
Supplementary Materials: A Solvable High-Dimensional Model of GAN

Chuang Wang^{1,2}
wangchuang@ia.ac.cn

Hong Hu²
honghu@g.harvard.edu

Yue M. Lu²
yuelu@seas.harvard.edu

1. National 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

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 I_d to denote the $d \times d$ identity matrix. Depending on the context, $\|\cdot\|$ denotes either the ℓ_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 τ and $\tilde{\tau}$ with the fixed model parameters η_T, η_G, Λ , and $\tilde{\Lambda}$.

Phase diagram for $d = 1$. By analyzing the local stabilities of these fixed points as illustrated in Figure 1(a), we obtain the phase diagram as shown in Figure 1(b). For simplicity, we only present the result when $\eta_T = \eta_G = 1$, and $\Lambda = \tilde{\Lambda}$, which is denoted by Λ 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 1(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 V has no correlation with the true feature vector U . In Figure 1(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

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 $\frac{\tilde{\tau}}{\tau} \leq 1 - \frac{\tau}{2}$	$\max\{2, \frac{\tau + \Lambda}{\Lambda}\} \leq \frac{\tilde{\tau}}{\tau} \leq \beta(\tau)$	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

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 1(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 2 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 2 show a case in the Info-1 phase. The bottom red dot in Figure 1.(b) represents this configuration of the step sizes, where $\tilde{\tau}/\tau$ is small. The top figure of Figure 2.(a) shows the dynamics of P_t, q_t and r_t , and the bottom figure shows the phase portrait on $P - q$ plane. Top figure of Figure 2.(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 τ/τ 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 2 show the two types of limiting cycles that can emerge in the oscillating phase. The middle two red dots in Figure 1.(b) represents these configurations of the step sizes. The last column of Figure 2 shows another stable phase in Info-2. In this phase, τ/τ is relatively large. The two time-scale dynamics are mixed, and another type of stable fixed points emerges.

Phase diagram for $d = 2$. Figure 3 shows the phase diagram when $d = 2$. In particular, the two red lines between Info-1 and Noninfo-1 in Figure 3 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.

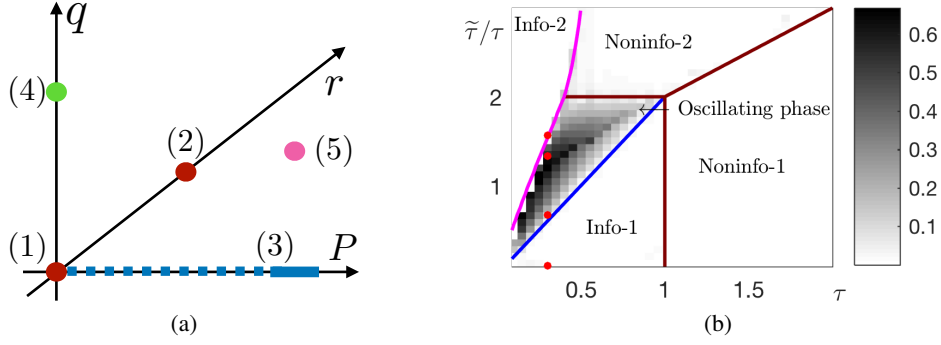


Figure 1: (a): The locations of the five types of fixed points of the ODE (13). Their properties are listed in Table 1. (b): The phase diagram for the stationary state of the ODE (13). The colored lines illustrate the theoretical prediction of the boundaries between the different phases. Simulations 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 2), whereas the variation is close to zero in all other phases.

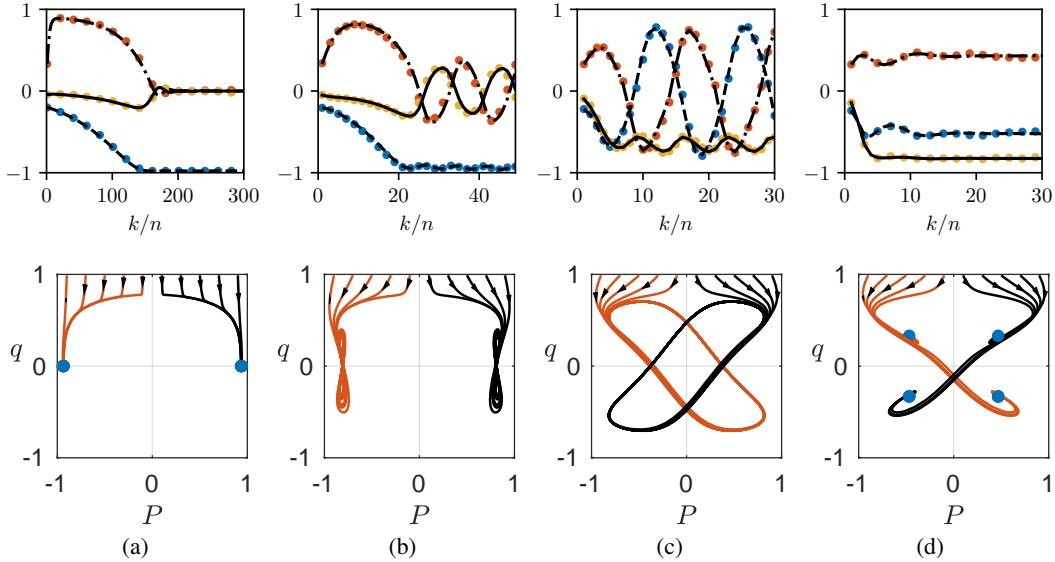


Figure 2: Macroscopic dynamics of Example 1 with $d = 1$. In the first row, the red, blue and yellow dots represent P_t , q_t , and r_t respectively of the experimental results of a single trial. The black curves under the dots are theoretical predictions given by the ODE (13). We set a fix the discriminator's 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 1. 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.

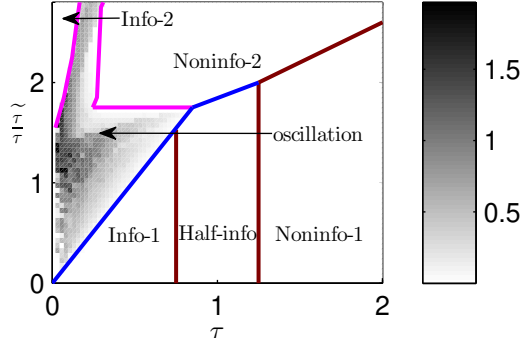


Figure 3: The phase diagram for the stationary states of the ODE (13) when $d = 2$. This phase diagram is generated by numerically computing the fixed points and eigenvalues of the Jacobian of the ODE (13).

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.

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 the i 'th row of the matrices $\sqrt{n}\mathbf{U}$ and $\sqrt{n}\mathbf{V}_k$ 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

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

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

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

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

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

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

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

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

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

where

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

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

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{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 :

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

Similarly, we have

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

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

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

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

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

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

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

$$h_k = \eta_{\Gamma} \left\langle f'(\mathbf{q}_k^\top \mathbf{c} + \sqrt{z_k \eta_{\Gamma}} e) \right\rangle_{\mathbf{c}, e} - \eta_G \left\langle \widetilde{f}'(\mathbf{r}_k^\top \widetilde{\mathbf{c}} + \sqrt{z_k \eta_G} e) \right\rangle_{\widetilde{\mathbf{c}}, e} - \lambda H'(z_k). \tag{S-15}$$

The second moment is

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

where

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

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 φ 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 φ 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 ℓ' th canonical basis (*i.e.*, all elements in $\mathbf{s}_{\ell'}$ are zeros, except that ℓ' 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 ℓ th row of \mathbf{P}_t and ℓ' 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

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

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 $\varphi = [\hat{\mathbf{v}}]_\ell [\hat{\mathbf{v}}]_{\ell'}$. If $\ell \neq \ell'$, we have $\nabla_{\hat{\mathbf{v}}} \varphi = [\hat{\mathbf{v}}]_\ell \mathbf{s}_{\ell'} + [\hat{\mathbf{v}}]_{\ell'} \mathbf{s}_\ell$, where $\mathbf{s}_{\ell'}$ is the ℓ' th canonical basis. Then

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

and

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

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

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

and

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

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

The rest two ODEs can be obtained in the similar way by letting φ to be each distinct component of $\hat{\mathbf{v}} \hat{w}$ and \hat{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 [1, 2]. We here build the whole proof on Lemma 2 in the supplementary materials of [3]. For reader's convenient, we present that lemma below.

Lemma 1 (Lemma 2 in the supplementary materials of [3]). *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 ϵ_1 ;*

(C.2) *$\mathbb{E} \|\boldsymbol{\delta}_k^{(n)}\| \leq C(T)/n^{1+\epsilon_2}$ for some positive ϵ_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 ϵ_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 denote \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 $\langle cf(\mathbf{c}^\top \mathbf{q} + e\sqrt{z}) \rangle_{c,e}$ is a Lipschitz function, because

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

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 $\mathbf{y}_k = \mathbf{U}\mathbf{c}_k + \sqrt{\eta_\Gamma}\mathbf{a}_k$ 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

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

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

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

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

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

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

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

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

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

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

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

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*

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

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

In the proof of this lemma and Lemma 2, we omit the two constants η_T and η_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

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

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

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

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

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

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

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

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

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

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

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

Using the above inequality iteratively, we have

$$\mathbb{E} \left(w_{k,i}^4 + \sum_{\ell=1}^d [V_k]_{i,\ell}^4 \right) \leq \left(n^{-2} + w_{0,i}^4 + \sum_{\ell=1}^d [V_0]_{i,\ell}^4 \right) e^{\frac{k}{n} C}.$$

Since $\mathbb{E} (w_{0,i}^4 + \sum_{\ell=1}^d [V_0]_{i,\ell}^4)$ 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

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

Correspondingly, ODE in (8) becomes:

$$\begin{aligned} \frac{d}{dt} P_t &= \tilde{\tau} (q_t \tilde{r}_t^{\top} \tilde{\Lambda} + P_t L_t) \\ \frac{d}{dt} q_t &= \tau (\Lambda q_t - P_t \tilde{\Lambda} r_t + q_t h_t) \\ \frac{d}{dt} r_t &= \tau (P_t^{\top} \Lambda q_t - S_t \tilde{\Lambda} r_t + r_t h_t) + \tilde{\tau} (\tilde{\Lambda} r_t + L_t r_t) \\ \frac{d}{dt} S_t &= \tilde{\tau} (r_t r_t^{\top} \tilde{\Lambda}^{\top} + \tilde{\Lambda} r_t r_t^{\top} + S_t L_t + L_t S_t) \\ \frac{d}{dt} z_t &= 2\tau (q_t^{\top} \Lambda q_t - r_t^{\top} \tilde{\Lambda} r_t + z_t h_t) \\ &\quad + \tau^2 [\eta_{\Gamma} (q_t^{\top} \Lambda q_t + z_t \eta_{\Gamma}) + \eta_{\text{G}} (r_t^{\top} \tilde{\Lambda} r_t + z_t \eta_{\text{G}})] \end{aligned} \quad (\text{S-48})$$

The first four equations are exactly (13). From last two equations of (S-48), by setting $\frac{d}{dt} \text{diag}\{S_t\} = \mathbf{0}$, $\frac{d}{dt} z_t = 0$, $\text{diag}(S_t) = \mathbf{I}$ 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 $\eta_{\Gamma} = \eta_{\text{G}} = 1$, 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 $g_1 = \tilde{\tau}\Lambda r(q - rP)$, $g_2 = \tau \left[\Lambda - \tau - \Lambda \left(1 + \frac{\tau}{2}\right) q^2 \right] q - \tau\Lambda \left[P + \left(\frac{\tau}{2} - 1\right) r q \right] r$ and $g_5 = \tau\Lambda P q + \left[\Lambda(\tilde{\tau} - \tau) - \frac{(\eta_I + \eta_G)\tau^2}{2} \right] r + \Lambda \left(\tau - \tilde{\tau} - \frac{\tau^2 \eta_G}{2} \right) r^3 - \tau\Lambda \left(1 + \frac{\tau \eta_I}{2} \right) r q^2$.

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

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

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

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 2×2 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 $\tilde{q} = q = 0$,

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

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 2×2 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 τ 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

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

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

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 2×2 matrix

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

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

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

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

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

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

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:

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

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.

References

- [1] P. Billingsley, *Convergence of probability measures*. John Wiley & Sons, 2013.
- [2] H. Kushner and G. G. Yin, *Stochastic approximation and recursive algorithms and applications*. Springer Science & Business Media, 2003, vol. 35.
- [3] 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.