

# High-Dimensional Limit of Stochastic Gradient Flow via Dynamical Mean-Field Theory

author names withheld

Under Review for the Workshop on High-dimensional Learning Dynamics, 2025

## Abstract

We study the high-dimensional dynamics of multi-pass stochastic gradient descent (SGD) with small batch sizes. Existing asymptotic descriptions of SGD primarily cover online learning, large-batch regimes, or linear models. We address this gap by analyzing a continuous-time stochastic differential equation, called *stochastic gradient flow* (SGF), which approximates small-batch multi-pass SGD. In the proportional limit  $n, d \rightarrow \infty$  with  $n/d \rightarrow \delta$ , we derive a closed dynamical mean-field theory (DMFT) system and prove that it characterizes the asymptotic empirical distribution of the SGF parameters. The framework applies to nonlinear models including generalized linear models and two-layer neural networks, and recovers several existing high-dimensional descriptions, including noiseless gradient flow, online SGD, and high-dimensional linear regression.

## 1. Introduction

Stochastic gradient descent (SGD) [47] is the standard training algorithm for modern machine learning models. Its stochasticity can qualitatively change the training dynamics relative to noiseless gradient descent or gradient flow, affecting both optimization and the learned solution [33, 34]. Understanding these effects in high-dimensional regimes is therefore a central problem in the theory of learning dynamics.

A common approach is to derive low-dimensional equations that characterize the macroscopic behavior of high-dimensional training dynamics. For noiseless gradient descent and gradient flow, random matrix theory and dynamical mean-field theory (DMFT) have led to precise asymptotic descriptions [1, 10, 11, 14, 21, 39]. For SGD, however, existing frameworks remain restricted. Online or one-pass SGD has been extensively studied [7, 28, 49], but does not capture repeated reuse of a finite dataset. DMFT approaches for multi-pass SGD typically require batch sizes proportional to the dataset size [26, 39], while analyses of multi-pass small-batch SGD are largely limited to least-squares linear models through homogenized SGD and random matrix methods [42, 43]. Thus, a high-dimensional theory for multi-pass small-batch SGD in nonlinear models is still missing. A detailed review of related work is deferred to Appendix A.2.

In this work, we address this gap by studying a stochastic differential equation called *stochastic gradient flow* (SGF), which approximates multi-pass SGD with small batch sizes. We derive a DMFT system that characterizes the high-dimensional limit of SGF and clarifies its relationship to several existing descriptions of SGD dynamics. Table 1 compares our framework with prior approaches.

Our contributions are as follows.

**DMFT for stochastic gradient flow.** We derive a closed system of low-dimensional stochastic processes that characterizes the empirical distribution of SGF parameters in the proportional

asymptotic regime  $n, d \rightarrow \infty$  with  $n/d \rightarrow \delta$ . We prove existence and uniqueness of the DMFT solution and show convergence of the empirical parameter and prediction distributions to this solution. The framework applies to a broad class of nonlinear models, including generalized linear models and two-layer neural networks.

**Unification of existing analyses.** We show that the DMFT equation for the SGF reduces to known descriptions in several limits: noiseless gradient flow when  $\tau = 0$ , online SGD dynamics in the infinite data limit  $\delta \rightarrow \infty$ , and the Volterra equations for high-dimensional linear regression [42, 43].

Concurrent with our work, Fan and Wang [22] derived a related DMFT characterization for high-dimensional SGD and SGF. Our work was developed independently and gives a complementary formulation; see the discussion in Appendix A.2.

Table 1: Comparison of frameworks for high-dimensional analysis of SGD.

	Online SGD [7]	Previous DMFT [39]	HSGD [43]	Ours
Multi-pass	✗	✓	✓	✓
Small batch sizes	✓	✗	✓	✓
Nonlinear models	✓	✓	✗	✓

## 2. Setup

Notations are summarized in Appendix A.1.

We study the high-dimensional dynamics of a *stochastic gradient flow* (SGF), a continuous-time approximation of small-batch multi-pass SGD. Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  be the data matrix,  $\mathbf{z} \in \mathbb{R}^n$  an independent noise vector, and  $\boldsymbol{\theta}^t \in \mathbb{R}^{d \times m}$  the parameter. For functions  $h_t : \mathbb{R}^m \rightarrow \mathbb{R}^m$  and  $\ell_t : \mathbb{R}^m \times \mathbb{R} \rightarrow \mathbb{R}^m$  applied row-wise, define

$$d\boldsymbol{\theta}^t = -\left(h_t(\boldsymbol{\theta}^t) + \frac{1}{\delta} \mathbf{X}^\top \ell_t(\mathbf{r}^t; \mathbf{z})\right) dt + \sqrt{\frac{\tau}{\delta}} \sum_{i=1}^n \mathbf{x}_i \ell_t(r_i^t; z_i)^\top dB_i^t, \quad \mathbf{r}^t = \mathbf{X} \boldsymbol{\theta}^t, \quad (1)$$

where  $\mathbf{B}^t = (B_i^t)_{i \in [n]}$  is a Brownian motion in  $\mathbb{R}^n$ ,  $\delta > 0$  is the sample-to-dimension ratio  $\delta := n/d$ , and  $\tau > 0$  is the *temperature* parameter controlling the noise magnitude. The SGF (1) can be seen as a continuous-time approximation of mini-batch SGD with learning rate  $\eta$  and batch size  $B$  in the regime  $n, d \rightarrow \infty$ ,  $B = o(n)$ , where  $\tau = \eta/B$ . The time variable  $t$  corresponds to the number of SGD steps  $k$  via  $t = \eta k/d$ , so that  $t = O(1)$  time corresponds to  $k = O(d/\eta)$  steps. The noiseless case  $\tau = 0$  reduces to the gradient flow dynamics studied in Celentano et al. [14].

The SGF (1) is obtained by matching the leading drift and covariance of one step of mini-batch SGD. The precise SGD recursion and moment calculation are deferred to Appendix A.3.

We analyze SGF as a tractable continuous model of small-batch SGD; the use of SDE approximations for SGD is standard in the literature and discussed in Appendix A.2.

We note that the case of planted models where the data label is generated by a planted signal  $\boldsymbol{\theta}^*$  can be handled by augmenting the parameter  $\boldsymbol{\theta}$  with  $\boldsymbol{\theta}^*$ ; see Appendix E.1.

**Assumptions.** We assume independent centered sub-Gaussian covariates  $\mathbf{X}$  with variance  $1/d$ , proportional asymptotics  $n/d \rightarrow \delta \in (0, \infty)$ , and empirical convergence of the initialization and noise. We further assume Lipschitz regularity of  $h_t, \ell_t$  and their relevant derivatives, with polynomial growth of  $\ell_t$ . The precise assumptions are stated in Appendix A.4 as Assumptions A.1 and A.2.

### 3. Main Result

We now state the DMFT characterization of the SGF (1). The DMFT equation describes the high-dimensional dynamics through low-dimensional effective processes  $\theta^t, r^t \in \mathbb{R}^m$  together with correlation and response kernels  $C_\theta, C_\ell, R_\theta, R_\ell : \mathbb{R}_{\geq 0}^2 \rightarrow \mathbb{R}^{m \times m}$  and  $\Gamma : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^{m \times m}$ . Informally, the system is

$$\begin{aligned} \frac{d}{dt} \theta^t &= u^t - (h_t(\theta^t) + \Gamma(t)\theta^t) - \int_0^t R_\ell(t, s) \theta^s ds, \quad u \sim \text{GP}(0, C_\ell/\delta), \\ r^t &= w^t - \frac{1}{\delta} \int_0^t R_\theta(t, s) \ell_s(r^s; z) (ds + \sqrt{\tau\delta} dB^s), \quad w \sim \text{GP}(0, C_\theta), \\ C_\theta(t, t') &= \mathbb{E}[\theta^t \theta^{t'\top}], \quad R_\theta(t, t') = \mathbb{E} \left[ \frac{\partial \theta^t}{\partial w^{t'}} \right], \\ C_\ell(t, t') &= \mathbb{E}[\ell_t(r^t; z) (1 + \sqrt{\tau\delta} \dot{B}^t) \ell_{t'}(r^{t'}; z)^\top (1 + \sqrt{\tau\delta} \dot{B}^{t'})], \\ R_\ell(t, t') &= \mathbb{E} \left[ \frac{\partial \ell_t(r^t; z)}{\partial w^{t'}} \right], \quad \Gamma(t) = \mathbb{E}[\nabla_r \ell_t(r^t; z)], \end{aligned} \tag{2}$$

with  $R_\theta(t, t') = R_\ell(t, t') = 0$  for  $t < t'$ . The Brownian motion  $B^t \in \mathbb{R}$  is one-dimensional and captures the stochasticity inherited from SGF. The equation is self-consistent: the law of  $(\theta^t, r^t)$  determines the kernels, which in turn determine the law of  $(\theta^t, r^t)$ . The display above is informal because it uses functional derivatives and the white noise process  $\dot{B}^t$ ; the rigorous formulation, replacing  $C_\ell$  by a well-defined covariance object  $\Sigma_\ell$ , is given in Appendix B.

When  $\tau = 0$ , the Brownian term vanishes and (2) reduces to the DMFT equation for noiseless gradient flow in Celentano et al. [14]. Thus, the present equation extends gradient-flow DMFT to stochastic dynamics.

**Theorem 3.1 (Existence and uniqueness of the DMFT equation)** *Suppose Assumptions A.1 and A.2 hold. Then there exists  $T_* > 0$  such that, for every  $T \in [0, T_*]$ , the DMFT system admits a unique bounded fixed point  $(C_\theta, \Sigma_\ell, R_\theta, R_\ell, \Gamma)$  on  $[0, T]$ . Moreover, the effective processes  $\theta^t$  and  $r^t$  have continuous sample paths.*

*If either  $\tau = 0$  or  $\nabla_r^2 \ell_t(r; z) = 0$  for all  $t, r, z$ , then the result holds globally for all  $T > 0$ .*

The global cases include noiseless gradient flow and SGF under quadratic losses, the latter covering least-squares linear regression. For general nonlinear  $\ell_t$  with  $\tau > 0$ , the theorem gives local well-posedness; a quantitative lower bound on  $T_*$  and the proof are deferred to Appendix C.

**Theorem 3.2 (DMFT characterization of SGF)** *Suppose Assumptions A.1 and A.2 hold, and let  $T_*$  be as in Theorem 3.1. For any  $T \in [0, T_*]$ ,  $L \in \mathbb{N}$ , and  $0 \leq t_1 < \dots < t_L \leq T$ ,*

$$\text{p-lim}_{n, d \rightarrow \infty} W_2 \left( \hat{\mathbb{P}}(\theta^{t_1}, \dots, \theta^{t_L}), \mathbb{P}(\theta^{t_1}, \dots, \theta^{t_L}) \right) = 0, \tag{3}$$

$$\text{p-lim}_{n,d \rightarrow \infty} W_2\left(\hat{\mathbb{P}}(\mathbf{r}^{t_1}, \dots, \mathbf{r}^{t_L}, \mathbf{z}), \mathbb{P}(r^{t_1}, \dots, r^{t_L}, z)\right) = 0. \quad (4)$$

Theorem 3.2 states that the empirical distribution of SGF parameters and predictions converges to the law of the effective DMFT process at any finite collection of times. Equivalently, empirical averages of test functions of the coordinates of  $\boldsymbol{\theta}^t$  and  $\mathbf{r}^t$  converge to expectations under the DMFT law. This extends the gradient-flow result of Celentano et al. [14] to stochastic gradient flow.

The proof follows the AMP-based strategy of Celentano et al. [14]: discretize the SGF, map the resulting iteration to an AMP recursion, identify its state evolution with the discretized DMFT system, and pass to the continuous-time limit. The stochastic setting requires an additional truncation argument for Brownian increments and Stein’s lemma to derive the stochastic correction terms. Details are given in Appendix D.

#### 4. Applications and Special Cases

We highlight two consequences of the DMFT characterization: its online SGD limit and its reduction to known Volterra equations in linear regression.

**Infinite-data limit and online SGD.** In the infinite-data limit  $\delta \rightarrow \infty$ , the response terms vanish and the DMFT equation simplifies to a low-dimensional SDE given in Equation (316) in Appendix E.2. This limit corresponds to online SGD: when  $n$  is much larger than the number of updates proportional to  $d$ , a constant-size mini-batch effectively consists of fresh samples. In Appendix E.2, we show that (316) recovers the online SGD characterization for linear regression derived in Wang et al. [54].

**Linear regression.** As a tractable special case, consider high-dimensional linear regression

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2, \quad \mathbf{y} = \mathbf{X}\boldsymbol{\theta}^* + \mathbf{z}, \quad (5)$$

with zero initialization. This corresponds to the planted model with  $m = 1$ ,  $h_t = 0$ , and  $\ell_t(r, r^*; z) = r - r^* - z$ . Since  $\ell_t$  is linear in  $r$ , Theorem 3.1 gives global well-posedness of the DMFT equation.

Let the training and test errors of the parameter  $\boldsymbol{\theta}$  be given by

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \boldsymbol{\theta} - y_i)^2, \quad \mathcal{R}(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x}, y)} [(\mathbf{x}^\top \boldsymbol{\theta} - y)^2]. \quad (6)$$

Solving the DMFT equation yields a closed Volterra system for the limiting train and test errors:

$$\mathcal{L}(t) = \mathcal{L}_0(t) + \tau \int_0^t H_2(t-s) \mathcal{L}(s) ds, \quad \mathcal{R}(t) = \mathcal{R}_0(t) + \tau \int_0^t H_1(t-s) \mathcal{L}(s) ds, \quad (7)$$

where  $H_i(t) = \int x^i e^{-2xt} d\mu_{\text{MP}}(x)$ , and  $\mu_{\text{MP}}$  is the Marchenko–Pastur law with aspect ratio  $\delta$ . Here  $\mathcal{L}_0$  and  $\mathcal{R}_0$  are the noiseless limiting train and test errors; their explicit forms are given in Appendix E.3.

More precisely, for any finite collection of times  $t_1, \dots, t_L$ , we have the following convergence in probability as  $n, d \rightarrow \infty$ :

$$\text{p-lim}_{n,d \rightarrow \infty} \max_{l=1, \dots, L} |\mathcal{L}(\boldsymbol{\theta}^{t_l}) - \mathcal{L}(t_l)| = 0, \quad \text{p-lim}_{n,d \rightarrow \infty} \max_{l=1, \dots, L} |\mathcal{R}(\boldsymbol{\theta}^{t_l}) - \mathcal{R}(t_l)| = 0, \quad (8)$$

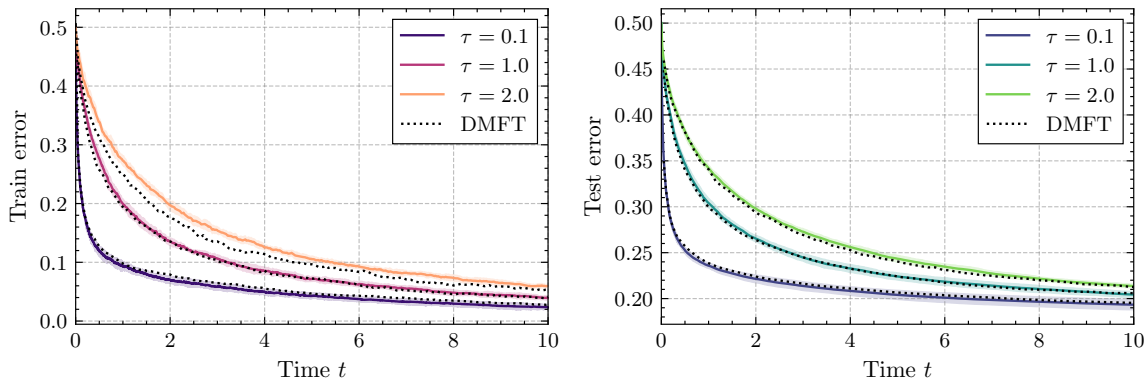


Figure 1: Train (left) and test (right) error dynamics for logistic regression at different temperatures  $\tau = \eta/B$ . Solid curves show the average over 10 SGD trials with  $d = 1024$  and  $n = 2048$ ; shaded regions show one standard deviation. Dotted curves show the DMFT predictions.

The scalar equation for  $\mathcal{L}$  can be solved independently, after which  $\mathcal{R}$  is obtained from the second equation in (7). These equations coincide, up to time rescaling, with the homogenized SGD Volterra equations of Paquette et al. [42, 43]. Thus, our DMFT framework recovers the known exact linear regression theory while extending naturally to nonlinear models.

## 5. Numerical Simulations and Discussion

We illustrate the DMFT prediction on a numerical simulation with a nonlinear model: logistic regression trained by multi-pass small-batch SGD. The data are sampled as  $\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_d/d)$  with  $d = 1024$ ,  $n = 2048$  ( $\delta = 2$ ), and labels  $y_i = \text{sign}(\boldsymbol{\theta}^{*\top} \mathbf{x}_i + z_i)$ ,  $\boldsymbol{\theta}^* \sim \mathcal{N}(0, \mathbf{I}_d)$ ,  $z_i \sim \mathcal{N}(0, 0.01)$ . We train with the regularized logistic objective

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{x}_i^\top \boldsymbol{\theta})) + \frac{\lambda}{2d} \|\boldsymbol{\theta}\|_2^2, \quad \lambda = 0.01. \quad (9)$$

SGD is run with batch size  $B = 10$  and varying learning rates, corresponding to different temperatures  $\tau = \eta/B$ . We compare the average train and test zero-one errors over 10 independent trials with the numerical solution of the DMFT equation. Details of the time discretization used to solve the DMFT equation are deferred to Appendix F.

Figure 1 shows that the DMFT prediction closely tracks the SGD train and test error dynamics across temperatures. Although the logistic model used here is not formally covered by our assumptions because of the non-differentiability induced by  $y = \text{sign}(r^* + z)$ , the agreement suggests that the DMFT equations remain predictive beyond the strictly covered setting.

**Discussion.** In summary, we derived a DMFT equation characterizing the high-dimensional limit of SGF, a continuous-time approximation of multi-pass small-batch SGD. The resulting theory covers nonlinear models, recovers known limits such as noiseless gradient flow and linear-regression Volterra equations, and gives accurate predictions in the logistic-regression experiment above. Future directions include proving sharper SGF–SGD approximation results, analyzing the long-time behavior of the DMFT equations, and applying the framework to more detailed nonlinear models.

## References

- [1] Madhu S. Advani, Andrew M. Saxe, and Haim Sompolinsky. High-dimensional dynamics of generalization error in neural networks. *Neural Networks*, 132:428–446, 2020. doi: 10.1016/j.neunet.2020.08.022.
- [2] Elisabeth Agoritsas, Giulio Biroli, Pierfrancesco Urbani, and Francesco Zamponi. Out-of-equilibrium dynamical mean-field equations for the perceptron model. *Journal of Physics A: Mathematical and Theoretical*, 51(8):085002, 2018. doi: 10.1088/1751-8121/aaa68d.
- [3] Alnur Ali, Edgar Dobriban, and Ryan Tibshirani. The implicit regularization of stochastic gradient flow for least squares. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pages 233–244, 2020.
- [4] Luca Arnaboldi, Ludovic Stephan, Florent Krzakala, and Bruno Loureiro. From high-dimensional & mean-field dynamics to dimensionless ODEs: a unifying approach to SGD in two-layers networks. In *Proceedings of Thirty Sixth Conference on Learning Theory*, volume 195, pages 1199–1227, 2023.
- [5] Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. Online stochastic gradient descent on non-convex losses from high-dimensional inference. *Journal of Machine Learning Research*, 22(106):1–51, 2021.
- [6] Mohsen Bayati and Andrea Montanari. The dynamics of message passing on dense graphs, with applications to compressed sensing. *IEEE Transactions on Information Theory*, 57(2): 764–785, 2011. doi: 10.1109/TIT.2010.2094817.
- [7] Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. High-dimensional limit theorems for SGD: effective dynamics and critical scaling. In *Advances in Neural Information Processing Systems*, volume 35, pages 25349–25362, 2022.
- [8] Marc A. Berger and Victor J. Mizel. Volterra equations with Itô integrals—I. *Journal of Integral Equations*, 2(3):187–245, 1980.
- [9] M Biehl and H Schwarze. Learning by on-line gradient descent. *Journal of Physics A: Mathematical and General*, 28(3):643–656, 1995. doi: 10.1088/0305-4470/28/3/018.
- [10] Antoine Bodin and Nicolas Macris. Model, sample, and epoch-wise descents: exact solution of gradient flow in the random feature model. In *Advances in Neural Information Processing Systems*, volume 34, pages 21605–21617, 2021.
- [11] Blake Bordelon, Alexander Atanasov, and Cengiz Pehlevan. A dynamical model of neural scaling laws. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235, pages 4345–4382, 2024.
- [12] Blake Bordelon, Alexander Atanasov, and Cengiz Pehlevan. How feature learning can improve neural scaling laws. *Journal of Statistical Mechanics: Theory and Experiment*, 2025(8):084002, 2025. doi: 10.1088/1742-5468/adeb1.

- [13] Michael Celentano, Andrea Montanari, and Yuchen Wu. The estimation error of general first order methods. In *Proceedings of Thirty Third Conference on Learning Theory*, volume 125, pages 1078–1141, 2020.
- [14] Michael Celentano, Chen Cheng, and Andrea Montanari. The high-dimensional asymptotics of first order methods with random data. *arXiv preprint arXiv:2112.07572*, 2021.
- [15] Pratik Chaudhari and Stefano Soatto. Stochastic gradient descent performs variational inference, converges to limit cycles for deep networks. In *2018 Information Theory and Applications Workshop (ITA)*, pages 1–10, 2018. doi: 10.1109/ITA.2018.8503224.
- [16] Xiang Cheng, Dong Yin, Peter Bartlett, and Michael Jordan. Stochastic gradient and Langevin processes. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pages 1810–1819, 2020.
- [17] Elizabeth Collins-Woodfin, Courtney Paquette, Elliot Paquette, and Inbar Seroussi. Hitting the high-dimensional notes: an ODE for SGD learning dynamics on GLMs and multi-index models. *Information and Inference: A Journal of the IMA*, 13(4):iaae028, 2024. doi: 10.1093/imaia/iaae028.
- [18] Andrea Crisanti, Heinz Horner, and H-J Sommers. The spherical p-spin interaction spin-glass model: the dynamics. *Zeitschrift für Physik B Condensed Matter*, 92(2):257–271, 1993.
- [19] L. F. Cugliandolo and J. Kurchan. Analytical solution of the off-equilibrium dynamics of a long-range spin-glass model. *Physical Review Letters*, 71(1):173–176, 1993. doi: 10.1103/PhysRevLett.71.173.
- [20] Leticia F. Cugliandolo. Recent applications of dynamical mean-field methods. *Annual Review of Condensed Matter Physics*, 15(1):177–213, 2024. doi: 10.1146/annurev-conmatphys-040721-022848.
- [21] Yatin Dandi, Emanuele Troiani, Luca Arnaboldi, Luca Pesce, Lenka Zdeborova, and Florent Krzakala. The benefits of reusing batches for gradient descent in two-layer networks: breaking the curse of information and leap exponents. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235, pages 9991–10016, 2024.
- [22] Zhou Fan and Leda Wang. High-dimensional learning dynamics of multi-pass stochastic gradient descent in multi-index models. *arXiv preprint arXiv:2601.21093*, 2026.
- [23] Zhou Fan, Justin Ko, Bruno Loureiro, Yue M. Lu, and Yandi Shen. Dynamical mean-field analysis of adaptive Langevin diffusions: propagation-of-chaos and convergence of the linear response. *arXiv preprint arXiv:2504.15556*, 2025.
- [24] Zhou Fan, Justin Ko, Bruno Loureiro, Yue M. Lu, and Yandi Shen. Dynamical mean-field analysis of adaptive Langevin diffusions: replica-symmetric fixed point and empirical Bayes. *arXiv preprint arXiv:2504.15558*, 2025.
- [25] Oliver Y. Feng, Ramji Venkataramanan, Cynthia Rush, and Richard J. Samworth. A unifying tutorial on approximate message passing. *Foundations and Trends in Machine Learning*, 15(4): 335–536, 2022. doi: 10.1561/22000000092.

- [26] Cédric Gerbelot, Emanuele Troiani, Francesca Mignacco, Florent Krzakala, and Lenka Zdeborová. Rigorous dynamical mean-field theory for stochastic gradient descent methods. *SIAM Journal on Mathematics of Data Science*, 6(2):400–427, 2024. doi: 10.1137/23M1594388.
- [27] Margalit Glasgow. SGD finds then tunes features in two-layer neural networks with near-optimal sample complexity: a case study in the XOR problem. In *International Conference on Learning Representations*, 2024.
- [28] Sebastian Goldt, Madhu Advani, Andrew M Saxe, Florent Krzakala, and Lenka Zdeborová. Dynamics of stochastic gradient descent for two-layer neural networks in the teacher-student setup. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [29] Qiyang Han. Entrywise dynamics and universality of general first order methods. *The Annals of Statistics*, 53(4), 2025. doi: 10.1214/25-AOS2544.
- [30] Qiyang Han and Masaaki Imaizumi. Precise gradient descent training dynamics for finite-width multi-layer neural networks. *arXiv preprint arXiv:2505.04898*, 2025.
- [31] Yuma Ichikawa and Koji Hukushima. Learning dynamics in linear VAE: posterior collapse threshold, superfluous latent space pitfalls, and speedup with KL annealing. In *Proceedings of the 27th International Conference on Artificial Intelligence and Statistics*, volume 238, pages 1936–1944, 2024.
- [32] Yuma Ichikawa, Shuhei Kashiwamura, and Ayaka Sakata. High-dimensional learning dynamics of quantized models with straight-through estimator. *arXiv preprint arXiv:2510.10693*, 2025.
- [33] Stanisław Jastrzębski, Zachary Kenton, Devansh Arpit, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amos Storkey. Three factors influencing minima in SGD. *arXiv preprint arXiv:1711.04623*, 2017.
- [34] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: generalization gap and sharp minima. In *International Conference on Learning Representations*, 2017.
- [35] Peter E. Kloeden and Eckhard Platen. *Numerical Solution of Stochastic Differential Equations*. Springer Berlin Heidelberg, 1992. doi: 10.1007/978-3-662-12616-5.
- [36] Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and adaptive stochastic gradient algorithms. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 2101–2110, 2017.
- [37] Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and dynamics of stochastic gradient algorithms I: mathematical foundations. *Journal of Machine Learning Research*, 20(40):1–47, 2019.
- [38] Stephan Mandt, Matthew D. Hoffman, and David M. Blei. Stochastic gradient descent as approximate Bayesian inference. *Journal of Machine Learning Research*, 18(134):1–35, 2017.
- [39] Francesca Mignacco, Florent Krzakala, Pierfrancesco Urbani, and Lenka Zdeborová. Dynamical mean-field theory for stochastic gradient descent in Gaussian mixture classification. In *Advances in Neural Information Processing Systems*, volume 33, pages 9540–9550, 2020.

- [40] Andrea Montanari and Pierfrancesco Urbani. Dynamical decoupling of generalization and overfitting in large two-layer networks. In *Advances in Neural Information Processing Systems*, 2025.
- [41] Sota Nishiyama and Masaaki Imaizumi. Precise dynamics of diagonal linear networks: a unifying analysis by dynamical mean-field theory. In *The 29th International Conference on Artificial Intelligence and Statistics*, 2026.
- [42] Courtney Paquette, Kiwon Lee, Fabian Pedregosa, and Elliot Paquette. SGD in the large: average-case analysis, asymptotics, and stepsize criticality. In *Proceedings of Thirty Fourth Conference on Learning Theory*, volume 134, pages 3548–3626, 2021.
- [43] Courtney Paquette, Elliot Paquette, Ben Adlam, and Jeffrey Pennington. Homogenization of SGD in high-dimensions: exact dynamics and generalization properties. *Mathematical Programming*, 214(1-2):1–90, 2025. doi: 10.1007/s10107-024-02171-3.
- [44] Scott Pesme, Loucas Pillaud-Vivien, and Nicolas Flammarion. Implicit bias of SGD for diagonal linear networks: a provable benefit of stochasticity. In *Advances in Neural Information Processing Systems*, volume 34, pages 29218–29230, 2021.
- [45] Yunwei Ren, Eshaan Nichani, Denny Wu, and Jason D. Lee. Emergence and scaling laws in SGD learning of shallow neural networks. In *Advances in Neural Information Processing Systems*, 2025.
- [46] P Riegler and M Biehl. On-line backpropagation in two-layered neural networks. *Journal of Physics A: Mathematical and General*, 28(20):L507–L513, 1995. doi: 10.1088/0305-4470/28/20/002.
- [47] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, 22(3):400–407, 1951.
- [48] David Saad and Sara A. Solla. Exact solution for on-line learning in multilayer neural networks. *Physical Review Letters*, 74(21):4337–4340, 1995. doi: 10.1103/PhysRevLett.74.4337.
- [49] David Saad and Sara A. Solla. On-line learning in soft committee machines. *Physical Review E*, 52(4):4225–4243, 1995. doi: 10.1103/PhysRevE.52.4225.
- [50] Stefano Sarao Mannelli, Giulio Biroli, Chiara Cammarota, Florent Krzakala, Pierfrancesco Urbani, and Lenka Zdeborová. Marvels and pitfalls of the Langevin algorithm in noisy high-dimensional inference. *Physical Review X*, 10(1):011057, 2020. doi: 10.1103/PhysRevX.10.011057.
- [51] H. Sompolinsky and Annette Zippelius. Dynamic theory of the spin-glass phase. *Physical Review Letters*, 47(5):359–362, 1981. doi: 10.1103/PhysRevLett.47.359.
- [52] H. Sompolinsky and Annette Zippelius. Relaxational dynamics of the Edwards-Anderson model and the mean-field theory of spin-glasses. *Physical Review B*, 25(11):6860–6875, 1982. doi: 10.1103/PhysRevB.25.6860.

- [53] Rodrigo Veiga, Ludovic Stephan, Bruno Loureiro, Florent Krzakala, and Lenka Zdeborová. Phase diagram of stochastic gradient descent in high-dimensional two-layer neural networks. In *Advances in Neural Information Processing Systems*, volume 35, pages 23244–23255, 2022.
- [54] Chuang Wang, Jonathan Mattingly, and Yue 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.
- [55] Chuang Wang, Hong Hu, and Yue Lu. A solvable high-dimensional model of GAN. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [56] Tianhao Wang, Xinyi Zhong, and Zhou Fan. Universality of approximate message passing algorithms and tensor networks. *The Annals of Applied Probability*, 34(4), 2024. doi: 10.1214/24-AAP2056.

**Contents**

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Setup</b>	<b>2</b>
<b>3</b>	<b>Main Result</b>	<b>3</b>
<b>4</b>	<b>Applications and Special Cases</b>	<b>4</b>
<b>5</b>	<b>Numerical Simulations and Discussion</b>	<b>5</b>
<b>A</b>	<b>Preliminaries</b>	<b>12</b>
A.1	Notation . . . . .	12
A.2	Related Works . . . . .	12
A.3	Derivation of the SGF as an Approximation of SGD . . . . .	14
A.4	Formal Statement of the Assumptions . . . . .	15
<b>B</b>	<b>Definition of the DMFT Equation</b>	<b>15</b>
B.1	Rigorous Definition of the DMFT System . . . . .	15
B.2	Discretized DMFT System . . . . .	17
B.3	Correspondence to the Informal Definition . . . . .	18
<b>C</b>	<b>Proof of Theorem 3.1</b>	<b>20</b>
C.1	Admissible Spaces $\mathcal{S}_\theta(T)$ and $\mathcal{S}_\ell(T)$ . . . . .	20
C.2	Equipping Metrics on $\mathcal{S}_\theta(T)$ and $\mathcal{S}_\ell(T)$ . . . . .	23
C.3	Proof of Lemma C.4 . . . . .	24
C.3.1	Construction of the Admissible Spaces . . . . .	24
C.3.2	$\mathcal{T}_{\theta \rightarrow \ell}$ maps $\mathcal{S}_\theta$ into $\mathcal{S}_\ell$ . . . . .	25
C.3.3	$\mathcal{T}_{\ell \rightarrow \theta}$ maps $\mathcal{S}_\ell$ into $\mathcal{S}_\theta$ . . . . .	29
C.3.4	A Rough Estimate of $T_*$ . . . . .	30
C.4	Proof of Lemma C.6 . . . . .	31
C.5	Proof of Lemma C.7 . . . . .	35
<b>D</b>	<b>Proof of Theorem 3.2</b>	<b>37</b>
D.1	Proof of Lemma D.1 . . . . .	39
D.2	Proof of Lemma D.2 . . . . .	41
D.2.1	Reduction to AMP . . . . .	41
D.2.2	Mapping the state evolution to DMFT . . . . .	44
D.3	Proof of Lemma D.3 . . . . .	46
D.3.1	Proof of Lemma D.6 . . . . .	48
D.3.2	Proof of Lemma D.7 . . . . .	53
<b>E</b>	<b>Details of Applications and Special Cases</b>	<b>53</b>
E.1	Planted Models . . . . .	53
E.1.1	Setup for Planted Models . . . . .	53
E.1.2	The DMFT Equation . . . . .	54

E.1.3	Proof of Corollary E.1 . . . . .	55
E.2	Infinite Data Limit . . . . .	57
E.2.1	Derivation of the Reduced DMFT Equation . . . . .	57
E.2.2	Example: Linear Regression . . . . .	58
E.3	Linear Regression . . . . .	58
E.3.1	Simplifying the DMFT Equations . . . . .	60
E.3.2	Solving the DMFT Equations . . . . .	64
<b>F Details of Numerical Simulations</b>		<b>70</b>
<b>G Discretization of SDEs in High Dimensions</b>		<b>71</b>

## Appendix A. Preliminaries

### A.1. Notation

$I_d \in \mathbb{R}^{d \times d}$  denotes the  $d \times d$  identity matrix.  $\mathbf{1}_d \in \mathbb{R}^d$  denotes the all-ones vector  $\mathbf{1}_d = (1, \dots, 1)^\top$ .  $\text{GP}(0, Q)$  denotes a centered Gaussian process with covariance kernel  $Q$ .  $\mathbb{I}(\cdot)$  denotes the indicator function that returns 1 if the argument is true and 0 otherwise.  $\|\cdot\|_2$  denotes the  $\ell_2$  norm for vectors and the 2-operator norm for matrices.  $\|\cdot\|_F$  denotes the Frobenius norm for matrices.  $W_p$  denotes the  $p$ -Wasserstein distance between probability measures.  $\text{p-lim}$  denotes convergence in probability.  $\text{P}(X)$  denotes the law of a random variable  $X$ .  $\text{P}(X, Y)$  denotes the joint law of random variables  $X$  and  $Y$ . For a matrix  $\mathbf{x} \in \mathbb{R}^{d \times m}$ , we denote by  $\hat{\text{P}}(\mathbf{x})$  the empirical distribution of its rows, i.e.,  $\hat{\text{P}}(\mathbf{x}) = \frac{1}{d} \sum_{i=1}^d \delta_{x_i}$ . Similarly,  $\hat{\text{P}}(\mathbf{x}, \mathbf{y})$  denotes the empirical joint distribution of the rows of  $\mathbf{x}$  and  $\mathbf{y}$ , i.e.,  $\hat{\text{P}}(\mathbf{x}, \mathbf{y}) = \frac{1}{d} \sum_{i=1}^d \delta_{x_i, y_i}$ .  $\|\cdot\|_{\psi_2}$  denotes the sub-Gaussian norm ( $\psi_2$ -Orlicz norm).

### A.2. Related Works

**One-pass SGD.** The study of online SGD using low-dimensional ODEs was pioneered in the statistical physics literature on learning in perceptrons and two-layer neural networks [9, 46, 48, 49]. They derived a closed system of low-dimensional ODEs for macroscopic parameters such as the correlation between the student and teacher weights and analyzed their dynamics, which closely approximates the original online SGD dynamics in high dimensions. The ODEs typically involve correction terms that account for the stochasticity of the dynamics, and analyzing these ODEs provides insights into how the noise affects the training dynamics.

Recently, these works have been put on a rigorous footing by Goldt et al. [28] using techniques developed in Wang et al. [54, 55]. Ben Arous et al. [7] extended these techniques to more general models and general scaling of learning rates. Analysis of online SGD in high dimensions using similar techniques has been applied to a wide range of models due to its versatility and simplicity [4, 5, 17, 27, 31, 32, 45, 53].

**Multi-pass SGD.** DMFT has recently gained attention as a powerful framework for analyzing high-dimensional random dynamics, including multi-pass GD and SGD, by reducing them to low-dimensional effective processes. DMFT was originally developed in spin glass theory [18, 19, 51, 52] and has been applied to analyzing various high-dimensional optimization dynamics [2, 20, 50].

In the context of SGD dynamics, Mignacco et al. [39] derived DMFT equations for multi-pass gradient flow and SGD in shallow neural networks, heuristically using statistical physics techniques.

To avoid the problem of vanishing stochasticity in the continuous-time limit  $\eta \rightarrow 0$ , they considered a variant of SGD called *persistent SGD* to retain nontrivial noise in the continuous-time limit. Their analysis depends on the batch size growing proportionally to the number of samples. In contrast, we work with a stochastic gradient flow which approximates mini-batch SGD with small (sublinear) batch sizes compared to the number of samples, which is a common setting in practice.

There are several rigorous works that derived DMFT equations for GD/SGD. Celentano et al. [14] rigorously derived DMFT equations for gradient flow dynamics in shallow neural networks by using time discretization and mapping to approximate message passing [6, 25]. We build upon their proof technique to analyze SGD dynamics in this work. Gerbelot et al. [26] derived DMFT equations for discrete-time GD and SGD for shallow neural networks with batch sizes proportional to the number of samples and a constant number of updates. More recently, Fan et al. [23, 24] derived DMFT equations for Langevin dynamics of Bayesian linear regression. A closely related line of work is the study of *general first order methods* (GFOMs), which provides a framework for analyzing a broad class of iterative algorithms, including GD, using a low-dimensional recursion similar to DMFT [13, 29, 30].

DMFT equations have been used for analyzing long-time behavior of optimization dynamics and provided insights into deep learning phenomena such as scaling laws and timescale separation [11, 12, 21, 40, 41].

For linear models, there is a framework that analyzes multi-pass SGD with small batch sizes and proportionally many updates using continuous-time equations. Paquette et al. [42] derived a low-dimensional and continuous-time Volterra equation characterizing the training loss dynamics of SGD in high-dimensional linear regression models. Paquette et al. [43] extended this work and introduced an SDE called *homogenized SGD* (HSGD) as a high-dimensional equivalent of SGD dynamics in linear regression. The HSGD framework allows deriving equations for macroscopic quantities of SGD dynamics, such as training and test errors. Our work provides a similar framework for broader settings, including generalized linear models and shallow neural networks.

**SDE approximations of SGD.** The SDE approximation of SGD has been used extensively in the literature as a continuous model of discrete-time SGD [3, 15, 33, 38, 44]. Tools from the Itô stochastic calculus can be used to analyze the dynamics of the SDE in detail, and this approach has led to fruitful insights into the behavior of SGD.

Although we do not prove that the high-dimensional behaviors of the SGD (10) and the SGF (1) match (and indeed they do not match exactly, as shown in the concurrent work [22]), we take the SGF (1) as an approximate continuous model of the SGD (10) and analyze its high-dimensional limiting behavior. It is shown that in the fixed-dimensional setting, when the learning rate is small ( $\eta \rightarrow 0$ ), the dynamics of SGD is well approximated by that of the SGF (1) [16, 36, 37]. Hence, we expect that the SGF (1) approximates the SGD (10) well when  $\tau$  is sufficiently small, which we empirically confirm by numerical experiments in Section 5.

**Comparison with the concurrent work [22].** A concurrent work by Fan and Wang [22] also derives a DMFT equation characterizing the high-dimensional dynamics of SGD and SGF (which they call SME). Their DMFT equation for SGF is equivalent to ours, although defined differently; they define the response function  $R_\ell$  as a linear operator, while we define it as a continuous function given by the expectation of a stochastic process. Another important difference is that they consider bounded  $\ell_t$ , while we allow for unbounded, Lipschitz continuous  $\ell_t$ . This allows us to rigorously apply our theory to linear regression settings which involve unbounded loss gradients.

### A.3. Derivation of the SGF as an Approximation of SGD

**Stochastic gradient descent.** Fix  $m \in \mathbb{N}$ , a learning rate  $\eta > 0$ , a batch size  $B \in \mathbb{N}$ , a data matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , a noise vector  $\mathbf{z} \in \mathbb{R}^n$ , the initial parameter  $\boldsymbol{\theta}^0 \in \mathbb{R}^{d \times m}$ , and functions  $h: \mathbb{R}^m \times \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^m; (\boldsymbol{\theta}, t) \mapsto h_t(\boldsymbol{\theta})$  and  $\ell: \mathbb{R}^m \times \mathbb{R} \times \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^m; (r, z, t) \mapsto \ell_t(r; z)$ . We consider the following stochastic process for  $\hat{\boldsymbol{\theta}}^k \in \mathbb{R}^{d \times m}$  for  $k = 0, 1, \dots$ , initialized with  $\hat{\boldsymbol{\theta}}^0 = \boldsymbol{\theta}^0$ :

$$\hat{\boldsymbol{\theta}}^{k+1} = \hat{\boldsymbol{\theta}}^k - \eta \cdot \left( \frac{1}{d} h_{t_k}(\hat{\boldsymbol{\theta}}^k) + \frac{1}{B} \sum_{i \in \mathcal{B}^k} \mathbf{x}_i \ell_{t_k}(\hat{r}_i^k; z_i)^\top \right), \quad \hat{r}^k = \mathbf{X} \hat{\boldsymbol{\theta}}^k \in \mathbb{R}^{n \times m}, \quad (10)$$

where  $t_k := \eta k / d$  and  $\mathbf{x}_i \in \mathbb{R}^d$  and  $\hat{r}_i^k \in \mathbb{R}^m$  are the  $i$ -th row of  $\mathbf{X}$  and  $\hat{r}^k$ , respectively. Here,  $h_t$  is applied row-wise to  $\hat{\boldsymbol{\theta}}^k$ . In each update  $k$ , the mini-batch  $\mathcal{B}^k$  is sampled uniformly at random from all subsets of  $[n] := \{1, 2, \dots, n\}$  with size  $B$ .

This stochastic process includes the mini-batch stochastic gradient descent on the following training objective with the learning rate  $\eta$  and the batch size  $B$ :

$$\mathcal{L}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n L(\boldsymbol{\theta}^\top \mathbf{x}_i; z_i) + \frac{1}{d} \sum_{i=1}^d H(\theta_i), \quad (11)$$

where  $\theta_i \in \mathbb{R}^m$  is the  $i$ -th row of  $\boldsymbol{\theta} \in \mathbb{R}^{d \times m}$ ,  $L: \mathbb{R}^m \times \mathbb{R} \rightarrow \mathbb{R}$  is a loss function, and  $H: \mathbb{R}^m \rightarrow \mathbb{R}$  is a regularization function. In this case, we have  $h_t(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} H(\boldsymbol{\theta})$  and  $\ell_t(r; z) = \nabla_r L(r; z)$ . This setting includes generalized linear models ( $m = 1$ ) and two-layer neural networks with width  $m$ . In the following, we refer to the general process (10) as an ‘SGD’ for simplicity, although it may not correspond to the gradient descent on any objective function.

**Derivation of SGF.** To see that the SGF (1) approximates the SGD (10), consider the one-step increment of the SGD (10) given by  $\Delta \hat{\boldsymbol{\theta}}^k = \hat{\boldsymbol{\theta}}^{k+1} - \hat{\boldsymbol{\theta}}^k$  and compute its first and second moments conditioned on  $\hat{\boldsymbol{\theta}}^k$ :

$$\mathbb{E}[\Delta \hat{\boldsymbol{\theta}}^k \mid \hat{\boldsymbol{\theta}}^k] = -\frac{\eta}{d} \left( h_{t_k}(\hat{\boldsymbol{\theta}}^k) + \frac{1}{\delta} \mathbf{X}^\top \ell_{t_k}(\hat{r}^k; \mathbf{z}) \right), \quad (12)$$

$$\begin{aligned} \text{Cov}(\Delta \hat{\boldsymbol{\theta}}^k \mid \hat{\boldsymbol{\theta}}^k) &= \frac{\eta^2(n-B)}{Bn(n-1)} \sum_{i=1}^n \mathbf{x}_i \ell_{t_k}(\hat{r}_i^k; z_i)^\top \otimes \mathbf{x}_i \ell_{t_k}(\hat{r}_i^k; z_i) \\ &\quad + \frac{\eta^2(n-B)}{Bn^2(n-1)} \mathbf{X}^\top \ell_{t_k}(\mathbf{X} \hat{\boldsymbol{\theta}}^k; \mathbf{z}) \otimes \mathbf{X}^\top \ell_{t_k}(\mathbf{X} \hat{\boldsymbol{\theta}}^k; \mathbf{z}). \end{aligned} \quad (13)$$

Here, the outer product  $\mathbf{A} \otimes \mathbf{B}$  for matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times m}$  is interpreted as  $\text{vec}(\mathbf{A}) \text{vec}(\mathbf{B})^\top \in \mathbb{R}^{dm \times dm}$ , where  $\text{vec}$  is the vectorization operator. In the proportional high-dimensional limit where  $n, d \rightarrow \infty$  with  $n/d \rightarrow \delta \in (0, \infty)$  and  $B = o(n)$ , the leading term of the covariance simplifies to

$$\text{Cov}(\Delta \hat{\boldsymbol{\theta}}^k \mid \hat{\boldsymbol{\theta}}^k) = \frac{\eta}{d} \cdot \frac{\tau}{\delta} \sum_{i=1}^n \mathbf{x}_i \ell_{t_k}(\hat{r}_i^k; z_i)^\top \otimes \mathbf{x}_i \ell_{t_k}(\hat{r}_i^k; z_i) + (\text{sub-leading terms}). \quad (14)$$

Thus, the first and second moments of the increment  $\Delta \hat{\boldsymbol{\theta}}^k$  match those of the increment of the SGF (1) up to the leading order in  $d$  over the time interval of length  $\Delta t = \eta/d$ .

#### A.4. Formal Statement of the Assumptions

##### Assumption A.1 (Data distribution)

- The entries  $\mathbf{X} = (x_{ij})_{i \in [n], j \in [d]}$  are independent, satisfying  $\mathbb{E} x_{ij} = 0$ ,  $\mathbb{E} x_{ij}^2 = 1/d$ , and  $\|x_{ij}\|_{\psi_2} \leq C/\sqrt{d}$  for some constant  $C > 0$ .
- Proportional high-dimensional asymptotics:  $n, d \rightarrow \infty$ ,  $n/d \rightarrow \delta \in (0, \infty)$ .
- $\mathbf{z} \in \mathbb{R}^n$  and  $\boldsymbol{\theta}^0 \in \mathbb{R}^{d \times m}$  are independent of  $\mathbf{X}$ , and for all  $p \geq 1$ , their empirical distributions  $\hat{\mathbb{P}}(\boldsymbol{\theta}^0)$  and  $\hat{\mathbb{P}}(\mathbf{z})$  converge in  $p$ -Wasserstein distance to  $\mathbb{P}(\boldsymbol{\theta}^0)$  and  $\mathbb{P}(\mathbf{z})$  respectively, almost surely as  $d \rightarrow \infty$ .

The distribution of the data  $\mathbf{X}$  is not restricted to the Gaussian distribution, and hence our analysis is *universal* with respect to the data distribution. The last condition on  $\mathbf{z}$  and  $\boldsymbol{\theta}^0$  is satisfied, for example, when their entries are i.i.d. samples from distributions with bounded moments of all orders.

The following assumption is used to guarantee that the SGF solution does not grow too fast and that the DMFT equation to be introduced is well-defined.

**Assumption A.2 (Function regularity)** *There exists a constant  $M > 0$  such that the following holds.*

- $h_t(\theta)$  and its Jacobian  $Dh = (\nabla_\theta h, \partial_t h)$  are Lipschitz continuous in  $t$  and  $\theta$ , i.e., for  $t_1, t_2 \geq 0$  and  $\theta_1, \theta_2 \in \mathbb{R}^m$ ,

$$\|h_{t_1}(\theta_1) - h_{t_2}(\theta_2)\|_2 + \|Dh_{t_1}(\theta_1) - Dh_{t_2}(\theta_2)\|_2 \leq M(\|\theta_1 - \theta_2\|_2 + |t_1 - t_2|). \quad (15)$$

- $\ell_t(r; z)$ , its Jacobian  $D\ell = (\nabla_r \ell, \partial_t \ell)$ , and its Hessian  $D^2 \ell$  are Lipschitz continuous in  $t$  and  $r$  for any  $z \in \mathbb{R}$ , i.e., for  $t_1, t_2 \geq 0$  and  $r_1, r_2 \in \mathbb{R}^m$ ,

$$\begin{aligned} \|\ell_{t_1}(r_1; z) - \ell_{t_2}(r_2; z)\|_2 + \|D\ell_{t_1}(r_1; z) - D\ell_{t_2}(r_2; z)\|_2 + \|D^2 \ell_{t_1}(r_1; z) - D^2 \ell_{t_2}(r_2; z)\|_2 \\ \leq M(\|r_1 - r_2\|_2 + |t_1 - t_2|). \end{aligned} \quad (16)$$

- $\ell_t(r; z)$  has polynomial growth for any  $t \geq 0$ , i.e., there exists some  $p \geq 1$  such that  $\|\ell_t(r; z)\|_2 \leq M(1 + \|r\|_2 + |z|)^p$ .

## Appendix B. Definition of the DMFT Equation

In this section, we provide a rigorous definition of the DMFT system  $\mathfrak{S}$  introduced informally in Equation (2). The key idea is to define all objects through well-defined auxiliary stochastic processes, avoiding functional derivatives and the formal derivative of the Brownian motion. We also introduce a discretized DMFT equation for which the connection between the informal definition and the rigorous definition is more transparent. The discretized DMFT equation is also used as an intermediate step in the proof of Theorem 3.2 and as a numerical method for solving the DMFT equation (see Appendix F).

### B.1. Rigorous Definition of the DMFT System

We rigorously define the DMFT system  $\mathfrak{S}$  for functions  $C_\theta, \Sigma_\ell, R_\theta, R_\ell: \mathbb{R}_{\geq 0}^2 \rightarrow \mathbb{R}^{m \times m}$  and  $\Gamma: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^{m \times m}$  self-consistently as follows. First, given  $\Sigma_\ell, R_\ell, \Gamma$ , define stochastic processes

$\{\theta^t \in \mathbb{R}^m\}_{t \geq 0}$  and  $\{\rho_\theta^{t,t'} \in \mathbb{R}^{m \times m}\}_{t \geq t' \geq 0}$  by the following equations.

$$\theta^t = \theta^0 + U^t - \int_0^t \left( h_s(\theta^s) + \Gamma(s)\theta^s + \int_0^s R_\ell(s, s')\theta^{s'} ds' \right) ds, \quad U \sim \text{GP}(0, \Sigma_\ell/\delta), \quad (17)$$

$$\rho_\theta^{t,t'} = I_m - \int_{t'}^t \left( (\nabla_\theta h_s(\theta^s) + \Gamma(s))\rho_\theta^{s,t'} + \int_{t'}^s R_\ell(s, s')\rho_\theta^{s',t'} ds' \right) ds. \quad (18)$$

Then, set  $C_\theta, R_\theta$  as

$$C_\theta(t, t') = \mathbb{E}[\theta^t \theta^{t'\top}], \quad (19)$$

$$R_\theta(t, t') = \mathbb{E}[\rho_\theta^{t,t'}] \quad (t \geq t'), \quad (20)$$

and  $R_\theta(t, t') = 0$  for  $t < t'$ . Here, the expectation is with respect to the randomness of  $\theta^0 \sim \text{P}(\theta^0)$  and the Gaussian process  $U$ .

Next, given  $C_\theta, R_\theta$ , define stochastic processes  $\{r^t \in \mathbb{R}^m\}_{t \geq 0}$  and  $\{\rho_\ell^{t,t'}, D_\ell^{t,t'} \in \mathbb{R}^{m \times m}\}_{t \geq t' \geq 0}$  by the following equations.

$$r^t = w^t - \frac{1}{\delta} \int_0^t R_\theta(t, s) \ell_s(r^s; z) (ds + \sqrt{\tau\delta} dB^s), \quad w \sim \text{GP}(0, C_\theta), \quad (21)$$

$$\rho_\ell^{t,t'} = \nabla_r \ell_t(r^t; z) \rho_r^{t,t'}, \quad (22)$$

$$D_\ell^{t,t'} = \nabla_r \ell_t(r^t; z) \left( -\frac{1}{\delta} \int_{t'}^t R_\theta(t, s) D_\ell^{s,t'} (ds + \sqrt{\tau\delta} dB^s) \right) + \nabla_r^2 \ell_t(r^t; z) [D_r^{t,t'}] \cdot \rho_r^{t,t'}, \quad (23)$$

where  $B^t$  is a Brownian motion in  $\mathbb{R}$ , and we defined the auxiliary processes  $\rho_r^{t,t'} \in \mathbb{R}^{m \times m}$  and  $D_r^{t,t'} \in \mathbb{R}^m$  as

$$\rho_r^{t,t'} = -\frac{1}{\delta} \int_{t'}^t R_\theta(t, s) \rho_\ell^{s,t'} (ds + \sqrt{\tau\delta} dB^s) - \frac{1}{\delta} R_\theta(t, t') \nabla_r \ell_{t'}(r^{t'}; z), \quad (24)$$

$$D_r^{t,t'} = -\frac{1}{\delta} \int_{t'}^t R_\theta(t, s) \nabla_r \ell_s(r^s; z) D_r^{s,t'} (ds + \sqrt{\tau\delta} dB^s) - \sqrt{\frac{\tau}{\delta}} R_\theta(t, t') \ell_{t'}(r^{t'}; z), \quad (25)$$

and the notation  $\nabla_r^2 \ell_t(r^t; z)[v]$  for  $v \in \mathbb{R}^m$  in Equation (23) denotes the product

$$\sum_{i=1}^m \partial_{r_i} \nabla_r \ell_t(r^t; z) v_i \in \mathbb{R}^{m \times m}. \quad (26)$$

Then, set  $\Sigma_\ell, R_\ell, \Gamma$  as

$$\Sigma_\ell(t, t') = \mathbb{E}[L^t L^{t'\top}], \quad L^t := \int_0^t \ell_s(r^s; z) (ds + \sqrt{\tau\delta} dB^s), \quad (27)$$

$$R_\ell(t, t') = \mathbb{E}[\rho_\ell^{t,t'}] + \sqrt{\tau\delta} \mathbb{E}[D_\ell^{t,t'}] \quad (t \geq t'), \quad (28)$$

$$\Gamma(t) = \mathbb{E}[\nabla_r \ell_t(r^t; z)], \quad (29)$$

and  $R_\ell(t, t') = 0$  for  $t < t'$ . Here, the expectation is with respect to the randomness of  $z \sim \text{P}(z)$ , the Gaussian process  $w$ , and the Brownian motion  $B$ .

Then, the solution of the DMFT system  $\mathfrak{S}$  is defined as a fixed point of the above two mappings.

Note that the above definition is written purely in terms of standard Itô integrals and is thus amenable to rigorous analysis.

## B.2. Discretized DMFT System

We present a discretized version of the DMFT system defined above, in which the time variable  $t$  is discretized with step size  $\gamma > 0$ .

Let  $t_i := i\gamma$  for  $i = 0, 1, 2, \dots$ . We define the discretized DMFT system  $\mathfrak{S}^\gamma$  by the following equations.

$$\theta_\gamma^{t_i} = \theta^0 + U_\gamma^{t_i} - \gamma \sum_{j=0}^{i-1} \left( h_{t_j}(\theta_\gamma^{t_j}) + \Gamma^\gamma(t_j) \theta_\gamma^{t_j} + \gamma \sum_{k=0}^{j-1} R_\ell^\gamma(t_j, t_k) \theta_\gamma^{t_k} \right), \quad (30a)$$

$$\rho_{\theta, \gamma}^{t_i, t_j} = I_m - \gamma \sum_{k=j+1}^{i-1} \left( (\nabla_\theta h_{t_k}(\theta_\gamma^{t_k}) + \Gamma^\gamma(t_k)) \rho_{\theta, \gamma}^{t_k, t_j} + \gamma \sum_{l=j+1}^{k-1} R_\ell^\gamma(t_k, t_l) \rho_{\theta, \gamma}^{t_l, t_j} \right), \quad (30b)$$

$$r_\gamma^{t_i} = w_\gamma^{t_i} - \frac{1}{\delta} \sum_{j=0}^{i-1} R_\theta^\gamma(t_i, t_j) \ell_{t_j}(r_\gamma^{t_j}; z) (\gamma + \sqrt{\tau} \delta (B^{t_{j+1}} - B^{t_j})), \quad (30c)$$

$$\rho_{\ell, \gamma}^{t_i, t_j} = \nabla_r \ell_{t_i}(r_\gamma^{t_i}; z) \rho_{r, \gamma}^{t_i, t_j}, \quad (30d)$$

$$\rho_{r, \gamma}^{t_i, t_j} = -\frac{1}{\delta} \sum_{k=j+1}^{i-1} R_\theta^\gamma(t_i, t_k) \rho_{\ell, \gamma}^{t_k, t_j} (\gamma + \sqrt{\tau} \delta (B^{t_{k+1}} - B^{t_k})) - \frac{1}{\delta} R_\theta^\gamma(t_i, t_j) \nabla_r \ell_{t_j}(r_\gamma^{t_j}; z), \quad (30e)$$

$$D_{\ell, \gamma}^{t_i, t_j} = \nabla_r \ell_{t_i}(r_\gamma^{t_i}; z) \left( -\frac{1}{\delta} \sum_{k=j+1}^{i-1} R_\theta^\gamma(t_i, t_k) D_{\ell, \gamma}^{t_k, t_j} (\gamma + \sqrt{\tau} \delta (B^{t_{k+1}} - B^{t_k})) \right) + \nabla_r^2 \ell_{t_i}(r_\gamma^{t_i}; z) [D_{r, \gamma}^{t_i, t_j}] \cdot \rho_{r, \gamma}^{t_i, t_j}, \quad (30f)$$

$$D_{r, \gamma}^{t_i, t_j} = -\frac{1}{\delta} \sum_{k=j+1}^{i-1} R_\theta^\gamma(t_i, t_k) \nabla_r \ell_{t_k}(r_\gamma^{t_k}; z) D_{r, \gamma}^{t_k, t_j} (\gamma + \sqrt{\tau} \delta (B^{t_{k+1}} - B^{t_k})) - \sqrt{\frac{\tau}{\delta}} R_\theta^\gamma(t_i, t_j) \ell_{t_j}(r_\gamma^{t_j}; z), \quad (30g)$$

where  $(U_\gamma^{t_i}, w_\gamma^{t_i})_{i \geq 0}$  satisfies

$$\mathbb{E}[U_\gamma^{t_i} U_\gamma^{t_j \top}] = \Sigma_\ell^\gamma(t_i, t_j) / \delta, \quad \mathbb{E}[w_\gamma^{t_i} w_\gamma^{t_j \top}] = C_\theta^\gamma(t_i, t_j). \quad (31)$$

Then, set  $C_\theta^\gamma, R_\theta^\gamma, \Sigma_\ell^\gamma, R_\ell^\gamma, \Gamma^\gamma$  as

$$\begin{aligned} C_\theta^\gamma(t_i, t_j) &= \mathbb{E}[\theta_\gamma^{t_i} \theta_\gamma^{t_j \top}], & R_\theta^\gamma(t_i, t_j) &= \mathbb{E}[\rho_{\theta, \gamma}^{t_i, t_j}], \\ \Sigma_\ell^\gamma(t_i, t_j) &= \mathbb{E}[L_\gamma^{t_i} L_\gamma^{t_j \top}], & L_\gamma^{t_i} &:= \sum_{k=0}^{i-1} \ell_{t_k}(r_\gamma^{t_k}; z) (\gamma + \sqrt{\tau} \delta (B^{t_{k+1}} - B^{t_k})), \\ R_\ell^\gamma(t_i, t_j) &= \mathbb{E}[\rho_{\ell, \gamma}^{t_i, t_j}] + \sqrt{\tau} \delta \mathbb{E}[D_{\ell, \gamma}^{t_i, t_j}], & \Gamma^\gamma(t_i) &= \mathbb{E}[\nabla_r \ell_{t_i}(r_\gamma^{t_i}; z)], \end{aligned} \quad (32)$$

where we set  $R_\theta^\gamma(t_i, t_j) = R_\ell^\gamma(t_i, t_j) = 0$  for  $i \leq j$ .

We will show that the solution of the discretized DMFT equation  $\mathfrak{S}^\gamma$  converges to the unique solution of the continuous-time DMFT equation  $\mathfrak{S}$  as  $\gamma \rightarrow 0$  in Lemma D.3.

### B.3. Correspondence to the Informal Definition

Once discretized, it is easy to see the correspondence between the rigorous definition of the DMFT system given above and the informal definition given in Equation (2). We distinguish the variables in the two definitions by writing bars over the variables in the informal definition, e.g.,  $\bar{\theta}^t$ ,  $\bar{r}^t$ , etc.

We first discretize the informal definition in Equation (2) with step size  $\gamma > 0$  in the same manner as in the previous section. We obtain

$$\begin{aligned}
 \frac{\bar{\theta}^{t_{i+1}} - \bar{\theta}^{t_i}}{\gamma} &= \bar{u}^{t_i} - (h_{t_i}(\bar{\theta}^{t_i}) + \bar{\Gamma}(t_i)\bar{\theta}^{t_i}) - \gamma \sum_{j=0}^{i-1} \bar{R}_\ell(t_i, t_j)\bar{\theta}^{t_j}, \quad \bar{u} \sim \text{GP}(0, \bar{C}_\ell/\delta), \\
 \bar{r}^{t_i} &= \bar{w}^{t_i} - \frac{1}{\delta} \sum_{j=1}^{i-1} \bar{R}_\theta(t_i, t_j)\ell_{t_j}(\bar{r}^{t_j}; z)(\gamma + \sqrt{\tau\delta}(\bar{B}^{t_{j+1}} - \bar{B}^{t_j})), \quad \bar{w} \sim \text{GP}(0, \bar{C}_\theta), \\
 \bar{C}_\theta(t_i, t_j) &= \mathbb{E}[\bar{\theta}^{t_i}\bar{\theta}^{t_j\top}], \quad \bar{R}_\theta(t_i, t_j) = \frac{1}{\gamma} \mathbb{E}\left[\frac{\partial \bar{\theta}^{t_i}}{\partial \bar{u}^{t_j}}\right] \quad (i > j), \\
 \bar{C}_\ell(t_i, t_j) &= \mathbb{E}\left[\ell_{t_i}(\bar{r}^{t_i}; z) \left(1 + \sqrt{\tau\delta} \frac{\bar{B}^{t_{i+1}} - \bar{B}^{t_i}}{\gamma}\right) \ell_{t_j}(\bar{r}^{t_j}; z)^\top \left(1 + \sqrt{\tau\delta} \frac{\bar{B}^{t_{j+1}} - \bar{B}^{t_j}}{\gamma}\right)\right], \\
 \bar{R}_\ell(t_i, t_j) &= \frac{1}{\gamma} \mathbb{E}\left[\frac{\partial \ell_{t_i}(\bar{r}^{t_i}; z)}{\partial \bar{w}^{t_j}}\right], \quad \bar{\Gamma}(t_i) = \mathbb{E}[\nabla_r \ell_{t_i}(\bar{r}^{t_i}; z)] \quad (i > j).
 \end{aligned} \tag{33}$$

Then, we transform the above equations to show their correspondence to the discretized DMFT system  $\mathfrak{S}^\gamma$ . The equations for  $\bar{r}^{t_i}$ ,  $\bar{C}_\theta$ , and  $\bar{\Gamma}$  directly correspond to definitions of  $r_\gamma^{t_i}$ ,  $C_\theta^\gamma$ , and  $\Gamma^\gamma$  in  $\mathfrak{S}^\gamma$ . Next, we show correspondence for  $\bar{\theta}^{t_i}$  and  $\bar{C}_\ell$ . Summing the equation for  $\bar{\theta}^{t_i}$  in (33) over  $i$  and multiplying by  $\gamma$ , we obtain

$$\bar{\theta}^{t_i} - \bar{\theta}^0 = \gamma \sum_{j=0}^{i-1} \bar{u}^{t_j} - \gamma \sum_{j=0}^{i-1} \left( h_{t_j}(\bar{\theta}^{t_j}) + \bar{\Gamma}(t_j)\bar{\theta}^{t_j} + \gamma \sum_{k=0}^{j-1} \bar{R}_\ell(t_j, t_k)\bar{\theta}^{t_k} \right), \quad \bar{u} \sim \text{GP}(0, \bar{C}_\ell/\delta). \tag{34}$$

Let  $\bar{U}^{t_i} := \gamma \sum_{j=0}^{i-1} \bar{u}^{t_j}$ .  $\bar{U}^{t_i}$  is a Gaussian process with covariance given by

$$\frac{1}{\delta} \bar{\Sigma}_\ell(t_i, t_j) := \mathbb{E}[\bar{U}^{t_i}\bar{U}^{t_j\top}] = \gamma^2 \sum_{k=0}^{i-1} \sum_{l=0}^{j-1} \mathbb{E}[\bar{u}^{t_k}\bar{u}^{t_l\top}] = \frac{\gamma^2}{\delta} \sum_{k=0}^{i-1} \sum_{l=0}^{j-1} \bar{C}_\ell(t_k, t_l) = \frac{1}{\delta} \mathbb{E}[\bar{L}^{t_i}\bar{L}^{t_j\top}], \tag{35}$$

where we set

$$\bar{L}^{t_i} := \sum_{j=0}^{i-1} \ell_{t_j}(\bar{r}^{t_j}; z)(\gamma + \sqrt{\tau\delta}(\bar{B}^{t_{j+1}} - \bar{B}^{t_j})). \tag{36}$$

Thus,  $\bar{\theta}^{t_i}$  and  $\bar{\Sigma}_\ell(t_i, t_j) = \gamma^2 \sum_{k=0}^{i-1} \sum_{l=0}^{j-1} \bar{C}_\ell(t_k, t_l)$  correspond to  $\theta_\gamma^{t_i}$  and  $\Sigma_\ell^\gamma$  in  $\mathfrak{S}^\gamma$ . Furthermore, differentiating  $\bar{\theta}^{t_i}$  with respect to  $\bar{u}^{t_j}$  ( $j < i$ ), we obtain

$$\frac{\partial \bar{\theta}^{t_i}}{\partial \bar{u}^{t_j}} = \gamma I_m - \sum_{k=j+1}^{i-1} \left( (\nabla_\theta h_{t_k}(\bar{\theta}^{t_k}) + \bar{\Gamma}(t_k)) \frac{\partial \bar{\theta}^{t_k}}{\partial \bar{u}^{t_j}} + \gamma \sum_{l=j+1}^{k-1} \bar{R}_\ell(t_k, t_l) \frac{\partial \bar{\theta}^{t_l}}{\partial \bar{u}^{t_j}} \right). \tag{37}$$

This shows correspondence for  $\gamma^{-1}\partial\bar{\theta}^{t_i}/\partial\bar{w}^{t_j}$  and  $\rho_{\theta,\gamma}^{t_i,t_j}$  and thus  $\bar{R}_\theta$  and  $R_\theta^\gamma$ . Finally, we show correspondence for  $\bar{R}_\ell$ . Differentiating  $\ell_{t_i}(\bar{r}^{t_i}; z)$  with respect to  $\bar{w}^{t_j}$  ( $j < i$ ), we obtain

$$\frac{\partial\ell_{t_i}(\bar{r}^{t_i}; z)}{\partial\bar{w}^{t_j}} = \nabla_r \ell_{t_i}(\bar{r}^{t_i}; z) \frac{\partial\bar{r}^{t_i}}{\partial\bar{w}^{t_j}}, \quad (38)$$

$$\begin{aligned} \frac{\partial\bar{r}^{t_i}}{\partial\bar{w}^{t_j}} &= -\frac{1}{\delta} \sum_{k=j+1}^{i-1} \bar{R}_\theta(t_i, t_k) \frac{\partial\ell_{t_k}(\bar{r}^{t_k}; z)}{\partial\bar{w}^{t_j}} (\gamma + \sqrt{\tau\delta}(\bar{B}^{t_{k+1}} - \bar{B}^{t_k})) \\ &\quad - \frac{1}{\delta} \bar{R}_\theta(t_i, t_j) \nabla_r \ell_{t_j}(\bar{r}^{t_j}; z) (\gamma + \sqrt{\tau\delta}(\bar{B}^{t_{j+1}} - \bar{B}^{t_j})). \end{aligned} \quad (39)$$

Let  $\bar{\rho}_\ell^{t_i,t_j}$  and  $\bar{\rho}_r^{t_i,t_j}$  be the solution of the following equation:

$$\bar{\rho}_\ell^{t_i,t_j} = \nabla_r \ell_{t_i}(\bar{r}^{t_i}; z) \bar{\rho}_r^{t_i,t_j}, \quad (40)$$

$$\bar{\rho}_r^{t_i,t_j} = -\frac{1}{\delta} \sum_{k=j+1}^{i-1} \bar{R}_\theta(t_i, t_k) \bar{\rho}_\ell^{t_k,t_j} (\gamma + \sqrt{\tau\delta}(\bar{B}^{t_{k+1}} - \bar{B}^{t_k})) - \frac{1}{\delta} \bar{R}_\theta(t_i, t_j) \nabla_r \ell_{t_j}(\bar{r}^{t_j}; z). \quad (41)$$

These equations correspond to definitions of  $\rho_{\ell,\gamma}^{t_i,t_j}$  and  $\rho_{r,\gamma}^{t_i,t_j}$ . By the linearity of the above equations, we see that

$$\frac{\partial\ell_{t_i}(\bar{r}^{t_i}; z)}{\partial\bar{w}^{t_j}} = \bar{\rho}_\ell^{t_i,t_j} (\gamma + \sqrt{\tau\delta}(\bar{B}^{t_{j+1}} - \bar{B}^{t_j})). \quad (42)$$

Let  $\bar{G}^j := (\bar{B}^{t_{j+1}} - \bar{B}^{t_j})/\sqrt{\gamma}$ . Then, we have  $\bar{G}^j \sim \mathcal{N}(0, 1)$  i.i.d. for  $j = 0, 1, 2, \dots$ . Taking the expectation of the above equation, we obtain

$$\bar{R}_\ell(t_i, t_j) = \frac{1}{\gamma} \mathbb{E} \left[ \frac{\partial\ell_{t_i}(\bar{r}^{t_i}; z)}{\partial\bar{w}^{t_j}} \right] = \mathbb{E}[\bar{\rho}_\ell^{t_i,t_j}] + \sqrt{\frac{\tau\delta}{\gamma}} \mathbb{E}[\bar{\rho}_\ell^{t_i,t_j} \bar{G}^j]. \quad (43)$$

By Stein's lemma (Gaussian integration by parts), we obtain

$$\mathbb{E}[\bar{\rho}_\ell^{t_i,t_j} \bar{G}^j] = \mathbb{E} \left[ \frac{\partial\bar{\rho}_\ell^{t_i,t_j}}{\partial\bar{G}^j} \right]. \quad (44)$$

Differentiating  $\bar{\rho}_\ell^{t_i,t_j}$  with respect to  $\bar{G}^j$  and using independence of  $\bar{r}^{t_k}$  and  $\bar{G}^j$  for  $k \leq j$ , we obtain

$$\frac{\partial\bar{\rho}_\ell^{t_i,t_j}}{\partial\bar{G}^j} = \nabla_r \ell_{t_i}(\bar{r}^{t_i}; z) \frac{\partial\bar{\rho}_r^{t_i,t_j}}{\partial\bar{G}^j} + \nabla_r^2 \ell_{t_i}(\bar{r}^{t_i}; z) \left[ \frac{\partial\bar{r}^{t_i}}{\partial\bar{G}^j} \right] \cdot \bar{\rho}_r^{t_i,t_j}, \quad (45)$$

$$\frac{\partial\bar{\rho}_r^{t_i,t_j}}{\partial\bar{G}^j} = -\frac{1}{\delta} \sum_{k=j+1}^{i-1} \bar{R}_\theta(t_i, t_k) \frac{\partial\bar{\rho}_\ell^{t_k,t_j}}{\partial\bar{G}^j} (\gamma + \sqrt{\tau\delta\gamma} \bar{G}^k), \quad (46)$$

$$\frac{\partial\bar{r}^{t_i}}{\partial\bar{G}^j} = -\frac{1}{\delta} \sum_{k=j+1}^{i-1} \bar{R}_\theta(t_i, t_k) \nabla_r \ell_{t_k}(\bar{r}^{t_k}; z) \frac{\partial\bar{r}^{t_k}}{\partial\bar{G}^j} (\gamma + \sqrt{\tau\delta\gamma} \bar{G}^k) - \sqrt{\frac{\tau\gamma}{\delta}} \bar{R}_\theta(t_i, t_j) \ell_{t_j}(\bar{r}^{t_j}; z). \quad (47)$$

This shows that  $\gamma^{-1/2} \partial \bar{\rho}_\ell^{t_i, t_j} / \partial \bar{G}^j$  and  $\gamma^{-1/2} \partial \bar{r}^{t_i} / \partial \bar{G}^j$  correspond to  $D_{\ell, \gamma}^{t_i, t_j}$  and  $D_{r, \gamma}^{t_i, t_j}$ . Thus, we have correspondence for  $\bar{R}_\ell$  and  $R_\ell^\gamma$ .

Summarizing, the informal definition corresponds to the formal definition by the following correspondence (informal definition on the left, rigorous definition on the right):

$$\begin{aligned}
 \bar{\theta}^{t_i} &\Leftrightarrow \theta_\gamma^{t_i}, & \frac{1}{\gamma} \frac{\partial \bar{\theta}^{t_i}}{\partial \bar{u}^{t_j}} &\Leftrightarrow \rho_{\theta, \gamma}^{t_i, t_j}, & \bar{C}_\theta(t_i, t_j) &\Leftrightarrow C_\theta^\gamma(t_i, t_j), & \bar{R}_\theta(t_i, t_j) &\Leftrightarrow R_\theta^\gamma(t_i, t_j), \\
 \bar{r}^{t_i} &\Leftrightarrow r_\gamma^{t_i}, & \frac{\partial \bar{r}^{t_i, t_j}}{\partial \bar{w}^{t_j}} &\Leftrightarrow \rho_{\ell, \gamma}^{t_i, t_j}, & \frac{1}{\sqrt{\gamma}} \frac{\partial \bar{\rho}_\ell^{t_i, t_j}}{\partial \bar{G}^{t_j}} &\Leftrightarrow D_{\ell, \gamma}^{t_i, t_j}, & \frac{1}{\gamma} \mathbb{E} \left[ \frac{\partial \ell_{t_i}(\bar{r}^{t_i}; z)}{\partial \bar{w}^{t_j}} \right] &\Leftrightarrow \mathbb{E}[\rho_{\ell, \gamma}^{t_i, t_j}] + \sqrt{\tau \delta} \mathbb{E}[D_{\ell, \gamma}^{t_i, t_j}] \\
 \gamma^2 \sum_{k=0}^{i-1} \sum_{l=0}^{j-1} \bar{C}_\ell(t_k, t_l) &\Leftrightarrow \Sigma_\ell^\gamma(t_i, t_j), & \bar{R}_\ell(t_i, t_j) &\Leftrightarrow R_\ell^\gamma(t_i, t_j), & \bar{\Gamma}(t_i) &\Leftrightarrow \Gamma^\gamma(t_i).
 \end{aligned} \tag{48}$$

Although these two definitions are equivalent in discrete time, the informal definition does not have a well-defined continuous-time limit as  $\gamma \rightarrow 0$ , while the rigorous definition does. Thus, for theoretical purposes, we work with the rigorous definition. In numerics, however, we work with the informal definition after time discretization, as it leads to simpler numerical schemes (see Appendix F for details).

### Appendix C. Proof of Theorem 3.1

We prove Theorem 3.1 using a contraction mapping argument similar to that of Celentano et al. [14, Theorem 1] and Fan et al. [23, Theorem 2.4]. It proceeds as follows.

1. For  $T > 0$ , we define *admissible spaces*  $\mathcal{S}_\theta(T)$  and  $\mathcal{S}_\ell(T)$  for the DMFT objects  $(C_\theta, R_\theta)$  and  $(\Sigma_\ell, R_\ell, \Gamma)$ , respectively. We define mappings  $\mathcal{T}_{\theta \rightarrow \ell}: \mathcal{S}_\theta(T) \rightarrow \mathcal{S}_\ell(T)$  and  $\mathcal{T}_{\ell \rightarrow \theta}: \mathcal{S}_\ell(T) \rightarrow \mathcal{S}_\theta(T)$  such that the fixed point of their composition  $\mathcal{T} := \mathcal{T}_{\ell \rightarrow \theta} \circ \mathcal{T}_{\theta \rightarrow \ell}$  solves the DMFT system. We show that for sufficiently small  $T > 0$ , these mappings are well-defined.
2. Next, we construct a metric on the function spaces  $\mathcal{S}_\theta(T)$  and  $\mathcal{S}_\ell(T)$  such that  $\mathcal{T}$  is a contraction.
3. Finally, we apply Banach's fixed point theorem to show the uniqueness and existence of the fixed point of  $\mathcal{T}$ .

#### C.1. Admissible Spaces $\mathcal{S}_\theta(T)$ and $\mathcal{S}_\ell(T)$

Since the following quantities are bounded by assumptions, we take  $M > 0$  sufficiently large so that we have

$$\max \left\{ \mathbb{E} \|\theta^0\|_2^2, \sup_{t \in [0, T]} h_t(0), \sup_{t \in [0, T]} \mathbb{E} \|\ell_t(0; z)\|_2^{2p} \right\} \leq M, \tag{49}$$

for  $p = 1, 2$ .

For  $T > 0$ , we define admissible spaces  $\mathcal{S}_\theta(T)$  and  $\mathcal{S}_\ell(T)$  as follows.

**Definition C.1 (Admissible space  $\mathcal{S}_\theta(T)$ )** Let  $D \subset (0, T)$  be a finite set. Let  $\mathcal{S}_\theta(T)$  be a set of function pairs  $(C_\theta, R_\theta)$  defined on  $[0, T]^2$ . We say that  $\mathcal{S}_\theta(T)$  is admissible if there exist constants  $\Phi_\theta, M_\theta > 0$  such that every  $(C_\theta, R_\theta) \in \mathcal{S}_\theta(T)$  satisfies the following.

- $C_\theta$  is a covariance kernel (in particular, it satisfies  $C_\theta(t, t') = C_\theta(t', t)^\top$ ) and satisfies  $\|C_\theta(t, t)\|_2 \leq \Phi_\theta$  for all  $t \in [0, T]$  and  $C_\theta(0, 0) = \mathbb{E}[\theta^0 \theta^{0\top}]$ . Furthermore,  $C_\theta(t, t')$  is uniformly continuous over  $t, t' \in I$  for each maximal interval  $I$  of  $[0, T] \setminus D$  and satisfies

$$\|C_\theta(t, t) - 2C_\theta(t, t') + C_\theta(t', t')\|_2 \leq M_\theta |t - t'|, \quad (50)$$

for any  $t, t' \in I$ .

- $R_\theta$  satisfies  $R_\theta(t, t') = 0$  for  $0 \leq t < t' \leq T$  and  $\|R_\theta(t, t')\|_2 \leq \Phi_\theta$  for  $0 \leq t' \leq t \leq T$ . Furthermore,  $R_\theta(t, t')$  is uniformly continuous over  $t \in I$  and  $t' \in I'$  for any two maximal intervals  $I, I'$  of  $[0, T] \setminus D$ .

We define  $\mathcal{S}_\theta^{\text{cont}}(T)$  as the subset of  $\mathcal{S}_\theta(T)$  with  $D = \emptyset$  in the above definition.

**Definition C.2 (Admissible space  $\mathcal{S}_\ell(T)$ )** Let  $D \subset (0, T)$  be a finite set. Let  $\mathcal{S}_\ell(T)$  be a set of function triples  $(\Sigma_\ell, R_\ell, \Gamma)$  defined on  $[0, T]^2$  and  $[0, T]$ . We say that  $\mathcal{S}_\ell(T)$  is admissible if there exist constants  $\Phi_\ell, M_\ell > 0$  such that every  $(\Sigma_\ell, R_\ell, \Gamma) \in \mathcal{S}_\ell(T)$  satisfies the following.

- $\Sigma_\ell(t, s)$  is a covariance kernel (in particular, it satisfies  $\Sigma_\ell(t, s) = \Sigma_\ell(s, t)^\top$ ) and satisfies  $\|\Sigma_\ell(t, t)\|_2 \leq \Phi_\ell$  for  $t \in [0, T]$  and  $\Sigma_\ell(0, 0) = \mathbb{E}[\ell_0(r^0; z)\ell_0(r^0; z)^\top]$  for  $r^0 \sim \mathcal{N}(0, \mathbb{E}[\theta^0 \theta^{0\top}])$ . Furthermore,  $\Sigma_\ell(t, t')$  is uniformly continuous over  $t, t' \in I$  for each maximal interval  $I$  of  $[0, T] \setminus D$  and satisfies

$$\|\Sigma_\ell(t, t) - 2\Sigma_\ell(t, t') + \Sigma_\ell(t', t')\|_2 \leq M_\ell |t - t'|, \quad (51)$$

for any  $t, t' \in I$ .

- $R_\ell(t, t')$  satisfies  $R_\ell(t, t') = 0$  for  $0 \leq t < t' \leq T$  and  $\|R_\ell(t, t')\|_2 \leq \Phi_\ell$  for  $0 \leq t' \leq t \leq T$ . Furthermore,  $R_\ell(t, t')$  is uniformly continuous over  $t \in I$  and  $t' \in I'$  for any two maximal intervals  $I, I'$  of  $[0, T] \setminus D$ .
- $\Gamma(t)$  satisfies  $\|\Gamma(t)\|_2 \leq M$  for  $t \in [0, T]$  and  $\Gamma(0) = \mathbb{E}[\nabla_r \ell_0(r^0; z)]$  for  $r^0 \sim \mathcal{N}(0, \mathbb{E}[\theta^0 \theta^{0\top}])$ . Furthermore,  $\Gamma(t)$  is uniformly continuous over  $t \in I$  for each maximal interval  $I$  of  $[0, T] \setminus D$ .

We define  $\mathcal{S}_\ell^{\text{cont}}(T)$  as the subset of  $\mathcal{S}_\ell(T)$  with  $D = \emptyset$  in the above definition.

In the above definitions, we allow for discontinuities at a finite set of time points  $D$  to handle the discretized DMFT system later in the proof of Theorem 3.2 in Appendix D.

We now show that the stochastic processes are uniquely defined given functions in admissible spaces.

**Lemma C.3** Given an admissible space  $\mathcal{S}_\theta(T)$  and any element  $(C_\theta, R_\theta) \in \mathcal{S}_\theta(T)$ , there exists a unique tuple of stochastic processes  $\{r^t, \rho_\ell^{t, t'}, D_\ell^{t, t'}\}_{0 \leq t' \leq t \leq T}$  satisfying Equations (21) to (23). Furthermore, for any  $(C_\theta, R_\theta) \in \mathcal{S}_\theta^{\text{cont}}(T)$ , the processes  $\{r^t, \rho_\ell^{t, t'}, D_\ell^{t, t'}\}_{0 \leq t' \leq t \leq T}$  have continuous sample paths.

Similarly, given an admissible space  $\mathcal{S}_\ell(T)$  and any element  $(\Sigma_\ell, R_\ell, \Gamma) \in \mathcal{S}_\ell(T)$ , there exists a unique pair of stochastic processes  $\{\theta^t, \rho_\theta^{t,t'}\}_{0 \leq t' \leq t \leq T}$  satisfying Equations (17) and (18). Furthermore, for any  $(\Sigma_\ell, R_\ell, \Gamma) \in \mathcal{S}_\ell^{\text{cont}}(T)$ , the processes  $\{\theta^t, \rho_\theta^{t,t'}\}_{0 \leq t' \leq t \leq T}$  have continuous sample paths.

**Proof** First, we show that  $r^t$  is uniquely defined. Let  $\{w^t\}_{t \in [0, T]}$  be a centered Gaussian process with covariance kernel  $C_\theta$ . Then, for any maximal interval  $I$  of  $[0, T] \setminus D$  and any  $t, t' \in I$ , we have

$$\mathbb{E}\|w^t - w^{t'}\|_2^4 \leq 3m^2 \|C_\theta(t, t) - 2C_\theta(t, t') + C_\theta(t', t')\|_2^2 \leq 3m^2 M_\theta^2 (t - t')^2. \quad (52)$$

By the Kolmogorov continuity theorem, there exists a modification of  $w^t$  that is locally Hölder continuous on  $I$ . Then, for each maximal interval  $I$ ,  $r^t$  follows a nonlinear Volterra stochastic integral equation of the second kind with a Lipschitz nonlinearity, a continuous kernel, and a continuous forcing term. By Berger and Mizel [8, Theorem 3.A], it has a unique continuous solution adapted to the filtration  $\mathcal{F}_t^\ell$  generated by  $(\mathbf{B}^s)_{s \leq t}$ . Applying this argument inductively over the maximal intervals of  $[0, T] \setminus D$ , we conclude that  $r^t$  is uniquely defined over  $[0, T]$ . The well-posedness of  $\rho_\ell^{t,t'}$  and  $D_\ell^{t,t'}$  can be shown similarly using the continuity of  $r^t$ .

Next, we show that  $\theta^t$  is uniquely defined. Let  $\{U^t\}_{t \in [0, T]}$  be a centered Gaussian process with covariance kernel  $\Sigma_\ell/\delta$ . Again, by the Kolmogorov continuity theorem, there exists a modification of  $U^t$  that is locally Hölder continuous on each  $I$ . From Equation (17),  $\theta^t$  satisfies the following equation:

$$\theta^t = \theta^0 + U^t - \int_0^t \left( h_s(\theta^s) + \Gamma(s)\theta^s + \left( \int_s^t R_\ell(t', s) dt' \right) \theta^s \right) ds. \quad (53)$$

For each  $I$ , this is a nonlinear Volterra integral equation of the second kind with a continuous kernel and a continuous forcing term. Again, by Berger and Mizel [8, Theorem 3.A], it has a unique continuous solution. Applying this argument inductively over the maximal intervals of  $[0, T] \setminus D$ , we conclude that  $\theta^t$  is uniquely defined over  $[0, T]$ . The well-posedness of  $\rho_\theta^{t,t'}$  can be shown similarly using the continuity of  $\theta^t$ .  $\blacksquare$

Next, we define mappings between the admissible spaces.

First, we define the map  $\mathcal{T}_{\theta \rightarrow \ell}: (C_\theta, R_\theta) \mapsto (\Sigma_\ell, R_\ell, \Gamma)$ . Given  $(C_\theta, R_\theta) \in \mathcal{S}_\theta$ , take the unique processes  $r^t, \rho_\ell^{t,t'}, D_\ell^{t,t'}$  satisfying Equations (21) to (23) whose existence is guaranteed by Lemma C.3. Then, we define  $(\Sigma_\ell, R_\ell, \Gamma)$  by Equations (27) to (29).

Next, we define the map  $\mathcal{T}_{\ell \rightarrow \theta}: (\Sigma_\ell, R_\ell, \Gamma) \mapsto (C_\theta, R_\theta)$ . Given  $(\Sigma_\ell, R_\ell, \Gamma) \in \mathcal{S}_\ell$ , take the unique processes  $\theta^t, \rho_\theta^{t,t'}$  satisfying Equations (17) and (18) whose existence is guaranteed by Lemma C.3. Then, we define  $(C_\theta, R_\theta)$  by Equations (19) and (20).

Finally, we define the composite map  $\mathcal{T} = \mathcal{T}_{\ell \rightarrow \theta} \circ \mathcal{T}_{\theta \rightarrow \ell}$ .

In the following lemma, we show that for sufficiently small  $T > 0$ ,  $\mathcal{T}_{\theta \rightarrow \ell}$  and  $\mathcal{T}_{\ell \rightarrow \theta}$  map  $\mathcal{S}_\theta(T)$  into  $\mathcal{S}_\ell(T)$  and  $\mathcal{S}_\ell(T)$  into  $\mathcal{S}_\theta(T)$ , respectively. We defer the proof to Appendix C.3.

#### Lemma C.4

1. There exists some  $T_* > 0$  such that, for any  $0 < T \leq T_*$ , there exist admissible spaces  $\mathcal{S}_\theta(T)$  and  $\mathcal{S}_\ell(T)$  such that  $\mathcal{T}_{\theta \rightarrow \ell}$  maps  $\mathcal{S}_\theta(T)$  into  $\mathcal{S}_\ell^{\text{cont}}(T)$  and  $\mathcal{T}_{\ell \rightarrow \theta}$  maps  $\mathcal{S}_\ell(T)$  to  $\mathcal{S}_\theta^{\text{cont}}(T)$ .
2. If either  $\tau = 0$  or  $\nabla^2 \ell_i(r; z) = 0$ ,  $T$  can be taken arbitrarily large (thus  $T_* = \infty$ ).

## C.2. Equipping Metrics on $\mathcal{S}_\theta(T)$ and $\mathcal{S}_\ell(T)$

In the following, we fix  $T > 0$  such that Lemma C.4 holds and fix admissible spaces  $\mathcal{S}_\theta := \mathcal{S}_\theta(T)$  and  $\mathcal{S}_\ell := \mathcal{S}_\ell(T)$ .

We equip the spaces  $\mathcal{S}_\theta$  and  $\mathcal{S}_\ell$  with metrics. For a constant  $\lambda > 0$ , we define

$$\text{dist}_\lambda(C_\theta^1, C_\theta^2) := \inf_{w_1 \sim \text{GP}(0, C_\theta^1), w_2 \sim \text{GP}(0, C_\theta^2)} \sup_{0 \leq t \leq T} e^{-\lambda t} \sqrt{\mathbb{E} \|w_1^t - w_2^t\|_2^2}, \quad (54a)$$

$$\text{dist}_\lambda(\Sigma_\ell^1, \Sigma_\ell^2) := \inf_{U_1 \sim \text{GP}(0, \Sigma_\ell^1/\delta), U_2 \sim \text{GP}(0, \Sigma_\ell^2/\delta)} \sup_{0 \leq t \leq T} e^{-\lambda t} \sqrt{\mathbb{E} \|U_1^t - U_2^t\|_2^2}, \quad (54b)$$

$$\text{dist}_\lambda(R_\theta^1, R_\theta^2) := \sup_{0 \leq s \leq t \leq T} e^{-\lambda t} \|R_\theta^1(t, s) - R_\theta^2(t, s)\|, \quad (54c)$$

$$\text{dist}_\lambda(R_\ell^1, R_\ell^2) := \sup_{0 \leq s \leq t \leq T} e^{-\lambda t} \|R_\ell^1(t, s) - R_\ell^2(t, s)\|, \quad (54d)$$

$$\text{dist}_\lambda(\Gamma^1, \Gamma^2) := \sup_{0 \leq t \leq T} e^{-\lambda t} \|\Gamma^1(t) - \Gamma^2(t)\|. \quad (54e)$$

In the first two definitions, the infima are taken over all couplings of the Gaussian processes with given marginal covariances. Finally, for  $X^i = (C_\theta^i, R_\theta^i) \in \mathcal{S}_\theta$  and  $Y^i = (\Sigma_\ell^i, R_\ell^i, \Gamma^i) \in \mathcal{S}_\ell$ , we define the distances

$$\text{dist}_\lambda(X^1, X^2) := \text{dist}_\lambda(C_\theta^1, C_\theta^2) + \text{dist}_\lambda(R_\theta^1, R_\theta^2), \quad (55)$$

$$\text{dist}_\lambda(Y^1, Y^2) := \sqrt{\lambda} \text{dist}_\lambda(\Sigma_\ell^1, \Sigma_\ell^2) + \text{dist}_\lambda(R_\ell^1, R_\ell^2) + \text{dist}_\lambda(\Gamma^1, \Gamma^2). \quad (56)$$

Notice the  $\sqrt{\lambda}$  factor in front of  $\text{dist}_\lambda(\Sigma_\ell^1, \Sigma_\ell^2)$ .

We show several properties of the metric spaces  $(\mathcal{S}_\theta, \text{dist}_\lambda)$  and  $(\mathcal{S}_\ell, \text{dist}_\lambda)$ .

**Lemma C.5** *The metric spaces  $(\mathcal{S}_\theta, \text{dist}_\lambda)$  and  $(\mathcal{S}_\ell, \text{dist}_\lambda)$  are complete.*

**Proof** For finite  $T$ , the distance  $\text{dist}_\lambda$  for  $R_\theta$ ,  $R_\ell$ , and  $\Gamma$  are equivalent to  $L^\infty$  distance and hence complete. Completeness for  $C_\theta$  and  $\Sigma_\ell$  are shown in the proof of Fan et al. [23, Theorem 2.4(b)]. ■

**Lemma C.6** *Let  $X^i = (C_\theta^i, R_\theta^i) \in \mathcal{S}_\theta$  and  $Y^i = \mathcal{T}_{\theta \rightarrow \ell}(X^i) = (\Sigma_\ell^i, R_\ell^i, \Gamma^i) \in \mathcal{S}_\ell$  for  $i = 1, 2$ . Then, there exists a constant  $K > 0$  such that for any sufficiently large  $\lambda > 0$  in the definitions of the metrics (54), we have*

$$\text{dist}_\lambda(Y^1, Y^2) \leq K \cdot \text{dist}_\lambda(X^1, X^2). \quad (57)$$

**Lemma C.7** *Let  $Y^i = (C_\ell^i, R_\ell^i, \Gamma^i) \in \mathcal{S}_\ell$  and  $X^i = \mathcal{T}_{\ell \rightarrow \theta}(Y^i) = (C_\theta^i, R_\theta^i) \in \mathcal{S}_\theta$  for  $i = 1, 2$ . Then, for any  $\varepsilon > 0$ , for sufficiently large  $\lambda > 0$ , we have*

$$\text{dist}_\lambda(X^1, X^2) \leq \varepsilon \cdot \text{dist}_\lambda(Y^1, Y^2). \quad (58)$$

We defer the proof of the last two lemmas to Appendices C.4 and C.5.

Finally, we show that  $\mathcal{T}$  is a contraction mapping under the above metrics, finishing the proof of Theorem 3.1. Take  $T \in [0, T_*]$  where  $T_*$  is as in Lemma C.4, and take admissible spaces  $\mathcal{S}_\theta := \mathcal{S}_\theta(T)$  and  $\mathcal{S}_\ell := \mathcal{S}_\ell(T)$ . By Lemmas C.6 and C.7, we can choose  $\varepsilon < 1/K$  and  $\lambda$  sufficiently large such

that  $\mathcal{T} = \mathcal{T}_{\ell \rightarrow \theta} \circ \mathcal{T}_{\theta \rightarrow \ell}$  is a contraction mapping on the metric space  $(\mathcal{S}_\theta, \text{dist}_\lambda)$  which is complete by Lemma C.5. By the Banach fixed-point theorem, there exists a unique fixed point  $(C_\theta, R_\theta) \in \mathcal{S}_\theta^{\text{cont}}$  such that  $\mathcal{T}(C_\theta, R_\theta) = (C_\theta, R_\theta)$ . Thus, this  $(C_\theta, R_\theta)$  and  $(\Sigma_\ell, R_\ell, \Gamma) = \mathcal{T}_{\theta \rightarrow \ell}(C_\theta, R_\theta) \in \mathcal{S}_\ell^{\text{cont}}$  together form a unique pair of fixed points satisfying the DMFT equation  $\mathfrak{S}$ .

The continuity of the sample paths follows from Lemma C.3.

### C.3. Proof of Lemma C.4

We will use the following bounds repeatedly in the proof.

**Lemma C.8** *Let  $X_1$  be a random variable and  $X_2^s, X_3^s$  be stochastic processes in  $\mathbb{R}^m$  adapted to the Brownian motion  $(B^{s'})_{0 \leq s' < s}$ . Then, for any  $0 \leq t' \leq t$  and any integer  $p \geq 1$ , there exists a constant  $C_p > 0$  such that we have*

$$\begin{aligned} & \mathbb{E} \left\| X_1 + \int_{t'}^t X_2^s ds + \int_{t'}^t X_3^s dB^s \right\|_2^{2p} \\ & \leq 3^{2p-1} \left( \mathbb{E} \|X_1\|_2^{2p} + (t-t')^{2p-1} \int_{t'}^t \mathbb{E} \|X_2^s\|_2^{2p} ds + C_p (t-t')^{p-1} \int_{t'}^t \mathbb{E} \|X_3^s\|_2^{2p} ds \right). \end{aligned} \quad (59)$$

**Proof** By Jensen's inequality, we have

$$\begin{aligned} & \mathbb{E} \left\| X_1 + \int_{t'}^t X_2^s ds + \int_{t'}^t X_3^s dB^s \right\|_2^{2p} \\ & \leq 3^{2p-1} \left( \mathbb{E} \|X_1\|_2^{2p} + (t-t')^{2p-1} \int_{t'}^t \mathbb{E} \|X_2^s\|_2^{2p} ds + \mathbb{E} \left\| \int_{t'}^t X_3^s dB^s \right\|_2^{2p} \right). \end{aligned} \quad (60)$$

By the Burkholder–Davis–Gundy inequality and Jensen's inequality again, we have

$$\mathbb{E} \left\| \int_{t'}^t X_3^s dB^s \right\|_2^{2p} \leq C_p \mathbb{E} \left( \int_{t'}^t \|X_3^s\|_2^2 ds \right)^p \leq C_p (t-t')^{p-1} \int_{t'}^t \mathbb{E} \|X_3^s\|_2^{2p} ds. \quad (61)$$

This proves the claim. Note that we can set  $C_1 = 1$  by the Itô isometry.  $\blacksquare$

#### C.3.1. CONSTRUCTION OF THE ADMISSIBLE SPACES

Take constants  $\Phi_{C_\theta} \geq 7(M + 4T^2M^2)$  and  $\Phi_{R_\theta} \geq 2$ . Define constants  $\Phi_{C_\ell}, \Phi_{\Sigma_\ell}, \Phi_{R_\ell}, \Phi_i$  ( $i = 1, \dots, 5$ ),  $M_\ell, M_\theta$ , and  $\bar{\lambda}$  as follows:

$$\Phi_{C_\ell} := (2M + 6M^2m\Phi_{C_\theta}) \exp \left( 6M^2 \left( \frac{T}{\delta^2} + \frac{\tau}{\delta} \right) \Phi_{R_\theta}^2 T \right), \quad (62)$$

$$\Phi_{\Sigma_\ell} := (T + \tau\delta) \Phi_{C_\ell}, \quad (63)$$

$$\Phi_{R_\ell} := \sqrt{2\Phi_1 + 2\tau\delta\Phi_5}, \quad (64)$$

$$\Phi_1 := \frac{3M^4}{\delta^2} \Phi_{R_\theta}^2 \exp \left( 3M^2 \left( \frac{T}{\delta^2} + \frac{\tau}{\delta} \right) \Phi_{R_\theta}^2 T \right), \quad (65)$$

$$\Phi_2 := \frac{27M^4\Phi_{R_\theta}^4}{\delta^4} \exp\left(27M^4\left(\frac{T^3}{\delta^4} + \frac{C_2\tau^2T}{\delta^2}\right)\Phi_{R_\theta}^4 T\right), \quad (66)$$

$$\Phi_3 := (8M + 648M^4m^2\Phi_{C_\theta}^2) \exp\left(216M^4\left(\frac{T^3}{\delta^4} + \frac{C_2\tau^2T}{\delta^2}\right)\Phi_{R_\theta}^4 T\right), \quad (67)$$

$$\Phi_4 := \frac{27\tau^2\Phi_{R_\theta}^4\Phi_3e^4}{\delta^2} \exp\left(27M^4\left(\frac{T^3}{\delta^4} + \frac{C_2\tau^2T}{\delta^2}\right)\Phi_{R_\theta}^4 T\right), \quad (68)$$

$$\Phi_5 := 2M^2\sqrt{\Phi_2\Phi_4}e^2 \exp\left(4M^2\left(\frac{T}{\delta^2} + \frac{\tau}{\delta}\right)\Phi_{R_\theta}^2 T\right), \quad (69)$$

$$M_\ell := 2(T + \tau\delta)\Phi_{C_\ell}e^{2\bar{\lambda}T}, \quad (70)$$

$$M_\theta := 2\left(\frac{mM_\ell}{\delta} + 2T\left(2M^2 + 3M^2\Phi_{C_\theta}e^{2\bar{\lambda}T} + T^2\Phi_{R_\ell}^2\Phi_{C_\theta}e^{4\bar{\lambda}T}\right)\right), \quad (71)$$

$$\bar{\lambda} := \max\left\{2T(3M^2 + T\Phi_{R_\ell}^2), \frac{6m\Phi_{\Sigma_\ell}}{\delta M}, 2(2M + \Phi_{R_\ell}T)\right\}, \quad (72)$$

where  $C_2 > 0$  is the constant in Lemma C.8 for  $p = 2$ .

For  $T > 0$ , we define the function spaces  $\mathcal{S}_\theta := \mathcal{S}_\theta(T)$  and  $\mathcal{S}_\ell := \mathcal{S}_\ell(T)$  as follows. We define  $\mathcal{S}_\theta$  as the space of pairs of functions  $(C_\theta, R_\theta)$  satisfying the continuity conditions and initial conditions in Definition C.1 and the following bounds for  $0 \leq t' \leq t \leq T$ :

$$\|C_\theta(t, t)\|_2 \leq \Phi_{C_\theta}e^{2\bar{\lambda}t}, \quad \|R_\theta(t, t')\|_2 \leq \Phi_{R_\theta}e^{\bar{\lambda}(t-t')}. \quad (73)$$

Then,  $\mathcal{S}_\theta$  is admissible with parameters  $\Phi_\theta = \max\{\Phi_{C_\theta}e^{2\bar{\lambda}T}, \Phi_{R_\theta}e^{\bar{\lambda}T}\}$  and  $M_\theta$ .

We define  $\mathcal{S}_\ell$  as the space of triples of functions  $(\Sigma_\ell, R_\ell, \Gamma)$  satisfying the continuity conditions and initial conditions in Definition C.2 and the following bounds for  $0 \leq t' \leq t \leq T$ :

$$\|\Sigma_\ell(t, t)\|_2 \leq \frac{\Phi_{\Sigma_\ell}e^{2\bar{\lambda}t}}{\bar{\lambda}}, \quad \|R_\ell(t, t')\|_2 \leq \Phi_{R_\ell}e^{\bar{\lambda}(t-t')}, \quad \|\Gamma(t)\|_2 \leq M. \quad (74)$$

Then,  $\mathcal{S}_\ell$  is admissible with parameters  $\Phi_\ell = \max\{\Phi_{\Sigma_\ell}e^{2\bar{\lambda}T}/\bar{\lambda}, \Phi_{R_\ell}e^{\bar{\lambda}T}\}$  and  $M_\ell$ .

In the following, we show that

1. For sufficiently small  $T$  with  $\bar{\lambda}T \leq 1$ , the mappings  $\mathcal{T}_{\theta \rightarrow \ell}$  and  $\mathcal{T}_{\ell \rightarrow \theta}$  map  $\mathcal{S}_\theta$  and  $\mathcal{S}_\ell$  into each other, respectively.
2. If either  $\tau = 0$  or  $\nabla_r^2 \ell_t(r; z) = 0$ , the above holds for any  $T > 0$ .

Note that it is possible to take  $\bar{\lambda}T \leq 1$  since  $\bar{\lambda}$  is monotonically increasing in  $T$ . Then, we can take  $T_*$  as the supremum of such  $T$ . We provide a rough estimate of  $T_* \gtrsim \delta^2/(\tau^2M^7m^2)$  where  $\gtrsim$  hides subleading terms in  $M, m, \tau, 1/\delta$  and the constant factor in Appendix C.3.4.

The proof is almost identical for both cases; the only difference lies in bounding  $R_\ell$ .

### C.3.2. $\mathcal{T}_{\theta \rightarrow \ell}$ MAPS $\mathcal{S}_\theta$ INTO $\mathcal{S}_\ell$ .

**Condition for  $\Sigma_\ell$ .** We have by the assumptions and the triangle inequality that

$$\mathbb{E}\|\ell_t(r^t; z)\|_2^2 \leq \mathbb{E}(\|\ell_t(0; z)\|_2 + M\|r^t\|_2)^2 \leq 2\mathbb{E}\|\ell_t(0; z)\|_2^2 + 2M^2\mathbb{E}\|r^t\|_2^2 \leq 2M + 2M^2\mathbb{E}\|r^t\|_2^2. \quad (75)$$

We apply Lemma C.8 to  $r^t$  in Equation (21) with

$$\mathbb{E}\|X_1\|_2^2 := \mathbb{E}\|w^t\|_2^2 = \text{tr}(C_\theta(t, t)) \leq m\|C_\theta(t, t)\|_2 \leq m\Phi_{C_\theta}e^{2\bar{\lambda}t}, \quad (76)$$

$$\mathbb{E}\|X_2^s\|_2^2 := \mathbb{E}\left\|\frac{1}{\delta}R_\theta(t, s)\ell_s(r^s; z)\right\|_2^2 \leq \frac{1}{\delta^2}\Phi_{R_\theta}^2e^{2\bar{\lambda}(t-s)}\mathbb{E}\|\ell_s(r^s; z)\|_2^2, \quad (77)$$

$$\mathbb{E}\|X_3^s\|_2^2 := \mathbb{E}\left\|\sqrt{\frac{\tau}{\delta}}R_\theta(t, s)\ell_s(r^s; z)\right\|_2^2 \leq \frac{\tau}{\delta}\Phi_{R_\theta}^2e^{2\bar{\lambda}(t-s)}\mathbb{E}\|\ell_s(r^s; z)\|_2^2, \quad (78)$$

to obtain

$$\mathbb{E}\|r^t\|_2^2 \leq 3\left(m\Phi_{C_\theta}e^{2\bar{\lambda}t} + \left(\frac{T}{\delta^2} + \frac{\tau}{\delta}\right)\Phi_{R_\theta}^2 \int_0^t e^{2\bar{\lambda}(t-s)}\mathbb{E}\|\ell_s(r^s; z)\|_2^2 ds\right). \quad (79)$$

Thus, we have

$$e^{-2\bar{\lambda}t}\mathbb{E}\|\ell_t(r^t; z)\|_2^2 \leq 2M + 6M^2m\Phi_{C_\theta} + 6M^2\left(\frac{T}{\delta^2} + \frac{\tau}{\delta}\right)\Phi_{R_\theta}^2 \int_0^t e^{-2\bar{\lambda}s}\mathbb{E}\|\ell_s(r^s; z)\|_2^2 ds. \quad (80)$$

By Grönwall's inequality, we have

$$\mathbb{E}\|\ell_t(r^t; z)\|_2^2 \leq (2M + 6M^2m\Phi_{C_\theta}) \exp\left(6M^2\left(\frac{T}{\delta^2} + \frac{\tau}{\delta}\right)\Phi_{R_\theta}^2 T\right) e^{2\bar{\lambda}t} = \Phi_{C_\ell}e^{2\bar{\lambda}t}, \quad (81)$$

where we used the definition of  $\Phi_{C_\ell}$  in Equation (62).

Now, we check the condition  $\|\Sigma_\ell(t, t)\|_2 \leq \Phi_{\Sigma_\ell}e^{2\bar{\lambda}t}/\bar{\lambda}$ . Using Lemma C.8, we have

$$\begin{aligned} \|\Sigma_\ell(t, t)\|_2 &= \|\mathbb{E}[L^t L^{t\top}]\|_2 \leq \mathbb{E}\|L^t\|_2^2 \leq 2(T + \tau\delta) \int_0^t \mathbb{E}\|\ell_s(r^s; z)\|_2^2 ds \\ &\leq 2(T + \tau\delta) \int_0^t \Phi_{C_\ell}e^{2\bar{\lambda}s} ds = \frac{(T + \tau\delta)\Phi_{C_\ell}}{\bar{\lambda}}(e^{2\bar{\lambda}t} - 1) \leq \frac{\Phi_{\Sigma_\ell}}{\bar{\lambda}}e^{2\bar{\lambda}t}, \end{aligned} \quad (82)$$

where we used the definition of  $\Phi_{\Sigma_\ell}$  in Equation (63).

Next, we check the condition (51). We have

$$\begin{aligned} \|\Sigma_\ell(t, t) - 2\Sigma_\ell(t, t') + \Sigma_\ell(t', t')\|_2 &= \|\mathbb{E}[(L^t - L^{t'})(L^t - L^{t'})^\top]\|_2 \leq \mathbb{E}\|L^t - L^{t'}\|_2^2 \\ &= \mathbb{E}\left\|\int_{t'}^t \ell_s(r^s; z)(ds + \sqrt{\tau\delta} dB^s)\right\|_2^2 \leq 2(T + \tau\delta) \int_{t'}^t \mathbb{E}\|\ell_s(r^s; z)\|_2^2 ds \\ &\leq 2(T + \tau\delta)\Phi_{C_\ell}e^{2\bar{\lambda}T}|t - t'| \leq M_\ell|t - t'|. \end{aligned} \quad (83)$$

**Condition for  $R_\ell$ : case (1).** We have

$$\|R_\ell(t, t')\|_2^2 \leq 2\mathbb{E}\|\rho_\ell^{t, t'}\|_2^2 + 2\tau\delta\mathbb{E}\|D_\ell^{t, t'}\|_2^2. \quad (84)$$

First, we bound  $\mathbb{E}\|\rho_\ell^{t,t'}\|_2^2$ . By the Lipschitz continuity of  $\ell_t(r; z)$  in  $r$ , we have  $\|\nabla_r \ell_t(r^t; z)\|_2 \leq M$  for all  $t$ . Therefore, we have

$$\mathbb{E}\|\rho_\ell^{t,t'}\|_2^2 = \mathbb{E}\|\nabla_r \ell_t(r^t; z) \rho_r^{t,t'}\|_2^2 \leq M^2 \mathbb{E}\|\rho_r^{t,t'}\|_2^2. \quad (85)$$

Applying Lemma C.8 to  $\rho_r^{t,t'}$  in Equation (22), we have

$$\mathbb{E}\|\rho_r^{t,t'}\|_2^2 \leq 3 \left( \frac{M^2}{\delta^2} \Phi_{R_\theta}^2 e^{2\bar{\lambda}(t-t')} + \left( \frac{T}{\delta^2} + \frac{\tau}{\delta} \right) \Phi_{R_\theta}^2 \int_{t'}^t e^{2\bar{\lambda}(t-s)} \mathbb{E}\|\rho_\ell^{s,t'}\|_2^2 ds \right). \quad (86)$$

Thus, we have

$$e^{-2\bar{\lambda}(t-t')} \mathbb{E}\|\rho_\ell^{t,t'}\|_2^2 \leq \frac{3M^4}{\delta^2} \Phi_{R_\theta}^2 + 3M^2 \left( \frac{T}{\delta^2} + \frac{\tau}{\delta} \right) \Phi_{R_\theta}^2 \int_{t'}^t e^{-2\bar{\lambda}(s-t')} \mathbb{E}\|\rho_\ell^{s,t'}\|_2^2 ds. \quad (87)$$

By Grönwall's inequality, we have

$$e^{-2\bar{\lambda}(t-t')} \mathbb{E}\|\rho_\ell^{t,t'}\|_2^2 \leq \frac{3M^4}{\delta^2} \Phi_{R_\theta}^2 \exp \left( 3M^2 \left( \frac{T}{\delta^2} + \frac{\tau}{\delta} \right) \Phi_{R_\theta}^2 T \right) \leq \Phi_1. \quad (88)$$

Next, we bound  $\mathbb{E}\|D_\ell^{t,t'}\|_2^2$ . We have

$$\begin{aligned} \mathbb{E}\|D_\ell^{t,t'}\|_2^2 &= 4M^2 \left( \frac{T}{\delta^2} + \frac{\tau}{\delta} \right) \int_{t'}^t \|R_\theta(t, s)\|_2^2 \mathbb{E}\|D_\ell^{s,t'}\|_2^2 ds + 2 \mathbb{E}[\|\nabla_r^2 \ell_t(r^t; z)\|_2^2 \|D_r^{t,t'}\|_2^2 \|\rho_r^{t,t'}\|_2^2] \\ &\leq 4M^2 \left( \frac{T}{\delta^2} + \frac{\tau}{\delta} \right) \Phi_{R_\theta}^2 \int_{t'}^t e^{2\bar{\lambda}(t-s)} \mathbb{E}\|D_\ell^{s,t'}\|_2^2 ds + 2M^2 \sqrt{\mathbb{E}\|D_r^{t,t'}\|_2^4} \sqrt{\mathbb{E}\|\rho_r^{t,t'}\|_2^4}, \end{aligned} \quad (89)$$

where we used that  $\|\nabla_r^2 \ell_t(r^t; z)\|_2 \leq M$  by assumption. We first bound  $\mathbb{E}\|\rho_r^{t,t'}\|_2^4$ . Applying Lemma C.8, we obtain

$$\begin{aligned} \mathbb{E}\|\rho_r^{t,t'}\|_2^4 &\leq 27 \left( \left( \frac{T^3}{\delta^4} + \frac{C_2 \tau^2 T}{\delta^2} \right) \int_{t'}^t \|R_\theta(t, s)\|_2^4 \mathbb{E}[\|\nabla_r \ell_s(r^s; z)\|_2^4 \|\rho_r^{s,t'}\|_2^4] ds \right. \\ &\quad \left. + \frac{1}{\delta^4} \|R_\theta(t, t')\|_2^4 \mathbb{E}\|\nabla_r \ell_{t'}(r^{t'}; z)\|_2^4 \right) \\ &\leq 27 \left( \left( \frac{T^3}{\delta^4} + \frac{C_2 \tau^2 T}{\delta^2} \right) \Phi_{R_\theta}^4 M^4 \int_{t'}^t e^{4\bar{\lambda}(t-s)} \mathbb{E}\|\rho_r^{s,t'}\|_2^4 ds + \frac{M^4}{\delta^4} \Phi_{R_\theta}^4 e^{4\bar{\lambda}(t-t')} \right). \end{aligned} \quad (90)$$

By Grönwall's inequality, we have

$$\mathbb{E}\|\rho_r^{t,t'}\|_2^4 \leq \frac{27M^4 \Phi_{R_\theta}^4}{\delta^4} \exp \left( 27M^4 \left( \frac{T^3}{\delta^4} + \frac{C_2 \tau^2 T}{\delta^2} \right) \Phi_{R_\theta}^4 T \right) e^{4\bar{\lambda}(t-t')} \leq \Phi_2 e^{4\bar{\lambda}(t-t')}. \quad (91)$$

We next bound  $\mathbb{E}\|D_r^{t,t'}\|_2^4$ . Applying Lemma C.8, we have

$$\mathbb{E}\|D_r^{t,t'}\|_2^4 \leq 27 \left( \left( \frac{T^3}{\delta^4} + \frac{C_2 \tau^2 T}{\delta^2} \right) \int_{t'}^t \|R_\theta(t, s)\|_2^4 \mathbb{E}[\|\nabla_r \ell_s(r^s; z)\|_2^4 \|D_r^{s,t'}\|_2^4] ds \right)$$

$$\begin{aligned}
 & + \frac{\tau^2}{\delta^2} \|R_\theta(t, t')\|_2^4 \mathbb{E} \|\ell_{t'}(r^{t'}; z)\|_2^4 \Big) \\
 & \leq 27 \left( \left( \frac{T^3}{\delta^4} + \frac{C_2 \tau^2 T}{\delta^2} \right) \Phi_{R_\theta}^4 M^4 \int_{t'}^t e^{4\bar{\lambda}(t-s)} \mathbb{E} \|D_r^{s,t'}\|_2^4 ds + \frac{\tau^2}{\delta^2} \Phi_{R_\theta}^4 e^{4\bar{\lambda}(t-t')} \mathbb{E} \|\ell_{t'}(r^{t'}; z)\|_2^4 \right). \tag{92}
 \end{aligned}$$

$\mathbb{E} \|\ell_{t'}(r^{t'}; z)\|_2^4$  can be bounded as

$$\begin{aligned}
 & \mathbb{E} \|\ell_t(r^t; z)\|_2^4 \leq 8 \mathbb{E} \|\ell_t(0; z)\|_2^4 + 8M^4 \mathbb{E} \|r^t\|_2^4 \\
 & \leq 8M + 8M^4 \cdot 27 \left( \mathbb{E} \|w^t\|_2^4 + \left( \frac{T^3}{\delta^4} + \frac{C_2 \tau^2 T}{\delta^2} \right) \int_0^t \|R_\theta(t, s)\|_2^4 \mathbb{E} \|\ell_s(r^s; z)\|_2^4 ds \right) \\
 & \leq 8M + 216M^4 \left( 3m^2 \Phi_{C_\theta}^2 e^{4\bar{\lambda}t} + \left( \frac{T^3}{\delta^4} + \frac{C_2 \tau^2 T}{\delta^2} \right) \Phi_{R_\theta}^4 \int_0^t e^{4\bar{\lambda}(t-s)} \mathbb{E} \|\ell_s(r^s; z)\|_2^4 ds \right). \tag{93}
 \end{aligned}$$

By Grönwall's inequality, we have

$$\mathbb{E} \|\ell_t(r^t; z)\|_2^4 \leq (8M + 648M^4 m^2 \Phi_{C_\theta}^2) \exp \left( 216M^4 \left( \frac{T^3}{\delta^4} + \frac{C_2 \tau^2 T}{\delta^2} \right) \Phi_{R_\theta}^4 T \right) e^{4\bar{\lambda}t} \leq \Phi_3 e^{4\bar{\lambda}t}. \tag{94}$$

Taking  $T$  small enough such that  $\bar{\lambda}T \leq 1$ , we have  $e^{4\bar{\lambda}t} \leq e^4$ . Thus, we have

$$\mathbb{E} \|D_r^{t,t'}\|_2^4 \leq 27 \left( \left( \frac{T^3}{\delta^4} + \frac{C_2 \tau^2 T}{\delta^2} \right) \Phi_{R_\theta}^4 M^4 \int_{t'}^t e^{4\bar{\lambda}(t-s)} \mathbb{E} \|D_r^{s,t'}\|_2^4 ds + \frac{\tau^2}{\delta^2} \Phi_{R_\theta}^4 \Phi_3 e^4 e^{4\bar{\lambda}(t-t')} \right). \tag{95}$$

By Grönwall's inequality, we have

$$\mathbb{E} \|D_r^{t,t'}\|_2^4 \leq \frac{27\tau^2 \Phi_{R_\theta}^4 \Phi_3 e^4}{\delta^2} \exp \left( 27M^4 \left( \frac{T^3}{\delta^4} + \frac{C_2 \tau^2 T}{\delta^2} \right) \Phi_{R_\theta}^4 T \right) e^{4\bar{\lambda}(t-t')} \leq \Phi_4 e^4. \tag{96}$$

Combining the above bounds, we have

$$\mathbb{E} \|D_\ell^{t,t'}\|_2^2 \leq 4M^2 \left( \frac{T}{\delta^2} + \frac{\tau}{\delta} \right) \Phi_{R_\theta}^2 \int_{t'}^t e^{2\bar{\lambda}(t-s)} \mathbb{E} \|D_\ell^{s,t'}\|_2^2 ds + 2M^2 \sqrt{\Phi_4 e^4} \sqrt{\Phi_2 e^{4\bar{\lambda}(t-t')}}. \tag{97}$$

By Grönwall's inequality, we have

$$\mathbb{E} \|D_\ell^{t,t'}\|_2^2 \leq 2M^2 \sqrt{\Phi_2 \Phi_4} e^2 \exp \left( 4M^2 \left( \frac{T}{\delta^2} + \frac{\tau}{\delta} \right) \Phi_{R_\theta}^2 T \right) e^{2\bar{\lambda}(t-t')} \leq \Phi_5 e^{2\bar{\lambda}(t-t')}. \tag{98}$$

Thus, we have

$$\|R_\ell(t, t')\|_2^2 \leq (2\Phi_1 + 2\tau\delta\Phi_5) e^{2\bar{\lambda}(t-t')} = \Phi_{R_\ell}^2 e^{2\bar{\lambda}(t-t')}. \tag{99}$$

**Condition for  $R_\ell$ : case (2).** When  $\nabla_r^2 \ell_t(r; z) = 0$ , we have  $D_\ell^{t,t'} = 0$ . Thus, when  $\tau = 0$  or  $\nabla^2 \ell_t(r; z) = 0$ , we have

$$\|R_\ell(t, t')\|_2^2 \leq \mathbb{E} \|\rho_\ell^{t,t'}\|_2^2 \leq \Phi_1 e^{-2\bar{\lambda}(t-t')} \leq \Phi_{R_\ell}^2 e^{2\bar{\lambda}(t-t')}. \tag{100}$$

This holds without taking  $T$  small.

**Condition for  $\Gamma$ .** We have

$$\|\Gamma(t)\|_2 \leq \mathbb{E}\|\nabla_r \ell_t(r^t; z)\|_2 \leq M. \quad (101)$$

C.3.3.  $\mathcal{T}_{\ell \rightarrow \theta}$  MAPS  $\mathcal{S}_\ell$  INTO  $\mathcal{S}_\theta$ .

**Condition for  $C_\theta$ .** We have

$$\mathbb{E}\|U^t\|_2^2 = \frac{1}{\delta} \text{tr}(\Sigma_\ell(t, t)) \leq \frac{m}{\delta} \|\Sigma_\ell(t, t)\|_2 \leq \frac{m\Phi_{\Sigma_\ell}}{\delta\bar{\lambda}} e^{2\bar{\lambda}t}. \quad (102)$$

By the Lipschitz continuity of  $h$ , we have  $\mathbb{E}\|h_t(\theta^t)\|_2^2 \leq 2\|h_t(0)\|_2^2 + 2\mathbb{E}\|h_t(\theta^t) - h_t(0)\|_2^2 \leq 2M^2 + 2M^2 \mathbb{E}\|\theta^t\|_2^2$ . Thus, we have

$$\begin{aligned} & e^{-2\bar{\lambda}t} \mathbb{E}\|\theta^t\|_2^2 \\ & \leq 3e^{-2\bar{\lambda}t} \left( \mathbb{E}\|\theta^0\|_2^2 + \mathbb{E}\|U^t\|_2^2 \right. \\ & \quad \left. + 2T \int_0^t \left( \mathbb{E}\|h_s(\theta^s)\|_2^2 + \|\Gamma(s)\|_2^2 \mathbb{E}\|\theta^s\|_2^2 + T \int_0^s \|R_\ell(s, s')\|_2^2 \mathbb{E}\|\theta^{s'}\|_2^2 ds' \right) ds \right) \\ & \leq 3 \left( e^{-2\bar{\lambda}t} (M + 4T^2M^2) + \frac{m\Phi_{\Sigma_\ell}}{\delta\bar{\lambda}} \right. \\ & \quad \left. + 2T \int_0^t e^{-2\bar{\lambda}(t-s)} \left( 3M^2 e^{-2\bar{\lambda}s} \mathbb{E}\|\theta^s\|_2^2 + T\Phi_{R_\ell}^2 \int_0^s e^{-2\bar{\lambda}s'} \mathbb{E}\|\theta^{s'}\|_2^2 ds' \right) ds \right) \\ & \leq 3 \left( M + 4T^2M^2 + \frac{m\Phi_{\Sigma_\ell}}{\delta\bar{\lambda}} + 2T(3M^2 + T^2\Phi_{R_\ell}^2) \left( \int_0^t e^{-2\bar{\lambda}(t-s)} ds \right) \sup_{s \in [0, T]} e^{-2\bar{\lambda}s} \mathbb{E}\|\theta^s\|_2^2 \right) \\ & \leq 3 \left( M + 4T^2M^2 + \frac{m\Phi_{\Sigma_\ell}}{\delta\bar{\lambda}} + \frac{T(3M^2 + T^2\Phi_{R_\ell}^2)}{\bar{\lambda}} \sup_{s \in [0, T]} e^{-2\bar{\lambda}s} \mathbb{E}\|\theta^s\|_2^2 \right) \end{aligned} \quad (103)$$

By  $\bar{\lambda} \geq 2T(3M^2 + T^2\Phi_{R_\ell}^2)$ , we have

$$\sup_{t \in [0, T]} e^{-2\bar{\lambda}t} \mathbb{E}\|\theta^t\|_2^2 \leq 6 \left( M + 4T^2M^2 + \frac{m\Phi_{\Sigma_\ell}}{\delta\bar{\lambda}} \right). \quad (104)$$

By  $\bar{\lambda} \geq 6m\Phi_{\Sigma_\ell}/(\delta(M + 4T^2M^2))$ , the right-hand side is bounded by  $7(M + 4T^2M^2) \leq \Phi_{C_\theta}$ . Therefore, we have

$$\|C_\theta(t, t)\|_2 \leq \mathbb{E}\|\theta^t\|_2^2 \leq \Phi_{C_\theta} e^{2\bar{\lambda}t}. \quad (105)$$

Next, we check the condition (50). We have

$$\begin{aligned} & \|C_\theta(t, t) - 2C_\theta(t, t') + C_\theta(t', t')\|_2 = \|\mathbb{E}[(\theta^t - \theta^{t'}) (\theta^t - \theta^{t'})^\top]\|_2 \leq \mathbb{E}\|\theta^t - \theta^{t'}\|_2^2 \\ & \leq \mathbb{E} \left\| U^t - U^{t'} - \int_{t'}^t \left( h_s(\theta^s) + \Gamma(s)\theta^s + \int_0^s R_\ell(s, s')\theta^{s'} ds' \right) ds \right\|_2^2 \end{aligned}$$

$$\begin{aligned}
 &\leq 2 \left( \mathbb{E} \|U^t - U^{t'}\|_2^2 + 2(t-t')^2 \sup_{s \in [0, T]} \left\{ 2M^2 + 3M^2 \mathbb{E} \|\theta^s\|_2^2 + T \int_0^s \|R_\ell(s, s')\|_2^2 \mathbb{E} \|\theta^{s'}\|_2^2 ds' \right\} \right) \\
 &\leq 2 \left( \frac{mM_\ell}{\delta} |t-t'| + 2(t-t')^2 \left( 2M^2 + 3M^2 \Phi_{C_\theta} e^{2\bar{\lambda}T} + T^2 \Phi_{R_\ell}^2 \Phi_{C_\theta} e^{4\bar{\lambda}T} \right) \right) \\
 &\leq M_\theta |t-t'|, \tag{106}
 \end{aligned}$$

where we used

$$\mathbb{E} \|U^t - U^{t'}\|_2^2 \leq \frac{m}{\delta} \|\Sigma_\ell(t, t) - 2\Sigma_\ell(t, t') + \Sigma_\ell(t', t')\|_2 \leq \frac{mM_\ell}{\delta} |t-t'|. \tag{107}$$

**Condition for  $R_\theta$ .** By the Lipschitz continuity of  $h$ , we have  $\|\nabla_\theta h_t(\theta)\|_2 \leq M$  for all  $t$  and  $\theta$ . Thus, we have

$$\begin{aligned}
 \|\rho_\theta^{t, t'}\|_2 &\leq 1 + \int_{t'}^t \left( (\|\nabla_\theta h_s(\theta^s)\|_2 + \|\Gamma(s)\|_2) \|\rho_\theta^{s, t'}\|_2 + \int_{t'}^s \|R_\ell(s, s')\|_2 \|\rho_\theta^{s', t'}\|_2 ds' \right) ds \\
 &\leq 1 + \int_{t'}^t \left( 2M \|\rho_\theta^{s, t'}\|_2 + \Phi_{R_\ell} \int_{t'}^s e^{\bar{\lambda}(s-s')} \|\rho_\theta^{s', t'}\|_2 ds' \right) ds, \tag{108}
 \end{aligned}$$

and thus

$$\begin{aligned}
 &e^{-\bar{\lambda}(t-t')} \|\rho_\theta^{t, t'}\|_2 \\
 &\leq 1 + \int_{t'}^t e^{-\bar{\lambda}(t-s)} \left( 2M e^{-\bar{\lambda}(s-t')} \|\rho_\theta^{s, t'}\|_2 + \Phi_{R_\ell} \int_{t'}^s e^{\bar{\lambda}(s'-t')} \|\rho_\theta^{s', t'}\|_2 ds' \right) ds \\
 &\leq 1 + (2M + \Phi_{R_\ell} T) \left( \int_{t'}^t e^{-\bar{\lambda}(t-s)} ds \right) \sup_{s \in [t', t]} e^{-\bar{\lambda}(s-t')} \|\rho_\theta^{s, t'}\|_2 \\
 &\leq 1 + \frac{2M + \Phi_{R_\ell} T}{\bar{\lambda}} \sup_{s \in [t', t]} e^{-\bar{\lambda}(s-t')} \|\rho_\theta^{s, t'}\|_2. \tag{109}
 \end{aligned}$$

By  $\bar{\lambda} \geq 2(2M + \Phi_{R_\ell} T)$ , we have

$$\|\rho_\theta^{t, t'}\|_2 \leq 2e^{\bar{\lambda}(t-t')}, \tag{110}$$

and thus

$$\|R_\theta(t, t')\|_2 \leq \mathbb{E} \|\rho_\theta^{t, t'}\|_2 \leq 2e^{\bar{\lambda}(t-t')} \leq \Phi_{R_\theta} e^{\bar{\lambda}(t-t')}. \tag{111}$$

#### C.3.4. A ROUGH ESTIMATE OF $T_*$

We derive a rough lower bound on  $T_*$  up to the leading dependencies on  $M, m, \tau, 1/\delta$  and ignoring constant factors. Take  $T \leq \min\{1, (M^2(1/\delta^2 + \tau/\delta)\Phi_{R_\theta}^2)^{-1}\}$ . Then, the exponents in the definitions of  $\Phi_{C_\ell}$  and  $\Phi_i$  ( $i = 1, \dots, 5$ ) are all bounded by constants since

$$M^2 \left( \frac{T}{\delta^2} + \frac{\tau}{\delta} \right) \Phi_{R_\theta}^2 T \leq 1, \quad M^4 \left( \frac{T^2}{\delta^4} + \frac{C_2 \tau^2}{\delta^2} \right) \Phi_{R_\theta}^4 T^2 \lesssim 1. \tag{112}$$

Then, the  $\Phi$  quantities can be bounded as

$$\begin{aligned} \Phi_{C_\theta} \lesssim M, \quad \Phi_{R_\theta} \lesssim 1, \quad \Phi_{C_\ell} \lesssim M^3 m, \quad \Phi_{\Sigma_\ell} \lesssim \tau M^3 m, \quad \Phi_1 \lesssim \frac{M^4}{\delta^2}, \quad \Phi_2 \lesssim \frac{M^4}{\delta^4}, \\ \Phi_3 \lesssim M^6 m^2, \quad \Phi_4 \lesssim \frac{\tau^2 M^6 m^2}{\delta^2}, \quad \Phi_5 \lesssim \frac{\tau M^7 m}{\delta^3}, \quad \Phi_{R_\ell} \leq \frac{\tau M^{7/2} m^{1/2}}{\delta} \end{aligned} \quad (113)$$

Then,  $\bar{\lambda}$  is bounded as

$$\bar{\lambda} \lesssim (M^2 + \Phi_{R_\ell}^2) + \frac{m\Phi_{\Sigma_\ell}}{\delta M} + (M + \Phi_{R_\ell}) \lesssim \frac{\tau^2 M^7 m^2}{\delta^2}. \quad (114)$$

Thus, further taking  $T \lesssim \delta^2/(\tau^2 M^7 m^2)$ , we have  $\bar{\lambda}T \lesssim 1$ .

#### C.4. Proof of Lemma C.6

Set  $\Phi := \max\{\Phi_\theta, \Phi_\ell\}$ . In the following,  $K$  denotes a positive constant that may depend on  $M, m, T, \delta, \tau$ , and  $\Phi$ , but not on other variables, and may change from line to line.

**Bound of  $\text{dist}_\lambda(\Sigma_\ell^1, \Sigma_\ell^2)$ .** Let  $w_1 \sim \text{GP}(0, C_\theta^1)$  and  $w_2 \sim \text{GP}(0, C_\theta^2)$  be Gaussian processes coupled such that

$$\sup_{t \in [0, T]} e^{-\lambda t} \sqrt{\mathbb{E}\|w_1^t - w_2^t\|_2^2} \leq 2 \cdot \text{dist}_\lambda(C_\theta^1, C_\theta^2). \quad (115)$$

For  $i = 1, 2$ , let  $r_i$  be the solution of

$$r_i^t = w_i^t - \frac{1}{\delta} \int_0^t R_\theta^i(t, s) \ell_s(r_i^s; z) (ds + \sqrt{\tau \delta} dB^s). \quad (116)$$

Note that we use the same Brownian motion  $B^t$  for  $i = 1, 2$ . Applying Lemma C.8 to  $r_1^t - r_2^t$  with

$$\mathbb{E}\|w_1^t - w_2^t\|_2^2 \leq 4e^{2\lambda t} \cdot \text{dist}_\lambda(C_\theta^1, C_\theta^2)^2, \quad (117)$$

$$\begin{aligned} \mathbb{E}\|R_\theta^1(t, s) \ell_s(r_1^s; z) - R_\theta^2(t, s) \ell_s(r_2^s; z)\|_2^2 \\ \leq 2(\|R_\theta^1(t, s) - R_\theta^2(t, s)\|_2^2 \mathbb{E}\|\ell_s(r_1^s; z)\|_2^2 + \|R_\theta^2(t, s)\|_2^2 \mathbb{E}\|\ell_s(r_1^s; z) - \ell_s(r_2^s; z)\|_2^2) \\ \leq 2\left(\Phi e^{2\lambda t} \cdot \text{dist}_\lambda(R_\theta^1, R_\theta^2)^2 + \Phi^2 M^2 \mathbb{E}\|r_1^s - r_2^s\|_2^2\right), \end{aligned} \quad (118)$$

we have

$$\begin{aligned} \mathbb{E}\|r_1^t - r_2^t\|_2^2 \leq 3\left(4e^{2\lambda t} \cdot \text{dist}_\lambda(C_\theta^1, C_\theta^2)^2 + 2\left(\frac{T}{\delta^2} + \frac{\tau}{\delta}\right) T \Phi e^{2\lambda t} \cdot \text{dist}_\lambda(R_\theta^1, R_\theta^2)^2\right. \\ \left. + 2\left(\frac{T}{\delta^2} + \frac{\tau}{\delta}\right) \Phi^2 M^2 \int_0^t \mathbb{E}\|r_1^s - r_2^s\|_2^2 ds\right), \end{aligned} \quad (119)$$

and thus

$$e^{-2\lambda t} \mathbb{E}\|r_1^t - r_2^t\|_2^2 \leq K \cdot \text{dist}_\lambda(X^1, X^2)^2 + K \int_0^t e^{-2\lambda s} \mathbb{E}\|r_1^s - r_2^s\|_2^2 ds. \quad (120)$$

By Grönwall's inequality, we have

$$e^{-2\lambda t} \mathbb{E} \|r_1^t - r_2^t\|_2^2 \leq K e^{KT} \cdot \text{dist}_\lambda(X^1, X^2)^2 \leq K \cdot \text{dist}_\lambda(X^1, X^2)^2. \quad (121)$$

For  $i = 1, 2$ , let

$$L_i^t := \int_0^t \ell_s(r_i^s; z) (ds + \sqrt{\tau\delta} dB^s). \quad (122)$$

Then, we have

$$\begin{aligned} e^{-2\lambda t} \mathbb{E} \|L_1^t - L_2^t\|_2^2 &\leq 2(T + \tau\delta) \int_0^t e^{-2\lambda(t-s)} \cdot e^{-2\lambda s} \mathbb{E} \|\ell_s(r_1^s; z) - \ell_s(r_2^s; z)\|_2^2 ds \\ &\leq \frac{(T + \tau\delta)M^2}{\lambda} \sup_{s \in [0, T]} e^{-2\lambda s} \mathbb{E} \|r_1^s - r_2^s\|_2^2 \leq \frac{K}{\lambda} \cdot \text{dist}_\lambda(X^1, X^2)^2. \end{aligned} \quad (123)$$

Let  $\{(U_1^t, U_2^t)\}_{t \in [0, T]}$  be a centered Gaussian process with covariance  $\mathbb{E} \left[ \begin{pmatrix} L_1^t \\ L_2^t \end{pmatrix} \begin{pmatrix} L_1^t \\ L_2^t \end{pmatrix}^\top \right] / \delta$ .

Since  $U_1$  and  $U_2$  have covariance kernels  $\Sigma_\ell^1 / \delta$  and  $\Sigma_\ell^2 / \delta$  respectively, we have

$$\text{dist}_\lambda(\Sigma_\ell^1, \Sigma_\ell^2) \leq \sup_{t \in [0, T]} e^{-\lambda t} \sqrt{\mathbb{E} \|U_1^t - U_2^t\|_2^2} = \sup_{t \in [0, T]} e^{-\lambda t} \sqrt{\mathbb{E} \|L_1^t - L_2^t\|_2^2 / \delta} \leq \sqrt{\frac{K}{\lambda}} \cdot \text{dist}_\lambda(X^1, X^2). \quad (124)$$

**Bound of  $\text{dist}_\lambda(\Gamma^1, \Gamma^2)$ .** By Equation (121), we have

$$\begin{aligned} \text{dist}_\lambda(\Gamma^1, \Gamma^2) &= \sup_{t \in [0, T]} e^{-\lambda t} \|\Gamma^1(t) - \Gamma^2(t)\|_2 \leq \sup_{t \in [0, T]} e^{-\lambda t} \sqrt{\mathbb{E} \|\nabla_r \ell_t(r_1^t; z) - \nabla_r \ell_t(r_2^t; z)\|_2^2} \\ &\leq M \sup_{t \in [0, T]} e^{-\lambda t} \sqrt{\mathbb{E} \|r_1^t - r_2^t\|_2^2} \leq K \cdot \text{dist}_\lambda(X^1, X^2). \end{aligned} \quad (125)$$

**Bound of  $\text{dist}_\lambda(R_\ell^1, R_\ell^2)$ .** For  $i = 1, 2$ , let

$$\rho_{\ell, i}^{t, t'} = \nabla_r \ell_t(r_i^t; z) \rho_{r, i}^{t, t'}, \quad \rho_{r, i}^{t, t'} = -\frac{1}{\delta} \int_{t'}^t R_\theta^i(t, s) \rho_{\ell, i}^{s, t'} (ds + \sqrt{\tau\delta} dB^s) - \frac{1}{\delta} R_\theta^i(t, t') \nabla_r \ell_{t'}(r_i^{t'}; z). \quad (126)$$

Then, we have

$$\|R_\ell^1(t, t') - R_\ell^2(t, t')\|_2^2 \leq 2 \mathbb{E} \|\rho_{\ell, 1}^{t, t'} - \rho_{\ell, 2}^{t, t'}\|_2^2 + 2\tau\delta \mathbb{E} \|D_{\ell, 1}^{t, t'} - D_{\ell, 2}^{t, t'}\|_2^2. \quad (127)$$

We first bound  $\mathbb{E} \|\rho_{\ell, 1}^{t, t'} - \rho_{\ell, 2}^{t, t'}\|_2^2$ . We have

$$\begin{aligned} \mathbb{E} \|\rho_{\ell, 1}^{t, t'} - \rho_{\ell, 2}^{t, t'}\|_2^2 &= \mathbb{E} \|\nabla_r \ell_t(r_1^t; z) \rho_{r, 1}^{t, t'} - \nabla_r \ell_t(r_2^t; z) \rho_{r, 2}^{t, t'}\|_2^2 \\ &\leq 2 \mathbb{E} \|(\nabla_r \ell_t(r_1^t; z) - \nabla_r \ell_t(r_2^t; z)) \rho_{r, 1}^{t, t'}\|_2^2 + 2 \mathbb{E} \|\nabla_r \ell_t(r_2^t; z) (\rho_{r, 1}^{t, t'} - \rho_{r, 2}^{t, t'})\|_2^2 \\ &\leq 2M^2 \sqrt{\mathbb{E} \|\rho_{r, 1}^{t, t'}\|_2^4} \sqrt{\mathbb{E} \|r_1^t - r_2^t\|_2^4} + 2M^2 \mathbb{E} \|\rho_{r, 1}^{t, t'} - \rho_{r, 2}^{t, t'}\|_2^2. \end{aligned} \quad (128)$$

First, we bound  $\mathbb{E}\|\rho_{r,1}^{t,t'}\|_2^4$ . Applying Lemma C.8 to  $\rho_{r,1}^{t,t'}$ , we have

$$\mathbb{E}\|\rho_{r,1}^{t,t'}\|_2^4 \leq 27 \left( \frac{\Phi^4 M^4}{\delta^4} + \left( \frac{T^3}{\delta^4} + \frac{3T\tau^2}{\delta^2} \right) \Phi^4 M^4 \int_{t'}^t \mathbb{E}\|\rho_{r,1}^{s,t'}\|_2^4 ds \right) \leq K + K \int_{t'}^t \mathbb{E}\|\rho_{r,1}^{s,t'}\|_2^4 ds, \quad (129)$$

By Grönwall's inequality, we have

$$\mathbb{E}\|\rho_{r,1}^{t,t'}\|_2^4 \leq K e^{KT} \leq K. \quad (130)$$

Next, we bound  $\mathbb{E}\|r_1^t - r_2^t\|_2^4$  by applying Lemma C.8 with

$$\begin{aligned} \mathbb{E}\|w_1^t - w_2^t\|_2^4 &\leq 3(\mathbb{E}\|w_1^t - w_2^t\|_2^2)^2 \leq 12e^{4\lambda t} \cdot \text{dist}_\lambda(C_\theta^1, C_\theta^2)^4, \\ \mathbb{E}\|R_\theta^1(t, s)\ell_s(r_1^s; z) - R_\theta^2(t, s)\ell_s(r_2^s; z)\|_2^4 &\leq 8(\|R_\theta^1(t, s) - R_\theta^2(t, s)\|_2^4 \mathbb{E}\|\ell_s(r_1^s; z)\|_2^4 + \|R_\theta^2(t, s)\|_2^4 \mathbb{E}\|\ell_s(r_1^s; z) - \ell_s(r_2^s; z)\|_2^4) \\ &\leq 8 \left( K e^{4\lambda t} \cdot \text{dist}_\lambda(R_\theta^1, R_\theta^2)^4 + \Phi^4 M^4 \mathbb{E}\|r_1^s - r_2^s\|_2^4 \right), \end{aligned} \quad (131)$$

Here, we used that  $\mathbb{E}\|\ell_t(r^t; z)\|_2^4$  is uniformly bounded by some constant  $K > 0$  by a similar argument as the bound on  $\mathbb{E}\|\rho_r^{t,t'}\|_2^4$ . Then, we have

$$\mathbb{E}\|r_1^t - r_2^t\|_2^4 \leq K e^{4\lambda t} \cdot \text{dist}_\lambda(X^1, X^2)^4 + K \int_0^t \mathbb{E}\|r_1^s - r_2^s\|_2^4 ds. \quad (132)$$

By Grönwall's inequality, we have

$$e^{-4\lambda t} \mathbb{E}\|r_1^t - r_2^t\|_2^4 \leq K e^{KT} \cdot \text{dist}_\lambda(X^1, X^2)^4 \leq K \cdot \text{dist}_\lambda(X^1, X^2)^4. \quad (133)$$

Finally, we bound  $\mathbb{E}\|\rho_{r,1}^{t,t'} - \rho_{r,2}^{t,t'}\|_2^2$ . We apply Lemma C.8 to  $\rho_{r,1}^{t,t'} - \rho_{r,2}^{t,t'}$  with

$$\begin{aligned} &\mathbb{E}\|R_\theta^1(t, t')\nabla_r \ell_{t'}(r_1^{t'}; z) - R_\theta^2(t, t')\nabla_r \ell_{t'}(r_2^{t'}; z)\|_2^2 \\ &\leq 2 \left( \|R_\theta^1(t, t') - R_\theta^2(t, t')\|_2^2 \mathbb{E}\|\nabla_r \ell_{t'}(r_1^{t'}; z)\|_2^2 + \|R_\theta^2(t, t')\|_2^2 \mathbb{E}\|\nabla_r \ell_{t'}(r_1^{t'}; z) - \nabla_r \ell_{t'}(r_2^{t'}; z)\|_2^2 \right) \\ &\leq 2 \left( M^2 e^{2\lambda t} \cdot \text{dist}_\lambda(R_\theta^1, R_\theta^2)^2 + \Phi^2 M^2 \mathbb{E}\|r_1^{t'} - r_2^{t'}\|_2^2 \right) \leq K e^{2\lambda t} \cdot \text{dist}_\lambda(X^1, X^2)^2, \end{aligned} \quad (134)$$

$$\begin{aligned} &\mathbb{E}\|R_\theta^1(t, s)\rho_{\ell,1}^{s,t'} - R_\theta^2(t, s)\rho_{\ell,2}^{s,t'}\|_2^2 \\ &\leq 2 \left( \|R_\theta^1(t, s) - R_\theta^2(t, s)\|_2^2 \mathbb{E}\|\rho_{\ell,1}^{s,t'}\|_2^2 + \|R_\theta^2(t, s)\|_2^2 \mathbb{E}\|\rho_{\ell,1}^{s,t'} - \rho_{\ell,2}^{s,t'}\|_2^2 \right) \\ &\leq 2 \left( K e^{2\lambda t} \cdot \text{dist}_\lambda(R_\theta^1, R_\theta^2)^2 + \Phi^2 \mathbb{E}\|\rho_{\ell,1}^{s,t'} - \rho_{\ell,2}^{s,t'}\|_2^2 \right). \end{aligned} \quad (135)$$

Here, we used that  $\mathbb{E}\|\rho_{\ell,i}^{s,t'}\|_2^2$  is uniformly bounded by some constant  $K > 0$  by a similar argument as the bound on  $\mathbb{E}\|\rho_r^{t,t'}\|_2^2$ . Then, we have

$$\mathbb{E}\|\rho_{r,1}^{t,t'} - \rho_{r,2}^{t,t'}\|_2^2 \leq K e^{2\lambda t} \cdot \text{dist}_\lambda(X^1, X^2)^2 + K \int_{t'}^t \mathbb{E}\|\rho_{\ell,1}^{s,t'} - \rho_{\ell,2}^{s,t'}\|_2^2 ds. \quad (136)$$

Plugging Equations (130), (134) and (137) into Equation (128), we have

$$\mathbb{E}\|\rho_{\ell,1}^{t,t'} - \rho_{\ell,2}^{t,t'}\|_2^2 \leq Ke^{2\lambda t} \cdot \text{dist}_\lambda(X^1, X^2)^2 + K \int_{t'}^t \mathbb{E}\|\rho_{\ell,1}^{s,t'} - \rho_{\ell,2}^{s,t'}\|_2^2 ds. \quad (138)$$

Applying Grönwall's inequality, we have

$$e^{-2\lambda t} \mathbb{E}\|\rho_{\ell,1}^{t,t'} - \rho_{\ell,2}^{t,t'}\|_2^2 \leq K \cdot \text{dist}_\lambda(X^1, X^2)^2. \quad (139)$$

Next, we bound  $\mathbb{E}\|D_{\ell,1}^{t,t'} - D_{\ell,2}^{t,t'}\|_2^2$ . We have

$$\begin{aligned} & \|D_{\ell,1}^{t,t'} - D_{\ell,2}^{t,t'}\|_2 \\ & \leq \frac{1}{\delta} \|\nabla_r \ell_t(r_1^t; z) - \nabla_r \ell_t(r_2^t; z)\|_2 \left\| \int_{t'}^t R_\theta^1(t, s) D_{\ell,1}^{s,t'} (ds + \sqrt{\tau\delta} dB^s) \right\|_2 \\ & \quad + \frac{1}{\delta} \|\nabla_r \ell_t(r_2^t; z)\|_2 \left\| \int_{t'}^t (R_\theta^1(t, s) D_{\ell,1}^{s,t'} - R_\theta^2(t, s) D_{\ell,2}^{s,t'}) (ds + \sqrt{\tau\delta} dB^s) \right\|_2 \\ & \quad + \|\nabla_r^2 \ell_t(r_1^t; z) - \nabla_r^2 \ell_t(r_2^t; z)\|_2 \|D_{r,1}^{t,t'} \rho_{r,1}^{t,t'}\|_2 + \|\nabla_r^2 \ell_t(r_2^t; z)\|_2 \|D_{r,1}^{t,t'} \rho_{r,1}^{t,t'} - D_{r,2}^{t,t'} \rho_{r,2}^{t,t'}\|_2 \\ & \leq \frac{M}{\delta} \|r_1^t - r_2^t\|_2 \left\| \int_{t'}^t R_\theta^1(t, s) D_{\ell,1}^{s,t'} (ds + \sqrt{\tau\delta} dB^s) \right\|_2 \\ & \quad + \frac{M}{\delta} \left\| \int_{t'}^t (R_\theta^1(t, s) D_{\ell,1}^{s,t'} - R_\theta^2(t, s) D_{\ell,2}^{s,t'}) (ds + \sqrt{\tau\delta} dB^s) \right\|_2 \\ & \quad + M \|r_1^t - r_2^t\|_2 \|D_{r,1}^{t,t'}\|_2 \|\rho_{r,1}^{t,t'}\|_2 + M \|D_{r,1}^{t,t'} - D_{r,2}^{t,t'}\|_2 \|\rho_{r,1}^{t,t'}\|_2 + M \|D_{r,2}^{t,t'}\|_2 \|\rho_{r,1}^{t,t'} - \rho_{r,2}^{t,t'}\|_2. \end{aligned} \quad (140)$$

Squaring and taking the expectation, we have

$$\begin{aligned} & \mathbb{E}\|D_{\ell,1}^{t,t'} - D_{\ell,2}^{t,t'}\|_2^2 \\ & \leq K \mathbb{E}\|r_1^t - r_2^t\|_2^2 \left\| \int_{t'}^t R_\theta^1(t, s) D_{\ell,1}^{s,t'} (ds + \sqrt{\tau\delta} dB^s) \right\|_2^2 \\ & \quad + K \mathbb{E} \left\| \int_{t'}^t (R_\theta^1(t, s) D_{\ell,1}^{s,t'} - R_\theta^2(t, s) D_{\ell,2}^{s,t'}) (ds + \sqrt{\tau\delta} dB^s) \right\|_2^2 \\ & \quad + K \mathbb{E}\|r_1^t - r_2^t\|_2^2 \|D_{r,1}^{t,t'}\|_2^2 \|\rho_{r,1}^{t,t'}\|_2^2 + K \mathbb{E}\|D_{r,1}^{t,t'} - D_{r,2}^{t,t'}\|_2^2 \|\rho_{r,1}^{t,t'}\|_2^2 + K \mathbb{E}\|D_{r,2}^{t,t'}\|_2^2 \|\rho_{r,1}^{t,t'} - \rho_{r,2}^{t,t'}\|_2^2 \\ & \leq K \sqrt{\mathbb{E}\|r_1^t - r_2^t\|_2^4} \sqrt{\mathbb{E} \left\| \int_{t'}^t R_\theta^1(t, s) D_{\ell,1}^{s,t'} (ds + \sqrt{\tau\delta} dB^s) \right\|_2^4} \\ & \quad + K \int_{t'}^t \mathbb{E} \|R_\theta^1(t, s) D_{\ell,1}^{s,t'} - R_\theta^2(t, s) D_{\ell,2}^{s,t'}\|_2^2 ds \\ & \quad + K \sqrt{\mathbb{E}\|r_1^t - r_2^t\|_2^4} \sqrt{\mathbb{E}\|D_{r,1}^{t,t'}\|_2^8} \sqrt{\mathbb{E}\|\rho_{r,1}^{t,t'}\|_2^8} + K \sqrt{\mathbb{E}\|D_{r,1}^{t,t'} - D_{r,2}^{t,t'}\|_2^4} \sqrt{\mathbb{E}\|\rho_{r,1}^{t,t'}\|_2^4} \\ & \quad + K \sqrt{\mathbb{E}\|D_{r,2}^{t,t'}\|_2^4} \sqrt{\mathbb{E}\|\rho_{r,1}^{t,t'} - \rho_{r,2}^{t,t'}\|_2^4} \\ & \leq Ke^{2\lambda t} \cdot \text{dist}_\lambda(X^1, X^2)^2 \cdot \left( \sqrt{\sup_{s \in [t', t]} \mathbb{E}\|D_{\ell,1}^{s,t'}\|_2^2} + \sqrt{\mathbb{E}\|D_{r,1}^{t,t'}\|_2^8} \sqrt{\mathbb{E}\|\rho_{r,1}^{t,t'}\|_2^8} \right) \end{aligned}$$

$$\begin{aligned}
 & + K \int_{t'}^t \left( \|R_\theta^1(t, s) - R_\theta^2(t, s)\|_2^2 \mathbb{E} \|D_{\ell,1}^{s,t'}\|_2^2 + \|R_\theta^2(t, s)\|_2^2 \mathbb{E} \|D_{\ell,1}^{s,t'} - D_{\ell,2}^{s,t'}\|_2^2 \right) ds \\
 & + K \sqrt{\mathbb{E} \|D_{r,1}^{t,t'} - D_{r,2}^{t,t'}\|_2^4} \sqrt{\mathbb{E} \|\rho_{r,1}^{t,t'}\|_2^4} + K \sqrt{\mathbb{E} \|D_{r,2}^{t,t'}\|_2^4} \sqrt{\mathbb{E} \|\rho_{r,1}^{t,t'} - \rho_{r,2}^{t,t'}\|_2^4}. \tag{141}
 \end{aligned}$$

We do not repeat all the details, but by a similar argument as the bounds on  $\mathbb{E} \|\rho_{r,1}^{t,t'}\|_2^4$ , we can show that  $\mathbb{E} \|\rho_{r,i}^{t,t'}\|_2^p$ ,  $\mathbb{E} \|D_{\ell,i}^{t,t'}\|_2^p$ ,  $\mathbb{E} \|D_{r,i}^{t,t'}\|_2^p$ , and  $\mathbb{E} \|D_{\ell,i}^{t,t'}\|_2^4$  are uniformly bounded by some constant  $K > 0$  for  $p \geq 2$  and  $i = 1, 2$ . Moreover, by a similar argument as the bound on  $\mathbb{E} \|\rho_{r,1}^{t,t'} - \rho_{r,2}^{t,t'}\|_2^2$ , we can show that  $\mathbb{E} \|D_{r,1}^{t,t'} - D_{r,2}^{t,t'}\|_2^4$  and  $\mathbb{E} \|\rho_{r,1}^{t,t'} - \rho_{r,2}^{t,t'}\|_2^4$  are bounded by  $Ke^{4\lambda t} \cdot \text{dist}_\lambda(X^1, X^2)^4$ . Thus, using these bounds and Grönwall's inequality, we have

$$e^{-2\lambda t} \mathbb{E} \|D_{\ell,1}^{t,t'} - D_{\ell,2}^{t,t'}\|_2^2 \leq K \cdot \text{dist}_\lambda(X^1, X^2)^2. \tag{142}$$

Therefore, we have

$$\text{dist}_\lambda(R_\ell^1, R_\ell^2) \leq K \cdot \text{dist}_\lambda(X^1, X^2). \tag{143}$$

Collecting the above bounds, we have

$$\text{dist}_\lambda(Y^1, Y^2) = \sqrt{\lambda} \text{dist}_\lambda(\Sigma_\ell^1, \Sigma_\ell^2) + \text{dist}_\lambda(R_\ell^1, R_\ell^2) + \text{dist}_\lambda(\Gamma^1, \Gamma^2) \leq K \cdot \text{dist}_\lambda(X^1, X^2). \tag{144}$$

### C.5. Proof of Lemma C.7

Again,  $K$  denotes a positive constant that may depend on  $M, m, T, \delta, \tau$ , and  $\Phi$ , but not on other variables, and may change from line to line.

**Bound of  $\text{dist}_\lambda(C_\theta^1, C_\theta^2)$ .** Let  $U_1 \sim \text{GP}(0, \Sigma_\ell^1/\delta)$  and  $U_2 \sim \text{GP}(0, \Sigma_\ell^2/\delta)$  be Gaussian processes coupled such that

$$\sup_{t \in [0, T]} e^{-\lambda t} \sqrt{\mathbb{E} \|U_1^t - U_2^t\|_2^2} \leq 2 \cdot \text{dist}_\lambda(\Sigma_\ell^1, \Sigma_\ell^2). \tag{145}$$

For  $i = 1, 2$ , let  $\theta_i$  be the solution of

$$\theta_i^t = \theta^0 + U_i^t - \int_0^t \left( h_s(\theta_i^s) + \Gamma^i(s) \theta_i^s + \int_0^s R_\ell^i(s, s') \theta_i^{s'} ds' \right) ds. \tag{146}$$

Then, we have

$$\begin{aligned}
 & e^{-2\lambda t} \mathbb{E} \|\theta_1^t - \theta_2^t\|_2^2 \\
 & \leq 6e^{-2\lambda t} \left( T \int_0^t \mathbb{E} \|h_s(\theta_1^s) - h_s(\theta_2^s)\|_2^2 ds + T \int_0^t \|\Gamma^1(s)\|_2^2 \mathbb{E} \|\theta_1^s - \theta_2^s\|_2^2 ds \right. \\
 & \quad + T^2 \int_0^t \int_0^s \|R_\ell^1(s, s')\|_2^2 \mathbb{E} \|\theta_1^{s'} - \theta_2^{s'}\|_2^2 ds' ds + \mathbb{E} \|U_1^t - U_2^t\|_2^2 \\
 & \quad \left. + T \int_0^t \|\Gamma^1(s) - \Gamma^2(s)\|_2^2 \mathbb{E} \|\theta_2^s\|_2^2 ds + T^2 \int_0^t \int_0^s \|R_\ell^1(s, s') - R_\ell^2(s, s')\|_2^2 \mathbb{E} \|\theta_2^{s'}\|_2^2 ds' ds \right)
 \end{aligned}$$

$$\begin{aligned}
 &\leq 6 \left( 2M^2T \int_0^t e^{-2\lambda(t-s)} \cdot e^{-2\lambda s} \mathbb{E} \|\theta_1^s - \theta_2^s\|_2^2 ds \right. \\
 &\quad + T^2 \Phi^2 \int_0^t e^{-2\lambda(t-s)} \int_0^s e^{-2\lambda s} \mathbb{E} \|\theta_1^{s'} - \theta_2^{s'}\|_2^2 ds' ds + 4 \cdot \text{dist}_\lambda(\Sigma_\ell^1, \Sigma_\ell^2)^2 \\
 &\quad + Tm\Phi \int_0^t e^{-2\lambda(t-s)} \cdot e^{-2\lambda s} \|\Gamma^1(s) - \Gamma^2(s)\|_2^2 \\
 &\quad \left. + T^2m\Phi \int_0^t e^{-2\lambda(t-s)} \int_0^s e^{-2\lambda s} \|R_\ell^1(s, s') - R_\ell^2(s, s')\|_2^2 ds' ds \right) \\
 &\leq 6 \left( \frac{M^2T}{\lambda} \sup_{s \in [0, t]} e^{-2\lambda s} \mathbb{E} \|\theta_1^s - \theta_2^s\|_2^2 + \frac{T^3\Phi^2}{2\lambda} \sup_{s \in [0, t]} e^{-2\lambda s} \mathbb{E} \|\theta_1^s - \theta_2^s\|_2^2 \right. \\
 &\quad \left. + \frac{4}{\lambda} \cdot \lambda \text{dist}_\lambda(\Sigma_\ell^1, \Sigma_\ell^2)^2 + \frac{Tm\Phi}{2\lambda} \cdot \text{dist}_\lambda(\Gamma^1, \Gamma^2)^2 + \frac{T^3m\Phi}{2\lambda} \cdot \text{dist}_\lambda(R_\ell^1, R_\ell^2)^2 \right) \\
 &\leq \frac{K}{\lambda} \sup_{s \in [0, t]} e^{-2\lambda s} \mathbb{E} \|\theta_1^s - \theta_2^s\|_2^2 + \frac{K}{\lambda} \cdot \text{dist}_\lambda(Y^1, Y^2)^2. \tag{147}
 \end{aligned}$$

Since  $K$  is independent of  $\lambda$ , we can take  $\lambda > 2K$  so that

$$\sup_{t \in [0, T]} e^{-2\lambda t} \mathbb{E} \|\theta_1^t - \theta_2^t\|_2^2 \leq \frac{2K}{\lambda} \cdot \text{dist}_\lambda(Y^1, Y^2)^2. \tag{148}$$

Thus, for any  $\varepsilon > 0$ , we can take  $\lambda$  large enough such that

$$\sup_{t \in [0, T]} e^{-2\lambda t} \mathbb{E} \|\theta_1^t - \theta_2^t\|_2^2 \leq \varepsilon^2 \cdot \text{dist}_\lambda(Y^1, Y^2)^2. \tag{149}$$

Let  $\{(w_1^t, w_2^t)\}_{t \in [0, T]}$  be a centered Gaussian process with covariance  $\mathbb{E} \begin{bmatrix} \left( \begin{smallmatrix} \theta_1^t \\ \theta_2^t \end{smallmatrix} \right) \left( \begin{smallmatrix} \theta_1^t \\ \theta_2^t \end{smallmatrix} \right)^\top \end{bmatrix}$ . Since  $w_1$  and  $w_2$  have covariance kernels  $C_\theta^1$  and  $C_\theta^2$  respectively, we have

$$\text{dist}_\lambda(C_\theta^1, C_\theta^2) \leq \sup_{t \in [0, T]} e^{-\lambda t} \sqrt{\mathbb{E} \|w_1^t - w_2^t\|_2^2} = \sup_{t \in [0, T]} e^{-\lambda t} \sqrt{\mathbb{E} \|\theta_1^t - \theta_2^t\|_2^2} \leq \varepsilon \cdot \text{dist}_\lambda(Y^1, Y^2). \tag{150}$$

**Bound of  $\text{dist}_\lambda(R_\theta^1, R_\theta^2)$ .** For  $i = 1, 2$ , let  $\rho_{\theta, i}$  be the solution of

$$\rho_{\theta, i}^{t, t'} = R_\theta^i(t, t') - \int_{t'}^t (\nabla_\theta h_s(\theta_i^s) + \Gamma^i(s)) \rho_{\theta, i}^{s, t'} + \int_{t'}^s R_\ell^i(s, s') \rho_{\theta, i}^{s', t'} ds'. \tag{151}$$

Thus, we have

$$\begin{aligned}
 &e^{-\lambda t} \|\rho_{\theta, 1}^{t, t'} - \rho_{\theta, 2}^{t, t'}\|_2 \\
 &\leq e^{-\lambda t} \int_0^t (\|\nabla_\theta h_s(\theta_1^s)\|_2 + \|\Gamma^1(s)\|_2) \|\rho_{\theta, 1}^{s, t'} - \rho_{\theta, 2}^{s, t'}\|_2 ds
 \end{aligned}$$

$$\begin{aligned}
 & + e^{-\lambda t} \int_0^t \left( \|\nabla_{\theta} h_s(\theta_1^s) - \nabla_{\theta} h_s(\theta_2^s)\|_2 + \|\Gamma^1(s) - \Gamma^2(s)\|_2 \right) \|\rho_{\theta,2}^{s,t'}\|_2 ds \\
 & + e^{-\lambda t} \int_0^t \int_{t'}^s \left( \|R_{\ell}^1(s, s')\|_2 \|\rho_{\theta,1}^{s',t'} - \rho_{\theta,2}^{s',t'}\|_2 + \|R_{\ell}^1(s, s') - R_{\ell}^2(s, s')\|_2 \|\rho_{\theta,2}^{s',t'}\|_2 \right) ds' ds \\
 \leq & \int_0^t e^{-\lambda(t-s)} \left( 2M e^{-\lambda s} \|\rho_{\theta,1}^{s,t'} - \rho_{\theta,2}^{s,t'}\|_2 + \Phi \int_{t'}^s e^{-\lambda s'} \|\rho_{\theta,1}^{s',t'} - \rho_{\theta,2}^{s',t'}\|_2 ds' \right) ds \\
 & + \int_0^t e^{-\lambda(t-s)} \left( M e^{-\lambda s} \|\theta_1^s - \theta_2^s\|_2 \|\rho_{\theta,2}^{s,t'}\|_2 + e^{-\lambda s} \|\Gamma^1(s) - \Gamma^2(s)\|_2 \|\rho_{\theta,2}^{s,t'}\|_2 \right) ds \\
 & + \int_0^t e^{-\lambda(t-s)} \int_{t'}^s e^{-\lambda s} \|R_{\ell}^1(s, s') - R_{\ell}^2(s, s')\|_2 \|\rho_{\theta,2}^{s',t'}\|_2 ds' ds. \tag{152}
 \end{aligned}$$

By Equation (110), we have  $\|\rho_{\theta,2}^{s',t'}\|_2 \leq \Phi$ . Thus, we have

$$\begin{aligned}
 & e^{-\lambda t} \mathbb{E} \|\rho_{\theta,1}^{t,t'} - \rho_{\theta,2}^{t,t'}\|_2 \\
 \leq & \int_0^t e^{-\lambda(t-s)} \left( 2M e^{-\lambda s} \|\rho_{\theta,1}^{s,t'} - \rho_{\theta,2}^{s,t'}\|_2 + \Phi \int_{t'}^s e^{-\lambda s'} \|\rho_{\theta,1}^{s',t'} - \rho_{\theta,2}^{s',t'}\|_2 ds' \right) ds \\
 & + \Phi \int_0^t e^{-\lambda(t-s)} \left( M e^{-\lambda s} \mathbb{E} \|\theta_1^s - \theta_2^s\|_2 + e^{-\lambda s} \|\Gamma^1(s) - \Gamma^2(s)\|_2 \right) ds \\
 & + \Phi \int_0^t e^{-\lambda(t-s)} \int_{t'}^s e^{-\lambda s} \|R_{\ell}^1(s, s') - R_{\ell}^2(s, s')\|_2 ds' ds \\
 \leq & \left( \int_0^t e^{-\lambda(t-s)} ds \right) \left( (2M + \Phi T) \sup_{0 \leq t' \leq s \leq t} e^{-\lambda s} \|\rho_{\theta,1}^{s,t'} - \rho_{\theta,2}^{s,t'}\|_2 + \Phi M \sup_{0 \leq s \leq t} \sqrt{\mathbb{E} \|\theta_1^s - \theta_2^s\|_2} \right. \\
 & \left. + \Phi \sup_{0 \leq s \leq t} e^{-\lambda s} \|\Gamma^1(s) - \Gamma^2(s)\|_2 + \Phi T \sup_{0 \leq t' \leq s' \leq s \leq t} e^{-\lambda s} \|R_{\ell}^1(s, s') - R_{\ell}^2(s, s')\|_2 \right) \\
 \leq & \frac{K}{\lambda} \sup_{0 \leq t' \leq s \leq t} e^{-\lambda s} \|\rho_{\theta,1}^{s,t'} - \rho_{\theta,2}^{s,t'}\|_2 + \frac{K}{\lambda} \cdot \text{dist}_{\lambda}(Y^1, Y^2). \tag{153}
 \end{aligned}$$

Taking  $\lambda$  large enough such that  $\lambda > 2K$ , we have

$$\sup_{0 \leq s \leq t \leq T} e^{-\lambda t} \|R_{\theta}^1(t, s) - R_{\theta}^2(t, s)\|_2 \leq \sup_{0 \leq s \leq t \leq T} e^{-\lambda t} \mathbb{E} \|\rho_{\theta,1}^{t,s} - \rho_{\theta,2}^{t,s}\|_2 \leq \frac{2K}{\lambda} \cdot \text{dist}_{\lambda}(Y^1, Y^2). \tag{154}$$

Thus, for any  $\varepsilon > 0$ , we can take  $\lambda$  large enough such that

$$\text{dist}_{\lambda}(R_{\theta}^1, R_{\theta}^2) \leq \varepsilon \cdot \text{dist}_{\lambda}(Y^1, Y^2). \tag{155}$$

Collecting the above bounds, we have that for any  $\varepsilon > 0$ , we can take  $\lambda$  large enough such that

$$\text{dist}_{\lambda}(X^1, X^2) = \text{dist}_{\lambda}(C_{\theta}^1, C_{\theta}^2) + \text{dist}_{\lambda}(R_{\theta}^1, R_{\theta}^2) \leq \varepsilon \cdot \text{dist}_{\lambda}(Y^1, Y^2). \tag{156}$$

## Appendix D. Proof of Theorem 3.2

We follow the three-step strategy outlined in the main text, which we repeat here for clarity.

1. We discretize the SGF (1) with time step  $\gamma > 0$  and analyze the discretization error (Lemma D.1).
2. We apply the AMP theory to characterize the asymptotic behavior of the discretized SGF using a low-dimensional state evolution recursion (Lemma D.2).
3. We take the continuous-time limit  $\gamma \rightarrow 0$  and show that the state evolution converges to the unique solution of the DMFT equation  $\mathfrak{S}$  (Lemma D.3).

We consider the discretization of the SGF (1) with time step  $\gamma > 0$ . Let  $[t] := \max\{k\gamma : k\gamma \leq t, k \in \mathbb{N}\}$ . We define  $\boldsymbol{\theta}_\gamma^t$  and  $\mathbf{r}_\gamma^t$  as the solution of the following equations:

$$\frac{d}{dt} \boldsymbol{\theta}_\gamma^t = - \left( h_{[t]}(\boldsymbol{\theta}_\gamma^{[t]}) + \frac{1}{\delta} \mathbf{X}^\top \ell_{[t]}(\mathbf{r}_\gamma^{[t]}; \mathbf{z}) \right) dt + \sqrt{\frac{\tau}{\delta}} \sum_{i=1}^n \mathbf{x}_i \ell_{[t]}(r_{\gamma,i}^{[t]}; z_i)^\top dB_i^t, \quad (157)$$

and  $\mathbf{r}_\gamma^t = \mathbf{X} \boldsymbol{\theta}_\gamma^t$  with initial condition  $\boldsymbol{\theta}_\gamma^0 = \boldsymbol{\theta}^0$ . At discrete time points  $t_1, t_2, \dots$  where we define  $t_k = k\gamma$  for  $k \in \mathbb{N}$ ,  $\boldsymbol{\theta}_\gamma^{t_k}$  satisfies the following recursion:

$$\boldsymbol{\theta}_\gamma^{t_{k+1}} = \boldsymbol{\theta}_\gamma^{t_k} - \gamma \left( h_{t_k}(\boldsymbol{\theta}_\gamma^{t_k}) + \frac{1}{\delta} \mathbf{X}^\top \ell_{t_k}(\mathbf{r}_\gamma^{t_k}; \mathbf{z}) \right) + \sqrt{\frac{\tau}{\delta}} \sum_{i=1}^n \mathbf{x}_i \ell_{t_k}(r_{\gamma,i}^{t_k}; z_i)^\top (B_i^{t_{k+1}} - B_i^{t_k}). \quad (158)$$

First, we control the discretization error between the SGF (1) and its time-discretized version (158). We prove it in Appendix D.1.

**Lemma D.1** *Under the assumptions of Theorem 3.2, for any  $T > 0$ , there exists a constant  $C > 0$  such that we have, almost surely over the randomness of  $\mathbf{X}, \mathbf{z}, \boldsymbol{\theta}^0$ ,*

$$\limsup_{n,d \rightarrow \infty} \frac{1}{d} \mathbb{E}_{\mathbf{B}} \left[ \sup_{0 \leq t \leq T} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t\|_{\mathbb{F}}^2 \right] \leq C\gamma, \quad \limsup_{n,d \rightarrow \infty} \frac{1}{n} \mathbb{E}_{\mathbf{B}} \left[ \sup_{0 \leq t \leq T} \|\mathbf{r}^t - \mathbf{r}_\gamma^t\|_{\mathbb{F}}^2 \right] \leq C\gamma. \quad (159)$$

Furthermore, for any  $L \in \mathbb{N}$  and  $0 \leq t_1, \dots, t_L \leq T$ , we have, almost surely,

$$\lim_{\gamma \rightarrow 0} \limsup_{n,d \rightarrow \infty} \mathbb{E}_{\mathbf{B}} W_2 \left( \hat{\mathbb{P}}(\boldsymbol{\theta}^{t_1}, \dots, \boldsymbol{\theta}^{t_L}), \hat{\mathbb{P}}(\boldsymbol{\theta}_\gamma^{t_1}, \dots, \boldsymbol{\theta}_\gamma^{t_L}) \right)^2 = 0, \quad (160)$$

$$\lim_{\gamma \rightarrow 0} \limsup_{n,d \rightarrow \infty} \mathbb{E}_{\mathbf{B}} W_2 \left( \hat{\mathbb{P}}(\mathbf{r}^{t_1}, \dots, \mathbf{r}^{t_L}, \mathbf{z}), \hat{\mathbb{P}}(\mathbf{r}_\gamma^{t_1}, \dots, \mathbf{r}_\gamma^{t_L}, \mathbf{z}) \right)^2 = 0. \quad (161)$$

Next, we relate the discretized SGF to the discretized DMFT equation  $\mathfrak{S}^\gamma$  defined in Appendix B.2. The following lemma shows that the unique solution of  $\mathfrak{S}^\gamma$  characterizes the asymptotic behavior of the discretized SGF (158). We prove it in Appendix D.2.

**Lemma D.2** *Under the assumptions of Theorem 3.2, for any  $T > 0, L \in \mathbb{N}$  and  $0 \leq t_1 < \dots < t_L \leq T$ , we have*

$$\text{p-lim}_{n,d \rightarrow \infty} W_2 \left( \hat{\mathbb{P}}(\boldsymbol{\theta}_\gamma^{t_1}, \dots, \boldsymbol{\theta}_\gamma^{t_L}), \mathbb{P}(\boldsymbol{\theta}_\gamma^{t_1}, \dots, \boldsymbol{\theta}_\gamma^{t_L}) \right) = 0, \quad (162)$$

$$\text{p-lim}_{n,d \rightarrow \infty} W_2 \left( \hat{\mathbb{P}}(\mathbf{r}_\gamma^{t_1}, \dots, \mathbf{r}_\gamma^{t_L}, \mathbf{z}), \mathbb{P}(\mathbf{r}_\gamma^{t_1}, \dots, \mathbf{r}_\gamma^{t_L}, \mathbf{z}) \right) = 0. \quad (163)$$

Finally, we establish the convergence of the discretized DMFT equation  $\mathfrak{S}^\gamma$  to the original DMFT equation  $\mathfrak{S}$  as  $\gamma \rightarrow 0$ . We prove it in Appendix D.3.

**Lemma D.3** *Under the assumptions of Theorem 3.2, for any  $T > 0, L \in \mathbb{N}$  and  $t_1 < \dots < t_L \in [0, T]$ , we have*

$$\lim_{\gamma \rightarrow 0} W_2(\mathbb{P}(\theta_\gamma^{t_1}, \dots, \theta_\gamma^{t_L}), \mathbb{P}(\theta^{t_1}, \dots, \theta^{t_L})) = 0, \quad (164)$$

$$\lim_{\gamma \rightarrow 0} W_2(\mathbb{P}(r_\gamma^{t_1}, \dots, r_\gamma^{t_L}, z), \mathbb{P}(r^{t_1}, \dots, r^{t_L}, z)) = 0. \quad (165)$$

We are now ready to prove Theorem 3.2.

**Proof** [Proof of Theorem 3.2] We prove for  $\theta^t$ ; the proof for  $r^t$  is similar.

For any  $0 \leq t_1 < \dots < t_L \leq T$  and  $\gamma > 0$ , by the triangle inequality, we have

$$E_d := W_2(\hat{\mathbb{P}}(\theta^{t_1}, \dots, \theta^{t_L}), \mathbb{P}(\theta^{t_1}, \dots, \theta^{t_L})) \leq E_{d,\gamma}^{(1)} + E_{d,\gamma}^{(2)} + E_\gamma^{(3)}, \quad (166)$$

where we defined

$$E_{d,\gamma}^{(1)} := W_2(\hat{\mathbb{P}}(\theta^{t_1}, \dots, \theta^{t_L}), \hat{\mathbb{P}}(\theta_\gamma^{t_1}, \dots, \theta_\gamma^{t_L})), \quad (167)$$

$$E_{d,\gamma}^{(2)} := W_2(\hat{\mathbb{P}}(\theta_\gamma^{t_1}, \dots, \theta_\gamma^{t_L}), \mathbb{P}(\theta_\gamma^{t_1}, \dots, \theta_\gamma^{t_L})), \quad (168)$$

$$E_\gamma^{(3)} := W_2(\mathbb{P}(\theta_\gamma^{t_1}, \dots, \theta_\gamma^{t_L}), \mathbb{P}(\theta^{t_1}, \dots, \theta^{t_L})). \quad (169)$$

By the union bound, we have, for any  $\varepsilon > 0$ ,

$$\mathbb{P}\{E_d \geq \varepsilon\} \leq \mathbb{P}\{E_{d,\gamma}^{(1)} \geq \varepsilon/3\} + \mathbb{P}\{E_{d,\gamma}^{(2)} \geq \varepsilon/3\} + \mathbb{P}\{E_\gamma^{(3)} \geq \varepsilon/3\}. \quad (170)$$

Taking the limit  $n, d \rightarrow \infty$  and applying Lemma D.2, the second term vanishes. Furthermore, by Markov's inequality, we have

$$\limsup_{n,d \rightarrow \infty} \mathbb{P}\{E_d \geq \varepsilon\} \leq \frac{9}{\varepsilon^2} \limsup_{n,d \rightarrow \infty} \mathbb{E}[(E_{d,\gamma}^{(1)})^2] + \mathbb{P}\{E_\gamma^{(3)} \geq \varepsilon/3\}. \quad (171)$$

Since the left-hand side does not depend on  $\gamma$ , we can take the limit  $\gamma \rightarrow 0$  and apply Lemmas D.1 and D.3 to obtain

$$\lim_{n,d \rightarrow \infty} \mathbb{P}\{E_d \geq \varepsilon\} = 0, \quad (172)$$

and thus  $E_d \rightarrow 0$  in probability as  $n, d \rightarrow \infty$ . ■

### D.1. Proof of Lemma D.1

In the following,  $C$  denotes a constant independent of  $n, d, \gamma$ , which may change from line to line.

We utilize the general results in Appendix G. We check that the SDE (1) satisfies Theorem G.1.

**Drift term.** For the drift coefficient  $\mathbf{b}(t, \boldsymbol{\theta}) = h_t(\boldsymbol{\theta}) + \frac{1}{\delta} \mathbf{X}^\top \ell_t(\mathbf{X}\boldsymbol{\theta}; \mathbf{z})$ , we use the Lipschitz continuity of  $h$  and  $\ell$  and that  $\|\mathbf{X}\|_2 \leq C$  and  $\|\ell_t(0; \mathbf{z})\|_F \leq Cn$  almost surely for sufficiently large  $n, d$  by assumption to obtain

$$\begin{aligned} \|\mathbf{b}(t, \boldsymbol{\theta})\|_F^2 &\leq 2\|h_t(\boldsymbol{\theta})\|_F^2 + \frac{2}{\delta^2} \|\mathbf{X}\|_2^2 \|\ell_t(\mathbf{X}\boldsymbol{\theta}; \mathbf{z})\|_F^2 \\ &\leq C(\|h_t(0)\|_F^2 + \|\boldsymbol{\theta}\|_F^2) + C(\|\ell_t(0; \mathbf{z})\|_F^2 + \|\mathbf{X}\boldsymbol{\theta}\|_F^2) \\ &\leq C(d + \|\boldsymbol{\theta}\|_F^2), \end{aligned} \quad (173)$$

$$\begin{aligned} \|\mathbf{b}(t_1, \boldsymbol{\theta}_1) - \mathbf{b}(t_2, \boldsymbol{\theta}_2)\|_F^2 &\leq 2\|h_{t_1}(\boldsymbol{\theta}_1) - h_{t_2}(\boldsymbol{\theta}_2)\|_F^2 + \frac{2}{\delta^2} \|\mathbf{X}\|_2^2 \|\ell_{t_1}(\mathbf{X}\boldsymbol{\theta}_1; \mathbf{z}) - \ell_{t_2}(\mathbf{X}\boldsymbol{\theta}_2; \mathbf{z})\|_F^2 \\ &\leq C(d|t_1 - t_2|^2 + \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_F^2) + C(n|t_1 - t_2|^2 + \|\mathbf{X}(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2)\|_F^2) \\ &\leq C(d|t_1 - t_2|^2 + \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_F^2). \end{aligned} \quad (174)$$

**Diffusion term.** For the diffusion coefficient  $\boldsymbol{\sigma}_i(t, \boldsymbol{\theta}) = \sqrt{\tau/\delta} \mathbf{x}_i \ell_t(\boldsymbol{\theta}^\top \mathbf{x}_i; z_i)^\top$ , we proceed similarly as above to obtain

$$\begin{aligned} \sum_{i=1}^n \|\boldsymbol{\sigma}_i(t, \boldsymbol{\theta})\|_F^2 &= \frac{\tau}{\delta} \sum_{i=1}^n \|\mathbf{x}_i\|_2^2 \|\ell_t(\boldsymbol{\theta}^\top \mathbf{x}_i; z_i)\|_2^2 \leq C(n + \|\mathbf{X}\boldsymbol{\theta}\|_F^2) \\ &\leq C(d + \|\boldsymbol{\theta}\|_F^2), \end{aligned} \quad (175)$$

$$\begin{aligned} \sum_{i=1}^n \|\boldsymbol{\sigma}_i(t_1, \boldsymbol{\theta}_1) - \boldsymbol{\sigma}_i(t_2, \boldsymbol{\theta}_2)\|_F^2 &= \frac{\tau}{\delta} \sum_{i=1}^n \|\mathbf{x}_i\|_2^2 \|\ell_{t_1}(\boldsymbol{\theta}_1^\top \mathbf{x}_i; z_i) - \ell_{t_2}(\boldsymbol{\theta}_2^\top \mathbf{x}_i; z_i)\|_2^2 \\ &\leq C(d|t_1 - t_2|^2 + \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_F^2). \end{aligned} \quad (176)$$

Therefore, by Lemma G.4, we have that, for sufficiently large  $d$ ,

$$\mathbb{E} \left[ \sup_{0 \leq t \leq T} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t\|_F^2 \right] \leq C\gamma(d + \|\boldsymbol{\theta}^0\|_F^2). \quad (177)$$

By Assumption A.1,  $\|\boldsymbol{\theta}^0\|_F^2 < Cd$  holds almost surely for sufficiently large  $d$ . Therefore, there exists a constant  $C$  independent of  $n, d, \gamma$  such that

$$\limsup_{n, d \rightarrow \infty} \frac{1}{d} \mathbb{E}_{\mathbf{B}} \left[ \sup_{0 \leq t \leq T} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t\|_F^2 \right] \leq C\gamma. \quad (178)$$

This shows the first claim. Furthermore, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{B}} W_2 \left( \hat{\mathbf{P}}(\boldsymbol{\theta}^{t_1}, \dots, \boldsymbol{\theta}^{t_L}), \hat{\mathbf{P}}(\boldsymbol{\theta}_\gamma^{t_1}, \dots, \boldsymbol{\theta}_\gamma^{t_L}) \right)^2 &\leq \mathbb{E}_{\mathbf{B}} \left[ \frac{1}{d} \sum_{l=1}^L \|\boldsymbol{\theta}^{t_l} - \boldsymbol{\theta}_\gamma^{t_l}\|_F^2 \right] \\ &\leq \frac{L}{d} \mathbb{E}_{\mathbf{B}} \left[ \sup_{0 \leq t \leq T} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t\|_F^2 \right]. \end{aligned} \quad (179)$$

Taking the limit  $n, d \rightarrow \infty$  followed by  $\gamma \rightarrow 0$  shows the second claim. The claim for  $\mathbf{r}^t$  follows from

$$\|\mathbf{r}^t - \mathbf{r}_\gamma^t\|_F = \|\mathbf{X}(\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t)\|_F \leq \|\mathbf{X}\|_2 \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t\|_F. \quad (180)$$

## D.2. Proof of Lemma D.2

### D.2.1. REDUCTION TO AMP

For notational simplicity, we omit the subscript  $\gamma$  and denote  $\boldsymbol{\theta}_\gamma^t, \mathbf{r}_\gamma^t$  by  $\boldsymbol{\theta}^t, \mathbf{r}^t$ . As we only work in discrete time, there is no risk of confusion.

Let  $\mathbf{G}^k = (\mathbf{B}^{t_{k+1}} - \mathbf{B}^{t_k})/\sqrt{\gamma} \sim \mathbf{N}(0, \mathbf{I}_n)$ . Then the recursion (158) can be rewritten as

$$\boldsymbol{\theta}^{t_{k+1}} = \boldsymbol{\theta}^{t_k} - \gamma h_{t_k}(\boldsymbol{\theta}^{t_k}) - \frac{\gamma}{\delta} \mathbf{X}^\top ((\mathbf{1}_n + \sqrt{\tau\delta/\gamma} \mathbf{G}^k) \odot \ell_{t_k}(\mathbf{r}^{t_k}; \mathbf{z})), \quad \mathbf{r}^{t_k} = \mathbf{X} \boldsymbol{\theta}^{t_k}. \quad (181)$$

Let  $M > 0$  be a constant and let  $[\cdot]_M : x \mapsto \max\{-M, \min\{x, M\}\}$  be the clipping function. We clip the Gaussian vector in (181) entry-wise as

$$\check{\boldsymbol{\theta}}^{t_{k+1}} = \check{\boldsymbol{\theta}}^{t_k} - \gamma h_{t_k}(\check{\boldsymbol{\theta}}^{t_k}) - \frac{\gamma}{\delta} \mathbf{X}^\top ((\mathbf{1}_n + \sqrt{\tau\delta/\gamma} [\mathbf{G}^k]_M) \odot \ell_{t_k}(\check{\mathbf{r}}^{t_k}; \mathbf{z})), \quad \check{\mathbf{r}}^{t_k} = \mathbf{X} \check{\boldsymbol{\theta}}^{t_k}. \quad (182)$$

We first control the difference between (181) and (182).

**Lemma D.4** *We have, almost surely over the randomness of  $\mathbf{X}, \mathbf{z}, \boldsymbol{\theta}^0$ ,*

$$\lim_{M \rightarrow \infty} \limsup_{n, d \rightarrow \infty} \frac{1}{d} \mathbb{E}_{\mathbf{B}} \left[ \max_{0 \leq k \leq T/\gamma} \|\boldsymbol{\theta}^{t_k} - \check{\boldsymbol{\theta}}^{t_k}\|_{\mathbb{F}}^2 \right] = 0, \quad (183)$$

and thus

$$\lim_{M \rightarrow \infty} \limsup_{n, d \rightarrow \infty} \mathbb{E}_{\mathbf{B}} W_2 \left( \hat{\mathbb{P}}(\boldsymbol{\theta}^{t_1}, \dots, \boldsymbol{\theta}^{t_L}), \hat{\mathbb{P}}(\check{\boldsymbol{\theta}}^{t_1}, \dots, \check{\boldsymbol{\theta}}^{t_L}) \right)^2 = 0, \quad (184)$$

$$\lim_{M \rightarrow \infty} \limsup_{n, d \rightarrow \infty} \mathbb{E}_{\mathbf{B}} W_2 \left( \hat{\mathbb{P}}(\mathbf{r}^{t_1}, \dots, \mathbf{r}^{t_L}, \mathbf{z}), \hat{\mathbb{P}}(\check{\mathbf{r}}^{t_1}, \dots, \check{\mathbf{r}}^{t_L}, \mathbf{z}) \right)^2 = 0, \quad (185)$$

**Proof** We have

$$\begin{aligned} & \|\boldsymbol{\theta}^{t_{k+1}} - \check{\boldsymbol{\theta}}^{t_{k+1}}\|_{\mathbb{F}} \\ & \leq \|\boldsymbol{\theta}^{t_k} - \check{\boldsymbol{\theta}}^{t_k}\|_{\mathbb{F}} + \gamma \|h_{t_k}(\boldsymbol{\theta}^{t_k}) - h_{t_k}(\check{\boldsymbol{\theta}}^{t_k})\|_{\mathbb{F}} \\ & \quad + \frac{\gamma}{\delta} \left\| \mathbf{X}^\top \left( (\mathbf{1}_n + \sqrt{\tau\delta/\gamma} \mathbf{G}^k) \odot \ell_{t_k}(\mathbf{r}^{t_k}; \mathbf{z}) - (\mathbf{1}_n + \sqrt{\tau\delta/\gamma} [\mathbf{G}^k]_M) \odot \ell_{t_k}(\check{\mathbf{r}}^{t_k}; \mathbf{z}) \right) \right\|_{\mathbb{F}} \\ & \leq C \|\boldsymbol{\theta}^{t_k} - \check{\boldsymbol{\theta}}^{t_k}\|_{\mathbb{F}} + C \|(\mathbf{G}^k - [\mathbf{G}^k]_M) \odot \ell_{t_k}(\mathbf{r}^{t_k}; \mathbf{z})\|_{\mathbb{F}} + C \|[\mathbf{G}^k]_M \odot (\ell_{t_k}(\mathbf{r}^{t_k}; \mathbf{z}) - \ell_{t_k}(\check{\mathbf{r}}^{t_k}; \mathbf{z}))\|_{\mathbb{F}}. \end{aligned} \quad (186)$$

Since  $\mathbf{G}^k$  is independent of  $\mathbf{r}^{t_k}$  and  $\check{\mathbf{r}}^{t_k}$ , we have

$$\begin{aligned} & \mathbb{E} \|(\mathbf{G}^k - [\mathbf{G}^k]_M) \odot \ell_{t_k}(\mathbf{r}^{t_k}; \mathbf{z})\|_{\mathbb{F}}^2 = \sum_{i=1}^n \mathbb{E} [(G_i^k - [G_i^k]_M)^2 \|\ell_{t_k}(\mathbf{r}_i^{t_k}; z_i)\|_2^2] \\ & = f(M) \sum_{i=1}^n \mathbb{E} [\|\ell_{t_k}(\mathbf{r}_i^{t_k}; z_i)\|_2^2] = f(M) \mathbb{E} \|\ell_{t_k}(\mathbf{r}^{t_k}; \mathbf{z})\|_{\mathbb{F}}^2 \leq C f(M) \mathbb{E} \|\boldsymbol{\theta}^{t_k}\|_{\mathbb{F}}^2, \end{aligned} \quad (187)$$

$$\mathbb{E} \|[\mathbf{G}^k]_M \odot (\ell_{t_k}(\mathbf{r}^{t_k}; \mathbf{z}) - \ell_{t_k}(\check{\mathbf{r}}^{t_k}; \mathbf{z}))\|_{\mathbb{F}}^2 = \sum_{i=1}^n \mathbb{E} [[G_i^k]_M^2 \|\ell_{t_k}(\mathbf{r}_i^{t_k}; z_i) - \ell_{t_k}(\check{\mathbf{r}}_i^{t_k}; z_i)\|_2^2]$$

$$= \mathbb{E}[[G]_M^2] \sum_{i=1}^n \mathbb{E}[\|\ell_{t_k}(\mathbf{r}_i^{t_k}; z_i) - \ell_{t_k}(\check{\mathbf{r}}_i^{t_k}; z_i)\|_2^2] \leq C \mathbb{E}[G^2] \mathbb{E}\|\mathbf{r}^{t_k} - \check{\mathbf{r}}^{t_k}\|_F^2 \leq C \mathbb{E}\|\boldsymbol{\theta}^{t_k} - \check{\boldsymbol{\theta}}^{t_k}\|_F^2, \quad (188)$$

where  $f(M) = \mathbb{E}[(G - [G]_M)^2]$  for  $G \sim \mathcal{N}(0, 1)$ . Therefore, we have

$$\frac{1}{d} \mathbb{E}\|\boldsymbol{\theta}^{t_{k+1}} - \check{\boldsymbol{\theta}}^{t_{k+1}}\|_F^2 \leq \frac{C}{d} \mathbb{E}\|\boldsymbol{\theta}^{t_k} - \check{\boldsymbol{\theta}}^{t_k}\|_F^2 + C f(M) \cdot \frac{1}{d} \mathbb{E}\|\boldsymbol{\theta}^{t_k}\|_F^2. \quad (189)$$

Iterating this inequality yields

$$\frac{1}{d} \mathbb{E}\|\boldsymbol{\theta}^{t_k} - \check{\boldsymbol{\theta}}^{t_k}\|_F^2 \leq C f(M) \cdot \frac{1}{d} \sup_{k \leq T/\gamma} \mathbb{E}\|\boldsymbol{\theta}^{t_k}\|_F^2. \quad (190)$$

As  $n, d \rightarrow \infty$ ,  $\mathbb{E}\|\boldsymbol{\theta}^{t_k}\|_F^2/d$  is uniformly bounded in  $t_k$  almost surely for large  $d$  by Lemma G.2. As  $M \rightarrow \infty$ , we have  $f(M) \rightarrow 0$  by the dominated convergence theorem, and the first claim follows. Convergence of the 2-Wasserstein distances follows from the first claim by the same argument as in the proof of Lemma D.1.  $\blacksquare$

Consider the following AMP iteration. Given sequences of functions  $f_i: \mathbb{R}^{(i+1)m+i+2} \rightarrow \mathbb{R}^m$  and  $g_i: \mathbb{R}^{im+m} \rightarrow \mathbb{R}^m$  ( $i \geq 0$ ) that are Lipschitz in the first  $(i+1)m$  and  $im$  arguments, respectively, we generate sequences of matrices  $\mathbf{a}^{i+1} \in \mathbb{R}^{d \times m}$  and  $\mathbf{b}^i \in \mathbb{R}^{n \times m}$  ( $i \geq 0$ ) as follows.

$$\mathbf{a}^{i+1} = -\frac{1}{\delta} \mathbf{X}^\top f_i(\mathbf{b}^0, \dots, \mathbf{b}^i; \mathbf{z}, [\mathbf{G}^0]_M, \dots, [\mathbf{G}^i]_M) + \sum_{j=0}^i g_j(\mathbf{a}^1, \dots, \mathbf{a}^j; \boldsymbol{\theta}^0) \xi_{i,j}^\top, \quad (191)$$

$$\mathbf{b}^i = \mathbf{X} g_i(\mathbf{a}^1, \dots, \mathbf{a}^i; \boldsymbol{\theta}^0) + \frac{1}{\delta} \sum_{j=0}^{i-1} f_j(\mathbf{b}^0, \dots, \mathbf{b}^j; \mathbf{z}, [\mathbf{G}^0]_M, \dots, [\mathbf{G}^j]_M) \zeta_{i,j}^\top, \quad (192)$$

with initial values  $g_0(\boldsymbol{\theta}^0) = \boldsymbol{\theta}^0$ ,  $\mathbf{b}^0 = \mathbf{X} \boldsymbol{\theta}^0$ . Here,  $f_i$  and  $g_i$  are applied row-wise, and  $\{\xi_{i,j}\}_{0 \leq j \leq i}$ ,  $\{\zeta_{i,j}\}_{0 \leq j \leq i-1} \subset \mathbb{R}^{m \times m}$  are defined as follows. We define a sequence of centered Gaussian random variables  $\{\bar{u}^{i+1}, \bar{w}^i\}_{i \geq 0}$  recursively as

$$\mathbb{E}[\bar{w}^i (\bar{w}^j)^\top] = \mathbb{E}[g_i(\bar{u}^1, \dots, \bar{u}^i; \boldsymbol{\theta}^0) g_j(\bar{u}^1, \dots, \bar{u}^j; \boldsymbol{\theta}^0)^\top], \quad (193)$$

$$\mathbb{E}[\bar{u}^{i+1} (\bar{u}^{j+1})^\top] = \frac{1}{\delta} \mathbb{E}[f_i(\bar{w}^0, \dots, \bar{w}^i; \mathbf{z}, [\bar{\mathbf{G}}^0]_M, \dots, [\bar{\mathbf{G}}^i]_M) f_j(\bar{w}^0, \dots, \bar{w}^j; \mathbf{z}, [\bar{\mathbf{G}}^0]_M, \dots, [\bar{\mathbf{G}}^j]_M)^\top], \quad (194)$$

for  $0 \leq j \leq i$ , and set  $\zeta_{i,j}, \xi_{i,j}$  as

$$\zeta_{i,j} = \mathbb{E} \left[ \frac{\partial}{\partial \bar{w}^{j+1}} g_i(\bar{u}^1, \dots, \bar{u}^i; \boldsymbol{\theta}^0) \right], \quad 0 \leq j \leq i-1, \quad (195)$$

$$\xi_{i,j} = \mathbb{E} \left[ \frac{\partial}{\partial \bar{w}^j} f_i(\bar{w}^0, \dots, \bar{w}^i; \mathbf{z}, [\bar{\mathbf{G}}^0]_M, \dots, [\bar{\mathbf{G}}^i]_M) \right], \quad 0 \leq j \leq i, \quad (196)$$

where the expectations are taken over  $\bar{u}^i, \bar{w}^i, \boldsymbol{\theta}^0 \sim \mathcal{P}(\boldsymbol{\theta}^0)$ ,  $\mathbf{z} \sim \mathcal{P}(\mathbf{z})$ , and  $\bar{\mathbf{G}}^i \sim \mathcal{N}(0, 1)$ .

This AMP iteration can be mapped to the recursion (181) by considering the specific choice of  $f_i$  and  $g_i$  as follows.

$$g_i(\mathbf{a}^1, \dots, \mathbf{a}^i; \boldsymbol{\theta}^0) = \check{\boldsymbol{\theta}}^{t_i}, \quad (197)$$

$$f_i(\mathbf{b}^0, \dots, \mathbf{b}^i; \mathbf{z}, [\mathbf{G}^0]_M, \dots, [\mathbf{G}^i]_M) = (\mathbf{1}_n + \sqrt{\tau\delta/\gamma}[\mathbf{G}^i]_M) \odot \ell_{t_i}(\mathbf{r}^{t_i}; \mathbf{z}). \quad (198)$$

We show that  $g_i$  is indeed a function of  $\mathbf{a}^1, \dots, \mathbf{a}^i, \boldsymbol{\theta}^0$  and Lipschitz in  $\mathbf{a}^j$ , and that  $f_i$  is a function of  $\mathbf{b}^0, \dots, \mathbf{b}^i, \mathbf{z}, [\mathbf{G}^0]_M, \dots, [\mathbf{G}^i]_M$  and Lipschitz in  $\mathbf{b}^j$ . It can be shown by the Lipschitz continuity of  $\ell$ , boundedness of  $[G]_M$ , and induction over  $i$  as follows.

$$\begin{aligned} \check{\boldsymbol{\theta}}^{t_{i+1}} &= \check{\boldsymbol{\theta}}^{t_i} - \gamma h_{t_i}(\check{\boldsymbol{\theta}}^{t_i}) - \frac{\gamma}{\delta} \mathbf{X}^\top f_i(\mathbf{b}^0, \dots, \mathbf{b}^i; \mathbf{z}, [\mathbf{G}^0]_M, \dots, [\mathbf{G}^i]_M) \\ &= \check{\boldsymbol{\theta}}^{t_i} - \gamma h_{t_i}(\check{\boldsymbol{\theta}}^{t_i}) + \gamma \left( \mathbf{a}^{i+1} - \sum_{j=0}^i g_j(\mathbf{a}^1, \dots, \mathbf{a}^j; \boldsymbol{\theta}^0) \xi_{i,j}^\top \right), \end{aligned} \quad (199)$$

$$\begin{aligned} \mathbf{r}^{t_i} &= \mathbf{X} \boldsymbol{\theta}^{t_i} = \mathbf{X} g_i(\mathbf{a}^1, \dots, \mathbf{a}^i; \boldsymbol{\theta}^0) \\ &= \mathbf{b}^i - \frac{1}{\delta} \sum_{j=0}^{i-1} f_j(\mathbf{b}^0, \dots, \mathbf{b}^j; \mathbf{z}, [\mathbf{G}^0]_M, \dots, [\mathbf{G}^j]_M) \zeta_{i,j}^\top. \end{aligned} \quad (200)$$

By Wang et al. [56, Theorem 2.21], for any second order pseudo-Lipschitz functions  $\psi: \mathbb{R}^{im+m} \rightarrow \mathbb{R}$  and  $\tilde{\psi}: \mathbb{R}^{(i+1)m+i+2} \rightarrow \mathbb{R}$ , we have almost surely

$$\lim_{n,d \rightarrow \infty} \frac{1}{d} \sum_{j=1}^d \psi(a_j^1, \dots, a_j^i; \theta_j^0) = \mathbb{E}[\psi(\bar{u}^1, \dots, \bar{u}^i; \theta^0)], \quad (201)$$

$$\lim_{n,d \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n \tilde{\psi}(b_j^0, \dots, b_j^i; z_j, G_j^0, \dots, G_j^i) = \mathbb{E}[\tilde{\psi}(\bar{w}^0, \dots, \bar{w}^i; z, \bar{G}^0, \dots, \bar{G}^i)]. \quad (202)$$

Since  $\check{\boldsymbol{\theta}}^{t_i}$  is a Lipschitz function of  $\mathbf{a}^1, \dots, \mathbf{a}^i, \boldsymbol{\theta}^0$ , we can take a Lipschitz function  $h_\theta$  such that  $\check{\boldsymbol{\theta}}^{t_i} = h_\theta(\mathbf{a}^1, \dots, \mathbf{a}^i; \boldsymbol{\theta}^0)$  and define  $\bar{\boldsymbol{\theta}}^i := h_\theta(\bar{u}^1, \dots, \bar{u}^i; \theta^0)$ . Similarly, we can take a Lipschitz function  $h_r$  such that  $\mathbf{r}^{t_i} = h_r(\mathbf{b}^0, \dots, \mathbf{b}^i; \mathbf{z}, [\mathbf{G}^0]_M, \dots, [\mathbf{G}^i]_M)$  and define  $\bar{r}^i := h_r(\bar{w}^0, \dots, \bar{w}^i; z, [\bar{G}^0]_M, \dots, [\bar{G}^i]_M)$ . Considering the composition of  $\psi, \tilde{\psi}$  with  $h_\theta, h_r$ , we have almost surely

$$\lim_{n,d \rightarrow \infty} \frac{1}{d} \sum_{j=1}^d \psi(\check{\boldsymbol{\theta}}_j^{t_1}, \dots, \check{\boldsymbol{\theta}}_j^{t_i}; \theta_j^0) = \mathbb{E}[\psi(\bar{\boldsymbol{\theta}}^1, \dots, \bar{\boldsymbol{\theta}}^i; \theta^0)], \quad (203)$$

$$\lim_{n,d \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n \tilde{\psi}(r_j^0, \dots, r_j^{t_i}; z_j, [G_j^0]_M, \dots, [G_j^i]_M) = \mathbb{E}[\tilde{\psi}(\bar{r}^0, \dots, \bar{r}^i; z, [\bar{G}^0]_M, \dots, [\bar{G}^i]_M)]. \quad (204)$$

As  $M \rightarrow \infty$ , by the dominated convergence theorem, we have  $W_2([G]_M, G) \rightarrow 0$  for  $G \sim \mathbf{N}(0, 1)$ . Combining this with the above and Lemma D.4, we have

$$\text{p-lim}_{n,d \rightarrow \infty} W_2(\hat{\mathbf{P}}(\boldsymbol{\theta}^{t_0}, \dots, \boldsymbol{\theta}^{t_i}), \mathbf{P}(\boldsymbol{\theta}^0, \bar{\boldsymbol{\theta}}^1, \dots, \bar{\boldsymbol{\theta}}^i)) = 0, \quad (205)$$

$$\text{p-lim}_{n,d \rightarrow \infty} W_2(\hat{\mathbf{P}}(\mathbf{r}^{t_0}, \dots, \mathbf{r}^{t_i}, \mathbf{z}), \mathbf{P}(\bar{r}^0, \dots, \bar{r}^i; z, \bar{G}^0, \dots, \bar{G}^i)) = 0. \quad (206)$$

## D.2.2. MAPPING THE STATE EVOLUTION TO DMFT

It remains to show that the state evolution process  $(\bar{\theta}^i, \bar{r}^i)_{i \geq 0}$  defined above satisfies the discretized DMFT equations (30)–(32).

By Equations (193) and (194), we have

$$\begin{aligned} \mathbb{E}[\bar{w}^i(\bar{w}^j)^\top] &= \mathbb{E}[\bar{\theta}^i(\bar{\theta}^j)^\top], \\ \mathbb{E}[\bar{u}^{i+1}(\bar{u}^{j+1})^\top] &= \frac{1}{\delta} \mathbb{E}[(1 + \sqrt{\tau\delta/\gamma G^i})\ell_{t_i}(\bar{r}^i; z)(1 + \sqrt{\tau\delta/\gamma G^j})\ell_{t_j}(\bar{r}^j; z)^\top]. \end{aligned} \quad (207)$$

Define  $\bar{U}^i$  and  $\bar{L}^i$  as

$$\bar{U}^i := \gamma \sum_{j=1}^i \bar{u}^j, \quad \bar{L}^i := \gamma \sum_{j=0}^{i-1} (1 + \sqrt{\tau\delta/\gamma G^j})\ell_{t_j}(\bar{r}^j; z). \quad (208)$$

Then, we have

$$\mathbb{E}[\bar{U}^i(\bar{U}^j)^\top] = \frac{1}{\delta} \mathbb{E}[\bar{L}^i(\bar{L}^j)^\top]. \quad (209)$$

By Equation (199),  $\bar{\theta}^i$  follows the following recursion.

$$\begin{aligned} \bar{\theta}^{i+1} &= \bar{\theta}^i - \gamma h_{t_i}(\bar{\theta}^i) + \gamma \left( \bar{u}^{i+1} - \sum_{j=0}^i \xi_{i,j} \bar{\theta}^j \right) \\ &= \bar{\theta}^i + \gamma \left( \bar{u}^{i+1} - h_{t_i}(\bar{\theta}^i) - \xi_{i,i} \bar{\theta}^i - \sum_{j=0}^{i-1} \xi_{i,j} \bar{\theta}^j \right). \end{aligned} \quad (210)$$

Thus,

$$\begin{aligned} \bar{\theta}^i &= \bar{\theta}^0 + \sum_{j=0}^{i-1} (\bar{\theta}^{j+1} - \bar{\theta}^j) = \bar{\theta}^0 + \gamma \sum_{j=0}^{i-1} \bar{u}^{j+1} - \gamma \sum_{j=0}^{i-1} \left( h_{t_j}(\bar{\theta}^j) + \xi_{j,j} \bar{\theta}^j + \sum_{k=0}^{j-1} \xi_{j,k} \bar{\theta}^k \right) \\ &= \bar{\theta}^0 + \bar{U}^i - \gamma \sum_{j=0}^{i-1} \left( h_{t_j}(\bar{\theta}^j) + \xi_{j,j} \bar{\theta}^j + \sum_{k=0}^{j-1} \xi_{j,k} \bar{\theta}^k \right). \end{aligned} \quad (211)$$

Furthermore,  $\partial \bar{\theta}^i / \partial \bar{u}^{j+1}$  satisfies

$$\frac{\partial \bar{\theta}^i}{\partial \bar{u}^{j+1}} = \gamma I_m - \gamma \sum_{k=j+1}^{i-1} \left( (\nabla_{\theta} h_{t_k}(\bar{\theta}^k) + \xi_{k,k}) \frac{\partial \bar{\theta}^k}{\partial \bar{u}^{j+1}} + \sum_{l=j+1}^{k-1} \xi_{k,l} \frac{\partial \bar{\theta}^l}{\partial \bar{u}^{j+1}} \right). \quad (212)$$

By Equation (200),  $\bar{r}^i$  follows the following recursions.

$$\bar{r}^i = \bar{w}^i - \frac{1}{\delta} \sum_{j=0}^{i-1} \zeta_{i,j} \ell_{t_j}(\bar{r}^j; z) (1 + \sqrt{\tau\delta/\gamma G^j}). \quad (213)$$

Furthermore,  $\partial \ell_{t_i}(\bar{r}^i; z) / \partial \bar{w}^j$  satisfies

$$\frac{\partial \ell_{t_i}(\bar{r}^i; z)}{\partial \bar{w}^j} = \nabla_r \ell_{t_i}(\bar{r}^i; z) \frac{\partial \bar{r}^i}{\partial \bar{w}^j}, \quad (214)$$

$$\frac{\partial \bar{r}^i}{\partial \bar{w}^j} = -\frac{1}{\delta} \sum_{k=j+1}^{i-1} \zeta_{i,k} \frac{\partial \ell_{t_k}(\bar{r}^k; z)}{\partial \bar{w}^j} (1 + \sqrt{\tau \delta / \gamma \bar{G}^k}) - \frac{1}{\delta} \zeta_{i,j} \nabla_r \ell_{t_j}(\bar{r}^j; z) (1 + \sqrt{\tau \delta / \gamma \bar{G}^j}). \quad (215)$$

Let  $\{\bar{\rho}_\ell^{i,j}\}_{0 \leq j < i}$  be the stochastic process satisfying

$$\bar{\rho}_\ell^{i,j} = \nabla_r \ell_{t_i}(\bar{r}^i; z) \bar{\rho}_r^{i,j}, \quad \bar{\rho}_r^{i,j} = -\frac{1}{\delta} \sum_{k=j+1}^{i-1} \zeta_{i,k} \bar{\rho}_\ell^{k,j} (1 + \sqrt{\tau \delta / \gamma \bar{G}^k}) - \frac{1}{\delta} \zeta_{i,j} \nabla_r \ell_{t_j}(\bar{r}^j; z). \quad (216)$$

Using the linearity of Equation (214), we have

$$\frac{\partial \ell_{t_i}(\bar{r}^i; z)}{\partial \bar{w}^j} = \bar{\rho}_\ell^{i,j} (1 + \sqrt{\tau \delta / \gamma \bar{G}^j}). \quad (217)$$

Then,  $\xi_{i,j}$  satisfies

$$\xi_{i,j} = \mathbb{E} \left[ \frac{\partial \ell_{t_i}(\bar{r}^i; z)}{\partial \bar{w}^j} \right] = \mathbb{E} \left[ \bar{\rho}_\ell^{i,j} (1 + \sqrt{\tau \delta / \gamma \bar{G}^j}) \right] = \mathbb{E}[\bar{\rho}_\ell^{i,j}] + \sqrt{\frac{\tau \delta}{\gamma}} \mathbb{E} \left[ \frac{\partial \bar{\rho}_\ell^{i,j}}{\partial \bar{G}^j} \right], \quad (218)$$

where we used Stein's lemma (Gaussian integration by parts) in the last equality. By Equation (216),  $\partial \bar{\rho}_\ell^{i,j} / \partial \bar{G}^j$  satisfies

$$\begin{aligned} \frac{\partial \bar{\rho}_\ell^{i,j}}{\partial \bar{G}^j} &= \nabla_r \ell_{t_i}(\bar{r}^i; z) \frac{\partial \bar{\rho}_r^{i,j}}{\partial \bar{G}^j} + \frac{\partial \nabla_r \ell_{t_i}(\bar{r}^i; z)}{\partial \bar{G}^j} \bar{\rho}_r^{i,j} \\ &= \nabla_r \ell_{t_i}(\bar{r}^i; z) \left( -\frac{1}{\delta} \sum_{k=j+1}^{i-1} \zeta_{i,k} \frac{\partial \bar{\rho}_\ell^{k,j}}{\partial \bar{G}^j} (1 + \sqrt{\tau \delta / \gamma \bar{G}^k}) \right) + \nabla_r^2 \ell_{t_i}(\bar{r}^i; z) \left[ \frac{\partial \bar{r}^i}{\partial \bar{G}^j} \right] \bar{\rho}_r^{i,j}, \end{aligned} \quad (219)$$

where  $\partial \bar{r}^i / \partial \bar{G}^j$  satisfies

$$\frac{\partial \bar{r}^i}{\partial \bar{G}^j} = -\frac{1}{\delta} \sum_{k=j+1}^{i-1} \zeta_{i,k} \nabla_r \ell_{t_k}(\bar{r}^k; z) \frac{\partial \bar{r}^k}{\partial \bar{G}^j} (1 + \sqrt{\tau \delta / \gamma \bar{G}^k}) - \sqrt{\frac{\tau}{\delta \gamma}} \zeta_{i,j} \ell_{t_j}(\bar{r}^j; z). \quad (220)$$

These state evolution recursions exactly correspond to the discrete DMFT equations  $\mathfrak{S}^\gamma$  shown in Equations (30) and (32) by the following mappings:

$$\theta_\gamma^{t_i} \stackrel{d}{=} \bar{\theta}^i, \quad i \geq 0, \quad (\text{Compare (30a) with (211)})$$

$$\rho_{\theta, \gamma}^{t_i, t_j} \stackrel{d}{=} \frac{1}{\gamma} \frac{\partial \bar{\theta}^i}{\partial \bar{w}^{j+1}}, \quad 0 \leq j < i, \quad (\text{Compare (30b) with (212)})$$

$$r_\gamma^{t_i} \stackrel{d}{=} \bar{r}^i, \quad i \geq 0, \quad (\text{Compare (30c) with (213)})$$

$$\rho_{\ell,\gamma}^{t_i,t_j} \stackrel{\text{d}}{=} \frac{1}{\gamma} \bar{\rho}_{\ell}^{i,j}, \quad 0 \leq j < i, \quad (\text{Compare (30d), (30e) with (216)})$$

$$D_{\ell,\gamma}^{t_i,t_j} \stackrel{\text{d}}{=} \frac{1}{\sqrt{\gamma}} \frac{\partial \bar{\rho}_{\ell}^{i,j}}{\partial \bar{G}^j}, \quad 0 \leq j < i, \quad (\text{Compare (30f), (30g) with (219), (220)})$$

$$U_{\gamma}^{t_i} \stackrel{\text{d}}{=} \bar{U}^i, \quad i \geq 0, \quad (221)$$

$$w_{\gamma}^{t_i} \stackrel{\text{d}}{=} \bar{w}^i, \quad i \geq 0, \quad (222)$$

$$B^{t_{i+1}} - B^{t_i} \stackrel{\text{d}}{=} \sqrt{\gamma} \bar{G}^i, \quad i \geq 0, \quad (223)$$

and

$$C_{\theta}^{\gamma}(t_i, t_j) = \mathbb{E}[\bar{\theta}^i (\bar{\theta}^j)^{\top}], \quad i, j \geq 0, \quad (224)$$

$$\Sigma_{\ell}^{\gamma}(t_i, t_j) = \mathbb{E}[\bar{L}^i (\bar{L}^j)^{\top}], \quad i, j \geq 0, \quad (225)$$

$$R_{\theta}^{\gamma}(t_i, t_j) = \zeta_{i,j}/\gamma, \quad 0 \leq j < i, \quad (226)$$

$$R_{\ell}^{\gamma}(t_i, t_j) = \xi_{i,j}/\gamma, \quad 0 \leq j < i, \quad (227)$$

$$\Gamma^{\gamma}(t_i) = \xi_{i,i}, \quad 0 \leq i. \quad (228)$$

### D.3. Proof of Lemma D.3

We first embed the discretized DMFT equation  $\mathfrak{S}^{\gamma}$  defined for discrete time knots  $t_k = k\gamma$  ( $k \geq 0$ ) into continuous time  $t \in [0, T]$  in a piecewise constant manner. We define the stochastic processes  $\{\theta_{\gamma}^t, r_{\gamma}^t\}_{t \in [0, T]}$  and  $\{\rho_{\theta,\gamma}^{t,t'}, \rho_{\ell,\gamma}^{t,t'}, D_{\ell,\gamma}^{t,t'}\}_{t \geq t' \geq 0}$  by the following equations.

$$\theta_{\gamma}^t = \theta^0 + U_{\gamma}^t - \int_0^{\lfloor t \rfloor} (h_{\lfloor s \rfloor}(\theta_{\gamma}^s) + \Gamma^{\gamma}(s) \theta_{\gamma}^s) + \int_0^{\lfloor s \rfloor} R_{\ell}^{\gamma}(s, s') \theta_{\gamma}^{s'} ds' ds, \quad (229)$$

$$\rho_{\theta,\gamma}^{t,t'} = I_m - \mathbb{I}(\lceil t' \rceil \leq \lfloor t \rfloor) \int_{\lceil t' \rceil}^{\lfloor t \rfloor} \left( (\nabla_{\theta} h_{\lfloor s \rfloor}(\theta_{\gamma}^s) + \Gamma^{\gamma}(s)) \rho_{\theta,\gamma}^{s,t'} + \int_{\lceil t' \rceil}^{\lfloor s \rfloor} R_{\ell}^{\gamma}(s, s') \rho_{\theta,\gamma}^{s',t'} ds' \right) ds, \quad (230)$$

$$r_{\gamma}^t = w_{\gamma}^t - \frac{1}{\delta} \int_0^{\lfloor t \rfloor} R_{\theta}^{\gamma}(t, s) \ell_{\lfloor s \rfloor}(r_{\gamma}^s; z) (ds + \sqrt{\tau \delta} dB^s), \quad (231)$$

$$\rho_{\ell,\gamma}^{t,t'} = \nabla_r \ell_{\lfloor t \rfloor}(r_{\gamma}^t; z) \rho_{r,\gamma}^{t,t'}, \quad (232)$$

$$\rho_{r,\gamma}^{t,t'} = -\frac{1}{\delta} \mathbb{I}(\lceil t' \rceil \leq \lfloor t \rfloor) \int_{\lceil t' \rceil}^{\lfloor t \rfloor} R_{\theta}^{\gamma}(t, s) \rho_{\ell,\gamma}^{s,t'} (ds + \sqrt{\tau \delta} dB^s) - \frac{1}{\delta} R_{\theta}^{\gamma}(t, t') \nabla_r \ell_{\lceil t' \rceil}(r_{\gamma}^{t'}; z), \quad (233)$$

$$D_{\ell,\gamma}^{t,t'} = \nabla_r \ell_{\lfloor t \rfloor}(r_{\gamma}^t; z) \left( -\frac{1}{\delta} \mathbb{I}(\lceil t' \rceil \leq \lfloor t \rfloor) \int_{\lceil t' \rceil}^{\lfloor t \rfloor} R_{\theta}^{\gamma}(t, s) D_{\ell,\gamma}^{s,t'} (ds + \sqrt{\tau \delta} dB^s) \right) + \nabla_r^2 \ell_{\lfloor t \rfloor}(r_{\gamma}^t; z) [D_{r,\gamma}^{t,t'}] \rho_{r,\gamma}^{t,t'}, \quad (234)$$

$$D_{r,\gamma}^{t,t'} = -\frac{1}{\delta} \mathbb{I}(\lceil t' \rceil \leq \lfloor t \rfloor) \int_{\lceil t' \rceil}^{\lfloor t \rfloor} R_{\theta}^{\gamma}(t, s) \nabla_r \ell_{\lfloor s \rfloor}(r_{\gamma}^s; z) D_{r,\gamma}^{s,t'} (ds + \sqrt{\tau \delta} dB^s) - \sqrt{\frac{\tau}{\delta}} R_{\theta}^{\gamma}(t, t') \ell_{\lceil t' \rceil}(r_{\gamma}^{t'}; z), \quad (235)$$

where  $U_\gamma^t, w_\gamma^t$  are centered Gaussian processes with covariance kernels  $\Sigma_\ell^\gamma/\delta$  and  $C_\theta^\gamma$  respectively. Then, set  $C_\theta^\gamma, R_\theta^\gamma, \Sigma_\ell^\gamma, R_\ell^\gamma, \Gamma^\gamma$  as

$$\begin{aligned} C_\theta^\gamma(t, t') &= \mathbb{E}[\theta_\gamma^t \theta_\gamma^{t'\top}], & R_\theta^\gamma(t, t') &= \mathbb{E}[\rho_{\theta, \gamma}^{t, t'}], \\ \Sigma_\ell^\gamma(t, t') &= \mathbb{E}[L_\gamma^t L_\gamma^{t'\top}], & L_\gamma^t &:= \int_0^{[t]} \ell_{[s]}(r_\gamma^s; z)(ds + \sqrt{\tau\delta} dB^s), \\ R_\ell^\gamma(t, t') &= \mathbb{E}[\rho_{\ell, \gamma}^{t, t'}] + \sqrt{\tau\delta} \mathbb{E}[D_{\ell, \gamma}^{t, t'}], & \Gamma^\gamma(t) &= \mathbb{E}[\nabla_r \ell_{[t]}(r_\gamma^t; z)], \end{aligned} \quad (236)$$

where we set  $R_\theta^\gamma(t, t') = R_\ell^\gamma(t, t') = 0$  for  $[t] < [t']$ .

The above equation agrees with the discretized DMFT equation  $\mathfrak{S}^\gamma$  at discrete time points  $t = t_i = i\gamma$  ( $i \geq 0$ ), and has a unique solution since it is piecewise constant in each interval  $[t_i, t_{i+1})$ .

We then define mappings  $\mathcal{T}_{\theta \rightarrow \ell}^\gamma: (C_\theta, R_\theta) \mapsto (\Sigma_\ell, R_\ell, \Gamma)$  and  $\mathcal{T}_{\ell \rightarrow \theta}^\gamma: (\Sigma_\ell, R_\ell, \Gamma) \mapsto (C_\theta, R_\theta)$  similarly to  $\mathcal{T}_{\theta \rightarrow \ell}$  and  $\mathcal{T}_{\ell \rightarrow \theta}$  in Appendix C but with the DMFT equation  $\mathfrak{S}$  replaced by its discretized version  $\mathfrak{S}^\gamma$ . We also define their composition  $\mathcal{T}^\gamma := \mathcal{T}_{\ell \rightarrow \theta}^\gamma \circ \mathcal{T}_{\theta \rightarrow \ell}^\gamma$ .

Since the solution to the discretized DMFT equation  $\mathfrak{S}^\gamma$  is determined by the values at discrete time points  $\{t_i = i\gamma : i \geq 0\}$ , the solution exists uniquely by induction. Let  $X^\gamma = (C_\theta^\gamma, R_\theta^\gamma)$  and  $Y^\gamma = (\Sigma_\ell^\gamma, R_\ell^\gamma, \Gamma^\gamma)$  be the solution to  $\mathfrak{S}^\gamma$ . Then,  $X^\gamma$  is the unique fixed point of  $\mathcal{T}^\gamma$ , i.e.,  $\mathcal{T}^\gamma(X^\gamma) = X^\gamma$ .

We show that the unique solutions  $X^\gamma$  and  $Y^\gamma$  belong to some admissible spaces defined in Appendix C.1.

**Lemma D.5** *There exist admissible spaces  $\mathcal{S}_\theta$  and  $\mathcal{S}_\ell$  such that  $X^\gamma \in \mathcal{S}_\theta$  and  $Y^\gamma \in \mathcal{S}_\ell$ .*

**Proof** Since  $C_\theta^\gamma, \Sigma_\ell^\gamma, R_\theta^\gamma, R_\ell^\gamma$ , and  $\Gamma^\gamma$  are piecewise constant, the continuity conditions are automatically satisfied. Since the solutions are bounded, we can take the spaces  $\mathcal{S}_\theta$  and  $\mathcal{S}_\ell$  large enough so that the boundedness conditions are also satisfied.  $\blacksquare$

Let  $X = (C_\theta, R_\theta)$  be the unique fixed point of  $\mathcal{T}$  shown in Theorem 3.1. We control their distance as follows.

$$\text{dist}_\lambda(X, X^\gamma) = \text{dist}_\lambda(\mathcal{T}(X), \mathcal{T}^\gamma(X^\gamma)) \leq \text{dist}_\lambda(\mathcal{T}(X), \mathcal{T}(X^\gamma)) + \text{dist}_\lambda(\mathcal{T}(X^\gamma), \mathcal{T}^\gamma(X^\gamma)). \quad (237)$$

As in the proof of Theorem 3.1, we can take  $\lambda$  large enough so that  $\mathcal{T}$  is a contraction and the first term is bounded by  $(1/2)\text{dist}_\lambda(X, X^\gamma)$ . Then, we have

$$\text{dist}_\lambda(X, X^\gamma) \leq 2 \cdot \text{dist}_\lambda(\mathcal{T}(X^\gamma), \mathcal{T}^\gamma(X^\gamma)). \quad (238)$$

The next lemma bounds the right-hand side. We prove it in Appendix D.3.1.

**Lemma D.6** *For any  $\lambda > 0$ , there exists a constant  $K > 0$  independent of  $\gamma$  such that*

$$\text{dist}_\lambda(\mathcal{T}(X^\gamma), \mathcal{T}^\gamma(X^\gamma)) \leq K\sqrt{\gamma}. \quad (239)$$

Therefore, for sufficiently large  $\lambda$ , we have

$$\text{dist}_\lambda(X, X^\gamma) \leq K\sqrt{\gamma}, \quad \text{dist}_\lambda(Y, Y^\gamma) \leq K\sqrt{\gamma}. \quad (240)$$

Finally, we couple  $\theta^t$  and  $\theta_\gamma^t$  so that they are close. We prove it in Appendix D.3.2.

**Lemma D.7** *There exist a constant  $K > 0$  independent of  $\gamma$  and a coupling of the processes  $\theta^t$  and  $\theta_\gamma^t$  such that*

$$\sup_{0 \leq t \leq T} \sqrt{\mathbb{E} \|\theta^t - \theta_\gamma^t\|_2^2} \leq K(\sqrt{\gamma} + \text{dist}_\lambda(Y, Y^\gamma)), \quad (241)$$

$$\sup_{0 \leq t \leq T} \sqrt{\mathbb{E} \|r^t - r_\gamma^t\|_2^2} \leq K(\sqrt{\gamma} + \text{dist}_\lambda(X, X^\gamma)). \quad (242)$$

By the above Lemmas, we have

$$W_2(\mathbb{P}(\theta^{t_1}, \dots, \theta^{t_L}), \mathbb{P}(\theta_\gamma^{t_1}, \dots, \theta_\gamma^{t_L})) \leq \sqrt{\sum_{i=1}^L \mathbb{E} \|\theta^{t_i} - \theta_\gamma^{t_i}\|_2^2} \leq \sqrt{KL} \cdot \sqrt{\gamma}, \quad (243)$$

$$W_2(\mathbb{P}(r^{t_1}, \dots, r^{t_L}, z), \mathbb{P}(r_\gamma^{t_1}, \dots, r_\gamma^{t_L}, z)) \leq \sqrt{\sum_{i=1}^L \mathbb{E} \|r^{t_i} - r_\gamma^{t_i}\|_2^2} \leq \sqrt{KL} \cdot \sqrt{\gamma}. \quad (244)$$

Sending  $\gamma \rightarrow 0$  completes the proof of Lemma D.3.

### D.3.1. PROOF OF LEMMA D.6

Let  $\bar{Y}^\gamma = (\bar{\Sigma}_\ell^\gamma, \bar{R}_\ell^\gamma, \bar{\Gamma}^\gamma) = \mathcal{T}_{\theta \rightarrow \ell}(X^\gamma)$  and  $Y^\gamma = (\Sigma_\ell^\gamma, R_\ell^\gamma, \Gamma^\gamma) = \mathcal{T}_{\theta \rightarrow \ell}(X^\gamma)$ . Also, let  $\bar{X}^\gamma = (\bar{C}_\theta^\gamma, \bar{R}_\theta^\gamma) = \mathcal{T}_{\ell \rightarrow \theta}(\bar{Y}^\gamma)$ .

In the following,  $K$  denotes a positive constant independent of  $\gamma$  whose value may change from line to line. Note that  $K$  can depend on  $\lambda$  since we fix  $\lambda$  and send  $\gamma \rightarrow 0$ .

The proof proceeds as follows.

1. We show that  $\text{dist}_\lambda(\bar{Y}^\gamma, Y^\gamma) \leq K\sqrt{\gamma}$  for some constant  $K > 0$ .
2. We show that  $\text{dist}_\lambda(\bar{X}^\gamma, X^\gamma) \leq K(\sqrt{\gamma} + \text{dist}_\lambda(\bar{Y}^\gamma, Y^\gamma))$  for some constant  $K > 0$ .

Together they prove the Lemma D.6.

**Bound of  $\text{dist}_\lambda(\bar{\Sigma}_\ell^\gamma, \Sigma_\ell^\gamma)$**  Let  $w_\gamma \sim \text{GP}(0, C_\theta^\gamma)$ . Let  $\bar{r}_\gamma^t$  and  $r_\gamma^t$  be the solutions to the following equations:

$$\bar{r}_\gamma^t = w_\gamma^t - \frac{1}{\delta} \int_0^t R_\theta^\gamma(t, s) \ell_s(\bar{r}_\gamma^s; z) (ds + \sqrt{\tau\delta} dB^s), \quad (245)$$

$$r_\gamma^t = w_\gamma^t - \frac{1}{\delta} \int_0^{\lfloor t \rfloor} R_\theta^\gamma(t, s) \ell_{\lfloor s \rfloor}(r_\gamma^s; z) (ds + \sqrt{\tau\delta} dB^s). \quad (246)$$

We have

$$\begin{aligned} \bar{r}_\gamma^t - r_\gamma^t &= -\frac{1}{\delta} \int_0^{\lfloor t \rfloor} R_\theta^\gamma(t, s) (\ell_s(\bar{r}_\gamma^s; z) - \ell_{\lfloor s \rfloor}(r_\gamma^s; z)) (ds + \sqrt{\tau\delta} dB^s) \\ &\quad - \frac{1}{\delta} \int_{\lfloor t \rfloor}^t R_\theta^\gamma(t, s) \ell_s(\bar{r}_\gamma^s; z) (ds + \sqrt{\tau\delta} dB^s), \end{aligned} \quad (247)$$

and thus

$$\begin{aligned}
 \mathbb{E}\|\bar{r}_\gamma^t - r_\gamma^t\|_2^2 &\leq 2\left(\frac{T}{\delta^2} + \frac{\tau}{\delta}\right) \int_0^{\lfloor t \rfloor} \|R_\theta^\gamma(t, s)\|_2^2 \mathbb{E}\|\ell_s(\bar{r}_\gamma^s; z) - \ell_{\lfloor s \rfloor}(r_\gamma^s; z)\|_2^2 ds \\
 &\quad + 2\left(\frac{T}{\delta^2} + \frac{\tau}{\delta}\right) \int_{\lfloor t \rfloor}^t \|R_\theta^\gamma(t, s)\|_2^2 \mathbb{E}\|\ell_s(\bar{r}_\gamma^s; z)\|_2^2 ds \\
 &\leq 2\left(\frac{T}{\delta^2} + \frac{\tau}{\delta}\right) \Phi^2 M^2 \int_0^{\lfloor t \rfloor} \mathbb{E}[ (|s - \lfloor s \rfloor| + \|\bar{r}_\gamma^s - r_\gamma^s\|_2)^2 ] ds + 2\left(\frac{T}{\delta^2} + \frac{\tau}{\delta}\right) \Phi^3 (t - \lfloor t \rfloor) \\
 &\leq 2\left(\frac{T}{\delta^2} + \frac{\tau}{\delta}\right) \Phi^2 M^2 \int_0^{\lfloor t \rfloor} (2\gamma^2 + 2\mathbb{E}\|\bar{r}_\gamma^s - r_\gamma^s\|_2^2) ds + 2\left(\frac{T}{\delta^2} + \frac{\tau}{\delta}\right) \Phi^3 \gamma \\
 &\leq K\gamma + K \int_0^t \mathbb{E}\|\bar{r}_\gamma^s - r_\gamma^s\|_2^2 ds. \tag{248}
 \end{aligned}$$

By Grönwall's inequality, we have

$$\sup_{t \in [0, T]} \mathbb{E}\|\bar{r}_\gamma^t - r_\gamma^t\|_2^2 \leq K\gamma e^{KT} \leq K\gamma. \tag{249}$$

Let

$$\bar{L}_\gamma^t := \int_0^t \ell_s(\bar{r}_\gamma^s; z)(ds + \sqrt{\tau\delta} dB^s), \quad L_\gamma^t := \int_0^{\lfloor t \rfloor} \ell_{\lfloor s \rfloor}(r_\gamma^s; z)(ds + \sqrt{\tau\delta} dB^s). \tag{250}$$

Then, we have

$$\begin{aligned}
 \mathbb{E}\|\bar{L}_\gamma^t - L_\gamma^t\|_2^2 &\leq 2(T + \tau\delta) \int_0^{\lfloor t \rfloor} \mathbb{E}\|\ell_s(\bar{r}_\gamma^s; z) - \ell_{\lfloor s \rfloor}(r_\gamma^s; z)\|_2^2 ds + 2(T + \tau\delta) \int_{\lfloor t \rfloor}^t \mathbb{E}\|\ell_s(\bar{r}_\gamma^s; z)\|_2^2 ds \\
 &\leq 2(T + \tau\delta) M^2 \int_0^{\lfloor t \rfloor} \mathbb{E}[ (|s - \lfloor s \rfloor| + \|\bar{r}_\gamma^s - r_\gamma^s\|_2)^2 ] ds + 2(T + \tau\delta) \Phi (t - \lfloor t \rfloor) \\
 &\leq K\gamma. \tag{251}
 \end{aligned}$$

Let  $\{(\bar{U}_\gamma^t, U_\gamma^t)\}_{t \in [0, T]}$  be a centered Gaussian process with covariance  $\mathbb{E}\left[\begin{pmatrix} \bar{L}_\gamma^t \\ L_\gamma^t \end{pmatrix} \begin{pmatrix} L_\gamma^t \\ \bar{L}_\gamma^t \end{pmatrix}^\top\right] / \delta$ . Since  $\bar{U}_\gamma$  and  $U_\gamma$  have covariance kernels  $\bar{\Sigma}_\ell^\gamma / \delta$  and  $\Sigma_\ell^\gamma / \delta$  respectively, we have

$$\text{dist}_\lambda(\bar{\Sigma}_\ell^\gamma, \Sigma_\ell^\gamma) \leq \sup_{t \in [0, T]} e^{-\lambda t} \sqrt{\mathbb{E}\|\bar{U}_\gamma^t - U_\gamma^t\|_2^2} = \sup_{t \in [0, T]} e^{-\lambda t} \sqrt{\mathbb{E}\|\bar{L}_\gamma^t - L_\gamma^t\|_2^2 / \delta} \leq K\sqrt{\gamma}. \tag{252}$$

**Bound of  $\text{dist}_\lambda(\bar{\Gamma}^\gamma, \Gamma^\gamma)$ .** By Equation (249), we have

$$\begin{aligned}
 \text{dist}_\lambda(\bar{\Gamma}^\gamma, \Gamma^\gamma) &= \sup_{t \in [0, T]} e^{-\lambda t} \|\bar{\Gamma}^\gamma(t) - \Gamma^\gamma(t)\|_2 \leq \sup_{t \in [0, T]} \sqrt{\mathbb{E}\|\nabla_r \ell_t(\bar{r}_\gamma^t; z) - \nabla_r \ell_{\lfloor t \rfloor}(r_\gamma^t; z)\|_2^2} \\
 &\leq M \sup_{t \in [0, T]} \sqrt{2(t - \lfloor t \rfloor)^2 + 2\mathbb{E}\|\bar{r}_\gamma^t - r_\gamma^t\|_2^2} \leq K\sqrt{\gamma}. \tag{253}
 \end{aligned}$$

**Bound of  $\text{dist}_\lambda(R_\ell^1, R_\ell^2)$ .** Let  $\bar{\rho}_{r,\gamma}^{t,t'}$  and  $\rho_{r,\gamma}^{t,t'}$  be the stochastic processes satisfying

$$\bar{\rho}_{r,\gamma}^{t,t'} = -\frac{1}{\delta} \int_{t'}^t R_\theta^\gamma(t, s) \bar{\rho}_{\ell,\gamma}^{s,t'} (ds + \sqrt{\tau\delta} dB^s) - \frac{1}{\delta} R_\theta^\gamma(t, t') \nabla_r \ell_{t'}(\bar{r}_\gamma^{t'}; z), \quad (254)$$

$$\rho_{r,\gamma}^{t,t'} = -\frac{1}{\delta} \mathbb{I}(\lceil t' \rceil \leq \lfloor t \rfloor) \int_{\lceil t' \rceil}^{\lfloor t \rfloor} R_\theta^\gamma(t, s) \rho_{\ell,\gamma}^{s,t'} (ds + \sqrt{\tau\delta} dB^s) - \frac{1}{\delta} R_\theta^\gamma(t, t') \nabla_r \ell_{\lfloor t \rfloor}(r_\gamma^{t'}; z). \quad (255)$$

Let  $\bar{\rho}_{\ell,\gamma}^{t,t'} = \nabla_r \ell_t(\bar{r}_\gamma^t; z) \bar{\rho}_{r,\gamma}^{t,t'}$  and  $\rho_{\ell,\gamma}^{t,t'} = \nabla_r \ell_{\lfloor t \rfloor}(r_\gamma^t; z) \rho_{r,\gamma}^{t,t'}$ . Similarly, define  $\bar{D}_{\ell,\gamma}^{t,t'}$  and  $D_{\ell,\gamma}^{t,t'}$  for processes defined by  $\mathcal{T}_{\theta \rightarrow \ell}$  and  $\mathcal{T}_\theta^\gamma$ , respectively. Then, we have

$$\|\bar{R}_\ell^\gamma(t, t') - R_\ell^\gamma(t, t')\|_2^2 \leq \mathbb{E} \|\bar{\rho}_{\ell,\gamma}^{t,t'} - \rho_{\ell,\gamma}^{t,t'}\|_2^2 + \tau\delta \mathbb{E} \|\bar{D}_{\ell,\gamma}^{t,t'} - D_{\ell,\gamma}^{t,t'}\|_2^2. \quad (256)$$

We first bound  $\mathbb{E} \|\bar{\rho}_{\ell,\gamma}^{t,t'} - \rho_{\ell,\gamma}^{t,t'}\|_2^2$ . We have

$$\begin{aligned} \mathbb{E} \|\bar{\rho}_{\ell,\gamma}^{t,t'} - \rho_{\ell,\gamma}^{t,t'}\|_2^2 &= \mathbb{E} \|\nabla_r \ell_t(\bar{r}_\gamma^t; z) \bar{\rho}_{r,\gamma}^{t,t'} - \nabla_r \ell_{\lfloor t \rfloor}(r_\gamma^t; z) \rho_{r,\gamma}^{t,t'}\|_2^2 \\ &\leq 2\sqrt{\mathbb{E} \|\bar{\rho}_{r,\gamma}^{t,t'}\|_2^4} \sqrt{\mathbb{E} \|\nabla_r \ell_t(\bar{r}_\gamma^t; z) - \nabla_r \ell_{\lfloor t \rfloor}(r_\gamma^t; z)\|_2^4} + 2M^2 \mathbb{E} \|\bar{\rho}_{r,\gamma}^{t,t'} - \rho_{r,\gamma}^{t,t'}\|_2^2 \\ &\leq K\sqrt{\gamma^4 + \mathbb{E} \|\bar{r}_\gamma^t - r_\gamma^t\|_2^4} + K \mathbb{E} \|\bar{\rho}_{r,\gamma}^{t,t'} - \rho_{r,\gamma}^{t,t'}\|_2^2. \end{aligned} \quad (257)$$

First, we bound  $\mathbb{E} \|\bar{r}_\gamma^t - r_\gamma^t\|_2^4$ . We have

$$\begin{aligned} &\mathbb{E} \|\bar{r}_\gamma^t - r_\gamma^t\|_2^4 \\ &\leq K \int_0^{\lfloor t \rfloor} \|R_\theta^\gamma(t, s)\|_2^4 \mathbb{E} \|\ell_s(\bar{r}_\gamma^s; z) - \ell_{\lfloor s \rfloor}(r_\gamma^s; z)\|_2^4 ds + K(t - \lfloor t \rfloor) \int_{\lfloor t \rfloor}^t \|R_\theta^\gamma(t, s)\|_2^4 \mathbb{E} \|\ell_s(\bar{r}_\gamma^s; z)\|_2^4 ds \\ &\leq K \int_0^{\lfloor t \rfloor} (\gamma^4 + \mathbb{E} \|\bar{r}_\gamma^s - r_\gamma^s\|_2^4) ds + K\gamma^2 \leq K\gamma^2 + K \int_0^t \mathbb{E} \|\bar{r}_\gamma^s - r_\gamma^s\|_2^4 ds. \end{aligned} \quad (258)$$

By Grönwall's inequality, we have

$$\sup_{t \in [0, T]} \mathbb{E} \|r_1^t - r_2^t\|_2^4 \leq K\gamma^2. \quad (259)$$

Next, we bound  $\mathbb{E} \|\bar{\rho}_{r,\gamma}^{t,t'} - \rho_{r,\gamma}^{t,t'}\|_2^2$ . We have

$$\begin{aligned} \|\rho_{r,\gamma}^{t,t'} - \bar{\rho}_{r,\gamma}^{t,t'}\|_2 &\leq \frac{1}{\delta} \left\| \int_{t'}^{\lceil t' \rceil} R_\theta^\gamma(t, s) \bar{\rho}_{\ell,\gamma}^{s,t'} (ds + \sqrt{\tau\delta} dB^s) \right\|_2 + \frac{1}{\delta} \left\| \int_{\lfloor t \rfloor}^t R_\theta^\gamma(t, s) \bar{\rho}_{\ell,\gamma}^{s,t'} (ds + \sqrt{\tau\delta} dB^s) \right\|_2 \\ &\quad + \frac{1}{\delta} \left\| \int_{\lceil t' \rceil}^{\lfloor t \rfloor} R_\theta^\gamma(t, s) (\rho_{\ell,\gamma}^{s,t'} - \bar{\rho}_{\ell,\gamma}^{s,t'}) (ds + \sqrt{\tau\delta} dB^s) \right\|_2 \\ &\quad + \frac{1}{\delta} \|R_\theta^\gamma(t, t') (\nabla_r \ell_t(\bar{r}_\gamma^{t'}; z) - \nabla_r \ell_{\lceil t' \rceil}(r_\gamma^{t'}; z))\|_2, \end{aligned} \quad (260)$$

Thus, we have

$$\mathbb{E} \|\rho_{r,\gamma}^{t,t'} - \bar{\rho}_{r,\gamma}^{t,t'}\|_2^2 \leq K \int_{t'}^{\lceil t' \rceil} \|R_\theta^\gamma(t, s)\|_2^2 \mathbb{E} \|\bar{\rho}_{\ell,\gamma}^{s,t'}\|_2^2 ds + K \int_{\lfloor t \rfloor}^t \|R_\theta^\gamma(t, s)\|_2^2 \mathbb{E} \|\bar{\rho}_{\ell,\gamma}^{s,t'}\|_2^2 ds$$

$$\begin{aligned}
 & + K \int_{\lceil t' \rceil}^{\lfloor t \rfloor} \|R_\theta^\gamma(t, s)\|_2^2 \mathbb{E} \|\rho_{\ell, \gamma}^{s, t'} - \bar{\rho}_{\ell, \gamma}^{s, t'}\|_2^2 ds \\
 & + K \|R_\theta^\gamma(t, t')\|_2^2 \mathbb{E} \|\nabla_r \ell_t(\bar{r}_\gamma^{t'}; z) - \nabla_r \ell_{\lfloor t' \rfloor}(r_\gamma^{t'}; z)\|_2^2 \\
 & \leq K\gamma + K \int_{t'}^t \mathbb{E} \|\rho_{\ell, \gamma}^{s, t'} - \bar{\rho}_{\ell, \gamma}^{s, t'}\|_2^2 ds.
 \end{aligned} \tag{261}$$

By Grönwall's inequality, we have

$$\sup_{0 \leq t' \leq t \leq T} \mathbb{E} \|\rho_{r, \gamma}^{t, t'} - \bar{\rho}_{r, \gamma}^{t, t'}\|_2^2 \leq K\gamma. \tag{262}$$

Next, we bound  $\mathbb{E} \|\bar{D}_{\ell, \gamma}^{t, t'} - D_{\ell, \gamma}^{t, t'}\|_2^2$  similarly. We do not repeat the details, but following the same argument, we have

$$\mathbb{E} \|\bar{D}_{\ell, \gamma}^{t, t'} - D_{\ell, \gamma}^{t, t'}\|_2^2 \leq K\gamma. \tag{263}$$

Therefore, we have

$$\|\bar{R}_\ell^\gamma(t, t') - R_\ell^\gamma(t, t')\|_2^2 \leq K\gamma. \tag{264}$$

By Equations (252), (253) and (264), we have

$$\text{dist}_\lambda(\bar{Y}^\gamma, Y^\gamma) \leq K\sqrt{\gamma}. \tag{265}$$

**Bound of  $\text{dist}_\lambda(\bar{C}_\theta^\gamma, C_\theta^\gamma)$ .** Let  $\bar{U}_\gamma \sim \text{GP}(0, \bar{\Sigma}_\ell^\gamma/\delta)$  and  $U_\gamma \sim \text{GP}(0, \Sigma_\ell^\gamma/\delta)$  be Gaussian processes coupled such that

$$\sup_{t \in [0, T]} e^{-\lambda t} \sqrt{\mathbb{E} \|\bar{U}_\gamma^t - U_\gamma^t\|_2^2} \leq 2 \cdot \text{dist}_\lambda(\bar{\Sigma}_\ell^\gamma, \Sigma_\ell^\gamma). \tag{266}$$

Let  $\bar{\theta}_\gamma^t$  and  $\theta_\gamma^t$  be the solution of

$$\bar{\theta}_\gamma^t = \theta^0 + \bar{U}_\gamma^t - \int_0^t \left( h_s(\bar{\theta}_\gamma^s) + \bar{\Gamma}^\gamma(s) \bar{\theta}_\gamma^s + \int_0^s \bar{R}_\ell^\gamma(s, s') \bar{\theta}_\gamma^{s'} ds' \right) ds, \tag{267}$$

$$\theta_\gamma^t = \theta^0 + U_\gamma^t - \int_0^{\lfloor t \rfloor} \left( h_{\lfloor s \rfloor}(\theta_\gamma^s) + \Gamma^\gamma(s) \theta_\gamma^s + \int_0^{\lfloor s \rfloor} R_\ell^\gamma(s, s') \theta_\gamma^{s'} ds' \right) ds. \tag{268}$$

Then, we have

$$\begin{aligned}
 \mathbb{E} \|\bar{\theta}_\gamma^t - \theta_\gamma^t\|_2^2 & \leq K \left( \gamma + \mathbb{E} \|\bar{U}_\gamma^t - U_\gamma^t\|_2^2 + \sup_{t \in [0, T]} \|\bar{\Gamma}^\gamma(t) - \Gamma^\gamma(t)\|_2^2 + \sup_{0 \leq t' \leq t \leq T} \|\bar{R}_\ell^\gamma(t, t') - R_\ell^\gamma(t, t')\|_2^2 \right) \\
 & \quad + K \int_0^t \mathbb{E} \|\bar{\theta}_\gamma^s - \theta_\gamma^s\|_2^2 ds \\
 & \leq K(\gamma + \text{dist}_\lambda(\bar{Y}^\gamma, Y^\gamma)^2) + K \int_0^t \mathbb{E} \|\bar{\theta}_\gamma^s - \theta_\gamma^s\|_2^2 ds.
 \end{aligned} \tag{269}$$

By Grönwall's inequality, we have

$$\sup_{t \in [0, T]} \mathbb{E} \|\bar{\theta}_\gamma^t - \theta_\gamma^t\|_2^2 \leq K(\gamma + \text{dist}_\lambda(\bar{Y}^\gamma, Y^\gamma)^2). \quad (270)$$

Let  $\{(\bar{w}_\gamma^t, w_\gamma^t)\}_{t \in [0, T]}$  be a centered Gaussian process with covariance  $\mathbb{E} \begin{bmatrix} \bar{\theta}_\gamma^t & \bar{\theta}_\gamma^t \\ \theta_\gamma^t & \theta_\gamma^t \end{bmatrix}^\top$ . Since  $\bar{w}_\gamma$  and  $w_\gamma$  have covariance kernels  $\bar{C}_\theta^\gamma$  and  $C_\theta^\gamma$  respectively, we have

$$\text{dist}_\lambda(\bar{C}_\theta^\gamma, C_\theta^\gamma) \leq \sup_{t \in [0, T]} e^{-\lambda t} \sqrt{\mathbb{E} \|\bar{w}_\gamma^t - w_\gamma^t\|_2^2} = \sup_{t \in [0, T]} e^{-\lambda t} \sqrt{\mathbb{E} \|\bar{\theta}_\gamma^t - \theta_\gamma^t\|_2^2} \leq K(\sqrt{\gamma} + \text{dist}_\lambda(\bar{Y}^\gamma, Y^\gamma)). \quad (271)$$

**Bound of  $\text{dist}_\lambda(R_\theta^1, R_\theta^2)$ .** Let  $\bar{\rho}_{\theta, \gamma}^{t, t'}$  and  $\rho_{\theta, \gamma}^{t, t'}$  be the stochastic processes satisfying

$$\bar{\rho}_{\theta, \gamma}^{t, t'} = I_m - \int_{t'}^t \left( (\nabla_\theta h_s(\bar{\theta}_\gamma^s) + \bar{\Gamma}^\gamma(s)) \bar{\rho}_{\theta, \gamma}^{s, t'} + \int_{t'}^s \bar{R}_\ell^\gamma(s, s') \bar{\rho}_{\theta, \gamma}^{s', t'} ds' \right) ds, \quad (272)$$

$$\rho_{\theta, \gamma}^{t, t'} = I_m - \mathbb{I}(\lceil t' \rceil \leq \lfloor t \rfloor) \int_{\lceil t' \rceil}^{\lfloor t \rfloor} \left( (\nabla_\theta h_{\lfloor s \rfloor}(\theta_\gamma^s) + \Gamma^\gamma(s)) \rho_{\theta, \gamma}^{s, t'} + \int_{\lceil t' \rceil}^{\lfloor s \rfloor} R_\ell^\gamma(s, s') \rho_{\theta, \gamma}^{s', t'} ds' \right) ds, \quad (273)$$

and set  $\bar{R}_\theta^\gamma(t, t') = \mathbb{E}[\bar{\rho}_{\theta, \gamma}^{t, t'}]$  and  $R_\theta^\gamma(t, t') = \mathbb{E}[\rho_{\theta, \gamma}^{t, t'}]$ . We have

$$\begin{aligned} \|\bar{\rho}_{\theta, \gamma}^{t, t'} - \rho_{\theta, \gamma}^{t, t'}\|_2 &\leq K \left( \gamma + \sup_{t \in [0, T]} \|\bar{\Gamma}^\gamma(t) - \Gamma^\gamma(t)\|_2 + \sup_{0 \leq t' \leq t \leq T} \|\bar{R}_\ell^\gamma(t, t') - R_\ell^\gamma(t, t')\|_2 \right) \\ &\quad + K \int_{t'}^t \|\bar{\rho}_{\theta, \gamma}^{s, t'} - \rho_{\theta, \gamma}^{s, t'}\|_2 ds \\ &\leq K(\gamma + \text{dist}_\lambda(\bar{Y}^\gamma, Y^\gamma)) + K \int_{t'}^t \|\bar{\rho}_{\theta, \gamma}^{s, t'} - \rho_{\theta, \gamma}^{s, t'}\|_2 ds. \end{aligned} \quad (274)$$

By Grönwall's inequality, we have

$$\sup_{0 \leq t' \leq t \leq T} \|\bar{\rho}_{\theta, \gamma}^{t, t'} - \rho_{\theta, \gamma}^{t, t'}\|_2 \leq K(\gamma + \text{dist}_\lambda(\bar{Y}^\gamma, Y^\gamma)). \quad (275)$$

Thus, we have

$$\begin{aligned} \text{dist}_\lambda(\bar{R}_\theta^\gamma, R_\theta^\gamma) &\leq \sup_{0 \leq t' \leq t \leq T} \|\bar{R}_\theta^\gamma(t, t') - R_\theta^\gamma(t, t')\|_2 \leq \sup_{0 \leq t' \leq t \leq T} \mathbb{E} \|\bar{\rho}_{\theta, \gamma}^{t, t'} - \rho_{\theta, \gamma}^{t, t'}\|_2 \\ &\leq K(\gamma + \text{dist}_\lambda(\bar{Y}^\gamma, Y^\gamma)). \end{aligned} \quad (276)$$

By Equations (271) and (276), we have

$$\text{dist}_\lambda(\bar{X}^\gamma, X^\gamma) \leq K(\sqrt{\gamma} + \text{dist}_\lambda(\bar{Y}^\gamma, Y^\gamma)). \quad (277)$$

### D.3.2. PROOF OF LEMMA D.7

Following the same calculation that led to Equation (249), we can show that there exists a coupling of the processes  $r^t$  and  $r_\gamma^t$  such that

$$\sup_{t \in [0, T]} \sqrt{\mathbb{E} \|r^t - r_\gamma^t\|_2^2} \leq K(\sqrt{\gamma} + \text{dist}_\lambda(X, X^\gamma)). \quad (278)$$

Following the same calculation that led to Equation (270), we can show that there exists a coupling of the processes  $\theta^t$  and  $\theta_\gamma^t$  such that

$$\sup_{t \in [0, T]} \sqrt{\mathbb{E} \|\theta^t - \theta_\gamma^t\|_2^2} \leq K(\sqrt{\gamma} + \text{dist}_\lambda(Y, Y^\gamma)). \quad (279)$$

## Appendix E. Details of Applications and Special Cases

### E.1. Planted Models

#### E.1.1. SETUP FOR PLANTED MODELS

The setting of supervised learning with (noisy) target labels  $y_i = f_*(\boldsymbol{\theta}^{*\top} \mathbf{x}_i) + z_i$  for some target parameter  $\boldsymbol{\theta}^* \in \mathbb{R}^{d \times m}$  and target function  $f_*: \mathbb{R}^m \rightarrow \mathbb{R}$  can be handled within our framework. Consider the following stochastic process with a planted signal  $\boldsymbol{\theta}^* \in \mathbb{R}^{d \times m}$ :

$$\hat{\boldsymbol{\theta}}^{k+1} = \hat{\boldsymbol{\theta}}^k - \eta \cdot \left( \frac{1}{d} h_{t_k}(\hat{\boldsymbol{\theta}}^k) + \frac{1}{B} \sum_{i \in \mathcal{B}^k} \mathbf{x}_i \ell_{t_k}(\hat{r}_i^k, r_i^*; z_i)^\top \right), \quad (280)$$

with  $\hat{r}^k = \mathbf{X} \hat{\boldsymbol{\theta}}^k$  and  $r^* = \mathbf{X} \boldsymbol{\theta}^*$ . We consider the SGF with the planted signal that approximates the above SGD dynamics:

$$d\boldsymbol{\theta}^t = - \left( h_t(\boldsymbol{\theta}^t) + \frac{1}{\delta} \mathbf{X}^\top \ell_t(\mathbf{r}^t, \mathbf{r}^*; \mathbf{z}) \right) dt + \sqrt{\frac{\tau}{\delta}} \sum_{i=1}^n \mathbf{x}_i \ell_t(r_i^t, r_i^*; z_i)^\top dB_i^t, \quad (281)$$

with  $\mathbf{r}^t = \mathbf{X} \boldsymbol{\theta}^t$ . This can be mapped to our general SGF setup (1) by concatenating the parameters  $\boldsymbol{\theta}^t$  and  $\boldsymbol{\theta}^*$  as follows.

$$\begin{aligned} d(\boldsymbol{\theta}^t, \boldsymbol{\theta}^*) &= - \left( (h_t(\boldsymbol{\theta}^t), 0) + \frac{1}{\delta} \mathbf{X}^\top (\ell_t(\mathbf{X}(\boldsymbol{\theta}^t, \boldsymbol{\theta}^*); \mathbf{z}), 0) \right) dt \\ &\quad + \sqrt{\frac{\tau}{\delta}} \sum_{i=1}^n \mathbf{x}_i \left( \ell_t((\boldsymbol{\theta}^t, \boldsymbol{\theta}^*)^\top \mathbf{x}_i; z_i)^\top, 0 \right) dB_i^t. \end{aligned} \quad (282)$$

Here,  $\boldsymbol{\theta}^*$  is constant over time. By applying Theorem 3.1 and Theorem 3.2 to the concatenated parameter of shape  $\mathbb{R}^{d \times 2m}$ , we obtain a DMFT characterization of the SGF with a planted signal (281) as a corollary.

**Corollary E.1 (DMFT characterization of SGF with a planted signal)** *Suppose Assumptions A.1 and A.2 hold. Take  $T_* > 0$  in Theorem 3.1. Furthermore, assume that  $\boldsymbol{\theta}^* \in \mathbb{R}^{d \times m}$  is independent of*

$\mathbf{X}, \mathbf{z}, \theta^0$  and for all  $p \geq 1$ , its empirical distribution converge in  $p$ -Wasserstein distance to  $\mathbb{P}(\theta^*)$  almost surely as  $d \rightarrow \infty$ . Then, for any  $T \in [0, T_*]$ ,  $L \in \mathbb{N}$ , and  $0 \leq t_1 < \dots < t_L \leq T$ , we have

$$\text{p-lim}_{n,d \rightarrow \infty} W_2 \left( \hat{\mathbb{P}}(\theta^{t_1}, \dots, \theta^{t_L}, \theta^*), \mathbb{P}(\theta^{t_1}, \dots, \theta^{t_L}, \theta^*) \right)^2 = 0, \quad (283)$$

$$\text{p-lim}_{n,d \rightarrow \infty} W_2 \left( \hat{\mathbb{P}}(\mathbf{r}^{t_1}, \dots, \mathbf{r}^{t_L}, \mathbf{r}^*, \mathbf{z}), \mathbb{P}(\mathbf{r}^{t_1}, \dots, \mathbf{r}^{t_L}, \mathbf{r}^*, \mathbf{z}) \right)^2 = 0, \quad (284)$$

where  $(\theta^t, \theta^*, \mathbf{r}^t, \mathbf{r}^*, \mathbf{z})$  is the unique solution of the DMFT equation given in Appendix E.1.2.

The proof appears in Appendix E.1.3. This result is an extension of Celentano et al. [14, Corollary 4.1] to the SGF dynamics.

### E.1.2. THE DMFT EQUATION

We state the DMFT equation  $\mathfrak{S}^*$  for planted model (281) for functions  $R_\ell, C_\theta: (\mathbb{R}_{\geq 0} \cup \{*\})^2 \rightarrow \mathbb{R}^{m \times m}$ ,  $R_\theta, \Sigma_\ell: \mathbb{R}_{\geq 0}^2 \rightarrow \mathbb{R}^{m \times m}$ , and  $\Gamma: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^{m \times m}$  self-consistently as follows. First, given  $\Sigma_\ell, R_\ell, \Gamma$ , define stochastic processes  $\{\theta^t \in \mathbb{R}^m\}_{t \geq 0}$  and  $\{\rho_\theta^{t,t'} \in \mathbb{R}^{m \times m}\}_{t \geq t' \geq 0}$  by the following equations.

$$\theta^t = \theta^0 + U^t - \int_0^t \left( (h_s(\theta^s) + \Gamma(s)\theta^s) + \int_0^s R_\ell(s, s')\theta^{s'} ds' + R_\ell(t, *)\theta^* \right) ds, \quad (285)$$

$$\rho_\theta^{t,t'} = I_m - \int_{t'}^t \left( (\nabla_\theta h_s(\theta^s) + \Gamma(s))\rho_\theta^{s,t'} + \int_{t'}^s R_\ell(s, s')\rho_\theta^{s',t'} ds' \right) ds, \quad (286)$$

where  $U \sim \text{GP}(0, \Sigma_\ell/\delta)$ . Then, set  $C_\theta, R_\theta$  as

$$C_\theta(t, t') = \mathbb{E}[\theta^t \theta^{t'\top}] \quad (s, t \in [0, T] \cup \{*\}), \quad R_\theta(t, t') = \mathbb{E}[\rho_\theta^{t,t'}] \quad (t \geq t'), \quad (287)$$

and  $R_\theta(t, t') = 0$  for  $t < t'$ .

Next, given  $C_\theta, R_\theta$ , define stochastic processes  $\{r^t \in \mathbb{R}^m\}_{t \geq 0}$ ,  $r^* \in \mathbb{R}^m$ , and  $\{\rho_\ell^{t,t'}, \rho_{\ell^*}^{t,t'}, D_\ell^{t,t'}, D_{\ell^*}^{t,t'} \in \mathbb{R}^{m \times m}\}_{t \geq t' \geq 0}$  by the following equations.

$$r^t = w^t - \frac{1}{\delta} \int_0^t R_\theta(t, s)\ell_s(r^s, r^*; z)(ds + \sqrt{\tau\delta} dB^s), \quad r^* = w^*, \quad w \sim \text{GP}(0, C_\theta), \quad (288)$$

$$\rho_\ell^{t,t'} = \nabla_r \ell_t(r^t, r^*; z)\rho_r^{t,t'}, \quad (289)$$

$$\rho_{\ell^*}^{t,t'} = \nabla_r \ell_t(r^t, r^*; z)\rho_{r^*}^{t,t'}, \quad (290)$$

$$D_\ell^{t,t'} = \nabla_r \ell_t(r^t, r^*; z) \left( -\frac{1}{\delta} \int_{t'}^t R_\theta(t, s)D_\ell^{s,t'}(ds + \sqrt{\tau\delta} dB^s) \right) + \nabla_r^2 \ell_t(r^t, r^*; z)[D_r^{t,t'}]\rho_r^{t,t'}, \quad (291)$$

$$D_{\ell^*}^{t,t'} = \nabla_r \ell_t(r^t, r^*; z) \left( -\frac{1}{\delta} \int_{t'}^t R_\theta(t, s)D_{\ell^*}^{s,t'}(ds + \sqrt{\tau\delta} dB^s) \right) + \nabla_r^2 \ell_t(r^t, r^*; z)[D_r^{t,t'}]\rho_{r^*}^{t,t'}, \quad (292)$$

where  $B^t$  is a Brownian motion in  $\mathbb{R}$ , and we defined the auxiliary processes  $\rho_r^{t,t'}, \rho_{r^*}^{t,t'} \in \mathbb{R}^{m \times m}$ , and  $D_r^{t,t'} \in \mathbb{R}^m$  as

$$\rho_r^{t,t'} = -\frac{1}{\delta} \int_{t'}^t R_\theta(t, s) \rho_\ell^{s,t'} (ds + \sqrt{\tau\delta} dB^s) - \frac{1}{\delta} R_\theta(t, t') \nabla_r \ell_{t'}(r^{t'}, r^*; z), \quad (293)$$

$$\rho_{r^*}^{t,t'} = -\frac{1}{\delta} \int_{t'}^t R_\theta(t, s) \rho_{\ell^*}^{s,t'} (ds + \sqrt{\tau\delta} dB^s) - \frac{1}{\delta} R_\theta(t, t') \nabla_{r^*} \ell_{t'}(r^{t'}, r^*; z), \quad (294)$$

$$D_r^{t,t'} = -\frac{1}{\delta} \int_{t'}^t R_\theta(t, s) \nabla_r \ell_s(r^s, r^*; z) D_r^{s,t'} (ds + \sqrt{\tau\delta} dB^s) - \sqrt{\frac{\tau}{\delta}} R_\theta(t, t') \ell_{t'}(r^{t'}, r^*; z). \quad (295)$$

Then, set  $\Sigma_\ell, R_\ell, \Gamma$  as

$$\Sigma_\ell(t, t') = \mathbb{E}[L^t L^{t'\top}], \quad L^t := \int_0^t \ell_s(r^s, r^*; z) (ds + \sqrt{\tau\delta} dB^s), \quad (296)$$

$$R_\ell(t, t') = \mathbb{E}[\rho_\ell^{t,t'}] + \sqrt{\tau\delta} \mathbb{E}[D_\ell^{t,t'}] \quad (t \geq t'), \quad (297)$$

$$R_\ell(t, *) = \mathbb{E}[\nabla_{r^*} \ell_t(r^t, r^*; z)] + \int_0^t \mathbb{E}[\rho_{\ell^*}^{t,t'} + \sqrt{\tau\delta} D_{\ell^*}^{t,t'}] dt' \quad (t \geq 0), \quad (298)$$

$$\Gamma(t) = \mathbb{E}[\nabla_r \ell_t(r^t, r^*; z)], \quad (299)$$

and  $R_\ell(t, t') = 0$  for  $t < t'$ .

Then, the solution of the DMFT system  $\mathfrak{S}^*$  is defined as a fixed point of the above two mappings.

In the informal notation, the DMFT equation for the planted model can be written as

$$\frac{d}{dt} \theta^t = u^t - (h_t(\theta^t) + \Gamma(t)\theta^t) - \int_0^t R_\ell(t, s) \theta^s ds - R_\ell(t, *) \theta^*, \quad u \sim \text{GP}(0, C_\ell/\delta),$$

$$r^t = w^t - \frac{1}{\delta} \int_0^t R_\theta(t, s) \ell_s(r^s, r^*; z) (ds + \sqrt{\tau\delta} dB^s), \quad r^* = w^*, \quad w \sim \text{GP}(0, C_\theta),$$

$$C_\theta(t, t') = \mathbb{E}[\theta^t \theta^{t'\top}] \quad (t, t' \in [0, T] \cup \{*\}), \quad R_\theta(t, t') = \mathbb{E}\left[\frac{\partial \theta^t}{\partial u^{t'}}\right],$$

$$C_\ell(t, t') = \mathbb{E}[\ell_t(r^t, r^*; z) (1 + \sqrt{\tau\delta} \dot{B}^t) \ell_{t'}(r^{t'}, r^*; z)^\top (1 + \sqrt{\tau\delta} \dot{B}^{t'})],$$

$$R_\ell(t, t') = \mathbb{E}\left[\frac{\partial \ell_t(r^t, r^*; z)}{\partial w^{t'}}\right], \quad R_\ell(t, *) = \mathbb{E}\left[\frac{\partial \ell_t(r^t, r^*; z)}{\partial w^*}\right], \quad \Gamma(t) = \mathbb{E}[\nabla_r \ell_t(r^t, r^*; z)]. \quad (300)$$

### E.1.3. PROOF OF COROLLARY E.1

We transform the DMFT equation  $\mathfrak{S}$  applied to the planted model (281) to the DMFT equation  $\mathfrak{S}^*$  defined above. We distinguish the variables in  $\mathfrak{S}$  by adding a bar over them, e.g.,  $\bar{\theta}^t$ . The variables in  $\mathfrak{S}$  have dimensions  $2m$ . Identifying components that are trivially zero, we see that the solution to  $\mathfrak{S}$

is of the form

$$\begin{aligned}
 \bar{\theta}^t &= \begin{pmatrix} \bar{\theta}_1^t \\ \theta^* \end{pmatrix}, \quad \bar{U}^t = \begin{pmatrix} \bar{U}_1^t \\ 0 \end{pmatrix}, \quad \bar{\rho}_\theta^{t,t'} = \begin{pmatrix} \bar{\rho}_{\theta,1}^{t,t'} & - \\ 0 & I_m \end{pmatrix}, \\
 \bar{C}_\theta(t, t') &= \begin{pmatrix} \bar{C}_\theta^{11}(t, t') & \bar{C}_\theta^{12}(t, t') \\ \bar{C}_\theta^{12}(t, t')^\top & \bar{C}_\theta^{22}(t, t') \end{pmatrix}, \quad \bar{R}_\theta(t, t') = \begin{pmatrix} \bar{R}_\theta^1(t, t') & - \\ 0 & I_m \end{pmatrix}, \\
 \bar{r}^t &= \begin{pmatrix} \bar{r}_1^t \\ \bar{r}_2^t \end{pmatrix}, \quad \bar{w}^t = \begin{pmatrix} \bar{w}_1^t \\ \bar{w}_2^t \end{pmatrix}, \quad \bar{\rho}_\ell^{t,t'} = \begin{pmatrix} \bar{\rho}_{\ell,1}^{t,t'} & \bar{\rho}_{\ell,2}^{t,t'} \\ 0 & 0 \end{pmatrix}, \quad \bar{D}_\ell^{t,t'} = \begin{pmatrix} \bar{D}_{\ell,1}^{t,t'} & \bar{D}_{\ell,2}^{t,t'} \\ 0 & 0 \end{pmatrix}, \\
 \bar{\rho}_r^{t,t'} &= \begin{pmatrix} \bar{\rho}_{r,1}^{t,t'} & \bar{\rho}_{r,2}^{t,t'} \\ 0 & 0 \end{pmatrix}, \quad \bar{D}_r^{t,t'} = \begin{pmatrix} \bar{D}_{r,1}^{t,t'} \\ 0 \end{pmatrix}, \\
 \bar{\Sigma}_\ell(t, t') &= \begin{pmatrix} \bar{\Sigma}_\ell^1(t, t') & 0 \\ 0 & 0 \end{pmatrix}, \quad \bar{R}_\ell(t, t') = \begin{pmatrix} \bar{R}_\ell^1(t, t') & \bar{R}_\ell^2(t, t') \\ 0 & 0 \end{pmatrix}, \quad \bar{\Gamma}(t) = \begin{pmatrix} \bar{\Gamma}_1(t) & \bar{\Gamma}_2(t) \\ 0 & 0 \end{pmatrix}.
 \end{aligned} \tag{301}$$

Here, we indicated by  $-$  the irrelevant variables. These variables satisfy

$$\begin{aligned}
 \bar{\theta}_1^t &= \theta^0 + \bar{U}_1^t - \int_0^t \left( h_s(\bar{\theta}_1^s) + \bar{\Gamma}_1(s) \bar{\theta}_1^s + \int_0^s \bar{R}_\ell^1(s, s') \bar{\theta}_1^{s'} ds' \right) ds \\
 &\quad - \int_0^t \left( \bar{\Gamma}_2(s) + \int_0^s \bar{R}_\ell^2(s, s') ds' \right) \theta^* ds, \quad \bar{U}_1 \sim \text{GP}(0, \bar{\Sigma}_\ell^{11} / \delta),
 \end{aligned} \tag{302}$$

$$\bar{\rho}_{\theta,1}^{t,t'} = I_m - \int_{t'}^t \left( (\nabla_\theta h_s(\bar{\theta}_1^s) + \bar{\Gamma}_1(s)) \bar{\rho}_{\theta,1}^{s,t'} + \int_{t'}^s \bar{R}_\ell^1(s, s') \bar{\rho}_{\theta,1}^{s',t'} ds' \right) ds, \tag{303}$$

$$\bar{C}_\theta^{11}(t, t') = \mathbb{E}[\bar{\theta}_1^t \bar{\theta}_1^{t'\top}], \quad \bar{C}_\theta^{12}(t, t') = \mathbb{E}[\bar{\theta}_1^t \theta^{*\top}], \quad \bar{C}_\theta^{22}(t, t') = \mathbb{E}[\theta^* \theta^{*\top}], \quad \bar{R}_\theta^1(t, t') = \mathbb{E}[\bar{\rho}_{\theta,1}^{t,t'}], \tag{304}$$

$$\bar{r}_1^t = \bar{w}_1^t - \frac{1}{\delta} \int_0^t \bar{R}_\theta^1(t, s) \ell_s(\bar{r}_1^s, \bar{r}_2^s; z) (ds + \sqrt{\tau\delta} dB^s), \quad \bar{r}_2^t = \bar{w}_2^t, \quad \begin{pmatrix} \bar{w}_1 \\ \bar{w}_2 \end{pmatrix} \sim \text{GP}(0, \bar{C}_\theta), \tag{305}$$

$$\bar{\rho}_{\ell,1}^{t,t'} = \nabla_r \ell_t(\bar{r}_1^t, \bar{r}_2^t; z) \bar{\rho}_{r,1}^{t,t'}, \quad \bar{\rho}_{\ell,2}^{t,t'} = \nabla_r \ell_t(\bar{r}_1^t, \bar{r}_2^t; z) \bar{\rho}_{r,2}^{t,t'}, \tag{306}$$

$$\bar{D}_{\ell,1}^{t,t'} = \nabla_r \ell_t(\bar{r}_1^t, \bar{r}_2^t; z) \left( -\frac{1}{\delta} \int_{t'}^t \bar{R}_\theta^1(t, s) \bar{D}_{\ell,1}^{s,t'} (ds + \sqrt{\tau\delta} dB^s) \right) + \nabla_{rr}^2 \ell_t(\bar{r}_1^t, \bar{r}_2^t; z) [\bar{D}_{r,1}^{t,t'}] \bar{\rho}_{r,1}^{t,t'}, \tag{307}$$

$$\bar{D}_{\ell,2}^{t,t'} = \nabla_r \ell_t(\bar{r}_1^t, \bar{r}_2^t; z) \left( -\frac{1}{\delta} \int_{t'}^t \bar{R}_\theta^1(t, s) \bar{D}_{\ell,2}^{s,t'} (ds + \sqrt{\tau\delta} dB^s) \right) + \nabla_{rr}^2 \ell_t(\bar{r}_1^t, \bar{r}_2^t; z) [\bar{D}_{r,1}^{t,t'}] \bar{\rho}_{r,2}^{t,t'}, \tag{308}$$

$$\bar{\rho}_{r,1}^{t,t'} = -\frac{1}{\delta} \int_{t'}^t \bar{R}_\theta^1(t, s) \bar{\rho}_{\ell,1}^{s,t'} (ds + \sqrt{\tau\delta} dB^s) - \frac{1}{\delta} \bar{R}_\theta^1(t, t') \nabla_r \ell_{t'}(\bar{r}_1^{t'}, \bar{r}_2^{t'}; z), \tag{309}$$

$$\bar{\rho}_{r,2}^{t,t'} = -\frac{1}{\delta} \int_{t'}^t \bar{R}_\theta^1(t, s) \bar{\rho}_{\ell,2}^{s,t'} (ds + \sqrt{\tau\delta} dB^s) - \frac{1}{\delta} \bar{R}_\theta^1(t, t') \nabla_w \ell_{t'}(\bar{r}_1^{t'}, \bar{r}_2^{t'}; z), \tag{310}$$

$$\bar{D}_{r,1}^{t,t'} = -\frac{1}{\delta} \int_{t'}^t \bar{R}_\theta^1(t, s) \nabla_r \ell_s(\bar{r}_1^s, \bar{r}_2^s; z) \bar{D}_{r,1}^{s,t'} (ds + \sqrt{\tau\delta} dB^s) - \sqrt{\frac{\tau}{\delta}} \bar{R}_\theta^1(t, t') \ell_{t'}(\bar{r}_1^{t'}, \bar{r}_2^{t'}; z), \tag{311}$$

$$\bar{\Sigma}_\ell^1(t, t') = \mathbb{E}[\bar{L}_1^t \bar{L}_1^{t'\top}], \quad \bar{L}_1^t := \int_0^t \ell_s(\bar{r}_1^s, \bar{r}_2^s; z) (ds + \sqrt{\tau\delta} dB^s), \quad (312)$$

$$\bar{R}_\ell^1(t, t') = \mathbb{E}[\bar{\rho}_{\ell,1}^{t,t'}] + \sqrt{\tau\delta} \mathbb{E}[\bar{D}_{\ell,1}^{t,t'}], \quad \bar{R}_\ell^2(t, t') = \mathbb{E}[\bar{\rho}_{\ell,2}^{t,t'}] + \sqrt{\tau\delta} \mathbb{E}[\bar{D}_{\ell,2}^{t,t'}], \quad (313)$$

$$\bar{\Gamma}_1(t) = \mathbb{E}[\nabla_r \ell_t(\bar{r}_1^t, \bar{r}_2^t; z)], \quad \bar{\Gamma}_2(t) = \mathbb{E}[\nabla_{r^*} \ell_t(\bar{r}_1^t, \bar{r}_2^t; z)]. \quad (314)$$

Since  $\bar{C}_\theta^{22}(t, t')$  is constant,  $\bar{r}_2^t = \bar{w}_2^t$  is constant. Then,  $\mathfrak{S}$  reduces to  $\mathfrak{S}^*$  by identifying

$$\begin{aligned} \bar{\theta}^t &= \theta^t, \quad \bar{U}_1^t = U^t, \quad \bar{\rho}_{\theta,1}^{t,t'} = \rho_\theta^{t,t'}, \\ \bar{C}_\theta^{11}(t, t') &= C_\theta(t, t'), \quad \bar{C}_\theta^{12}(t, t') = C_\theta(t, *), \quad \bar{C}_\theta^{22}(t, t') = C_\theta(*, *), \quad \bar{R}_\theta^1(t, t') = R_\theta(t, t'), \\ \bar{r}_1^t &= r^t, \quad \bar{w}_1^t = w^t, \quad \bar{r}_2^t = \bar{w}_2^t = w^*, \quad \bar{\rho}_{\ell,1}^{t,t'} = \rho_{\ell^*}^{t,t'}, \quad \bar{\rho}_{\ell,2}^{t,t'} = \rho_{\ell^*}^{t,t'}, \quad \bar{D}_{\ell,1}^{t,t'} = D_{\ell^*}^{t,t'}, \quad \bar{D}_{\ell,2}^{t,t'} = D_{\ell^*}^{t,t'}, \\ \bar{\rho}_{r,1}^{t,t'} &= \rho_r^{t,t'}, \quad \bar{\rho}_{r,2}^{t,t'} = \rho_{r^*}^{t,t'}, \quad \bar{D}_{r,1}^{t,t'} = D_r^{t,t'}, \\ \bar{\Sigma}_\ell^1(t, t') &= \Sigma_\ell(t, t'), \quad \bar{R}_\ell^1(t, t') = R_\ell(t, t'), \quad \bar{\Gamma}_1(t) = \Gamma(t), \quad \bar{\Gamma}_2(t) + \int_0^t \bar{R}_\ell^2(t, t') dt' = R_\ell(t, *). \end{aligned} \quad (315)$$

## E.2. Infinite Data Limit

### E.2.1. DERIVATION OF THE REDUCED DMFT EQUATION

We show that in the infinite data limit  $\delta \rightarrow \infty$ , the DMFT equations reduce to the following simple form.

$$\begin{aligned} d\theta^t &= -(h_t(\theta^t) + \Gamma(t)\theta^t) dt + \sqrt{\tau C_\ell(t)} dW^t, \quad C_\theta(t, t') = \mathbb{E}[\theta^t \theta^{t'\top}], \\ C_\ell(t) &= \mathbb{E}[\ell_t(w^t; z) \ell_t(w^t; z)^\top], \quad \Gamma(t) = \mathbb{E}[\nabla_r \ell_t(w^t; z)], \quad w \sim \text{GP}(0, C_\theta), \end{aligned} \quad (316)$$

where  $W^t$  is a Brownian motion in  $\mathbb{R}^m$ .

It is easy to see that as  $\delta \rightarrow \infty$ ,  $R_\ell(t, t') \rightarrow 0$  and hence  $R_\theta(t, t') \rightarrow I_m$ . Then, the equations for  $\theta^t$  and  $r^t$  reduce to

$$\theta^t = \theta^0 + U^t - \int_0^t (h_s(\theta^s) + \Gamma(s)\theta^s) ds, \quad U \sim \text{GP}(0, \tilde{\Sigma}_\ell), \quad r^t = w^t, \quad w \sim \text{GP}(0, C_\theta), \quad (317)$$

where we introduced the rescaled variable  $\tilde{\Sigma}_\ell := \Sigma_\ell/\delta$ . In addition, the equation for  $\tilde{\Sigma}_\ell$  reduces to

$$\tilde{\Sigma}_\ell(t, t') = \tau \int_0^{t \wedge t'} \mathbb{E}[\ell_s(w^s; z) \ell_s(w^s; z)^\top] ds. \quad (318)$$

Defining  $C_\ell(t) = \mathbb{E}[\ell_t(w^t; z) \ell_t(w^t; z)^\top]$ ,  $\tilde{\Sigma}_\ell$  is a covariance kernel of an integrated Brownian motion:

$$U^t = \int_0^t \sqrt{\tau C_\ell(s)} dW^s, \quad (319)$$

where  $W^t$  is a standard Brownian motion in  $\mathbb{R}^m$ . Thus,  $\theta^t$  follows the following SDE:

$$d\theta^t = -(h_t(\theta^t) + \Gamma(t)\theta^t) dt + \sqrt{\tau C_\ell(t)} dW^t. \quad (320)$$

These equations recover Equation (316).

For planted models of the form in Equation (281), the DMFT equation in the infinite data limit (316) becomes

$$\begin{aligned} d\theta^t &= -(h_t(\theta^t) + \Gamma(t)\theta^t + \Gamma^*(t)\theta^*) dt + \sqrt{\tau C_\ell(t)} dW^t, \\ C_\theta(t, t') &= \mathbb{E}[\theta^t \theta^{t'\top}] \quad (t, t' \in [0, T] \cup \{*\}), \quad C_\ell(t) = \mathbb{E}[\ell_t(w^t, w^*; z) \ell_t(w^t, w^*; z)^\top], \\ \Gamma(t) &= \mathbb{E}[\nabla_w \ell_t(w^t, w^*; z)], \quad \Gamma^*(t) = \mathbb{E}[\nabla_{w^*} \ell_t(w^t, w^*; z)], \quad w \sim \text{GP}(0, C_\theta). \end{aligned} \quad (321)$$

### E.2.2. EXAMPLE: LINEAR REGRESSION

Consider the linear regression setting described in Section 4. This corresponds to the choice  $m = 1$ ,  $h_t = 0$ , and  $\ell_t(r, r^*; z) = r - r^* - z$ . Then, Equation (321) reduces to

$$\begin{aligned} d\theta^t &= -(\theta^t - \theta^*) dt + \sqrt{\tau C_\ell(t)} dW^t, \\ C_\theta(t, t') &= \mathbb{E}[\theta^t \theta^{t'}] \quad (t, t' \in [0, T] \cup \{*\}), \quad C_\ell(t) = \mathbb{E}[(w^t - w^* - z)^2], \quad w \sim \text{GP}(0, C_\theta). \end{aligned} \quad (322)$$

Using the Fokker–Planck equation, the density of  $\theta^t$  given by the above SDE is given by the following partial differential equation (PDE):

$$\partial_t \mu(t, \theta) = \partial_\theta((\theta - \theta^*)\mu(t, \theta)) + \frac{\tau C_\ell(t)}{2} \partial_\theta^2 \mu(t, \theta). \quad (323)$$

This equation coincides with the PDE derived in Wang et al. [54].

Furthermore, the training dynamics for test errors  $\mathcal{L}(t) = \mathbb{E}[(\theta - \theta^*)^2] + \sigma^2$  can be obtained in closed form, where  $\sigma^2 := \mathbb{E}[z^2]$ . Let  $\rho^2 := \mathbb{E}[(\theta^*)^2]$ . Then, we have  $C_\ell(t) = C_\theta(t, t) - 2C_\theta(t, *) + \rho^2 + \sigma^2$ . By using Itô's lemma on  $C_\theta(t, *) = \mathbb{E}[\theta^t \theta^*]$  and  $C_\theta(t, t) = \mathbb{E}[(\theta^t)^2]$ , we obtain the following system of ODEs for  $C_\theta(t, *)$  and  $C_\theta(t, t)$ :

$$\frac{d}{dt} C_\theta(t, *) = -C_\theta(t, *) + \rho^2, \quad (324)$$

$$\frac{d}{dt} C_\theta(t, t) = -2C_\theta(t, t) + 2C_\theta(t, *) + \tau(C_\theta(t, t) - 2C_\theta(t, *) + \rho^2 + \sigma^2). \quad (325)$$

These are linear ODEs and can be solved in closed form as follows (assuming zero initialization, i.e.,  $\theta^0 = C_\theta(0, 0) = C_\theta(0, *) = 0$ ):

$$C_\theta(t, *) = \rho^2(1 - e^{-t}), \quad C_\theta(t, t) = \rho^2(1 - 2e^{-t} + e^{-(2-\tau)t}) + \frac{\tau\sigma^2}{2-\tau}(1 - e^{-(2-\tau)t}). \quad (326)$$

### E.3. Linear Regression

In this section, we derive the DMFT equation for the SGF for the ridge regression (regularized linear regression) problem. The loss for the ridge regression problem is defined as

$$\mathcal{L}(\theta) = \frac{1}{2n} \|\mathbf{X}\theta - \mathbf{y}\|_2^2 + \frac{\lambda}{2d} \|\theta\|_2^2, \quad \mathbf{y} = \mathbf{X}\theta^* + \mathbf{z}, \quad (327)$$

with regularization parameter  $\lambda \geq 0$ . The SGF for this problem is given by

$$d\theta^t = -\left(\lambda\theta^t + \frac{1}{\delta} \mathbf{X}^\top(\mathbf{r}^t - \mathbf{r}^* - \mathbf{z})\right) dt + \sqrt{\frac{\tau}{\delta}} \sum_{i=1}^n \mathbf{x}_i (r_i^t - r_i^* - z_i) dB_i^t. \quad (328)$$

where  $\mathbf{r}^t = \mathbf{X}\boldsymbol{\theta}^t$  and  $\mathbf{r}^* = \mathbf{X}\boldsymbol{\theta}^*$ . For simplicity, we consider the case of zero initialization  $\boldsymbol{\theta}^0 = 0$ . This is a special case of the SGF for planted models (281) with  $m = 1$ ,  $h_t(\boldsymbol{\theta}) = \lambda\boldsymbol{\theta}$ , and  $\ell_t(r, r^*; z) = r - r^* - z$ .

We are particularly interested in the training and test errors of the parameter  $\boldsymbol{\theta}^t$  given by

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \boldsymbol{\theta} - y_i)^2, \quad \mathcal{R}(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x}, y)} [(\mathbf{x}^\top \boldsymbol{\theta} - y)^2] = \frac{1}{d} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2^2 + \sigma^2. \quad (329)$$

Applying Corollary E.1, we obtain a DMFT characterization of the SGF for ridge regression (328). Note that since we have  $\nabla_r^2 \ell_t(r, r^*; z) = 0$ , the DMFT equation is guaranteed to have a unique solution for all  $T > 0$  by Theorem 3.1. Solving the DMFT equation, we can characterize the training and test errors of SGD for linear regression by a simple set of linear Volterra equations.

**Proposition E.2 (DMFT characterization of SGF for ridge regression)** *Assume that the noise  $z \in \mathbb{R}^n$  and the target parameter  $\boldsymbol{\theta}^* \in \mathbb{R}^d$  satisfy the same assumptions as in Corollary E.1. Let  $\rho^2 := \mathbb{E}[(\boldsymbol{\theta}^*)^2]$  and  $\sigma^2 := \mathbb{E}[z^2]$ . Let  $\mu_{\text{MP}}$  be the Marchenko–Pastur distribution with parameter  $\delta$  given by*

$$\mu_{\text{MP}}(x) = \frac{\delta \sqrt{(\lambda_+ - x)(x - \lambda_-)}}{2\pi x} + (1 - \delta)\delta(x)\mathbb{I}(\delta < 1), \quad x \in [\lambda_-, \lambda_+], \quad (330)$$

where  $\lambda_{\pm} = (1 \pm 1/\sqrt{\delta})^2$  and  $\delta$  is the Dirac delta function (an upright  $\delta$  is used to distinguish it from  $\delta = n/d$ ).

Then, for any  $0 \leq t_1, \dots, t_L < \infty$ , we have

$$\text{p-lim}_{n, d \rightarrow \infty} \max_{l=1, \dots, L} |\mathcal{L}(\boldsymbol{\theta}^{t_l}) - \mathcal{L}(t_l)| = 0, \quad \text{p-lim}_{n, d \rightarrow \infty} \max_{l=1, \dots, L} |\mathcal{R}(\boldsymbol{\theta}^{t_l}) - \mathcal{R}(t_l)| = 0, \quad (331)$$

where  $\mathcal{L}(t)$  and  $\mathcal{R}(t)$  solve the following system of linear Volterra equations:

$$\mathcal{L}(t) = \mathcal{L}_0(t) + \tau \int_0^t H_2(t-s)\mathcal{L}(s) \, ds, \quad \mathcal{R}(t) = \mathcal{R}_0(t) + \tau \int_0^t H_1(t-s)\mathcal{L}(s) \, ds, \quad (332)$$

where  $H_i(t) := \int x^i e^{-2(x+\lambda)t} \, d\mu_{\text{MP}}(x)$ , and  $\mathcal{L}_0(t)$  and  $\mathcal{R}_0(t)$  are the asymptotic train and test errors for the noiseless case  $\tau = 0$ , which are given by

$$\mathcal{L}_0(t) = \rho^2 \int \frac{x(\lambda + xe^{-(x+\lambda)t})^2}{(x+\lambda)^2} \, d\mu_{\text{MP}}(x) + \frac{\sigma^2}{\delta} \int \frac{(\lambda + xe^{-(x+\lambda)t})^2}{(x+\lambda)^2} \, d\mu_{\text{MP}}(x) + \frac{\delta - 1}{\delta} \sigma^2, \quad (333)$$

$$\mathcal{R}_0(t) = \rho^2 \int \frac{(\lambda + xe^{-(x+\lambda)t})^2}{(x+\lambda)^2} \, d\mu_{\text{MP}}(x) + \frac{\sigma^2}{\delta} \int \frac{x}{(x+\lambda)^2} (1 - e^{-(x+\lambda)t})^2 \, d\mu_{\text{MP}}(x) + \sigma^2. \quad (334)$$

Setting  $\lambda = 0$ , we obtain the DMFT characterization in the main text (Section 4).

The equation for  $\mathcal{L}$  in (332) is a scalar linear Volterra integral equation, which can be solved numerically efficiently. Once we obtain  $\mathcal{L}$ , we can compute  $\mathcal{R}$  using the second equation in (332).

The equations (332) are equivalent to those derived in Paquette et al. [42, 43] using the theory of homogenized SGD and random matrix theory (up to time rescaling to match their settings). Our result provides an alternative derivation of these equations using the DMFT framework (although our framework does not directly apply to SGD rigorously). For further analysis of these equations, such as their exact solutions and long-time behaviors, see Paquette et al. [42, 43].

## E.3.1. SIMPLIFYING THE DMFT EQUATIONS

We can derive the DMFT equation for ridge regression by specializing the DMFT equation  $\mathfrak{S}^*$  for planted models in Appendix E.1.2. Since the loss is quadratic, i.e.,  $\partial_r \ell_t(r, r^*; z) = 1$ ,  $\partial_{r^*} \ell_t(r, r^*; z) = -1$ , and  $\partial^2 \ell_t(r, r^*; z) = 0$ , the DMFT equation simplifies significantly. We have  $\Gamma(t) = 1$ ,  $D_\ell^{t,t'} = D_{\ell^*}^{t,t'} = 0$ ,  $\rho_r^{t,t'} = -\rho_{r^*}^{t,t'}$ , and  $\rho_\ell^{t,t'} = -\rho_{\ell^*}^{t,t'}$ . The DMFT equation for ridge regression reduces to

$$\theta^t = U^t - \int_0^t \left( (1 + \lambda)\theta^s - \theta^* + \int_0^s R_\ell(s, s')(\theta^{s'} - \theta^*) ds' \right) ds, \quad U \sim \text{GP}(0, \Sigma_\ell/\delta), \quad (335a)$$

$$\rho_\theta^{t,t'} = 1 - \int_{t'}^t \left( (1 + \lambda)\rho_\theta^{s,t'} + \int_{t'}^s R_\ell(s, s')\rho_\theta^{s',t'} ds' \right) ds, \quad (335b)$$

$$C_\theta(t, t') = \mathbb{E}[\theta^t \theta^{t'}] \quad (s, t \in [0, T] \cup \{*\}), \quad R_\theta(t, t') = \mathbb{E}[\rho_\theta^{t,t'}] \quad (t \geq t'), \quad (335c)$$

$$r^t = w^t - \frac{1}{\delta} \int_0^t R_\theta(t, s)(r^s - r^* - z)(ds + \sqrt{\tau\delta} dB^s), \quad w \sim \text{GP}(0, C_\theta), \quad (335d)$$

$$\rho_\ell^{t,t'} = -\frac{1}{\delta} \int_{t'}^t R_\theta(t, s)\rho_\ell^{s,t'}(ds + \sqrt{\tau\delta} dB^s) - \frac{1}{\delta} R_\theta(t, t'), \quad (335e)$$

$$\Sigma_\ell(t, t') = \mathbb{E}[L^t L^{t'}], \quad L^t := \int_0^t (r^s - r^* - z)(ds + \sqrt{\tau\delta} dB^s), \quad R_\ell(t, t') = \mathbb{E}[\rho_\ell^{t,t'}] \quad (t \geq t'). \quad (335f)$$

Next, we eliminate the stochastic processes to close the system in terms of the correlation and response functions. To this end, we discretize time with step size  $\gamma$  as in Appendix B.2, reduce the equations, and then take the continuous-time limit  $\gamma \rightarrow 0$ . This operation can be justified along the lines of the proof of Theorem 3.2.

Let  $t_i = i\gamma$  for  $i = 0, 1, \dots, \lfloor T/\gamma \rfloor$ . The discretized version of Equation (335) is given by

$$\theta^{t_i} = \gamma \sum_{j=0}^{i-1} \left( u^{t_j} - (1 + \lambda)\theta^{t_j} + \theta^* - \gamma \sum_{k=0}^{j-1} R_\ell(t_j, t_k)(\theta^{t_k} - \theta^*) \right), \quad U \sim \text{GP}(0, C_\ell/\delta), \quad (336a)$$

$$C_\theta(t_i, t_j) = \mathbb{E}[\theta^{t_i} \theta^{t_j}], \quad C_\theta(t_i, *) = \mathbb{E}[\theta^{t_i} \theta^*], \quad R_\theta(t_i, t_j) = \gamma^{-1} \mathbb{E} \left[ \frac{\partial \theta^{t_i}}{\partial u^{t_j}} \right] \quad (i > j), \quad (336b)$$

$$r^{t_i} = w^{t_i} - \frac{1}{\delta} \sum_{j=0}^{i-1} R_\theta(t_i, t_j)(r^{t_j} - r^* - z)(\gamma + \sqrt{\tau\delta}(B^{t_{j+1}} - B^{t_j})), \quad w \sim \text{GP}(0, C_\theta), \quad (336c)$$

$$C_\ell(t_i, t_j) = \mathbb{E} \left[ (r^{t_i} - r^* - z) \left( 1 + \sqrt{\tau\delta} \frac{B^{t_{i+1}} - B^{t_i}}{\gamma} \right) (r^{t_j} - r^* - z) \left( 1 + \sqrt{\tau\delta} \frac{B^{t_{j+1}} - B^{t_j}}{\gamma} \right) \right], \quad (336d)$$

$$R_\ell(t_i, t_j) = \gamma^{-1} \mathbb{E} \left[ \frac{\partial r^{t_i}}{\partial w^{t_j}} \right] \quad (i > j). \quad (336e)$$

Let  $K = \lfloor T/\gamma \rfloor + 1$  and let  $\mathbf{C}_\theta, \mathbf{R}_\theta, \mathbf{C}_\ell, \mathbf{R}_\ell \in \mathbb{R}^{K \times K}$  be the matrices with entries  $(\mathbf{C}_\theta)_{ij} = C_\theta(t_i, t_j)$ ,  $(\mathbf{R}_\theta)_{ij} = R_\theta(t_i, t_j)$ ,  $(\mathbf{C}_\ell)_{ij} = C_\ell(t_i, t_j)$ , and  $(\mathbf{R}_\ell)_{ij} = R_\ell(t_i, t_j)$  for  $i, j = 0, 1, \dots, K-1$ . Let  $\mathbf{c}^* \in \mathbb{R}^K$  be the vector with entries  $c_i^* = C_\theta(t_i, *)$  for  $i = 0, 1, \dots, K-1$ . Let  $\boldsymbol{\theta}, \mathbf{r}, \mathbf{u}, \mathbf{w}, \mathbf{G} \in \mathbb{R}^K$  be the vectors with entries  $\theta^{t_i}, r^{t_i}, u^{t_i}, w^{t_i}, G^{t_i} := (B^{t_{i+1}} - B^{t_i})/\sqrt{\gamma}$  for  $i = 0, 1, \dots, K-1$ .

**Equation for  $R_\theta$ .** Differentiating Equation (336a) with respect to  $u^{t_j}$  for  $j < i$  and taking the expectation, we have

$$R_\theta(t_i, t_j) = 1 - \gamma \sum_{k=j}^{i-1} \left( (1 + \lambda) R_\theta(t_k, t_j) + \gamma \sum_{l=0}^{k-1} R_\ell(t_k, t_l) R_\theta(t_l, t_j) \right). \quad (337)$$

Let  $\mathbf{T} := \mathbb{R}^{K \times K}$  be the lower-triangular matrix with entries  $T_{ij} = 1$  for  $i > j$  and  $T_{ij} = 0$  for  $i \geq j$ . Equation (337) can be written in matrix form as

$$\mathbf{R}_\theta = \mathbf{T} - \gamma \mathbf{T}((1 + \lambda)\mathbf{R}_\theta + \gamma \mathbf{R}_\ell \mathbf{R}_\theta). \quad (338)$$

Solving for  $\mathbf{R}_\theta$ , we obtain

$$\mathbf{R}_\theta = (\mathbf{I} + \gamma \mathbf{T}((1 + \lambda)\mathbf{I} + \gamma \mathbf{R}_\ell))^{-1} \mathbf{T}. \quad (339)$$

**Equation for  $C_\theta$ .** Multiplying Equation (336a) by  $\theta^*$  and taking the expectation, we have

$$C_\theta(t_i, *) = -\gamma \sum_{j=0}^{i-1} \left( (1 + \lambda) C_\theta(t_j, *) - \rho^2 + \gamma \sum_{k=0}^{j-1} R_\ell(t_j, t_k) (C_\theta(t_k, *) - \rho^2) \right). \quad (340)$$

In matrix form, this can be written as

$$\mathbf{c}^* = -\gamma \mathbf{T}((1 + \lambda)\mathbf{c}^* - \rho^2 \mathbf{1} + \gamma \mathbf{R}_\ell(\mathbf{c}^* - \rho^2 \mathbf{1})). \quad (341)$$

Solving for  $\mathbf{c}^*$ , we obtain

$$\mathbf{c}^* = \rho^2 \gamma (\mathbf{I} + \gamma \mathbf{T}((1 + \lambda)\mathbf{I} + \gamma \mathbf{R}_\ell))^{-1} \mathbf{T}(\mathbf{I} + \gamma \mathbf{R}_\ell) \mathbf{1} = \rho^2 \gamma \mathbf{R}_\theta (\mathbf{I} + \gamma \mathbf{R}_\ell) \mathbf{1}. \quad (342)$$

Multiplying Equation (336a) by  $\theta^{t_j}$  and taking the expectation, we have

$$\begin{aligned} & C_\theta(t_i, t_j) \\ &= \gamma \sum_{k=0}^{i-1} \left( \mathbb{E}[u^{t_k} \theta^{t_j}] - (1 + \lambda) C_\theta(t_k, t_j) + C_\theta(t_j, *) - \gamma \sum_{l=0}^{j-1} R_\ell(t_k, t_l) (C_\theta(t_l, t_j) - C_\theta(t_j, *)) \right). \end{aligned} \quad (343)$$

By Stein's lemma (the Gaussian integration by parts), we have

$$\mathbb{E}[u^{t_k} \theta^{t_j}] = \sum_{l=0}^{j-1} \text{Cov}(u^{t_k}, u^{t_l}) \mathbb{E} \left[ \frac{\partial \theta^{t_j}}{\partial u^{t_l}} \right] = \frac{\gamma}{\delta} \sum_{l=0}^{j-1} C_\ell(t_k, t_l) R_\theta(t_j, t_l). \quad (344)$$

Therefore, Equation (343) can be written in matrix form as

$$\mathbf{C}_\theta = \gamma \mathbf{T} \left( \frac{\gamma}{\delta} \mathbf{C}_\ell \mathbf{R}_\theta^\top - (1 + \lambda) \mathbf{C}_\theta + \mathbf{1} \mathbf{c}^{*\top} - \gamma \mathbf{R}_\ell (\mathbf{C}_\theta - \mathbf{1} \mathbf{c}^{*\top}) \right). \quad (345)$$

Solving for  $\mathbf{C}_\theta$ , we obtain

$$\begin{aligned} \mathbf{C}_\theta &= \gamma (\mathbf{I} + \gamma \mathbf{T} ((1 + \lambda) \mathbf{I} + \gamma \mathbf{R}_\ell))^{-1} \mathbf{T} (\gamma \mathbf{C}_\ell \mathbf{R}_\theta^\top + (\mathbf{I} + \gamma \mathbf{R}_\ell) \mathbf{1} \mathbf{c}^{*\top}) \\ &= \frac{\gamma^2}{\delta} \mathbf{R}_\theta \mathbf{C}_\ell \mathbf{R}_\theta^\top + \frac{1}{\rho^2} \mathbf{c}^* \mathbf{c}^{*\top}. \end{aligned} \quad (346)$$

Let  $C_\theta^*(t_i, t_j) := \mathbb{E}[(\theta^{t_i} - \theta^*)(\theta^{t_j} - \theta^*)] = C_\theta(t_i, t_j) - C_\theta(t_i, *) - C_\theta(t_j, *) + \rho^2$ . Let  $\mathbf{C}_\theta^* \in \mathbb{R}^{K \times K}$  be the matrix with entries  $(\mathbf{C}_\theta^*)_{ij} = C_\theta^*(t_i, t_j)$  for  $i, j = 0, 1, \dots, K - 1$ . From Equation (346), we have

$$\mathbf{C}_\theta^* = \mathbf{C}_\theta - \mathbf{c}^* \mathbf{1}^\top - \mathbf{1} \mathbf{c}^{*\top} + \rho^2 \mathbf{1} \mathbf{1}^\top = \frac{\gamma^2}{\delta} \mathbf{R}_\theta \mathbf{C}_\ell \mathbf{R}_\theta^\top + \frac{1}{\rho^2} (\mathbf{c}^* - \rho^2 \mathbf{1})(\mathbf{c}^* - \rho^2 \mathbf{1})^\top. \quad (347)$$

**Equation for  $\mathbf{R}_\ell$ .** Differentiating Equation (336c) with respect to  $w^{t_j}$  for  $j < i$  and taking the expectation, we have

$$R_\ell(t_i, t_j) = -\frac{1}{\delta} \sum_{k=j}^{i-1} R_\theta(t_i, t_k) R_\ell(t_k, t_j) - \frac{1}{\delta} R_\theta(t_i, t_j). \quad (348)$$

This can be written in matrix form as

$$\mathbf{R}_\ell = -\frac{1}{\delta} \mathbf{R}_\theta \mathbf{R}_\ell - \frac{1}{\delta} \mathbf{R}_\theta. \quad (349)$$

Solving for  $\mathbf{R}_\ell$ , we obtain

$$\mathbf{R}_\ell = -(\delta \mathbf{I} + \gamma \mathbf{R}_\theta)^{-1} \mathbf{R}_\theta. \quad (350)$$

**Equation for  $\mathbf{C}_\ell$ .** Since  $G^{t_i} \sim \mathcal{N}(0, 1)$  are i.i.d. standard normal variables, we have

$$\begin{aligned} C_\ell(t_i, t_j) &= \mathbb{E}[(r^{t_i} - r^* - z)(r^{t_j} - r^* - z)] + \frac{\tau \delta}{\gamma} \mathbb{E}[(r^{t_i} - r^* - z)^2] \mathbb{I}(i = j) \\ &\quad + \sqrt{\frac{\tau \delta}{\gamma}} (\mathbb{E}[(r^{t_i} - r^* - z)(r^{t_j} - r^* - z) G^{t_j}] + \mathbb{E}[(r^{t_i} - r^* - z)(r^{t_j} - r^* - z) G^{t_i}]) \end{aligned} \quad (351)$$

$$= L(t_i, t_j) + \frac{\tau \delta}{\gamma} L(t_i, t_i) \mathbb{I}(i = j) + \sqrt{\tau \delta} (M(t_i, t_j) + M(t_j, t_i)), \quad (352)$$

where we defined

$$L(t_i, t_j) = \mathbb{E}[(r^{t_i} - r^* - z)(r^{t_j} - r^* - z)], \quad (353)$$

$$M(t_i, t_j) = \gamma^{-1/2} \mathbb{E}[(r^{t_i} - r^* - z)(r^{t_j} - r^* - z) G^{t_j}] = \gamma^{-1/2} \mathbb{E} \left[ \frac{\partial r^{t_i}}{\partial G^{t_j}} (r^{t_j} - r^* - z) \right]. \quad (354)$$

In the definition of  $M(t_i, t_j)$ , we used Stein's lemma.

We first derive a closed equation for  $M(t_i, t_j)$ . Differentiating Equation (336c) with respect to  $G^{t_j}$ , multiplying by  $\gamma^{-1/2}(r^{t_j} - r^* - z)$ , and taking the expectation, we have

$$M(t_i, t_j) = -\frac{\gamma}{\delta} \sum_{k=j}^{i-1} R_\theta(t_i, t_k) M(t_k, t_j) - \sqrt{\frac{\tau}{\delta}} R_\theta(t_i, t_j) L(t_j, t_j). \quad (355)$$

In matrix form, this can be written as

$$\mathbf{M} = -\frac{\gamma}{\delta} \mathbf{R}_\theta \mathbf{M} - \sqrt{\frac{\tau}{\delta}} \mathbf{R}_\theta \text{diag}(\mathbf{L}), \quad (356)$$

where  $\mathbf{M}, \mathbf{L} \in \mathbb{R}^{K \times K}$  are the matrices with entries  $M_{ij} = M(t_i, t_j)$  and  $L_{ij} = L(t_i, t_j)$ , and  $\text{diag}(\mathbf{L})$  is the diagonal matrix with diagonal entries equal to those of  $\mathbf{L}$ . Solving for  $\mathbf{M}$ , we obtain

$$\mathbf{M} = -\sqrt{\tau\delta}(\delta\mathbf{I} + \gamma\mathbf{R}_\theta)^{-1} \mathbf{R}_\theta \text{diag}(\mathbf{L}) = \sqrt{\tau\delta} \mathbf{R}_\ell \text{diag}(\mathbf{L}). \quad (357)$$

Next, we derive a closed equation for  $L(t_i, t_j)$ . Multiplying Equation (336c) by  $r^{t_j} - r^* - z$  and taking the expectation, we have

$$L(t_i, t_j) = \mathbb{E}[(w^{t_i} - r^* - z)(r^{t_j} - r^* - z)] - \frac{\gamma}{\delta} \sum_{k=0}^{i-1} R_\theta(t_i, t_k) (L(t_k, t_j) + \sqrt{\tau\delta} M(t_j, t_k)). \quad (358)$$

By Stein's lemma, we have

$$\begin{aligned} & \mathbb{E}[(w^{t_i} - r^* - z)(r^{t_j} - r^* - z)] \\ &= \sum_{k=0}^j \text{Cov}(w^{t_i} - r^* - z, w^{t_k}) \mathbb{E}\left[\frac{\partial r^{t_j}}{\partial w^{t_k}}\right] + \text{Cov}(w^{t_i} - r^* - z, r^* + z) \left(\mathbb{E}\left[\frac{\partial r^{t_j}}{\partial (r^* + z)}\right] - 1\right) \\ &= \gamma \sum_{k=0}^{j-1} (C_\theta(t_i, t_k) - C_\theta(t_k, *)) R_\ell(t_j, t_k) + C_\theta(t_i, t_j) - C_\theta(t_j, *) \\ & \quad - (C_\theta(t_i, *) - \rho^2 - \sigma^2) \left(\gamma \sum_{k=0}^{j-1} R_\ell(t_j, t_k) + 1\right) \\ &= C_\theta(t_i, t_j) - C_\theta(t_i, *) - C_\theta(t_j, *) + \rho^2 + \sigma^2 \\ & \quad + \gamma \sum_{k=0}^{j-1} (C_\theta(t_i, t_k) - C_\theta(t_i, *) - C_\theta(t_k, *) + \rho^2 + \sigma^2) R_\ell(t_j, t_k) \\ &= C_\theta^*(t_i, t_j) + \sigma^2 + \gamma \sum_{k=0}^{j-1} (C_\theta^*(t_i, t_k) + \sigma^2) R_\ell(t_j, t_k). \end{aligned} \quad (359)$$

Therefore, Equation (358) can be written in matrix form as

$$\mathbf{L} = (\mathbf{C}_\theta^* + \sigma^2 \mathbf{1}\mathbf{1}^\top)(\mathbf{I} + \gamma \mathbf{R}_\ell^\top) - \frac{\gamma}{\delta} \mathbf{R}_\theta (\mathbf{L} + \sqrt{\tau\delta} \mathbf{M}^\top). \quad (360)$$

Solving for  $\mathbf{L}$ , we obtain

$$\begin{aligned} \mathbf{L} &= \delta(\delta\mathbf{I} + \gamma\mathbf{R}_\theta)^{-1}(\mathbf{C}_\theta^* + \sigma^2\mathbf{1}\mathbf{1}^\top)(\mathbf{I} + \gamma\mathbf{R}_\ell^\top) - \sqrt{\tau\delta}\gamma(\delta\mathbf{I} + \gamma\mathbf{R}_\theta)^{-1}\mathbf{R}_\theta\mathbf{M}^\top \\ &= (\mathbf{I} + \gamma\mathbf{R}_\ell)(\mathbf{C}_\theta^* + \sigma^2\mathbf{1}\mathbf{1}^\top)(\mathbf{I} + \gamma\mathbf{R}_\ell^\top) + \tau\delta\gamma\mathbf{R}_\ell \text{diag}(\mathbf{L})\mathbf{R}_\ell^\top, \end{aligned} \quad (361)$$

where we used that

$$\begin{aligned} \mathbf{I} + \gamma\mathbf{R}_\ell &= \mathbf{I} - \gamma(\delta\mathbf{I} + \gamma\mathbf{R}_\theta)^{-1}\mathbf{R}_\theta = \mathbf{I} - (\delta\mathbf{I} + \gamma\mathbf{R}_\theta)^{-1}(\delta\mathbf{I} + \gamma\mathbf{R}_\theta - \delta\mathbf{I}) \\ &= \delta(\delta\mathbf{I} + \gamma\mathbf{R}_\theta)^{-1}. \end{aligned} \quad (362)$$

**Continuous-time limit.** Taking the continuous-time limit  $\gamma \rightarrow 0$  in Equations (337), (342), and (347), we obtain

$$R_\theta(t, t') = 1 - \int_{t'}^t \left( (1 + \lambda)R_\theta(s, t') + \int_{t'}^s R_\ell(s, s')R_\theta(s', t') ds' \right) ds, \quad (363)$$

$$C_\theta(t, *) = \rho^2 \int_0^t R_\theta(t, s) \left( 1 + \int_0^s R_\ell(s, s') ds' \right) ds, \quad (364)$$

$$C_\theta^*(t, t') = \frac{1}{\delta} \int_0^t \int_0^{t'} R_\theta(t, s) R_\theta(t', s') C_\ell(s, s') ds' ds + \frac{1}{\rho^2} (C_\theta(t, *) - \rho^2)(C_\theta(t', *) - \rho^2). \quad (365)$$

Taking the continuous-time limit  $\gamma \rightarrow 0$  in Equations (348), (352), and (361), we obtain

$$R_\ell(t, t') = -\frac{1}{\delta} \int_{t'}^t R_\theta(t, s) R_\ell(s, t') ds - \frac{1}{\delta} R_\theta(t, t'), \quad (366)$$

$$C_\ell(t, t') = L(t, t') + \tau\delta(R_\ell(t', t)L(t', t') + R_\ell(t, t')L(t, t) + L(t, t)\delta(t - t')), \quad (367)$$

$$\begin{aligned} L(t, t') &= \int_0^t \int_0^{t'} (\delta(t - s) + R_\ell(t, s))(\delta(t' - s') + R_\ell(t', s'))(C_\theta^*(s, s') + \sigma^2) ds' ds \\ &\quad + \tau\delta \int_0^{t \wedge t'} R_\ell(t, s) R_\ell(t', s) L(s, s) ds. \end{aligned} \quad (368)$$

Here,  $\delta(\cdot)$  is the Dirac delta function.

### E.3.2. SOLVING THE DMFT EQUATIONS

We now solve the DMFT equations derived in the previous section for ridge regression.

**Lemma E.3** *Define the function  $K_i(t)$  for  $i \geq 0$  as*

$$K_i(t) := \int x^i e^{-(x+\lambda)t} d\mu_{\text{MP}}(x), \quad (369)$$

where  $\mu_{\text{MP}}$  is the Marchenko–Pastur law, whose density is given by Equation (330).

The solution of the DMFT equations for ridge regression (363)–(368) is given by

$$R_\theta(t, t') = K_0(t - t'), \quad R_\ell(t, t') = -\frac{1}{\delta} K_1(t - t'), \quad (370)$$

$$C_{\theta}^*(t, t') = C_{\theta 0}^*(t, t') + \tau \int_0^{t \wedge t'} K_1(t + t' - 2s)L(s, s) ds, \quad (371)$$

$$L(t, t') = L_0(t, t') + \tau \int_0^{t \wedge t'} K_2(t + t' - 2s)L(s, s) ds, \quad (372)$$

where  $C_{\theta 0}^*$  and  $L_0$  are given by

$$\begin{aligned} C_{\theta 0}^*(t, t') &= \rho^2 \int \frac{(\lambda + xe^{-(x+\lambda)t})(\lambda + xe^{-(x+\lambda)t'})}{(x + \lambda)^2} d\mu_{\text{MP}}(x) \\ &\quad + \frac{\sigma^2}{\delta} \int \frac{x}{(x + \lambda)^2} (1 - e^{-(x+\lambda)t})(1 - e^{-(x+\lambda)t'}) d\mu_{\text{MP}}(x), \end{aligned} \quad (373)$$

$$\begin{aligned} L_0(t, t') &= \rho^2 \int x \frac{(\lambda + xe^{-(x+\lambda)t})(\lambda + xe^{-(x+\lambda)t'})}{(x + \lambda)^2} d\mu_{\text{MP}}(x) \\ &\quad + \frac{\sigma^2}{\delta} \int \frac{(\lambda + xe^{-(x+\lambda)t})(\lambda + xe^{-(x+\lambda)t'})}{(x + \lambda)^2} d\mu_{\text{MP}}(x) + \frac{\delta - 1}{\delta} \sigma^2. \end{aligned} \quad (374)$$

Before proving Lemma E.3, we summarize the necessary background on the *Laplace transform*, which is a useful technique for analyzing linear differential equations and will be used extensively in the proof. Given a function  $f: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$ , its Laplace transform  $\mathcal{L}[f] = \bar{f}$  is defined as

$$\bar{f}(p) := \int_0^{\infty} f(t)e^{-pt} dt, \quad (375)$$

for  $p \in \mathbb{C}$  with sufficiently large real part for the integral to be convergent.

We state several of its basic properties used in the proof.

- *Linearity:* For  $f: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$  and  $a, b \in \mathbb{R}$ , we have  $\mathcal{L}[af + b] = a\mathcal{L}[f] + b$ .
- *Laplace transforms of derivatives, integrals, and convolutions:* For  $f, g: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$ , we have

$$\mathcal{L}[f'(t)](p) = p\bar{f}(p) - f(0), \quad (376)$$

$$\mathcal{L}\left[\int_0^t f(s) ds\right](p) = \frac{\bar{f}(p)}{p}, \quad (377)$$

$$\mathcal{L}\left[\int_0^t f(t-s)g(s) ds\right](p) = \bar{f}(p)\bar{g}(p). \quad (378)$$

- *Laplace transform of the Dirac delta function:* We have  $\mathcal{L}[\delta(t)](p) = 1$ , where  $\delta(t)$  is the Dirac delta function.

We also utilize a two-dimensional version of the Laplace transform, which is defined for  $g: \mathbb{R}_{\geq 0}^2 \rightarrow \mathbb{R}$  as

$$\bar{g}(p, q) := \int_0^{\infty} \int_0^{\infty} g(t, t')e^{-pt}e^{-qt'} dt dt', \quad (379)$$

for  $p, q \in \mathbb{C}$  with sufficiently large real parts for the integral to be convergent. Its properties are similar to the one-dimensional case.

**Proof** We proceed as follows. First, we solve the equation for the noiseless case  $\tau = 0$  in the frequency domain using the Laplace transform. Next, we perform the inverse Laplace transform to obtain the time-domain solution for  $\tau = 0$ . Along the way, we use ideas and techniques from the random matrix theory. Finally, we solve the full equations for  $\tau > 0$ .

**Solving in the frequency domain for  $\tau = 0$ .** Note that the equations for  $R_\theta$ ,  $R_\ell$ , and  $C_\theta(\cdot, *)$  do not depend on  $\tau$  and are thus the same for  $\tau = 0$  and  $\tau > 0$ . Let  $C_{\theta 0}^*$  and  $L_0$  be the solutions of Equations (365) and (368) for  $\tau = 0$ .

Since the equations for  $R_\theta(t, t')$  (363) and  $R_\ell(t, t')$  (366) depend on time only through the time difference  $t - t'$ , they are time-translation invariant, i.e.,  $R_\theta(t, t') = R_\theta(t - t')$  and  $R_\ell(t, t') = R_\ell(t - t')$ . Therefore, they satisfy the following one-dimensional integral equations:

$$R_\theta(t) = 1 - \int_0^t \left( (1 + \lambda)R_\theta(s) + \int_0^s R_\ell(s - s')R_\theta(s') ds' \right) ds, \quad (380)$$

$$R_\ell(t) = -\frac{1}{\delta}R_\theta(t) - \frac{1}{\delta} \int_0^t R_\theta(t - s)R_\ell(s) ds. \quad (381)$$

Taking the Laplace transforms of these equations, we have

$$\bar{R}_\theta(p) = \frac{1}{p} \left( 1 - (1 + \lambda)\bar{R}_\theta(p) - \bar{R}_\ell(p)\bar{R}_\theta(p) \right), \quad (382)$$

$$\bar{R}_\ell(p) = -\frac{1}{\delta}\bar{R}_\theta(p) - \frac{1}{\delta}\bar{R}_\theta(p)\bar{R}_\ell(p). \quad (383)$$

Rearranging, we obtain

$$\bar{R}_\theta(p) = -\frac{\delta(1 + p + \lambda) - 1 - \sqrt{(\delta(1 + p + \lambda) - 1)^2 + 4\delta(p + \lambda)}}{2(p + \lambda)}, \quad (384)$$

$$\bar{R}_\ell(p) = -\frac{\bar{R}_\theta(p)}{\delta + \bar{R}_\theta(p)} = \frac{(p + \lambda)\bar{R}_\theta(p) - 1}{\delta}. \quad (385)$$

Taking the Laplace transforms of the correlation functions in Equations (364), (365), and (368), we have

$$\bar{C}_\theta(p, *) = \rho^2 \bar{R}_\theta(p) \cdot \frac{1 + \bar{R}_\ell(p)}{p}, \quad (386)$$

$$\bar{C}_{\theta 0}^*(p, q) = \frac{1}{\delta} \bar{R}_\theta(p) \bar{R}_\theta(q) \bar{L}_0(p, q) + \frac{1}{\rho^2} \left( \bar{C}_\theta(p, *) - \frac{\rho^2}{p} \right) \left( \bar{C}_\theta(q, *) - \frac{\rho^2}{q} \right), \quad (387)$$

$$\bar{L}_0(p, q) = (1 + \bar{R}_\ell(p))(1 + \bar{R}_\ell(q)) \left( \bar{C}_{\theta 0}^*(p, q) + \frac{\sigma^2}{pq} \right). \quad (388)$$

Simplifying these equations, we obtain

$$\bar{C}_{\theta 0}^*(p, q) = \frac{\rho^2(p + \lambda)(q + \lambda)}{pq} \frac{\bar{R}_\theta(p)\bar{R}_\theta(q)}{1 - \delta\bar{R}_\ell(p)\bar{R}_\ell(q)} + \frac{\sigma^2}{pq} \frac{\delta\bar{R}_\ell(p)\bar{R}_\ell(q)}{1 - \delta\bar{R}_\ell(p)\bar{R}_\ell(q)}, \quad (389)$$

$$\begin{aligned} \bar{L}_0(p, q) &= \frac{\rho^2\delta(p + \lambda)(q + \lambda)}{pq} \frac{\delta\bar{R}_\ell(p)\bar{R}_\ell(q)}{1 - \delta\bar{R}_\ell(p)\bar{R}_\ell(q)} \\ &\quad + \frac{\sigma^2(p + \lambda)(q + \lambda)}{\delta pq} \frac{\bar{R}_\theta(p)\bar{R}_\theta(q)}{1 - \delta\bar{R}_\ell(p)\bar{R}_\ell(q)} + \frac{\delta - 1}{\delta} \frac{\sigma^2}{pq}. \end{aligned} \quad (390)$$

Using the relation

$$1 - \delta\bar{R}_\ell(p)\bar{R}_\ell(q) = \frac{(q - p)\bar{R}_\theta(p)\bar{R}_\theta(q)}{\bar{R}_\theta(p) - \bar{R}_\theta(q)}, \quad (391)$$

we have

$$\frac{\overline{R}_\theta(p)\overline{R}_\theta(q)}{1 - \delta\overline{R}_\ell(p)\overline{R}_\ell(q)} = \frac{\overline{R}_\theta(p) - \overline{R}_\theta(q)}{q - p}, \quad \frac{\delta\overline{R}_\ell(p)\overline{R}_\ell(q)}{1 - \delta\overline{R}_\ell(p)\overline{R}_\ell(q)} = \frac{\overline{R}_\ell(q) - \overline{R}_\ell(p)}{q - p}. \quad (392)$$

Thus, we obtain

$$\overline{C}_{\theta 0}^*(p, q) = \frac{\rho^2(p + \lambda)(q + \lambda)}{pq} \frac{\overline{R}_\theta(p) - \overline{R}_\theta(q)}{q - p} + \frac{\sigma^2}{pq} \frac{\overline{R}_\ell(q) - \overline{R}_\ell(p)}{q - p}, \quad (393)$$

$$\overline{L}_0(p, q) = \frac{\rho^2\delta(p + \lambda)(q + \lambda)}{pq} \frac{\overline{R}_\ell(q) - \overline{R}_\ell(p)}{q - p} + \frac{\sigma^2(p + \lambda)(q + \lambda)}{\delta pq} \frac{\overline{R}_\theta(p) - \overline{R}_\theta(q)}{q - p} + \frac{\delta - 1}{\delta} \frac{\sigma^2}{pq}. \quad (394)$$

**Solving in the time domain for  $\tau = 0$ .** We now perform the inverse Laplace transform to obtain the time-domain solution. Before proceeding, we introduce *Stieltjes transform*. The Stieltjes transform  $S : \mathbb{C} \setminus I \rightarrow \mathbb{C}$  of a (signed) measure  $\mu$  on an interval  $I \subseteq \mathbb{R}$  is defined as follows:

$$S(z) = \int_I \frac{1}{x - z} d\mu(x). \quad (395)$$

We check that the time-domain solution stated in Lemma E.3 has the same Laplace transform as the frequency-domain solution obtained above. First, we take the Laplace transform of  $R_\theta$  to obtain

$$\begin{aligned} \int_0^\infty R_\theta(t) e^{-pt} dt &= \int_0^\infty \left( \int e^{-(x+\lambda+p)t} d\mu_{\text{MP}}(x) \right) dt = \int \left( \int_0^\infty e^{-(x+\lambda+p)t} dt \right) d\mu_{\text{MP}}(x) \\ &= \int \frac{1}{x + \lambda + p} d\mu_{\text{MP}}(x) = S_{\text{MP}}(-(p + \lambda)). \end{aligned} \quad (396)$$

Here,  $S_{\text{MP}}(z)$  is the Stieltjes transform of the Marchenko–Pastur law which is given by

$$S_{\text{MP}}(z) = \frac{\delta(1 - z) - 1 - \sqrt{(\delta(1 - z) - 1)^2 - 4\delta z}}{2z}. \quad (397)$$

Setting  $z = -(p + \lambda)$ , we obtain  $\overline{R}_\theta(p)$  given in Equation (384).

Next, we check  $R_\ell$ .

$$\int_0^\infty R_\ell(t) e^{-pt} dt = -\frac{1}{\delta} \int \frac{x}{x + p + \lambda} d\mu_{\text{MP}}(x) = -\frac{1}{\delta} (-(p + \lambda)S_{\text{MP}}(-(p + \lambda)) + 1), \quad (398)$$

which is equal to  $\overline{R}_\ell(p)$  in Equation (385).

Finally, we check  $C_{\theta 0}^*$  and  $L_0$ . Define  $F_1(t, t')$ ,  $F_2(t, t')$  as follows:

$$F_1(t, t') = K_0(t + t'), \quad F_2(t, t') = \frac{1}{\delta} K_1(t + t'). \quad (399)$$

The Laplace transforms of these functions are

$$\int_0^\infty \int_0^\infty F_1(t, t') e^{-pt - qt'} dt dt' = \int \frac{1}{(x + p + \lambda)(x + q + \lambda)} d\mu_{\text{MP}}(x) = \frac{\overline{R}_\theta(p) - \overline{R}_\theta(q)}{q - p}, \quad (400)$$

$$\int_0^\infty \int_0^\infty F_2(t, t') e^{-pt - qt'} dt dt' = \frac{1}{\delta} \int \frac{x}{(x+p+\lambda)(x+q+\lambda)} d\mu_{\text{MP}}(x) = \frac{\bar{R}_\ell(q) - \bar{R}_\ell(p)}{q-p}. \quad (401)$$

Furthermore, we have

$$\begin{aligned} \mathcal{L}^{-1} \left[ \frac{(p+\lambda)(q+\lambda)}{pq} \frac{\bar{R}_\theta(p) - \bar{R}_\theta(q)}{q-p} \right] &= \mathcal{L}^{-1} \left[ \frac{(p+\lambda)(q+\lambda)}{pq} \bar{F}_1(p, q) \right] \\ &= \int_0^t \int_0^{t'} (\delta(t-s) + \lambda)(\delta(t'-s') + \lambda) K_0(s+s') ds' dt' \\ &= \int \frac{(\lambda + xe^{-(x+\lambda)t})(\lambda + xe^{-(x+\lambda)t'})}{(x+\lambda)^2} d\mu_{\text{MP}}(x), \quad (402) \\ \mathcal{L}^{-1} \left[ \frac{1}{pq} \frac{\bar{R}_\ell(q) - \bar{R}_\ell(p)}{q-p} \right] &= \mathcal{L}^{-1} \left[ \frac{\bar{F}_2(p, q)}{pq} \right] = \frac{1}{\delta} \int_0^t \int_0^{t'} K_1(s, s') ds' ds \\ &= \frac{1}{\delta} \int \frac{x}{(x+\lambda)^2} (1 - e^{-(x+\lambda)t})(1 - e^{-(x+\lambda)t'}) d\mu_{\text{MP}}(x), \quad (403) \end{aligned}$$

$$\begin{aligned} \mathcal{L}^{-1} \left[ \frac{(p+\lambda)(q+\lambda)}{pq} \frac{\bar{R}_\ell(q) - \bar{R}_\ell(p)}{q-p} \right] &= \mathcal{L}^{-1} \left[ \frac{(p+\lambda)(q+\lambda)}{pq} \frac{\bar{F}_2(p, q)}{pq} \right] \\ &= \frac{1}{\delta} \int_0^t \int_0^{t'} (\delta(t-s) + \lambda)(\delta(t'-s') + \lambda) K_1(s, s') ds' dt' \\ &= \frac{1}{\delta} \int x \frac{(\lambda + xe^{-(x+\lambda)t})(\lambda + xe^{-(x+\lambda)s})}{(x+\lambda)^2} d\mu_{\text{MP}}(x). \quad (404) \end{aligned}$$

Thus, by Equations (393) and (394), we have the desired expressions for  $C_{\theta 0}^*$  and  $L_0$ .

**Solving for  $\tau > 0$ .** Equations for  $\Delta C_\theta^*(t, t') := C_\theta^*(t, t') - C_{\theta 0}^*(t, t')$  and  $\Delta L(t, t') := L(t, t') - L_0(t, t')$  are given by

$$\begin{aligned} \Delta C_\theta^*(t, t') &= \frac{1}{\delta} \int_0^t \int_0^{t'} R_\theta(t-s) R_\theta(t'-s') \Delta L(s, s') ds' ds \\ &+ \tau \int_0^t \int_0^{t'} R_\theta(t-s) R_\theta(t'-s) (R_\ell(s-s') L(s', s') + R_\ell(s'-s) L(s, s) + L(s, s') \delta(s-s')) ds' ds, \quad (405) \end{aligned}$$

and

$$\begin{aligned} \Delta L(t, t') &= \frac{1}{\delta} \int_0^t \int_0^{t'} (\delta(t-s) + R_\ell(t-s)) (\delta(t'-s') + R_\ell(t'-s')) \Delta C_\theta^*(s, s') ds' ds \\ &+ \tau \int_0^t \int_0^{t'} R_\ell(t-s) R_\ell(t'-s) L(s, s) ds. \quad (406) \end{aligned}$$

Define  $L'(t, t') := L(t, t)\delta(t - t')$ . Taking the Laplace transform of these equations, we obtain

$$\Delta \bar{C}_\theta^*(p, q) = \frac{1}{\delta} \bar{R}_\theta(p) \bar{R}_\theta(q) \Delta \bar{L}(p, q) + \tau \bar{R}_\theta(p) \bar{R}_\theta(q) (1 + \bar{R}_\ell(p) + \bar{R}_\ell(q)) \bar{L}'(p, q), \quad (407)$$

$$\Delta \bar{L}(p, q) = (1 + \bar{R}_\ell(p))(1 + \bar{R}_\ell(q)) \Delta \bar{C}_\theta^*(p, q) + \tau \delta \bar{R}_\ell(p) \bar{R}_\ell(q) \bar{L}'(p, q). \quad (408)$$

Further simplification gives

$$\Delta \bar{C}_\theta^*(p, q) = \frac{\delta^2 \bar{R}_\ell(p) \bar{R}_\ell(q)}{1 - \delta \bar{R}_\ell(p) \bar{R}_\ell(q)} \tau \bar{L}'(p, q) = \frac{\delta(\bar{R}_\ell(q) - \bar{R}_\ell(p))}{q - p} \tau \bar{L}'(p, q), \quad (409)$$

$$\begin{aligned} \Delta \bar{L}(p, q) &= \frac{\delta(\delta(\bar{R}_\ell(p) + \bar{R}_\ell(q)) + \delta + 1) R_\ell(p) R_\ell(q)}{1 - \delta \bar{R}_\ell(p) \bar{R}_\ell(q)} \tau \bar{L}'(p, q) \\ &= \frac{(\delta R_\ell(q)^2 + (\delta + 1) R_\ell(q)) - (\delta R_\ell(p)^2 + (\delta + 1) R_\ell(p))}{q - p} \tau \bar{L}'(p, q). \end{aligned} \quad (410)$$

We perform the inverse Laplace transform. For  $\Delta \bar{C}_\theta^*(p, q)$ , we have

$$\Delta C_\theta^*(t, t') = \tau \delta \int_0^t \int_0^{t'} F_2(t - s, t' - s') L'(s, s') ds' ds = \tau \int_0^{t \wedge t'} K_1(t + t' - 2s) L(s, s) ds, \quad (411)$$

and we obtain Equation (371). For  $\Delta \bar{L}(p, q)$ , we use that

$$\begin{aligned} \int_0^\infty \int_0^\infty K_2(t + t') e^{-pt - qt'} dt' dt &= \int \frac{x^2}{(x + \lambda + p)(x + \lambda + q)} d\mu_{\text{MP}}(x) \\ &= \frac{(\delta R_\ell(q)^2 + (\delta + 1) R_\ell(q)) - (\delta R_\ell(p)^2 + (\delta + 1) R_\ell(p))}{q - p}, \end{aligned} \quad (412)$$

and proceeding similarly, we obtain Equation (372). ■

Finally, we prove Proposition E.2. The asymptotic train and test errors can be expressed in terms of the DMFT solution as

$$\mathcal{L}(\boldsymbol{\theta}^t) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \boldsymbol{\theta}^t - y_i)^2 \rightarrow \mathbb{E}[(r^t - r^* - z)^2] = L(t, t), \quad (413)$$

$$\mathcal{R}(\boldsymbol{\theta}^t) = \frac{1}{d} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|_2^2 + \sigma^2 \rightarrow \mathbb{E}[(\theta^t - \theta^*)^2] + \sigma^2 = C_\theta^*(t, t) + \sigma^2. \quad (414)$$

By Lemma E.3, the asymptotic train and test errors  $\mathcal{L}(t) := L(t, t)$  and  $\mathcal{R}(t) := C_\theta^*(t, t) + \sigma^2$  satisfy

$$\mathcal{L}(t) = \mathcal{L}_0(t) + \tau \int_0^t K_2(2(t - s)) \mathcal{L}(s) ds, \quad \mathcal{R}(t) = \mathcal{R}_0(t) + \tau \int_0^t K_1(2(t - s)) \mathcal{L}(s) ds, \quad (415)$$

where  $\mathcal{L}_0(t) := L_0(t, t)$  and  $\mathcal{R}_0(t) := C_{\theta_0}^*(t, t) + \sigma^2$  are the asymptotic train and test errors for the noiseless case  $\tau = 0$ , which can be written explicitly as

$$\mathcal{L}_0(t) = \rho^2 \int \frac{x(\lambda + xe^{-(x+\lambda)t})^2}{(x+\lambda)^2} d\mu_{\text{MP}}(x) + \frac{\sigma^2}{\delta} \int \frac{(\lambda + xe^{-(x+\lambda)t})^2}{(x+\lambda)^2} d\mu_{\text{MP}}(x) + \frac{\delta - 1}{\delta} \sigma^2, \quad (416)$$

$$\mathcal{R}_0(t) = \rho^2 \int \frac{(\lambda + xe^{-(x+\lambda)t})^2}{(x+\lambda)^2} d\mu_{\text{MP}}(x) + \frac{\sigma^2}{\delta} \int \frac{x}{(x+\lambda)^2} (1 - e^{-(x+\lambda)t})^2 d\mu_{\text{MP}}(x) + \sigma^2. \quad (417)$$

Setting  $H_i(t) := K_i(2t)$ , we obtain Equations (332)–(334). This concludes the proof of Proposition E.2.

## Appendix F. Details of Numerical Simulations

We numerically solve the discretized DMFT equation given in Appendix B.2 using Monte Carlo sampling. Instead of directly working with the system  $\mathfrak{S}^\gamma$  given in Appendix B.2, we work with the equivalent system (33), as it is simpler to implement. We solve Equation (33) by iterating the following steps until convergence:

0. Start with a random guess of the DMFT solution  $(C_\theta, R_\theta)$ .
1. Given the current estimate of  $(C_\theta, R_\theta)$ , sample  $M$  instances of the stochastic processes  $r^{t_i}$  and  $\partial \ell_{t_i}(r^{t_i}; z) / \partial w^{t_j}$ , and compute the functions  $(C_\ell, R_\ell, \Gamma)$  by averaging over the samples.
2. Given the current estimate of  $(C_\ell, R_\ell, \Gamma)$ , sample  $M$  instances of the stochastic processes  $\theta^{t_i}$  and  $\partial \theta^{t_i} / \partial w^{t_j}$ , and compute the functions  $\tilde{C}_\theta$  and  $\tilde{R}_\theta$  by averaging over the samples.
3. Update the DMFT solution as  $(C_\theta, R_\theta) \leftarrow (1 - \alpha)(C_\theta, R_\theta) + \alpha(\tilde{C}_\theta, \tilde{R}_\theta)$  where  $\alpha \in (0, 1]$  is a damping factor.

In our experiments, we set the number of samples to  $M = 8000$ , the damping factor to  $\alpha = 0.8$ , and the time step to  $\gamma = 0.05$ . We observe that the above iteration converges in around 10 iterations for the settings considered in this paper.

For the logistic regression setting with  $\ell_t(r, r^*; z) = -y / (1 + \exp(yr))$  where  $y = \text{sign}(r^* + z)$ , there is an issue with computing the function  $R_\ell(t_i, *) = \mathbb{E}[\partial \ell_{t_i}(r^{t_i}, r^*; z) / \partial r^*]$  due to the non-differentiability of  $\ell$  with respect to  $r^*$  at  $r^* = -z$ . However, we can avoid the differentiation by  $r^*$  entirely by (heuristically) using Stein's lemma:

$$\begin{aligned} \mathbb{E}[r^* \ell_{t_i}(r^{t_i}, r^*; z)] &= \sum_{j=0}^{i-1} \text{Cov}(r^*, w^{t_j}) \mathbb{E}\left[\frac{\partial \ell_{t_i}(r^{t_i}, r^*; z)}{\partial w^{t_j}}\right] + \text{Cov}(r^*, r^*) \mathbb{E}\left[\frac{\partial \ell_{t_i}(r^{t_i}, r^*; z)}{\partial r^*}\right] \\ &= \gamma \sum_{j=0}^{i-1} C_\ell(t_j, *) R_\ell(t_i, t_j) + \rho^2 R_\ell(t_i, *), \end{aligned} \quad (418)$$

and computing  $R_\ell(t_i, *)$  as

$$R_\ell(t_i, *) = \frac{1}{\rho^2} \left( \mathbb{E}[r^* \ell_{t_i}(r^{t_i}, r^*; z)] - \gamma \sum_{j=0}^{i-1} R_\ell(t_i, t_j) C_\ell(t_j, *) \right), \quad (419)$$

and calculating the expectation using Monte Carlo integration.

### Appendix G. Discretization of SDEs in High Dimensions

In this section, we analyze the discretization error of a general stochastic differential equation in high dimensions. The results presented here are used in the proof of Lemma D.1.

Let  $T > 0$ ,  $\mathbf{b}: [0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ , and  $\boldsymbol{\sigma}: [0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times n}$ . Consider the following SDE in  $\mathbb{R}^d$ :

$$d\boldsymbol{\theta}^t = \mathbf{b}(t, \boldsymbol{\theta}^t) dt + \boldsymbol{\sigma}(t, \boldsymbol{\theta}^t) d\mathbf{W}^t, \quad (420)$$

with given initial condition  $\boldsymbol{\theta}^0 \in \mathbb{R}^d$ . Here,  $\mathbf{W}^t \in \mathbb{R}^n$  is a Brownian motion.

We discretize the SDE using the Euler–Maruyama method with step size  $\gamma > 0$  as follows:

$$\boldsymbol{\theta}_\gamma^{t_k+1} = \boldsymbol{\theta}_\gamma^{t_k} + \mathbf{b}(t_k, \boldsymbol{\theta}_\gamma^{t_k})\gamma + \boldsymbol{\sigma}(t_k, \boldsymbol{\theta}_\gamma^{t_k})(\mathbf{W}^{t_{k+1}} - \mathbf{W}^{t_k}), \quad \boldsymbol{\theta}_\gamma^0 = \boldsymbol{\theta}^0, \quad (421)$$

where  $t_k := k\gamma$ .

There is a standard result [35, Theorem 10.2.2] for bounding the difference between  $\boldsymbol{\theta}^t$  and  $\boldsymbol{\theta}_\gamma^t$ . However, they treat  $n, d$  as constants, and the dependence of the bound on  $n, d$  is not obvious. Here, we show a bound that explicitly tracks the dependence on  $n, d$ .

**Assumption G.1** *There exists a constant  $L > 0$  independent of  $n$  and  $d$  such that the following hold.*

1. (Proportional asymptotics):  $1/L \leq n/d \leq L$ .
2. (Lipschitz continuity):  $\|\mathbf{b}(t_1, \mathbf{x}_1) - \mathbf{b}(t_2, \mathbf{x}_2)\|_2^2 + \|\boldsymbol{\sigma}(t_1, \mathbf{x}_1) - \boldsymbol{\sigma}(t_2, \mathbf{x}_2)\|_F^2 \leq L(d(t_1 - t_2)^2 + \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2)$  for  $t_1, t_2 \in [0, T]$ ,  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$ .
3. (Linear growth):  $\|\mathbf{b}(t, \mathbf{x})\|_2^2 + \|\boldsymbol{\sigma}(t, \mathbf{x})\|_F^2 \leq L(d + dt^2 + \|\mathbf{x}\|_2^2)$  for  $t \in [0, T]$ ,  $\mathbf{x} \in \mathbb{R}^d$ .

**Lemma G.2 (Norm bound on the discretized SDE iterates)** *Under Assumption G.1 and  $\gamma < 1$ , there exists a constant  $C := C(L) > 0$  independent of  $n, d, T, \gamma$  such that the following holds.*

$$\mathbb{E} \left[ \max_{0 \leq k \leq \lfloor T/\gamma \rfloor} \|\boldsymbol{\theta}_\gamma^{t_k}\|_2^2 \right] \leq e^{CT} (Td + \|\boldsymbol{\theta}^0\|_2^2). \quad (422)$$

**Proof** For notational simplicity, we denote  $K := \lfloor T/\gamma \rfloor$ ,  $\boldsymbol{\theta}^k := \boldsymbol{\theta}_\gamma^{t_k}$ ,  $\mathbf{b}^k := \mathbf{b}(t_k, \boldsymbol{\theta}_\gamma^{t_k})$ ,  $\boldsymbol{\sigma}^k := \boldsymbol{\sigma}(t_k, \boldsymbol{\theta}_\gamma^{t_k})$ , and  $\boldsymbol{\xi}^k := \mathbf{W}^{t_{k+1}} - \mathbf{W}^{t_k}$ . From Equation (421), we have

$$\|\boldsymbol{\theta}^{k+1}\|_2^2 = \|\boldsymbol{\theta}^k\|_2^2 + 2\gamma \boldsymbol{\theta}^{k\top} \mathbf{b}^k + \gamma^2 \|\mathbf{b}^k\|_2^2 + \|\boldsymbol{\sigma}^k \boldsymbol{\xi}^k\|_2^2 + 2(\boldsymbol{\theta}^k + \gamma \mathbf{b}^k)^\top \boldsymbol{\sigma}^k \boldsymbol{\xi}^k. \quad (423)$$

Summing over  $k$  steps, we have

$$\|\boldsymbol{\theta}^k\|_2^2 = \|\boldsymbol{\theta}^0\|_2^2 + \underbrace{\sum_{j=0}^{k-1} \left( 2\gamma \boldsymbol{\theta}^j \top \mathbf{b}^j + \gamma^2 \|\mathbf{b}^j\|_2^2 + \|\boldsymbol{\sigma}^j \boldsymbol{\xi}^j\|_2^2 \right)}_{=: A^k} + \underbrace{\sum_{j=0}^{k-1} 2(\boldsymbol{\theta}^j + \gamma \mathbf{b}^j)^\top \boldsymbol{\sigma}^j \boldsymbol{\xi}^j}_{=: M^k}. \quad (424)$$

We first bound  $A^k$ . We have

$$\max_{0 \leq k \leq K} |A^k| \leq \sum_{j=0}^{K-1} \left( 2\gamma |\boldsymbol{\theta}^j \mathbf{b}^j| + \gamma^2 \|\mathbf{b}^j\|_2^2 + \|\boldsymbol{\sigma}^j \boldsymbol{\xi}^j\|_2^2 \right). \quad (425)$$

By the linear growth condition, we have

$$2\gamma |\boldsymbol{\theta}^j \mathbf{b}^j| + \gamma^2 \|\mathbf{b}^j\|_2^2 \leq 2\gamma \|\boldsymbol{\theta}^j\|_2^2 + (2\gamma + \gamma^2) \|\mathbf{b}^j\|_2^2 \leq C\gamma((1 + T^2)d + \|\boldsymbol{\theta}^j\|_2^2). \quad (426)$$

Using the covariance of  $\boldsymbol{\xi}^j \sim \mathbf{N}(0, \gamma \mathbf{I})$  and the linear growth condition, we have

$$\mathbb{E}[\|\boldsymbol{\sigma}^j \boldsymbol{\xi}^j\|_2^2 \mid \mathcal{F}^j] = \gamma \|\boldsymbol{\sigma}^j\|_{\mathbb{F}}^2 \leq C\gamma((1 + T^2)d + \|\boldsymbol{\theta}^j\|_2^2), \quad (427)$$

where  $\mathcal{F}^j$  is the filtration generated by  $\{\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^{j-1}\}$ . Thus, taking the expectation, we have

$$\mathbb{E} \left[ \max_{0 \leq k \leq K} |A^k| \right] \leq C\gamma \sum_{j=0}^{K-1} ((1 + T^2)d + \mathbb{E}\|\boldsymbol{\theta}^j\|_2^2). \quad (428)$$

Next, we bound  $M^k$ . Notice that  $M^k$  is a martingale with respect to the filtration  $\mathcal{F}^k$  since  $\mathbb{E}[\boldsymbol{\xi}^j \mid \mathcal{F}^j] = 0$ . Thus, by the Burkholder–Davis–Gundy inequality, we have

$$\begin{aligned} \mathbb{E} \left[ \max_{0 \leq k \leq K} |M^k| \right] &\leq C \mathbb{E} \left( \sum_{j=0}^{K-1} \mathbb{E} \left[ \left( 2(\boldsymbol{\theta}^j + \gamma \mathbf{b}^j)^\top \boldsymbol{\sigma}^j \boldsymbol{\xi}^j \right)^2 \mid \mathcal{F}^j \right] \right)^{1/2} \\ &= C \mathbb{E} \left( \sum_{j=0}^{K-1} 4\gamma \|\boldsymbol{\theta}^j + \gamma \mathbf{b}^j\|_2^2 \|\boldsymbol{\sigma}^j\|_2^2 \right)^{1/2} \\ &= C \mathbb{E} \left( \gamma \sum_{j=0}^{K-1} \|\boldsymbol{\theta}^j + \gamma \mathbf{b}^j\|_2^2 \|\boldsymbol{\sigma}^j\|_{\mathbb{F}}^2 \right)^{1/2} \\ &\leq C \mathbb{E} \left[ \left( \max_{0 \leq k \leq K} \|\boldsymbol{\theta}^k + \gamma \mathbf{b}^k\|_2^2 \right)^{1/2} \left( \gamma \sum_{j=0}^{K-1} \|\boldsymbol{\sigma}^j\|_{\mathbb{F}}^2 \right)^{1/2} \right] \\ &\leq C \mathbb{E} \left[ \left( \gamma^2(1 + T^2)d + \max_{0 \leq k \leq K} \|\boldsymbol{\theta}^k\|_2^2 \right)^{1/2} \left( \gamma \sum_{j=0}^{K-1} ((1 + T^2)d + \|\boldsymbol{\theta}^j\|_2^2) \right)^{1/2} \right], \end{aligned} \quad (429)$$

where we used the linear growth condition in the last line. By Young's inequality, we have

$$\mathbb{E} \left[ \max_{0 \leq k \leq K} |M^k| \right] \leq \frac{1}{2} \left( \gamma^2(1 + T^2)d + \mathbb{E} \left[ \max_{0 \leq k \leq K} \|\boldsymbol{\theta}^k\|_2^2 \right] \right) + C\gamma \sum_{j=0}^{K-1} ((1 + T^2)d + \mathbb{E}\|\boldsymbol{\theta}^j\|_2^2). \quad (430)$$

Combining Equations (428) and (430), we have

$$\mathbb{E} \left[ \max_{0 \leq k \leq K} \|\boldsymbol{\theta}^k\|_2^2 \right] \leq \|\boldsymbol{\theta}^0\|_2^2 + C\gamma^2(1+T^2)d + C\gamma \sum_{j=0}^{K-1} ((1+T^2)d + \mathbb{E}\|\boldsymbol{\theta}^j\|_2^2) + \frac{1}{2} \mathbb{E} \left[ \max_{0 \leq k \leq K} \|\boldsymbol{\theta}^k\|_2^2 \right]. \quad (431)$$

Rearranging the terms, we have

$$\mathbb{E} \left[ \max_{0 \leq k \leq K} \|\boldsymbol{\theta}^k\|_2^2 \right] \leq C(T(1+T^2)d + \|\boldsymbol{\theta}^0\|_2^2) + C\gamma \sum_{j=0}^{K-1} \mathbb{E} \left[ \max_{0 \leq j \leq k} \|\boldsymbol{\theta}^j\|_2^2 \right]. \quad (432)$$

The bound (422) follows by applying Grönwall's inequality and absorbing  $1+T^2$  into the exponential factor  $e^{CT}$ .  $\blacksquare$

**Remark G.3** Lemma G.2 holds verbatim to the case where the Gaussian increments  $\mathbf{W}^{t_{k+1}} - \mathbf{W}^{t_k}$  are replaced by independent random vectors  $\boldsymbol{\xi}^k \in \mathbb{R}^n$  with  $\mathbb{E}[\boldsymbol{\xi}^k] = 0$  and  $\mathbb{E}[\boldsymbol{\xi}^k \boldsymbol{\xi}^{k\top}] = \gamma \mathbf{I}_n$ , as the proof only uses up to the second moment of the increments.

**Lemma G.4 (Strong approximation of SDE)** Under Assumption G.1 and  $\gamma < 1$ , there exists a constant  $C := C(L) > 0$  that does not depend on  $n, d, T, \gamma$  such that the following holds.

$$\mathbb{E} \left[ \sup_{0 \leq t \leq T} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t\|_2^2 \right] \leq e^{CT} \gamma (Td + \|\boldsymbol{\theta}^0\|_2^2). \quad (433)$$

**Proof** Equation (420) can be written as

$$\boldsymbol{\theta}^t = \boldsymbol{\theta}^0 + \int_0^t \mathbf{b}(s, \boldsymbol{\theta}^s) ds + \int_0^t \boldsymbol{\sigma}(s, \boldsymbol{\theta}^s) d\mathbf{W}^s. \quad (434)$$

Let  $\lfloor t \rfloor := \max\{k\gamma \mid k\gamma \leq t, k \in \mathbb{N}\}$  and consider the following stochastic process that embeds Equation (421) into continuous time.

$$\boldsymbol{\theta}_\gamma^t = \boldsymbol{\theta}^0 + \int_0^t \mathbf{b}(\lfloor s \rfloor, \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor}) ds + \int_0^t \boldsymbol{\sigma}(\lfloor s \rfloor, \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor}) d\mathbf{W}^s. \quad (435)$$

Using Itô's lemma on  $\|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t\|_2^2$ , we have

$$\begin{aligned} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t\|_2^2 &= 2 \int_0^t (\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^s)^\top (\mathbf{b}(s, \boldsymbol{\theta}^s) - \mathbf{b}(\lfloor s \rfloor, \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor})) ds \\ &\quad + \int_0^t \|\boldsymbol{\sigma}(s, \boldsymbol{\theta}^s) - \boldsymbol{\sigma}(\lfloor s \rfloor, \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor})\|_{\mathbb{F}}^2 ds + M^t, \end{aligned} \quad (436)$$

$$M^t := 2 \int_0^t (\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^s)^\top (\boldsymbol{\sigma}(s, \boldsymbol{\theta}^s) - \boldsymbol{\sigma}(\lfloor s \rfloor, \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor})) d\mathbf{W}^s. \quad (437)$$

By the Lipschitz assumption, we have

$$2(\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^s)^\top (\mathbf{b}(s, \boldsymbol{\theta}^s) - \mathbf{b}(\lfloor s \rfloor, \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor})) \leq \|\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^s\|_2^2 + \|\mathbf{b}(s, \boldsymbol{\theta}^s) - \mathbf{b}(\lfloor s \rfloor, \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor})\|_2^2$$

$$\leq C(\gamma^2 d + \|\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor}\|_2^2), \quad (438)$$

$$\|\boldsymbol{\sigma}(s, \boldsymbol{\theta}^s) - \boldsymbol{\sigma}(\lfloor s \rfloor, \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor})\|_{\mathbb{F}}^2 \leq C(\gamma^2 d + \|\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor}\|_2^2). \quad (439)$$

Therefore, we have

$$\mathbb{E} \left[ \sup_{0 \leq t \leq T} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t\|_2^2 \right] \leq C \int_0^T (\gamma^2 d + \mathbb{E} \|\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor}\|_2^2) ds + \mathbb{E} \left[ \sup_{0 \leq t \leq T} |M^t| \right]. \quad (440)$$

We next bound the martingale term  $M^t$ . By the Burkholder–Davis–Gundy inequality,

$$\begin{aligned} \mathbb{E} \left[ \sup_{0 \leq t \leq T} |M^t| \right] &\leq C \mathbb{E} \left[ \left( 4 \int_0^T \|(\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^s)^\top (\boldsymbol{\sigma}(s, \boldsymbol{\theta}^s) - \boldsymbol{\sigma}(\lfloor s \rfloor, \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor}))\|_2^2 ds \right)^{1/2} \right] \\ &\leq C \mathbb{E} \left[ \left( \int_0^T \|\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^s\|_2^2 \|\boldsymbol{\sigma}(s, \boldsymbol{\theta}^s) - \boldsymbol{\sigma}(\lfloor s \rfloor, \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor})\|_{\mathbb{F}}^2 ds \right)^{1/2} \right] \\ &\leq C \mathbb{E} \left[ \left( \sup_{0 \leq t \leq T} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t\|_2^2 \right)^{1/2} \left( \int_0^T (\gamma^2 d + \|\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor}\|_2^2) ds \right)^{1/2} \right] \\ &\leq \frac{1}{2} \mathbb{E} \left[ \sup_{0 \leq t \leq T} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t\|_2^2 \right] + C \int_0^T (\gamma^2 d + \mathbb{E} \|\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor}\|_2^2) ds. \end{aligned} \quad (441)$$

In the last line, we used Young’s inequality. Therefore, we have

$$\mathbb{E} \left[ \sup_{0 \leq t \leq T} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t\|_2^2 \right] \leq C \int_0^T (\gamma^2 d + \mathbb{E} \|\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor}\|_2^2) ds + \frac{1}{2} \mathbb{E} \left[ \sup_{0 \leq t \leq T} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t\|_2^2 \right]. \quad (442)$$

By rearranging the terms, we have

$$\mathbb{E} \left[ \sup_{0 \leq t \leq T} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t\|_2^2 \right] \leq C\gamma^2 Td + C \int_0^T \mathbb{E} \|\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor}\|_2^2 ds. \quad (443)$$

Finally, we have  $\|\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor}\|_2^2 \leq 2\|\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^s\|_2^2 + 2\|\boldsymbol{\theta}_\gamma^s - \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor}\|_2^2$  and

$$\boldsymbol{\theta}_\gamma^s - \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor} = \mathbf{b}(\lfloor s \rfloor, \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor})(s - \lfloor s \rfloor) + \boldsymbol{\sigma}(\lfloor s \rfloor, \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor})(\mathbf{W}^s - \mathbf{W}^{\lfloor s \rfloor}). \quad (444)$$

By the linear growth condition and Lemma G.2, we have

$$\begin{aligned} \mathbb{E} \|\boldsymbol{\theta}_\gamma^s - \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor}\|_2^2 &\leq C\gamma^2 \mathbb{E} \|\mathbf{b}(\lfloor s \rfloor, \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor})\|_2^2 + C\gamma \mathbb{E} \|\boldsymbol{\sigma}(\lfloor s \rfloor, \boldsymbol{\theta}_\gamma^{\lfloor s \rfloor})\|_{\mathbb{F}}^2 \\ &\leq C\gamma(d + \mathbb{E} \|\boldsymbol{\theta}_\gamma^{\lfloor s \rfloor}\|_2^2) \leq e^{CT} \gamma(d + \|\boldsymbol{\theta}^0\|_2^2). \end{aligned} \quad (445)$$

Combining the above bounds, we have

$$\mathbb{E} \left[ \sup_{0 \leq t \leq T} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\gamma^t\|_2^2 \right] \leq e^{CT} \gamma(Td + \|\boldsymbol{\theta}^0\|_2^2) + C \int_0^T \mathbb{E} \left[ \sup_{0 \leq s \leq t} \|\boldsymbol{\theta}^s - \boldsymbol{\theta}_\gamma^s\|_2^2 \right] dt. \quad (446)$$

Applying Grönwall’s inequality, we obtain the bound (433). ■