}{n}} \sum_{j \neq i} a_{k,j} \hat{\mathbf{w}}_{k,j} + \sqrt{\frac{\eta_I}{n}} a_{k,i} \hat{\mathbf{w}}_{k,i}\right) \\ &= f(\mathbf{q}_k^\top \mathbf{c}_k + \sqrt{\frac{\eta_I}{n}} \sum_{j \neq i} a_{k,j} \hat{\mathbf{w}}_{k,j}) + f'(\mathbf{q}_k^\top \mathbf{c}_k + \sqrt{\frac{\eta_I}{n}} \sum_{j \neq i} a_{k,j} \hat{\mathbf{w}}_{k,j}) \sqrt{\frac{\eta_I}{n}} a_{k,i} \hat{\mathbf{w}}_{k,i} + \mathcal{O}(\frac{1}{n}). \end{aligned} \quad (\text{S-10})$$

Similarly, we have

$$\begin{aligned} \tilde{f}_{2k} &= \tilde{f}(\mathbf{r}_k^\top \tilde{\mathbf{c}}_{2k} + \sqrt{\frac{\eta_G}{n}} \sum_{j \neq i} \tilde{a}_{k,j} \hat{\mathbf{w}}_{k,j} + \sqrt{\frac{\eta_G}{n}} \tilde{a}_{2k,i} \hat{\mathbf{w}}_{k,i}) \\ &= \tilde{f}(\mathbf{r}_k^\top \tilde{\mathbf{c}}_{2k} + \sqrt{\frac{\eta_G}{n}} \sum_{j \neq i} \tilde{a}_{k,j} \hat{\mathbf{w}}_{k,j}) + \tilde{f}'(\mathbf{r}_k^\top \tilde{\mathbf{c}}_{2k} + \sqrt{\frac{\eta_G}{n}} \sum_{j \neq i} \tilde{a}_{k,j} \hat{\mathbf{w}}_{k,j}) \sqrt{\frac{\eta_G}{n}} \tilde{a}_{2k,i} \hat{\mathbf{w}}_{k,i} + \mathcal{O}(\frac{1}{n}) \end{aligned} \quad (\text{S-11})$$

Substituting (S-10) and (S-11) into (S-3), we have

$$\begin{aligned} &\frac{\hat{\mathbf{w}}_{k+1,i} - \hat{\mathbf{w}}_{k,i}}{\tau/n} \\ &= \hat{\mathbf{u}}_i^\top \mathbf{c}_k f(\mathbf{q}_k^\top \mathbf{c}_k + \sqrt{\frac{\eta_I}{n}} \sum_{j \neq i} a_{k,j} \hat{\mathbf{w}}_{k,j}) - \hat{\mathbf{v}}_{k,i}^\top \tilde{\mathbf{c}}_{2k} \tilde{f}(\mathbf{r}_k^\top \tilde{\mathbf{c}}_{2k} + \sqrt{\frac{\eta_G}{n}} \sum_{j \neq i} \tilde{a}_{k,j} \hat{\mathbf{w}}_{k,j}) \\ &+ \hat{\mathbf{w}}_{k,i} \left[ a_{k,i}^2 f'(\mathbf{q}_k^\top \mathbf{c}_k + \sqrt{\frac{\eta_I}{n}} \sum_{j \neq i} a_{k,j} \hat{\mathbf{w}}_{k,j}) - \tilde{a}_{2k,i}^2 \tilde{f}'(\mathbf{r}_k^\top \tilde{\mathbf{c}}_{2k} + \sqrt{\frac{\eta_G}{n}} \sum_{j \neq i} \tilde{a}_{2k,j} \hat{\mathbf{w}}_{k,j}) - \lambda H'(z_k) \right] \\ &+ \sqrt{n} \left[ a_{k,i} f(\mathbf{q}_k^\top \mathbf{c}_k + \sqrt{\frac{\eta_I}{n}} \sum_{j \neq i} a_{k,j} \hat{\mathbf{w}}_{k,j}) + \tilde{a}_{2k,i} \tilde{f}(\mathbf{r}_k^\top \tilde{\mathbf{c}}_{2k} + \sqrt{\frac{\eta_G}{n}} \sum_{j \neq i} \tilde{a}_{2k,j} \hat{\mathbf{w}}_{k,j}) \right] + \delta_{k,i}, \end{aligned} \quad (\text{S-12})$$

where  $\delta_{k,i}$  collects all higher-order terms whose contributions will vanish as  $n \rightarrow \infty$ . From this equation, we can already infer the SDE (10). Specifically, on the right hand side of (S-12), the terms in the first two lines correspond to the drift term in the SDE. Furthermore, the first term in the third line in (S-12) contributes to the SDE as a Brownian motion. More precisely, we can derive the third equation of the SDE (10) in the main text by the moments estimations. Specifically, the first-order moment is

$$\langle \hat{\mathbf{w}}_{k+1,i} - \hat{\mathbf{w}}_{k,i} \rangle_{\mathbf{c}_k, \mathbf{a}_k, \tilde{\mathbf{c}}_{2k}, \tilde{\mathbf{a}}_{2k}} = \frac{\tau}{n} [\hat{\mathbf{u}}_i^\top \mathbf{g}_k - \hat{\mathbf{v}}_{k,i}^\top \tilde{\mathbf{g}}_k + \hat{\mathbf{w}}_{k,i} h_k] + \mathcal{O}(n^{-\frac{3}{2}}) \quad (\text{S-13})$$

where  $\tilde{\mathbf{g}}_k$  is defined in (S-7), and

$$\mathbf{g}_k = \left\langle \mathbf{c} f(\mathbf{q}_k^\top \mathbf{c} + \sqrt{z_k \eta_I} e) \right\rangle_{\mathbf{c}, e} \quad (\text{S-14})$$

$$h_k = \eta_I \left\langle f'(\mathbf{q}_k^\top \mathbf{c} + \sqrt{z_k \eta_I} e) \right\rangle_{\mathbf{c}, e} - \tilde{\eta}_G \left\langle \tilde{f}'(\mathbf{r}_k^\top \tilde{\mathbf{c}} + \sqrt{z_k \eta_G} e) \right\rangle_{\tilde{\mathbf{c}}, e} - \lambda H'(z_k). \quad (\text{S-15})$$

The second moment is

$$\left\langle (\widehat{w}_{k+1,i} - \widehat{w}_{k,i})^2 \right\rangle_{\mathbf{c}_k, \mathbf{a}_k, \widetilde{\mathbf{c}}_{2k}, \widetilde{\mathbf{a}}_{2k}} = \frac{\tau^2}{n} b_k + \mathcal{O}(n^{-\frac{3}{2}}), \quad (\text{S-16})$$

where

$$b_k = \eta_{\text{T}} \left\langle f^2(\mathbf{q}_k^\top \mathbf{c} + \sqrt{z_k \eta_{\text{T}} e}) \right\rangle_{\mathbf{c}, e} + \eta_{\text{G}} \left\langle \widetilde{f}^2(\mathbf{r}_k^\top \widetilde{\mathbf{c}} + \sqrt{z_k \eta_{\text{G}} e}) \right\rangle_{\widetilde{\mathbf{c}}, e}. \quad (\text{S-17})$$

From the (S-13) and (S-16), we derive the SDE for  $\widehat{w}_t$  in (10) in the main text.

### S-III Derive the ODE in Theorem 1 from the weak formulation of the PDE

In this section, we show how to derive the ODE (8) from the weak formulation of the PDE (12). Choosing the test function  $\varphi$  being each element of  $\widehat{\mathbf{u}}\widehat{\mathbf{v}}^\top$ ,  $\widehat{\mathbf{u}}\widehat{w}$ ,  $\widehat{\mathbf{v}}\widehat{w}$ ,  $\widehat{\mathbf{v}}\widehat{\mathbf{v}}^\top$ ,  $\widehat{w}^2$ , and substituting those  $\varphi$  into the weak formulation of the PDE (12), we will get the ODE (8) as presented in Theorem 1. In what follows, we provide additional details of this derivation.

We first derive the first ODE  $\frac{d}{dt} \mathbf{P}_t = \dots$  in (8). Let  $\varphi = [\widehat{\mathbf{u}}]_\ell [\widehat{\mathbf{v}}]_{\ell'}$ ,  $\ell, \ell' = 1, 2, \dots, d$ , we have  $\nabla_{\widehat{\mathbf{v}}} \varphi = [\widehat{\mathbf{u}}]_\ell \mathbf{s}_{\ell'}$ , where  $\mathbf{s}_{\ell'}$  is the  $\ell'$ th canonical basis (i.e., all elements in  $\mathbf{s}_{\ell'}$  are zeros, except that  $\ell'$ th element is 1). From the PDE (12) in the main text, we have  $\forall \ell, \ell' = 1, 2, \dots, d$ :

$$\langle \mu_t, \varphi(\widehat{\mathbf{u}}, \widehat{\mathbf{v}}, \widehat{w}) \rangle = \langle \mu_t, [\widehat{\mathbf{u}}]_\ell [\widehat{\mathbf{v}}]_{\ell'} \rangle = [\mathbf{P}_t]_{\ell, \ell'},$$

$$\begin{aligned} \left\langle \mu_t, (\widehat{w} \widetilde{\mathbf{g}}_t^\top + \widehat{\mathbf{v}}^\top \mathbf{L}_t) \nabla_{\widehat{\mathbf{v}}} \varphi \right\rangle &= \left\langle \mu_t, ([\widehat{\mathbf{u}}]_\ell \widehat{w}) [\widetilde{\mathbf{g}}_t]_{\ell'} + ([\widehat{\mathbf{u}}]_\ell \widehat{\mathbf{v}}^\top) [\mathbf{L}_t]_{:, \ell'} \right\rangle \\ &= [\mathbf{q}_t]_{\ell} [\widetilde{\mathbf{g}}_t]_{\ell'} + [\mathbf{P}_t]_{\ell, :} [\mathbf{L}_t]_{:, \ell'}, \end{aligned}$$

where  $[\mathbf{P}_t]_{\ell, :}$  and  $[\mathbf{L}_t]_{:, \ell'}$  are  $\ell$ th row of  $\mathbf{P}_t$  and  $\ell'$ th column of  $\mathbf{L}$ , respectively. In addition, we know that  $\frac{\partial}{\partial \widehat{w}} \varphi = \frac{\partial^2}{\partial \widehat{w}^2} \varphi = 0$ . Combining above results, we can recover the first ODE in (8).

Next, we derive the second ODE  $\frac{d\mathbf{q}_t}{dt} = \dots$  in (8). Let  $\varphi = [\widehat{\mathbf{u}}]_\ell \widehat{w}$ ,  $\ell = 1, 2, \dots, d$ . We have  $\nabla_{\widehat{\mathbf{v}}} \varphi = 0$ ,  $\frac{\partial}{\partial \widehat{w}} \varphi = [\widehat{\mathbf{u}}]_\ell$  and  $\frac{\partial^2}{\partial \widehat{w}^2} \varphi = 0$ . Then  $\forall \ell = 1, 2, \dots, d$ ,

$$\langle \mu_t, \varphi(\widehat{\mathbf{u}}, \widehat{\mathbf{v}}, \widehat{w}) \rangle = \langle \mu_t, [\widehat{\mathbf{u}}]_\ell \widehat{w} \rangle = [\mathbf{q}_t]_\ell$$

and

$$\begin{aligned} \left\langle \mu_t, (\widehat{\mathbf{u}}^\top \mathbf{g}_t - \widehat{\mathbf{v}}^\top \widetilde{\mathbf{g}}_t + h_t \widehat{w}) \frac{\partial}{\partial \widehat{w}} \varphi \right\rangle &= \left\langle \mu_t, (\widehat{\mathbf{u}}^\top \mathbf{g}_t - \widehat{\mathbf{v}}^\top \widetilde{\mathbf{g}}_t + h_t \widehat{w}) [\widehat{\mathbf{u}}]_\ell \right\rangle \\ &= [\mathbf{g}_t]_\ell - [\mathbf{P}_t]_{\ell} \widetilde{\mathbf{g}}_t + [\mathbf{q}_t]_\ell h_t. \end{aligned}$$

With above results, we can obtain the second ODE in (8).

Next, let's derive the ODE for  $\frac{d\mathbf{S}_t}{dt}$ . We set  $\varphi = [\widehat{\mathbf{v}}]_\ell [\widehat{\mathbf{v}}]_{\ell'}$ . If  $\ell \neq \ell'$ , we have  $\nabla_{\widehat{\mathbf{v}}} \varphi = [\widehat{\mathbf{v}}]_\ell \mathbf{s}_{\ell'} + [\widehat{\mathbf{v}}]_{\ell'} \mathbf{s}_\ell$ , where  $\mathbf{s}_{\ell'}$  is the  $\ell'$ th canonical basis. Then

$$\langle \mu_t, \varphi(\widehat{\mathbf{u}}, \widehat{\mathbf{v}}, \widehat{w}) \rangle = [\mathbf{S}_t]_{\ell, \ell'}$$

and

$$\begin{aligned} \left\langle \mu_t, (\widehat{w} \widetilde{\mathbf{g}}_t^\top + \widehat{\mathbf{v}}^\top \mathbf{L}_t) \nabla_{\widehat{\mathbf{v}}} \varphi \right\rangle &= \left\langle \mu_t, ([\widehat{\mathbf{v}}]_\ell \widehat{w}) [\widetilde{\mathbf{g}}_t]_{\ell'} + ([\widehat{\mathbf{v}}]_\ell \widehat{\mathbf{v}}^\top) [\mathbf{L}_t]_{:, \ell'} \right\rangle \\ &\quad + \left\langle \mu_t, ([\widehat{\mathbf{v}}]_{\ell'} \widehat{w}) [\widetilde{\mathbf{g}}_t]_\ell + ([\widehat{\mathbf{v}}]_{\ell'} \widehat{\mathbf{v}}^\top) [\mathbf{L}_t]_{:, \ell} \right\rangle \\ &= [\mathbf{r}_t]_\ell [\widetilde{\mathbf{g}}_t]_{\ell'} + [\widetilde{\mathbf{g}}_t]_\ell [\mathbf{r}_t]_{\ell'} + [\mathbf{S}_t]_{\ell, :} [\mathbf{L}_t]_{:, \ell'} + [\mathbf{L}_t]_{\ell, :} [\mathbf{S}_t]_{:, \ell'} \end{aligned}$$

If  $\ell = \ell'$ , we have  $\nabla_{\widehat{\mathbf{v}}} \varphi = 2[\widehat{\mathbf{v}}]_\ell \mathbf{s}_\ell$ , then

$$\langle \mu_t, \varphi(\widehat{\mathbf{u}}, \widehat{\mathbf{v}}, \widehat{w}) \rangle = [\mathbf{S}_t]_{\ell, \ell}$$

and

$$\left\langle \mu_t, (\widehat{w} \widetilde{\mathbf{g}}_t^\top + \widehat{\mathbf{v}}^\top \mathbf{L}_t) \nabla_{\widehat{\mathbf{v}}} \varphi \right\rangle = 2([\mathbf{r}_t]_\ell [\widetilde{\mathbf{g}}_t]_\ell + [\mathbf{S}_t]_{\ell, :} [\mathbf{L}_t]_{:, \ell})$$

Plugging back the above two equations and combining the fact that  $\frac{\partial}{\partial \widehat{w}} \varphi = \frac{\partial^2}{\partial \widehat{w}^2} \varphi = 0$ , we recover the ODE of  $\frac{d\mathbf{S}_t}{dt}$ .

The rest two ODEs can be obtained in the similar way by letting  $\varphi$  to be each distinct component of  $\widehat{\mathbf{v}}\widehat{w}$  and  $\widehat{w}^2$ .



## S-IV Proof of Theorem 1

In this section, we prove Theorem 1 shown in the main text. In the previous section, we have already provided a derivation of the ODE in Theorem 1 from the weak formulation of the PDE for the microscopic states. In this section, we follow a different path to prove the theorem without referencing the PDE, because it is easier to establish the rigorous bound of the convergence rate. Thus, the proof itself also provides another derivation of the ODE, where the most relevant part is Lemma 5.

### S-IV.1 Sketch of the proof

The proof follows the standard procedure of the convergence of stochastic processes [29, 30]. We here build the whole proof on Lemma 2 in the supplementary materials of [22]. For reader's convenient, we present that lemma below.

**Lemma 1** (Lemma 2 in the supplementary materials of [22]). *Consider a sequence of stochastic process  $\{\mathbf{x}_k^{(n)}, k = 0, 1, 2, \dots, \lfloor nT \rfloor\}_{n=1,2,\dots}$ , with some constant  $T > 0$ . If  $\mathbf{x}_k^{(n)}$  can be decomposed into three parts*

$$\mathbf{x}_{k+1}^{(n)} - \mathbf{x}_k^{(n)} = \frac{1}{n}\phi(\mathbf{x}_k^{(n)}) + \boldsymbol{\rho}_k^{(n)} + \boldsymbol{\delta}_k^{(n)} \quad (\text{S-18})$$

such that

(C.1) *The process  $\sum_{k'=0}^k \boldsymbol{\rho}_{k'}^{(n)}$  is a martingale, and  $\mathbb{E} \|\boldsymbol{\rho}_k^{(n)}\|^2 \leq C(T)/n^{1+\epsilon_1}$  for some positive  $\epsilon_1$ ;*

(C.2)  $\mathbb{E} \|\boldsymbol{\delta}_k^{(n)}\| \leq C(T)/n^{1+\epsilon_2}$  *for some positive  $\epsilon_2$ ;*

(C.3)  $\phi(\mathbf{x})$  *is a Lipschitz function, i.e.,  $\|\phi(\mathbf{x}) - \phi(\tilde{\mathbf{x}})\| \leq C\|\mathbf{x} - \tilde{\mathbf{x}}\|$ ;*

(C.4)  $\mathbb{E} \|\mathbf{x}_k^{(n)}\|^2 \leq C$  *for all  $k \leq \lfloor nT \rfloor$ ;*

(C.5)  $\mathbb{E} \|\mathbf{x}_0^{(n)} - \mathbf{x}_0^*\| \leq C/n^{\epsilon_3}$  *for some positive  $\epsilon_3$  and a deterministic vector  $\mathbf{x}_0^*$ ,*

then we have

$$\|\mathbf{x}_k^{(n)} - \mathbf{x}(\frac{k}{n})\| \leq C(T)n^{-\min\{\frac{1}{2}\epsilon_1, \epsilon_2, \epsilon_3\}},$$

where  $\mathbf{x}(t)$  is the solution of the ODE

$$\frac{d}{dt}\mathbf{x}(t) = \phi(\mathbf{x}(t)), \quad \text{with } \mathbf{x}(0) = \mathbf{x}_0^*.$$

In Theorem 1, the stochastic process is the macroscopic states  $\{\mathbf{M}_k, k = 0, 1, \dots\}$ , where  $\mathbf{M}_k$  is a symmetric matrix consists of 5 non-trivial parts  $\mathbf{P}_k, \mathbf{q}_k, \mathbf{r}_k, \mathbf{S}_k$ , and  $z_k$  as shown in (6) in the main text. Following (S-18), we have the following decomposition for  $\mathbf{M}_k$

$$\mathbf{M}_{k+1} - \mathbf{M}_k = \frac{1}{n}\phi(\mathbf{M}_k) + (\mathbf{M}_{k+1} - \mathbb{E}_k \mathbf{M}_{k+1}) + [\mathbb{E}_k \mathbf{M}_{k+1} - \mathbf{M}_k - \frac{1}{n}\phi(\mathbf{M}_k)], \quad (\text{S-19})$$

in which the matrix-valued function  $\phi(\mathbf{M})$  represents the functions on the right hand sides of the ODE (8), and  $\mathbb{E}_k$  denotes the conditional expectation given the state of the Markov chain  $\mathbf{X}_k$ . Note that the stochastic process of the macroscopic state  $\mathbf{M}_k$  is driven by the Markov chain of the microscopic state  $\mathbf{X}_k$ . Thus,  $\mathbb{E}_k$  is well-defined. For future reference, we denotes  $\mathbb{E}$  the unconditional expectation of all the randomness of the Markov chain  $\mathbf{X}_k$ , i.e., the initial state  $\mathbf{U}, \mathbf{V}_0, \mathbf{w}_0$  and  $\{\mathbf{a}_k, \mathbf{c}_k, \tilde{\mathbf{a}}_k, \tilde{\mathbf{c}}_k | k = 0, 1, 2, \dots\}$ . By definition,  $\sum_{k'=0}^k (\mathbf{M}_{k'+1} - \mathbb{E}_{k'} \mathbf{M}_{k'})$  is a Martingale.

### S-IV.2 Check the conditions provided in Lemma 1

In this subsection, we check the condition (C.1)–(C.5) for the decomposition of (S-19). Once all conditions are proved to be satisfied, Theorem 1 will be proved.

We first note that (C.5) is the assumption (A.5) in the main text. Thus, (C.5) is satisfied. Before proving other conditions, we declare a lemma.

**Lemma 2.** *Under the same setting as Theorem 1, given  $T > 0$ , then*

$$\mathbb{E} \left( \sum_{\ell=1}^d [\mathbf{V}_k]_{i,\ell}^4 + [\mathbf{w}_k]_i^4 \right) \leq C(T)n^{-2}, \quad \forall i = 1, 2, \dots, n, \text{ and } k = 0, 1, \dots, \lfloor nT \rfloor, \quad (\text{S-20})$$

The proof can be founded in Section S-IV.3.

#### Check Condition (C.4)

**Lemma 3.** *Under the same setting as Theorem 1, for all  $k = 0, 1, \dots, \lfloor nT \rfloor$  with a given  $T > 0$ , then*

$$\begin{aligned}\mathbb{E} \|\mathbf{P}_k\|^2 &\leq C(T), & \mathbb{E} \|\mathbf{q}_k\|^2 &\leq C(T), \\ \mathbb{E} \|\mathbf{S}_k\|^2 &\leq C(T), & \mathbb{E} z_k^2 &\leq C(T), \\ \mathbb{E} \|\mathbf{r}_k\|^2 &\leq C(T).\end{aligned}$$

*Proof.* It's a direct consequence of Lemma 2. We first verify  $\mathbb{E} z_k^2 \leq C(T)$ . Using Holder's inequality, we have

$$\mathbb{E} z_k^2 = \mathbb{E} \left( \sum_{i=1}^n w_{k,i}^2 \right)^2 \leq n \mathbb{E} \sum_{i=1}^n w_{k,i}^4 \leq C(T)$$

For  $[\mathbf{S}_k]_{\ell,\ell}$ ,  $\ell = 1, \dots, d$ , similarly, we have

$$\mathbb{E} [\mathbf{S}_k]_{\ell,\ell}^2 = \mathbb{E} \left( \sum_{i=1}^n [\mathbf{V}_k]_{i,\ell}^2 \right)^2 \leq C(T).$$

and for  $\mathbb{E} [\mathbf{S}_k]_{\ell,\ell'}, \ell \neq \ell'$ , we have:

$$\begin{aligned}\mathbb{E} [\mathbf{S}_k]_{\ell,\ell'}^2 &= \mathbb{E} \left( \sum_{i=1}^n [\mathbf{V}_k]_{i,\ell} [\mathbf{V}_k]_{i,\ell'} \right)^2 \\ &\leq \mathbb{E} \left( \sum_{i=1}^n [\mathbf{V}_k]_{i,\ell}^2 \right) \left( \sum_{i=1}^n [\mathbf{V}_k]_{i,\ell'}^2 \right) \\ &\leq \sqrt{\mathbb{E} \left( \sum_{i=1}^n [\mathbf{V}_k]_{i,\ell}^2 \right)^2 \mathbb{E} \left( \sum_{i=1}^n [\mathbf{V}_k]_{i,\ell'}^2 \right)^2} \\ &\leq C(T)\end{aligned}$$

where in reaching the third and last line, we used the Cauchy-Schwartz inequality. Now, we get  $\mathbb{E} \|\mathbf{S}_k\|^2 \leq C(T)$ . The rest bounds of  $\mathbb{E} \|\mathbf{P}_k\|^2$ ,  $\mathbb{E} \|\mathbf{q}_k\|^2$  and  $\mathbb{E} \|\mathbf{r}_k\|^2$  in Lemma 3 can also be directly verified using the Cauchy-Schwartz inequality.  $\square$

#### Check Condition (C.3)

**Lemma 4.** *If Assumption (A.3) hold,  $\phi(\mathbf{M})$  is a Lipschitz function.*

*Proof.* It suffices to verify each component of gradient  $\nabla \phi(\mathbf{M})$  is bounded. Assumption (A.3) ensures that  $H'$  is Lipschitz and the derivatives up to fourth order of the functions  $f$ ,  $\tilde{f}$  exists and uniformly bounded. These conditions guarantee that the partial derivatives of  $\phi(\mathbf{M})$  w.r.t.  $\mathbf{P}$ ,  $\mathbf{q}$ ,  $\mathbf{S}$  and  $\mathbf{r}$  are bounded. The remaining thing is to show that  $\frac{\partial \phi(\mathbf{M})}{\partial z}$  is also bounded. Since there is a  $\sqrt{z}$  term in  $\phi(\mathbf{M})$ , the boundness can be potentially broken at  $z = 0$ . However, we can show that it is not the case. For example, we can show that  $\langle \mathbf{c} f(\mathbf{c}^\top \mathbf{q} + e\sqrt{z}) \rangle_{\mathbf{c},e}$  is a Lipschitz function, because

$$\begin{aligned}\frac{\partial}{\partial z} \langle \mathbf{c} f(\mathbf{c}^\top \mathbf{q} + e\sqrt{z}) \rangle_{\mathbf{c},e} &= \frac{1}{2} z^{-\frac{1}{2}} \langle e c f'(c\mathbf{q} + e\sqrt{z}) \rangle_{\mathbf{c},e} \\ &= \frac{1}{2} \langle c f''(c\mathbf{q} + e\sqrt{z}) \rangle_{\mathbf{c},e}\end{aligned}$$

is always a well-defined bounded function. In reaching the first line, we here interchanged the expectation and derivative, which is valid because of the boundness of  $f(\cdot)$ , and in reaching the second line, we used the Stein's lemma. Finally, other terms in (9) involving  $\sqrt{z}$  can be treated in the same way. Thus,  $\phi(\mathbf{M})$  is a Lipschitz function.  $\square$

#### Check Condition (C.2)

**Lemma 5.** *Under the same setting as Theorem 1, for all  $k = 0, 1, \dots, \lfloor nT \rfloor$  with a given  $T > 0$ , then*

$$\mathbb{E} \|\mathbb{E}_k \mathbf{M}_{k+1} - \mathbf{M}_k - \frac{1}{n} \phi(\mathbf{M}_k)\| \leq C(T) n^{-\frac{3}{2}}.$$

*Proof.* The above inequality can be split into 5 parts

$$\mathbb{E} \|\mathbb{E}_k \mathbf{P}_{k+1} - \mathbf{P}_k - \frac{\tilde{\tau}}{n}(\mathbf{q}_k \tilde{\mathbf{g}}_k^\top + \mathbf{P}_k \mathbf{L}_k)\| \leq C(T)n^{-\frac{3}{2}} \quad (\text{S-21})$$

$$\mathbb{E} \|\mathbb{E}_k \mathbf{q}_{k+1} - \mathbf{q}_k - \frac{\tau}{n}(\mathbf{g}_k - \mathbf{P}_k \tilde{\mathbf{g}}_k + \mathbf{q}_k h_k)\| \leq C(T)n^{-\frac{3}{2}} \quad (\text{S-22})$$

$$\mathbb{E} \|\mathbb{E}_k \mathbf{S}_{k+1} - \mathbf{S}_k - \frac{\tilde{\tau}}{n}(\mathbf{r}_k \tilde{\mathbf{g}}_k^\top + \tilde{\mathbf{g}}_k \mathbf{r}_k^\top + \mathbf{S}_k \mathbf{L}_k + \mathbf{L}_k \mathbf{S}_k)\| \leq C(T)n^{-\frac{3}{2}} \quad (\text{S-23})$$

$$\mathbb{E} \|\mathbb{E}_k z_{k+1} - z_k - \frac{2\tau}{n}(\mathbf{q}_k^\top \mathbf{g}_k - \mathbf{r}_k^\top \tilde{\mathbf{g}}_k + z_k h_k) - \frac{\tau^2}{n} b_k\| \leq C(T)n^{-\frac{3}{2}}, \quad (\text{S-24})$$

$$\mathbb{E} \|\mathbb{E}_k \mathbf{r}_{k+1} - \mathbf{r}_k - \frac{\tau}{n}(\mathbf{P}_k^\top \mathbf{g}_k - \mathbf{S}_k \tilde{\mathbf{g}}_k + \mathbf{r}_k h_k) - \frac{\tilde{\tau}}{n}(z_k \tilde{\mathbf{g}}_k + \mathbf{L}_k \mathbf{r}_k)\| \leq C(T)n^{-\frac{3}{2}} \quad (\text{S-25})$$

where  $\tilde{\mathbf{g}}_k$ ,  $\mathbf{L}_k$ ,  $\mathbf{g}_k$ ,  $h_k$ ,  $b_k$  are defined in (S-7), (S-8), (S-14), (S-15) and (S-17), respectively.

We first prove (S-21). From (S-2), we have

$$\mathbf{V}_{k+1} - \mathbf{V}_k = \frac{\tilde{\tau}}{n}[\mathbf{w}_k \tilde{\mathbf{c}}_{2k+1}^\top \tilde{f}(\tilde{\mathbf{c}}_{2k+1}^\top \mathbf{V}_k^\top \mathbf{w}_k + \eta_G \tilde{\mathbf{a}}_{2k+1}^\top \mathbf{w}_k) - \lambda \mathbf{V}_k \text{diag}(H'(\mathbf{S}_k))]. \quad (\text{S-26})$$

Averaging both sides of the above equation over  $\tilde{\mathbf{c}}_{2k+1}$  and  $\tilde{\mathbf{a}}_{2k+1}$ , we have

$$\mathbb{E}_k \mathbf{V}_{k+1} - \mathbf{V}_k = \frac{\tilde{\tau}}{n}[\mathbf{w}_k \tilde{\mathbf{g}}_k^\top + \mathbf{V}_k \mathbf{L}_k], \quad (\text{S-27})$$

where  $\tilde{\mathbf{g}}_k$  and  $\mathbf{L}_k$  are defined in (S-7) and (S-8), respectively. Multiplying  $\mathbf{U}^\top$  from the left on the both sides of the above equation, we have

$$\mathbb{E}_k \mathbf{P}_{k+1} - \mathbf{P}_k = \frac{\tilde{\tau}}{n}[\mathbf{q}_k \tilde{\mathbf{g}}_k^\top + \mathbf{P}_k \mathbf{L}_k],$$

which implies (S-21). In fact, there is no higher-order term in (S-21), and the left hand side of (S-21) is exactly zero.

Then, we prove (S-22). From (S-1), we have

$$\mathbf{w}_{k+1} - \mathbf{w}_k = \frac{\tau}{n}[\mathbf{y}_k f(\mathbf{y}_k^\top \mathbf{w}_k) - \tilde{\mathbf{y}}_{2k} \tilde{f}(\tilde{\mathbf{y}}_{2k}^\top \mathbf{w}_k) - \lambda \mathbf{w}_k \text{diag}(H'(z_k))], \quad (\text{S-28})$$

where  $\mathbf{y}_k = \mathbf{U} \mathbf{c}_k + \sqrt{\eta_T} \mathbf{a}_k$  and  $\tilde{\mathbf{y}}_{2k} = \mathbf{V}_k \tilde{\mathbf{c}}_{2k} + \sqrt{\eta_G} \tilde{\mathbf{a}}_{2k}$ . Averaging both sides of the above equation over  $\mathbf{c}_k$ ,  $\mathbf{a}_k \tilde{\mathbf{c}}_{2k}$  and  $\tilde{\mathbf{a}}_{2k}$ , we have

$$\begin{aligned} \mathbb{E}_k \mathbf{w}_{k+1} - \mathbf{w}_k = & \frac{\tau}{n} \left[ \mathbf{U} \mathbf{g}_k + \left\langle \mathbf{a}_k f(\mathbf{c}_k^\top \mathbf{q}_k + \sqrt{\eta_T} \mathbf{a}_k^\top \mathbf{w}_k) \right\rangle \right. \\ & \left. - \mathbf{V}_k \tilde{\mathbf{g}}_k - \left\langle \tilde{\mathbf{a}}_{2k} \tilde{f}(\tilde{\mathbf{c}}_{2k}^\top \mathbf{r}_k + \sqrt{\eta_G} \tilde{\mathbf{a}}_{2k}^\top \mathbf{w}_k) \right\rangle - \lambda \mathbf{w}_k \text{diag}(H'(z_k)) \right]. \end{aligned}$$

Multiplying  $\mathbf{U}^\top$  from the left on the both sides of the above equation, we have

$$\begin{aligned} \mathbb{E}_k \mathbf{q}_{k+1} - \mathbf{q}_k = & \frac{\tau}{n} \left[ \mathbf{g}_k - \mathbf{P}_k \tilde{\mathbf{g}}_k + \sqrt{\eta_T} \left\langle \mathbf{U}^\top \mathbf{a}_k f(\mathbf{c}_k^\top \mathbf{q}_k + \sqrt{\eta_T} \mathbf{a}_k^\top \mathbf{w}_k) \right\rangle_{\mathbf{c}, \mathbf{a}} \right. \\ & \left. - \sqrt{\eta_G} \left\langle \mathbf{U}^\top \tilde{\mathbf{a}} \tilde{f}(\tilde{\mathbf{c}}^\top \mathbf{r}_k + \sqrt{\eta_G} \tilde{\mathbf{a}}^\top \mathbf{w}_k) \right\rangle_{\tilde{\mathbf{c}}, \tilde{\mathbf{a}}} - \lambda \mathbf{q}_k \text{diag}(H'(z_k)) \right] \quad (\text{S-29}) \end{aligned}$$

We note that  $\begin{bmatrix} \mathbf{U}^\top \mathbf{a}_k \\ \mathbf{w}_k^\top \mathbf{a}_k \end{bmatrix}$  are Gaussian random vector with zero-mean and covariance matrix  $\begin{bmatrix} \mathbf{I} & \mathbf{q}_k \\ \mathbf{q}_k^\top & z_k \end{bmatrix}$ .

We can rewrite

$$\begin{aligned} \left\langle \mathbf{U}^\top \mathbf{a}_k f(\mathbf{c}_k^\top \mathbf{q}_k + \sqrt{\eta_T} \mathbf{a}_k^\top \mathbf{w}_k) \right\rangle_{\mathbf{c}, \mathbf{a}} &= z_k^{-1/2} \mathbf{U}^\top \mathbf{w}_k \left\langle e f(\mathbf{c}_k^\top \mathbf{q}_k + \sqrt{z_k \eta_T} e) \right\rangle_{\mathbf{c}, e} \\ &= \sqrt{\eta_T} \mathbf{q}_k \left\langle f'(\mathbf{c}_k^\top \mathbf{q}_k + \sqrt{z_k \eta_T} e) \right\rangle_{\mathbf{c}, e}, \end{aligned} \quad (\text{S-30})$$

where the second line is due to Stein's lemma (i.e., integral by part for Gaussian random variable.) Similarly, we have

$$\left\langle \mathbf{U}^\top \tilde{\mathbf{a}} \tilde{f}(\tilde{\mathbf{c}}^\top \mathbf{r}_k + \sqrt{\eta_G} \tilde{\mathbf{a}}^\top \mathbf{w}_k) \right\rangle_{\tilde{\mathbf{c}}, \tilde{\mathbf{a}}} = \sqrt{\eta_G} \mathbf{q}_k \left\langle \tilde{f}'(\tilde{\mathbf{c}}^\top \mathbf{r}_k + \sqrt{z_k \eta_G} e) \right\rangle_{\tilde{\mathbf{c}}, e}. \quad (\text{S-31})$$

Substituting (S-30) and (S-31) into (S-29), we get

$$\mathbb{E}_k \mathbf{q}_{k+1} - \mathbf{q}_k = \frac{\tau}{n}[\mathbf{g}_k - \mathbf{P}_k \tilde{\mathbf{g}}_k + \mathbf{q}_k h_k],$$

where  $\tilde{\mathbf{g}}_k$ ,  $\mathbf{g}_k$ , and  $h_k$  are defined in (S-7), (S-14), and (S-15), respectively. Now, we proved (S-22), which again has no higher-order term.

We next prove (S-23). Note that

$$\begin{aligned}\mathbf{S}_{k+1} - \mathbf{S}_k &= (\mathbf{V}_k + \mathbf{V}_{k+1} - \mathbf{V}_k)^\top (\mathbf{V}_k + \mathbf{V}_{k+1} - \mathbf{V}_k) - \mathbf{S}_k \\ &= \mathbf{V}_k^\top (\mathbf{V}_{k+1} - \mathbf{V}_k) + (\mathbf{V}_{k+1} - \mathbf{V}_k)^\top \mathbf{V}_k + (\mathbf{V}_{k+1} - \mathbf{V}_k)^\top (\mathbf{V}_{k+1} - \mathbf{V}_k).\end{aligned}$$

Averaging both sides of the above equation over  $\tilde{\mathbf{c}}_{2k+1}$  and  $\tilde{\mathbf{a}}_{2k+1}$  and substituting (S-27) into above equation, we have

$$\mathbb{E}_k \mathbf{S}_{k+1} - \mathbf{S}_k = \frac{\tau}{n} [\mathbf{r}_k \tilde{\mathbf{g}}_k^\top + \mathbf{S}_k \mathbf{L}_k + \tilde{\mathbf{g}}_k \mathbf{r}_k^\top + \mathbf{L}_k \mathbf{S}_k] + \frac{\tau^2}{n^2} [\mathbf{w}_k \tilde{\mathbf{g}}_k^\top + \mathbf{V}_k \mathbf{L}_k]^\top [\mathbf{w}_k \tilde{\mathbf{g}}_k^\top + \mathbf{V}_k \mathbf{L}_k]. \quad (\text{S-32})$$

We know that

$$\begin{aligned}\mathbb{E} \|\mathbf{w}_k \tilde{\mathbf{g}}_k^\top + \mathbf{V}_k \mathbf{L}_k\|^\top [\mathbf{w}_k \tilde{\mathbf{g}}_k^\top + \mathbf{V}_k \mathbf{L}_k] &\leq \mathbb{E} \|\mathbf{w}_k \tilde{\mathbf{g}}_k^\top + \mathbf{V}_k \mathbf{L}_k\|^2 \\ &\leq 2z_k \|\tilde{\mathbf{g}}_k\|^2 + 2\|\mathbf{S}_k\| \|\mathbf{L}_k\|^2 \\ &\leq C \mathbb{E} [z_k + \|\mathbf{S}_k\|] \\ &\leq C(T),\end{aligned} \quad (\text{S-33})$$

where  $\tilde{\mathbf{g}}_k$ ,  $\mathbf{L}_k$  are defined in (S-7) and (S-8), respectively. The third line of the above inequalities is due to the fact that  $\tilde{f}$  and  $H'$  are uniformly bounded, and in reaching the last line, we used Lemma 3. Combining (S-32) and (S-33), we reach (S-23).

The other two inequalities (S-24) and (S-25) can be proved in a similar way. We omit the details here.  $\square$

### Check Condition (C.1)

**Lemma 6.** *Under the same setting as Theorem 1, for all  $k = 0, 1, \dots, \lfloor nT \rfloor$  with a given  $T > 0$ , then*

$$\mathbb{E} \|\mathbf{M}_{k+1} - \mathbb{E}_k \mathbf{M}_{k+1}\|^2 \leq C(T)n^{-2}.$$

*Proof.* Note that  $\mathbb{E} \|\mathbf{M}_{k+1} - \mathbb{E}_k \mathbf{M}_{k+1}\|^2 = \mathbb{E} \|\mathbf{M}_{k+1} - \mathbf{M}_k - \mathbb{E}_k (\mathbf{M}_{k+1} - \mathbf{M}_k)\|^2 \leq \mathbb{E} \|\mathbf{M}_{k+1} - \mathbf{M}_k\|^2$ . It is sufficient to prove

$$\mathbb{E} \|\mathbf{M}_{k+1} - \mathbf{M}_k\|^2 \leq C(T)n^{-2}. \quad (\text{S-34})$$

In what follows, we are going to bound the second-order moment of each element in  $\mathbf{M}_{k+1} - \mathbf{M}_k$ . In particular, we bound the 5 blocks  $\mathbf{P}_k$ ,  $\mathbf{S}_k$ ,  $\mathbf{q}_k$ ,  $z_k$  and  $\mathbf{r}_k$  of  $\mathbf{M}_k$  separately.

We first bound  $\mathbb{E} \|\mathbf{P}_{k+1} - \mathbf{P}_k\|^2$ . Multiplying  $\mathbf{U}^\top$  from left on both sides of (S-26), we have

$$\mathbf{P}_{k+1} - \mathbf{P}_k = \frac{\tau}{n} [\mathbf{q}_k \tilde{\mathbf{c}}_{2k+1}^\top \tilde{f}(\tilde{\mathbf{c}}_{2k+1}^\top \mathbf{V}_k^\top \mathbf{w}_k + \eta_G \tilde{\mathbf{a}}_{2k+1}^\top \mathbf{w}_k) - \lambda \mathbf{P}_k \text{diag}(H'(\mathbf{V}_k^\top \mathbf{V}_k))]$$

We then get

$$\begin{aligned}\mathbb{E} \|\mathbf{P}_{k+1} - \mathbf{P}_k\|^2 &\leq Cn^{-2} \mathbb{E} [\|\mathbf{q}_k\|^2 \mathbb{E}_k \|\tilde{\mathbf{c}}_{2k+1}\|^2 + \|\mathbf{P}_k\|^2] \\ &\leq Cn^{-2} \mathbb{E} [1 + \|\mathbf{q}_k\|^2 + \|\mathbf{P}_k\|^2] \\ &\leq C(T)n^{-2}.\end{aligned} \quad (\text{S-35})$$

Here the last line is due to Lemma 3.

We next bound  $\mathbb{E} \|\mathbf{q}_{k+1} - \mathbf{q}_k\|^2$  in the same way. Specifically, multiplying  $\mathbf{U}^\top$  from the left on both sides of (S-28), we get

$$\mathbf{q}_{k+1} - \mathbf{q}_k = \frac{\tau}{n} [\mathbf{U}^\top \mathbf{y}_k f(\mathbf{y}_k^\top \mathbf{w}_k) - \mathbf{U}^\top \tilde{\mathbf{y}}_{2k} \tilde{f}(\tilde{\mathbf{y}}_{2k}^\top \mathbf{w}_k) - \lambda \mathbf{q}_k \text{diag}(H'(\mathbf{w}_k^\top \mathbf{w}_k))].$$

We then have

$$\begin{aligned}
& \mathbb{E} \|\mathbf{q}_{k+1} - \mathbf{q}_k\|^2 \\
& \leq \frac{\tau^2}{n^2} \mathbb{E} [\|\mathbf{c}_k\|^2 f_k^2 + \|\mathbf{U}^\top \mathbf{a}_k\|^2 f_k^2 + \|\mathbf{P}_k\|^2 \|\tilde{\mathbf{c}}_{2k}\|^2 \tilde{f}_{2k}^2 + \|\mathbf{U}^\top \tilde{\mathbf{a}}_{2k}\|^2 \tilde{f}_{2k}^2 + \|\mathbf{q}_k\|^2 h_k^2] \\
& \leq Cn^{-2} [1 + \sqrt{\mathbb{E} \|\mathbf{U}^\top \mathbf{a}_k\|^4} \sqrt{\mathbb{E} f_k^4} + \sqrt{\mathbb{E} \|\mathbf{U}^\top \tilde{\mathbf{a}}_{2k}\|^4} \sqrt{\mathbb{E} \tilde{f}_{2k}^4} + \mathbb{E} z_k^2 + \mathbb{E} \|\mathbf{S}_k\|^2] \\
& \leq Cn^{-2} [1 + \mathbb{E} z_k^2 + \mathbb{E} \|\mathbf{S}_k\|^2] \\
& \leq C(T)n^{-2},
\end{aligned} \tag{S-36}$$

where  $f_k$  and  $\tilde{f}_{2k}$  are shorthands for  $f(\mathbf{y}_k^\top \mathbf{w}_k)$  and  $\tilde{f}(\tilde{\mathbf{y}}_{2k}^\top \mathbf{w}_k)$  respectively. In reaching the last line, we used Lemma 3 again.

Similarly, we can also prove that

$$\begin{aligned}
\mathbb{E} \|\mathbf{S}_{k+1} - \mathbf{S}_k\|^2 & \leq C(T)n^{-2} \\
\mathbb{E} (z_{k+1} - z_k)^2 & \leq C(T)n^{-2} \\
\mathbb{E} \|\mathbf{r}_{k+1} - \mathbf{r}_k\|^2 & \leq C(T)n^{-2}.
\end{aligned} \tag{S-37}$$

Combining (S-35), (S-36) and (S-37), we can prove (S-34), which concludes the whole proof.  $\square$

### S-IV.3 Proof of Lemma 2

Before proving Lemma 2, we first present and prove the following lemma. Let  $\mathbf{u}_i$  and  $\mathbf{v}_{k,i}$  denote the  $i$ th row vectors of  $\mathbf{U}$  and  $\mathbf{V}_k$  in column view, respectively, and let  $w_{k,i}$  be the  $i$ th element of the vector  $\mathbf{w}_k$ .

**Lemma 7.** *Under the same setting as Theorem 1, for all  $k = 0, 1, \dots, \lfloor nT \rfloor$  with a given  $T > 0$ , then*

$$\|\mathbb{E}_k \mathbf{v}_{k+1,i} - \mathbf{v}_{k,i}\| \leq Cn^{-1} (\|\mathbf{v}_{k,i}\| + |w_{k,i}|) \tag{S-38}$$

$$|\mathbb{E}_k w_{k,i} - w_{k,i}| \leq Cn^{-1} (\|\mathbf{u}_i\| + \|\mathbf{v}_{k,i}\| + |w_{k,i}|). \tag{S-39}$$

In the proof of this lemma and Lemma 2, we omit the two constants  $\eta_T$  and  $\eta_G$  for simplicity.

*Proof.* From (S-2) and knowing that the function  $\tilde{f}$  and  $H'$  are uniformly bounded, we can immediately prove (S-38).

Next, we are going to prove (S-39). From (S-1), we know

$$\begin{aligned}
& |\mathbb{E}_k w_{k+1,i} - w_{k,i}| \\
& \leq \frac{\tau}{n} \left( \left| \mathbf{u}_i^\top \left\langle \mathbf{c}_k f(\mathbf{y}_k^\top \mathbf{w}_k) \right\rangle_{\mathbf{c}_k, \mathbf{a}_k} \right| + \left| \left\langle a_{k,i} f(\mathbf{y}_k^\top \mathbf{w}_k) \right\rangle_{\mathbf{c}_k, \mathbf{a}_k} \right| \\
& \quad + \left| \mathbf{v}_{k,i}^\top \left\langle \tilde{\mathbf{c}}_{2k} \tilde{f}(\tilde{\mathbf{y}}_{2k}^\top \mathbf{w}_k) \right\rangle_{\tilde{\mathbf{c}}_{2k}, \tilde{\mathbf{a}}_{2k}} \right| + \left| \left\langle \tilde{a}_{2k,i} \tilde{f}(\tilde{\mathbf{y}}_{2k}^\top \mathbf{w}_k) \right\rangle_{\tilde{\mathbf{c}}_{2k}, \tilde{\mathbf{a}}_{2k}} \right| + \lambda |w_{k,i} H'(\mathbf{w}_k^\top \mathbf{w}_k)| \Big) \\
& \leq Cn^{-1} \left( \|\mathbf{u}_i\| + \|\mathbf{v}_{k,i}\| + |w_{k,i}| + \left| \left\langle a_{k,i} f(\mathbf{y}_k^\top \mathbf{w}_k) \right\rangle_{\mathbf{c}_k, \mathbf{a}_k} \right| + \left| \left\langle \tilde{a}_{2k,i} \tilde{f}(\tilde{\mathbf{y}}_{2k}^\top \mathbf{w}_k) \right\rangle_{\tilde{\mathbf{c}}_{2k}, \tilde{\mathbf{a}}_{2k}} \right| \right),
\end{aligned} \tag{S-40}$$

where the last is due to the fact that  $H'$ ,  $f$  and  $\tilde{f}$  are uniformly bounded. Using Taylor's expansion up-to zero-order

$$\begin{aligned}
f(\mathbf{y}_k^\top \mathbf{w}_k) &= f(\mathbf{q}_k^\top \mathbf{c}_k + \sum_{j \neq i} w_{k,j} a_{k,j} + w_{k,i} a_{k,i}) \\
&= f(\mathbf{q}_k^\top \mathbf{c}_k + \sum_{j \neq i} w_{k,j} a_{k,j}) + f'(\mathbf{q}_k^\top \mathbf{c}_k + \sum_{j \neq i} w_{k,j} a_{k,j} + \chi_{k,i}) w_{k,i} a_{k,i},
\end{aligned}$$

with  $\chi_{k,i}$  being some number such that  $|\chi_{k,i}| \leq |w_{k,i} a_{k,i}|$ , we have

$$\begin{aligned}
& \left| \left\langle a_{k,i} f(\mathbf{y}_k^\top \mathbf{w}_k) \right\rangle_{\mathbf{c}_k, \mathbf{a}_k} \right| \\
& \leq \left| \left\langle f(\mathbf{q}_k^\top \mathbf{c}_k + \sum_{j \neq i} w_{k,j} a_{k,j}) a_{k,i} \right\rangle_{\mathbf{c}_k, \mathbf{a}_k} \right| + \left| \left\langle f'(\mathbf{q}_k^\top \mathbf{c}_k + \sum_{j \neq i} w_{k,j} a_{k,j} + \chi_{k,i}) w_{k,i} a_{k,i}^2 \right\rangle_{\mathbf{c}_k, \mathbf{a}_k} \right| \\
& = \left| \left\langle f'(\mathbf{q}_k^\top \mathbf{c}_k + \sum_{j \neq i} w_{k,j} a_{k,j} + \chi_{k,i}) w_{k,i} a_{k,i}^2 \right\rangle_{\mathbf{c}_k, \mathbf{a}_k} \right| \\
& \leq C |w_{k,i}|.
\end{aligned} \tag{S-41}$$

The second line is due to the fact  $a_{k,i}$  is zero-mean, and in reaching the last line, we used the boundness of  $f'$ . Similarly, we can get

$$\left| \left\langle \tilde{a}_{2k,i} \tilde{f}(\tilde{\mathbf{y}}_{2k}^\top \mathbf{w}_k) \right\rangle_{\tilde{\mathbf{c}}_{2k}, \tilde{\mathbf{a}}_{2k}} \right| \leq C |w_{k,i}|. \tag{S-42}$$

Substituting (S-41) and (S-42) into (S-40), we prove (S-40).  $\square$

Now we are in the position to prove Lemma 2.

*Proof of Lemma 2.* Because of the exchangeability,  $\mathbb{E} w_{k,i}^4 = \mathbb{E} w_{k,j}^4$ , and  $\mathbb{E} [\mathbf{V}_k]_{i,\ell}^4 = \mathbb{E} [\mathbf{V}_k]_{j,\ell}^4$  for all  $i, j = 1, 2, \dots, n$  and  $\ell = 1, 2, \dots, d$ . Thus, we only need to prove (S-20) for any specific  $i$ .

We first prove  $\mathbb{E} w_{k,i}^4 \leq C(T)n^{-2}$ . We know that

$$\begin{aligned}
\mathbb{E} w_{k+1,i}^4 - \mathbb{E} w_{k,i}^4 &= 4\mathbb{E} \left[ w_{k,i}^3 \mathbb{E}_k (w_{k+1,i} - w_{k,i}) \right] + 6\mathbb{E} \left[ w_{k,i}^2 \mathbb{E}_k (w_{k+1,i} - w_{k,i})^2 \right] \\
&\quad + 4\mathbb{E} \left[ w_{k,i} \mathbb{E}_k (w_{k+1,i} - w_{k,i})^3 \right] + \mathbb{E} \mathbb{E}_k (w_{k+1,i} - w_{k,i})^4.
\end{aligned} \tag{S-43}$$

From (S-1) and knowing that  $h, f$  and  $\tilde{f}$  are uniformly bounded, we have

$$\mathbb{E}_k (w_{k+1,i} - w_{k,i})^\gamma \leq \frac{C}{n^\gamma} \left( 1 + \|\mathbf{u}_i\|^\gamma + \|\mathbf{v}_{k,i}\|^\gamma + |w_{k,i}|^\gamma \right) \quad \text{for } \gamma = 2, 3, 4. \tag{S-44}$$

Substituting (S-39) and (S-44) into (S-43) and using the Young's inequality, we have

$$\begin{aligned}
\mathbb{E} w_{k+1,i}^4 - \mathbb{E} w_{k,i}^4 &\leq \frac{C}{n} \left( n^{-2} + \mathbb{E} \|\mathbf{u}_i\|^4 + \mathbb{E} \|\mathbf{v}_{k,i}\|^4 + \mathbb{E} w_{k,i}^4 \right) \\
&\leq \frac{C}{n} \mathbb{E} \left( n^{-2} + \sum_{\ell=1}^d [\mathbf{V}_k]_{i,\ell}^4 + w_{k,i}^4 \right),
\end{aligned} \tag{S-45}$$

where the last line is due to Assumption A.4), which implies  $\sum_{\ell} [\mathbf{U}]_{i,\ell}^4 \leq C$ . Similarly, we can prove

$$\sum_{\ell=1}^d \mathbb{E} \left( [\mathbf{V}_{k+1}]_{i,\ell}^4 - [\mathbf{V}_k]_{i,\ell}^4 \right) \leq \frac{C}{n} \mathbb{E} \left( n^{-2} + \sum_{\ell=1}^d [\mathbf{V}_k]_{i,\ell}^4 + w_{k,i}^4 \right). \tag{S-46}$$

Combining (S-45) and (S-46), we have

$$\mathbb{E} (w_{k+1,i}^4 + \sum_{\ell=1}^d [\mathbf{V}_{k+1}]_{i,\ell}^4) - \mathbb{E} (w_{k,i}^4 + \sum_{\ell=1}^d [\mathbf{V}_k]_{i,\ell}^4) \leq \frac{C}{n} \left[ n^{-2} + \mathbb{E} (w_{k,i}^4 + \sum_{\ell=1}^d [\mathbf{V}_k]_{i,\ell}^4) \right].$$

Using the above inequality iteratively, we have

$$\mathbb{E} \left( w_{k,i}^4 + \sum_{\ell=1}^d [\mathbf{V}_k]_{i,\ell}^4 \right) \leq \left( n^{-2} + w_{0,i}^4 + \sum_{\ell=1}^d [\mathbf{V}_0]_{i,\ell}^4 \right) e^{\frac{k}{n} C}.$$

Since  $\mathbb{E} (w_{0,i}^4 + \sum_{\ell=1}^d [\mathbf{V}_0]_{i,\ell}^4)$  are bounded in Assumption A.4), we now reach (S-20).  $\square$

## S-V Local stability analysis of the fixed points of the ODE

In this section, we provide additional details on the local stability analysis of the ODE for Example 1. We first its simplified ODE (13) in the main text. Then, we provide the derivation of the local stability analysis when  $d = 1$ , where the main results are summarized in Section S-I. Finally, we establish the proof of Claim 1 in the main text.

### S-V.1 Derive the reduced ODE for Example 1 when $\lambda \rightarrow \infty$

In Example 1,  $f(x) = \tilde{f}(x) = x$ . Plugging back to (9), we obtain that

$$\begin{aligned} \mathbf{g}_t &= \Lambda \mathbf{q}_t \\ \tilde{\mathbf{g}}_t &= \tilde{\Lambda} \mathbf{r}_t \\ b_t &= \eta_{\Gamma}(\mathbf{q}_t^{\top} \Lambda \mathbf{q}_t + \eta_{\Gamma} z_t) + \eta_{\mathrm{G}}(\mathbf{r}_t^{\top} \tilde{\Lambda} \mathbf{r}_t + \eta_{\mathrm{G}} z_t). \end{aligned} \quad (\text{S-47})$$

Correspondingly, ODE in (8) becomes:

$$\begin{aligned} \frac{d}{dt} \mathbf{P}_t &= \tilde{\tau}(\mathbf{q}_t \tilde{\mathbf{r}}_t^{\top} \tilde{\Lambda} + \mathbf{P}_t \mathbf{L}_t) \\ \frac{d}{dt} \mathbf{q}_t &= \tau(\Lambda \mathbf{q}_t - \mathbf{P}_t \tilde{\Lambda} \mathbf{r}_t + \mathbf{q}_t h_t) \\ \frac{d}{dt} \mathbf{r}_t &= \tau(\mathbf{P}_t^{\top} \Lambda \mathbf{q}_t - \mathbf{S}_t \tilde{\Lambda} \mathbf{r}_t + \mathbf{r}_t h_t) + \tilde{\tau}(\tilde{\Lambda} \mathbf{r}_t + \mathbf{L}_t \mathbf{r}_t) \\ \frac{d}{dt} \mathbf{S}_t &= \tilde{\tau}(\mathbf{r}_t \mathbf{r}_t^{\top} \tilde{\Lambda}^{\top} + \tilde{\Lambda} \mathbf{r}_t \mathbf{r}_t^{\top} + \mathbf{S}_t \mathbf{L}_t + \mathbf{L}_t^{\top} \mathbf{S}_t) \\ \frac{d}{dt} z_t &= 2\tau(\mathbf{q}_t^{\top} \Lambda \mathbf{q}_t - \mathbf{r}_t^{\top} \tilde{\Lambda} \mathbf{r}_t + z_t h_t) \\ &\quad + \tau^2[\eta_{\Gamma}(\mathbf{q}_t^{\top} \Lambda \mathbf{q}_t + z_t \eta_{\Gamma}) + \eta_{\mathrm{G}}(\mathbf{r}_t^{\top} \tilde{\Lambda} \mathbf{r}_t + z_t \eta_{\mathrm{G}})] \end{aligned} \quad (\text{S-48})$$

The first four equations are exactly (13). From last two equations of (S-48), by setting  $\frac{d}{dt} \text{diag}\{\mathbf{S}_t\} = \mathbf{0}$ ,  $\frac{d}{dt} z_t = 0$ ,  $\text{diag}(\mathbf{S}_t) = \mathbf{I}$  and  $z_t = 1$ , we can get (14).

### S-V.2 A complete study of all fixed points when $d = 1$

We next provide the local stability analysis of the fixed points of the ODE (13). When  $d = 1$  and  $\lambda \rightarrow \infty$ , the macroscopic state is described by only 3 scalars,  $P_t$ ,  $q_t$  and  $r_t$ . The result is summarized in Table 1. For the sake of simplicity, we only consider the case  $\Lambda = \tilde{\Lambda}$ , and set  $\eta_{\Gamma} = \eta_{\mathrm{G}} = 1$ , but all analysis can be extended to general cases.

The fixed points are given by the condition  $\frac{d}{dt} P_t = \frac{d}{dt} q_t = \frac{d}{dt} r_t = 0$ . From (13), we get

$$\begin{cases} \tilde{\tau} \Lambda r (q - r P) = 0 \\ \tau [\Lambda - \tau - \Lambda (1 + \frac{\tau}{2}) q^2] q - \tau \Lambda [P + (\frac{\tau}{2} - 1) r q] r = 0 \\ \tau \Lambda P q + [\Lambda (\tilde{\tau} - \tau) - \tau^2] r + \Lambda (\tau - \tilde{\tau} - \frac{\tau^2}{2}) r^3 - \tau \Lambda (1 + \frac{\tau}{2}) r q^2 = 0, \end{cases} \quad (\text{S-49})$$

where  $P, q, r$  are the stationary macroscopic state. The local stability of a fixed point is identified by whether the Jacobian matrix

$$J(P, q, r) \stackrel{\text{def}}{=} \begin{bmatrix} \frac{\partial}{\partial P} g_1 & \frac{\partial}{\partial q} g_1 & \frac{\partial}{\partial r} g_1 \\ \frac{\partial}{\partial P} g_3 & \frac{\partial}{\partial q} g_3 & \frac{\partial}{\partial r} g_3 \\ \frac{\partial}{\partial P} g_5 & \frac{\partial}{\partial q} g_5 & \frac{\partial}{\partial r} g_5 \end{bmatrix}$$

has eigenvalue with non-negative real part or not, where  $g_1 = \tilde{\tau} \Lambda r (q - r P)$ ,  $g_2 = \tau [\Lambda - \tau - \Lambda (1 + \frac{\tau}{2}) q^2] q - \tau \Lambda [P + (\frac{\tau}{2} - 1) r q] r$  and  $g_5 = \tau \Lambda P q + [\Lambda (\tilde{\tau} - \tau) - \frac{(\eta_{\Gamma} + \eta_{\mathrm{G}}) \tau^2}{2}] r + \Lambda (\tau - \tilde{\tau} - \frac{\tau^2}{2}) r^3 - \tau \Lambda (1 + \frac{\tau}{2}) r q^2$ .

#### Type (1) fixed point at $P = q = r = 0$

It is easy to verify that  $q = r = 0$  and any  $P \in [-1, 1]$  is a solution of (S-49), but we first consider  $P = 0$ .

The Jacobian at  $P = q = r = 0$  is

$$J(0, 0, 0) = \begin{bmatrix} 0 & 0 & 0 \\ 0 & \tau(\Lambda - \tau) & 0 \\ 0 & 0 & \Lambda(\tilde{\tau} - \tau) - \tau^2 \end{bmatrix}.$$

Thus, type (1) fixed point is stable if and only if

$$\tau \geq \Lambda \quad \text{and} \quad \frac{\tilde{\tau}}{\tau} \leq \frac{\tau + \Lambda}{\Lambda}.$$

**Type (2) fixed points at  $P = q = 0, r = \pm r^* \neq 0$**

We first analyze when such fixed point exists and then study its local stability.

If  $P = q = 0$ , the first two equations in (S-49) trivially hold. The third equation becomes

$$\tau[\Lambda(r^2 - 1) - \frac{\tau}{2}(\Lambda r^2 + 2)] - \tilde{\tau}\Lambda(r^2 - 1) = 0.$$

The solution is

$$r^2 = \frac{\tau - \tilde{\tau} + \tau^2/\Lambda}{\tau - \tilde{\tau} - \tau^2/2}. \quad (\text{S-50})$$

Since only the positive solution corresponds a fixed one. Thus, type (2) fixed point exists if

$$\frac{\tilde{\tau}}{\tau} \leq 1 - \frac{\tau}{2} \quad (\text{S-51})$$

$$\text{or } \frac{\tilde{\tau}}{\tau} \geq \frac{\tau + \Lambda}{\Lambda}. \quad (\text{S-52})$$

Next, we investigate the local stability of this fixed point. The Jacobian at  $\tilde{q} = q = 0$  for a given  $r$  is

$$J(0, 0, r) = \begin{bmatrix} -\tilde{\tau}\Lambda r^2 & \tilde{\tau}\Lambda r & 0 \\ -\tau\Lambda r & \tau(\Lambda - \tau) - \Lambda\tau(\frac{\tau}{2} - 1)r^2 & 0 \\ 0 & 0 & 3r^2\Lambda(\tau - \frac{\tau^2}{2} - \tilde{\tau}) - \tau^2 + \Lambda(\tilde{\tau} - \tau) \end{bmatrix} \quad (\text{S-53})$$

Plugging (S-50) into  $[J(0, 0, r)]_{3,3}$  of (S-53), then  $[J(0, 0, r)]_{3,3} \leq 0$  implies

$$\frac{\tilde{\tau}}{\tau} \geq \frac{\tau}{\Lambda} + 1.$$

It indicates that the stationary points at the region (S-51) are always unstable. Thus, we only need to consider the second region specified by (S-52).

For the upper-left  $2 \times 2$  sub-matrix of (S-53), the eigenvalues are non-positive if and only if

$$-\tilde{\tau}\Lambda r^2 + \tau(\Lambda - \tau) - \Lambda\tau(\frac{\tau}{2} - 1)r^2 \leq 0 \quad (\text{S-54})$$

$$\tau + \Lambda(\frac{\tau}{2} - 1)r^2 + \Lambda - \Lambda \geq 0. \quad (\text{S-55})$$

Plugging (S-50) into (S-54), we can get

$$\frac{\tilde{\tau}}{\tau} \geq 2. \quad (\text{S-56})$$

Plugging (S-50) into (S-55) and combining (S-52), we can get

$$[\tau + \Lambda(\frac{\tau}{2} - 1)]\tilde{\tau} \geq \tau\Lambda(\frac{\tau}{2} - 1).$$

Solving this inequality implies that

$$\frac{\tilde{\tau}}{\tau} \leq \frac{(\frac{\tau}{2} - 1)\Lambda}{(\frac{\tau}{2} - 1)\Lambda + \tau}, \text{ when } \tau < \frac{2\Lambda}{\Lambda + 2} \quad (\text{S-57})$$

and

$$\frac{\tilde{\tau}}{\tau} \geq \frac{(\frac{\tau}{2} - 1)\Lambda}{(\frac{\tau}{2} - 1)\Lambda + \tau}, \text{ when } \tau > \frac{2\Lambda}{\Lambda + 2}. \quad (\text{S-58})$$

Note that (S-58) is included by (S-56), as  $\frac{(\frac{\tau}{2} - 1)\Lambda}{(\frac{\tau}{2} - 1)\Lambda + \tau} \leq 2$  when  $\tau > \frac{2\Lambda}{\Lambda + 2}$ .

Then, combining (S-52), (S-56), and (S-57) we obtain the stability region for  $\tilde{q} = q = 0$ ,

$$\frac{\tilde{\tau}}{\tau} \geq 1 + \frac{\tau}{\Lambda}, \frac{\tilde{\tau}}{\tau} \geq 2, \text{ and } \frac{\tilde{\tau}}{\tau} \leq \beta(\tau),$$

where  $\beta(\tau)$  is defined as

$$\beta(\tau) \stackrel{\text{def}}{=} \begin{cases} \frac{(\frac{\tau}{2} - 1)\Lambda}{(\frac{\tau}{2} - 1)\Lambda + \tau} & \text{if } \tau \leq \frac{2\Lambda}{\Lambda + 2} \\ +\infty & \text{otherwise.} \end{cases}$$



**Type (3) fixed points at  $q = r = 0$  and  $|P| \in (0, 1]$**

As mentioned, we can check that  $q = r = 0$  and any  $P \in [-1, 1]$  is a solution of (S-49). We next investigate the stable region for the fixed point  $P = \pm 1$  and  $q = r = 0$ , which represents the perfect recovery state. For general  $P$ , we can analyze its fixed point similarly.

The Jacobian at  $q = r = 0$  for any given  $P$  is

$$J(1, 0, 0) = \begin{bmatrix} 0 & 0 & 0 \\ 0 & \tau(\Lambda - \tau) & -\tau\Lambda \\ 0 & \tau\Lambda & \Lambda(\tilde{\tau} - \tau) - \tau^2 \end{bmatrix}.$$

In this case,  $J(1, 0, 0)$  always has an eigenvalue 0 and to calculate the rest two eigenvalues, we only need to analyze the bottom-right  $2 \times 2$  sub-matrix of  $J(\tilde{q})$ . The characteristic polynomial of this sub-matrix is  $f(\lambda) = \lambda^2 - (a + d)\lambda + ad - bc$ , where  $a = \tau(\Lambda - \tau)$ ,  $b = -\tau\Lambda$ ,  $c = \tau\Lambda$ , and  $d = \Lambda(\tilde{\tau} - \tau) - \tau^2$ . The roots of  $f(\lambda) = 0$  both have non-positive real part if and only if  $a + d \leq 0$ ,  $ad - bc \geq 0$ , which implies

$$\frac{\tilde{\tau}}{\tau} \leq \frac{2\tau}{\Lambda} \quad \text{and} \quad \frac{\tilde{\tau}}{\tau}(\tau - \Lambda) \leq \frac{\tau^2}{\Lambda}. \quad (\text{S-59})$$

Noting that when  $\tau < \Lambda$ , the second inequality always hold, and when  $\tau > \Lambda$ ,  $\frac{\tau^2}{\Lambda(\tau - \Lambda)} \geq 4$ , we can combine the two inequalities in (S-59) into compact form

$$\frac{\tilde{\tau}}{\tau} \leq \min\left\{\frac{2\tau}{\Lambda}, \max\left\{\frac{\tau^2}{\Lambda|\tau - \Lambda|}, 4\right\}\right\}.$$

The stable regions of the fixed points for  $q = r = 0$  and  $|P| < 1$  can be derived in a similar way, which turns out to be a subset of the stable region for  $P = \pm 1$ .

**Type (4) fixed point at  $P = r = 0$  and  $q \neq 0$ .**

From (S-49), we know when at fixed point,  $\tilde{q} = r = 0$ , then  $q^2 = \frac{\Lambda - \tau}{\Lambda(1 + \tau/2)}$ , so  $\tau$  must satisfy  $\tau \leq \Lambda$ . The corresponding Jacobian is:

$$J(0, 0, q) = \begin{bmatrix} 0 & 0 & \tilde{\tau}\Lambda q \\ 0 & \tau(\Lambda - \tau) - 3\tau\Lambda q^2(1 + \frac{\tau}{2}) & 0 \\ \tau\Lambda q & 0 & (\tilde{\tau} - \tau)\Lambda - \tau^2 - \tau\Lambda q^2(1 + \frac{\tau}{2}) \end{bmatrix}.$$

After plugging in  $q^2 = \frac{\Lambda - \tau}{\Lambda(1 + \tau/2)}$ , we can obtain that the characteristic function  $\det(\lambda \mathbf{I} - J(0, 0, q))$  is equal to:

$$\det(\lambda \mathbf{I} - J(0, 0, q)) = [\lambda + 2\tau(\Lambda - \tau)][\lambda(\lambda + (2\tau - \tilde{\tau})\Lambda) - \tau\tilde{\tau}\Lambda^2 q^2]$$

Clearly,  $\det(\lambda \mathbf{I} - J(0, 0, q)) = 0$  has a non-negative root, so  $J(0, 0, q)$  always has a non-negative eigenvalue. This means type (4) fixed points are always unstable.

**Type (5) fixed points at  $P, q, r \neq 0$**

The fixed points equation (S-49) can also have solutions that none of  $P$ ,  $q$  and  $r$  is zero. In what follows, we derive the analytical expression of this type of solutions. It turns out that there can be maximum 8 solutions, which are symmetric by flipping the signs. We are unable to derive the analytical expression for their stable region, but it can be computed numerically.

If  $P, q, r \neq 0$ , (S-49) yields

$$r = \frac{q}{P} \quad (\text{S-60})$$

$$\Lambda - \tau - \Lambda(1 + \frac{\tau}{2})q^2 - \Lambda[\frac{P}{q} + (\frac{\tau}{2} - 1)r]r = 0 \quad (\text{S-61})$$

$$\tau\Lambda\tilde{P}q + r[\Lambda(\tilde{\tau} - \tau) - \tau^2] + r^3\Lambda(\tau - \tilde{\tau} - \frac{\tau^2}{2}) - rq^2\tau\Lambda(1 + \frac{\tau}{2}) = 0. \quad (\text{S-62})$$

Plugging (S-60) into (S-61), we can get

$$q^{-2} = -\frac{1}{\tau}[\Lambda(\frac{\tau}{2} - 1)P^{-2} + \Lambda(1 + \frac{\tau}{2})]. \quad (\text{S-63})$$

Then combining (S-60) (S-63) and (S-62), we can obtain the following equations:

$$AP^{-4} + BP^{-2} + C = 0 \quad (\text{S-64})$$

where  $A = \Lambda(\tilde{\tau} - \tau)(\frac{1}{2} - \frac{1}{\tau}) + \tilde{\tau}$ ,  $B = \Lambda[\frac{\tilde{\tau}}{\tau}(1 + \frac{\tau}{2}) - 2]$ ,  $C = \Lambda(1 + \frac{\tau}{2})$ . We can find that (S-64) is an equation of  $P^{-2}$  with at most two roots. Combining (S-63), we know there are at most 2 solutions for the pair  $(q^{-2}, P^{-2})$  and hence there are at most 8 solutions for  $(q, P, r)$ , where  $r = P/q$ .

### S-V.3 Proof of Claim 1

*Proof of Claim 1.* We first compute the Jacobian  $\partial\{\frac{d}{dt}\mathbf{P}_t, \frac{d}{dt}\mathbf{q}_t, \frac{d}{dt}\mathbf{r}_t\}/\partial\{\mathbf{P}_t, \mathbf{q}_t, \mathbf{r}_t\}$  of the ODE (13) when  $\mathbf{q}_t = \mathbf{r}_t = \mathbf{0}$ . In the Jacobian, the  $d \times d$  matrix  $\mathbf{P}_t$  is considered as a  $d^2$  vector. In fact, all elements in the Jacobian matrix related to  $\mathbf{P}_t$  are 0. Specifically, the Jacobian for any  $\mathbf{P}$  and  $\mathbf{q}_t = \mathbf{r}_t = \mathbf{0}$  is

$$\mathbf{J}(\mathbf{P}) = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \tau(\Lambda - \tau\bar{\eta}^2\mathbf{I}_d) & -\tau\mathbf{P}\tilde{\Lambda} \\ \mathbf{0} & \tau\mathbf{P}^\top\Lambda & \tilde{\Lambda}(\tilde{\tau} - \tau) - \tau^2\bar{\eta}^2 \end{bmatrix}, \quad (\text{S-65})$$

where  $\bar{\eta}^2 = (\eta_T^2 + \eta_G^2)/2$ .

When  $\mathbf{P}$  is diagonal, under a suitable column-row permutation, the  $\mathbf{J}(\mathbf{P})$  in (S-65) becomes a block diagonal matrix, where each non-zero block is a  $2 \times 2$  matrix

$$\begin{bmatrix} \tau([\Lambda]_{\ell,\ell} - \tau\bar{\eta}^2) & -\tau[\mathbf{P}]_{\ell,\ell}[\tilde{\Lambda}]_{\ell,\ell} \\ \tau[\mathbf{P}]_{\ell,\ell}[\Lambda]_{\ell,\ell} & [\tilde{\Lambda}]_{\ell,\ell}(\tilde{\tau} - \tau) - \tau^2\bar{\eta}^2 \end{bmatrix} \quad (\text{S-66})$$

for  $\ell = 1, 2, \dots, d$ . Intuitively, the above matrix is the Jacobian matrix of  $\partial\{\frac{d}{dt}[\mathbf{q}_t]_\ell, \frac{d}{dt}[\mathbf{r}_t]_\ell\}/\partial\{[\mathbf{q}_t]_\ell, [\mathbf{r}_t]_\ell\}$ , and the Jacobian  $\partial\{\frac{d}{dt}[\mathbf{q}_t]_\ell, \frac{d}{dt}[\mathbf{r}_t]_\ell\}/\partial\{[\mathbf{q}_t]_{\ell'}, [\mathbf{r}_t]_{\ell'}\}$  is zero for  $\ell \neq \ell'$ .

Now the problem reduces into investigate eigenvalues of  $n$  2-by-2 matrices. For any given  $\ell = 1, 2, \dots, n$ , we have studied this problem in Section S-V.2 (type (1) and type (3) fixed points).

Specifically, the perfect recovery point  $\mathbf{P} = \mathbf{I}$ ,  $\mathbf{q} = \mathbf{r} = \mathbf{0}$  is stable if and only if  $\lambda_{\max}(\mathbf{J}(\mathbf{P})) \leq 0$ , where  $\mathbf{J}(\mathbf{P})$  is defined in (S-65). Similar to the analysis of the type (3) fixed points in Section S-V.2, the condition that both eigenvalues of the matrix in (S-66) is non-positive implies

$$\frac{1}{2}([\Lambda]_{\ell,\ell} - [\tilde{\Lambda}]_{\ell,\ell} + \alpha[\tilde{\Lambda}]_{\ell,\ell}) \leq \tau\bar{\eta}^2 \quad (\text{S-67})$$

$$\text{and } \alpha(\tau\bar{\eta}^2 - [\Lambda]_{\ell,\ell}) \leq \frac{\tau\bar{\eta}^2}{[\Lambda]_{\ell,\ell}}(\tau\bar{\eta}^2 - [\Lambda]_{\ell,\ell} + [\tilde{\Lambda}]_{\ell,\ell}), \quad (\text{S-68})$$

for all  $\ell = 1, 2, \dots, n$ . The inequality (S-67) is the first inequality of (15) in Claim 1 in the main text.

Next, we investigate the condition when the trivial fixed point of the origin  $\mathbf{P} = \mathbf{0}$  and  $\mathbf{q} = \mathbf{r} = \mathbf{0}$  is unstable. Put  $\mathbf{P} = \mathbf{0}$  into (S-66), we get a diagonal matrix

$$\begin{bmatrix} \tau([\Lambda]_{\ell,\ell} - \tau\bar{\eta}^2) & 0 \\ 0 & [\tilde{\Lambda}]_{\ell,\ell}(\tilde{\tau} - \tau) - \tau^2\bar{\eta}^2 \end{bmatrix}.$$

When any eigenvalue of the above matrices for  $\ell = 1, 2, \dots, n$  is positive, this trivial fixed point will be unstable. A sufficient condition is the first eigenvalues of all matrices are positive:

$$\tau\bar{\eta}^2 < [\Lambda]_{\ell,\ell} \text{ for all } \ell = 1, 2, \dots, n. \quad (\text{S-69})$$

The above inequality is the second inequality of (15) in the main text. In addition, (S-69) implies (S-68) hold as the left hand side of (S-68) is negative. Now, we prove that (15) is a sufficient condition that the perfect fixed point is stable and the trivial fixed point is unstable.

□

We further note that (15) is not a necessary condition. There may be a region that (S-69) does not hold, but the origin is still unstable, and the perfect recovery point is stable. Such region is hard to characterize analytically, and numerically, we found the training algorithms always converge to other bad fixed points (e.g. mode collapsing state, or a state that  $\mathbf{P}$  and  $\mathbf{q}$  are still zero, but  $\mathbf{r}$  is non-zero. The situation of the latter is similar to the noninfo-2 phase in the  $d = 1$  case, which converges to the type (2) fixed point). Further study on those bad fixed points will be established in future works under a more general model.