

---

# Three-stage Evolution and Fast Equilibrium for SGD with Non-degenerate Critical Points

---

Yi Wang<sup>1</sup> Zhiren Wang<sup>2</sup>

## Abstract

We justify the fast equilibrium conjecture on stochastic gradient descent from (Li et al., 2020) under the assumptions that critical points are non-degenerate and the stochastic noise is a standard Gaussian. In this case, we prove an scaling invariant SGD with constant effective learning rate consists of three stages: descent, diffusion and tunneling, and explicitly identify temporary equilibrium states that can be observed within practical training time. This interprets the gap between the mixing time in the fast equilibrium conjecture and the previously known upper bound. While our assumptions do not represent typical implementations of SGD of neural networks in practice, this is the first description of the three-stage mechanism in any case. The main finding in this mechanism is that a temporary equilibrium of local nature is quickly achieved after polynomial time (in term of the reciprocal of the intrinsic learning rate) and then stabilizes within observable time scales; and that the temporary equilibrium is in general different from the global Gibbs equilibrium, which will only appear after an exponentially long period beyond typical training limits. Our experiments support that this mechanism may extend to the general case.

## 1. Introduction

### 1.1. Background and motivation

Stochastic gradient descent (SGD) has been an indispensable tool in the training of neural networks and is known to work in both convex and non-convex settings. One theme that has recently attracted much attention is the descrip-

<sup>1</sup>Department of Mathematics, Johns Hopkins University  
<sup>2</sup>Department of Mathematics, Pennsylvania State University. Correspondence to: Yi Wang <ywang@math.jhu.edu>, Zhiren Wang <zhirenw@psu.edu>.

tion of the asymptotic behavior of the SGD, especially how learning rate schemes affect the distribution of random trajectories. A useful approach is to model the SGD

$$\mathbf{w}_{k+1} \leftarrow (1 - \lambda)\mathbf{w}_k - \eta \nabla \mathcal{L}_{\mathcal{B}_k}(\mathbf{w}_k), \quad (1)$$

where  $\mathcal{L}_{\mathcal{B}_k}$  is the loss function over a stochastic batch  $\mathcal{B}_k$  and  $\lambda$  and  $\eta$  are respectively the weight decay and the learning rate, by stochastic differential equations (SDE). In this approach, the random noise in  $\nabla \mathcal{L}_{\mathcal{B}_k}(\mathbf{w}_k)$  is regarded, by taking a continuous limit on time increments, as a Brownian motion with covariance matrix  $\Sigma(\mathbf{w}_k)$ . A widely adopted approach is to approximate the SGD (1) by the value at  $t = k\eta$  of the SDE:

$$d\mathbf{W}_t = -\eta(\nabla \mathcal{L}(\mathbf{W}_t))dt + \Sigma(\mathbf{W}_t)^{\frac{1}{2}}d\mathbf{B}_t^d - \lambda_e \mathbf{W}_t dt, \quad (2)$$

starting from  $\mathbf{W}_0 = \mathbf{w}_0$ , where  $\lambda_e = \lambda\eta$ , and  $\mathbf{B}_t^d$  stands for the standard Brownian motion in  $\mathbb{R}^d$  (see §2.1 for more details). The quantity  $\lambda_e$  is called the **intrinsic learning rate** and is known to decide the limit behavior of the dynamics.

Modern neural networks often contain various normalization steps, such as batch normalization, weight normalization and layer normalization. Normalization makes the loss function  $\mathcal{L}$  **scaling invariant** and typically non-smooth near the origin. A priori, it may take a long time for the SDE (2) to reach an equilibrium as  $|\mathbf{W}_t|$  may remain too large or too small for a long period. Recently (Li et al., 2020) studied how  $\lambda_e$  affects the distribution of the solutions through another quantity  $\gamma_t^{-\frac{1}{2}} := |\mathbf{W}_t|^{-2}\eta$ , called the **effective learning rate**. Under mild assumptions, they proved that  $\gamma_t^{-\frac{1}{2}}$  stabilizes to the magnitude  $O(\lambda_e^{\frac{1}{2}})$  after  $O(\frac{1}{\lambda_e})$  time. They further conjectured:

**Conjecture 1.1. (Fast Equilibrium Conjecture)** (Li et al., 2020) Suppose  $F(\mathbf{W}, \mathbf{x})$  is the output of the neural network with parameter  $\mathbf{W}$  and input data  $\mathbf{x}$ , then the distribution of  $F(\mathbf{W}_t, \mathbf{x})$  stabilizes in total variation distance for all  $\mathbf{x}$  after  $O(\frac{1}{\lambda_e})$  time to an equilibrium state. Moreover, this distribution is independent of the initial parameter  $\mathbf{W}_0$ .

Remark that the conjecture concerns the speed of convergence of a sequence of distributions towards an equilibrium state in terms of the total variation distance between probability measures. An interesting but different question in

similar settings, about convergence of a typical trajectory towards a local minima in terms of the gap in the loss function, has been studied in (Raginsky et al., 2017; Zhang et al., 2017; Xu et al., 2018; Huang & Becker, 2021).

While the conjecture is supported by numerical experiments, it currently lacks theoretical explanation. The relation between the convergence time of the SDE model and the learning rate  $\lambda_e$  has been studied in (Bovier et al., 2004; Shi et al., 2020) and the best known upper bound for mixing time is  $O(e^{C\lambda_e^{-1}})$  for networks without normalization. For systems with normalization, after adapting the stabilized value of effective learning rate  $O(\lambda_e^{\frac{1}{2}})$  in (Li et al., 2020), in lieu of  $\lambda_e$ , this bound becomes  $O(e^{C\lambda_e^{-\frac{1}{2}}})$ , which is much larger than  $O(\frac{1}{\lambda_e})$ . This gap has not been theoretically explained and will be the main focus of this paper.

Thanks to scaling-invariance, the distribution of  $F(\mathbf{W}_t, \mathbf{x})$  is determined by the distribution of  $\frac{\mathbf{W}_t}{|\mathbf{W}_t|}$ . We will study this later distribution on the unit sphere  $\mathbb{S}^{d-1}$  and prove that it displays fast convergence to certain temporary equilibria. It turns out that, unlike in Conjecture 1.1, the temporary equilibrium of  $\frac{\mathbf{W}_t}{|\mathbf{W}_t|}$  after such fast convergence does depend on the initial value  $\mathbf{W}_0$ . For possible interpretations that reconcile such dependence with the initial parameter independence in Conjecture 1.1, see the discussions in §6.1.

We also note that several earlier works, such as (Mandt et al., 2017; Izmailov et al., 2018) investigated local mixing in the convex optimization case, when there is only one local minimum and the loss function is assumed to be quadratic. Compared to these works, the current paper deals with other obstructions that arise in a non-convex setting. Namely, we will mainly focus on the separation between stages, which is not an issue in the convex case since that case does not have a final stage during which trajectories move across basins.

## 1.2. Restriction of mathematical tools and our goal

It is worth noting that both the works (Bovier et al., 2004; Shi et al., 2020) assumed two hypothesis: (i) the noise  $\Sigma$  is the standard Gaussian (see Assumption 4.1), (ii)  $\mathcal{L}$  is a Morse function, i.e. all critical points are non-degenerate and hence isolated (see Assumption 4.2). The reason lies in the limitation of mathematical tools: the only currently available mathematical theory that applies to the study of convergence of distributions towards equilibria around local minima is the seminal work of Barry Simon (Simon, 1983) on semiclassical analysis of low lying eigenvalues. However this theory is limited to the two assumption above. It is unclear whether similar conclusions can be achieved beyond these two assumptions without a major update to the underlying mathematical theory.

As no such updates exist to date, the goal of this paper is

to work within the framework of Simon’s theory and to explain the huge gap between the experimentally supported upper bound in Conjecture 1.1 and the previously known exponentially large upper bounds under these assumptions.

## 1.3. Our contributions

The main contributions of this paper are:

(1) We derive a spherical SDE model (Definition 3.2) of the SGD with constant effective learning rate. This model focuses on the normalized parameter  $\frac{\mathbf{W}_t}{|\mathbf{W}_t|}$  and uses intrinsic differential operators of  $\mathbb{S}^{d-1}$  (instead of those in  $\mathbb{R}^d$ ). Since the output of the neural network depends only on  $\frac{\mathbf{W}_t}{|\mathbf{W}_t|}$ , this does not affect Conjecture 1.1. This spherical SDE is also the mathematical model of the Riemannian realization of batch normalization in (Cho & Lee, 2017).

(2) We introduce, for the first time to the best of our knowledge, the three-stage description of the SGD: **descent**, **diffusion** and **tunneling**. The descent stage sends, with high probability, a point near the bottom of attracting basin containing it. The diffusion stages stabilizes the distribution towards a temporary Gibbs equilibrium that is locally Gibbs within each individual attracting basin. The tunneling stage allows mass to slowly leak between basins to achieve the unique global Gibbs state. The three stages respectively take at most  $O(\frac{1}{\lambda_e^{\frac{1}{2}}})$ ,  $O(\frac{1}{\lambda_e})$ ,  $O(e^{\frac{C}{\lambda_e}})$  in time. It was previously known (Shi et al., 2020) that convergence towards the global Gibbs state is exponentially slow. But to our best knowledge our result for the first time identifies the distinction between the three stages, especially the fast-slow contrast between the diffusion and tunneling stages.

(3) Our proof is based on a completely new strategy. Instead of fully relying on Simon’s semiclassical analysis, we complement this theory with a probabilistic argument (Lemma E.6) that characterizes microscopically the difficulty for an individual trajectory to escape from an attracting basin. This is the key to materialize the gap between small and large non-zero eigenvalues discovered in (Simon, 1983), which leads to the aforementioned separation between the diffusion and tunneling stages. Previous works (Bovier et al., 2004; Shi et al., 2020) only used the gap between 0 and the smallest non-zero eigenvalue.

(4) We derive an explicit formula of the temporary equilibrium achieved by the diffusion stage, which is the equilibrium observed in real word trainings: it is a linear combination of local Gibbs states  $\mu_k$  in the attracting basins  $U_k$ . And the weight of each basin is approximately the same as the initial distribution of parameter among basins.

(5) Finally, we remark that our analysis also works in the setting of (Shi et al., 2020), where scaling invariance is not assumed. Instead of the compact parametric space  $\mathbb{S}^{d-1}$ ,

one assumes  $\mathcal{L}$  is a loss function on  $\mathbb{R}^d$  that grows fast as  $|\mathbf{x}| \rightarrow \infty$ .  $\mathbf{y}$  (which is the case when the neural network has regularization such as weight decay). In fact, the main mathematical theories from (Simon, 1983; Freidlin & Wentzell, 2012), that our proof relies on, work in both the  $\mathbb{S}^{d-1}$  and the fast growing functions in  $\mathbb{R}^d$  settings.

#### 1.4. Limitations

The main limitation of our study is the adoption of Assumptions 4.1 and Assumption 4.2. As mentioned above, the reason for these assumptions is their indispensability for the application of (Simon, 1983). On the one hand, the isotropic Gaussian is a popular assumption to use study SGD via SDE, and Morse functions are mathematically generic among  $C^2$ -functions. On the other hand, a serious limitation of these assumptions arises from the fact that most modern neural networks are overparametrized, which forces local minima to form connected regions, instead of being isolated like in the case of Morse functions, see (Garipov et al., 2018; Kudipudi et al., 2019; Maddox et al., 2020; Benton et al., 2021; Cooper, 2021). The local dynamics near such regions has more recently been studied in (Li et al., 2021b).

Experimental evidences suggest that the temporary equilibria at constant effective learning rates should still be localized in the general setting. We hope the study of this phenomenon may shed light on better understanding the dynamics of SGD in schemes not covered by Simon’s theory, and formulate it as Conjecture 6.3.

## 2. Preliminaries

### 2.1. SDE model

The SDE model of SGD has been extensively studied in recent years. Given a dataset  $\mathcal{S} = \{\mathbf{x}_i\}_{i=1}^N$ , at the  $k$ -th step of the SGD, a subset  $\mathcal{B}_k \subset \mathcal{S}$  of fixed size  $n$  is randomly drawn to train a neural network whose parameters are denoted by  $\mathbf{w} \in \mathbb{R}^d$ . For  $\mathbf{x}$  from a mini-batch  $\mathcal{B}$  and a parameter  $\mathbf{w}$ , the neural network outputs a loss function  $\ell_{\mathcal{B}}(\mathbf{w}, \mathbf{x})$ , which is assumed to be differentiable in  $\mathbf{w}$ . The loss function depends on  $\mathcal{B}$  because of batch normalization steps inside the neural network. The loss functions over  $\mathcal{B}$  is  $\mathcal{L}_{\mathcal{B}}(\mathbf{w}) := \mathbb{E}_{\mathbf{x} \in \mathcal{B}} \ell_{\mathcal{B}}(\mathbf{w}, \mathbf{x})$ . Also define the average loss function by  $\mathcal{L}(\mathbf{w}) := \mathbb{E}_{\mathbf{x} \in \mathcal{S}} \ell_{\mathcal{B}}(\mathbf{w}, \mathbf{x}) = \mathbb{E}_{\substack{\mathcal{B} \subset \mathcal{S} \\ |\mathcal{B}|=n}} \mathcal{L}_{\mathcal{B}}(\mathbf{w})$  (averaging all subsets  $\mathcal{B}$  of size  $n$  randomly drawn from  $\mathcal{S}$ ). In the  $k$ -th step, the parameter  $\mathbf{w}_k$  is updated by (1).

The function  $\mathcal{L}$  now defines the gradient vector field  $\nabla \mathcal{L}$  at every point  $\mathbf{w} \in \mathbb{R}^d$ . We also define a non-negative definite

symmetric matrix  $\Sigma(\mathbf{w}) \in \text{Mat}(d, d)$  at every  $\mathbf{w} \in \mathbb{R}^d$  by

$$\Sigma(\mathbf{w}) := \mathbb{E}_{\substack{\mathcal{B} \subset \mathcal{S} \\ |\mathcal{B}|=n}} \left( (\nabla \mathcal{L}_{\mathcal{B}}(\mathbf{w}) - \nabla \mathcal{L}(\mathbf{w})) (\nabla \mathcal{L}_{\mathcal{B}}(\mathbf{w}) - \nabla \mathcal{L}(\mathbf{w}))^\top \right). \quad (3)$$

The matrix  $\Sigma(\mathbf{w})$  is called the gradient noise, and will take the role as a diffusion matrix.

The popular approach of approximating (1) by (2) has been studied in (Jastrzebski et al., 2017; Goyal et al., 2017; Smith & Le, 2018; Smith et al., 2018; Chaudhari & Soatto, 2018; Shi et al., 2020; Li et al., 2019, Li et al., 2020). For example, Li et al. (2019) proved that, for small  $\eta$ , starting from the same initial position the probability distributions of the random variables  $\mathbf{W}_{k\eta}$  in (2) and  $\mathbf{w}_k$  in (1) are close to each other for fixed  $K$  and all  $k \leq \frac{K}{\eta}$ . See also (Yaida, 2019; Smith et al., 2021; Li et al., 2021a) for discussions on the deficiencies of SDE view of SGD as well as conditions that guarantee validate this view.

Li et al. (2020, Theorem 5.1) proved that the SDE (2) is in turn equivalent to

$$\begin{aligned} d\overline{\mathbf{W}}_t &= -\gamma_t^{-\frac{1}{2}} \left( \nabla \mathcal{L}(\overline{\mathbf{W}}_t) dt + \Sigma(\overline{\mathbf{W}}_t)^{\frac{1}{2}} d\mathbf{B}_t^d \right) \\ &\quad - \frac{1}{2} \gamma_t^{-1} \text{Tr} \Sigma(\overline{\mathbf{W}}_t) \overline{\mathbf{W}}_t dt \end{aligned}; \quad (4)$$

$$\frac{d\gamma_t}{dt} = -4\lambda_e \gamma_t + 2\text{Tr} \Sigma(\overline{\mathbf{W}}_t), \quad (5)$$

where  $\overline{\mathbf{W}}_t = \frac{\mathbf{w}_t}{|\mathbf{w}_t|}$  and  $\gamma_t = |\mathbf{w}_t|^4 \eta^{-2}$ . The quantity  $\gamma_t^{-\frac{1}{2}}$  is called the *effective learning rate*. In particular,  $\overline{\mathbf{W}}_t$  is a random process on the unit sphere  $\mathbb{S}^{d-1}$ . In addition, experimental results from (Li et al., 2020) suggest that both  $\text{Tr} \Sigma(\overline{\mathbf{W}}_t)$  and  $\gamma_t$  stabilizes quickly near constant values, and the learning performance has little dependence on the initial parameter  $\gamma_0$ .

### 2.2. Normalization and scaling invariance within components

The neural network is assumed to be batch normalized, which guarantees that  $\ell_{\mathcal{B}}(\mathbf{w}, \mathbf{x}) = \ell_{\mathcal{B}}(c\mathbf{w}, \mathbf{x})$  for all  $c > 0$ . In consequence (Li et al., 2020, Lemma B.1),  $\mathbf{w}^\top \nabla \ell_{\mathcal{B}}(\mathbf{w}, \mathbf{x}) = 0$ . The same orthogonality holds for  $\nabla \mathcal{L}_{\mathcal{B}}(\mathbf{w})$  and  $\nabla \mathcal{L}(\mathbf{w})$ , which are linear combinations of the  $\nabla \ell_{\mathcal{B}}$ ’s. Thus,

$$\mathbf{w}^\top \nabla \mathcal{L}(\mathbf{w}) = 0, \Sigma(\mathbf{w})\mathbf{w} = 0. \quad (6)$$

For a vector  $\mathbf{w}$ , denote by  $V_{\mathbf{w}}^\perp$  its  $(d-1)$ -dimensional orthogonal complement vector space. As  $\Sigma(\mathbf{w})$  is symmetric and has mutually orthogonal eigenspaces,  $\Sigma(\mathbf{w})$  preserves  $V_{\mathbf{w}}^\perp$ . Therefore  $\Sigma(\mathbf{w})^{\frac{1}{2}}$  annihilates  $\mathbf{w}$  and preserves  $V_{\mathbf{w}}^\perp$  as well. Recall  $\Sigma(\mathbf{w})^{\frac{1}{2}}$  is a uniquely defined positive semi-definite matrix (Horn & Johnson, 2013, Thm. 7.2.6).

### 3. The spherical model

By scaling invariance, the parameters  $\mathbf{W}_t$  and  $\overline{\mathbf{W}}_t$  yield the same outcome. So it suffices to understand the distribution of  $\overline{\mathbf{W}}_t$  to consider Conjecture 1.1. We shall ignore (5) and focus on (4).

#### 3.1. Description of model

Denote by  $\overline{\nabla}$  the gradient on  $\mathbb{S}^{d-1}$  with respect to the standard sphere metric. Then  $\overline{\nabla}\mathcal{L} : \mathbb{S}^{d-1} \rightarrow T\mathbb{S}^{d-1}$  is a vector field on  $\mathbb{S}^{d-1}$ , whose value coincides with  $\nabla\mathcal{L}$  by (6). For later use, we also write  $\Delta$  for the Laplacian on  $\mathbb{S}^{d-1}$ .

Write  $\overline{\Sigma}(\mathbf{w})$  and  $\overline{\Sigma}(\mathbf{w})^{\frac{1}{2}}$  for the restrictions of  $\Sigma(\mathbf{w})$  and  $\Sigma(\mathbf{w})^{\frac{1}{2}}$  to  $\mathbb{S}^{d-1}$ , viewed as tensor fields that send  $\mathbb{R}^d$  to  $T_{\mathbf{w}}\mathbb{S}^{d-1} = V_{\mathbf{w}}^{\perp}$  for  $\mathbf{w} \in \mathbb{S}^{d-1}$ . In particular, given the Brownian motion  $d\mathbf{B}_t^d$  on  $\mathbb{R}^d$ ,  $\overline{\Sigma}(\overline{\mathbf{W}}_t)^{\frac{1}{2}}d\mathbf{B}_t^d$  is a random infinitesimal vector along the tangent space  $T_{\mathbf{w}}\mathbb{S}^{d-1}$ .

**Theorem 3.1.** *Starting from an initial value  $\overline{\mathbf{W}}_0 \in \mathbb{S}^{d-1}$ , the SDE (4) is equivalent to the following SDE on  $\mathbb{S}^{d-1}$ :*

$$d\overline{\mathbf{W}}_t = -\gamma_t^{-\frac{1}{2}} \left( \overline{\nabla}\mathcal{L}(\overline{\mathbf{W}}_t)dt + \overline{\Sigma}(\overline{\mathbf{W}}_t)^{\frac{1}{2}}d\mathbf{B}_t^d \right). \quad (7)$$

**Difference between the SDE's (4) and (7).** We now explain the meaning of Theorem 3.1, as it may at first glance seem unnatural to remove the last term  $-\frac{1}{2}\gamma_t^{-1}\text{Tr}\Sigma(\overline{\mathbf{W}}_t)\overline{\mathbf{W}}_t dt$  from (4) without destroying the equality. The difference is as follows: (7) is a differential equation defined on the manifold  $\mathbb{S}^{d-1}$  with respect to the intrinsic geometry of this manifold, where the terms  $\overline{\nabla}\mathcal{L}(\overline{\mathbf{W}}_t)dt$  and  $\overline{\Sigma}(\overline{\mathbf{W}}_t)^{\frac{1}{2}}d\mathbf{B}_t^d$  are infinitesimal tangent vectors of  $\mathbb{S}^{d-1}$ . While the parameter  $\overline{\mathbf{W}}_t$  flows along these vector fields according to (7), it by construction stays inside  $\mathbb{S}^{d-1}$ . Indeed, in the construction the diffusion process (7), the differential geometry of  $\mathbb{S}^{d-1}$  is used (Hsu, 2002, §1.3) in addition to the values of vector fields. But (4) is a differential equation defined on the larger ambient manifold  $\mathbb{R}^d$ , and  $\nabla\mathcal{L}(\overline{\mathbf{W}}_t)dt$  and  $\Sigma(\overline{\mathbf{W}}_t)^{\frac{1}{2}}d\mathbf{B}_t^d$  are only vector fields on  $\mathbb{R}^d$ . A priori, (4) moves  $\overline{\mathbf{W}}_t$  around in the entire  $\mathbb{R}^d$  but not necessarily within  $\mathbb{S}^{d-1}$ . However, it turns out that given an initial position  $\overline{\mathbf{W}}_0 \in \mathbb{S}^{d-1}$ , almost all random trajectories of (4) must remain in  $\mathbb{S}^{d-1}$  since (4) arises from (2) with  $\overline{\mathbf{W}}_t = \frac{\mathbf{W}_t}{|\overline{\mathbf{W}}_t|}$ . Theorem 3.1 then asserts the probability distributions of  $\overline{\mathbf{W}}_t$  are the same for (4) and for (7).

In most of this paper,  $\Sigma(\mathbf{w})|_{V_{\mathbf{w}}^{\perp}}$  will be assumed to be the constant matrix  $\sigma^2\text{Id}|_{V_{\mathbf{w}}^{\perp}}$  for  $\mathbf{w} \in \mathbb{S}^{d-1}$ . In this case, (Li et al., 2020, Lemma 5.2) proved  $\gamma_t = \gamma + (\gamma_0 - \gamma)e^{-4\lambda_e t}$  converges exponentially fast to  $\gamma = \frac{\sigma^2}{2\lambda_e}$ . In light of this, we will focus on the SDE model assuming  $\gamma_t$  is the constant  $\gamma$ :

**Definition 3.2.** The spherical model of SGD with constant effective learning rate  $\zeta := \gamma^{-\frac{1}{2}} = \frac{\sqrt{2\lambda_e}}{\sigma}$  is the following

SDE on  $\mathbb{S}^{d-1}$ :

$$d\overline{\mathbf{W}}_t = -\zeta \left( \overline{\nabla}\mathcal{L}(\overline{\mathbf{W}}_t)dt + \overline{\Sigma}(\overline{\mathbf{W}}_t)^{\frac{1}{2}}d\mathbf{B}_t^d \right). \quad (8)$$

We will denote by  $\mathcal{P}_{\overline{\mathbf{W}}_0=\mathbf{w}}(\overline{\mathbf{W}}_t)$  and  $\mathcal{P}_{\overline{\mathbf{W}}_0\sim\nu}(\overline{\mathbf{W}}_t)$  the probability distributions of the random solution  $\overline{\mathbf{W}}_t$  to the equation (8), respectively under initial conditions  $\overline{\mathbf{W}}_0 = \mathbf{w}$  and  $\overline{\mathbf{W}}_0 \sim \nu$ , where  $\mathbf{w} \in \mathbb{S}^{d-1}$  and  $\nu$  is a probability measure on  $\mathbb{S}^{d-1}$ .

#### 3.2. Relation to Riemannian implement of BN

Besides simulating the stabilizing behavior of  $\gamma_t$ , the SDE (8) is also a model to the Riemannian approach to batch normalization introduced by (Cho & Lee, 2017). This approach aims to eliminate the ambiguity in scaling by perform SGD on the Riemannian manifold  $\mathbb{S}^{d-1}$ .

The basic algorithm from (Cho & Lee, 2017) can be simply stated as follows: the sequence of parameters  $\mathbf{w}_k$  will always remain in  $\mathbb{S}^{d-1}$  and be updated by

$$\mathbf{w}_{k+1} \leftarrow \exp_{\mathbf{w}_k}(-\zeta\nabla\mathcal{L}_{\mathcal{B}_k}(\mathbf{w}_k)), \quad (9)$$

where  $\zeta$  is the learning rate and  $\exp_{\mathbf{w}}(\mathbf{v}) := \mathbf{w} \cos|\mathbf{v}| + \mathbf{v} \frac{\sin|\mathbf{v}|}{|\mathbf{v}|}$  is the exponential map on  $\mathbb{S}^{d-1}$  for  $\mathbf{w} \in \mathbb{S}^{d-1}$  and  $\mathbf{v} \in T_{\mathbf{w}}\mathbb{S}^{d-1}$ . Note that  $\nabla\mathcal{L}_{\mathcal{B}_k}(\mathbf{w}_k) \in T_{\mathbf{w}}\mathbb{S}^{d-1}$  as  $\mathbf{w}^{\top}\nabla\mathcal{L}_{\mathcal{B}}(\mathbf{w}, \mathbf{x}) = 0$ . An alternative to (9) is

$$\mathbf{w}_{k+1} \leftarrow \frac{\mathbf{w}_k - \zeta\nabla\mathcal{L}_{\mathcal{B}_k}(\mathbf{w}_k)}{|\mathbf{w}_k - \zeta\nabla\mathcal{L}_{\mathcal{B}_k}(\mathbf{w}_k)|}. \quad (10)$$

The updating methods (9) and (10) differ only by  $O(\zeta^2)$  and hence are asymptotic to each other for small  $\zeta$ , as demonstrated by Theorem 3.3 below. However (10) has the advantage that it is computationally cheaper and the trajectories are guaranteed to stay in  $\mathbb{S}^{d-1}$  despite of numerical errors.

**Theorem 3.3.** *If the vector  $\nabla\mathcal{L}_{\mathcal{B}}(\mathbf{w}) - \nabla\mathcal{L}(\mathbf{w})$  is considered as a Gaussian noise  $\xi(\mathbf{w})$  with covariance matrix  $\Sigma(\mathbf{w})$  defined by (3), then (8) is the continuous time limit of both (9) and (10).*

### 4. Three-stage evolution

We will from now on assume following conditions like in (Bovier et al., 2004; Shi et al., 2020) :

**Assumption 4.1.** (Standard noise) *On  $\mathbb{S}^{d-1}$ , the covariance matrix  $\Sigma$  coincides with a constant multiple of the Riemannian metric on  $\mathbb{S}^{d-1}$ . More precisely, there exists  $\sigma > 0$  such that for all  $\mathbf{w} \in \mathbb{S}^{d-1}$  and  $\mathbf{v} \in V_{\mathbf{w}}^{\perp} = T_{\mathbf{w}}\mathbb{S}^{d-1}$ ,  $\mathbf{v}^{\top}\Sigma(\mathbf{w})\mathbf{v} = \sigma^2|\mathbf{v}|^2$ .*

**Assumption 4.2.** (Morse loss function) *The restriction of  $\mathcal{L}$  to  $\mathbb{S}^{d-1}$  is a function in  $C^2(\mathbb{S}^{d-1})$ . Furthermore, it is a Morse function: every critical point  $\mathbf{z}$  is non-degenerate, i.e.  $\det \overline{\nabla}^2\mathcal{L}(\mathbf{z}) \neq 0$ .*

Note that the Morse function condition is generically true in the class  $C^2(\mathbb{S}^{d-1})$ .

#### 4.1. Fokker-Planck equation and Gibbs density

Because the diffusion matrix  $\bar{\Sigma}(\mathbf{w})^{\frac{1}{2}}(\bar{\Sigma}(\mathbf{w})^{\frac{1}{2}})^{\top} = \bar{\Sigma}(\mathbf{w})$  coincides with  $\sigma^2 \text{Id}|_{T_{\mathbf{w}}\mathbb{S}^{d-1}}$  on  $T_{\mathbf{w}}\mathbb{S}^{d-1}$ , the marginal distribution of the stochastic process (8) is absolutely continuous on  $\mathbb{S}^{d-1}$  for  $t > 0$  and its density function  $u_t(\mathbf{w})$  is a solution for  $t \in (0, \infty)$  and  $\mathbf{w} \in \mathbb{S}^{d-1}$  to the following Fokker-Planck equation by standard arguments on diffusion process (see §B for a proof).

**Proposition 4.3.** *The density function  $u : (0, \infty) \times \mathbb{S}^{d-1} \rightarrow \mathbb{R}$  of the distribution of the stochastic process (8) satisfies*

$$\partial_t u = \zeta \bar{\nabla} \cdot (u \bar{\nabla} \mathcal{L}) + \frac{1}{2} \zeta^2 \sigma^2 \bar{\Delta} u. \quad (11)$$

Define

$$\beta := \frac{1}{2} \zeta \sigma^2 = \sqrt{\frac{\lambda_e}{2}} \sigma. \quad (12)$$

It is known that the Gibbs density

$$\mu^{(\beta)}(\mathbf{w}) := \frac{e^{-\frac{\mathcal{L}(\mathbf{w})}{\beta}}}{\int_{\mathbb{S}^{d-1}} e^{-\frac{\mathcal{L}}{\beta}} \mathbf{d}\mathbf{m}} \quad (13)$$

is a stationary solution to (11).

#### 4.2. Attracting basins

In view of Assumption 4.2, denote the set of critical points of the restriction of  $\mathcal{L}$  on  $\mathbb{S}^{d-1}$  by  $Z \subset \mathbb{S}^{d-1}$ , and the local minima among them by  $z_1, \dots, z_m \in Z$ . Because the critical points are non-degenerate, they must be isolated. Hence  $Z$  is finite.

**Definition 4.4.** The **attracting basin**  $U_z$  of a critical point  $z \in Z$  is the set of all points  $\mathbf{w}_0 \in \mathbb{S}^{d-1}$  such that the unique solution  $\bar{\mathbf{W}}_t$  to the ODE

$$\bar{\mathbf{W}}_t = -\zeta \bar{\nabla} \mathcal{L}(\bar{\mathbf{W}}_t) dt \quad (14)$$

defined on  $\mathbb{S}^{d-1}$ , subject to the initial condition  $\bar{\mathbf{W}}_0 = \mathbf{w}_0$ , converges to  $z$  as  $t \rightarrow \infty$ . For simplicity, when  $z = z_i$  is a local minima, we will write  $U_i := U_{z_i}$ .

For  $Q > 0$ , let  $U_{i,Q} := \{\mathbf{w} \in U_i : \mathcal{L}(\mathbf{w}) - \mathcal{L}(z_i) < Q\}$ , which is a neighborhood of  $z_i$  in  $U_i$ .

#### 4.3. Main theorems

In the analysis below, fix a parameter  $\epsilon > 0$  as error tolerance, and let  $\lambda_e$  be sufficiently small. All choices of parameters, as well as implicit constants in  $O(\cdot)$  notations, are supposed to be dependent of the loss function  $\mathcal{L}$ .

The dynamics of (8) consists of three different stages: **descent**, **diffusion** and **tunneling**.

**Stage 1: Descent.** In this stage, the trajectory, with probability close to 1, takes  $O(\lambda_e^{-\frac{1}{2}})$  time, to descend to  $U_{i,Q_1}$  in each  $U_i$ .

**Theorem 4.5.** *Under Assumptions 4.1 and 4.2, for all  $\epsilon > 0$  and  $Q_1 > 0$ , there exist  $C_{\text{des}} > 0, \lambda_{\text{des}} > 0$  and a set  $\Lambda_\epsilon$  of volume  $\mathbf{m}(\Lambda_\epsilon) > 1 - \epsilon$ , such that for all  $\lambda_e < \lambda_{\text{des}}$ , and all  $\mathbf{w}_0 \in \Lambda_\epsilon$ , the random solution  $\bar{\mathbf{W}}_t$  to (8) starting at  $\mathbf{w}_0$  satisfies  $\mathbb{P}_{\bar{\mathbf{W}}_0 = \mathbf{w}_0} \left( \bar{\mathbf{W}}_{C_{\text{des}} \lambda_e^{-\frac{1}{2}}} \in U_{k,Q_1} \right) > 1 - \epsilon$ , where  $U_k$  is the unique attracting basin that contains  $\mathbf{w}_0$ .*

**Stage 2: Diffusion.** The diffusion stage takes at most  $O(\lambda_e^{-1})$  time in terms of  $t$ . During this period, the distribution converges to a **temporary equilibrium** given by the conditional measure of the Gibbs equilibrium inside each basin and remains stable for exponentially long time after that. The weights assigned to basins correspond to the initial distribution of mass among them. The next theorem is the central one in this paper.

**Theorem 4.6.** *Under Assumptions 4.1 and 4.2, for all  $\epsilon > 0$ , there exist constants  $C_{\text{dif}}, c_{\text{dif}}, \lambda_{\text{dif}} > 0$ , and a set  $\Lambda_\epsilon$  of volume  $\mathbf{m}(\Lambda_\epsilon) > 1 - \epsilon$ , such that:*

With  $\beta = \sqrt{\frac{\lambda_e}{2}} \sigma$ , for all  $\lambda_e < \lambda_{\text{dif}}$ , the random solution  $\bar{\mathbf{W}}_t$  to (8) satisfies: for all  $t \in [\frac{C_{\text{dif}}}{\lambda_e}, e^{\frac{c_{\text{dif}}}{\sqrt{\lambda_e}}}]$ :

(i) For all initial positions  $\mathbf{w}_0 \in \Lambda_\epsilon$ ,

$$\text{dist}_{\text{TV}} \left( \mathcal{P}_{\bar{\mathbf{W}}_0 = \mathbf{w}_0}(\bar{\mathbf{W}}_t), \frac{\mu^{(\beta)}|_{U_k}}{\int_{U_k} \mu^{(\beta)} \mathbf{d}\mathbf{m}} \mathbf{d}\mathbf{m} \right) \leq \epsilon,$$

where  $U_k$  is the unique attracting basin of the gradient flow of  $\mathcal{L}$  that contains  $\mathbf{w}_0$ .

(ii) For all initial probability distribution  $\nu_0$ ,

$$\text{dist}_{\text{TV}} \left( \mathcal{P}_{\bar{\mathbf{W}}_0 \sim \nu_0}(\bar{\mathbf{W}}_t), \sum_{i=1}^m \nu_0(U_i) \frac{\mu^{(\beta)}|_{U_i}}{\int_{U_i} \mu^{(\beta)} \mathbf{d}\mathbf{m}} \mathbf{d}\mathbf{m} \right) \leq \epsilon + \nu_0(\Lambda_\epsilon^c).$$

Here  $\text{dist}_{\text{TV}}$  is the total variation distance  $\text{dist}_{\text{TV}}(\mu, \nu) = \sup_A |\mu(A) - \nu(A)|$  between measures, and  $\Lambda_\epsilon^c$  is the complement of  $\Lambda_\epsilon$ .

Theorem 4.6(ii) identifies a temporary equilibrium  $\sum_{i=1}^m \nu_0(U_i) \frac{\mu^{(\beta)}|_{U_i}}{\int_{U_i} \mu^{(\beta)} \mathbf{d}\mathbf{m}} \mathbf{d}\mathbf{m}$ . The corollary below claims that in the case of non-convex optimization, this temporary equilibrium is often different from the eventual Gibbs equilibrium  $\mu^{(\beta)} \mathbf{d}\mathbf{m}$ .

**Corollary 4.7.** *In the setting of Theorem 4.6, if in addition  $m \geq 2$ , then there exists a constant  $\kappa_{\text{dif}} > 0$  and a subset  $\Omega_{\text{dif}}$  with  $\mathbf{m}(\Omega_{\text{dif}}) > \kappa_{\text{dif}}$ , such that for all sufficiently small*

$\lambda_e, \mathbf{w}_0 \in \Omega_{\text{dif}}$  and  $t \in [\frac{C_{\text{dif}}}{\lambda_e}, e^{\frac{C_{\text{dif}}}{\lambda_e}}]$ , the random solution  $\overline{\mathbf{W}}_t$  to (8) satisfies

$$\text{dist}_{\text{TV}}\left(\mathcal{P}_{\overline{\mathbf{W}}_0=\mathbf{w}_0}(\overline{\mathbf{W}}_t), \mu^{(\beta)} \mathbf{d}\mathbf{m}\right) > \kappa_{\text{dif}}.$$

Moreover, if  $\mathcal{L}$  attains its global minimum value on  $\mathbb{S}^{d-1}$  at more than one local minima, then given any small  $\epsilon > 0$ , one can replace  $\mathbf{m}(\Omega_{\text{dif}}) > \kappa_{\text{dif}}$  by  $\mathbf{m}(\Omega_{\text{dif}}) > 1 - \epsilon$  for sufficiently small  $\lambda_e$ .

The exceptional set  $\Lambda_\epsilon^c$  in Theorems 4.5 and 4.6 is a small neighborhood of the set of points whose deterministic gradient descent trajectory converge to a stationary point of  $\mathcal{L}$  that is not a local minimum. This later set is a submanifold of strictly lower dimension under Assumption 4.2, and should be viewed as the boundary between attracting basins of local minima. A careful analysis of the arguments in Appendix §D would control the radius of the neighborhood in terms of  $\epsilon$ , the values of  $\mathcal{L}$  and  $\overline{\nabla}^2 \mathcal{L}$  at the stationary points of  $\mathcal{L}$  and the Lipschitz norm of  $\mathcal{L}$ .

**Stage 3: Tunneling.** The final stage takes  $O(e^{\frac{C}{\lambda_e}})$  time. During this stage, mass leaks slowly between basins, and eventually equidistributes towards the Gibbs equilibrium  $\mu^{(\beta)} \mathbf{d}\mathbf{m}$ . Because of the slow rate, this stage is not expected to be observed within typical training time.

The name ‘‘tunneling’’, following previous works, e.g. (Helffer & Sjöstrand, 1985; Hérau et al., 2011), comes from the quantum tunnel effect in solutions to the Schrödinger equation, which is related to our model through the Schrödinger operator  $\mathcal{D}^{\beta, \#} f = \beta \overline{\Delta} f - \left(\frac{|\overline{\nabla} \mathcal{L}|^2}{4\beta} - \frac{\overline{\Delta} \mathcal{L}}{2}\right) f$  (see §E.1).

**Theorem 4.8.** *Under Assumptions 4.1 and 4.2, and with  $\beta = \sqrt{\frac{\lambda_e}{2}} \sigma$ , there exist constants  $C_{\text{tun}}, \lambda_{\text{tun}} > 0$  such that for all  $\mathbf{w}_0 \in \mathbb{S}^{d-1}$ ,  $\lambda_e < \lambda_{\text{tun}}$  and  $t \geq 0$ , the random solution  $\overline{\mathbf{W}}_t$  to (8) satisfies :*

$$\text{dist}_{\text{TV}}\left(\mathcal{P}_{\overline{\mathbf{W}}_0=\mathbf{w}_0}(\overline{\mathbf{W}}_t), \mu^{(\beta)} \mathbf{d}\mathbf{m}\right) \leq O\left(e^{-e^{-\frac{C_{\text{tun}}}{\sqrt{\lambda_e} t}}}\right).$$

We don’t claim originality for Theorem 4.8 and only include it for a complete description of the stages. It was proved by (Shi et al., 2020) for fast growing loss functions  $\mathcal{L}$  on  $\mathbb{R}^d$ , and their argument also works here (see §F).

Instead, our main contribution is Theorem 4.6, which identifies the temporary equilibrium and proves the time needed to reach it is  $\leq O(\lambda_e^{-1})$ , as well as Corollary 4.7, which asserts that typically the temporary and eventual equilibria are different, and the time needed to deviate from the temporary equilibrium towards the eventual one is at least  $O(e^{\frac{C_{\text{dif}}}{\lambda_e}})$ , and thus usually beyond practical observable windows. The contrast between Corollary 4.7 and Theorem 4.8 establishes the separation between the diffusion and tunneling stages.

## 5. Experiments

**Model with our assumptions.** We implemented the SDE (8) on  $\mathbb{S}^{d-1}$  with  $d = 100$ ,  $\sigma = 0.1$ , and a scaling-invariant loss function  $\mathcal{L}(\mathbf{w}) = \sum_{i=1}^{200} \sin(\mathbf{a}_i \cdot \frac{\mathbf{w}}{|\mathbf{w}|})^2$  where  $\mathbf{a}_i$  are randomly chosen. For each  $\zeta$ , we ran 32 independent random instances of the discrete implementation

$$\mathbf{w}_{k+1} \leftarrow \frac{\mathbf{w}_k - \zeta (\nabla \mathcal{L}(\mathbf{w}_k) + \mathcal{N}_d(0, \sigma))}{|\mathbf{w}_k - \zeta (\nabla \mathcal{L}(\mathbf{w}_k) + \mathcal{N}_d(0, \sigma))|} \quad (15)$$

of (8). The same proof of Theorem 3.3 shows (8) is a continuous time limit of (15). All instances start at the same initial position and last  $8 \times 10^5$  iteration steps.

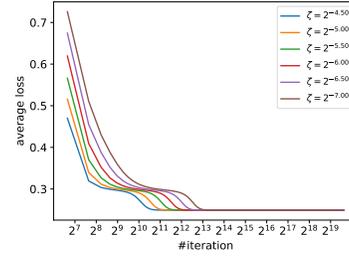


Figure 1. Train loss [Toy function experiment]

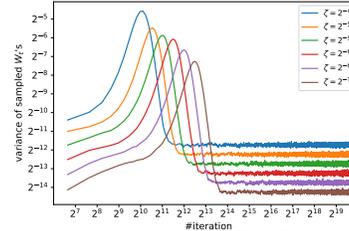


Figure 2. Variance among instances [Toy function experiment]

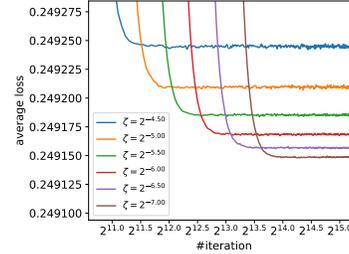


Figure 3. Train loss, zoom in view [Toy function experiment]

We choose two features to indicate the achievement of an equilibrium: the average loss among the instances, as well as the variance of the weights  $\mathbf{W}_t$  of all the instances. Figures 1-4 show that when  $\zeta$  is divided by  $2^{0.5}$ , the step at which these features become approximately constant is roughly multiplied by  $2^{0.5}$ . That is, equilibrium is observed in  $O(\zeta^{-1})$  time. As  $\zeta = \frac{\sqrt{2\lambda_e}}{\sigma}$ ,  $O(\zeta^{-1}) = O(\lambda_e^{-\frac{1}{2}}) < O(\lambda_e^{-1})$ , this is consistent with Theorem 4.6 and suggests there is room for further improvement (Question 6.2 below).

Experiments also support that the observed equilibrium is the one in Theorem 4.6. First of all, the observed equilib-

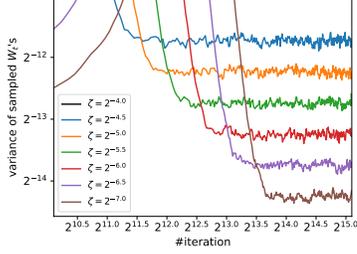
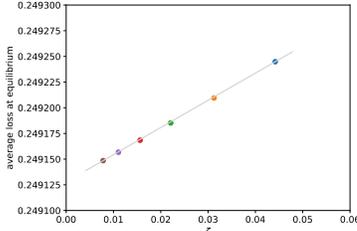
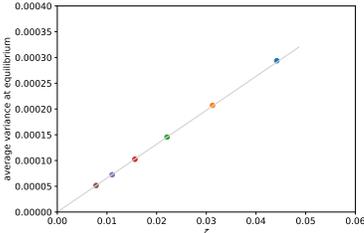


Figure 4. Variance, zoom in view [Toy function experiment]

rium is localized near a single local minimum point, as the variance among instances become very close to 0.

In addition, note  $\mu_k^{(\beta)} := \frac{\mu^{(\beta)}|_{U_k}}{\int_{U_k} \mu^{(\beta)} d\mathbf{m}}$  is concentrated near  $\mathbf{z}_k$  for small  $\beta$ . As  $\mathbf{z}_k$  is a non-degenerate local minimum, in suitable local linear coordinates  $\mathbf{y} \in T_{\mathbf{z}_k} \mathbb{S}^{d-1}$  for  $\mathbf{w}$  near  $\mathbf{z}_k$ ,  $\mathcal{L}(\mathbf{w}) = \mathcal{L}(\mathbf{z}_k) + \frac{1}{2}|\mathbf{y}|^2 + O(|\mathbf{y}|^3)$ . Then  $\mu_k^{(\beta)}$  is approximated by  $\frac{e^{-\frac{1}{2\beta}|\mathbf{y}|^2}}{\text{constant}} d\mathbf{y}$ , which is the normal distribution  $\mathcal{N}(0, \beta \text{Id})$  in terms of  $\mathbf{y}$ . Since  $\mathbf{w}$  is approximately an affine function of  $\mathbf{y}$ ,  $\text{Var}_{\mu_k^{(\beta)}}(\mathbf{w})$  should be approximately proportional to  $\beta = \frac{1}{2}\zeta\sigma^2$ , and thus to  $\zeta$  as well. In addition,  $\mathbb{E}_{\mu_k^{(\beta)}} \mathcal{L}(\mathbf{w}) - \mathcal{L}(\mathbf{z}_k) \approx \frac{1}{2} \mathbb{E}_{\mu_k^{(\beta)}} |\mathbf{y}|^2 = \frac{1}{2} \text{Var}(\mathbf{y}) \approx \frac{1}{2} \text{Tr} \beta \text{Id} = \frac{d-1}{2} \beta$  is also proportional to both  $\beta$  and  $\zeta$ .


 Figure 5. Train loss at equilibrium vs  $\zeta$  [Toy function experiment]

 Figure 6. Variance of equilibrium vs  $\zeta$  [Toy function experiment]

Figures 5 and 6 show that near the observed equilibrium, the training loss  $\mathbb{E}\mathcal{L}(\mathbf{W}_t)$  is approximately linear in  $\zeta$  and  $\text{Var}(\mathbf{W}_t)$  is approximately proportional to  $\zeta$ , which perfectly matches the predictions above.

**Underparametrized neural network.** We apply (10) to a tiny 4-layer CNN network of the same structure as above, but with on the MNIST dataset of only  $d = 332$  parameters. There are two convolution layer followed by two linear layers. All layers are batch normalized to ensure

scale-invariance. The number of parameters is chosen to be unpractically low for two reasons: (1) the underparametrization of the network guarantees that Assumption 4.2 still holds and local minima are isolated; (2) to allow faster training, as we need to run many independent instances and observe their distribution. However, Assumption 4.1 no longer holds in this case.

Due to the small size of the network, the achieved minimum of the cross entropy loss function is much worse than those in practical trainings. But our goal here is to empirically test the likelihood of whether our results extend to settings satisfying Assumption 4.2 but not Assumption 4.1.

For each learning rate  $\zeta$ , 16 independent random instances of (10) with the same initial position are performed, for 20,000 epochs. Each epoch has 54 iteration steps with a batch size of 1000.

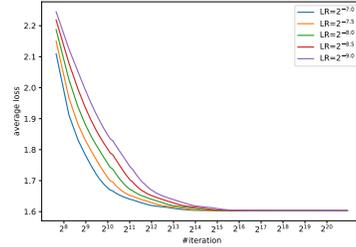


Figure 7. Train loss [Underparametrized NN]

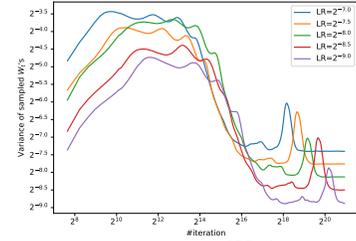


Figure 8. Variance among instances [Underparametrized NN]

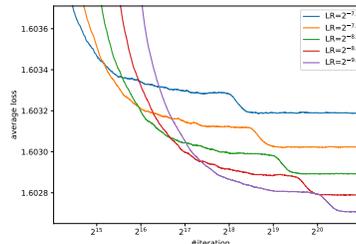


Figure 9. Train loss, zoom in view [Underparametrized NN]

As before, we use the average loss as well as the variance among parameters as indicators for reaching equilibrium. Figures 7-10<sup>1</sup> again show that random trajectories with the same initial position stabilize to a local equilibrium

<sup>1</sup>Curves in Figures 7-10, as well as in Figure 14 later, are logarithmically smoothed by displaying at  $n$  the average values from epoch  $0.9n$  to epoch  $n$ .

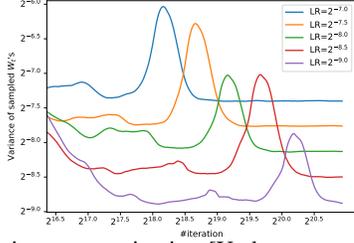
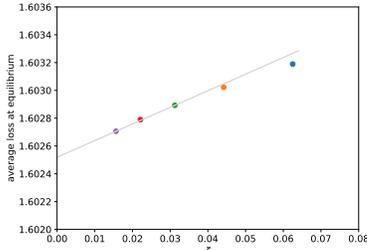
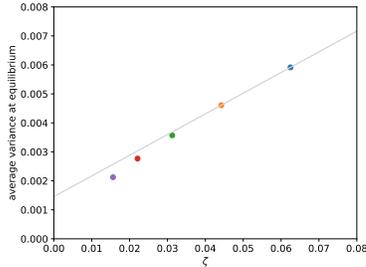


Figure 10. Variance, zoom in view [Underparametrized NN]

within  $O(\zeta^{-1}) = O(\lambda_e^{-\frac{1}{2}})$  time. This local minimum is again near a single critical point as  $\text{Var}_{\mu_k^{(\beta)}}(\mathbf{w})$  is extremely small. As  $\zeta$  decreases, the loss and variance decay in similar patterns as in the previous model. However, Figure 12 shows  $\text{Var}_{\mu_k^{(\beta)}}(\mathbf{w})$  is not exactly proportional to  $\alpha$ , but decreases to 0 faster than  $\eta$ , which suggests that with Assumption 4.2 but not Assumption 4.1, the local equilibrium in general doesn't take the form  $\frac{\mu^{(\beta)}|_{U_k}}{\int_{U_k} \mu^{(\beta)} d\mathbf{m}}$  than in the standard noise setting. See Conjecture 6.3 below.


 Figure 11. Train loss at equilibrium vs  $\zeta$  [Underparametrized NN]

 Figure 12. Variance of equilibrium vs  $\zeta$  [Underparametrized NN]

**Overparametrized neural network.** Finally we apply (10) to a 4-layer CNN network of the same structure as above, but with  $d \approx 100k$  parameters instead, on the MNIST. This general case is very different from the setting we studied as it is highly overparametrized and the local minima are likely submanifolds of large dimensions and sizes. Neither Assumptions 4.1 nor 4.2 is expected to hold. In particular, we no longer expect to see equilibria concentrated around single critical points.

In this setting it takes a huge number ( $> 10k$ ) of epochs before statistics stabilize, and each epoch runs for longer time. Hence it was infeasible within our training budget to carry out similar experiments with many independent instances

across several learning rates like in the two previous cases. However, in this case we want to know: Starting from a given initial position, is the fast equilibrium for (7) localized inside a single basin, and different from the global Gibbs equilibrium, like in Theorem 4.6? In other words, does the separation between the diffusion and tunneling stages extend to the overparametrized case?

Our experiment suggests that the answer is yes. At small effective learning rate, within observable time windows different initial positions lead to different equilibria that are far from each other. Since they cannot both be the unique Gibbs equilibrium, these equilibria are both local and temporary in nature. Starting from each of two randomly chosen initial positions  $\mathbf{w}_0^1, \mathbf{w}_0^2 \in \mathbb{S}^{d-1}$ , we ran 5 instances of the implementation (10) of (7) for 40000 epochs ( $\approx 2M$  iterations) at  $\zeta = 2^{-4.5}$ . All instances use independent random seeds. For  $1 \leq i \leq 5$  and  $p = 1, 2$ , let  $\mathbf{w}_k^{i,p} \in \mathbb{S}^{d-1}$  denote the parameter at  $k$ -th step of the  $i$ -th instance starting at  $\mathbf{w}_0^p$ . We tracked the average distance squares between pairs of parameters that originated from the same initial position, as well as from different initial positions, or more precisely the quantities  $\mathcal{V}_{11}(k) := \mathbb{E}_{i \neq j} |\mathbf{w}_k^{i,1} - \mathbf{w}_k^{j,1}|^2$ ;  $\mathcal{V}_{22}(k) := \mathbb{E}_{i \neq j} |\mathbf{w}_k^{i,2} - \mathbf{w}_k^{j,2}|^2$ ;  $\mathcal{V}_{12}(k) := \mathbb{E}_{i,j} |\mathbf{w}_k^{i,1} - \mathbf{w}_k^{j,2}|^2$ . where the parameters are regarded as vectors in  $\mathbb{R}^d$ .

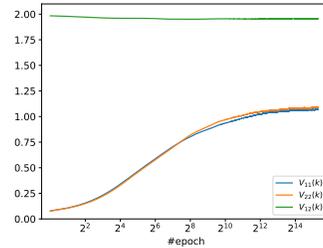


Figure 13. Different temporary equilibria [Overparametrized NN]

Figure 13 shows: (1) The equilibrium derived from an initial position is not concentrated near a single point, as  $\mathcal{V}_{11}$  and  $\mathcal{V}_{22}$  do not approach 0; (2) The two equilibria derived from  $\mathbf{w}_0^1$  and  $\mathbf{w}_0^2$  are far from each other, because  $\mathcal{V}_{12}$  is much larger than  $\mathcal{V}_{11}$  and  $\mathcal{V}_{22}$ . Indeed,  $\mathcal{V}_{12} \approx 2$  for all time. Since the distance between two uniformly chosen random points on  $\mathbb{S}^{d-1}$  is  $\approx \sqrt{2}$  when  $d$  is large, this suggests the two equilibria are independently located.

## 6. Discussions

### 6.1. Interpretations of initial parameter independence in Conjecture 1.1

Our Theorem 4.6 shows that, in the parameter space the equilibrium reached in  $O(\frac{1}{\lambda_e})$  time is mostly concentrated in the attracting basin of the initial parameter. While this

is supported by real world neural network implementations (Figure 13), one needs to explain such deviation from the independence on the initial parameter in Conjecture 1.1 of (Li et al., 2020). Two possible interpretations are sketched below.

**Initial mass distribution.** Consider SGD with batch normalization starting from an initial distribution  $\nu_0$  of parameters and let  $\nu_t$  be the distribution at moment  $t$ . Theorem 4.6 shows that equilibrium can be reached in  $t_1 \leq O(\frac{1}{\lambda_e})$  steps after  $\gamma_t^{-\frac{1}{2}}$  stabilizes. Moreover, the equilibrium is a linear combination  $\sum_k p_k \mu_k^{(\beta)}$  of local equilibria  $\mu_k^{(\beta)}$  in basins  $U_k$ . The allocation of weights  $p_k$  is approximately the same as the initial allocation for this stage, or in other words the allocation at the end of the preparatory phase when  $\gamma_t$  stabilizes. (Li et al., 2020) suggests this preparatory phase takes  $t_0 \approx O(\frac{1}{\lambda_e})$  steps, during which  $\gamma_t$  stabilizes exponentially fast (at least when  $\text{Tr}\Sigma$  is constant). Experiments in (Li et al., 2020, Figure 5) also show that the initial effective learning rate  $\gamma_0^{-\frac{1}{2}}$  is typically much larger than the limit value at  $t_0$ .

When  $t$  is small, the effective learning rate  $\gamma_t^{-\frac{1}{2}}$  is large for some time. Though this period is short, the large learning rate allows mixing in short time and reaches a global equilibrium (which we think of as  $\mu^{(\beta_0)}$  with some very large  $\beta_0$ ). This process fixes the allocation of mass  $p_k$  among basins  $U_k$ . Once the effective learning rate becomes small, mechanisms similar to Theorem 4.6 prevent mass from leaking between basins. Finally, after  $\gamma_t^{-\frac{1}{2}}$  stabilizes, the distribution is locally fine tuned to local Gibbs equilibria  $\mu_k^{(\beta)}$  with small  $\beta$  without changing the weights  $p_k$ .

In summary, in the fast equilibrium  $\sum_k p_k \mu_k^{(\beta)}$ , the components  $\mu_k^{(\beta)}$ 's are determined later at small effective learning rates, but the weights  $p_k$ 's are determined earlier at large effective learning rates.

**Similar landscapes among local minima.** We are grateful to an anonymous reviewer for suggesting to us that the use of an output function such as testing/training loss might be responsible for the independence in the initial parameter in Conjecture 1.1. On the one hand, for an arbitrary loss function  $\mathcal{L}$ , this is likely not the case because  $\mathcal{L}$  may have two attracting basins with different local minimum values  $\mathcal{L}(x_1)$ ,  $\mathcal{L}(x_2)$ . Different initial parameters in these basins will output the corresponding training losses. On the other hand, further experiments show that in realistic neural network implementations the use of an output function such as the loss  $\mathcal{L}$  might explain the initial parameter independence. For instance, for the same initial parameters  $w_0^1, w_0^2$  in Figure 13, though their stochastic trajectories converge in the diffusions stage to different temporary equilibria in the parameter space, Figure 14 shows that the loss function  $\mathcal{L}$  has similar

distribution over these temporary equilibria. More precisely, the quantities  $\mathcal{V}\mathcal{L}_{11}(k) := \mathbb{E}_{i \neq j} |\mathcal{L}(w_k^{i,1}) - \mathcal{L}(w_k^{j,1})|^2$ ,  $\mathcal{V}\mathcal{L}_{22}(k) := \mathbb{E}_{i \neq j} |\mathcal{L}(w_k^{i,2}) - \mathcal{L}(w_k^{j,2})|^2$  and  $\mathcal{V}\mathcal{L}_{12}(k) := \mathbb{E}_{i,j} |\mathcal{L}(w_k^{i,1}) - \mathcal{L}(w_k^{j,2})|^2$  are all distributed near 0 with the same pattern. This suggests the following possible interpretation of initial parameter independence: in an overparametrized scaling-invariant neural network, a majority portion of the parametric space  $\mathbb{S}^{d-1}$  might be covered by attracting basins on which the landscapes of  $\mathcal{L}$  are approximately the same.

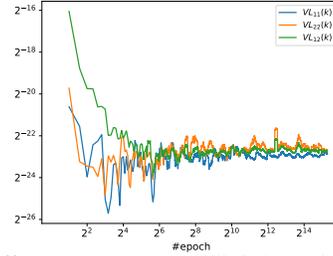


Figure 14. Different temporary equilibria have similar distributions of loss function [Overparametrized NN]

## 6.2. Open questions and future directions

The following questions might be of interest.

**Question 6.1.** Analyze either one or both of the interpretations in §6.1.

Though Theorem 4.6 shows the distribution stabilizes near temporary equilibria within  $O(\frac{1}{\lambda_e})$  time, in our experiments, most features only require  $O(\frac{1}{\sqrt{\lambda_e}})$  time to stabilize with constant  $\gamma_t^{-\frac{1}{2}}$ . Note that this wouldn't affect Conjecture 1.1 as the convergence of  $\gamma_t^{-\frac{1}{2}}$  itself needs  $O(\frac{1}{\lambda_e})$  time.

**Question 6.2.** With Assumptions 4.1 and 4.2, can the bound  $t \geq O(\frac{1}{\lambda_e})$  in Theorem 4.6 be improved to  $t \geq O(\frac{1}{\sqrt{\lambda_e}})$ ?

In light of the experiments in §5, we conjecture that

**Conjecture 6.3.** Without Assumptions 4.1 and 4.2, around each connected sets  $Z_k \subset \mathbb{S}^{d-1}$  of local minima of  $\mathcal{L}$ , there are probability measures  $\mu_k^{(\beta)}$  supported on attracting basins, such that  $\mathbb{E}_{w \sim \mu_k^{(\beta)}} \text{dist}(w, Z_k) \rightarrow 0$  as  $\beta \rightarrow 0$ ; and Theorem 4.6 and Corollary 4.7 hold after replacing

$$\frac{\mu^{(\beta)}|_{U_k}}{\int_{U_k} \mu^{(\beta)} dm} dm \text{ with } \mu_k^{(\beta)}.$$

## Acknowledgments

Y. W.'s research and the use of GPU units were supported by NSF grant DMS-1845033. Z.W.'s research was supported by NSF grant DMS-1753042. We thank the anonymous reviewers for their insightful comments and helpful suggestions.

## References

- Benton, G., Maddox, W., Lotfi, S., and Wilson, A. G. G. Loss surface simplexes for mode connecting volumes and fast ensembling. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 769–779. PMLR, 18–24 Jul 2021.
- Bovier, A., Eckhoff, M., Gaynard, V., and Klein, M. Metastability in reversible diffusion processes i: Sharp asymptotics for capacities and exit times. *Journal of the European Mathematical Society*, 006(4):399–424, 2004.
- Chaudhari, P. and Soatto, S. Stochastic gradient descent performs variational inference, converges to limit cycles for deep networks. In *International Conference on Learning Representations*, 2018.
- Cho, M. and Lee, J. Riemannian approach to batch normalization. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/file/3a0844cee4fcf57de0c71e9ad3035478-Paper.pdf>.
- Cooper, Y. Global minima of overparameterized neural networks. *SIAM Journal on Mathematics of Data Science*, 3(2):676–691, 2021.
- Freidlin, M. I. and Wentzell, A. D. *Random perturbations of dynamical systems*, volume 260 of *Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*. Springer, Heidelberg, third edition, 2012. ISBN 978-3-642-25846-6. doi: 10.1007/978-3-642-25847-3. URL <https://doi.org/10.1007/978-3-642-25847-3>. Translated from the 1979 Russian original by Joseph Szücs.
- Garipov, T., Izmailov, P., Podoprikin, D., Vetrov, D. P., and Wilson, A. G. Loss surfaces, mode connectivity, and fast ensembling of dnns. In *Advances in Neural Information Processing Systems*, pp. 8789–8798, 2018.
- Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., and He, K. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- Helfffer, B. and Sjöstrand, J. Effet tunnel pour l’opérateur de Schrödinger semi-classique. I. In *Proceedings of the conference on partial differential equations, Vol. 1, 2 (Saint Jean de Monts, 1985)*, pp. Exp. No. 13, 38. Soc. Math. France, Paris, 1985.
- Hérau, F., Hitrik, M., and Sjöstrand, J. Tunnel effect and symmetries for Kramers-Fokker-Planck type operators. *J. Inst. Math. Jussieu*, 10(3):567–634, 2011. ISSN 1474-7480. doi: 10.1017/S1474748011000028. URL <https://doi.org/10.1017/S1474748011000028>.
- Horn, R. A. and Johnson, C. R. *Matrix analysis*. Cambridge University Press, Cambridge, second edition, 2013. ISBN 978-0-521-54823-6.
- Hsu, E. P. *Stochastic analysis on manifolds*, volume 38 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, 2002. ISBN 0-8218-0802-8. doi: 10.1090/gsm/038. URL <https://doi.org/10.1090/gsm/038>.
- Huang, Z. and Becker, S. Stochastic gradient langevin dynamics with variance reduction. *arXiv:2102.06759*, 2021.
- Hwang, C.-R. Laplace’s method revisited: weak convergence of probability measures. *Ann. Probab.*, 8(6):1177–1182, 1980.
- Izmailov, P., Podoprikin, D., Garipov, T., Vetrov, D., and Wilson, A. Averaging weights leads to wider optima and better generalization. In Silva, R., Globerson, A., and Globerson, A. (eds.), *34th Conference on Uncertainty in Artificial Intelligence 2018, UAI 2018*, 34th Conference on Uncertainty in Artificial Intelligence 2018, UAI 2018, pp. 876–885. Association For Uncertainty in Artificial Intelligence (AUAI), January 2018. 34th Conference on Uncertainty in Artificial Intelligence 2018, UAI 2018 ; Conference date: 06-08-2018 Through 10-08-2018.
- Jastrzebski, S., Kenton, Z., Arpit, D., Ballas, N., Fischer, A., Bengio, Y., and Storkey, A. Three factors influencing minima in sgd. *arXiv preprint arXiv:1711.04623*, 2017.
- Kolokoltsov, V. N. *Semiclassical analysis for diffusions and stochastic processes*, volume 1724 of *Lecture Notes in Mathematics*. Springer-Verlag, Berlin, 2000. ISBN 3-540-66972-8. doi: 10.1007/BFb0112488. URL <https://doi.org/10.1007/BFb0112488>.
- Kuditipudi, R., Wang, X., Lee, H., Zhang, Y., Li, Z., Hu, W., Ge, R., and Arora, S. Explaining landscape connectivity of low-cost solutions for multilayer nets. In Wallach, H., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/46a4378f835dc8040c8057beb6a2da52-Paper.pdf>.

- Li, Q., Tai, C., and E, W. Stochastic modified equations and dynamics of stochastic gradient algorithms i: Mathematical foundations. *Journal of Machine Learning Research*, 20(40):1–47, 2019. URL <http://jmlr.org/papers/v20/17-526.html>.
- Li, Z., Lyu, K., and Arora, S. Reconciling modern deep learning with traditional optimization analyses: The intrinsic learning rate. *Advances in Neural Information Processing Systems 33 pre-proceedings (NeurIPS 2020)*, abs/2010.02916, 2020.
- Li, Z., Malladi, S., and Arora, S. On the validity of modeling SGD with stochastic differential equations (SDEs). In Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, 2021a. URL [https://openreview.net/forum?id=goEdyJ\\_nVQI](https://openreview.net/forum?id=goEdyJ_nVQI).
- Li, Z., Wang, T., and Arora, S. What happens after sgd reaches zero loss? –a mathematical framework, 2021b.
- Maddox, W. J., Benton, G., and Wilson, A. G. Rethinking parameter counting in deep models: Effective dimensionality revisited. *arXiv:2003.02139*, 2020.
- Mandt, S., Hoffman, M. D., and Blei, D. M. Stochastic gradient descent as approximate bayesian inference. *The Journal of Machine Learning Research*, 18(1):4873–4907, 2017.
- Michel, L. About small eigenvalues of the Witten Laplacian. *Pure Appl. Anal.*, 1(2):149–206, 2019. ISSN 2578-5885. doi: 10.2140/paa.2019.1.149. URL <https://doi.org/10.2140/paa.2019.1.149>.
- Raginsky, M., Rakhlin, A., and Telgarsky, M. Non-convex learning via stochastic gradient langevin dynamics: a nonasymptotic analysis. In Kale, S. and Shamir, O. (eds.), *Proceedings of the 2017 Conference on Learning Theory*, volume 65 of *Proceedings of Machine Learning Research*, pp. 1674–1703. PMLR, 07–10 Jul 2017. URL <https://proceedings.mlr.press/v65/raginsky17a.html>.
- Shi, B., Su, W. J., and Jordan, M. I. On learning rates and schrödinger operators. *arXiv preprint arXiv:2004.06977*, 2020.
- Simon, B. Semiclassical analysis of low lying eigenvalues. I. Nondegenerate minima: asymptotic expansions. *Ann. Inst. H. Poincaré Sect. A (N.S.)*, 38(3):295–308, 1983. ISSN 0246-0211.
- Smith, S. L. and Le, Q. V. A bayesian perspective on generalization and stochastic gradient descent. In *International Conference on Learning Representations*, 2018.
- Smith, S. L., Kindermans, P.-J., and Le, Q. V. Don’t decay the learning rate, increase the batch size. In *International Conference on Learning Representations*, 2018.
- Smith, S. L., Dherin, B., Barrett, D., and De, S. On the origin of implicit regularization in stochastic gradient descent. In *International Conference on Learning Representations*, 2021. URL [https://openreview.net/forum?id=rq\\_Qr0clHyo](https://openreview.net/forum?id=rq_Qr0clHyo).
- Xu, P., Chen, J., Zou, D., and Gu, Q. Global convergence of langevin dynamics based algorithms for nonconvex optimization. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pp. 3126?3137, 2018.
- Yaida, S. Fluctuation-dissipation relations for stochastic gradient descent. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=SkNksoRctQ>.
- Zhang, Y., Liang, P., and Charikar, M. A hitting time analysis of stochastic gradient langevin dynamics. In Kale, S. and Shamir, O. (eds.), *Proceedings of the 2017 Conference on Learning Theory*, volume 65 of *Proceedings of Machine Learning Research*, pp. 1980–2022. PMLR, 07–10 Jul 2017.

## A. Proofs of results in §3

*Proof of Theorem 3.1.* First remark that given any initial position  $\overline{\mathbf{W}}_0 \in \mathbb{S}^{d-1}$ , the solution to (4) almost surely stays in  $\mathbb{S}^{d-1}$  by (Li et al., 2020, Thm 5.1).

Write  $\Sigma(\mathbf{w})^{\frac{1}{2}} = (\mathbf{a}_1(\mathbf{w}), \dots, \mathbf{a}_d(\mathbf{w}))$ . Note that for  $\mathbf{w} \in \mathbb{S}^{d-1}$ ,  $\mathbf{a}_j(\mathbf{w}) \in T_{\mathbf{w}}\mathbb{S}^{d-1}$ . Then the random solution  $\overline{\mathbf{W}}_t$  to (7) satisfies the following property:

$$\frac{d}{dt}\mathbb{E}f(\overline{\mathbf{W}}_t) = \mathbb{E}(\overline{\Psi}f(\overline{\mathbf{W}}_t)), \quad (16)$$

where  $\overline{\Psi}$  is a differential operator on  $\mathbb{S}^{d-1}$  defined by

$$\overline{\Psi}f(\mathbf{w}) = -\gamma_t^{-\frac{1}{2}}\nabla\mathcal{L}(\mathbf{w}) \cdot \nabla f(\mathbf{w}) + \frac{1}{2}\gamma_t^{-1} \sum_{j=1}^d (\overline{\nabla}_{\mathbf{a}_j})^2 f(\mathbf{w}). \quad (17)$$

where  $\overline{\nabla}_{\mathbf{a}_j}$  is the covariant derivative along the vector field  $\mathbf{a}_j$  on  $\mathbb{S}^{d-1}$  with respect to the spherical metric (see (Hsu, 2002, Theorem 1.3.4)). The operator  $\overline{\Psi}$  is called the generator of (7). Similarly, the random solution  $\overline{\mathbf{W}}_t$  to (4) satisfies:  $\frac{d}{dt}\mathbb{E}f(\overline{\mathbf{W}}_t) = \mathbb{E}(\Psi f(\overline{\mathbf{W}}_t))$  where  $\Psi$  is a differential operator on  $\mathbb{R}^d$  defined by

$$\begin{aligned} \Psi f(\mathbf{w}) &= -\gamma_t^{-\frac{1}{2}}\nabla\mathcal{L}(\mathbf{w}) \cdot \nabla f(\mathbf{w}) \\ &+ \frac{1}{2}\gamma_t^{-1} \sum_{j=1}^d \nabla_{\mathbf{a}_j}^2 f(\mathbf{w}) \\ &- \frac{1}{2}\gamma_t^{-1} \text{Tr}\Sigma(\mathbf{w})\mathbf{w} \cdot \nabla f(\mathbf{w}). \end{aligned} \quad (18)$$

We emphasize that besides the removal of the last term,  $\overline{\Psi}$  differs from  $\Psi$  in that differential operators and the Riemannian metric on  $\mathbb{S}^{d-1}$  are used instead of those on  $\mathbb{R}^{d-1}$ .

We next analyze the difference between corresponding terms in (17) and (18) at  $\mathbf{w} \in \mathbb{S}^{d-1}$ .

Recall that  $T_{\mathbf{w}}\mathbb{S}^{d-1}$  is the orthogonal complement  $V_{\mathbf{w}}^{\perp}$  of  $\mathbf{w}$ , and the unit normal vector at  $\mathbf{w}$  is given by  $\mathbf{w}$  itself. Thus  $\overline{\nabla}f(\mathbf{w})$  is the  $V_{\mathbf{w}}^{\perp}$ -component of  $\nabla f(\mathbf{w})$  and  $\overline{\nabla}f(\mathbf{w}) = \nabla f(\mathbf{w}) - (\nabla_{\mathbf{w}}f(\mathbf{w}))\mathbf{w}$ . Similarly  $\overline{\nabla}\mathcal{L} = \nabla\mathcal{L} - (\nabla_{\mathbf{w}}\mathcal{L})\mathbf{w} = \nabla\mathcal{L}$ , where we used the fact that  $\mathcal{L}$  is scaling-invariant and  $\nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w}) = 0$ . Moreover, by (6),  $\nabla\mathcal{L}(\mathbf{w}) \cdot (\overline{\nabla}f(\mathbf{w}) - \nabla f(\mathbf{w})) = (-\nabla_{\mathbf{w}}f(\mathbf{w}))(\nabla\mathcal{L}(\mathbf{w}) \cdot \mathbf{w}) = 0$ . Therefore, (17) and (18) share the same first term.

For the second term, we use the fact that

$$(\overline{\nabla}_{\mathbf{a}_i})^2 - (\nabla_{\mathbf{a}_i})^2 = -\Pi_{\mathbf{w}}(\mathbf{a}_i, \mathbf{a}_i)\nabla_{\mathbf{w}}$$

at  $\mathbf{w} \in \mathbb{S}^{d-1}$ , where  $\Pi_{\mathbf{w}}$  is the second fundamental form of  $\mathbb{S}^{d-1}$  at  $\mathbf{w}$  and  $\nabla_{\mathbf{w}}$  is the directional derivative along the outward normal direction. Since  $\mathbf{w} \cdot \nabla f = \nabla_{\mathbf{w}}f$ , after comparing, we know that there exists a function  $G$  on  $\mathbb{S}^{d-1}$  such that  $\overline{\Psi}f(\mathbf{w}) - \Psi f(\mathbf{w}) = G(\mathbf{w})\nabla_{\mathbf{w}}f$ .

Because trajectories of both (4) and (7) remain in  $\mathbb{S}^{d-1}$ , for  $\mathbf{w} \in \mathbb{S}^{d-1}$  and two smooth functions  $f, g$  that coincide on  $\mathbb{S}^{d-1}$ ,  $\Psi f(\mathbf{w}) = \Psi g(\mathbf{w})$  and  $\overline{\Psi}f(\mathbf{w}) = \overline{\Psi}g(\mathbf{w})$ . Hence  $G(\mathbf{w})\nabla_{\mathbf{w}}(f - g)(\mathbf{w}) = 0$ . Because  $\nabla_{\mathbf{w}}(f - g)(\mathbf{w})$  can be chosen arbitrarily, we must have  $G(\mathbf{w}) = 0$ . Therefore,  $\Psi f = \overline{\Psi}f$  on  $\mathbb{S}^{d-1}$  for all functions  $f$ .

Therefore, restricted to  $\mathbb{S}^{d-1}$ , (4) and (7) have the same generator and are hence equivalent as diffusion processes. The proof is completed.  $\square$

*Proof of Theorem 3.3.* We first prove the theorem for the updating method (9). Under continuous limit, the surrogate for  $\nabla\mathcal{L}_{\mathcal{B}}(\mathbf{w}) = \nabla\mathcal{L}(\mathbf{w}) + (\nabla\mathcal{L}_{\mathcal{B}}(\mathbf{w}) - \nabla\mathcal{L}(\mathbf{w}))$  is  $\nabla\mathcal{L}(\mathbf{w})dt + d\xi_t$  where  $\xi_t$  is a Gaussian diffusion of covariance  $\Sigma(\mathbf{w})$ , such a diffusion can be taken to be  $\xi_t = \int \Sigma(\mathbf{w})^{\frac{1}{2}} dB_t^d$ .

Using the Taylor expansion

$$\begin{aligned} \exp_{\mathbf{w}}(\mathbf{v}) &= \mathbf{w} + \mathbf{w}(\cos|\mathbf{v}| - 1) + \mathbf{v} \frac{\sin|\mathbf{v}|}{|\mathbf{v}|} \\ &= \mathbf{w}(1 - \frac{1}{2}|\mathbf{v}|^2 + \frac{1}{4!}|\mathbf{v}|^4 - \dots) \\ &\quad + \mathbf{v}(1 - \frac{1}{3!}|\mathbf{v}|^2 + \dots), \end{aligned}$$

the surrogate dynamics for (9) is

$$\begin{aligned} d\overline{\mathbf{W}}_t &= \overline{\mathbf{W}}_t \left( -\frac{1}{2} \left| \zeta(\nabla\mathcal{L}(\overline{\mathbf{W}}_t)dt + \Sigma(\overline{\mathbf{W}}_t)^{\frac{1}{2}}dB_t^d \right|^2 \right. \\ &\quad \left. + \frac{1}{4!} \left| \zeta(\nabla\mathcal{L}(\overline{\mathbf{W}}_t)dt + \Sigma(\overline{\mathbf{W}}_t)^{\frac{1}{2}}dB_t^d \right|^4 - \dots \right) \\ &\quad - \zeta(\nabla\mathcal{L}(\overline{\mathbf{W}}_t)dt + \Sigma(\overline{\mathbf{W}}_t)^{\frac{1}{2}}dB_t^d) \\ &\quad \left( 1 - \frac{1}{3!} \left| \zeta(\nabla\mathcal{L}(\overline{\mathbf{W}}_t)dt + \Sigma(\overline{\mathbf{W}}_t)^{\frac{1}{2}}dB_t^d \right|^2 + \dots \right). \end{aligned}$$

Itô's calculus gives

$$\begin{aligned} &\left| \zeta(\nabla\mathcal{L}(\overline{\mathbf{W}}_t)dt + \Sigma(\overline{\mathbf{W}}_t)^{\frac{1}{2}}dB_t^d \right|^2 \\ &= \zeta^2 (dB_t^d)^\top (\Sigma(\overline{\mathbf{W}}_t)^{\frac{1}{2}})^\top \Sigma(\overline{\mathbf{W}}_t)^{\frac{1}{2}} dB_t^d \\ &= \zeta^2 \text{Tr}((\Sigma(\overline{\mathbf{W}}_t)^{\frac{1}{2}})^\top \Sigma(\overline{\mathbf{W}}_t)^{\frac{1}{2}}) dt \\ &= \zeta^2 \text{Tr}\Sigma(\overline{\mathbf{W}}_t) dt. \end{aligned}$$

Hence,

$$\begin{aligned} d\overline{\mathbf{W}}_t &= \overline{\mathbf{W}}_t \left( -\frac{1}{2}\zeta^2 \text{Tr}\Sigma(\overline{\mathbf{W}}_t)dt \right) \\ &\quad - \zeta(\nabla\mathcal{L}(\overline{\mathbf{W}}_t)dt + \Sigma(\overline{\mathbf{W}}_t)^{\frac{1}{2}}dB_t^d) \\ &\quad \left( 1 - \frac{1}{6}\zeta^2 \text{Tr}\Sigma(\overline{\mathbf{W}}_t)dt \right) \\ &= -\zeta(\nabla\mathcal{L}(\overline{\mathbf{W}}_t)dt + \Sigma(\overline{\mathbf{W}}_t)^{\frac{1}{2}}dB_t^d) \\ &\quad - \frac{1}{2}\zeta^2 \text{Tr}\Sigma(\overline{\mathbf{W}}_t)\overline{\mathbf{W}}_t dt. \end{aligned}$$

Note that only the terms  $\mathbf{w}(1 - \frac{1}{2}|\mathbf{v}|^2) + \mathbf{v}$  came into the final expression. Since this is exactly the same equation as (4) with  $\gamma_t^{-\frac{1}{2}} = \zeta$ , the statement now follows from Theorem 3.1.

For the updating method (9), it suffices to use, for  $\mathbf{w} \in \mathbb{S}^{d-1}$  and  $\mathbf{v} \perp \mathbf{w}$ , the Taylor expansion

$$\begin{aligned} \frac{\mathbf{w} + \mathbf{v}}{|\mathbf{w} + \mathbf{v}|} &= (\mathbf{w} + \mathbf{v})(1 + |\mathbf{v}|^2)^{-\frac{1}{2}} \\ &= \mathbf{w}(1 - \frac{1}{2}|\mathbf{v}|^2 + \frac{3}{4}|\mathbf{v}|^4 + \dots) \\ &\quad + \mathbf{v}(1 - \frac{1}{2}|\mathbf{v}|^2 + \frac{3}{4}|\mathbf{v}|^4 + \dots) \end{aligned}$$

Since we still have  $\mathbf{w}(1 - \frac{1}{2}|\mathbf{v}|^2) + \mathbf{v}$  as leading terms, the same argument as above applies and concludes the proof.  $\square$

## B. Formulation of the Fokker-Planck equation

In this section, we justify the Fokker-Planck equation (11), which is a standard fact.

*Proof of Proposition 4.3.* As in the proof of Theorem 3.1, write  $\Sigma(\mathbf{w})^{\frac{1}{2}} = (\mathbf{a}_1(\mathbf{w}), \dots, \mathbf{a}_d(\mathbf{w}))$ . Because for  $\mathbf{w} \in \mathbb{S}^{d-1}$

$$\begin{aligned} & \left( \sum_{j=1}^d \mathbf{a}_j(\mathbf{w})^\top \mathbf{a}_j(\mathbf{w}) \right) \Big|_{T_{\mathbf{w}} \mathbb{S}^{d-1}} \\ &= \left( (\Sigma(\mathbf{w})^{\frac{1}{2}})^\top \Sigma(\mathbf{w})^{\frac{1}{2}} \right) \Big|_{T_{\mathbf{w}} \mathbb{S}^{d-1}} \\ &= \Sigma(\mathbf{w}) \Big|_{T_{\mathbf{w}} \mathbb{S}^{d-1}} \\ &= \sigma^2 \text{Id} \Big|_{T_{\mathbf{w}} \mathbb{S}^{d-1}}, \end{aligned}$$

the operator  $\sum_{j=1}^d \bar{\nabla}_{\mathbf{a}_j}^2$  equals  $\sigma^2 \bar{\Delta}$ . (Recall  $\bar{\Delta}$  denotes the Laplacian on  $\mathbb{S}^{d-1}$ .)

As we are working with the effective learning rate  $\gamma_t^{-\frac{1}{2}}$  equal to a constant  $\zeta$ . The random solution  $\bar{\mathbf{W}}_t$  to (8) satisfies (16) with

$$\begin{aligned} \bar{\Psi}f(\mathbf{w}) &= -\zeta \bar{\nabla} \mathcal{L}(\mathbf{w}) \cdot \bar{\nabla} f(\mathbf{w}) + \frac{1}{2} \zeta^2 \sum_{j=1}^d (\bar{\nabla}_{\mathbf{a}_j})^2 f(\mathbf{w}) \\ &= -\zeta \bar{\nabla} \mathcal{L}(\mathbf{w}) \cdot \bar{\nabla} f(\mathbf{w}) + \frac{1}{2} \zeta^2 \sigma^2 \bar{\Delta} f(\mathbf{w}). \end{aligned} \tag{19}$$

Then for all smooth test functions  $f$ , (16) can be reformulated as

$$\begin{aligned} & \partial_t \int_{\mathbb{S}^{d-1}} u(t, \mathbf{w}) f(\mathbf{w}) \\ &= \int_{\mathbb{S}^{d-1}} u(t, \mathbf{w}) \left( -\zeta \bar{\nabla} \mathcal{L}(\mathbf{w}) \cdot \bar{\nabla} f(\mathbf{w}) + \frac{1}{2} \zeta^2 \sigma^2 \bar{\Delta} f(\mathbf{w}) \right). \end{aligned}$$

After integration by parts, this is equivalent to the desired Fokker-Planck equation (11).  $\square$

## C. Linear time change and notations

We first perform a time change. Let  $\beta$  and  $\mu^{(\beta)}$  be as in (12), (13). Write  $T = \zeta t$  and accordingly

$$\widetilde{\mathbf{W}}_T^{(\beta)} := \bar{\mathbf{W}}_{\frac{T}{\zeta}}; \quad \tilde{u}(T, \mathbf{w}) = u\left(\frac{T}{\zeta}, \mathbf{w}\right). \tag{20}$$

After this time change, (8) and (11) respectively become

$$d\widetilde{\mathbf{W}}_T^{(\beta)} = -\bar{\nabla} \mathcal{L}(\widetilde{\mathbf{W}}_T^{(\beta)}) dT - \beta^{\frac{1}{2}} \bar{\Sigma}(\widetilde{\mathbf{W}}_T^{(\beta)})^{\frac{1}{2}} d\mathbf{B}_T^d, \tag{21}$$

and

$$\partial_T \tilde{u} = \mathcal{D}^{(\beta)} \tilde{u}, \tag{22}$$

where

$$\mathcal{D}^{(\beta)} \tilde{u} := \bar{\nabla} \cdot (\tilde{u} \bar{\nabla} \mathcal{L}) + \beta \bar{\Delta} \tilde{u}. \tag{23}$$

Like in the notations for  $\bar{\mathbf{W}}_t$ , We will denote by  $\mathcal{P}_{\widetilde{\mathbf{W}}_0^{(\beta)} = \mathbf{w}}(\widetilde{\mathbf{W}}_T^{(\beta)})$  and  $\mathcal{P}_{\widetilde{\mathbf{W}}_0^{(\beta)} \sim \nu}(\widetilde{\mathbf{W}}_T^{(\beta)})$  the probability distribution of the random solution  $\widetilde{\mathbf{W}}_t$  to (21), respectively under initial conditions  $\widetilde{\mathbf{W}}_0^{(\beta)} = \mathbf{w}$  and  $\widetilde{\mathbf{W}}_0^{(\beta)} \sim \nu$  where  $\nu$  is a measure (which we allow to be a non-probability). In addition we define an operator  $\tilde{\mathcal{F}}_T^{(\beta)}$  between measures by:

$$\tilde{\mathcal{F}}_T^{(\beta)} \nu := \mathcal{P}_{\widetilde{\mathbf{W}}_0^{(\beta)} \sim \nu}(\widetilde{\mathbf{W}}_T^{(\beta)}), \tag{24}$$

in other words  $\tilde{\mathcal{F}}_T^{(\beta)}$  is the pushforward by the SDE (21) for time  $T$ .

Remark that  $\tilde{\mathcal{F}}_T^{(\beta)}$  is a linear operator between positive measures and preserves total mass. Moreover,  $\tilde{\mathcal{F}}_T^{(\beta)}$  forms a semigroup parametered by  $T$ :  $\tilde{\mathcal{F}}_T^{(\beta)} \circ \tilde{\mathcal{F}}_{T'}^{(\beta)} = \tilde{\mathcal{F}}_{T+T'}^{(\beta)}$ .

For convenience, we will occasionally denote by  $\widetilde{\mathbf{W}}_T^{(\beta)}(\mathbf{w})$  fo the value of  $\widetilde{\mathbf{W}}_T^{(\beta)}$  subject to initial condition  $\widetilde{\mathbf{W}}_0^{(\beta)} = \mathbf{w}$ . Note that this value is random.

## D. Proof of the descent stage

This section contains the proof of Theorem 4.5. After the time change, Theorem 4.5 is equivalent to:

**Theorem D.1.** *Under Assumptions 4.1 and 4.2, for all  $\epsilon > 0$  and  $Q_1 > 0$ , there exists  $C_{\text{des}}, \beta_{\text{des}} > 0$  and a set  $\Lambda_\epsilon$  of volume  $\mathbf{m}(\Lambda_\epsilon) > 1 - \epsilon$ , such that for all  $\beta < \beta_{\text{des}}$  and all  $\mathbf{w}_0 \in \Lambda_\epsilon$ , the solutions to (21) starting at  $\mathbf{w}_0$  satisfy  $\mathbb{P}_{\widetilde{\mathbf{W}}_0^{(\beta)} = \mathbf{w}_0}(\widetilde{\mathbf{W}}_{T_{\text{des}}}^{(\beta)} \in U_{k, Q_1}) > 1 - \epsilon$ , where  $U_k$  is the unique attracting basin in  $\mathbb{S}^{d-1}$  that contains  $\mathbf{w}_0$ .*

We will first focus on the gradient flow without an diffusion term and then prove the diffusion component is a small perturbation in the beginning of the dynamics.

### D.1. Study of the gradient flow

Note that by allowing to endow  $\beta$  with value 0, again via the time change  $T = \zeta t$ , the gradient flow (14) is equivalent to the flow  $\widetilde{\mathbf{W}}_T^{(0)}$  defined by (21) with the same initial condition. Hence we will think of the attracting basins  $U_z$  as attracting basins for the flow  $\widetilde{\mathbf{W}}_T^{(0)}$ . We emphasize that  $\widetilde{\mathbf{W}}_T^{(0)}$  is deterministic and has only one possible trajectory starting from a given initial position.

**Lemma D.2.** *Under Assumption 4.2,*

1.  $\mathbb{S}^{d-1}$  is the disjoint union of  $\{U_z\}_{z \in Z}$ ;
2.  $U_i$  is open for each  $i = 1, \dots, m$ , and the union  $\bigcup_{i=1}^m U_i$  is dense in  $\mathbb{S}^{d-1}$ ;
3. If  $z \in Z \setminus \{z_1, \dots, z_m\}$ , then  $U_z$  is a submanifold whose dimension is strictly less than  $d - 1$ .

This is an elementary fact and the proof is omitted.

**Lemma D.3.** *For all  $Q' > 0$  and  $\epsilon > 0$ , there exists  $\theta = \theta(Q', \epsilon) > 0$ , such that the volume of the set*

$$\Lambda_\theta^1 := \bigcup_{i=1}^m \left\{ \mathbf{w}_0 \in U_i : \forall T \geq 0, \right. \\ \left. \begin{aligned} &\text{either } \mathcal{L}(\widetilde{\mathbf{W}}_T^{(0)}(\mathbf{w}_0)) - \mathcal{L}(z_i) < Q' \\ &\text{or } |\nabla \mathcal{L}(\widetilde{\mathbf{W}}_T^{(0)}(\mathbf{w}_0))|^2 > \theta. \end{aligned} \right\}$$

satisfies  $\mathbf{m}(\Lambda_\theta^1) > 1 - \epsilon$ .

*Proof.* Since the set  $\Lambda_\theta^1$  is decreasing in  $\theta$ , it suffices to show that  $\mathbf{m}(\bigcup_{\theta > 0} \Lambda_\theta^1) = 1$ . We claim that  $\bigcup_{\theta > 0} \Lambda_\theta^1 \supseteq \bigcup_{i=1}^m U_i$ . The lemma would then follow from Lemma D.2.

Suppose  $\widetilde{\mathbf{W}}_0^{(0)} = \mathbf{w}_0 \in U_i$  and  $T \geq 0$ , then  $\widetilde{\mathbf{W}}_T^{(0)} \in U_i$  as well. It suffices to show there exists  $\theta$  that depends on  $\mathbf{w}_0$  but not on  $T$ , such that if  $\mathcal{L}(\widetilde{\mathbf{W}}_T^{(0)}) - \mathcal{L}(z_i) \geq Q'$  then  $|\nabla \mathcal{L}(\widetilde{\mathbf{W}}_T^{(0)})|^2 > \theta$ .

Indeed, as  $\lim_{T \rightarrow \infty} (\mathcal{L}(\widetilde{\mathbf{W}}_T^{(0)}) - \mathcal{L}(z_i)) = 0$ ,  $\mathcal{L}(\widetilde{\mathbf{W}}_T^{(0)}) - \mathcal{L}(z_i) \geq Q'$  only happens on a fixed interval  $[0, T_0]$ .

Moreover,  $\nabla \mathcal{L}(\widetilde{\mathbf{W}}_T^{(0)}) > 0$  for all  $T \geq 0$ , because otherwise  $\widetilde{\mathbf{W}}_T^{(0)} \in Z$  and hence  $\mathbf{w} = \widetilde{\mathbf{W}}_T^{(0)}$  is a critical point, which must be  $z_i$  since  $\mathbf{w} \in U_i$ . In this case  $\mathcal{L}(\widetilde{\mathbf{W}}_T^{(0)}) - \mathcal{L}(z_i) \geq Q'$  does not hold.

Therefore, as  $\nabla \mathcal{L}$  is continuous on the compact interval  $[0, T_0]$ , it suffices to choose  $\theta = \sup_{T \in [0, T_0]} |\nabla \mathcal{L}(\widetilde{\mathbf{W}}_T^{(0)})|^2 > 0$ .  $\square$

**Corollary D.4.** *Let  $\theta = \theta(Q', \epsilon)$  be as in the lemma above. If  $\mathbf{w}_0 \in \Lambda_\theta^1$  then  $\mathbf{w}_0 \in U_i$  for some  $1 \leq i \leq m$ . Moreover, if  $\widetilde{\mathbf{W}}_0^{(0)} = \mathbf{w}_0$ , then for  $T = \frac{\max \mathcal{L} - \min \mathcal{L}}{\theta}$ ,  $\widetilde{\mathbf{W}}_T^{(0)} \in U_i$  and  $\mathcal{L}(\widetilde{\mathbf{W}}_T^{(0)}) - \mathcal{L}(z_i) < Q'$ .*

*Proof.* By definition  $\mathbf{w} \in U_i$  for some  $i$ . Assume that  $\mathcal{L}(\widetilde{\mathbf{W}}_T^{(0)}) - \mathcal{L}(z_i) \geq Q'$ , then  $\mathcal{L}(\widetilde{\mathbf{W}}_s^{(0)}) - \mathcal{L}(z_i) \geq Q$  for all  $s \in [0, T]$  as the gradient flow decreases  $\mathcal{L}$ . Thus  $|\nabla \mathcal{L}(\widetilde{\mathbf{W}}_s^{(0)})|^2 \geq \theta$  and

$$\begin{aligned} &\mathcal{L}(\mathbf{w}_0) - \mathcal{L}(\widetilde{\mathbf{W}}_T^{(0)}) \\ &= - \int_0^T \frac{d}{ds} \mathcal{L}(\widetilde{\mathbf{W}}_s^{(0)}) \\ &= - \int_0^T \nabla \mathcal{L}(\widetilde{\mathbf{W}}_s^{(0)})^\top \cdot (-\nabla \mathcal{L}(\widetilde{\mathbf{W}}_s^{(0)})) ds \\ &= \int_0^T |\nabla \mathcal{L}(g_s \mathbf{w})|^2 ds \geq \theta T \geq \max \mathcal{L} - \min \mathcal{L}, \end{aligned}$$

which cannot be true as  $\mathcal{L}(\widetilde{\mathbf{W}}_T^{(0)}) > \mathcal{L}(z_i) \geq \min \mathcal{L}$ . This completes the proof.  $\square$

### D.2. Perturbative estimate

As the following lemma shows, the SDE (21) is a perturbation of the gradient flow on short time scales, during which the diffusion effect is weak and dominated by the speed to the gradient flow.

**Lemma D.5.** *Assume  $\widetilde{\mathbf{W}}_T^{(\beta)}$  and  $\widetilde{\mathbf{W}}_T^{(0)}$  start from the same initial position  $\mathbf{w}_0$  at  $T = 0$ , then*

$$\mathbb{E} \left( \text{dist}(\widetilde{\mathbf{W}}_T^{(\beta)}(\mathbf{w}_0), \widetilde{\mathbf{W}}_T^{(0)}(\mathbf{w}_0)) \right)^2 = O_T(\beta)$$

for a fixed  $T > 0$  and sufficiently small  $\beta$ .

*Proof.* By (Freidlin & Wentzell, 2012, p32), there exists a Lipschitz type constant  $C = C(\mathcal{L})$ , such that  $\mathbb{E} \left( \text{dist}(\widetilde{\mathbf{W}}_T^{(\beta)}(\mathbf{w}_0), \widetilde{\mathbf{W}}_T^{(0)}(\mathbf{w}_0)) \right)^2 = O \left( \beta C^2 e^{2CT} \int_0^T e^{(2C + \beta C^2)s} ds \right)$ . For small  $\beta$ , the right hand side is  $O_T(\beta)$ .  $\square$

**Corollary D.6.** *In the setting as above, for a fixed  $T$ , there exists a subset  $\Lambda_{T, \epsilon}^2$  with  $\mathbf{m}(\Lambda_{T, \epsilon}^2) > 1 - \epsilon$ , such that for sufficiently small  $\beta$  and all  $\mathbf{w}_0 \in \Lambda_{T, \epsilon}^2$ , the solutions to (21) starting at  $\mathbf{w}$  satisfy*

$$\begin{aligned} &\mathbb{P} \left( \widetilde{\mathbf{W}}_T^{(\beta)}(\mathbf{w}_0) \text{ and } \mathbf{w} \text{ belong to the same } U_i, \right. \\ &\quad \left. \text{and } \text{dist}(\widetilde{\mathbf{W}}_T^{(\beta)}(\mathbf{w}_0), \widetilde{\mathbf{W}}_T^{(0)}(\mathbf{w}_0)) < O_T(\beta^{\frac{1}{2}}) \right) \\ &> 1 - \epsilon. \end{aligned}$$

*Proof.* Lemma D.5 implies, by Chebyshev inequality, that  $\mathbb{P} \left( \text{dist}(\widetilde{\mathbf{W}}_T^{(\beta)}(\mathbf{w}_0), \widetilde{\mathbf{W}}_T^{(0)}(\mathbf{w}_0)) < O_T(\beta^{\frac{1}{2}}) \right) > 1 - \epsilon$  for all initial positions  $\mathbf{w}$ . Remember that  $\widetilde{\mathbf{W}}_T^{(0)}(\mathbf{w}_0)$  and  $\mathbf{w}_0$  belong to the same attracting basin. For  $\widetilde{\mathbf{W}}_T^{(\beta)}$  and  $\mathbf{w}$  to be

in the same  $U_i$ , we can take

$$\tilde{\Lambda}_{T,\beta}^2 := \bigcup_{i=1}^m \{w_0 \in U_i : B_{O_T(\beta^{\frac{1}{2}})}(\tilde{W}_T^{(0)}(w_0)) \subseteq U_i\}.$$

By Lemma D.2, if  $w_0 \notin \tilde{\Lambda}_{T,\beta}^2$ , then the neighborhood  $B_{O_T(\beta^{\frac{1}{2}})}(\tilde{W}_T^{(0)}(w_0))$  meets the boundary of  $U_i$ , which is contained in  $\bigcup_{z \in Z \setminus \{z_1, \dots, z_i\}} U_z$ , a union of finitely many proper submanifolds. In other words,  $\tilde{W}_T^{(0)}(w_0)$  is in the  $O_T(\beta^{\frac{1}{2}})$ -neighborhood  $B_{O_T(\beta^{\frac{1}{2}})}\left(\bigcup_{z \in Z \setminus \{z_1, \dots, z_i\}} U_z\right)$  of this union. It follows that  $(\tilde{\Lambda}_{T,\beta}^2)^c \subseteq \tilde{W}_{-T}^{(0)}\left(B_{O_T(\beta^{\frac{1}{2}})}\left(\bigcup_{z \in Z \setminus \{z_1, \dots, z_i\}} U_z\right)\right)$ , where  $\tilde{W}_{-T}^{(0)}$  is the time-reversed gradient flow. Because

$$\begin{aligned} & \bigcap_{\beta > 0} \tilde{W}_{-T}^{(0)}\left(B_{O_T(\beta^{\frac{1}{2}})}\left(\bigcup_{z \in Z \setminus \{z_1, \dots, z_i\}} U_z\right)\right) \\ &= \tilde{W}_{-T}^{(0)}\left(\bigcup_{z \in Z \setminus \{z_1, \dots, z_i\}} U_z\right) \\ &= \bigcup_{z \in Z \setminus \{z_1, \dots, z_i\}} U_z, \end{aligned}$$

we conclude  $\lim_{\beta \rightarrow 0} \mathbf{m}(\tilde{\Lambda}_{T,\beta}^2) = 1$ . Since the sets  $\tilde{\Lambda}_{T,\beta}^2$  are decreasing in  $\beta$ , one may fix a sufficiently small  $\beta_0 = \beta_0(T, \epsilon)$  and set  $\Lambda_{T,\epsilon}^2 = \tilde{\Lambda}_{T,\beta_0}^2$ , such that  $\mathbf{m}(\Lambda_\epsilon^2) > 1 - \epsilon$ .  $\square$

### D.3. Proof of Theorem 4.5

*Proof of Theorem 4.5.* It suffices to prove Theorem D.1.

Take  $T_{\text{des}} = \frac{\max \mathcal{L} - \min \mathcal{L}}{\theta(\frac{Q_1}{2}, \epsilon)}$  and  $\Lambda_\epsilon = \Lambda_\epsilon^1 \cap \Lambda_{T_{\text{des}}, \epsilon}^2$ . It follows from Corollaries D.4 and D.6 that

$$\begin{aligned} & \mathbb{P}_{\tilde{W}_0^{(\beta)}=w_0} \left( \tilde{W}_{T_{\text{des}}}^{(\beta)} \in U_k, \text{ and} \right. \\ & \quad \left. \mathcal{L}(\tilde{W}_{T_{\text{des}}}^{(\beta)}) - \mathcal{L}(z_k) < \frac{Q_1}{2} + \max |\bar{\nabla} \mathcal{L}| \cdot O_\epsilon(\beta^{\frac{1}{2}}) \right) \\ & > 1 - \epsilon. \end{aligned}$$

For sufficiently small  $\beta$ ,  $\frac{Q_1}{2} + \max |\bar{\nabla} \mathcal{L}| \cdot O_\epsilon(\beta^{\frac{1}{2}}) < Q_1$ . Moreover,  $\mathbf{m}(\Lambda_\epsilon) > 1 - 2\epsilon$ . To deduce the proposition, it suffices to rewrite  $2\epsilon$  as  $\epsilon$ .  $\square$

## E. Proof of the diffusion stage

This section contains the proof of Theorem 4.6. Using the time change  $T = \zeta t$ , Theorem 4.6 follows from:

**Theorem E.1.** *Under Assumptions 4.1 and 4.2, for all  $\epsilon > 0$ , there exist constants  $R_{\text{dif}}, r_{\text{dif}}, \lambda_{\text{dif}} > 0$ , and a set  $\Lambda_\epsilon$  of volume  $\mathbf{m}(\Lambda_\epsilon) > 1 - \epsilon$ , such that:*

With  $\beta = \sqrt{\frac{\lambda_\epsilon}{2}} \sigma$ , for all  $\lambda_\epsilon < \lambda_{\text{dif}}$ , the following is true for all  $T \in [\frac{R_{\text{dif}}}{\beta}, e^{\frac{r_{\text{dif}}}{\beta}}]$ :

(i) For all initial positions  $w_0 \in \Lambda_\epsilon$ ,  $\text{dist}_{\text{TV}}\left(\mathcal{P}_{\tilde{W}_0^{(\beta)}=w_0}(\tilde{W}_T^{(\beta)}), \frac{\mu^{(\beta)}|_{U_k} \mathbf{d}\mathbf{m}}{\int_{U_k} \mu^{(\beta)} \mathbf{d}\mathbf{m}}\right) \leq \epsilon$ , where  $U_k$  is the unique attracting basin of the gradient flow of  $\mathcal{L}$  that contains  $w_0$ .

(ii) More generally, for all initial probability distribution  $\nu_0$ ,

$$\text{dist}_{\text{TV}}\left(\mathcal{P}_{\tilde{W}_0^{(\beta)} \sim \nu_0}(\tilde{W}_T^{(\beta)}), \sum_{i=1}^m \nu_0(U_i) \frac{\mu^{(\beta)}|_{U_i} \mathbf{d}\mathbf{m}}{\int_{U_i} \mu^{(\beta)} \mathbf{d}\mathbf{m}}\right) \leq \epsilon + \nu_0(\Lambda_\epsilon^c).$$

### E.1. Relevant Hilbert spaces and operators

Consider the adjoint operator  $(\mathcal{D}^{(\beta)})^* = -\bar{\nabla} \mathcal{L} \cdot \bar{\nabla} + \beta \bar{\Delta}$ . It is known that  $(\mathcal{D}^{(\beta)})^*$  is self-adjoint on the Hilbert space  $L^2(\mu^{(\beta)}) := L^2(\mathbb{S}^{d-1}, \mu^{(\beta)} \mathbf{d}\mathbf{m})$  (Kolokoltsov, 2000, §8.5). More precisely, for smooth functions  $f, g$ ,

$$\int f((\mathcal{D}^{(\beta)})^* g) \mu^{(\beta)} \mathbf{d}\mathbf{m} = \int ((\mathcal{D}^{(\beta)})^* f) g \mu^{(\beta)} \mathbf{d}\mathbf{m}. \quad (25)$$

This is equivalent to

$$\int \mathcal{D}^{(\beta)}(f \mu^{(\beta)}) g \mathbf{d}\mathbf{m} = \int f \mathcal{D}^{(\beta)}(g \mu^{(\beta)}) \mathbf{d}\mathbf{m},$$

or

$$\begin{aligned} & \int \mathcal{D}^{(\beta)}(f \mu^{(\beta)}) g \mu^{(\beta)} \cdot \frac{1}{\mu^{(\beta)}} \mathbf{d}\mathbf{m} \\ &= \int f \mu^{(\beta)} \mathcal{D}^{(\beta)}(g \mu^{(\beta)}) \cdot \frac{1}{\mu^{(\beta)}} \mathbf{d}\mathbf{m}. \end{aligned}$$

This shows  $\mathcal{D}^{(\beta)}$  is self-adjoint for the Hilbert space  $L^2(\frac{1}{\mu^{(\beta)}}) := L^2(\mathbb{S}^{d-1}, \frac{1}{\mu^{(\beta)}} \mathbf{d}\mathbf{m})$ . We shall also write  $L^2(1) := L^2(\mathbb{S}^{d-1}, \mathbf{d}\mathbf{m})$  for the unweighted  $L^2$  space.

The equality (25) can also be written as  $\int ((\mathcal{D}^{(\beta)})^* f) g \mu^{(\beta)} \mathbf{d}\mathbf{m} = \int \mathcal{D}^{(\beta)}(f \mu^{(\beta)}) g \mathbf{d}\mathbf{m}$ , which implies

$$(\mathcal{D}^{(\beta)})^* f = \frac{1}{\mu^{(\beta)}} \mathcal{D}^{(\beta)}(f \mu^{(\beta)}). \quad (26)$$

On the other hand, one can directly check the following facts:

**Lemma E.2.** *Let  $f_T(w) = \mu^{(\beta)}(w)^{-\frac{1}{2}} \tilde{u}_T(w)$ . If  $u$  is a solution to (22) then  $f$  satisfies  $\partial_T f = \mathcal{D}^{\beta, \#} f$  where*

$$\mathcal{D}^{\beta, \#} f := \beta \bar{\Delta} f - \left( \frac{|\bar{\nabla} \mathcal{L}|^2}{4\beta} - \frac{\bar{\Delta} \mathcal{L}}{2} \right) f. \quad (27)$$

We also note

$$\mathcal{D}^{(\beta)} \tilde{u} = \mathcal{D}^{(\beta)}((\mu^{(\beta)})^{\frac{1}{2}} f) = (\mu^{(\beta)})^{\frac{1}{2}} \mathcal{D}^{\beta, \#} f. \quad (28)$$

The equalities (25) and (28) guarantee the following commutative diagram of operators:

$$\begin{array}{ccccc}
 L^2(\mu^{(\beta)}) & \xrightarrow{f \rightarrow (\mu^{(\beta)})^{\frac{1}{2}} f} & L^2(1) & \xrightarrow{f \rightarrow (\mu^{(\beta)})^{\frac{1}{2}} f} & L^2(\frac{1}{\mu^{(\beta)}}) \\
 \downarrow (\mathcal{D}^{(\beta)})^* & & \downarrow \mathcal{D}^{\beta, \#} & & \downarrow \mathcal{D}^{(\beta)} \\
 L^2(\mu^{(\beta)}) & \xrightarrow{f \rightarrow (\mu^{(\beta)})^{\frac{1}{2}} f} & L^2(1) & \xrightarrow{f \rightarrow (\mu^{(\beta)})^{\frac{1}{2}} f} & L^2(\frac{1}{\mu^{(\beta)}})
 \end{array} \tag{29}$$

In this diagram, every horizontal arrow is an isometry and every vertical arrow is a self-adjoint operator. In particular,  $f \rightarrow f(\mu^{(\beta)})^{\frac{1}{2}}$  is a bijection, from the eigenfunctions of  $(\mathcal{D}^{(\beta)})^*$  to those of  $\mathcal{D}^{\beta, \#}$  with the same eigenvalues, then again a bijection from the later ones to the eigenfunctions of  $\mathcal{D}^{(\beta)}$  with the same eigenvalues. In other words, the spectra of  $(\mathcal{D}^{(\beta)})^*$ ,  $\mathcal{D}^{\beta, \#}$  and  $\mathcal{D}^{(\beta)}$ , as self-adjoint operators in their corresponding spaces, are the same. By self-adjointness, this spectrum is actually contained in  $\mathbb{R}$ .

Recall that  $U_j$  is the attracting basin containing  $z_j$ . From now on, we will denote the indicator function of  $U_j$  by

$$\chi_j := \chi_{U_j}.$$

The low lying eigenvalues of  $-\mathcal{D}^{(\beta)}$  correspond to the local minima. This was first proved in (Simon, 1983). The precise version that we need can be found in (Kolokoltsov, 2000).

**Proposition E.3.** (Kolokoltsov, 2000, p248) *Given  $\mathcal{L}$ , there exist constants  $Q_0, \rho > 0$ , determined by  $\mathcal{L}$ , such that:*

1. *The first  $m$  eigenvalues (counted with multiplicity) of  $-\mathcal{D}^{(\beta)}$  in  $L^2(\frac{1}{\mu^{(\beta)}})$  are  $\leq O(e^{-\frac{Q_0}{\beta}})$ ;*
2. *With  $\Pi^{(\beta)}$  and  $(\Pi^{(\beta)})^\perp = \text{Id} - \Pi^{(\beta)}$  respectively denoting the orthogonal projections in  $L^2(\frac{1}{\mu^{(\beta)}})$  to the span of the first  $m$  eigenfunctions of  $-\mathcal{D}^{(\beta)}$  and to its orthogonal complement,  $\frac{\|(\Pi^{(\beta)})^\perp(\chi_j \mu^{(\beta)})\|_{L^2(\frac{1}{\mu^{(\beta)}})}}{\|\chi_j \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} = O(e^{-\frac{Q_0}{\beta}})$  holds for  $1 \leq j \leq m$ ;*
3. *All other eigenvalues of  $-\mathcal{D}^{(\beta)}$  are greater than or equal to  $\rho$ .*

More precisely,  $Q_0$  can be any positive value such that  $Q_0 < \min_i \sup_{\mathbf{w} \in U_i} ((\mathcal{L}(\mathbf{w}) - \mathcal{L}(z_i)))$ .

The original formulation in (Kolokoltsov, 2000) was for the spectral decomposition of  $-(\mathcal{D}^{(\beta)})^*$ . But in light of the correspondence  $\mathcal{D}^{(\beta)}(\mu^{(\beta)} f) = \mu^{(\beta)}(\mathcal{D}^{(\beta)})^* f$  (see §E.1), the translation to the  $\mathcal{D}^{(\beta)}$  setting is straightforward.

## E.2. Approximate spectral decomposition

For brevity, let  $\hat{\chi}_j^{(\beta)} = \frac{\chi_j \mu^{(\beta)}}{\|\chi_j \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}}$  and  $\psi_j^{(\beta)} =$

$\Pi^{(\beta)} \hat{\chi}_j^{(\beta)}$ . Then  $\|\hat{\chi}_j^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})} = 1$ . By Proposition E.3,

$$\|\psi_j^{(\beta)} - \hat{\chi}_j^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \leq O(e^{-\frac{Q_0}{\beta}}), \tag{30}$$

and the matrix

$$K = (K_{ij})_{i,j=1}^m := (\langle \psi_i^{(\beta)}, \psi_j^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})})_{i,j=1}^m$$

is within distance  $O(e^{-\frac{Q_0}{\beta}})$  from Id. In consequence,

$$\|K^{-1} - \text{Id}\| \leq O(e^{-\frac{Q_0}{\beta}}). \tag{31}$$

The entries of  $K^{-1}$  will be denoted by  $K^{-1} = (K^{ij})$ .

Suppose  $\tilde{u}_T(\cdot) = \tilde{u}(T, \cdot)$  is the density function of the distribution of trajectories  $\tilde{\mathbf{W}}_T$  that start from an initial point  $\mathbf{w}_0$ . We are interested in the spectral decomposition of  $\tilde{u}_T$ , or more precisely the projection

$$\Pi^{(\beta)} f = \sum_{i=1}^m \left( \sum_{j=1}^m K^{ij} \langle \tilde{u}_T, \psi_j^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} \right) \psi_j^{(\beta)}. \tag{32}$$

The strategy is to show that the coefficient  $\sum_{j=1}^m K^{ij} \langle \tilde{u}_T, \psi_j^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})}$  in the decomposition above has little dependence on the choice of  $\mathbf{w}_0$  as long as it is supported in  $U_{k, Q_1}$ .

## E.3. Regularity bound

From now on, suppose  $U_k$  is the attracting basin containing  $\mathbf{w}_0$ . By Theorem D.1, after running the SDE (21) for  $T_1 = O_{Q_1, \epsilon}(1)$  time starting from a randomly sampled point  $\mathbf{w}_0$ , the resulting point  $\tilde{\mathbf{W}}_{T_1}$  is with probability  $1 - \epsilon$  in the ‘‘well’’  $U_{k, Q_1}$  near the local minimum  $z_k$  at the bottom of  $U_k$ .

We will start the SDE again, on a longer time scale, from the restricted probability distribution  $(\tilde{\mathcal{F}}_{T_1}^{(\beta)} \delta_{\mathbf{w}_0})|_{U_{k, Q_1}}$ . In the next stage, gradient descent slows down because  $\nabla \mathcal{L}$  is small in  $U_{k, Q_1}$ , and diffusion behavior plays a more important role than in the previous stage.

It is noteworthy that  $(\tilde{\mathcal{F}}_{T_1}^{(\beta)} \delta_{\mathbf{w}_0})^{U_{k, Q_1}}$  is absolutely continuous because  $T_1 > 0$  and the parabolic PDE (22) is non-degenerate.

Instead of bounding  $\tilde{u}(T, \cdot)$ , we allow the stochastic process (21) to run for a duration of  $\frac{R}{\beta}$  before estimating the regularity of the new marginal distribution  $u_{\frac{R}{\beta}} \mathbf{d}\mathbf{m}$ . Let  $p : (0, \infty) \times \mathbb{S}^{d-1} \times \mathbb{S}^{d-1}$  be the heat kernel of the operator  $\mathcal{D}^{(\beta)}$ , i.e.  $p(T, \mathbf{x}, \mathbf{y}) \mathbf{d}\mathbf{m}(\mathbf{y}) = \mathbf{d}(\tilde{\mathcal{F}}_T^{(\beta)} \delta_{\mathbf{x}})(\mathbf{y})$ . In other words, for each fixed  $\mathbf{x}$ ,  $p(T, \mathbf{x}, \cdot)$  solves (22) and  $p(T, \mathbf{x}, \mathbf{y}) \mathbf{d}\mathbf{m}(\mathbf{y})$  converges weakly to  $\delta_{\mathbf{x}}$  as  $T \rightarrow 0$ .

**Lemma E.4.** *There exist constants  $C_1, C_2 > 0$  determined by  $\mathcal{L}$ , such that for every  $R \geq 1$  and sufficiently small  $\beta$ ,  $|p(\frac{R}{\beta}, \cdot, \cdot)|$  is uniformly bounded by  $O(e^{\frac{C_1 R + C_2}{\beta}})$ .*

*Proof.* First remark that, by (29),  $e^{-\frac{\mathcal{L}(x) - \mathcal{L}(y)}{2\beta}} p(T, \mathbf{x}, \mathbf{y}) = (\mu^{(\beta)})^{\frac{1}{2}}(\mathbf{x})(\mu^{(\beta)})^{-\frac{1}{2}}(\mathbf{y}) p(T, \mathbf{x}, \mathbf{y})$  is the heat kernel of the operator  $\mathcal{D}^{\beta, \#}$ .

We then renormalize  $\mathcal{D}^{\beta, \#}$  to  $\frac{\mathcal{D}^{\beta, \#}}{\beta} = \bar{\Delta} - V^{(\beta)}$  where  $V^{(\beta)} := \frac{|\bar{\Delta}\mathcal{L}|^2}{4\beta^2} - \frac{\bar{\Delta}\mathcal{L}}{2\beta}$ . The heat kernel of  $\frac{\mathcal{D}^{\beta, \#}}{\beta}$  at time  $R$  takes the form  $e^{-\frac{\mathcal{L}(x) - \mathcal{L}(y)}{2\beta}} p(\frac{R}{\beta}, \mathbf{x}, \mathbf{y})$ .

Note that  $V^{(\beta)}$  may take negative values, but  $\min V^{(\beta)} \geq -\frac{\|\bar{\Delta}\mathcal{L}\|_{L^\infty}}{2\beta}$ . Define the Schrödinger operator

$$\mathcal{S}^{(\beta)} := \bar{\Delta} - (V^{(\beta)} - \min V^{(\beta)}),$$

whose heat kernel at time  $R$  is  $e^{(\min V^{(\beta)})R} e^{-\frac{\mathcal{L}(x) - \mathcal{L}(y)}{2\beta}} p(\frac{R}{\beta}, \mathbf{x}, \mathbf{y})$ .

Because the new potential function  $V^{(\beta)} - \min V^{(\beta)}$  is non-negative, by a standard maximal principle argument, the heat kernel of  $\mathcal{S}^{(\beta)}$  is bounded by the Gaussian heat kernel, i.e. the one of the Laplacian  $\bar{\Delta}$  on  $\mathbb{S}^{d-1}$ . This in particular shows that

$$\begin{aligned} & e^{(\min V^{(\beta)})R} e^{-\frac{\mathcal{L}(x) - \mathcal{L}(y)}{2\beta}} p(\frac{R}{\beta}, \mathbf{x}, \mathbf{y}) \\ & \leq O\left(R^{-\frac{d-1}{2}} e^{-\frac{c \operatorname{dist}(\mathbf{x}, \mathbf{y})^2}{R}}\right) \end{aligned}$$

for some constant  $c > 0$ . Thus

$$p(\frac{R}{\beta}, \mathbf{x}, \mathbf{y}) \leq O\left(R^{-\frac{d-1}{2}} e^{-(\min V^{(\beta)})R} e^{\frac{\mathcal{L}(x) - \mathcal{L}(y)}{2\beta}}\right).$$

To establish the lemma, it suffices to take  $C_1 = \frac{\|\bar{\Delta}\mathcal{L}\|_{L^\infty}}{2}$  and  $C_2 = \max \mathcal{L} - \min \mathcal{L}$ .  $\square$

**Corollary E.5.** *For  $R \geq 1$ , small  $\beta$  and any initial measure  $\nu_0$  on  $\mathbb{S}^{d-1}$ ,  $\tilde{\mathcal{F}}_{\frac{R}{\beta}}^{(\beta)} \nu_0$  is absolutely continuous. Moreover,  $\tilde{\mathcal{F}}_{\frac{R}{\beta}}^{(\beta)} \nu_0 = q d\mathbf{m}$  for a function  $q$  with  $\|q\|_{L^\infty} \leq O(e^{\frac{C_1 R + C_2}{\beta}}) \nu_0(\mathbb{S}^{d-1})$ .*

The operator  $\tilde{\mathcal{F}}_{\frac{R}{\beta}}^{(\beta)}$  was defined in §C.

*Proof.* This follows from Lemma E.4 and the direct decomposition  $q(\mathbf{y}) = \int p(\frac{R}{\beta}, \mathbf{x}, \mathbf{y}) d\nu_0(\mathbf{x})$ .  $\square$

#### E.4. Non-escaping from well

We will make a choice of the parameter  $Q_1$  in Theorem D.1 that depends only on  $\mathcal{L}$ . The goal is to assert that the

a trajectory of the SDE (21) with starting point  $\mathbf{w} \in U_{k, Q_1}$  will be trapped in  $U_{k, Q}$  for an exponentially long period with high probability. The argument below is essentially due to (Freidlin & Wentzell, 2012, §4.4). However, while their proof works for all  $\mathbf{w} \in U_{k, Q_0}$ , the purpose of reproducing it here is to demonstrate the uniformity of the estimate for all  $\mathbf{w} \in U_{k, Q_1}$ .

For every  $\mathbf{w} \in U_{k, Q_0}$ , its first exit time with respect to  $U_{k, Q_0}$  is denoted by

$$\tau^{(\beta)}(\mathbf{w}) := \inf \left\{ T > 0 : \widetilde{\mathbf{W}}_T^{(\beta)}(\mathbf{w}) \in \partial U_{i, Q_0} \right\}. \quad (33)$$

**Lemma E.6.** *Given  $\mathcal{L}$  and  $Q_0$ , there exists a constant  $Q_1$ , such that the bound  $\mathbb{P}(\tau^{(\beta)}(\mathbf{w}) \leq e^{\frac{Q_0}{3\beta}}) \leq O(e^{-\frac{Q_0}{3\beta}})$  holds uniformly for all  $\mathbf{w} \in U_{k, Q_1}$ .*

*Proof.* Together with another constant  $r$ , also determined by  $\mathcal{L}$ , we can make  $Q_1$  satisfy the following condition:

$$\forall 1 \leq i \leq m, \overline{U_{i, Q_1}} \subset B_r(\mathbf{z}_i) \text{ and } \overline{B_{3r}(\mathbf{z}_i)} \subset U_{i, Q_0} \quad (34)$$

For every  $\mathbf{w} \in U_{i, Q_0} \cup \partial U_{i, Q_0}$ , and let  $\tau_1^{(\beta)}(\mathbf{w})$  be the following stopping time:

$$\begin{aligned} & \tau_1^{(\beta)}(\mathbf{w}) \\ & := \inf \left\{ T > 0 : \widetilde{\mathbf{W}}_T^{(\beta)}(\mathbf{w}) \in \partial B_r(\mathbf{z}_i) \cup \partial U_{i, Q_0} \right. \\ & \quad \left. \text{and } \exists T' \in [0, T) \text{ s.t. } \widetilde{\mathbf{W}}_{T'}^{(\beta)}(\mathbf{w}) \in \partial B_{2r}(\mathbf{z}_i) \right\}. \end{aligned} \quad (35)$$

Fix constants  $Q'_0, Q''_0$  such that  $0 < Q''_0 < Q'_0 < Q_0$ . In addition to (34), by (Freidlin & Wentzell, 2012, Ch. 4, (4.6)), one can choose  $r$  to be sufficiently small, such that for all sufficiently small  $\beta$ ,

$$\mathbb{P}\left(\widetilde{\mathbf{W}}_{\tau_1^{(\beta)}(\mathbf{w})}^{(\beta)}(\mathbf{w}) \in \partial U_{i, Q_0}\right) \leq e^{-\frac{Q'_0}{\beta}} \quad (36)$$

for all  $1 \leq i \leq m$  and  $\mathbf{w} \in \partial B_r(\mathbf{z}_i) \cup \partial U_{i, Q_0}$ . In addition, define for every  $\mathbf{w} \in \overline{U_{i, Q_0}}$  recursively the  $n$ -th stopping times for all  $n \geq 1$ :

$$\tau_n^{(\beta)}(\mathbf{w}) := \tau_{n-1}^{(\beta)}(\mathbf{w}) + \tau_1^{(\beta)}(\widetilde{\mathbf{W}}_{\tau_{n-1}^{(\beta)}(\mathbf{w})}^{(\beta)}). \quad (37)$$

Note that if  $\mathbf{w} \in B_{2r}(\mathbf{z}_i)$ , then  $\tau^{(\beta)}(\mathbf{w}) = \tau_{N^{(\beta)}(\mathbf{w})}^{(\beta)}(\mathbf{w})$  for a random variable  $N = N^{(\beta)}(\mathbf{w}) \in \mathbb{N}$ .

Once  $r$  is fixed, because the stochastic process (21) is a perturbation of the geodesic flow when  $\beta$  is very small and the separation condition (34) holds, there exists a constant  $\theta > 0$  such that  $\mathbb{P}(\tau_1^{(\beta)}(\mathbf{w}) > \theta) > \frac{1}{2}$  for all sufficiently small  $\beta$  and  $\mathbf{w} \in \partial B_{2r}(\mathbf{z}_k)$ . By construction of  $\tau_1^{(\beta)}$ , this inequality is also true for all  $\mathbf{w} \in \partial B_r(\mathbf{z}_k)$

From now on suppose  $\mathbf{w} \in U_{i,Q_1} \subset B_r(\mathbf{z}_i)$ . Set a target iteration number at  $M^{(\beta)} = \lfloor \frac{4}{\theta} e^{\frac{Q'_0}{\beta}} \rfloor$ . Then  $\tau_1^{(\beta)}(\mathbf{w}) = \tau_0^{(\beta)}(\mathbf{w}) + \tau_1^{(\beta)}(\widetilde{\mathbf{W}}_{\tau_0^{(\beta)}(\mathbf{w})}^{(\beta)}(\mathbf{w}))$  where  $\tau_0^{(\beta)}(\mathbf{w}) := \inf\{T > 0 : \widetilde{\mathbf{W}}_T^{(\beta)}(\mathbf{w}) \in \partial B_r(\mathbf{z}_i)\}$ . In particular,

$$N^{(\beta)}(\mathbf{w}) = N^{(\beta)}(\widetilde{\mathbf{W}}_{\tau_0^{(\beta)}(\mathbf{w})}^{(\beta)}(\mathbf{w})). \quad (38)$$

Since  $\widetilde{\mathbf{W}}_{\tau_0^{(\beta)}(\mathbf{w})}^{(\beta)}(\mathbf{w}) \in \partial B_r(\mathbf{z}_i) \cup \partial U_{i,Q_0}$ , by (36),

$$\begin{aligned} \mathbb{P}(N^{(\beta)}(\mathbf{w}) < M^{(\beta)}) &\leq 1 - (1 - e^{-\frac{Q'_0}{\beta}})^{\lfloor \frac{4}{\theta} e^{\frac{Q'_0}{\beta}} \rfloor} \\ &= O(e^{\frac{Q'_0 - Q'_0}{\beta}}), \end{aligned} \quad (39)$$

where the implied constant is uniform for  $\mathbf{w} \in U_{i,Q_1} \subset B_r(\mathbf{z}_i)$ .

On the other hand, since  $\tau_n^{(\beta)}(\mathbf{w}) = \tau_0^{(\beta)}(\mathbf{w}) + \sum_{l=1}^n \tau_l^{(\beta)}(\widetilde{\mathbf{W}}_{\tau_{l-1}^{(\beta)}(\mathbf{w})}^{(\beta)}(\mathbf{w}))$ , and each term in the summation, given all precedent terms, is greater than  $\theta$  with at least  $\frac{1}{2}$  probability, we know by large deviation principle that for an absolute constant  $c > 0$ ,

$$\mathbb{P}(\tau_{M^{(\beta)}}^{(\beta)}(\mathbf{w}) \leq \frac{\theta}{3} M^{(\beta)}) \leq e^{-cM^{(\beta)}} \quad (40)$$

uniformly for all  $\mathbf{w} \in U_{i,Q_1}$ . Combining (39) and (40) yields that for all sufficiently small  $\beta$  (depending only on  $\mathcal{L}$ ), we have uniformly for all  $\mathbf{w} \in U_{i,Q_1}$ ,

$$\begin{aligned} \mathbb{P}(\tau^{(\beta)}(\mathbf{w}) \leq e^{\frac{Q'_0}{\beta}}) &= \mathbb{P}(\tau_{N^{(\beta)}(\mathbf{w})}^{(\beta)}(\mathbf{w}) \leq e^{\frac{Q'_0}{\beta}}) \leq \mathbb{P}(\tau_{N^{(\beta)}(\mathbf{w})}^{(\beta)}(\mathbf{w}) \leq \frac{\theta}{3} M^{(\beta)}) \\ &\leq \mathbb{P}(N^{(\beta)}(\mathbf{w}) < M^{(\beta)}) + \mathbb{P}(\tau_{M^{(\beta)}}^{(\beta)}(\mathbf{w}) \leq \frac{\theta}{3} M^{(\beta)}) \\ &\leq O(e^{\frac{Q'_0 - Q'_0}{\beta}}) + e^{-c \lfloor \frac{4}{\theta} e^{\frac{Q'_0}{\beta}} \rfloor} \\ &\leq O(e^{\frac{Q'_0 - Q'_0}{\beta}}). \end{aligned} \quad (41)$$

The proof is completed by taking  $Q'_0 = \frac{Q_0}{3}$  and  $Q'_0 = \frac{2Q_0}{3}$ .  $\square$

### E.5. Short term convergence towards local equilibria

We are now ready to put the pieces together and understand the evolution during the current stage.

Let  $T_1 = T_{\text{des}}$  as in Theorem D.1. Fix an arbitrary parameter  $R' > 0$ , say  $R' = 1$ . Let

$$C'_2 = C_2 + \frac{1}{2}(\max \mathcal{L} - \min \mathcal{L}) = \frac{3}{2}(\max \mathcal{L} - \min \mathcal{L}),$$

which depends only on  $\mathcal{L}$ . Consider  $R \geq R'$  such that the following assumptions are satisfied:

$$R \geq \frac{\rho + C_1}{\rho} R' + \frac{C'_2 + Q_0}{\rho}, \quad (42)$$

$$\frac{R'}{\beta} \leq \frac{R}{\beta} \leq e^{\frac{Q_0}{3\beta}}. \quad (43)$$

We first write

$$\mathcal{P}_{\overline{\mathbf{W}}_0 = \mathbf{w}_0}(\overline{\mathbf{W}}_{\frac{R}{\beta}}) = \widetilde{\mathcal{F}}_{\frac{R}{\beta}}^{(\beta)} \delta_{\mathbf{w}_0} = \widetilde{\mathcal{F}}_{\frac{R-R'}{\beta}}^{(\beta)} \widetilde{\mathcal{F}}_{\frac{R'}{\beta}}^{(\beta)} \delta_{\mathbf{w}_0}. \quad (44)$$

Then decompose

$$\begin{aligned} \widetilde{\mathcal{F}}_{\frac{R'}{\beta}}^{(\beta)} \delta_{\mathbf{w}_0} &= \widetilde{\mathcal{F}}_{\frac{R'}{\beta} - T_1}^{(\beta)} ((\widetilde{\mathcal{F}}_{T_1}^{(\beta)} \delta_{\mathbf{w}_0})|_{U_{k,Q_1}}) \\ &\quad + \widetilde{\mathcal{F}}_{\frac{R'}{\beta} - T_1}^{(\beta)} ((\widetilde{\mathcal{F}}_{T_1}^{(\beta)} \delta_{\mathbf{w}_0})|_{U_{k,Q_1}^c}). \end{aligned}$$

By Theorem 4.5, the total mass of  $\widetilde{\mathcal{F}}_{\frac{R'}{\beta} - T_1}^{(\beta)} ((\widetilde{\mathcal{F}}_{T_1}^{(\beta)} \delta_{\mathbf{w}_0})|_{U_{k,Q_1}^c})$  is less than  $\epsilon$ .

For simplicity, denote the measure  $\widetilde{\mathcal{F}}_{\frac{R'}{\beta} - T_1}^{(\beta)} ((\widetilde{\mathcal{F}}_{T_1}^{(\beta)} \delta_{\mathbf{w}_0})|_{U_{k,Q_1}})$  by  $\gamma$ . Then its properties can be summarized as follows:

- (i)  $\text{dist}_{\text{TV}}(\widetilde{\mathcal{F}}_{\frac{R'}{\beta}}^{(\beta)} \delta_{\mathbf{w}_0}, \gamma) \leq \epsilon$ ;
- (ii)  $\gamma$  is absolutely continuous. Moreover,  $\gamma = h(\cdot) \mathbf{d}\mathbf{m}$  with a density function  $h$  satisfying  $\|h\| \leq O(e^{\frac{C_1 R' + C_2}{\beta}})$ .

Here the second property follows from Corollary E.5.

In particular, the total mass of  $\gamma$  satisfies  $1 - \epsilon \leq \gamma(\mathbb{S}^{d-1}) \leq 1$ .

We now run the Fokker-Planck equation (22) starting at initial time  $\frac{R'}{\beta}$  and initial data  $h(\cdot)$ , and denote the solution by  $\tilde{h}$ ; that is,  $\tilde{h}_T(\mathbf{w}) := \tilde{h}(T, \mathbf{w})$  is defined on  $T \geq \frac{R'}{\beta}$ ,  $\mathbf{w} \in \mathbb{S}^{d-1}$ , and solves

$$\begin{cases} \partial_T \tilde{h}_T = \mathcal{D}^{(\beta)} \tilde{h}_T \\ \tilde{h}_{\frac{R'}{\beta}}(\cdot) = h(\cdot). \end{cases} \quad (45)$$

In other words,  $\tilde{h}_{T - \frac{R'}{\beta}} = e^{(T - \frac{R'}{\beta}) \mathcal{D}^{(\beta)}} h$ . Note that as the Fokker-Planck equation preserves total mass,

$$\int \tilde{h}_T \mathbf{d}\mathbf{m} = \int h \mathbf{d}\mathbf{m} \in [1 - \epsilon, 1], \forall T \geq \frac{R'}{\beta}. \quad (46)$$

Since  $\tilde{h}$  and  $\tilde{u}$  both satisfy (22), for  $T \geq \frac{R'}{\beta}$ ,  $\tilde{h}_T$  will be viewed as an approximation of  $\tilde{u}_T$  in the sense that

$$\begin{aligned} &\text{dist}_{\text{TV}}(\widetilde{\mathcal{F}}_T^{(\beta)} \delta_{\mathbf{w}_0}, \tilde{h}_T \mathbf{d}\mathbf{m}) \\ &= \text{dist}_{\text{TV}}(\tilde{u}_T \mathbf{d}\mathbf{m}, \tilde{h}_T \mathbf{d}\mathbf{m}) \\ &= \text{dist}_{\text{TV}}(\widetilde{\mathcal{F}}_{T - \frac{R'}{\beta}}^{(\beta)}(\tilde{u}_{\frac{R'}{\beta}} \mathbf{d}\mathbf{m}), \widetilde{\mathcal{F}}_{T - \frac{R'}{\beta}}^{(\beta)}(\tilde{h}_{\frac{R'}{\beta}} \mathbf{d}\mathbf{m})) \\ &= \text{dist}_{\text{TV}}(\widetilde{\mathcal{F}}_{T - \frac{R'}{\beta}}^{(\beta)}(\tilde{u}_{\frac{R'}{\beta}} \mathbf{d}\mathbf{m}), \widetilde{\mathcal{F}}_{T - \frac{R'}{\beta}}^{(\beta)}(\gamma)) \\ &\leq \epsilon. \end{aligned} \quad (47)$$

Recall  $\psi_j^{(\beta)} = \Pi^{(\beta)} \hat{\chi}_j^{(\beta)}$  and  $\hat{\chi}_j^{(\beta)} - \psi_j^{(\beta)} = (\Pi^{(\beta)})^\perp \hat{\chi}_j^{(\beta)}$  are respectively the projections of  $\hat{\chi}_j^{(\beta)}$  to the span of the first  $m$  eigenvalues of  $-\mathcal{D}^{(\beta)}$  and its orthogonal complement in  $L^2(\frac{1}{\mu^{(\beta)}})$ . The conditions (42) and (43) will be assumed on  $R$  and  $R'$  without further notice.

**Lemma E.7.** For sufficiently small  $\beta$ ,

$$\|(\Pi^{(\beta)})^\perp \tilde{h}_{\frac{R}{\beta}}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \leq O(e^{-\frac{Q_0}{\beta}}).$$

The projections  $\Pi^{(\beta)}$  and  $(\Pi^{(\beta)})^\perp$  were defined in Proposition E.3.

*Proof.* The projection  $(\Pi^{(\beta)})^\perp \tilde{h}_{\frac{R}{\beta}}$  can be bounded by

$$\begin{aligned} & \|(\Pi^{(\beta)})^\perp \tilde{h}_{\frac{R}{\beta}}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \\ &= \|(\Pi^{(\beta)})^\perp e^{\frac{R-R'}{\beta} \mathcal{D}^{(\beta)}} \tilde{h}_{\frac{R'}{\beta}}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \\ &= \|e^{\frac{R-R'}{\beta} \mathcal{D}^{(\beta)}} (\Pi^{(\beta)})^\perp \tilde{h}_{\frac{R'}{\beta}}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \\ &\leq e^{-\rho(\frac{R-R'}{\beta})} \|\tilde{h}_{\frac{R'}{\beta}}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \end{aligned} \quad (48)$$

because by Proposition E.3.(3), the spectrum of  $-\mathcal{D}^{(\beta)}$  on the image of  $(\Pi^{(\beta)})^\perp$  is in  $[\rho, \infty)$ .

By the bound on  $\tilde{h}_{\frac{R'}{\beta}} = h$ , we have:

$$\begin{aligned} & \|\tilde{h}_{\frac{R'}{\beta}}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \\ &\leq O(e^{\frac{C_1 R' + C_2}{\beta}}) \|1\|_{L^2(\frac{1}{\mu^{(\beta)}})} \\ &= O(e^{\frac{C_1 R' + C_2}{\beta}}) \left( \int \frac{1}{\mu^{(\beta)}} \mathbf{d}\mathbf{m} \right)^{-\frac{1}{2}} \\ &\leq O(e^{\frac{C_1 R' + C_2}{\beta}}) (\min \mu^{(\beta)})^{-\frac{1}{2}} \\ &\leq O(e^{\frac{C_1 R' + C_2}{\beta}}) (e^{\frac{\min \mathcal{L} - \max \mathcal{L}}{\beta}})^{-\frac{1}{2}} \\ &= O(e^{\frac{C_1 R' + C_2'}{\beta}}). \end{aligned} \quad (49)$$

Combining (48) with (49), we obtain

$$\begin{aligned} & \|(\Pi^{(\beta)})^\perp \tilde{h}_{\frac{R}{\beta}}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \\ &\leq O(e^{-\rho(\frac{R-R'}{\beta})} e^{\frac{C_1 R' + C_2'}{\beta}}) \\ &\leq O(e^{-\frac{\rho R - (\rho + C_1) R' - C_2'}{\beta}}) \\ &\leq O(e^{-\frac{Q_0}{\beta}}) \end{aligned}$$

because of (42).  $\square$

**Notation E.8.** In addition to the usual big  $O(\cdot)$  notation, we will use  $\tilde{O}(\epsilon)$  to represent an error taking value in  $[-\epsilon, \epsilon]$ .

**Lemma E.9.** For each  $j$ ,

$$\int_{U_{j, Q_0}} \tilde{h}_{\frac{R}{\beta}} \mathbf{d}\mathbf{m} = \begin{cases} 1 + \tilde{O}(\epsilon) + O(e^{-\frac{Q_0}{3\beta}}) & j = k; \\ O(e^{-\frac{Q_0}{3\beta}}) & j \neq k. \end{cases}$$

*Proof.* By the fact that  $\tilde{h}_{\frac{R}{\beta}} \mathbf{d}\mathbf{m} = \tilde{\mathcal{F}}_{\frac{R}{\beta} - T_1}^{(\beta)}((\tilde{\mathcal{F}}_{T_1}^{(\beta)} \delta_{w_0})|_{U_{k, Q_1}^\epsilon})$ , condition (43), and Lemma E.6,  $\int_{U_{j, Q_0}} \tilde{h}_{\frac{R}{\beta}} \mathbf{d}\mathbf{m} = \delta_{kj} \int_{\mathbb{S}^{d-1}} \tilde{h}_{\frac{R}{\beta}} \mathbf{d}\mathbf{m} + O(e^{-\frac{Q_0}{3\beta}})$ . We then conclude with (46).  $\square$

**Corollary E.10.** For sufficiently small  $\beta$ ,

$$\begin{aligned} & \langle \tilde{h}_{\frac{R}{\beta}}, \psi_j^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} \\ &= \frac{\delta_{kj} + \tilde{O}(\epsilon) + O(e^{-\frac{Q_0}{3\beta}})}{\|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} + O(e^{-\frac{2Q_0}{\beta}}). \end{aligned}$$

*Proof.* It follows from (30) and the lemma above that

$$\begin{aligned} & \langle \tilde{h}_{\frac{R}{\beta}}, \psi_j^{(\beta)} - \hat{\chi}_j^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} \\ &= \langle (\Pi^{(\beta)})^\perp \tilde{h}_{\frac{R}{\beta}}, \psi_j^{(\beta)} - \hat{\chi}_j^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} \\ &\leq O(e^{-\frac{Q_0}{\beta}}) \cdot O(e^{-\frac{Q_0}{\beta}}) = O(e^{-\frac{2Q_0}{\beta}}). \end{aligned} \quad (50)$$

Together with Lemma E.9, this implies

$$\begin{aligned} & \langle \tilde{h}_{\frac{R}{\beta}}, \hat{\chi}_j^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} \\ &= \frac{\langle \tilde{h}_{\frac{R}{\beta}}, \chi_j \mu^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})}}{\|\chi_j \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} \\ &= \frac{\int_{U_{j, Q_0}} \tilde{h}_{\frac{R}{\beta}} \mu^{(\beta)} \cdot \frac{1}{\mu^{(\beta)}} \mathbf{d}\mathbf{m}}{\|\chi_j \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} \\ &= \frac{\int_{U_{j, Q_0}} \tilde{h}_{\frac{R}{\beta}} \mathbf{d}\mathbf{m}}{\|\chi_j \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} \\ &= \frac{\delta_{kj} + \tilde{O}(\epsilon) + O(e^{-\frac{Q_0}{3\beta}})}{\|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}}. \end{aligned} \quad (51)$$

The corollary is proved by adding (50) and (51).  $\square$

## E.6. Proof of Theorem 4.6

To prove Theorem 4.6, it suffices to prove Theorem E.1.

*Proof of Theorem E.1.* The constant  $R_{\text{dif}}$  will be the right hand in (42) and let  $r_{\text{dif}} = \frac{Q_0}{3}$ . For  $T \in [\frac{R_{\text{dif}}}{\beta}, e^{\frac{r_{\text{dif}}}{\beta}}]$ , write  $T = \frac{R}{\beta}$ . Then  $R$  satisfies both (42) and (43).

Part (i). Let  $U_k$  be the unique basin containing  $w_0$ . First assume that  $A \subseteq U_k$ . In this case, using (30)

$$\begin{aligned}
 & \langle \chi_A \mu^{(\beta)}, \psi_j^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} \\
 &= \langle \chi_A \mu^{(\beta)}, \hat{\chi}_j^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} \\
 & \quad + \langle \chi_A \mu^{(\beta)}, \psi_j^{(\beta)} - \hat{\chi}_j^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} \\
 &= \frac{\langle \chi_A \mu^{(\beta)}, \chi_j \mu^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})}}{\|\chi_j \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} \\
 & \quad + \|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \|\psi_j^{(\beta)} - \hat{\chi}_j^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \\
 &= \frac{\delta_{kj} \|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}^2}{\|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} + O(e^{-\frac{Q_0}{\beta}}) \|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \\
 &= \left( \frac{\delta_{kj} \|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}}{\|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} + O(e^{-\frac{Q_0}{\beta}}) \right) \|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \tag{52}
 \end{aligned}$$

We can now, from (31), (32), (52) and Corollary E.10, deduce

$$\begin{aligned}
 & \langle \Pi^{(\beta)} \tilde{h}_{\frac{R}{\beta}}, \chi_A \mu^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} \\
 &= \sum_{i=1}^m \sum_{j=1}^m K^{ij} \langle \tilde{h}_{\frac{R}{\beta}}, \psi_i^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} \\
 & \quad \cdot \langle \chi_A \mu^{(\beta)}, \psi_j^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} \\
 &= \sum_{i=1}^m \sum_{j=1}^m (\delta_{ij} + O(e^{-\frac{Q_0}{\beta}})) \\
 & \quad \cdot \left( \frac{\delta_{ki} + O(\epsilon + e^{-\frac{Q_0}{3\beta}})}{\|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} + O(e^{-\frac{2Q_0}{\beta}}) \right) \\
 & \quad \cdot \left( \frac{\delta_{kj} \|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}}{\|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} + O(e^{-\frac{Q_0}{\beta}}) \right) \\
 & \quad \cdot \|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \\
 &= \sum_{j=1}^m \left( \frac{\delta_{kj} + \tilde{O}(\epsilon) + O(e^{-\frac{Q_0}{3\beta}})}{\|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} + O(e^{-\frac{2Q_0}{\beta}}) \right) \\
 & \quad \cdot \left( \frac{\delta_{kj} \|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}}{\|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} + O(e^{-\frac{Q_0}{\beta}}) \right) \\
 & \quad \cdot \|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \tag{53}
 \end{aligned}$$

By separating the  $j = k$  and  $j \neq k$  terms, this becomes

$$\begin{aligned}
 & \langle \Pi^{(\beta)} \tilde{h}_{\frac{R}{\beta}}, \chi_A \mu^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} \\
 &= \left( \left( \frac{1 + \tilde{O}(\epsilon) + O(e^{-\frac{Q_0}{3\beta}})}{\|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} + O(e^{-\frac{2Q_0}{\beta}}) \right) \right. \\
 & \quad \cdot \left( \frac{\|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}}{\|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} + O(e^{-\frac{Q_0}{\beta}}) \right) \\
 & \quad \left. + \frac{O(e^{-\frac{Q_0}{\beta}})}{\|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} \right) \|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \\
 &= \left( \frac{\|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}}{\|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}^2} + \frac{\tilde{O}(\epsilon) + O(e^{-\frac{Q_0}{3\beta}})}{\|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} \right) \\
 & \quad \cdot \|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \\
 &= \frac{\|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}^2}{\|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}^2} \\
 & \quad + \frac{(\tilde{O}(\epsilon) + O(e^{-\frac{Q_0}{3\beta}})) \|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}}{\|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})}} \\
 &= \frac{\int_A \mu^{(\beta)} d\mathbf{m}}{\int_{U_k} \mu^{(\beta)} d\mathbf{m}} + \tilde{O}(\epsilon) + O(e^{-\frac{Q_0}{3\beta}}). \tag{54}
 \end{aligned}$$

Here we used the fact that  $\|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \leq \|\chi_k \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})} = (\int_{U_k} \mu^{(\beta)} d\mathbf{m})^{\frac{1}{2}} \leq 1$ .

On the other hand, by Lemma E.7,

$$\begin{aligned}
 & \langle (\Pi^{(\beta)})^\perp \tilde{h}_{\frac{R}{\beta}}, \chi_A \mu^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} \\
 & \leq O(e^{-\frac{Q_0}{\beta}}) \|\chi_A \mu^{(\beta)}\|_{L^2(\frac{1}{\mu^{(\beta)}})} \leq O(e^{-\frac{Q_0}{\beta}}) \tag{55}
 \end{aligned}$$

Adding the last two inequalities gives, for all subset  $A \subseteq U_k$ ,

$$\begin{aligned}
 & \langle \tilde{h}_{\frac{R}{\beta}}, \chi_A \mu^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} \\
 &= \frac{\int_A \mu^{(\beta)} d\mathbf{m}}{\int_{U_k} \mu^{(\beta)} d\mathbf{m}} + \tilde{O}(\epsilon) + O(e^{-\frac{Q_0}{3\beta}}) \tag{56}
 \end{aligned}$$

Assume now  $A \subseteq U_k^c$ . In this case, by Lemma E.9,

$$\langle \tilde{h}_{\frac{R}{\beta}}, \chi_A \mu^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} = \int_A \tilde{h}_{\frac{R}{\beta}} d\mathbf{m} = O(e^{-\frac{Q_0}{3\beta}}) \tag{57}$$

Finally, for a general subset  $A \subseteq \mathbb{S}^{d-1}$ , we decompose  $A$  into  $A \cap U_k$  and  $A \setminus U_k$  and apply (56) and (57), and conclude

that

$$\begin{aligned}
 & \int_A \tilde{h}_{\frac{R}{\beta}} \mathbf{d}\mathbf{m} \\
 &= \langle \tilde{h}_{\frac{R}{\beta}}, \chi_A \mu^{(\beta)} \rangle_{L^2(\frac{1}{\mu^{(\beta)}})} \\
 &= \frac{\int_{A \cap U_k} \mu^{(\beta)} \mathbf{d}\mathbf{m}}{\int_{U_k} \mu^{(\beta)} \mathbf{d}\mathbf{m}} + \tilde{O}(\epsilon) + O(e^{-\frac{Q_0}{3\beta}}) \\
 &= \frac{\int_A \mu^{(\beta)}|_{U_k} \mathbf{d}\mathbf{m}}{\int_{U_k} \mu^{(\beta)} \mathbf{d}\mathbf{m}} + \tilde{O}(\epsilon) + O(e^{-\frac{Q_0}{3\beta}})
 \end{aligned} \tag{58}$$

By (47), the relation (58) is equivalent to

$$\begin{aligned}
 & \text{dist}_{\text{TV}}\left(\tilde{\mathcal{F}}_{\frac{R}{\beta}}^{(\beta)} \delta_{\mathbf{w}_0}, \frac{\mu^{(\beta)}|_{U_k}}{\int_{U_k} \mu^{(\beta)} \mathbf{d}\mathbf{m}} \mathbf{d}\mathbf{m}\right) \\
 & \leq \epsilon + O(e^{-\frac{Q_0}{3\beta}}) \leq 2\epsilon.
 \end{aligned} \tag{59}$$

for sufficiently small  $\beta$ . Let  $T = \frac{R}{\beta}$ . After renaming  $2\epsilon$  as  $\epsilon$  and rewriting  $\tilde{\mathcal{F}}_{\frac{R}{\beta}}^{(\beta)} \delta_{\mathbf{w}_0} = \mathcal{P}_{\tilde{\mathbf{W}}_0^{(\beta)} = \mathbf{w}_0}(\tilde{\mathbf{W}}_T^{(\beta)})$ , in view of (43), we obtain Part (i) of Theorem E.1.

*Part (ii).* The second part of the theorem follows directly from Part (i) by disintegrating  $\nu_0$  as  $\sum_{j=1}^m \nu_0|_{U_j \cap \Lambda_\epsilon} + \nu_0|_{\Lambda_\epsilon^c}$ .  $\square$

Corollary 4.7 follows from Theorem 4.6 and the following lemma:

**Lemma E.11.** *In the setting of Theorem 4.6, given  $\epsilon > 0$ , for sufficiently small  $\beta$ ,*

(i) *For all  $i$  with  $\mathcal{L}(\mathbf{z}_i) > \min \mathcal{L}$ ,*

$$\text{dist}_{\text{TV}}\left(\frac{\mu^{(\beta)}|_{U_i}}{\int_{U_i} \mu^{(\beta)} \mathbf{d}\mathbf{m}}, \mu^{(\beta)} \mathbf{d}\mathbf{m}\right) \geq 1 - \epsilon;$$

(ii) *For all  $i$  with  $\mathcal{L}(\mathbf{z}_i) = \min \mathcal{L}$ ,*

$$\begin{aligned}
 & \text{dist}_{\text{TV}}\left(\frac{\mu^{(\beta)}|_{U_i}}{\int_{U_i} \mu^{(\beta)} \mathbf{d}\mathbf{m}}, \mu^{(\beta)} \mathbf{d}\mathbf{m}\right) \\
 & \geq 1 - \frac{(\det \bar{\nabla}^2 \mathcal{L}(\mathbf{z}_i))^{-\frac{1}{2}}}{\sum_{k: \mathcal{L}(\mathbf{z}_k) = \min \mathcal{L}} (\det \bar{\nabla}^2 \mathcal{L}(\mathbf{z}_k))^{-\frac{1}{2}}} - \epsilon.
 \end{aligned}$$

*Proof.* By (Hwang, 1980, Thm. 2.1), as  $\beta \rightarrow 0$ ,  $\mu^{(\beta)}$  converges to a probability measure  $\mu^{(0)}$  supported on  $\{\mathbf{z}_i : \mathcal{L}(\mathbf{z}_i) = \min \mathcal{L}\}$  and  $\mu^{(0)}(\{\mathbf{z}_i\}) = \frac{(\det \bar{\nabla}^2 \mathcal{L}(\mathbf{z}_i))^{-\frac{1}{2}}}{\sum_{k: \mathcal{L}(\mathbf{z}_k) = \min \mathcal{L}} (\det \bar{\nabla}^2 \mathcal{L}(\mathbf{z}_k))^{-\frac{1}{2}}}$  if  $\mathcal{L}(\mathbf{z}_i) = \min \mathcal{L}$ . Therefore:

For all  $i$  with  $\mathcal{L}(\mathbf{z}_i) > \min \mathcal{L}$ ,  $\mu^{(0)}$  is supported outside  $U_i$ , thus

$$\text{dist}_{\text{TV}}\left(\frac{\mu^{(\beta)}|_{U_i}}{\int_{U_i} \mu^{(\beta)} \mathbf{d}\mathbf{m}}, \mu^{(\beta)} \mathbf{d}\mathbf{m}\right) \geq 1;$$

For all  $i$  with  $\mathcal{L}(\mathbf{z}_i) = \min \mathcal{L}$ ,

$$\begin{aligned}
 & \text{dist}_{\text{TV}}\left(\frac{\mu^{(\beta)}|_{U_i}}{\int_{U_i} \mu^{(\beta)} \mathbf{d}\mathbf{m}}, \mu^{(0)} \mathbf{d}\mathbf{m}\right) \\
 & \geq 1 - \frac{(\det \bar{\nabla}^2 \mathcal{L}(\mathbf{z}_i))^{-\frac{1}{2}}}{\sum_{k: \mathcal{L}(\mathbf{z}_k) = \min \mathcal{L}} (\det \bar{\nabla}^2 \mathcal{L}(\mathbf{z}_k))^{-\frac{1}{2}}}.
 \end{aligned}$$

Both inequalities above are obtained by comparing the measures of  $U_i$ .

It now suffices to take sufficiently small  $\beta$  in the limit.  $\square$

*Proof of Corollary 4.7.* Suppose  $m \geq 2$ , choose  $\lambda_e < \lambda_{\text{dif}}$  sufficiently small such that  $\beta$  is small enough for Lemma E.11. Let  $\Omega_{\text{dif}} = \Lambda_\epsilon \cap \bigcup_{i: \mathcal{L}(\mathbf{z}_i) > \min \mathcal{L}} U_i$  if there is only one  $\mathbf{z}_k$  with  $\mathcal{L}(\mathbf{z}_k) = \min \mathcal{L}$ ; and  $\Omega_{\text{dif}} = \Lambda_\epsilon \cap \bigcup_i U_i$  otherwise. Note that  $\mathbf{m}(\bigcup_{i: \mathcal{L}(\mathbf{z}_i) > \min \mathcal{L}} U_i) > 0$  in the first case and  $\bigcup_i U_i = 1$  in the second case. So in the first case  $\mathbf{m}(\Omega_{\text{dif}}) > \kappa_1 - \epsilon$  for some constant  $\kappa_1 > 0$ , and in the second case  $\mathbf{m}(\Omega_{\text{dif}}) > 1 - \epsilon$ .

Suppose  $\mathbf{w} \in \Omega_{\text{dif}}$  and  $\mathbf{w} \in U_i$ . Notice that in the second case,  $1 - \frac{(\det \bar{\nabla}^2 \mathcal{L}(\mathbf{z}_i))^{-\frac{1}{2}}}{\sum_{k: \mathcal{L}(\mathbf{z}_k) = \min \mathcal{L}} (\det \bar{\nabla}^2 \mathcal{L}(\mathbf{z}_k))^{-\frac{1}{2}}} > 0$  by the non-degeneracy Assumption 4.2. By Lemma E.11, in both cases,

$$\text{dist}_{\text{TV}}\left(\frac{\mu^{(\beta)}|_{U_i}}{\int_{U_i} \mu^{(\beta)} \mathbf{d}\mathbf{m}}, \mu^{(\beta)} \mathbf{d}\mathbf{m}\right) > \kappa_2 - \epsilon$$

for some constant  $\kappa_2 > 0$ . By Theorem 4.6,

$$\text{dist}_{\text{TV}}\left(\mathcal{P}_{\tilde{\mathbf{W}}_0 = \mathbf{w}_0}(\tilde{\mathbf{W}}_t), \mu^{(\beta)} \mathbf{d}\mathbf{m}\right) > \kappa_2 - 2\epsilon.$$

The corollary follows by fixing a sufficiently small  $\epsilon < \frac{1}{4} \min(\kappa_1, \kappa_2)$  and choose  $\lambda_e$  accordingly.  $\square$

## F. Proof of the tunneling stage

Using again the time change  $T = \zeta t$ , Theorem 4.8 would follow from:

**Theorem F.1.** *Under the genericity Assumption 4.2, There exists a constant  $Q_{\text{tun}}$  such that for all  $T \geq 0$*

$$\begin{aligned}
 & \text{dist}_{\text{TV}}\left(\mathcal{P}_{\tilde{\mathbf{W}}_0^{(\beta)} = \mathbf{w}_0}(\tilde{\mathbf{W}}_T^{(\beta)}), \mu^{(\beta)} \mathbf{d}\mathbf{m}\right) \\
 & \leq O(e^{-(e^{-\frac{Q_{\text{tun}}}{\beta}})T})
 \end{aligned}$$

*holds uniformly for all initial position  $\mathbf{w}_0 \in \mathbb{S}^{d-1}$  and sufficiently small  $\beta$ .*

*Proof of Theorem F.1.* As noted in the discussion following Theorem 4.8, this part (or the  $\mathbb{R}^d$  version of it) has essentially been proved in (Shi et al., 2020). Their proof is for fast growing functions  $\mathcal{L}$  on  $\mathbb{R}^d$  and is based on the fact that

there exists  $Q$  such that for the second eigenvalue  $\lambda_2^{(\beta)}$  of  $-\mathcal{D}^{(\beta)}$ ,

$$\lambda_2^{(\beta)} \gtrsim e^{-\frac{Q}{\beta}}. \quad (60)$$

(Recall that the first eigenvalue is 0 as  $\mathcal{D}^{(\beta)}\mu^{(\beta)} = 0$ .) In our setting of the compact manifold  $\mathbb{S}^{d-1}$ , (60) was proved in (Michel, 2019, Thm 2.8), and the same argument as in (Shi et al., 2020) applies.  $\square$