1 We thank the reviewers for their useful comments, which help us improve the clarity and quality of our paper. We
2 address these comments point by point. The paper will also be revised accordingly, with additional details given there.

3 **Reviewer 1** *suggested us to provide a short "proof sketch".* We will present a brief summary of our proof strategy. It
4 generally follows the classical recipe for establishing the weak convergence of interactive particle systems. Our proof
5 sketch will also provide a roadmap to the SM, which contains all the technical details of the proof.

6 **Reviewer 1** *asked why we choose* $\log\cosh(x)$ *as the regularization function* $H(x)$ *in example 1.* In fact, any convex
7 function with its minimum reached at zero would be fine, for example $H(x) = |x|$ or $x^2$. The function $H(x) =$
8 $\log\cos(x)$ is just a convenient special case since its derivative $H'(x) = \tanh(x)$ is smooth and bounded.

9 **Reviewer 1** *asked what "a sufficiently strong regularizer" means in example 1.* It means to use a sufficiently large $\lambda$. In
10 our experiments, $\lambda > 1$ is sufficient. In such cases, the regularized variables, e.g. $\boldsymbol{w}^\top \boldsymbol{w} - 1$ and each diagonal term of
11 $\boldsymbol{V}^\top \boldsymbol{V} - \boldsymbol{I}$ in (4) are restricted to around 0 during the training process.

12 **Reviewer 1***'s additional comments:*
13 • We will emphasize that $P_{\tilde{c}}$ is fixed in our model. Meanwhile, a learnable $P_{\tilde{c}}$ is indeed an interesting idea to explore.
14 • By "analyzing the long-time behavior', we meant the investigation of the local stability of the ODE in Section 4. We
15 will clarify this point.
16 • The matrix $\boldsymbol{M}(t)$ is $(2d + 1) \times (2d + 1)$, where $d$ is the number of features in the true data model (1).
17 • Thank you for catching the various typos and other miscellaneous issues. We will fix these and post the code.

18 **Reviewer 2** *suggested avoiding using a single matrix M for the macroscopic states.* We agree with the reviewer that
19 $\boldsymbol{M}$ is simply a concatenation of different macroscopic states $\boldsymbol{P}$, $\boldsymbol{q}$, $\boldsymbol{S}$, $\boldsymbol{r}$ and $z$, all of which have concrete meanings.
20 However, we believe that there is indeed some benefit in introducing $\boldsymbol{M}$, which serves to streamline later presentations
21 (e.g. those in Sec. 3.1). In the revised paper, we will start our discussions by first defining $\boldsymbol{P}$, $\boldsymbol{q}$, $\boldsymbol{S}$, $\boldsymbol{r}$ and $z$ as the actual
22 macroscopic states we study. We will then introduce $\boldsymbol{M}$ and emphasize that it is just a convenient and compact notation.

23 **Reviewer 2** *remarked that the notions of the macroscopic and microscopic states are standard in physics.* It is indeed a
24 standard assumption in statistical physics that the macroscopic states of large systems tend to converge to deterministic
25 values due to self-averaging. However, we note that the mean-field regime in our work was not considered in previous
26 theoretical analysis of GAN. For example, a series of recent work [11-16] considers a different scaling regime where
27 the learning rate goes to zero but the system dimension $n$ stays fixed. In that regime, the microscopic dynamics are
28 deterministic even with the presence of the microscopic noise. In contrast, we study the regime where the learning
29 rate is fixed but the dimension $n \to \infty$. This setting allows us to quantify the effect of training noise in the learning
30 dynamics. Thus, we believe our work provides new and valuable insights into the theoretical understanding of GAN.

31 **Reviewer 2** *suggested adding more discussions on the learning dynamics.* In the revised paper, we will use the one
32 extra page available to us to provide more discussions on the hierarchical interactions of the macroscopic variables.
33 Specifically, we will consider the special setting where the ratio of the generator's and the discriminator's learning rates
34 $\tilde{\tau}/\tau \to 0$. This allows us to further simplify the ODE in (8) and to show that the learning process will be a combination
35 of two dynamics with two different time scales. The discriminator's process is associated with the faster time scale,
36 whereas the generator's process relates to the slower one.

37 **Reviewer 3** *suggested us to extend the discussion of the phase diagram to* $d \geq 2$ *cases.* After the initial submission
38 we have done some additional work on understanding the phase diagram for cases when $d \geq 2$. Accordingly, we will
39 expand the current discussions on line 45–50 in the SM. Specifically, we will present local stability analysis about
40 the phases of Info-1 and Noninfo-1 for $d = 2$. The analytical characterizations of the other phases are much more
41 challenging when $d \geq 2$. Instead of showing analytical results, we will characterize these phases numerically by directly
42 computing the eigenvalues of the Jacobian of the limiting ODEs.

43 **Reviewer 3** *asked how one can orthonormalize the columns of a general* $\boldsymbol{U}$. Suppose $\boldsymbol{U}$ in the data model (1) is not
44 orthogonal. We can always rewrite the product $\boldsymbol{U}\boldsymbol{c}$ as $(\boldsymbol{U}\boldsymbol{R})(\boldsymbol{R}^{-1}\boldsymbol{c})$, where $\boldsymbol{c}$ is the original random variable in the
45 feature space, and $\boldsymbol{R} \in \mathbb{R}^{d \times d}$ is a matrix that orthogonalizes and normalizes the columns of $\boldsymbol{U}$. We can then study an
46 equivalent system where the new feature vector is $\boldsymbol{R}^{-1}\boldsymbol{c}$.

47 **Reviewer 3** *asked how one can generalize our analysis to handle non-linear generators.* One possible extension is to
48 add a non-linear function, *e.g.* ReLU, to the linear data model (1) and the generator (2). Our analysis technique can be
49 extended to handle this case, where one can still obtain a differential equation in the scaling limit. Another possible
50 extension is a mixture of Gaussian model, where we need to treat $P_{\tilde{c}}$ as a learnable distribution rather than a fixed one.
51 (This was also mentioned by Reviewer 1.) We will comment on these possible extensions in the revised paper.

52 **Reviewer 3** *pointed out several typos and presentation issues.* Thank you. We will correct them in the revised version.
53 By "heuristic derivation", we meant that certain steps presented in that section were not fully rigorous as we directly
54 discard higher-order terms without any justification, in order to highlight the main ideas. In Section S-IV in the SM, we
55 rigorously justify these steps by providing bounds on those terms.