

# Eliminating Sharp Minima from SGD with Truncated Heavy-tailed Noise

Xingyu Wang

Department of Industrial Engineering and Management Sciences  
Northwestern University, Evanston, IL, 60208  
`xingyuwang2017@u.northwestern.edu`

Sewoong Oh

Allen School of Computer Science & Engineering  
University of Washington, Seattle, WA, 98195  
`sewoong@cs.washington.edu`

Chang-Han Rhee

Department of Industrial Engineering and Management Sciences  
Northwestern University, Evanston, IL, 60208  
`chang-han.rhee@northwestern.edu`

June 28, 2021

## Abstract

The empirical success of deep learning is often attributed to SGD’s mysterious ability to avoid sharp local minima in the loss landscape, as sharp minima are known to lead to poor generalization. Recently, empirical evidence of *heavy-tailed* gradient noise was reported in many deep learning tasks, and it was shown in [28, 29] that SGD can *escape* sharp local minima under the presence of such heavy-tailed gradient noise, providing a partial solution to the mystery. In this work, we analyze a popular variant of SGD where gradients are truncated above a fixed threshold. We show that it achieves a stronger notion of avoiding sharp minima: it can effectively *eliminate* sharp local minima entirely from its training trajectory. We characterize the dynamics of truncated SGD driven by heavy-tailed noises. First, we show that the truncation threshold and width of the attraction field dictate the order of the first exit time from the associated local minimum. Moreover, when the objective function satisfies appropriate structural conditions, we prove that as the learning rate decreases, the dynamics of the heavy-tailed truncated SGD closely resemble those of a continuous-time Markov chain that never visits any sharp minima. Real data experiments on deep learning confirm our theoretical prediction that heavy-tailed SGD with gradient clipping finds a *flatter* local minima and achieves better generalization.

## 1 Introduction

Stochastic gradient descent (SGD) and its variants have seen unprecedented empirical successes in training deep neural networks. The training of deep neural networks is typically posed as a non-convex optimization problem without explicit regularization, but the solutions obtained by SGD often perform surprisingly well on test data. Such an unexpected generalization performance of SGD in deep neural networks are often attributed to SGD’s ability to avoid sharp local minima in the loss

landscape, which is well known to lead to poor generalization [7, 11, 14]. Despite significant efforts to explain such phenomena theoretically, understanding how SGD avoids sharp local minima still remains as a central mystery of deep learning. Recently, the heavy-tailed dynamics of SGD received significant attention, and it was suggested that the heavy tails in the stochastic gradients may be a key ingredient that facilitates SGD’s escape from sharp local minima: for example, [28] and [29] report the empirical evidence of heavy-tails in stochastic gradient noise in popular deep learning architectures (see also [30, 3]) and show that SGD can escape sharp local minima in polynomial time under the presence of the heavy-tailed gradient noise. More specifically, they view heavy-tailed SGD as discrete approximations of Lévy driven Langevin equations and argue that the amount of time SGD spends in each local minimum is proportional to the width of the associated minimum according to the metastability theory [23, 9, 10] for such simplified heavy-tailed processes.

In this paper, we study the global dynamics and long-run behavior of heavy-tailed SGD and its practical variant in depth. In particular, we consider a popular version of SGD, where the stochastic gradient is truncated above a fixed threshold. Such truncation scheme is often called *gradient clipping* and employed as default in various contexts [2, 16, 5, 21, 34, 4]. We uncover a rich mathematical structure in the dynamics of SGD under this scheme and prove that the long-run behavior of such SGD is fundamentally different from that of the pure form of SGD: in particular, under a suitable structural condition on the geometry of the loss landscape, we prove that gradient clipping *completely eliminates sharp minima from the trajectory of SGDs*. Moreover, we rigorously establish that, under small learning rates, the dynamics of clipped heavy-tailed SGD closely resemble a continuous-time Markov chain (CTMC) that never visits sharp local minima in the loss landscape. Our theoretical results provide critical insights into how heavy-tailed dynamics of SGD can be utilized to find a local minimum that generalizes better.

Figure 1 (Left, Middle) clearly illustrates these points with the histograms of the sample trajectories of SGDs. Note first that SGDs with light-tailed gradient noise—(c) and (d) of Figure 1 (Left, Middle)—never manages to escape a (sharp) minimum regardless of gradient clipping. In contrast, SGDs with heavy-tailed gradient noise—(a) and (b) of Figure 1 (Left, Middle)—easily escapes from local minima. Moreover, there is a clear difference between SGDs with gradient clipping and without gradient clipping. In (a) of Figure 1 (Left), SGD without gradient clipping spends a significant amount of time at each of all four local minima ( $\{m_1, m_2, m_3, m_4\}$ ), although it spends more time around the wide ones ( $\{m_2, m_4\}$ ) than the sharp ones ( $\{m_1, m_3\}$ ). On the other hand, in (b) of Figure 1 (Left), SGD with gradient clipping not only escapes from local minima but also avoids sharp minima ( $\{m_1, m_3\}$ ) almost completely. This means that if we stop training SGD at an arbitrary time point, it is almost guaranteed that it won’t be at a sharp minimum, effectively eliminating sharp minima from its training trajectories. Behind these phenomena is the different scaling of first exit times from different attraction fields with different widths. Figure 1 (Right) demonstrates that Theorem 1 (the dashed lines) accurately predicts the first exit times.

We also propose a novel training strategy that takes advantage of the newly discovered global dynamics of truncated heavy-tailed SGDs. Despite the evidence of heavy tails reported in numerous statistical learning tasks [29, 28, 3, 6, 8, 19, 15, 30, 34], there seem to be deep learning contexts where the heavy-tails are lacking; see, for instance, [20, 33]. Given the critical role of the heavy-tails in our theory, the natural idea is to inflate the tail distribution by carefully injecting heavy-tailed noises. At first glance, the benefits of introducing heavy-tails may seem unclear, given that we truncate the stochastic gradient in the end. Nevertheless, in light of the global dynamics of SGD revealed in our theory (demonstrated in Figure 1), truncation of heavy-tailed noises produces a strong regularization effect. As detailed in Section 4, we investigate various image classification tasks and deep neural network architectures and compare the performance of different optimization methods with or without the tail inflation and gradient clipping mechanism. Results reported in Tables 2 and 3 illustrate that the tail-inflation strategy we propose here indeed consistently improves the generalization performance of the SGD and drives the SGD iterates to reach local minima with “flatter” geometry. In contrast, the pure form of (unclipped) SGD exhibits drastically deteriorated performance when heavy-tails are

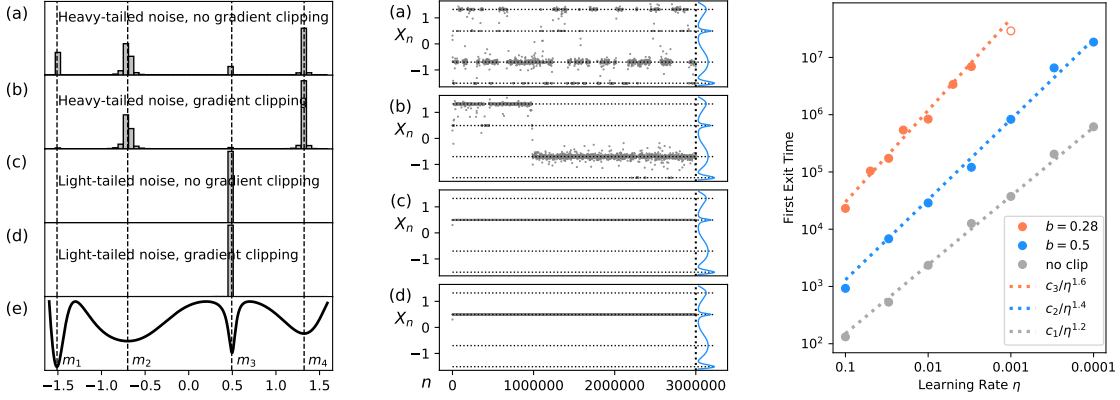


Figure 1: Elimination of Sharp Minima under Truncated Heavy-tailed Noises in  $\mathbb{R}^1$ . **(Left)** Histograms of the locations visited by SGD. With truncated heavy-tailed noises, SGD hardly ever visits the two sharp minima  $m_1$  and  $m_3$ . The objective function  $f$  is plotted at the bottom, and dashed lines are added as references for the locations of local minima. **(Middle)** Typical trajectories of SGD in different cases: (a) Heavy-tailed noises, no gradient clipping; (b) Heavy-tailed noises, gradient clipping at  $b = 0.5$ ; (c) Light-tailed noises, no gradient clipping; (d) Light-tailed noises, gradient clipping at  $b = 0.5$ . The objective function  $f$  is plotted at the right of each figure, and dashed lines are added as references for locations of the local minima. **(Right)** First Exit Time from  $\Omega_2 = (-1.3, 0.2)$ . Each dot represents the average of 20 samples of first exit time. Each dashed line shows a polynomial function  $c_i/\eta^\beta$  where  $\beta$  is predicted by Theorem 1 and  $c_i$  is chosen to fit the dots. The non-solid green dot indicates that for some of the 20 samples of the termination threshold  $5 \times 10^7$  was reached, and hence, it is an underestimation. Results in **(Left)** and **(Middle)** are obtained under learning rate  $\eta = 0.001$ .

injected. This clearly shows that injecting heavy-tailed noises alone may not be sufficient in real-world deep learning problems, and gradient clipping is indeed a key ingredient. To the best of our knowledge the idea of adding heavy-tailed noises has not been explored successfully in statistical learning context, let alone for the purpose of achieving better generalization performance for the test data. For analyses of injecting light-tailed (Gaussian) noises in deep learning tasks, see [31, 35].

The rest of the paper is organized as follows. Section 2 formulates the problem setting and characterizes the global dynamics of the SGD driven by heavy-tailed noises. Section 3 presents numerical experiments that confirm our theory. Section 4 proposes a new algorithm that artificially injects heavy tailed gradient noise in actual deep learning tasks and demonstrates the improved performance. Section 5 concludes with several potential future directions.

**Technical Contributions:** 1) We rigorously characterize the global behavior—Eyring-Kramer type formula and Markov chain model reduction for metastability—of the heavy-tailed SGD with gradient clipping. We focus on the case where the loss function is in  $\mathbb{R}^1$  with some (standard) simplifying assumptions on its geometry. Even with such assumptions, the proofs of our theorems involve substantial technical challenges since the traditional tools for analyzing SGD fail in our context due to the adaptive nature of its dynamics and non-Gaussian distributional assumptions. Moreover, we believe that the technical ideas we developed in this paper lay the foundation for a fully general theory—Freidlin-Wentzell type sample path large deviations and metastability of general heavy-tailed stochastic processes—that can address the high-dimensional loss landscapes. 2) We propose a novel computational strategy for improving the generalization performance of SGD by carefully injecting heavy-tailed noise. We test the proposed algorithm with deep learning tasks and confirm that it improves the generalization performance of SGD. This also suggests that the key phenomenon we characterize in our theory—elimination of sharp local minima—manifests itself in real deep learning

problems.

## 2 Global Dynamics of Truncated Heavy-tailed SGD

This section characterizes the global dynamics of SGD iterates with gradient clipping when applied to a non-convex objective function  $f$ . Specifically, we provide a sharp characterization of the time it takes for SGDs to exit from an attraction field: the first exit time is of order  $O(1/\eta^{1+l(\alpha-1)})$  as the learning rate  $\eta$  approaches 0, where  $\alpha > 1$  is the heavy-tailed index of noise distribution and the integer  $l$  indicates how “wide” the attraction field is when compared to the gradient clipping threshold. This characterization immediately reveals a *first-exit-time strata* on all attraction fields in the optimization landscape, where the sojourn times at wide attraction fields clearly dominate the sojourn times at other locations. Moreover, in terms of the global dynamics of clipped heavy-tailed SGD, we rigorously establish that the time-scaled version of the sample path of clipped heavy-tailed SGD converges in distribution to a continuous-time Markov chain (CTMC) as  $\eta$  approaches 0. In particular, this CTMC only visits flat minima in the optimization landscape and completely avoids any sharp local minima in narrow attraction fields. As a result, the truncated heavy-tailed noises exhibit strong regularization effect and eliminate all sharp local minima from SGD trajectories.

### 2.1 Problem setting

We make the following assumptions for the sake of the simplicity of analysis. However, as illustrated in Section 3 and 4, we believe that the gist of the phenomena we analyze—elimination of sharp local minima from the global dynamics of SGD—persists in general contexts where the domain of  $f$  is multi-dimensional, and the stationary points are not necessarily strict local optima separated from one another.

**Assumption 1.** *Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a  $C^2$  function. There exist a positive integer  $n_{\min}$ , a real number  $L \in (0, \infty)$ , and an ordered sequence of real numbers  $m_1, s_1, m_2, s_2, \dots, s_{n_{\min}-1}, m_{n_{\min}}$  such that (1)  $-L < m_1 < s_1 < m_2 < s_2 < \dots < s_{n_{\min}-1} < m_{n_{\min}} < L$ ; (2)  $f'(x) = 0$  iff  $x \in \{m_1, s_1, \dots, s_{n_{\min}-1}, m_{n_{\min}}\}$ ; (3) For any  $x \in \{m_1, m_2, \dots, m_{n_{\min}}\}$ ,  $f''(x) > 0$ ; (4) For any  $x \in \{s_1, s_2, \dots, s_{n_{\min}-1}\}$ ,  $f''(x) < 0$ .*

As illustrated in Figure 2 (Left), the assumption above requires that  $f$  has finitely many local minima (to be specific, the count is  $n_{\min}$ ) and they all lie in the interval  $[-L, L]$ . Moreover, the points  $s_1, \dots, s_{n_{\min}-1}$  naturally partition the entire real line into different regions  $\Omega_i = (s_{i-1}, s_i)$  (here we adopt the convention that  $s_0 = -\infty, s_{n_{\min}} = +\infty$ ). We call each region  $\Omega_i$  the **attraction field** of the local minimum  $m_i$ , as the gradient flow in  $\Omega_i$  always points to  $m_i$ .

Throughout the optimization procedure, given any location  $x \in \mathbb{R}$  we assume that we have access to the noisy estimator  $f'(x) - Z_n$  of the true gradient  $f'(x)$ , and  $f'(x)$  itself is difficult to evaluate. Specifically, in this work we are interested in the case where the iid sequence of noises  $(Z_n)_{n \geq 1}$  are heavy-tailed. Typically, the heavy-tailed phenomena are captured by the concept of regular variation: for a measurable function  $\phi : \mathbb{R}_+ \mapsto \mathbb{R}_+$ , we say that  $\phi$  is regularly varying at  $+\infty$  with index  $\beta$  (denoted as  $\phi \in \mathcal{RV}_\beta$ ) if  $\lim_{x \rightarrow \infty} \phi(tx)/\phi(x) = t^\beta$  for all  $t > 0$ . For details on the definition and properties of regularly varying functions, see, for example, chapter 2 of [25]. In this paper, we work with the following distributional assumption on the gradient noise. Let

$$H_+(x) \triangleq \mathbb{P}(Z_1 > x), \quad H_-(x) \triangleq \mathbb{P}(Z_1 < -x), \quad H(x) \triangleq H_+(x) + H_-(x) = \mathbb{P}(|Z_1| > x).$$

**Assumption 2.**  $\mathbb{E}Z_1 = 0$ . Furthermore, there exists some  $\alpha \in (1, \infty)$  such that function  $H(x)$  is regularly varying (at  $+\infty$ ) with index  $-\alpha$ . Besides, regarding the positive and negative tail for distribution of the noises, we have

$$\lim_{x \rightarrow \infty} \frac{H_+(x)}{H(x)} = p_+, \quad \lim_{x \rightarrow \infty} \frac{H_-(x)}{H(x)} = p_- = 1 - p_+$$

where  $p_+$  and  $p_-$  are constants in interval  $(0, 1)$ .

Roughly speaking, Assumption 2 means that the shape of the tail for the distribution of noises  $Z_n$  resembles a polynomial function  $x^{-\alpha}$ , which is much heavier than the exponential tail of Gaussian distributions. Therefore, large values of  $Z_n$  are much more likely to be observed under the Assumption 2 compared to the typical Gaussian assumption. Note that the index  $\alpha$  of regular variation encodes the heaviness of the tail—the smaller the heavier—and we are assuming that the left and right tails share the same index  $\alpha$ . We make this simplifying assumption for the purpose of clear presentation, but it is straightforward to extend this to the case where the left and right tails have different regular variation indices.

Our work concerns a popular variant of SGD where the stochastic gradient is truncated. Specifically, when updating the SGD iterates with a learning rate  $\eta > 0$ , rather than using the original noisy gradient descent step  $\eta(f'(X_n) - Z_n)$ , we will truncate it at a threshold  $b > 0$  and use  $\varphi_b(\eta(f'(X_n) - Z_n))$  instead. Here the truncation operator  $\varphi(\cdot)$  is defined as

$$\varphi_c(w) \triangleq \varphi(w, c) \triangleq (w \wedge c) \vee (-c) \quad \forall w \in \mathbb{R}, c > 0 \quad (1)$$

where  $u \wedge v = \min\{u, v\}$ ,  $u \vee v = \max\{u, v\}$ . Besides truncating the stochastic gradient, we also project the SGD into  $[-L, L]$  at each iteration; recall that  $L$  is the constant in Assumption 1. That is, the main object of our study is the stochastic process  $\{X_n^\eta\}_{n \geq 1}$  driven by the following recursion:

$$X_n^\eta \triangleq \varphi_L\left(X_n^\eta - \varphi_b(\eta(f'(X_n^\eta) - Z_n))\right). \quad (2)$$

The projection  $\varphi_L$  and truncation  $\varphi_b$  here are common practices in many statistical learning contexts as well as other optimization tasks for the purpose of ensuring that the SGD does not explode and drift to infinity. Besides, the projection also allows us to drop the sophisticated assumptions on the tail behaviors of  $f$  that are commonly seen in previous works; see, for instance, the dissipativity conditions in [19].

For technical reasons, we make the following assumption about the truncation threshold  $b > 0$ . Note that this assumption is a very mild one, as it is obviously satisfied by (Lebesgue) almost every  $b > 0$ .

**Assumption 3.** For each  $i = 1, 2, \dots, n_{\min}$ ,  $\min\{|s_i - m_i|, |s_{i-1} - m_i|\}/b$  is not an integer.

## 2.2 First exit times

Denote the SGD's first exit time from the attraction field  $\Omega_i$  with

$$\sigma_i(\eta) \triangleq \min\{n \geq 0 : X_n^\eta \notin \Omega_i\}.$$

In this section, we prove that  $\sigma_i(\eta)$  converges to an exponential distribution when scaled properly. To characterize such a scaling, we first introduce a few concepts. For each attraction field  $\Omega_i$ , define (note that  $\lceil x \rceil = \min\{n \in \mathbb{Z} : n \geq x\}$ ,  $\lfloor x \rfloor = \max\{n \in \mathbb{Z} : n \leq x\}$ )

$$r_i \triangleq \min\{|m_i - s_{i-1}|, |s_i - m_i|\}, \quad l_i^* \triangleq \lceil r_i/b \rceil. \quad (3)$$

Note that  $l_i^*$ 's in fact depend on the value of gradient clipping threshold  $b$  even though this dependency is not highlighted by the notation. Here  $r_i$  can be interpreted as the radius or the effective *width* of the attraction field, and  $l_i^*$  is the *minimum number of jumps* required to escape  $\Omega_i$  when starting from  $m_i$ . Indeed, the gradient clipping threshold  $b$  dictates that no single SGD update step can travel more than  $b$ , and to exit  $\Omega_i$  when starting from  $m_i$  (which requires the length of travel to be at least  $r_i$ ) we can see that at least  $\lceil r_i/b \rceil$  steps are required. We can interpret  $l_i^*$  as the minimum effort required to exit  $\Omega_i$ . Note that, (for a fixed  $b$ ) the minimum number of jumps  $l_i^*$

is an indicator of the width of the attraction field  $\Omega_i$ . Theorem 1 states that  $l_i^*$  dictates the order of magnitude of the first exit time as well as where the iterates  $X_n^\eta$  land on at the first exit time.

To stress the initial condition, we write  $\mathbb{P}_x$  for the probability law when conditioning on  $X_0^\eta = x$ , or simply write  $X_n^\eta(x)$ . For each  $\Omega_i$ , define a scaling function

$$\lambda_i(\eta) \triangleq H(1/\eta) \left( \frac{H(1/\eta)}{\eta} \right)^{l_i^* - 1}.$$

The distribution of the first exit time under this scaling is characterized in Theorem 1.

**Theorem 1.** *Under Assumptions 1-3, there exist constants  $q_i > 0 \forall i \in \{1, 2, \dots, n_{\min}\}$  and  $q_{i,j} \geq 0 \forall j \in \{1, 2, \dots, n_{\min}\} \setminus \{i\}$  such that*

(i) *Suppose that  $x \in \Omega_k$  for some  $k \in \{1, 2, \dots, n_{\min}\}$ . Under  $\mathbb{P}_x$ ,  $q_k \lambda_k(\eta) \sigma_k(\eta)$  converges in distribution to an exponential random variable with rate 1 as  $\eta \rightarrow 0$ .*

(ii) *For  $k, l \in \{1, 2, \dots, n_{\min}\}$  such that  $k \neq l$ , we have  $\lim_{\eta \rightarrow 0} \mathbb{P}_x(X_{\sigma_k(\eta)}^\eta \in \Omega_l) = q_{k,l}/q_k$ .*

See Section 2.5 for the proof of Theorem 1. We note here that Theorem 1 implies (i) for  $X_n^\eta$  to escape the current attraction field, say  $\Omega_i$ , it takes  $O(1/\lambda_i(\eta))$  time, and (ii) the destination is most likely to be reachable within  $l_i^*$  jumps from  $m_i$ . Before concluding this section, we state the definition of the constants  $q_i, q_{i,j}$  in Theorem 1 to show how these constants vary with the loss landscape  $f$  and the shape of gradient noise distributions (in particular  $\alpha, p_+, p_-$ ). Let  $\mathbf{Leb}_+$  denote the Lebesgue measure restricted on  $[0, \infty)$ , and define a (Borel) measure  $\nu_\alpha$  on  $\mathbb{R} \setminus \{0\}$  as

$$\nu_\alpha(dx) = \mathbb{1}\{x > 0\} \frac{\alpha p_+}{x^{\alpha+1}} + \mathbb{1}\{x < 0\} \frac{\alpha p_-}{|x|^{\alpha+1}} \quad (4)$$

where  $\alpha, p_-$ , and  $p_+$  are constants in Assumption 2. Define a Borel measure  $\mu_i$  on  $\mathbb{R}^{l_i^*} \times \left(\mathbb{R}_+\right)^{l_i^* - 1}$  as the product measure  $\mu_i = (\nu_\alpha)^{l_i^*} \times (\mathbf{Leb}_+)^{l_i^* - 1}$ . We also define mappings  $h_i$  as follows. For a real sequence  $\mathbf{w} = (w_1, w_2, \dots, w_{l_i^*})$  and a positive real number sequence  $\mathbf{t} = (t'_j)_{j=2}^{l_i^*}$ , define  $t_1 = t'_1 = 0$  and  $t_j = t'_1 + t'_2 + \dots + t'_j$  for  $j = 2, \dots, l_i^*$ . Now we define a path  $\hat{\mathbf{x}} : [0, \infty) \mapsto \mathbb{R}$  as the solution to the following ODE with jumps:

$$\hat{\mathbf{x}}(0) = \varphi_L(m_i + \varphi_b(w_1)); \quad (5)$$

$$\frac{d\hat{\mathbf{x}}(t)}{dt} = -f'(\hat{\mathbf{x}}(t)), \quad \forall t \in [t_{j-1}, t_j), \quad \forall j = 2, \dots, l_i^*; \quad (6)$$

$$\hat{\mathbf{x}}(t_j) = \varphi_L(\hat{\mathbf{x}}(t_{j-1}) + \varphi_b(w_j)), \quad \forall j = 2, \dots, l_i^*. \quad (7)$$

Now we define the mappings  $h_i : \mathbb{R}^{l_i^*} \times \left(\mathbb{R}_+\right)^{l_i^* - 1} \mapsto \mathbb{R}$  as  $h_i(\mathbf{w}, \mathbf{t}) = \hat{\mathbf{x}}(t_{l_i^*})$ . It is easy to see that  $h_i$ 's are continuous mappings. With these mappings, we define the following sets:

$$E_i \triangleq \{(\mathbf{w}, \mathbf{t}) \subseteq \mathbb{R}^{l_i^*} \times \mathbb{R}_+^{l_i^* - 1} : h_i(\mathbf{w}, \mathbf{t}) \notin \Omega_i\}; \quad (8)$$

$$E_{i,j} \triangleq \{(\mathbf{w}, \mathbf{t}) \subseteq \mathbb{R}^{l_i^*} \times \mathbb{R}_+^{l_i^* - 1} : h_i(\mathbf{w}, \mathbf{t}) \in \Omega_j\}. \quad (9)$$

Now we can define the constant  $q_i$  and  $q_{i,j}$  as follows:

$$q_i = \mu_i(E_i), \quad q_{i,j} = \mu_i(E_{i,j}) \quad \forall i \neq j. \quad (10)$$

We add a few concluding remarks regarding the intuition behind Theorem 1.

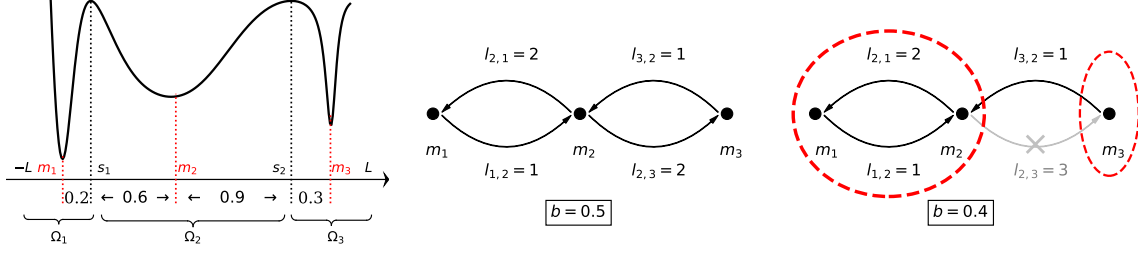


Figure 2: Typical transition graphs  $\mathcal{G}$  under different gradient clipping thresholds  $b$ . (Left) The function  $f$  illustrated here has 3 attraction fields. For the second one  $\Omega_2 = (s_1, s_2)$ , we have  $s_2 - m_2 = 0.9, m_2 - s_1 = 0.6$ . (Middle) The typical transition graph induced by  $b = 0.5$ . The entire graph  $\mathcal{G}$  is irreducible since all nodes communicate with each other. (Right) The typical transition graph induced by  $b = 0.4$ . When  $b = 0.4$ , since  $0.6 < 2b$  and  $0.9 > 2b$ , the SGD can only exit  $\Omega_2$  from the left with only 2 jumps if started from  $m_2$ . Therefore, on the graph  $\mathcal{G}$  there are two communication classes:  $G_1 = \{m_1, m_2\}, G_2 = \{m_3\}$ ;  $G_1$  is absorbing while  $G_2$  is transient.

- Suppose that  $X_n^\eta$  is started at the  $i^{\text{th}}$  local minimum  $m_i$  of  $f$ , and consider the behavior of  $X_n^\eta$  over the time period  $H_1 \triangleq \{1, \dots, \lceil t/\eta \rceil\}$  for a sufficiently large  $t$ . The heavy-tailed large deviations theory [26] and a heuristic application of the contraction principle implies that the path of  $X_{\lceil n/\eta \rceil}^\eta$  over this period will converge to the gradient flow of  $f$ , and the event that  $X_{\lceil n/\eta \rceil}^\eta$  escapes  $\Omega_i$  within this period is a (heavy-tailed) rare event. This means that the probability of such an event is of order  $(1/\eta)^{(\alpha-1)l_i^*}$ . Moreover, whenever it happens, it is almost always because  $X_n^\eta$  is shaken by *exactly*  $l_i^*$  large gradient noises of size  $\mathcal{O}(1/\eta)$ , which translates to  $l_i^*$  jumps in  $X_{\lceil n/\eta \rceil}^\eta$ 's path, while the rest of its path closely resemble the deterministic gradient flow. Moreover, conditional on the event that  $X_n^\eta$  fails to escape from  $\Omega_i$  within this period, the endpoint of the path is most likely to be close to the local minima, i.e.,  $X_{\lceil t/\eta \rceil}^\eta \approx m_i$ . This suggests that over the next time period  $H_2 \triangleq \{\lceil t/\eta \rceil + 1, \lceil t/\eta \rceil + 2, \dots, 2\lceil t/\eta \rceil\}$  of length  $\lceil t/\eta \rceil$ ,  $X_n^\eta$  will behave similarly to its behavior over the first period  $H_1$ . The same argument applies to the subsequent periods  $H_3, H_4, \dots$  as well. Therefore, over each time period of length  $\lceil t/\eta \rceil$ , there is  $(1/\eta)^{(\alpha-1)l_i^*}$  probability of exit. In view of this, the exit time should be of order  $(1/\eta)^{1+(\alpha-1)l_i^*}$  and resemble an exponential distribution when scaled properly.
- Part (i) of Theorem 1 builds on this intuition and rigorously prove that the first exit time is indeed roughly of order  $1/\lambda_i(\eta) \approx (1/\eta)^{1+(\alpha-1)l_i^*}$  and resembles an exponential distribution.
- Given this, one would expect that  $X_{\sigma_i(\eta)}^\eta$ , the location of SGD right at the time of exit, will hardly ever be farther than  $l_i^*b$  away from  $m_i$ : the length of each update is clipped by  $b$ , and there will most likely be only  $l_i^*$  large SGD steps during this successful attempt. Indeed, from the definition of  $q_{i,j}$ 's, one can see that  $q_{i,j} > 0$  if and only if  $\inf_{y \in \Omega_j} |y - m_i| < l_i^*b$ .

Summarizing the three bullet points here, we see that the minimum number of jumps  $l_i^*$  dictates *how* heavy-tailed SGD escapes an attraction field, *where* the SGD lands on upon its exit, and *when* the exit occurs.

### 2.3 Elimination of Small Attraction Fields

In this section, we show that, under proper structural assumptions on the geometry of  $f$ , the “sharp minima” of  $f$  can be effectively eliminated from the trajectory of heavy-tailed SGD. Before we state the result, we first introduce a few new concepts. Similar to the the minimum number of jumps  $l_i^*$

defined in 3, we can define the following quantities as the *minimum number of jumps to reach  $\Omega_j$  from  $m_i$*  for any  $j \neq i$ :

$$l_{i,j} = \begin{cases} \lceil (s_{j-1} - m_i)/b \rceil & \text{if } j > i, \\ \lceil (m_i - s_j)/b \rceil & \text{if } j < i. \end{cases} \quad (11)$$

Recall that Theorem 1 dictates that  $X_n^\eta$  is most likely to move out of the current attraction field, say  $\Omega_i$ , to somewhere else after  $O(1/\lambda_i(\eta))$  time steps, and the destination is most likely to be reachable within  $l_i^*$  jumps from  $m_i$ . Therefore, the transitions from  $\Omega_i$  to  $\Omega_j$  can be considered typical if  $\Omega_j$  can be reached from  $m_i$  with  $l_i^*$  jumps—that is,  $l_{i,j} = l_i^*$ . Now we define the following directed graph that only includes these typical transitions.

**Definition 1** (Typical Transition Graph). *Given a function  $f$  satisfying Assumption 1 and gradient clipping threshold  $b > 0$  satisfying Assumption 3, a directed graph  $\mathcal{G} = (V, E)$  is the corresponding typical transition graph if (1)  $V = \{m_1, \dots, m_{n_{\min}}\}$ ; (2) An edge  $(m_i \rightarrow m_j)$  is in  $E$  iff  $l_{i,j} = l_i^*$ .*

Naturally, the typical transition graph  $\mathcal{G}$  can be decomposed into different communication classes  $G_1, \dots, G_K$  that are mutually exclusive by considering the equivalence relation associated with the existence of the (two-way) paths between  $i$  and  $j$ . More specifically, for  $i \neq j$ , we say that  $i$  and  $j$  communicate if and only if there exists a path  $(m_i, m_{k_1}, \dots, m_{k_n}, m_j)$  as well as a path  $(m_j, m_{k'_1}, \dots, m_{k'_n}, m_i)$  in  $\mathcal{G}$ ; in other words, by travelling through edges on  $\mathcal{G}$ ,  $m_i$  can be reached from  $m_j$  and  $m_j$  can be reached from  $m_i$ .

We say that a communication class  $G$  is *absorbing* if there does not exist any edge  $(m_i \rightarrow m_j) \in E$  such that  $m_i \in G$  and  $m_j \notin G$ . Otherwise, we say that  $G$  is *transient*. In the case that all  $m_i$ 's communicate with each other on graph  $\mathcal{G}$ , we say  $\mathcal{G}$  is *irreducible*. See Figure 2 (Middle) for the illustration of an irreducible case. When  $\mathcal{G}$  is irreducible, we define the set of *largest* attraction fields  $M^{\text{large}} \triangleq \{m_i : i = 1, 2, \dots, n_{\min}, l_i^* = l^{\text{large}}\}$  where  $l^{\text{large}} = \max_j l_j^*$ ; recall that  $l_i^*$  characterizes the width of  $\Omega_i$ . Also, we define the longest time scale  $\lambda^{\text{large}}(\eta) = H(1/\eta) \left( \frac{H(1/\eta)}{\eta} \right)^{l^{\text{large}}-1}$ . Note that this corresponds exactly to the order of the first exit time of the largest attraction fields; see Theorem 1. The following theorem is the main result of this paper.

**Theorem 2.** *Let Assumptions 1-3 hold and assume that the graph  $\mathcal{G}$  is **irreducible**. Given  $t > 0$ ,  $\beta > 1 + (\alpha - 1)l^{\text{large}}$ , and  $x \in [-L, L]$ ,*

$$\frac{1}{\lfloor t/\eta^\beta \rfloor} \int_0^{\lfloor t/\eta^\beta \rfloor} \mathbb{1} \left\{ X_{\lfloor u \rfloor}^\eta(x) \in \bigcup_{j: m_j \notin M^{\text{large}}} \Omega_j \right\} du \rightarrow 0 \quad (12)$$

*in probability as  $\eta \rightarrow 0$ .*

We refer the readers to Section 2.6 for the proof. Here we briefly discuss the implication of the result. Suppose that we terminate the training after a reasonably long time, say,  $\lfloor t/\eta^\beta \rfloor$  iterations. Then the random variable that converges to zero in (12) is exactly the proportion of time that  $X_n^\eta$  spent in the attraction fields that are not wide. Therefore, by truncating the gradient noise of the heavy-tailed SGD, we can effectively eliminate small attraction fields from its training trajectory. In other words, it is almost guaranteed that SGD is in one of the widest attraction fields after sufficiently long training iterations.

Interestingly enough, Theorem 2 is merely a manifestation of the global dynamics of heavy-tailed SGD, and a lot more can be said even when  $\mathcal{G}$  is not irreducible. The main message can be summarized as follows: (a) SGD with truncated and heavy-tailed noise naturally partitions the entire training landscape into different regions; (b) In each region, the dynamics of  $X_n^\eta$  for small  $\eta$  closely resemble that of a continuous-time Markov chain that *only* visits local minima; (3) In particular, any sharp minima within each region is almost completely avoided by SGD.



When the typical transition graph (see Definition 1 in the main paper) is not irreducible, there will be multiple communication classes on the graph. Suppose that there are  $K$  communication classes  $G_1, \dots, G_K$ . From now on, we zoom in on a specific communication class  $G \in \{G_1, \dots, G_K\}$ . For this communication class  $G$ , define  $l_G^* \triangleq \max\{l_i^* : i = 1, 2, \dots, n_{\min}; m_i \in G\}$ . For each local minimum  $m_i \in G$ , we call its attraction field  $\Omega_i$  a *large* attraction field if  $l_i^* = l_G^*$ , and a *small* attraction field if  $l_i^* < l_G^*$ . We have thus classified all  $m_i$  in  $G$  into two groups: the ones in large attraction fields  $m_1^{\text{large}}, \dots, m_{i_G}^{\text{large}}$  and the ones in small attraction fields  $m_1^{\text{small}}, \dots, m_{i'_G}^{\text{small}}$ . Also, define a scaling function  $\lambda_G$  associated with  $G$  as  $\lambda_G(\eta) \triangleq H(1/\eta) \left( \frac{H(1/\eta)}{\eta} \right)^{l_G^* - 1}$ .

**Theorem 3.** *Under Assumptions 1-3, if  $G$  is **absorbing**, then there exists a continuous-time Markov chain  $Y$  on  $\{m_1^{\text{large}}, \dots, m_{i_G}^{\text{large}}\}$  such that for any  $x \in \Omega_i, |x| \leq L$  (where  $i \in \{1, 2, \dots, n_{\min}\}$ ) with  $m_i \in G$ , and*

$$X_{\lfloor t/\lambda_G(\eta) \rfloor}^\eta(x) \rightarrow Y_t(\pi_G(m_i)) \quad \text{as } \eta \downarrow 0$$

*in the following sense: for any positive integer  $k$  and any  $0 < t_1 < \dots < t_k$ , the random vector  $(X_{\lfloor \frac{t_k}{\lambda_G(\eta)} \rfloor}^\eta(x), \dots, X_{\lfloor \frac{t_1}{\lambda_G(\eta)} \rfloor}^\eta(x))$  converges in distribution to  $(Y_{t_1}(\pi_G(m_i)), \dots, Y_{t_k}(\pi_G(m_i)))$  as  $\eta \downarrow 0$ . Here  $\pi_G$  is a random mapping satisfying (1)  $\pi_G(m) \equiv m$  if  $m \in \{m_1^{\text{large}}, \dots, m_{i_G}^{\text{large}}\}$ ; (2)  $\pi_G(m)$  is a random variable that only takes value in  $\{m_1^{\text{large}}, \dots, m_{i_G}^{\text{large}}\}$  if  $m \in \{m_1^{\text{small}}, \dots, m_{i'_G}^{\text{small}}\}$ .*

The proof is provided in Section 2.6 and we add a few remarks here. Theorem 3 tells us that, when initialized on a absorbing class  $G$ , the dynamics of the clipped heavy-tailed SGD converge to a continuous-time Markov chain avoiding any local minima that is not in the largest attraction fields in  $G$ . Second, under small learning rate  $\eta > 0$ , if  $X_n^\eta(x)$  is initialized at  $x \in \Omega_i$  where  $\Omega_i$  is NOT a largest attraction field in  $G$ , then SGD will quickly escape  $\Omega_i$  and arrive at some  $\Omega_j$  that is indeed a largest one—i.e.,  $m_j \in M^{\text{large}}$ ; such a transition is so quick that, under time scaling  $\lambda^{\text{large}}(\eta)$ , it is almost instantaneous as if  $X_n^\eta(x)$  is actually initialized randomly at some of the largest attraction fields in  $G$ . This randomness is compressed in the random mapping  $\pi_G$ .

Next, to state the corresponding result for the *transient* case, we introduce a couple of extra definitions. We consider a version of  $X_n^\eta$  that is killed when  $X_n^\eta$  leaves  $G$ . Define stopping time

$$\tau_G(\eta) \triangleq \min\{n \geq 0 : X_n^\eta \notin \bigcup_{i: m_i \in G} \Omega_i\} \quad (13)$$

as the first time the SGD iterates leave all attraction fields in  $G$ , and we use a cemetery state  $\dagger$  to construct the following process  $X_n^{\dagger, \eta}$  as a version of  $X_n^\eta$  with killing at  $\tau_G$ :

$$X_n^{\dagger, \eta} = \begin{cases} X_n^\eta & \text{if } n < \tau_G(\eta), \\ \dagger & \text{if } n \geq \tau_G(\eta). \end{cases} \quad (14)$$

**Theorem 4.** *Under Assumptions 1-3, if  $G$  is **transient**, then there exists a continuous-time Markov chain  $Y$  **with killing** that has state space  $\{m_1^{\text{large}}, \dots, m_{i_G}^{\text{large}}, \dagger\}$  (we say the Markov chain  $Y$  is killed when it enters the absorbing cemetery state  $\dagger$ ) such that for any  $x \in \Omega_i, |x| \leq L$  (where  $i \in \{1, 2, \dots, n_{\min}\}$ ) with  $m_i \in G$ , and  $X_{\lfloor t/\lambda_G(\eta) \rfloor}^{\dagger, \eta}(x) \rightarrow Y_t(\pi_G(m_i))$  as  $\eta \downarrow 0$  in the sense of finite-dimensional distributions, where  $\pi_G$  is a random mapping satisfying (1)  $\pi_G(m) \equiv m$  if  $m \in \{m_1^{\text{large}}, \dots, m_{i_G}^{\text{large}}\}$ ; (2)  $\pi_G(m)$  is a random variable that only takes value in  $\{m_1^{\text{large}}, \dots, m_{i_G}^{\text{large}}, \dagger\}$  if  $m \in \{m_1^{\text{small}}, \dots, m_{i'_G}^{\text{small}}\}$ .*

See Section 2.6 for the proof. Theorem 3 and 4 tell us that, regardless of the initial location of SGD iterates, the elimination of small attraction fields can be observed on each communication class in the graph  $\mathcal{G}$ . In the case that  $\mathcal{G}$  is indeed irreducible, the following result follows immediately from Theorem 3 and guarantees the elimination of small attraction fields of the entire loss landscape.

**Theorem 5.** *Let Assumptions 1-2 hold. Let  $x \in \Omega_i \cap [-L, L]$  for some  $i = 1, 2, \dots, n_{\min}$ . If Assumptions 1-3 hold and  $\mathcal{G}$  is irreducible, then there exist a continuous-time Markov chain  $Y$  on  $M^{\text{large}}$  as well as a random mapping  $\pi$  such that the scaled process  $\{X_{[t/\lambda^{\text{large}}(\eta)]}^\eta(x) : t \geq 0\}$  converges to process  $\{Y_t(\pi(m_i)) : t \geq 0\}$  in the sense of finite-dimensional distributions.*

## 2.4 Implications of the theoretical results

**Systematic control of the exit times from attraction fields:** In light of the wide minima folklore, one may want to find techniques to modify the sojourn time of SGD at each attraction field. Theorem 1 suggests that the order of the first exit time (w.r.t. learning rate  $\eta$ ) is directly controlled by the gradient clipping threshold  $b$ . Recall that for an attraction field with minimum jump number  $l^*$ , Theorem 1 tells us the exit time from this attraction field is roughly of order  $(1/\eta)^{1+(\alpha-1)l^*}$ . Given the width of the attraction field, its minimum jump number  $l^*$  is dictated by gradient clipping threshold  $b$ . Therefore, gradient clipping provides us with a very systematic method to control the exit time of each attraction field. For instance, given clipping threshold  $b$ , the exit time from an attraction field with width less than  $b$  is of order  $(1/\eta)^\alpha$ , while the exit time from one larger than  $b$  is at least  $(1/\eta)^{2\alpha-1}$ , which dominates the exit time from smaller ones.

**The role of structural properties of  $\mathcal{G}$  and  $f$ :** Recall that in order for Theorem 5 to apply, the irreducibility of  $\mathcal{G}$  is required. Along with the choice of  $b$ , the geometry of function  $f$  is a deciding factor of the irreducibility. For instance, we say that  $\mathcal{G}$  is *symmetric* if for any attraction field  $\Omega_i$  such that  $i = 2, 3, \dots, n_{\min} - 1$  (so that  $\Omega_i$  is not the leftmost or rightmost one at the boundary), we have  $q_{i,i-1} > 0, q_{i,i+1} > 0$ . One can see that  $\mathcal{G}$  is symmetric if and only if, for any  $i = 2, 3, \dots, n_{\min} - 1$ ,  $|s_i - m_i| \vee |m_i - s_{i-1}| < l_i^* b$ , and symmetry is a sufficient condition for the irreducibility of  $\mathcal{G}$ . The graph illustrated in Figure 2(Middle) is symmetric, while the one in Figure 2(Right) is not. As the name suggests, in the  $\mathbb{R}^1$  case the symmetry of  $\mathcal{G}$  is more likely to hold if the shape of attraction fields in  $f$  is also nearly symmetric around its local minimum. If not, the symmetry (as well as irreducibility) of  $\mathcal{G}$  can be violated as illustrated in Figure 2, especially when a small gradient clipping threshold  $b$  is used.

Generally speaking, our results imply that, even with the truncated heavy-tailed noises, the function  $f$  needs to satisfy certain regularity conditions to ensure that SGD iterates avoid undesirable minima. This is consistent with the observations in [14]: the deep neural nets that are more trainable with SGD tend to have a much more regular structure in terms of the number and shape of local minima.

**Heavy-tailed SGD without gradient clipping:** It is worth mentioning that our results also characterize the dynamics of heavy-tailed SGDs without gradient clipping. For instance, recall that we restrict the iterates on the compact set  $[-L, L]$ . If we set truncation threshold  $b > 2L$ , then the gradient clipping at norm  $b$  becomes superfluous given the weight clipping at  $\pm L$  and the update recursion degenerates to

$$X_n^{\eta, \text{unclipped}} = \varphi_L \left( X_{n-1}^{\eta, \text{unclipped}} - \eta f'(X_{n-1}^{\eta, \text{unclipped}}) + \eta Z_n \right). \quad (15)$$

The next result follows immediately from Theorem 1 and 3.

**Corollary 6.** *There exist constants  $q_i > 0 \forall i$ ,  $q_{i,j} > 0 \forall j \neq i$  such that the following claims hold for any  $i$  and any  $x \in \Omega_i, |x| \leq L$ .*

- 1) Under  $\mathbb{P}_x$ ,  $q_i H(1/\eta) \sigma_i(\eta)$  converges weakly to an Exponential random variable with rate 1 as  $\eta \downarrow 0$ ;
- 2) For any  $j = 1, 2, \dots, n_{\min}$  such that  $j \neq i$ ,  $\lim_{\eta \downarrow 0} \mathbb{P}_x(X_{\sigma_i(\eta)}^\eta \in \Omega_j) = q_{i,j}/q_i$ .
- 3) Let  $Y$  be a continuous-time Markov chain on  $\{m_1, \dots, m_{n_{\min}}\}$  with generator matrix  $Q$  parametrized by  $Q_{i,i} = -q_i, Q_{i,j} = q_{i,j}$ . Then the scaled process  $\{X_{[t/H(1/\eta)]}^{\eta, \text{unclipped}}(x) : t > 0\}$  converges to  $\{Y_t(m_i) : t > 0\}$  in the sense of finite-dimensional distributions.

At first glance, Corollary 6 may seem similar to the results in [28] and [23]. However, the object studied in [28, 23] is different: they study the following Langevin-type stochastic differential equation (SDE) driven by an  $\alpha$ -stable Lévy process  $L_t$ :  $dY_t^\eta = -f'(Y_t^\eta)dt + \eta dL_t$ . In particular, [23] studies the metastability of  $Y_t^\eta$  and concludes that as the scaling factor  $\eta \downarrow 0$ , the first exit time and global dynamics of  $Y_t^\eta$  admit a similar characterization as described in our Theorem 6. Then Theorem 4 in [28] argues that when the learning rate  $\eta$  is sufficiently small, the distribution of the first exit time of the SGD  $X_n^\eta$  and that of the Lévy-driven Langevin SDE  $Y_t^\eta$  are similar. However, the analysis of [28] hinges critically on the assumption that  $L_t^\alpha$  is symmetric and  $\alpha$ -stable. While such an assumption is convenient for their analysis, it is a strong assumption. It implies that the gradient noise distribution belongs to a very specific parametric family and excludes all the other heavy-tailed distributions. In particular, the assumption precludes analysis of any heavy tails with finite variance. On the contrary, our work directly analyzes the SGD  $X_n^\eta$  and reveals the heavy-tailed SGD dynamics at a much greater level of generality. Specifically, we allow the noise to have general regularly varying distributions with arbitrary tail index—which includes  $\alpha$ -stable distributions as a (very) special case—and extend the characterization of global dynamics of heavy-tailed SGD to the adaptive versions of SGD where gradient clipping is applied.

## 2.5 Proof of Theorem 1

This section is organized as follows. We first introduce some key lemmas that analyze the dynamics of truncated heavy-tailed SGD before the first exit from an attraction field. Building upon these lemmas, we then present the proof of Theorem 1. We detail the proof of the technical lemmas in Appendix B.

As stated above, the proof of Theorem 1 hinges on the following two lemmas that characterize the behavior of  $X_n^\eta$  in two different phases respectively. Let  $k \in [n_{\min}]$  and  $x \in \Omega_k$ . We consider the SGD iterates initialized at  $X_0^\eta = x$ . In the first phase, the SGD iterates return to  $[m_k - 2\epsilon, m_k + 2\epsilon]$ , a small neighborhood of the local minimizer in attraction field  $\Omega_k$ ; in other words, it ends at

$$T_{\text{return}}^{(k)}(\eta, \epsilon) \triangleq \min\{n \geq 0 : X_n^\eta \in [m_k - 2\epsilon, m_k + 2\epsilon]\}. \quad (16)$$

During this phase, we show that for all learning rate  $\eta$  that is sufficiently small, it is almost always the case that  $X_n^\eta$  would quickly return to  $[m_k - 2\epsilon, m_k + 2\epsilon]$ , and it never leaves  $\Omega_k$  before  $T_{\text{return}}^{(k)}$ .

**Lemma 7.** *Under Assumptions 1-3, there exists some  $c \in (0, \infty)$  such that for any  $k \in [n_{\min}]$ , the following claim holds for all  $\epsilon > 0$  small enough:*

$$\lim_{\eta \downarrow 0} \inf_{y \in [-L, L]} \inf_{y \in (s_{k-1} + \epsilon, s_k - \epsilon)} \mathbb{P}_y \left( X_n^\eta \in \Omega_k \ \forall n \in [T_{\text{return}}^{(k)}(\eta, \epsilon)], \ T_{\text{return}}^{(k)}(\eta, \epsilon) \leq \frac{c \log(1/\epsilon)}{\eta} \right) = 1.$$

During the second phase,  $X_n^\eta$  starts from somewhere in  $[m_k - 2\epsilon, m_k + 2\epsilon]$  and tries to escape from  $\Omega_k$ , meaning that the phase ends at  $\sigma_k(\eta)$ . During this phase, we show that the distributions of the first exit time  $\sigma_k(\eta)$  and the location  $X_{\sigma_k(\eta)}^\eta$  do converge to the ones described in Theorem 1 as learning rate  $\eta$  tends to 0.

**Lemma 8.** *Let Assumptions 1-3 hold. Given  $C > 0, u > 0$  and  $k, l \in [n_{\min}]$  with  $k \neq l$ , the following claims*

$$\limsup_{\eta \downarrow 0} \sup_{x \in [-L, L], x \in (m_k - 2\epsilon, m_k + 2\epsilon)} \mathbb{P}_x \left( q_k \lambda_k(\eta) \sigma_k(\eta) > u \right) \leq C + \exp \left( - (1 - C)u \right) \quad (17)$$

$$\liminf_{\eta \downarrow 0} \inf_{x \in [-L, L], x \in (m_k - 2\epsilon, m_k + 2\epsilon)} \mathbb{P}_x \left( q_k \lambda_k(\eta) \sigma_k(\eta) > u \right) \geq -C + \exp \left( - (1 + C)u \right) \quad (18)$$

$$\limsup_{\eta \downarrow 0} \sup_{x \in [-L, L], x \in (m_k - 2\epsilon, m_k + 2\epsilon)} \mathbb{P}_x \left( X_{\sigma_k(\eta)}^\eta \in \Omega_l \right) \leq \frac{q_{k,l} + C}{q_k} \quad (19)$$

$$\liminf_{\eta \downarrow 0} \inf_{x \in [-L, L], x \in (m_k - 2\epsilon, m_k + 2\epsilon)} \mathbb{P}_x \left( X_{\sigma_k(\eta)}^\eta \in \Omega_l \right) \geq \frac{q_{k,l} - C}{q_k} \quad (20)$$

hold for all  $\epsilon > 0$  that are sufficiently small.

Now we are ready to show Theorem 1.

*Proof of Theorem 1.* Fix some  $k \in [n_{\min}]$  and  $x \in \Omega_k \cap [-L, L]$ . Let  $q_k$  and  $q_{k,l}$  be the constants in Lemma 8.

We first prove the weak convergence claim in Theorem 1(i). Arbitrarily choose some  $u > 0$  and  $C \in (0, 1)$ . It suffices to show that

$$\begin{aligned} \limsup_{\eta \downarrow 0} \mathbb{P}_x(q_k \lambda_k(\eta) \sigma_k(\eta) > u) &\leq 2C + \exp(-(1-C)u), \\ \liminf_{\eta \downarrow 0} \mathbb{P}_x(q_k \lambda_k(\eta) \sigma_k(\eta) > u) &\geq (1-C) \left( -C + \exp(-(1+C)u) \right). \end{aligned}$$

Recall the definition of the stopping time  $T_{\text{return}}^{(k)}$  in (16). Define event

$$A_k(\eta, \epsilon) \triangleq \left\{ X_n^\eta \in \Omega_k \ \forall n \in [T_{\text{return}}^{(k)}(\eta, \epsilon)], \ T_{\text{return}}^{(k)}(\eta, \epsilon) \leq \frac{c \log(1/\epsilon)}{\eta} \right\}$$

where  $c < \infty$  is the constant in Lemma 7. First, since  $x \in \Omega_k = (s_{k-1}, s_k)$ , it holds for all  $\epsilon > 0$  small enough that  $x \in (s_{k-1} + \epsilon, s_k - \epsilon)$ . Next, one can find some  $\epsilon > 0$  such that

- (Due to Lemma 7)

$$\mathbb{P}_x((A_k(\eta, \epsilon))^c) \leq C \ \forall \eta \text{ sufficiently small};$$

- (Due to (17)(18) and strong Markov property) For all  $\eta$  sufficiently small,

$$\begin{aligned} \mathbb{P}_x \left( q_k \lambda_k(\eta) (\sigma(\eta) - T_{\text{return}}^{(k)}(\eta, \epsilon)) > (1-C)u \mid A_k(\eta, \epsilon) \right) &\leq C + \exp(-(1-C)u), \\ \mathbb{P}_x \left( q_k \lambda_k(\eta) (\sigma(\eta) - T_{\text{return}}^{(k)}(\eta, \epsilon)) > u \mid A_k(\eta, \epsilon) \right) &\geq -C + \exp(-(1+C)u). \end{aligned}$$

Fix such  $\epsilon$ . Lastly, for this fixed  $\epsilon$ , due to  $\lambda_k \in \mathcal{RV}_{-1-l_k^*(\alpha-1)}(\eta)$  and  $\alpha > 1$ , we have  $q_k \lambda_k(\eta) \cdot \frac{c \log(1/\epsilon)}{\eta} < Cu$  for all  $\eta$  sufficiently small. In summary, for all  $\eta$  sufficiently small, we have

$$\begin{aligned} &\mathbb{P}_x(q_k \lambda_k(\eta) \sigma_k(\eta) > u) \\ &\leq \mathbb{P}_x((A_k(\eta, \epsilon))^c) + \mathbb{P}_x(\{q_k \lambda_k(\eta) \sigma_k(\eta) > u\} \cap A_k(\eta, \epsilon)) \\ &\leq C + \mathbb{P}_x(\{q_k \lambda_k(\eta) \sigma_k(\eta) > u\} \cap A_k(\eta, \epsilon)) \\ &= C + \mathbb{P}_x \left( q_k \lambda_k(\eta) (\sigma(\eta) - T_{\text{return}}^{(k)}(\eta, \epsilon)) > u - q_k \lambda_k(\eta) T_{\text{return}}^{(k)}(\eta, \epsilon) \mid A_k(\eta, \epsilon) \right) \cdot \mathbb{P}_x(A_k(\eta, \epsilon)) \\ &\leq C + \mathbb{P}_x \left( q_k \lambda_k(\eta) (\sigma(\eta) - T_{\text{return}}^{(k)}(\eta, \epsilon)) > (1-C)u \mid A_k(\eta, \epsilon) \right) \\ &\leq 2C + \exp(-(1-C)u) \end{aligned}$$

and

$$\begin{aligned} &\mathbb{P}_x(q_k \lambda_k(\eta) \sigma_k(\eta) > u) \\ &\geq \mathbb{P}_x(\{q_k \lambda_k(\eta) \sigma_k(\eta) > u\} \cap A_k(\eta, \epsilon)) \\ &= \mathbb{P}_x \left( q_k \lambda_k(\eta) (\sigma(\eta) - T_{\text{return}}^{(k)}(\eta, \epsilon)) > u - q_k \lambda_k(\eta) T_{\text{return}}^{(k)}(\eta, \epsilon) \mid A_k(\eta, \epsilon) \right) \cdot \mathbb{P}_x(A_k(\eta, \epsilon)) \\ &\geq \mathbb{P}_x \left( q_k \lambda_k(\eta) (\sigma(\eta) - T_{\text{return}}^{(k)}(\eta, \epsilon)) > u - q_k \lambda_k(\eta) T_{\text{return}}^{(k)}(\eta, \epsilon) \mid A_k(\eta, \epsilon) \right) \cdot (1-C) \end{aligned}$$

$$\begin{aligned}
&\geq \mathbb{P}_x \left( q_k \lambda_k(\eta) (\sigma(\eta) - T_{\text{return}}^{(k)}(\eta, \epsilon)) > u \mid A_k(\eta, \epsilon) \right) \cdot (1 - C) \\
&\geq (1 - C) \left( -C + \exp(-(1 + C)u) \right)
\end{aligned}$$

so this concludes the proof for Theorem 1(i).

In order to prove claims in Theorem 1(ii), we first observe that on event  $A_k(\eta, \epsilon)$  we must have  $\sigma(\eta) > T_{\text{return}}^{(k)}(\eta, \epsilon)$ . Next, arbitrarily choose some  $C \in (0, 1)$ , and note that it suffices to show that  $\mathbb{P}_x(X_{\sigma(\eta)}^\eta \in \Omega_l) \in \left( (1 - C) \frac{q_{k,l} - C}{q_k}, C + \frac{q_{k,l} + C}{q_k} \right)$  holds for all  $\eta$  sufficiently small. Again, we can find  $\epsilon > 0$  such that

- (Due to Lemma 7)

$$\mathbb{P}_x((A_k(\eta, \epsilon))^c) \leq C \quad \forall \eta \text{ sufficiently small};$$

- (Due to (19)(20) and strong Markov property) For all  $\eta$  sufficiently small,

$$\frac{q_{k,l} - C}{q_k} \leq \mathbb{P}_x(X_{\sigma_k(\eta)}^\eta \in \Omega_l \mid A_k(\eta, \epsilon)) \leq \frac{q_{k,l} + C}{q_k}.$$

In summary, for all  $\eta$  sufficiently small, we have

$$\begin{aligned}
\mathbb{P}_x(X_{\sigma_k(\eta)}^\eta \in \Omega_l) &\leq \mathbb{P}_x((A_k(\eta, \epsilon))^c) + \mathbb{P}_x(\{X_{\sigma_k(\eta)}^\eta \in \Omega_l\} \cap A_k(\eta, \epsilon)) \\
&\leq C + \mathbb{P}_x(X_{\sigma_k(\eta)}^\eta \in \Omega_l \mid A_k(\eta, \epsilon)) \mathbb{P}_x(A_k(\eta, \epsilon)) \\
&\leq C + \frac{q_{k,l} + C}{q_k}, \\
\mathbb{P}_x(X_{\sigma_k(\eta)}^\eta \in \Omega_l) &\geq \mathbb{P}_x(\{X_{\sigma_k(\eta)}^\eta \in \Omega_l\} \cap A_k(\eta, \epsilon)) \\
&= \mathbb{P}_x(X_{\sigma_k(\eta)}^\eta \in \Omega_l \mid A_k(\eta, \epsilon)) \mathbb{P}_x(A_k(\eta, \epsilon)) \\
&\geq (1 - C) \frac{q_{k,l} - C}{q_k}
\end{aligned}$$

and this concludes the proof.  $\square$

## 2.6 Proofs of Section 2.3

In this section, we first present some key lemmas and then use them to prove the Theorem 2-5. The proofs of the technical lemmas are deferred to Appendix C.

In order to prove Theorem 2, we will make use of the following lemma, where we show that the type of claim in Theorem 2 is indeed valid if we look at a much shorter time interval. Then when we move onto the proof of Theorem 2, it suffices to partition the entire horizon into pieces of these short time intervals, on each of which we analyze the dynamics of SGD respectively.

**Lemma 9.** *Let Assumptions 1-3 hold. Assume the graph  $\mathcal{G}$  is irreducible, and let  $\epsilon > 0, \delta > 0$  be any positive real numbers. For the following random variables (indexed by  $\eta$ )*

$$V^{\text{small}}(\eta, \epsilon, t) \triangleq \frac{1}{\lfloor t/\lambda^{\text{large}}(\eta) \rfloor} \int_0^{\lfloor t/\lambda^{\text{large}}(\eta) \rfloor} \mathbb{1}\{X_{\lfloor u \rfloor}^\eta \in \bigcup_{j: m_j \notin M^{\text{large}}} \Omega_j\} du, \quad (21)$$

the following claim holds for any sufficiently small  $t > 0$ :

$$\limsup_{\eta \downarrow 0} \sup_{x \in [-L, L]} \mathbb{P}_x(V^{\text{small}}(\eta, \epsilon, t) > \epsilon) \leq 5\delta$$

*Proof of Theorem 2.* It suffices to show that for any  $\epsilon > 0, \delta \in (0, \epsilon)$ , we have

$$\limsup_{\eta \downarrow 0} \mathbb{P}_x \left( V^*(\eta, t, \kappa) > 3\epsilon \right) < \delta.$$

Let us fix some  $\epsilon > 0, \delta \in (0, \epsilon)$ . First, let

$$N(\eta) = \lceil \frac{\lfloor t/\eta^\kappa \rfloor}{\lfloor t/\lambda^{\text{large}}(\eta) \rfloor} \rceil.$$

The regularly varying nature of  $H$  implies that  $\lambda^{\text{large}} \in \mathcal{RV}_{-(1+l^{\text{large}}(\alpha-1))}$ . Since  $\kappa > 1 + l^{\text{large}}(\alpha - 1)$ , we know that  $\lim_{\eta \downarrow 0} N(\eta) = \infty$ . Next, due to Lemma 9, we can find  $t_0 > 0$  and  $\bar{\eta} > 0$  such that for any  $\eta \in (0, \bar{\eta})$

$$\sup_{y \in [-L, L]} \mathbb{P}_y(V^{\text{small}}(\eta, \epsilon, t_0) > \epsilon) < \delta. \quad (22)$$

For any  $k \geq 1$ , define

$$V_k(\eta) \triangleq \frac{1}{\lfloor t_0/\lambda^{\text{large}}(\eta) \rfloor} \int_{(k-1)\lfloor t_0/\lambda^{\text{large}}(\eta) \rfloor}^{k\lfloor t_0/\lambda^{\text{large}}(\eta) \rfloor} \mathbb{1} \left\{ X_{\lfloor u \rfloor}^\eta \in \bigcup_{j: m_j \notin M^{\text{large}}} \Omega_j \right\} du. \quad (23)$$

It is clear from its definition that  $V_k$  stands for the proportion of time that the SGD iterates are outside of *large* attraction fields on the interval  $[(k-1)\lfloor \frac{t_0}{\lambda^{\text{large}}(\eta)} \rfloor, k\lfloor \frac{t_0}{\lambda^{\text{large}}(\eta)} \rfloor]$ . From (22) and Markov property, one can see that for any  $\eta \in (0, \bar{\eta})$

$$\sup_{x \in [-L, L]} \mathbb{P}_x(V_k(\eta) > \epsilon \mid X_0^\eta, \dots, X_{(k-1)\lfloor t_0/\lambda^{\text{large}}(\eta) \rfloor}^\eta) \leq \delta$$

uniformly for all  $k \geq 1$ . Now define  $K(\eta) \triangleq \#\{n = 1, 2, \dots, N(\eta) : V_k(\eta) > \epsilon\}$ . By a simple stochastic dominance argument, we have

$$\sup_{x \in [-L, L]} \mathbb{P}_x(K(\eta) \geq j) \leq \mathbb{P}(\text{Binomial}(N(\eta), \delta) \geq j) \quad \forall j = 1, 2, \dots.$$

Meanwhile, strong law of large numbers implies the existence of some  $\bar{\eta}_1 > 0$  such that  $\mathbb{P}(\frac{\text{Binomial}(N(\eta), \delta)}{N(\eta)} > 2\delta) < \delta$  for all  $\eta \in (0, \bar{\eta}_1)$ , thus

$$\sup_{|x| \leq L} \mathbb{P}_x(K(\eta)/N(\eta) > 2\delta) \leq \delta \quad \forall \eta \in (0, \bar{\eta}_1 \wedge \bar{\eta}).$$

Lastly, from the definition of  $K(\eta)$  and  $N(\eta)$ , we know that for all the  $N(\eta)$  intervals  $[(k-1)\lfloor \frac{t_0}{\lambda^{\text{large}}(\eta)} \rfloor, k\lfloor \frac{t_0}{\lambda^{\text{large}}(\eta)} \rfloor]$  with  $k \in [N(\eta)]$ , only on  $K(\eta)$  of them did the SGD iterates spent more than  $\epsilon$  proportion of time outside of the *large* attraction fields, hence

$$V^*(\eta, t, \kappa) \leq \epsilon + \frac{K(\eta)}{N(\eta)}.$$

In summary, we now have

$$\mathbb{P}_x(V^*(\eta, t, \kappa) > 3\epsilon) < \delta$$

for all  $\eta \in (0, \bar{\eta}_1 \wedge \bar{\eta})$ . This concludes the proof.  $\square$

To show Theorem 3 and 4, we introduce the following concepts. First, we consider the case where the SGD iterates  $X_n^\eta$  is initialized on the communication class  $G$  and  $G$  is absorbing. For some  $\Delta > 0, \eta > 0$ , define (let  $B(u, v) \triangleq [u - v, u + v]$ )

$$\sigma_0^G(\eta, \Delta) \triangleq \min\{n \geq 0 : X_n^\eta \in \bigcup_{i: m_i \in G} B(m_i, 2\Delta)\} \quad (24)$$

$$\tau_0^G(\eta, \Delta) \triangleq \min\{n \geq \sigma_0^G(\eta, \Delta) : X_n^\eta \in \bigcup_{i: m_i \notin G^{\text{small}}} B(m_i, 2\Delta)\} \quad (25)$$

$$I_0^G(\eta, \Delta) = j \iff X_{\tau_0^G}^\eta \in B(m_j, 2\Delta), \quad \tilde{I}_0^G(\eta, \Delta) = j \iff X_{\sigma_0^G}^\eta \in B(m_j, 2\Delta) \quad (26)$$

$$\sigma_k^G(\eta, \Delta) \triangleq \min\{n > \tau_{k-1}^G(\eta, \Delta) : X_n^\eta \in \bigcup_{i: m_i \in G, i \neq I_{k-1}^G} B(m_i, 2\Delta)\} \quad \forall k \geq 1 \quad (27)$$

$$\tau_k^G(\eta, \Delta) \triangleq \min\{n \geq \sigma_{k-1}^G(\eta, \Delta) : X_n^\eta \in \bigcup_{i: m_i \notin G^{\text{small}}} B(m_i, 2\Delta)\} \quad \forall k \geq 1 \quad (28)$$

$$I_k^G(\eta, \Delta) = j \iff X_{\tau_k^G}^\eta \in B(m_j, 2\Delta), \quad \tilde{I}_k^G(\eta, \Delta) = j \iff X_{\sigma_k^G}^\eta \in B(m_j, 2\Delta) \quad \forall k \geq 1. \quad (29)$$

Intuitively speaking, at each  $\tau_k^G$  the SGD iterates visits a minimizer that is not in a *small* attraction field on  $G$ , and we use  $I_k^G$  to mark the label of that large attraction field. Stopping time  $\sigma_k^G$  is the first time that SGD visits a minimizer that is *different* from the one visited at  $\tau_k^G$ , and  $\tau_{k+1}^G$  is the first time that a minimizer not in a small attraction field of  $G$  is visited again since  $\sigma_k^G$  (and including  $\sigma_k^G$ ). It is worth mentioning that, under this definition, we could have  $I_k^G = I_{k+1}^G$  for any  $k \geq 0$ . Meanwhile, define the following process that only keeps track of the updates on the labels  $(I_k^G)_{k \geq 0}$  instead of the information of the entire trajectory of  $(X_n^\eta)_{n \geq 0}$ :

$$\hat{X}_n^{\eta, \Delta} = \begin{cases} m_{I_k^G} & \text{if } \exists k \geq 0 \text{ such that } \tau_k^G \leq n < \tau_{k+1}^G \\ 0 & \text{otherwise} \end{cases} \quad (30)$$

In other words, when  $n < \tau_0^G$  we simply let  $\hat{X}_n^{\eta, \Delta} = 0$ , otherwise it is equal to the latest “marker” for the last visited wide minimum up until step  $n$ . This marker process  $\hat{X}$  jumps between the different minimizers of the large attractions in  $G$ . In particular, if for some  $n$  we have  $X_n^\eta \in B(m_j, 2\Delta)$  for some  $j$  with  $m_j \in G^{\text{large}}$ , then we must have  $\hat{X}_n^{\eta, \Delta} = m_j$ , which implies that, in this case,  $\hat{X}_n^{\eta, \Delta}$  indeed indicates the location of  $X_n^\eta$ .

Note that results in Theorem 3 and 4 concern a *scaled* version of  $X^\eta$ . Here we also define the corresponding *scaled* version of the processes

$$X_t^{*, \eta} \triangleq X_{\lfloor t/\lambda_G(\eta) \rfloor}^\eta \quad (31)$$

$$\hat{X}_t^{*, \eta, \Delta} \triangleq \hat{X}_{\lfloor t/\lambda_G(\eta) \rfloor}^{\eta, \Delta}, \quad (32)$$

a mapping  $\mathbf{T}^*(n, \eta) \triangleq n\lambda_G(\eta)$  that translates a step  $n$  to the corresponding timestamp for the scaled processes, and the following series of scaled stopping times

$$\tau_k^*(\eta, \Delta) = \mathbf{T}^*(\tau_k^G(\eta, \Delta), \eta), \quad \sigma_k^*(\eta, \Delta) = \mathbf{T}^*(\sigma_k^G(\eta, \Delta), \eta). \quad (33)$$

Before presenting the proof of Theorem 3 and 4, we make several preparations. First, our proof is inspired by ideas in [22] and here we provide a briefing. At any time  $t > 0$ , if we can show that  $X_t^{*, \eta}$  is almost always in set  $\bigcup_{i: m_i \in G^{\text{large}}} B(m_i, 2\Delta)$  (so the SGD iterates is almost always close to a minimizer in a large attraction field), then the marker process  $\hat{X}_t^{*, \eta, \Delta}$  is a pretty accurate indicator of the location of  $X_t^{*, \eta}$ , so it suffices to show that the marker process  $\hat{X}_t^{*, \eta, \Delta}$  converges to a continuous-time Markov chain  $Y$ .

Second, we construct the limiting process  $Y$  and the random mapping  $\pi_G$  before utilizing them in Theorem 3 and 4. As an important building block for this purpose, we start by considering the following discrete time Markov chain (DTMC) on the entire graph  $\mathcal{G} = (V, E)$ . Let  $\mathbf{P}^{DTMC}$  be a transition matrix with  $\mathbf{P}^{DTMC}(m_i, m_j) = \mu_i(E_{i,j})/\mu_i(E_i)$  for all  $j \neq i$ , and  $Y^{DTMC} = (Y_j^{DTMC})_{j \geq 0}$  be the DTMC induced by the said transition matrix. Let

$$T_G^{DTMC} \triangleq \min\{j \geq 0 : Y_j^{DTMC} \notin G^{\text{small}}\} \quad (34)$$

be the first time this DTMC visits a large attraction field on the communication class  $G$ , or escapes from  $G$ . Lastly, define (for any  $j$  such that  $m_j \notin G^{\text{small}}$ )

$$p_{i,j} \triangleq \mathbb{P}(Y_{T_G^{DTMC}}^{DTMC}(m_i) = m_j) \quad (35)$$

as the probability that the first large attraction field on  $G$  visited by  $Y^{DTMC}$  is  $m_j$  when initialized at  $m_i$ .

We add a comment regarding the stopping times  $T_G^{DTMC}$  and probabilities  $p_{i,j}$  defined above. In the case that  $G$  is absorbing, we have  $Y_j^{DTMC}(m_i) \in G$  for all  $j \geq 0$  if  $m_i \in G$ . Therefore, in this case, given any  $i$  with  $m_i \in G$ , we must have

$$T_G^{DTMC} = \min\{j \geq 0 : Y_j^{DTMC}(m_i) \in G^{\text{large}}\}, \quad \sum_{j: m_j \in G^{\text{large}}} p_{i,j} = 1.$$

On the contrary, when  $G$  is transient we may have  $\sum_{j: m_j \in G^{\text{large}}} p_{i,j} < 1$  and  $\sum_{j: m_j \notin G} p_{i,j} > 0$ . Lastly, whether  $G$  is absorbing or transient, we always have  $p_{i,j} = \mathbb{1}\{i = j\}$  if  $m_i \in G^{\text{large}}$ .

Next, consider the following definition of (continuous-time) jump processes.

**Definition 2.** A continuous-time process  $Y_t$  on  $\mathbb{R}$  is a  $((U_j)_{j \geq 0}, (V_j)_{j \geq 0})$  **jump process** if

$$Y_t = \begin{cases} 0 & \text{if } t < U_0 \\ \sum_{j \geq 0} V_j \mathbb{1}_{[U_0+U_1+\dots+U_j, U_0+U_1+\dots+U_{j+1})}(t) & \text{otherwise} \end{cases},$$

where  $(U_j)_{j \geq 0}$  is a sequence of non-negative random variables such that  $U_j > 0 \forall j \geq 1$  almost surely, and  $(V_j)_{j \geq 0}$  is a sequence of random variables in  $\mathbb{R}$ .

Obviously, the definition above implies that  $Y_t = V_j$  for any  $t \in [U_j, U_{j+1})$ .

Now we are ready to construct the limiting continuous-time Markov chain  $Y$ . To begin with, we address the case where  $G$  is absorbing. For any  $m' \in G^{\text{large}}$ , let  $Y(m')$  be a  $((S_k)_{k \geq 0}, (W_k)_{k \geq 0})$ -jump process where  $S_0 = 0, W_0 = m'$  and (for all  $k \geq 0$  and  $i, j$  with  $m_i \in G^{\text{large}}, m_j \notin G^{\text{small}}$ )

$$\mathbb{P}(W_{k+1} = m_j, S_{k+1} > t \mid W_k = m_i, (W_l)_{l=0}^{k-1}, (S_l)_{l=0}^k) \quad (36)$$

$$= \mathbb{P}(W_{k+1} = m_j, S_{k+1} > t \mid W_k = m_i) = \exp(-q_i t) \frac{q_{i,j}}{q_i} \forall t > 0 \quad (37)$$

where

$$q_i = \mu_i(E_i) \quad (38)$$

$$q_{i,j} = \mathbb{1}\{i \neq j\} \mu_i(E_{i,j}) + \sum_{k: m_k \in G^{\text{small}}} \mu_i(E_{i,k}) p_{k,j} \quad (39)$$

and  $p_{k,j}$  is defined in (35). In other words, conditioning on  $W_k = m_i$ , the time until next jump  $S_{k+1}$  and the jump location  $W_{k+1}$  are independent, where  $S_{k+1}$  is  $\text{Exp}(q_i)$  and  $W_{k+1} = m_j$  with probability  $q_{i,j}/q_i$ . First, it is easy to see that  $Y$  is a continuous-time Markov chain. Second, under this definition



$Y$  is allowed to have some *dummy* jumps where  $W_k = W_{k+1}$ : in this case the process  $Y_t$  does not move to a different minimizer after the  $k+1$ -th jump, and by inspecting the path of  $Y$  we cannot tell that this dummy jump has occurred. As a result, that generator  $Q$  of this Markov chain admits the form (for all  $i \neq j$  with  $m_i, m_j \in G^{\text{large}}$ )

$$Q_{i,i} = - \sum_{k: k \neq i, m_k \in G^{\text{large}}} q_{i,k}, \quad Q_{i,j} = q_{i,j}.$$

Moreover, define the following random function  $\pi_G(\cdot)$  such that for any  $m_i \in G$ ,

$$\pi_G(m_i) = \begin{cases} m_j & \text{with probability } q_{i,j}/q_i \text{ if } m_i \in G^{\text{small}} \\ m_i & \text{if } m_i \in G^{\text{large}} \end{cases} \quad (40)$$

By  $Y(\pi_G(m_i))$  we refer to the version of the Markov chain  $Y$  where we randomly initialize  $W_0 = \pi_G(m_i)$ . The following lemma is the key tool for proving Theorem 3.

**Lemma 10.** *Let Assumptions 1-3 hold. Assume that the communication class  $G$  is absorbing. Given any  $m_i \in G$ ,  $x \in \Omega_i$ , finitely many real numbers  $(t_l)_{l=1}^{k'}$  such that  $0 < t_1 < t_2 < \dots < t_{k'}$ , and a sequence of strictly positive real numbers  $(\eta_n)_{n \geq 1}$  with  $\lim_{n \rightarrow 0} \eta_n = 0$ , there exists a sequence of strictly positive real numbers  $(\Delta_n)_{n \geq 1}$  with  $\lim_n \Delta_n = 0$  such that*

- As  $n$  tends to  $\infty$ ,

$$(\hat{X}_{t_1}^{*,\eta_n,\Delta_n}(x), \dots, \hat{X}_{t_{k'}}^{*,\eta_n,\Delta_n}(x)) \Rightarrow (Y_{t_1}(\pi_G(m_i)), \dots, Y_{t_{k'}}(\pi_G(m_i))) \quad (41)$$

- For all  $k \in [k']$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}_x \left( X_{t_k}^{*,\eta_n} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, \Delta_n) \right) = 0. \quad (42)$$

Now we address the case where  $G$  is transient, let  $\dagger$  be a real number such that  $\dagger \notin [-L, L]$ , and we use  $\dagger$  as the cemetery state since the processes  $X_n^\eta$  or  $X_t^{*,\eta}$  are restricted on  $[-L, L]$ . Recall the definition of  $\tau_G$  defined in (13). Analogous to the process  $X^\dagger$  in (14), we can also define

$$X_t^{\dagger,*,\eta} = \begin{cases} X_t^{*,\eta} & \text{if } t < \mathbf{T}^*(\tau_G(\eta), \eta) \\ \dagger & \text{otherwise} \end{cases}, \quad \hat{X}_t^{\dagger,*,\eta,\Delta} = \begin{cases} \hat{X}_t^{*,\eta,\Delta} & \text{if } t < \mathbf{T}^*(\tau_G(\eta), \eta) \\ \dagger & \text{otherwise} \end{cases}. \quad (43)$$

Next, analogous to  $\tau_G$ , consider the stopping time

$$\tau_G^Y \triangleq \min\{t > 0 : Y_t \notin G\}.$$

When  $G$  is transient, due to the construction of  $Y$  we know that  $\tau_G^Y < \infty$  almost surely. The introduction of  $\tau_G^Y$  allows us to define

$$Y_t^\dagger = \begin{cases} Y_t & \text{if } t < \tau_G^Y \\ \dagger & \text{otherwise.} \end{cases} \quad (44)$$

The following Lemma will be used to prove Theorem 4.

**Lemma 11.** *Let Assumptions 1-3 hold. Assume that the communication class  $G$  is transient. Given any  $m_i \in G$ ,  $x \in \Omega_i$ , finitely many real numbers  $(t_l)_{l=1}^{k'}$  such that  $0 < t_1 < t_2 < \dots < t_{k'}$ , and a sequence of strictly positive real numbers  $(\eta_n)_{n \geq 1}$  with  $\lim_{n \rightarrow 0} \eta_n = 0$ , there exists a sequence of strictly positive real numbers  $(\Delta_n)_{n \geq 1}$  with  $\lim_n \Delta_n = 0$  such that*

- As  $n$  tends to  $\infty$ ,

$$(\hat{X}_{t_1}^{\dagger,*,\eta_n,\Delta_n}(x), \dots, \hat{X}_{t_{k'}}^{\dagger,*,\eta_n,\Delta_n}(x)) \Rightarrow (Y_{t_1}^{\dagger}(\pi_G(m_i)), \dots, Y_{t_{k'}}^{\dagger}(\pi_G(m_i))) \quad (45)$$

- For all  $k \in [k']$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}_x \left( X_{t_k}^{\dagger,*,\eta_n} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, \Delta_n) \text{ and } X_{t_k}^{\dagger,*,\eta_n} \neq \dagger \right) = 0. \quad (46)$$

*Proof of Theorem 3 and 4.* We first address the case where  $G$  is absorbing. Arbitrarily choose some  $\Delta > 0$ , a sequence of strictly positive real numbers  $(\eta_n)_{n \geq 1}$  with  $\lim_n \eta_n = 0$ , a positive integer  $k'$ , a series of real numbers  $(t_j)_{j=1}^{k'}$  with  $0 < t_1 < \dots < t_{k'}$ , and a sequence  $(w_j)_{j=1}^{k'}$  with  $w_j \in G^{\text{large}}$  for all  $j \in [k']$ . It suffices to show that

$$\lim_n \mathbb{P}_x \left( X_{t_k}^{*,\eta_n} \in B(w_k, \Delta) \ \forall k \in [k'] \right) = \mathbb{P} \left( Y_{t_k}(\pi_G(m_i)) = w_k \ \forall k \in [k'] \right).$$

Using Lemma 10, we can find a sequence of strictly positive real numbers  $(\Delta_n)_{n \geq 1}$  with  $\lim_n \Delta_n = 0$  such that (41) and (42) hold. From the weak convergence in (41), we only need to show

$$\lim_n \mathbb{P}_x \left( X_{t_k}^{*,\eta_n} \notin B(\hat{X}_{t_k}^{*,\eta_n,\Delta_n}, \Delta) \right) = 0 \ \forall k \in [k'].$$

For all  $n$  large enough, we have  $2\Delta_n < \Delta$ . For such large  $n$ , observe that

$$\begin{aligned} & \mathbb{P}_x \left( X_{t_k}^{*,\eta_n} \notin B(\hat{X}_{t_k}^{*,\eta_n,\Delta_n}, \Delta) \right) \\ & \leq \mathbb{P}_x \left( X_{t_k}^{*,\eta_n} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, 2\Delta_n) \right) \text{ due to definition of marker process } \hat{X} \text{ in (24)-(32),} \\ & \leq \mathbb{P}_x \left( X_{t_k}^{*,\eta_n} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, \Delta_n) \right) \end{aligned}$$

and by applying (42) we conclude the proof for Theorem 3.

The proof of Theorem 4 is almost identical, with the only modification being that we apply Lemma 11 instead of Lemma 10. In doing so, we are able to find a sequence of  $(\Delta_n)_{n \geq 1}$  with  $\lim_n \Delta_n = 0$  such that (45) and (46) hold. Given the weak convergence claim in (45), it suffices to show that

$$\lim_n \mathbb{P}_x \left( X_{t_k}^{\dagger,*,\eta_n} \notin B(\hat{X}_{t_k}^{\dagger,*,\eta_n,\Delta_n}, \Delta) \right) = 0 \ \forall k \in [k'].$$

If  $X_{t_k}^{\dagger,*,\eta_n} = \dagger$ , we must have  $\hat{X}_{t_k}^{\dagger,*,\eta_n,\Delta_n} = \dagger$  as well. Therefore, for all  $n$  large enough so that  $3\Delta_n < \Delta$ ,

$$\begin{aligned} & \mathbb{P}_x \left( X_{t_k}^{\dagger,*,\eta_n} \notin B(\hat{X}_{t_k}^{\dagger,*,\eta_n,\Delta_n}, \Delta) \right) \\ & = \mathbb{P}_x \left( X_{t_k}^{\dagger,*,\eta_n} \notin B(\hat{X}_{t_k}^{\dagger,*,\eta_n,\Delta_n}, \Delta), \ X_{t_k}^{\dagger,*,\eta_n} \neq \dagger \right) \\ & \leq \mathbb{P}_x \left( X_{t_k}^{\dagger,*,\eta_n} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, 2\Delta_n), \ X_{t_k}^{\dagger,*,\eta_n} \neq \dagger \right) \\ & \leq \mathbb{P}_x \left( X_{t_k}^{\dagger,*,\eta_n} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, \Delta_n), \ X_{t_k}^{\dagger,*,\eta_n} \neq \dagger \right). \end{aligned}$$

Apply (46) and we conclude the proof.  $\square$

### 3 Simulation Experiments

#### 3.1 $\mathbb{R}^1$ experiment

Our numerical experiments in this section demonstrate that, (a) as indicated by Theorem 1, the minimum jump number defined in (3) accurately characterizes the first exit times of the SGDs with clipped heavy-tailed gradient noises; (b) with clipped heavy-tailed noises, sharp minima can be effectively eliminated from SGD; (c) properties studied in this paper are exclusive to heavy-tailed noises; under light-tailed noises SGDs are trapped in sharp minima for extremely long time. The test function  $f \in \mathcal{C}^2(\mathbb{R})$  is the same one depicted in Figure 1 (Left,e). Note that  $m_1$  and  $m_3$  are sharp minima in narrow attraction fields, while  $m_2$  and  $m_4$  are flatter and located in larger attraction fields. For all experiments on this  $f$ , heavy-tailed noises are  $Z_n = 0.1U_nW_n$  where  $W_n$  are sampled from Pareto distributions with shape parameter  $\alpha = 1.2$  and the signs  $U_n$  are iid RVs such that  $\mathbb{P}(U_n = 1) = \mathbb{P}(U_n = -1) = 1/2$ ; All *light-tailed* noises are  $\mathcal{N}(0, 1)$ . See Appendix A for exact expression for test function  $f$  and experiment details.

In the first experiment, we compare the first exit time of heavy-tailed SGD (when initialized at  $-0.7$ ) from  $\Omega_2 = (-1.3, 0.2)$  under 3 different clipping mechanism: (1)  $b = 0.28$ , where the minimum jump number required to escape is  $l^* = 3$ ; (2)  $b = 0.5$ , where  $l^* = 2$ ; (3) no gradient clipping, where  $l^* = 1$  obviously. To prevent excessively long running time, each simulation run is forced to terminate after  $5 \times 10^7$  SGD iterations. According to Theorem 1, the first exit times for the aforementioned 3 clipping mechanism are of order  $(1/\eta)^{1.6}$ ,  $(1/\eta)^{1.4}$  and  $(1/\eta)^{1.2}$  respectively. As demonstrated in Figure 1 (Right), our results accurately predict how first exit time varies with learning rate and gradient clipping scheme.

Next, we investigate the global dynamics of heavy-tailed SGD. We compared the clipped case (with  $b = 0.5$ ) against the case without gradient clipping. For each scenario we simulate 10 SGD paths, each of length 10,000,000 iterations and initialized at  $X_0 = 0.3$ . Figure 1 (Left, a, b) show the histograms of the empirical distributions of SGD, and Figure 1(Middle, a,b) plots the SGD trajectories. We can clearly see that, without gradient clipping,  $X_n$  still visits the two sharp minima  $m_1, m_3$ ; under gradient clipping, the time spent at  $m_1, m_3$  is almost completely eliminated and is negligible compared to the time  $X_n$  spent at  $m_2, m_4$ , both of which are in larger attraction fields. This is exactly the dynamics predicted by Theorem 2-5: the elimination of sharp minima with truncated heavy-tailed noises. We stress that the said properties are exclusive to heavy-tailed SGD. As shown in Figure 1(Left,c,d) and Figure 1(Middle, c,d), light-tailed SGD are easily trapped at sharp minima for extremely long time.

#### 3.2 $\mathbb{R}^d$ experiment

Figure 3 illustrates that the same phenomena are observed in  $\mathbb{R}^2$ . The details of the experiment are provided in Appendix A, but here we point out that the objective function  $f$  in this experiment has several saddle points and infinitely many local minima—the local minima of  $\Omega_2$  form a line segment, which is an uncountably infinite set. Besides, under gradient clipping threshold  $b$ , attraction fields  $\Omega_1$  and  $\Omega_2$  are the *larger* ones since the escape from them requires at least two jumps. This suggests that the theoretical results proved in Section 2 holds under more general contexts than Assumptions 1-3.

### 4 Improving Generalization Performance in Deep Learning with Injected Heavy-Tailed Noises

In this section, we verify our theoretical results and demonstrate the effectiveness of clipped heavy-tailed noise in the deep learning context. Theorems in Section 2 suggest the possibility of modifying gradient noise distributions to drive SGD to “flatter” local minima and achieve better generalization performance. Meanwhile, recent experiments (such as [20, 33]) show that SGD noises in image classification tasks might be light-tailed in several cases. (As displayed in Appendix A, for tasks considered

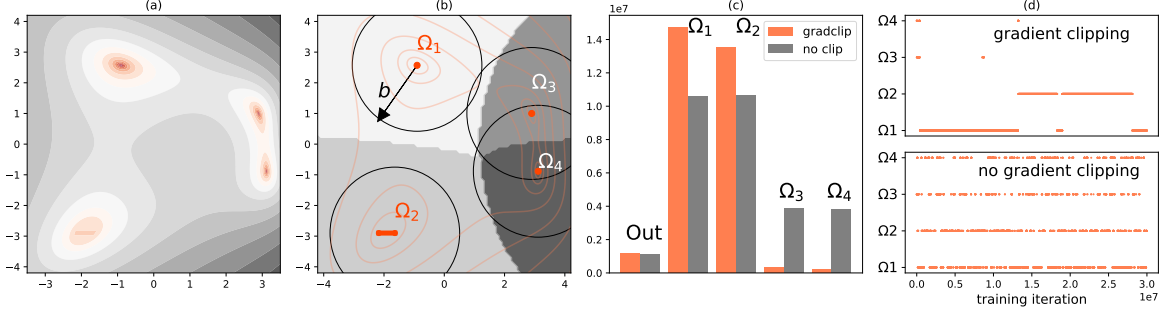


Figure 3: Experiment result of heavy-tailed SGD when optimizing the modified Himmelblau function. **(a)** Contour plot of the test function  $f$ . **(b)** Different shades of gray are used to indicate the area of the four different attraction fields  $\Omega_1, \Omega_2, \Omega_3, \Omega_4$  of  $f$ . We say that a point belongs to an attraction field  $\Omega_i$  if, when initializing at this point, the gradient descent iterates converge to the local minima in  $\Omega_i$ , which are indicated by the colored dots. The circles are added to imply whether the SGD iterates can escape from each  $\Omega_i$  with one large jump or not under clipping threshold  $b$ . **(c)** The time heavy-tailed SGD spent at different region. An iterate  $X_k$  is considered “visiting”  $\Omega_i$  if its distance to the local minimizer of  $\Omega_i$  is less than 0.5; otherwise we label  $X_k$  as “out”. **(d)** The transition trajectories of heavy-tailed SGD. The dots represent the attraction field each iteration lies in.

in this section, the gradient noise distributions are not heavy-tailed either when models are randomly initialized.) Inspired by these facts, we conduct experiments to show that SGD with gradient clipping exhibits improved performance and prefers solutions with a flatter geometry if we *make the SGD noise heavy-tailed*. Let  $\theta$  be the current model weight of a neural net during training,  $g_{SB}(\theta)$  be the typical small-batch gradient direction, and  $g_{GD}(\theta)$  be the true, deterministic gradient direction evaluated on the entire training dataset. Then by evaluating  $g_{SB}(\theta) - g_{GD}(\theta)$  we obtain a sample of gradient noise in SGD. Due to the prohibitive cost of evaluating the true gradient in real-world deep learning tasks, we consider using  $g_{SB}(\theta) - g_{LB}(\theta)$  as its approximation where  $g_{LB}$  denotes the gradient evaluated on a large batch of samples. This approximation is justified by the unbiasedness in  $\mathbb{E}g_{LB}(\theta) = g_{GD}(\theta)$ . For some heavy-tailed random variable  $Z$ , by multiplying  $Z$  with SGD noise, we obtain the following perturbed gradient direction

$$g_{heavy}(\theta) = g_{SB}(\theta) + Z(g_{SB*}(\theta) - g_{LB}(\theta)) \quad (47)$$

where  $SB$  and  $SB^*$  are two mini batches that may or may not be identical. In other words, we consider the following update recursion under gradient clipping threshold  $b$ :  $X_{k+1}^\eta = X_k^\eta - \varphi_b(\eta g_{heavy}(X_k^\eta))$  where  $\varphi_b$  is the truncation operator. In particular, we consider two different implementations: in *our method 1* (labeled as “our 1” in Table 2),  $SB$  and  $SB^*$  are chosen independently, while in *our method 2* (labeled as “our 2” in Table 2), we use the same batch of samples for  $SB$  and  $SB^*$ . In summary, by simply multiplying gradient noise with heavy-tailed random variables, we inject heavy-tailed noise into the optimization procedure.

To concretely demonstrate the effect of clipped heavy-tailed noise, we benchmark the proposed clipped heavy-tailed methods against the following optimization methods. *LB*: large-batch SGD with  $X_{k+1}^\eta = X_k^\eta - \eta g_{LB}(X_k^\eta)$ ; *SB*: small-batch SGD with  $X_{k+1}^\eta = X_k^\eta - \eta g_{SB}(X_k^\eta)$ ; *SB + Clip*: the update recursion is  $X_{k+1}^\eta = X_k^\eta - \varphi_b(\eta g_{SB}(X_k^\eta))$ ; *SB + Noise*: Our method 2 WITHOUT the gradient clipping mechanism, leading to the update recursion  $X_{k+1}^\eta = X_k^\eta - \eta g_{heavy}(X_k^\eta)$ .

The experiment setting and choice of hyperparameters are mostly adapted from [35]. We consider three different tasks: (1) LeNet [13] on corrupted FashionMNIST [32]; we use a 1200-sample subset of the original training dataset, and 200 samples points thereof are randomly assigned a label; (2)

Table 1: Hyperparameters for Training in Different Tasks

Hyperparameters	FashionMNIST, LeNet	SVHN, VGG11	CIFAR10, VGG11
learning rate	0.05	0.05	0.05
batch size for $g_{SB}$	100	100	100
training iterations	10,000	30,000	30,000
gradient clipping threshold	5	20	20
$c$	0.5	0.5	0.5
$\alpha$	1.4	1.4	1.4

Table 2: Test accuracy and expected sharpness of different methods across different tasks. The reported numbers are averaged over 5 replications. All the data points are provided in the appendix.

Test accuracy	LB	SB	SB + Clip	SB + Noise	Our 1	Our 2
FashionMNIST, LeNet	68.66%	69.20%	68.77%	64.43%	69.47%	<b>70.06%</b>
SVHN, VGG11	82.87%	85.92%	85.95%	38.85%	<b>88.42%</b>	88.37%
CIFAR10, VGG11	69.39%	74.42%	74.38%	40.50%	75.69%	<b>75.87%</b>
Expected Sharpness	LB	SB	SB + Clip	SB + Noise	Our 1	Our 2
FashionMNIST, LeNet	0.032	0.008	0.009	0.047	0.003	<b>0.002</b>
SVHN, VGG11	0.694	0.037	0.041	0.012	<b>0.002</b>	0.005
CIFAR10, VGG11	2.043	0.050	0.039	2.046	<b>0.024</b>	0.037

VGG11 [27] on SVHN [17], where we use a 25000-sample subset of the training dataset; (3) VGG11 on CIFAR10 [12] using the entire training dataset. For all tasks we use the entire test dataset when evaluating test accuracy. Whenever heavy-tailed noise is needed, the heavy-tailed multipliers used in the experiment are  $Z_n = cW_n$  where  $W_n$  are iid Pareto( $\alpha$ ) RVs. For details of the hyperparameters for training, see Table 1. Here we highlight a few points: First, within the same task, for all the 6 candidate methods will use the same  $\eta$ , batch size, training iteration, and (when needed) the same clipping threshold  $b$  and heavy-tailed multiplier  $Z_n$  for a fair comparison; the training duration is long enough so that  $LB$  and  $SB$  have attained 100% training accuracy and close-to-0 training loss long before the end of training (the exception here is “ $SB + Noise$ ” method; see Appendix A for the details); Second, to facilitate convergence to local minima for *our methods 1 and 2*, we remove heavy-tailed noise for last final 5,000 iterations and run  $LB$  instead<sup>1</sup>.

Table 2 shows that in all 3 tasks both *our method 1* and *our method 2* attain better test accuracy than the other candidate methods. Meanwhile, both methods exhibit similar test performance, implying that the implementation of the heavy-tailed method may not be a the deciding factor. We also report the *expected sharpness* metric  $\mathbb{E}_{\nu \sim \mathcal{N}(\mathbf{0}, \delta^2 \mathbf{I})} |L(\theta^* + \nu) - L(\theta^*)|$  used in [35, 18]: where  $\mathcal{N}(\mathbf{0}, \delta^2 \mathbf{I})$  is a Gaussian distribution,  $\theta^*$  is the trained model weight and  $L$  is training loss. In our experiment,

<sup>1</sup>Another interpretation of our proposed heavy-tailed methods is that it is a simplified version of GD + annealed heavy-tailed perturbation, where a detailed annealing is substituted by a two-phase training schedule.

Table 3: Our method’s gain on test accuracy persists even when applied with standard techniques for high-quality predictions: data augmentation and scheduled learning rates.

CIFAR10-VGG11	Rep 1	Rep 2	Rep 3	Rep 4	Rep 5	Average
SB+Clip	89.40%	89.41%	89.89%	89.52%	89.47%	89.54%
Our method	90.76%	90.57%	90.49%	90.85%	90.79%	<b>90.67%</b>
CIFAR100-VGG16	Rep 1	Rep 2	Rep 3	Rep 4	Rep 5	Average
SB+Clip	55.76%	56.8%	56.38%	56.35%	56.32%	56.32%
Our method	67.43%	65.12%	65.14%	65.96%	63.57%	<b>65.44%</b>

we use  $\delta = 0.01$  and the expectation is evaluated by averaging over 100 samples. We conduct 5 replications for each experiment scenario and report the averaged performance in Table 2. Smaller sharpness of clipped heavy-tailed SGD confirms its preference for minimizers that have a “flatter” geometry. In particular, the comparisons to *SB* confirm that, even in the deep learning context, clipped heavy-tailed noise can drive SGD to explore “wider” attraction fields, thus attaining better test performances.

Results in Table 2 shows that both heavy-tailed noise and gradient clipping are necessary to achieve better generalization, which is expected from our theoretical analyses. *SB* and *SB + Clip* achieve similar inferior performances, confirming that clipping does not help when noise is light-tailed. In *SB + Noise*, we inject heavy-tailed noise without gradient clipping, which still achieves an inferior performance. The poor performance of this method—even after extensive parameter tuning and engineering (see Appendix A for more details)—demonstrates the difficulty on the optimization front when heavy-tailed noise is present yet little effort is put into controlling the highly volatile gradient noises. This is aligned with the observations in [34, 4] where adaptive gradient clipping methods are proposed to improve convergence of SGD under heavy-tailed noises. This confirms that gradient clipping is crucial for heavy-tailed SGD.

Lastly, Table 3 shows that the gain of truncated heavy-tailed noise persists even when trained together with more sophisticated techniques—in particular, data augmentation and scheduled learning rates—to achieve higher performances. For experiment details, see Appendix A.

## 5 Conclusion

We characterized the global dynamics of SGD with truncated heavy-tailed gradient noise and illustrated our theoretical results with numerical experiments. Our characterizations provide key insights into the global dynamics of SGD and reveal the strong regularization effects of truncated heavy-tailed noises.

For the sake of analysis and presentation, we made some simplifying assumptions. However, as indicated by our numerical experiments, it is very likely that the same overall principle—elimination of sharp minima—persists in high-dimensional optimization landscapes that arise in deep learning contexts. We aim to develop the theories of heavy-tailed SGD at the full level of generality.

Also, the conditions—irreducibility and symmetry of graph  $\mathcal{G}$ —we impose on the geometry of  $f$  in Theorem 2 are likely to be stronger than necessary for eliminating sharp minima. In view of these, a natural next step is to investigate how the structures of state-of-the-art deep neural nets lend themselves to the SGD’s global dynamics under general conditions.

Lastly, the results of our deep learning experiments suggest the possibility of a full-fledged, rigorously justified large-batch SGD algorithm with heavy-tailed perturbations that can overcome the generalization gap between small-batch and large-batch optimization methods. For instance, it is worth exploring whether the training would benefit from a more delicate approach for inducing heavy-tailed noises or a detailed annealing of the noise magnitude.

## References

- [1] P. Billingsley. *Convergence of probability measures*. John Wiley & Sons, 2013.
- [2] L. Engstrom, A. Ilyas, S. Santurkar, D. Tsipras, F. Janoos, L. Rudolph, and A. Madry. Implementation matters in deep rl: A case study on ppo and trpo. In *International Conference on Learning Representations*, 2020.
- [3] S. Garg, J. Zhanson, E. Parisotto, A. Prasad, J. Z. Kolter, Z. C. Lipton, S. Balakrishnan, R. Salakhutdinov, and P. K. Ravikumar. On proximal policy optimization’s heavy-tailed gradients, 2021.

- [4] E. Gorbunov, M. Danilova, and A. Gasnikov. Stochastic optimization with heavy-tailed noise via accelerated gradient clipping. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 15042–15053. Curran Associates, Inc., 2020.
- [5] A. Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013.
- [6] M. Gurbuzbalaban, U. Simsekli, and L. Zhu. The heavy-tail phenomenon in sgd. *arXiv preprint arXiv:2006.04740*, 2020.
- [7] S. Hochreiter and J. Schmidhuber. Flat minima. *Neural computation*, 9(1):1–42, 1997.
- [8] L. Hodgkinson and M. W. Mahoney. Multiplicative noise and heavy tails in stochastic optimization. *arXiv preprint arXiv:2006.06293*, 2020.
- [9] P. Imkeller, I. Pavlyukevich, and M. Stauch. First exit times of non-linear dynamical systems in  $\mathbb{R}^d$  perturbed by multifractal Lévy noise. *Journal of Statistical Physics*, 141(1):94–119, 2010.
- [10] P. Imkeller, I. Pavlyukevich, and T. Wetzel. The hierarchy of exit times of Lévy-driven Langevin equations. *The European Physical Journal Special Topics*, 191(1):211–222, 2010.
- [11] N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016.
- [12] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [13] Y. LeCun, B. E. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. E. Hubbard, and L. D. Jackel. Handwritten digit recognition with a back-propagation network. In *Advances in neural information processing systems*, pages 396–404, 1990.
- [14] H. Li, Z. Xu, G. Taylor, C. Studer, and T. Goldstein. Visualizing the loss landscape of neural nets. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 6391–6401, 2018.
- [15] M. Mahoney and C. Martin. Traditional and heavy tailed self regularization in neural network models. In *International Conference on Machine Learning*, pages 4284–4293. PMLR, 2019.
- [16] S. Merity, N. S. Keskar, and R. Socher. Regularizing and optimizing LSTM language models. In *International Conference on Learning Representations*, 2018.
- [17] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- [18] B. Neyshabur, S. Bhojanapalli, D. Mcallester, and N. Srebro. Exploring generalization in deep learning. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [19] T. H. Nguyen, U. Simsekli, and G. Richard. Non-asymptotic analysis of fractional langevin monte carlo for non-convex optimization. In *International Conference on Machine Learning*, pages 4810–4819. PMLR, 2019.
- [20] A. Panigrahi, R. Somani, N. Goyal, and P. Netrapalli. Non-gaussianity of stochastic gradient noise. *arXiv preprint arXiv:1910.09626*, 2019.

- [21] R. Pascanu, T. Mikolov, and Y. Bengio. On the difficulty of training recurrent neural networks. In *International conference on machine learning*, pages 1310–1318. PMLR, 2013.
- [22] I. Pavlyukevich. Metastable behaviour of small noise lévy-driven diffusion. *arXiv preprint math/0601771*, 2005.
- [23] I. Pavlyukevich. Cooling down lévy flights. *Journal of Physics A: Mathematical and Theoretical*, 40(41):12299, 2007.
- [24] P. E. Protter. *Stochastic integration and differential equations*. Springer, 2005.
- [25] S. I. Resnick. *Heavy-tail phenomena: probabilistic and statistical modeling*. Springer Science & Business Media, 2007.
- [26] C.-H. Rhee, J. Blanchet, B. Zwart, et al. Sample path large deviations for lévy processes and random walks with regularly varying increments. *The Annals of Probability*, 47(6):3551–3605, 2019.
- [27] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [28] U. Şimşekli, M. Gürbüzbalaban, T. H. Nguyen, G. Richard, and L. Sagun. On the heavy-tailed theory of stochastic gradient descent for deep neural networks. *arXiv preprint arXiv:1912.00018*, 2019.
- [29] U. Şimşekli, L. Sagun, and M. Gurbuzbalaban. A tail-index analysis of stochastic gradient noise in deep neural networks. In *International Conference on Machine Learning*, pages 5827–5837. PMLR, 2019.
- [30] V. Srinivasan, A. Prasad, S. Balakrishnan, and P. K. Ravikumar. Efficient estimators for heavy-tailed machine learning, 2021.
- [31] Y. Wen, K. Luk, M. Gazeau, G. Zhang, H. Chan, and J. Ba. An empirical study of stochastic gradient descent with structured covariance noise. In S. Chiappa and R. Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 3621–3631. PMLR, 26–28 Aug 2020.
- [32] H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017.
- [33] Z. Xie, I. Sato, and M. Sugiyama. A diffusion theory for deep learning dynamics: Stochastic gradient descent exponentially favors flat minima. In *International Conference on Learning Representations*, 2021.
- [34] J. Zhang, T. He, S. Sra, and A. Jadbabaie. Why gradient clipping accelerates training: A theoretical justification for adaptivity. In *International Conference on Learning Representations*, 2020.
- [35] Z. Zhu, J. Wu, B. Yu, L. Wu, and J. Ma. The anisotropic noise in stochastic gradient descent: Its behavior of escaping from sharp minima and regularization effects. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7654–7663. PMLR, 09–15 Jun 2019.



## A Details of Numerical Experiments

### A.1 Details of the $\mathbb{R}^1$ simulation experiment

The function  $f$  used in the experiments is

$$f(x) = (x + 1.6)(x + 1.3)^2(x - 0.2)^2(x - 0.7)^2(x - 1.6)(0.05|1.65 - x|)^{0.6} \\ \left(1 + \frac{1}{0.01 + 4(x - 0.5)^2}\right) \left(1 + \frac{1}{0.1 + 4(x + 1.5)^2}\right) \left(1 - \frac{1}{4} \exp(-5(x + 0.8)(x + 0.8))\right).$$

The four isolated local minimizers of  $f$  are  $m_1 = -1.51, s_1 = -1.3, m_2 = -0.66, s_2 = 0.2, m_3 = 0.49, s_3 = 0.7, m_4 = 1.32$ , and in our experiment we restrict the iterates on  $[-L, L]$  with  $L = 1.6$ .

### A.2 Details of the $\mathbb{R}^d$ simulation experiment

As illustrated in the contour plot in Figure 3 (a), the function  $f$  in this experiment is a modified version of Himmelblau function, a commonly used test function for optimization algorithm. The modifications serve two purposes. First, as shown in Figure 3 (b), for the modified function the four attraction fields  $\Omega_1, \Omega_2, \Omega_3, \Omega_4$  have different sizes; in particular, under gradient clipping threshold  $b = 2.15$ , from the local minimizers of  $\Omega_1$  and  $\Omega_2$  (indicated by red dots in the corresponding area) at least two jumps are required to escape from the attraction field, while from the local minimizer in  $\Omega_3$  or  $\Omega_4$  it is possible to escape with one jump. Therefore, for the minimum jump number required to escape, we have  $l_1^* = l_2^* = 2 > l_3^* = l_4^* = 1$  in this case. Second, for the modified test function  $f$ , the local minimizer in  $\Omega_2$  is not a single point but a connected line segment, which is indicated by the dark line in bottom-left region in Figure 3 (a) and the red line segment in Figure 3 (b).

Now we describe the construction of the test function  $f$ . Let  $h$  be the Himmelblau function with expression  $h(x, y) = (x^2 + y - 11)^2 + (x + y^2 - 7)^2$ . Next, define the following transformation for coordinates:  $\phi(x, y) = (x(\exp(c_0(x - c_x) + 1)), y(\exp(c_0(x - c_x) + 1)))$ . Let the composition be  $h_\phi(x, y) = h(\phi(x - a_x, y))$ . To create the connected region of local minimizers, define the following locally “cut” version of  $h_\phi$ :

$$i(x, y) = \mathbb{1}\{x \in [b_l, b_r], |y - a_y| < b_y\}, \\ h^*(x, y) = (1 - i(x, y))h_\phi(x, y) + i(x, y) \min\{h_\phi(x, y), c_1|y - a_y|^{1.1}\}.$$

Lastly, the test function we use in the experiment is  $f = 0.1h^*$ , with  $a_x = 1.5, a_y = -2.9, b_l = -5.5, b_r = -0.5, b_y = 2.0, c_0 = 0.4, c_1 = 12$ .

In the experiment, we initialize the SGD iterates at  $X_0 = (2.9, 1.0)$ , which is very close to the local minimizer in  $\Omega_3$ . For both the clipped and unclipped SGD, we perform updates for  $3 \times 10^7$  steps under learning rate  $5 \times 10^{-4}$  and heavy-tailed noise  $Z_k = 0.75W_k$  where  $W_k$  are isotropic and the law of  $\|W_k\|$ , the size of the noise, is Pareto(1.2). For clipped SGD, we use threshold  $b = 2.15$ . To prevent the iterates from drifting to infinity, after each update  $X_k$  is projected back to the  $L_2$  ball centered at origin with radius 4.2.

### A.3 Details of the deep learning experiments comparing different methods

We first mention that the all experiments using neural networks are conducted on Nvidia GeForce GTX 1080 T, and the scripts are adapted from the ones in [35].<sup>2</sup> In Figure A.1 (Top), we display the gradient noise distribution in the three tasks after the model is randomly initialized.

The hyperparameters and training procedures for “*SB + Noise*” method is different from the other methods: we use learning rate  $\eta = 0.005$ ; for FashionMNIST task we train for 100,000 iterations and

<sup>2</sup><https://github.com/uuujf/SGDNoise>

the heavy-tailed noise is removed for the final 50,000 iterations; for SVHN and CIFAR10 tasks, we train for 150,000 iterations and heavy-tailed noise is removed for the last 70,000 iterations. Besides, for this method we always clip the model weights if its  $L_\infty$  norm exceeds 1. The reason for the extra tuning and extended training in “ $SB + Noise$ ” method is that, without the said modifications, in all three tasks we observed that the model weights quickly drift to infinity and explodes; even with the weight clipping implemented, the model performance stays at random level with no signs of improvements if we do not tune down learning rate. In Figure A.1 (Bottom), we plot the test accuracy of our method against that of the SGD for all 5 replications and 3 tasks.

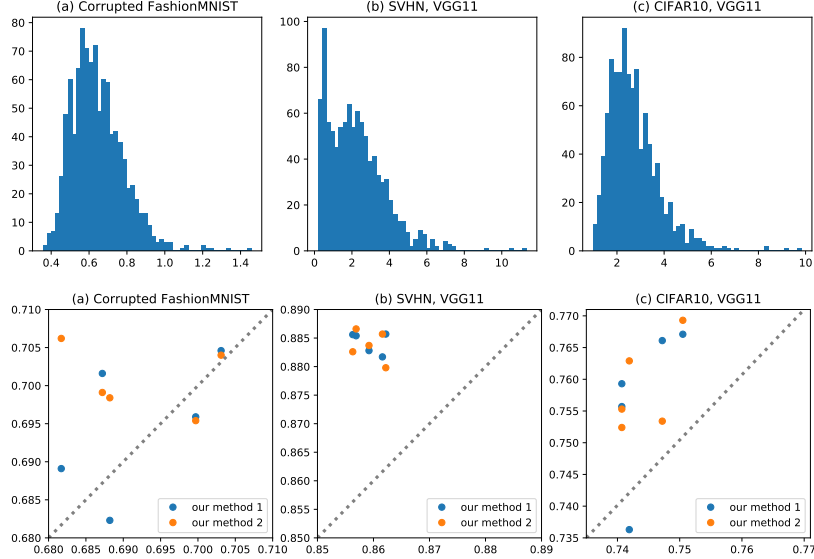


Figure A.1: **(Top)** Distribution of gradient noise in different tasks. **(Bottom)** Test accuracy of the proposed clipped heavy-tailed methods vs. test accuracy of vanilla SGD in different tasks.

#### A.4 Details of CIFAR10 experiments with data augmentation

For both methods, we train the model for 300 epochs. The initial learning rate is set at 0.1, and the training can be partitioned into two phases. In the first phase (the first 200 epochs), the learning rate is kept at a constant. In the second phase, for every 30 epoch we reduce the learning rate by half. Also, an  $L_2$  weight decaying with coefficient  $5 \times 10^{-4}$  is enforced. As for parameters for heavy-tailed noises in (47), we use  $c = 0.5$  and  $\alpha = 1.4$  in the first phase, and in the second phase we remove heavy-tailed noise and use  $SB$  to update weights. In both methods for the small-batch direction  $g_{SB}$  the batch size is 128, while for  $g_{LB}$  we evaluate the gradient on a large sample batch of size 1,024. Under the epoch number 300 and batch size 128, the count of total iterations performed during training is  $1.17 \times 10^5$ . To augment the dataset, random horizontal flipping and cropping with padding size 4 is applied for each training batch. Lastly, gradient clipping scheme is applied for both methods, and we fix  $b = 0.5$ . In other words, when the learning rate is  $\eta$  (note that due to the scheduling of learning rates,  $\eta$  will be changing throughout the training), the gradient is clipped if its  $L_2$  norm is larger than  $b/\eta$ . The scripts are adapted from the ones in <https://github.com/chengyangfu/pytorch-vgg-cifar10>.

## B Proofs of Lemma 7, 8

First, note that Assumption 1 implies the following:

- There exist  $c_0 > 0, \epsilon_0 \in (0, 1)$  such that for any  $x \in \{m_1, s_1, \dots, s_{n_{\min}-1}, m_{n_{\min}}\}, |y - x| < \epsilon_0$ ,

$$|f'(y)| > c_0|y - x|, \quad (\text{B.1})$$

and for any  $y \in [-L, L]$  such that  $|y - x| \geq \epsilon_0$  for all  $x \in \{m_1, s_1, \dots, s_{n_{\min}-1}, m_{n_{\min}}\}$ , we have

$$|f'(y)| > c_0; \quad (\text{B.2})$$

- There exist constants  $L \in (0, \infty), M \in (0, \infty)$  such that  $|m_0| < L, |m_{n_{\min}}| < L$ , and (for any  $x \in [-L, L]$ )

$$|f'(x)| \leq M, |f''(x)| \leq M. \quad (\text{B.3})$$

Recall that in (1) we define the truncation operator as  $\varphi_c(w) \triangleq \varphi(w, c) \triangleq (w \wedge c) \vee (-c)$ . Also, in (2) we defined the update recursion for SGD iterates under clipping threshold  $b$  as

$$X_n^\eta = \varphi_L \left( X_n^\eta - \varphi_b(\eta(f'(X_n^\eta) - Z_{n+1})) \right). \quad (\text{B.4})$$

Here  $\eta > 0$  is the learning rate (step length) and  $b > 0$  is the gradient clipping threshold. Also, recall that we let  $\sigma_k(\eta) \triangleq \min\{n \geq 0 : X_n^\eta \notin \Omega_k\}$  to be the first exit time of  $X_n^\eta$  from the  $k$ -th attraction field  $\Omega_k$ . Meanwhile, recall that in (3) we have defined

$$r_i \triangleq \min\{m_i - s_{i-1}, s_i - m_i\}, \quad (\text{B.5})$$

$$l_i^* \triangleq \lceil r_i/b \rceil. \quad (\text{B.6})$$

Intuitively speaking,  $l_i^*$  tells us the minimum number of jumps with size no larger than  $b$  required in order to escape the attraction field if we start from the local minimum of this attraction field  $\Omega_i$ . Lastly, recall the definition of  $H(\cdot) = \mathbb{P}(|Z_1| > \cdot)$  and the scaling function for the  $i$ -th attraction field

$$\lambda_i(\eta) \triangleq H(1/\eta) \left( \frac{H(1/\eta)}{\eta} \right)^{l_i^* - 1}.$$

The rest of this section is devoted to the proofs of Lemma 7 and Lemma 8. Specifically, Lemma 7 is an immediate Corollary of Lemma B.11, the proof of which will be provided below. The proof of Lemma 8 can be found at the end of this section.

The following three lemmas will be applied repeatedly throughout this section. The proofs are straightforward but provided in Appendix D for the sake of completeness.

**Lemma B.1.** *Given two real functions  $a : \mathbb{R}_+ \mapsto \mathbb{R}_+, b : \mathbb{R}_+ \mapsto \mathbb{R}_+$  such that  $a(\epsilon) \downarrow 0, b(\epsilon) \downarrow 0$  as  $\epsilon \downarrow 0$ , and a family of geometric RVs  $\{U(\epsilon) : \epsilon > 0\}$  with success rate  $a(\epsilon)$  (namely,  $\mathbb{P}(U(\epsilon) > k) = (1 - a(\epsilon))^k$  for  $k \in \mathbb{N}$ ), for any  $c > 1$ , there exists  $\epsilon_0 > 0$  such that for any  $\epsilon \in (0, \epsilon_0)$ ,*

$$\exp\left(-\frac{c \cdot a(\epsilon)}{b(\epsilon)}\right) \leq \mathbb{P}\left(U(\epsilon) > \frac{1}{b(\epsilon)}\right) \leq \exp\left(-\frac{a(\epsilon)}{c \cdot b(\epsilon)}\right).$$

**Lemma B.2.** *Given two real functions  $a : \mathbb{R}_+ \mapsto \mathbb{R}_+, b : \mathbb{R}_+ \mapsto \mathbb{R}_+$  such that  $a(\epsilon) \downarrow 0, b(\epsilon) \downarrow 0$  and*

$$a(\epsilon)/b(\epsilon) \rightarrow 0$$

*as  $\epsilon \downarrow 0$ , and a family of geometric RVs  $\{U(\epsilon) : \epsilon > 0\}$  with success rate  $a(\epsilon)$  (namely,  $\mathbb{P}(U(\epsilon) > k) = (1 - a(\epsilon))^k$  for  $k \in \mathbb{N}$ ), for any  $c > 1$  there exists some  $\epsilon_0 > 0$  such that for any  $\epsilon \in (0, \epsilon_0)$ ,*

$$a(\epsilon)/(c \cdot b(\epsilon)) \leq \mathbb{P}(U(\epsilon) \leq 1/b(\epsilon)) \leq c \cdot a(\epsilon)/b(\epsilon)$$

**Lemma B.3.** Suppose that a function  $g : I_g \mapsto \mathbb{R}$  (where  $I_g$  is an open interval of  $\mathbb{R}$ ) is  $g \in \mathcal{C}^2(I_g)$  and  $|g''(\cdot)| \leq C$  on its domain  $I_g$  for some constant  $C < \infty$ . For a finite integer  $n$ , a sequence of real numbers  $\{z_1, \dots, z_n\}$ , and real numbers  $x, \tilde{x} \in I_g, \eta > 0$ , consider two sequences  $\{x_k\}_{k=0, \dots, n}$  and  $\{\tilde{x}_k\}_{k=0, \dots, n}$  constructed as follows:

$$\begin{aligned} x_0 &= x \\ x_k &= x_{k-1} - \eta g'(x_{k-1}) + \eta z_k \quad \forall k = 1, 2, \dots, n \\ \tilde{x}_0 &= \tilde{x} \\ \tilde{x}_k &= \tilde{x}_{k-1} - \eta g'(\tilde{x}_{k-1}) \quad \forall k = 1, 2, \dots, n \end{aligned}$$

If we have  $x_k \in I_g, \tilde{x}_k \in I_g$ , and there exists  $\tilde{c} \in (0, \infty)$  such that  $\eta|z_1 + \dots + z_k| + |x - \tilde{x}| \leq \tilde{c}$  for  $k = 1, 2, \dots, n$ , then

$$|x_k - \tilde{x}_k| \leq \tilde{c} \cdot \exp(\eta C k) \quad \forall k = 1, 2, \dots, n.$$

To facilitate the analysis below, we introduce some additional notations. First, we will group the noises  $Z_n$  based on a threshold level  $\delta > 0$ : let us define

$$Z_n^{\leq \delta, \eta} \triangleq Z_n \mathbb{1}\{\eta|Z_n| \leq \delta\}, \quad (\text{B.7})$$

$$Z_n^{> \delta, \eta} \triangleq Z_n \mathbb{1}\{\eta|Z_n| > \delta\}. \quad (\text{B.8})$$

The former are viewed as *small noises* while the latter will be referred to as *large noises* or *large jumps*. Furthermore, for any  $j \geq 1$ , define the  $j^{\text{th}}$  arrival time and size of large jumps as

$$T_j^\eta(\delta) \triangleq \min\{n > T_{j-1}^\eta(\delta) : \eta|Z_n| > \delta\}, \quad T_0^\eta(\delta) = 0 \quad (\text{B.9})$$

$$W_j^\eta(\delta) \triangleq Z_{T_j^\eta(\delta)}. \quad (\text{B.10})$$

Next, for any  $\epsilon > 0$ , let  $\Omega_i(\epsilon) = [m_i - \epsilon, m_i + \epsilon]$  be an  $\epsilon$ -neighborhood of the local minimum  $m_i$ , and  $S_i(\epsilon) = [s_i - \epsilon, s_i + \epsilon]$  be an  $\epsilon$ -neighborhood of the local maximum  $s_i$ .

For most part of this section, we will zoom in on one of the local minima  $m_i$  and its attraction field  $\Omega_i = (s_{i-1}, s_i)$ . Without loss of generality, we assume  $m_i = 0$ , and denote the attraction field as  $\Omega = (s_-, s_+)$ . (If  $m_i$  happens to be the local minimum at the left or right boundary, then the attraction field is  $[-L, s_+)$  or  $(s_-, L]$  where the SGD iterates will be reflected at  $\pm L$ .) Henceforth we will drop the dependency on notation  $i$  when referring to this specific attraction field until the very end of this section. Throughout the proof, the following (deterministic) dynamic systems will be used frequently as benchmark processes to indicate the most likely location of the SGD iterates. Specifically, given any  $x \in \Omega$ , we use  $X_n(x)$  to indicate that the starting point is  $x$ , namely  $X_0(x) = x$ . Similarly, consider the following ODE  $\mathbf{x}^\eta(t; x)$  as

$$\mathbf{x}^\eta(0; x) = x; \quad (\text{B.11})$$

$$\frac{d\mathbf{x}^\eta(t; x)}{dt} = -\eta f'(\mathbf{x}^\eta(t; x)). \quad (\text{B.12})$$

When we use update rate  $\eta = 1$ , we will drop the dependency of  $\eta$  and simply use  $\mathbf{x}(t; x)$  to denote the process.

Based on Assumption 3, we know the existence of some constant  $\bar{\epsilon} \in (0, \epsilon_0)$  (note that  $\epsilon_0$  is the constant in (B.1)) such that

$$r \triangleq \min\{-s_-, s_+\}, \quad (\text{B.13})$$

$$l^* \triangleq \lceil r/b \rceil, \quad (\text{B.14})$$

$$(l^* - 1)b + 100l^*\bar{\epsilon} < r - 100l^*\bar{\epsilon} \quad (\text{B.15})$$

$$r + 100l^*\bar{\epsilon} < l^*b - 100l^*\bar{\epsilon}. \quad (\text{B.16})$$

Here  $r$  can be understood as the effective *radius* of the said attraction field. Also, we fix such  $\bar{\epsilon}$  small enough so that (let  $c_-^L = -f'(-L)$ ,  $c_+^L = -f'(-L)$ ), we have

$$0.9c_-^L \leq -f'(x) \leq 1.1c_-^L \quad \forall x \in [-L, -L + 100\bar{\epsilon}], \quad 0.9c_+^L \geq -f'(x) \geq 1.1c_+^L \quad \forall x \in [L - 100\bar{\epsilon}, L]. \quad (\text{B.17})$$

Similar to the definition of ODE  $\mathbf{x}^\eta$ , let us consider the following construction of ODE  $\tilde{\mathbf{x}}^\eta$  that can be understood as  $\mathbf{x}^\eta$  perturbed by  $l^*$  shocks. Specifically, consider a sequence of real numbers  $0 = t_1 < t_2 < t_3 < \dots < t_{l^*}$  and real numbers  $w_1, \dots, w_{l^*}$  where  $|w_j| \leq b$  for each  $j$ . Let  $\mathbf{t} = (t_1, \dots, t_{l^*})$ ,  $\mathbf{w} = (w_1, \dots, w_{l^*})$ . Based on these two sequences and rate  $\eta > 0$ , define  $\tilde{\mathbf{x}}^\eta(t; x)$  as

$$\tilde{\mathbf{x}}^\eta(0, x; \mathbf{t}, \mathbf{w}) = \varphi_L(x + \varphi_b(w_1)); \quad (\text{B.18})$$

$$\frac{d\tilde{\mathbf{x}}^\eta(t, x; \mathbf{t}, \mathbf{w})}{dt} = -\eta f'(\tilde{\mathbf{x}}^\eta(t, x; \mathbf{t}, \mathbf{w})) \quad \forall t \notin \{t_1, t_2, \dots, t_{l^*}\} \quad (\text{B.19})$$

$$\tilde{\mathbf{x}}^\eta(t_j, x; \mathbf{t}, \mathbf{w}) = \varphi_L(\tilde{\mathbf{x}}^\eta(t_{j-}, x; \mathbf{t}, \mathbf{w}) + \varphi_b(w_j)) \quad \forall j = 2, \dots, l^* \quad (\text{B.20})$$

Again, when  $\eta = 1$  we drop the notational dependency on  $\eta$  and use  $\tilde{\mathbf{x}}$  to denote the process. Now from (B.15)(B.16) one can easily see the following fact: there exist constants  $\bar{t}, \bar{\delta} > 0$  such that  $\tilde{\mathbf{x}}(t_{l^*}, 0; \mathbf{t}, \mathbf{w}) \notin \Omega$  (note that the starting point is 0, the local minimum) only if (under the condition that  $|w_j| \leq b \quad \forall j$ )

$$t_j - t_{j-1} \leq \bar{t} \quad \forall j = 2, 3, \dots, l^* \quad (\text{B.21})$$

$$|w_j| > \bar{\delta} \quad \forall j = 1, 2, \dots, l^*. \quad (\text{B.22})$$

The intuition is as follows: if the inter-arrival time between any of the  $l^*$  jumps is too long, then the path of  $\tilde{\mathbf{x}}^\eta(t; x)$  will drift back to the local minimum  $m_i$  so that the remaining  $l^* - 1$  shocks (whose sizes are bounded by  $b$ ) cannot overcome the radius  $r$  which is strictly larger than  $(l^* - 1)b$ ; similarly, if size of any of the shocks is too small, then since all other jumps have sizes bounded by  $b$ , the shock created by the  $l_i^*$  jumps will be smaller than  $(l^* - 1)b + 100\bar{\epsilon}$ , which is strictly less than  $r$ . We fix these constants  $\bar{t}, \bar{\delta}$  throughout the analysis, and stress again that their values are dictated by the geometry of the function  $f$ , thus *do not* vary with the accuracy parameters  $\epsilon$  and  $\delta$  mentioned earlier. In particular, choose  $\bar{\delta}$  such that  $\bar{\delta} < \bar{\epsilon}$ .

In our analysis below,  $\epsilon > 0$  will be a variable representing the level of *accuracy* in our analysis. For instance, for small  $\epsilon$ , the chance that SGD iterates will visit somewhere that is  $\epsilon$ -close to  $s_-$  or  $s_+$  (namely, the boundary of the attraction field) should be small. Consider some  $\epsilon \in (0, \epsilon_0)$  where  $\epsilon_0$  is the constant in Assumption 1. Due to (B.1)(B.2), one can see the existence of some  $g_0 > 0, c_1 < \infty$  such that

- $|f'(x)| \geq g_0$  for any  $x \in \Omega$  such that  $|x - s_-| > \epsilon_0, |x - s_+| > \epsilon_0$ ;
- Let  $\hat{t}_{\text{ODE}}(x, \eta) \triangleq \min\{t \geq 0 : \mathbf{x}^\eta(t, x) \in [-\epsilon, \epsilon]\}$  be the time that the ODE returns to a  $\epsilon$ -neighborhood of local minimum of  $\Omega$  when starting from  $x$ . As proved in Lemma 3.5 of [22], for any  $x \in \Omega$  such that  $|x - s_-| > \epsilon, |x - s_+| > \epsilon$ , we have

$$\hat{t}_{\text{ODE}}(x, \eta) \leq c_1 \frac{\log(1/\epsilon)}{\eta} \quad (\text{B.23})$$

and we define the function

$$\hat{t}(\epsilon) \triangleq c_1 \log(1/\epsilon). \quad (\text{B.24})$$

In short, given any accuracy level  $\epsilon$ , the results above give us an upper bound for how fast the ODE would return to a neighborhood of the local minimum, if the starting point is not too close to the boundary of this attraction field  $\Omega$ .

For the first few technical results established below, we show that, without large jumps, the SGD iterates  $X_n^\eta(x)$  are unlikely to show significant deviation from the deterministic gradient descent process  $\mathbf{y}_n^\eta(x)$  defined as

$$\mathbf{y}_0^\eta(x) = x, \quad (\text{B.25})$$

$$\mathbf{y}_n^\eta(x) = \mathbf{y}_{n-1}^\eta(x) - \eta f'(\mathbf{y}_{n-1}^\eta(x)). \quad (\text{B.26})$$

We are ready to state the first lemma, where we bound the distance between the gradient descent iterates  $\mathbf{y}_n^\eta(y)$  and the ODE  $\mathbf{x}^\eta(t, x)$  when the initial conditions  $x, y$  are close enough.

**Lemma B.4.** *The following claim holds for all  $\eta > 0$ : for any  $t > 0$ , we have*

$$\sup_{s \in [0, t]} |\mathbf{x}^\eta(s, x) - \mathbf{y}_{\lfloor s \rfloor}^\eta(y)| \leq (2\eta M + |x - y|) \exp(\eta M t)$$

where  $M \in (0, \infty)$  is the constant in (B.3).

*Proof.* Define a continuous-time process  $\mathbf{y}^\eta(s; y) \triangleq \mathbf{y}_{\lfloor s \rfloor}^\eta(y)$ , and note that

$$\begin{aligned} \mathbf{x}^\eta(s, x) &= \mathbf{x}^\eta(\lfloor s \rfloor, x) - \eta \int_{\lfloor s \rfloor}^s f'(\mathbf{x}^\eta(u, x)) du \\ \mathbf{x}^\eta(\lfloor s \rfloor, x) &= x - \eta \int_0^{\lfloor s \rfloor} f'(\mathbf{x}^\eta(u, x)) du \\ \mathbf{y}_{\lfloor s \rfloor}^\eta(y) &= \mathbf{y}^\eta(\lfloor s \rfloor, y) = y - \eta \int_0^{\lfloor s \rfloor} f'(\mathbf{y}^\eta(u, y)) du. \end{aligned}$$

Therefore, if we define function

$$b(u) = \mathbf{x}^\eta(u, x) - \mathbf{y}^\eta(u, y),$$

from the fact  $|f'(\cdot)| \leq M$ , one can see that  $|b(u)| \leq \eta M + |x - y|$  for any  $u \in [0, 1]$  and  $|b(1)| \leq 2\eta M + |x - y|$ . In case that  $s > 1$ , from the display above and the fact  $|f''(\cdot)| \leq M$ , we now have

$$\begin{aligned} |\mathbf{y}_{\lfloor s \rfloor}^\eta(x) - \mathbf{x}^\eta(s, x)| &\leq |b(\lfloor s \rfloor)| + \eta M; \\ |b(\lfloor s \rfloor)| &\leq \eta M \int_1^{\lfloor s \rfloor} |b(u)| du. \end{aligned}$$

From Gronwall's inequality (see Theorem 68, Chapter V of [24], where we let function  $\alpha(u)$  be  $\alpha(u) = |b(u + 1)|$ ), we have

$$|\mathbf{y}_{\lfloor s \rfloor}^\eta(x) - \mathbf{x}^\eta(s, x)| \leq (2\eta M + |x - y|) \exp(\eta M t).$$

This concludes the proof.  $\square$

Now we consider an extension of the previous Lemma in the following sense: we add perturbations to the gradient descent process and ODE, and show that, when both perturbed by  $l^*$  similar perturbations, the ODE and gradient descent process should still stay close enough. Analogous to the definition of the perturbed ODE  $\tilde{x}^\eta$  in (B.18)-(B.20), we can construct a process  $\tilde{\mathbf{Y}}^\eta$  as a perturbed gradient descent process as follows. For a sequence of integers  $0 = t_1 < t_2 < \dots < t_{l^*}$  (let  $\mathbf{t} = (t_j)_{j \geq 1}$ ) and a sequence of real numbers  $\tilde{w}_1, \dots, \tilde{w}_{l^*}$  (let  $\tilde{\mathbf{w}} = (\tilde{w}_j)_{j \geq 1}$ ) and  $y \in \mathbb{R}$ , define (for all  $n = 1, 2, \dots, t_{l^*}$ ) the perturbed gradient descent iterates with gradient clipping at  $b$  and reflection at  $\pm L$  as

$$\tilde{\mathbf{y}}_n^\eta(y; \mathbf{t}, \tilde{\mathbf{w}}) = \varphi_L \left( \tilde{\mathbf{y}}_{n-1}^\eta(y; \mathbf{t}, \tilde{\mathbf{w}}) + \varphi_b \left( -\eta f'(\tilde{\mathbf{y}}_{n-1}^\eta(y; \mathbf{t}, \tilde{\mathbf{w}})) + \sum_{j=2}^{l^*} \mathbb{1}\{n = t_j\} \tilde{w}_j \right) \right) \quad (\text{B.27})$$

with initial condition  $\tilde{\mathbf{y}}_0^\eta(y; \mathbf{t}, \tilde{\mathbf{w}}) = \varphi_L(y + \varphi_b(\tilde{w}_1))$ .

**Corollary B.5.** *Given any  $\epsilon > 0$ , the following claim holds for all sufficiently small  $\eta > 0$ : for any  $x, y \in \Omega$ , and sequence of integers  $\mathbf{t} = (t_j)_{j=1}^{l^*}$  and any two sequence of real numbers  $\mathbf{w} = (w_j)_{j=1}^{l^*}$ ,  $\tilde{\mathbf{w}} = (\tilde{w}_j)_{j=1}^{l^*}$  such that*

- $|x - y| < \epsilon$ ;
- $t_1 = 0$ , and  $t_j - t_{j-1} \leq 2\bar{t}/\eta$  for all  $j \geq 1$  where  $\bar{t}$  is the constant in (B.21);
- $|w_j - \tilde{w}_j| < \epsilon$  for all  $j \geq 1$ ;

then we have

$$\sup_{t \in [0, t_{l^*}]} |\tilde{\mathbf{x}}^\eta(t, x; \mathbf{t}, \mathbf{w}) - \tilde{\mathbf{y}}_{[t]}^\eta(y; \mathbf{t}, \tilde{\mathbf{w}})| \leq \bar{\rho}\epsilon$$

where the constant  $\bar{\rho} = (3 \exp(\eta M \bar{t}) + 3)^{l^*}$ .

*Proof.* Throughout this proof, fix some  $\eta \in (0, \epsilon/2M)$ . We will show that for any  $\eta$  in the range the claim would hold.

First, on interval  $[0, t_2]$ , from Lemma B.4, one can see that (since  $2M\eta < \epsilon$ )

$$\sup_{t \in [0, t_2]} |\tilde{\mathbf{x}}^\eta(t, x; \mathbf{t}, \mathbf{w}) - \tilde{\mathbf{y}}_{[t]}^\eta(y; \mathbf{t}, \tilde{\mathbf{w}})| \leq 3 \exp(\eta M \bar{t}) \cdot \epsilon.$$

Then at  $t = t_2$ , by considering the difference between  $w_2$  and  $\tilde{w}_2$ , and the possible change due to one more gradient descent step (which is bounded by  $\eta M < \epsilon$ ), we have

$$\sup_{t \in [0, t_2]} |\tilde{\mathbf{x}}^\eta(t, x; \mathbf{t}, \mathbf{w}) - \tilde{\mathbf{y}}_{[t]}^\eta(y; \mathbf{t}, \tilde{\mathbf{w}})| \leq (3 \exp(\eta M \bar{t}) + 2) \cdot \epsilon.$$

Now we proceed inductively. For any  $j = 2, 3, \dots, l^* - 1$ , assume that

$$\sup_{t \in [0, t_j]} |\tilde{\mathbf{x}}^\eta(t, x; \mathbf{t}, \mathbf{w}) - \tilde{\mathbf{y}}_{[t]}^\eta(y; \mathbf{t}, \tilde{\mathbf{w}})| \leq (3 \exp(\eta M \bar{t}) + 3)^{j-1} \cdot \epsilon.$$

Then by focusing on interval  $[t_j, t_{j+1}]$  and using Lemma B.4 again, one can show that

$$\begin{aligned} \sup_{t \in [t_j, t_{j+1}]} |\tilde{\mathbf{x}}^\eta(t, x; \mathbf{t}, \mathbf{w}) - \tilde{\mathbf{y}}_{[t]}^\eta(y; \mathbf{t}, \tilde{\mathbf{w}})| &\leq 2\epsilon + ((3 \exp(\eta M \bar{t}) + 3)^{j-1} + 1) \exp(\eta M \bar{t}) \epsilon \\ &\leq (3 \exp(\eta M \bar{t}) + 3)^j \cdot \epsilon. \end{aligned}$$

This concludes the proof.  $\square$

In the next few results, we show that the same can be said for gradient descent iterates  $\tilde{\mathbf{y}}_n$  and the SGD iterates  $X_n$ . Specifically, our first goal is to show that before any *large* jump (see the definition in B.8), it is unlikely that the gradient descent process  $\mathbf{y}_n^\eta$  would deviate too far from  $X_n^\eta$ . Define the event

$$A(n, \eta, \epsilon, \delta) = \left\{ \max_{k=1, 2, \dots, n \wedge (T_1^\eta(\delta) - 1)} \eta |Z_1 + \dots + Z_k| \leq \epsilon \right\} \quad (\text{B.28})$$

and recall that arrival times  $T_j^\eta(\delta)$  are defined in (B.9).

As a building block, we first study the case when the starting point  $x$  is close to the reflection boundary  $-L$ . The takeaway from the next result is that the reflection operator hardly comes into play, since the SGD iterates would most likely quickly move to somewhere far enough from  $\pm L$ ; besides, throughout this procedure the SGD iterates would most likely stay pretty close to the corresponding deterministic gradient descent process.

**Lemma B.6.** *Given  $\epsilon \in (0, \bar{\epsilon}/9)$ , it holds for any sufficiently small  $\epsilon, \delta, \eta > 0$  that, if  $x \in [-L, -L + \bar{\epsilon}]$  and  $\rho_0(|x - y| + 9\epsilon) < \bar{\epsilon}$ , then on event  $A(n, \eta, \epsilon, \delta)$  we have*

$$|X_k^\eta(x) - \mathbf{y}_k^\eta(y)| \leq \rho_0 \cdot (|x - y| + 9\epsilon) \quad \forall k = 1, 2, \dots, n \wedge (T_1^\eta(\delta) - 1) \wedge \tilde{T}_{\text{escape}}^\eta(x)$$

where  $\tilde{T}_{\text{escape}}^\eta(x) \triangleq \min\{n \geq 0 : X_n^\eta(x) > -L + \bar{\epsilon}\}$  and  $\rho_0 \triangleq \exp\left(\frac{2M\bar{\epsilon}}{0.9c_-^L}\right)$  is a constant that does not vary with our choice of  $\epsilon, \delta, \eta$ .

*Proof.* For any  $k < T_1^\eta(\delta)$ , we know that  $Z_k = Z_k^{\leq \delta}$  (thus  $\eta|Z_k| < \delta$ ). Also, recall that  $|f'(x)| \leq M$  for any  $x \in \Omega_i$ . Therefore, as long as  $\eta$  and  $\delta$  are small enough, we will have that

$$|\eta(-f'(X_n^\eta(x)) + Z_k^{\leq \delta})| \leq b \quad (\text{B.29})$$

so the gradient clipping operator in (B.4) has no effect when  $k < T_1^\eta(\delta)$ , and in fact the only possible time for the gradient clipping trick to work is at  $T_j^\eta(\delta)$ . Therefore, we can safely rewrite the SGD update as

$$X_k^\eta(x) = X_{k-1}^\eta(x) - \eta f'(X_{k-1}^\eta(x)) + \eta Z_k + R_k \quad \forall k < T_1^\eta(\delta)$$

where each  $R_k \geq 0$  and it represents the push caused by reflection at  $-L$ .

First, choose  $\epsilon$  small enough so that  $9\epsilon < \bar{\epsilon}$ . Next, based on (B.17) we have the following lower bound:

$$X_k^\eta(x) \geq x + 0.9c_-^L \eta k - \epsilon \quad \forall k < T_1^\eta(\delta).$$

Let  $\tilde{t}_0(x, \epsilon) \triangleq \min\{n \geq 0 : X_n^\eta(x) \geq -L + 2\epsilon\}$ . Due to the inequality above, we know that

$$\tilde{t}_0(x, \epsilon) \leq \frac{3\epsilon}{0.9c_-^L \eta}. \quad (\text{B.30})$$

One the other hand, given the current choice of  $\epsilon$ , if we choose  $\eta$  and  $\delta$  small enough, then using the same argument leading to (B.29), we will have

$$X_{\tilde{t}_0(x, \epsilon)}^\eta(x) \leq -L + 2.1\epsilon \leq x + 1.1c_-^L \eta k + 2.1\epsilon$$

if  $\tilde{t}_0(x, \epsilon) \geq 1$  (namely  $x < -L + 2\epsilon$ ).

Let us inspect the two scenarios separately. First, assume  $\tilde{t}_0(x, \epsilon) \geq 1$ . For the deterministic gradient descent process  $\mathbf{y}_n^\eta(y)$ , we have the following bounds:

$$y + 0.9c_-^L \eta k \leq \mathbf{y}_k^\eta(y) \leq y + 1.1c_-^L \eta k \quad \forall k \leq \tilde{t}_0(x, \epsilon) \wedge (T_1^\eta(\delta) - 1).$$

This gives us

$$|X_k^\eta(x) - \mathbf{y}_k^\eta(y)| \leq |x - y| + 0.2c_-^L \eta k + 2.1\epsilon \quad \forall k \leq \tilde{t}_0(x, \epsilon) \wedge (T_1^\eta(\delta) - 1).$$

At time  $k = \tilde{t}_0(x, \epsilon)$ , due to previous bound on  $\tilde{t}_0(x, \epsilon)$ , we know that  $|X_{\tilde{t}_0(x, \epsilon)}^\eta(x) - \mathbf{y}_{\tilde{t}_0(x, \epsilon)}^\eta(y)| \leq |x - y| + 7\epsilon$ . If  $n \wedge (T_1^\eta(\delta) - 1) \leq \tilde{t}_0(x, \epsilon)$  then we have already shown the desired claim. Otherwise, starting from time  $\tilde{t}_0(x, \epsilon)$ , due to the definition of event  $A(n, \eta, \epsilon, \delta)$  in (B.28), we know that the SGD iterates  $X_n^\eta(x)$  will not touch the boundary  $-L$  afterwards. Therefore, by directly applying Lemma B.3, and notice that  $|f''(x)| \leq M$  for any  $x \in [-L, L]$  and, we have

$$|X_k^\eta(x) - \mathbf{y}_k^\eta(y)| \leq (|x - y| + 9\epsilon) \cdot \exp\left(\frac{2M\bar{\epsilon}}{0.9c_-^L}\right) \quad \forall k \leq n \wedge (T_1^\eta(\delta) - 1) \wedge \tilde{T}_{\text{escape}}^\eta(x).$$



Indeed, it suffices to use Lemma B.3 for the next  $\lceil 2\bar{\epsilon}/(0.9\eta c_-^L) \rceil$  steps to show that  $|X_k^\eta(x) - \mathbf{y}_k^\eta(y)|$  should be smaller than  $\epsilon$  for the next  $\lceil 2\bar{\epsilon}/(0.9\eta c_-^L) \rceil$  steps, while  $\mathbf{y}_k^\eta(y)$  will reach some where in  $(-L + 2\bar{\epsilon}, -L + 3\bar{\epsilon})$  within  $\lceil 2\bar{\epsilon}/(0.9\eta c_-^L) \rceil$  steps so we must have

$$n \wedge (T_1^\eta(\delta) - 1) \wedge \tilde{T}_{\text{escape}}^\eta(x) \wedge \tilde{t}_0(x, \epsilon) - \tilde{t}_0(x, \epsilon) \leq 2\bar{\epsilon}/(0.9\eta c_-^L) \quad (\text{B.31})$$

Lastly, in the case  $\tilde{t}_0(x, \epsilon) = 0$  (which means  $x \geq -L + \epsilon$ ), we can use Lemma B.3 directly as we did above and establish the same bound. This concludes the proof.  $\square$

Obviously, a similar result can be shown if  $x$  is in the rightmost attraction field  $(s_{n_{\min}-1}, L]$  and the approach is identical. We omit the details here. In the next Lemma, we consider the scenario where the starting point  $x$  is far enough from the boundaries.

**Lemma B.7.** *Given any  $\epsilon > 0$ , the following holds for all sufficiently small  $\eta > 0$ : for any  $x, y \in \Omega$  and positive integer  $n$  such that  $|x-L| > 2\epsilon$ ,  $|x+L| > 2\epsilon$ ,  $|x-s_-| > 2\epsilon$ ,  $|x-s_+| > 2\epsilon$  and  $|x-y| < \frac{\epsilon}{2\exp(\eta Mn)}$ , on event*

$$A(n, \eta, \frac{\epsilon}{2\exp(\eta Mn)}, \delta) \cap \left\{ |\mathbf{y}_j^\eta(y)| \in \Omega, |X_j^\eta(x)| \in \Omega \quad \forall j = 1, 2, \dots, n \wedge (T_1^\eta(\delta) - 1) \right\}$$

we have

$$|\tilde{X}_m^\eta(y) - X_m^\eta(x)| \leq \epsilon \quad \forall m = 1, 2, \dots, n \wedge (T_1^\eta(\delta) - 1).$$

*Proof.* For sufficiently small  $\eta$ , we will have that the (deterministic) gradient descent iterates  $|\mathbf{y}_n^\eta|$  is monotonically decreasing in  $n$ , which ensures that  $\mathbf{y}_n^\eta$  always stays in the range that are at least  $\epsilon$ -away from  $\pm L$  or  $s_-, s_+$ . We now show that the claim holds for any such  $\eta$ .

On event  $\left\{ |\mathbf{y}_j^\eta(y)| \in \Omega, |X_j^\eta(x)| \in \Omega \quad \forall j = 1, 2, \dots, n \wedge (T_1^\eta(\delta) - 1) \right\}$ , we are able to apply Lemma B.3 inductively for any  $m \in [n]$  and obtain that

$$|\mathbf{y}_j^\eta(y) - X_j^\eta(x)| \leq (|x - y| + \frac{\epsilon}{2\exp(\eta Mn)}) \exp(\eta Mj) < \epsilon \quad \forall j = 1, 2, \dots, m$$

and conclude the proof. The reason to apply the Lemma inductively for  $m = 1, 2, \dots, n$ , instead of directly at step  $n$ , is to ensure that SGD iterates  $X_n^\eta$  would not hit the boundary  $\pm L$  (so the reflection operator would not come into play on the time interval we are currently interested in), thus ensuring that Lemma B.3 is applicable.  $\square$

Similar to the extension from Lemma B.4 to Corollary B.5, we can extend Lemma B.7 to show that, if we consider the a gradient descent process that is only perturbed by large noises, then it should stay pretty close to the SGD iterates  $X_n^\eta$ . To be specific, let

$$Y_0^\eta(x) = x \quad (\text{B.32})$$

$$Y_n^\eta(x) = \varphi_L \left( Y_{n-1}^\eta(x) - \varphi_b(-\eta f'(Y_{n-1}^\eta(x)) + \sum_{j \geq 1} \mathbb{1}\{n = T_j^\eta(\delta)\} \eta Z_n) \right). \quad (\text{B.33})$$

be a gradient descent process (with gradient clipping at threshold  $b$ ) that is only shocked by large noises in  $(Z_n)_{n \geq 1}$ . The next corollary can be shown by an approach that is identical to Corollary B.5 (namely, inductively repeating Lemma B.7 at each jump time) so we omit the details here.

**Corollary B.8.** *Given any  $\epsilon > 0$ , the following holds for any sufficiently small  $\eta > 0$ : For any  $|x| < 2\epsilon$ , on event  $A_0(\epsilon, \eta, \delta) \cap B_0(\epsilon, \eta, \delta)$ , we have*

$$|Y_n^\eta(x) - X_n^\eta(x)| < \tilde{\rho}\epsilon \quad \forall n = 1, 2, \dots, T_{l^*}^\eta(\delta)$$

where

$$A_0(\epsilon, \eta, \delta) \triangleq \left\{ \forall i = 1, \dots, l^*, \max_{j=T_{i-1}^\eta(\delta)+1, \dots, T_i^\eta(\delta)-1} \eta |Z_{T_{i-1}^\eta(\delta)+1} + \dots + Z_j| \leq \frac{\epsilon}{2 \exp(2\bar{t}M)} \right\};$$

$$B_0(\epsilon, \eta, \delta) \triangleq \left\{ \forall j = 2, \dots, l^*, T_j^\eta(\delta) - T_{j-1}^\eta(\delta) \leq 2\bar{t}/\eta \right\}$$

and  $\tilde{\rho} \in (0, \infty)$  is a constant that does not vary with  $\eta, \delta, \epsilon$ .

The next two results shows that the type of events  $A(n, \eta, \epsilon, \delta)$  defined in (B.28) is indeed very likely to occur, especially for small  $\epsilon$ . For clarity of the presentation, we introduce the following definitions that are slightly more general than the *small* and *large* jumps defined in (B.7)(B.8) (for any  $c > 0$ )

$$Z_n^{\leq c} \triangleq Z_n \mathbb{1}\{|Z_n| \leq c\},$$

$$Z_n^{> c} \triangleq Z_n \mathbb{1}\{|Z_n| > c\}.$$

**Lemma B.9.** Define functions  $u(\eta) = \delta/\eta^{1-\Delta}$ ,  $v(\eta) = \epsilon\eta^{\tilde{\Delta}}$  with  $\epsilon, \delta > 0$ . If real numbers  $\Delta, \tilde{\Delta}, \beta, \epsilon, \delta$  and positive integers  $j, N$  are such that the following conditions hold:

$$\Delta \in [0, (1 - \frac{1}{\alpha}) \wedge \frac{1}{2}), \quad (\text{B.34})$$

$$\beta \in (1, (2 - 2\Delta) \wedge \alpha(1 - \Delta)), \quad (\text{B.35})$$

$$\tilde{\Delta} \in [0, \frac{\Delta}{2}], \quad \tilde{\Delta} < \alpha(1 - \Delta) - \beta, \quad (\text{B.36})$$

$$N < (\alpha(1 - \Delta) - \beta)j, \quad (\text{B.37})$$

$$v(\eta) - j\eta u(\eta) \geq v(\eta)/2 \quad \text{for all } \eta > 0 \text{ sufficiently small}, \quad (\text{B.38})$$

then

$$\mathbb{P}\left(\max_{k=1,2,\dots,\lceil 1/\eta^\beta \rceil} \eta |Z_1^{\leq u(\eta)} + \dots + Z_k^{\leq u(\eta)}| > 3v(\eta)\right) = o(\eta^N)$$

as  $\eta \downarrow 0$ .

*Proof.* From the stated range of the parameters, we know that

$$\alpha(1 - \Delta) > \beta,$$

$$(\alpha(1 - \Delta) - \beta)j > N,$$

so we are able to find  $\gamma \in (0, 1)$  small enough such that

$$\alpha(1 - \Delta)(1 - 2\gamma) > \beta, \quad (\text{B.39})$$

$$(\alpha(1 - \Delta)(1 - 2\gamma) - \beta)j > N. \quad (\text{B.40})$$

Fix such  $\gamma \in (0, 1)$  for the rest of the proof, and let  $n(\eta) \triangleq \lceil (1/\eta)^\beta \rceil$ ,  $I \triangleq \#\{i \in [n(\eta)] : |Z_i^{\leq u(\eta)}| > u(\eta)^{1-\gamma}\}$ . Then

$$\begin{aligned} & \mathbb{P}\left(|Z_1^{\leq u(\eta)} + \dots + Z_{n(\eta)}^{\leq u(\eta)}| > v(\eta)\right) \\ &= \sum_{i=0}^{j-1} \underbrace{\mathbb{P}\left(|Z_1^{\leq u(\eta)} + \dots + Z_{n(\eta)}^{\leq u(\eta)}| > v(\eta), I = i\right)}_{\triangleq \text{(I)}} + \underbrace{\mathbb{P}\left(|Z_1^{\leq u(\eta)} + \dots + Z_{n(\eta)}^{\leq u(\eta)}| > v(\eta), I \geq j\right)}_{\triangleq \text{(II)}} \end{aligned}$$

Note that since  $|Z_i^{\leq u(\eta)}| < u(\eta)$ ,

$$\begin{aligned}
(\text{I}) &\leq \binom{n(\eta)}{i} \cdot \mathbb{P}\left(|Z_1^{\leq u(\eta)} + \dots + Z_{n(\eta)-i}^{\leq u(\eta)}| > \frac{v(\eta) - i\eta u(\eta)}{\eta}, |Z_i^{\leq u(\eta)}| \leq u(\eta)^{1-\gamma} \forall i \in [n(\eta) - i]\right) \\
&\leq n(\eta)^i \cdot \mathbb{P}\left(|Z_1^{\leq u(\eta)^{1-\gamma}} + \dots + Z_{n(\eta)-i}^{\leq u(\eta)^{1-\gamma}}| > \frac{v(\eta) - i\eta u(\eta)}{\eta}\right) \\
&\leq n(\eta)^i \cdot \mathbb{P}\left(|Z_1^{\leq u(\eta)^{1-\gamma}} + \dots + Z_{n(\eta)-i}^{\leq u(\eta)^{1-\gamma}}| > \frac{v(\eta)}{2\eta}\right)
\end{aligned} \tag{B.41}$$

where the last inequality follows from (B.38). First, since  $\mathbb{E}Z_1 = 0$ , we have

$$\begin{aligned}
|\mathbb{E}Z_1^{\leq u(\eta)^{1-\gamma}}| &= |\mathbb{E}Z_1^{> u(\eta)^{1-\gamma}}| \\
&= \int_{u(\eta)^{1-\gamma}}^{\infty} \mathbb{P}(|Z_1| > x) dx \in \mathcal{RV}_{(\alpha-1)(1-\gamma)(1-\Delta)}(\eta).
\end{aligned}$$

Therefore, for all  $\eta > 0$  that are sufficiently small,

$$\begin{aligned}
&|\mathbb{E}Z_1^{\leq u(\eta)^{1-\gamma}} + \dots + \mathbb{E}Z_{n(\eta)-i}^{\leq u(\eta)^{1-\gamma}}| \\
&\leq n(\eta) \cdot \eta^{(\alpha-1)(1-\Delta)(1-2\gamma)} \leq 2\eta^{(\alpha-1)(1-\Delta)(1-2\gamma)-\beta} \\
&\leq (1/\eta)^{(1-\Delta)(1-2\gamma)} \quad \text{due to (B.39)} \\
&\leq \frac{v(\eta)}{4\eta} \quad \text{due to } \tilde{\Delta}/2 \leq \Delta \text{ in (B.36) and } 1-\gamma < 1.
\end{aligned}$$

If we let  $Y_n = Z_n^{\leq u(\eta)^{1-\gamma}} - \mathbb{E}Z_n^{\leq u(\eta)^{1-\gamma}}$  and plug the bound above back into (B.41), then (for all  $\eta > 0$  that are sufficiently small)

$$\begin{aligned}
(\text{I}) &\leq n(\eta)^i \cdot \mathbb{P}(|Y_1 + \dots + Y_{n(\eta)-i}| > \frac{v(\eta)}{4\eta}) \\
&\leq n(\eta)^i \exp\left(-\frac{\frac{\epsilon^2}{16} \cdot 1/\eta^{2-2\tilde{\Delta}}}{2(n(\eta)-i)\mathbb{E}|Y_1|^2 + \frac{2}{3}\delta^{1-\gamma} \cdot (1/\eta)^{(1-\Delta)(1-\gamma)} \cdot \frac{\epsilon}{4}/\eta^{1-\tilde{\Delta}}}\right)
\end{aligned} \tag{B.42}$$

where the last inequality is obtained from Bernstein's inequality. Note that from Karamata's theorem,

$$\begin{aligned}
\mathbb{E}|Y_1|^2 &= \text{var}(Z_1^{\leq u(\eta)^{1-\gamma}}) \leq \mathbb{E}|Z_1^{\leq u(\eta)^{1-\gamma}}|^2 \\
&\leq \int_0^{u(\eta)^{1-\gamma}} 2x\mathbb{P}(|Z_1| > x) dx \in \mathcal{RV}_{-(1-\Delta)(1-\gamma)(2-\alpha)}(\eta).
\end{aligned}$$

Now note that

- In case that  $\alpha < 2$ , for all  $\eta > 0$  that are sufficiently small, we have (using (B.36))

$$\begin{aligned}
2(n(\eta) - i)\mathbb{E}|Y_1|^2 &\leq (1/\eta)^{\beta+(2-\alpha)(1-\Delta)} < (1/\eta)^{2(1-\Delta)} \\
&\Rightarrow \frac{1/\eta^{2-2\tilde{\Delta}}}{2(n(\eta) - i)\mathbb{E}|Y_1|^2} \geq 1/\eta^{\tilde{\Delta}};
\end{aligned}$$

- In case that  $\alpha \geq 2$ , for all  $\eta > 0$  that are sufficiently small,

$$2(n(\eta) - i)\mathbb{E}|Y_1|^2 < 1/\eta^{\beta+\frac{\Delta}{2}}$$

and we know that  $\beta + \frac{\Delta}{2} < 2 - 2\tilde{\Delta}$  due to  $2 - \beta > 2\Delta$  and  $2\tilde{\Delta} \leq \Delta$  (see (B.34)-(B.36));

- Since  $\gamma > 0$  and  $2\tilde{\Delta} \leq \Delta$ , we know that

$$(1 - \Delta)(1 - \gamma) + (1 - \tilde{\Delta}) < 2 - 2\tilde{\Delta}.$$

Therefore, it is easy to see that the R.H.S. of (B.42) decays at a geometric rate as  $\eta$  tends to zero, hence  $o(\eta^N)$ . On the other hand,

$$(II) \leq \mathbb{P}(I \geq j) \leq \binom{n(\eta)}{j} \cdot \mathbb{P}\left(|Z_i^{\leq u(\eta)}| > u(\eta)^{1-\gamma} \forall i = 1, \dots, j\right) \leq n(\eta)^j \cdot \mathbb{P}\left(|Z_1^{\leq u(\eta)}| > u(\eta)^{1-\gamma}\right)^j,$$

which is regularly varying w.r.t.  $\eta$  with index  $(\alpha(1 - \gamma)(1 - \Delta) - \beta)j$ . Therefore, for all  $\eta > 0$  sufficiently small,

$$(II) \leq \eta^{(\alpha(1-2\gamma)(1-\Delta)-\beta)j} < \eta^N \quad \text{due to (B.40).}$$

Collecting results above, we have established that

$$\mathbb{P}\left(\eta|Z_1^{\leq u(\eta)} + \dots + Z_{n(\eta)}^{\leq u(\eta)}| > v(\eta)\right) = o(\eta^N).$$

The conclusion of the lemma now follows from Etemadi's theorem.  $\square$

Now consider the following setting. Let us fix some positive integer  $N$  and  $\beta \in (1, 2 \wedge \alpha)$ . Then we can find some positive integer  $j$  such that  $(\alpha - \beta)j > N$ . Meanwhile, given any  $\epsilon > 0$ , we will have  $\epsilon - j\delta \geq \epsilon/2$  for all  $\delta > 0$  sufficiently small. Therefore, by applying Lemma B.9 with  $\Delta = \tilde{\Delta} = 0$  (hence  $u(\eta) = \delta/\eta$ ,  $v(\eta) = \epsilon$ ) and  $\beta, j, N, \epsilon, \delta$  as described here, we immediately get the following result.

**Lemma B.10.** *Given any  $\beta \in (1, \alpha \wedge 2)$ ,  $\epsilon > 0$ , and  $N > 0$ , the following holds for any sufficiently small  $\delta > 0$ :*

$$\mathbb{P}\left(\max_{j=1,2,\dots,\lceil(1/\eta)^\beta\rceil} \eta|Z_1^{\leq \delta/\eta} + \dots + Z_j^{\leq \delta/\eta}| > \epsilon\right) = o(\eta^N)$$

as  $\eta \downarrow 0$ .

Using results and arguments above, we are able to illustrate the typical behavior of the SGD iterates  $X_n^\eta$  in the following two scenarios. First, we show that, when starting from most parts in the attraction field  $\Omega$ , the SGD iterates  $X_n^\eta$  will most likely return to the neighborhood of the local minimum within a short period of time without exiting  $\Omega$ . Given that there are only finitely many attraction fields on  $f$ , it is easy to see that the key technical tool Lemma 7 follows immediately from the next result.

**Lemma B.11.** *For sufficiently small  $\epsilon > 0$ , the following claim holds:*

$$\lim_{\eta \downarrow 0} \sup_{x \in \Omega: |x-s_-| \wedge |x-s_+| > \epsilon} \mathbb{P}_x\left(X_n^\eta \in \Omega \forall n \leq T_{\text{return}}(\eta, \epsilon), \text{ and } T_{\text{return}}(\eta, \epsilon) \leq \rho(\epsilon)/\eta\right) = 1$$

where the stopping time involved is defined as

$$T_{\text{return}}(\eta, \epsilon) \triangleq \min\{n \geq 0 : X_n^\eta(x) \in [-2\epsilon, 2\epsilon]\}$$

the function  $\hat{t}(\epsilon)$  is defined in (B.24), and the function  $\rho(\cdot)$  is defined as  $\rho(\epsilon) = \frac{3\bar{\epsilon}}{0.9c_-^L \wedge c_+^L} + 2\hat{t}(\epsilon)$

*Proof.* Throughout, we only consider  $\epsilon$  small enough so that Lemma B.6 could hold. Also, fix some  $N > 0, \Delta_\alpha \in (0, \alpha - 1)$  and  $\beta \in (1, \alpha)$ . Let  $\sigma(x, \eta) \triangleq \min\{n \geq 0 : X_n^\eta \notin \Omega\}$ .

Without loss of generality, we assume  $\Omega = [-L, s_+)$  and  $x < 0$  (so reflection at  $-L$ ) is a possibility. Any other case can be addressed similarly as shown below.

From Lemma B.2 and the regular varying nature of  $H(\cdot)$ , we have, for any  $\epsilon, \delta > 0$ ,

$$\mathbb{P}(T_1^\eta(\delta) \leq \rho(\epsilon)/\eta) \leq \eta^{\alpha-1-\Delta_\alpha} \quad (\text{B.43})$$

for any sufficiently small  $\eta$ .

Let  $\tilde{T}_{\text{escape}}^\eta(x)$  be the stopping time defined in Lemma B.6. From (B.30), (B.31), (B.43) and Lemma B.10, we know that

$$\begin{aligned} \sup_{x \in [-L, -L+\bar{\epsilon}]} \mathbb{P}\left(\tilde{T}_{\text{escape}}^\eta(x) < \sigma(x, \eta), \tilde{T}_{\text{escape}}^\eta(x) \leq \frac{3\bar{\epsilon}}{0.9c_-^L \eta} \text{ and } X_{\tilde{T}_{\text{escape}}^\eta}^\eta(x) \in [-L + \bar{\epsilon}, -L + 2\bar{\epsilon}]\right) \\ \geq 1 - \eta^N - \eta^{\alpha-1-\Delta_\alpha} \end{aligned} \quad (\text{B.44})$$

for all sufficiently small  $\eta$ .

Next, we focus on  $x \in \Omega$  such that  $|x - s_-| \wedge |x - s_+| > \epsilon$  and  $x \geq -L + \bar{\epsilon}$ . We start by considering the time it took for the (deterministic) gradient descent process  $\mathbf{y}_n^\eta(x)$  to return to  $[-1.5\epsilon, 1.5\epsilon]$ . From the definition of  $\hat{t}(\epsilon)$  in (B.24) and Lemma B.4, we know that for  $\eta$  small enough such that  $\eta \exp(2M\hat{t}(\epsilon)) < 0.5\epsilon$ , we have

$$\min\{n \geq 0 : \mathbf{y}_n^\eta(x) \in [-1.5\epsilon, 1.5\epsilon]\} \leq 2\hat{t}(\epsilon)/\eta.$$

Now consider event  $A(\lceil(1/\eta)^\beta\rceil, \eta, \frac{\epsilon}{4\exp(2M\hat{t}(\epsilon))}, \delta)$  (see definition in (B.28)). From Lemma B.10, we know that for any sufficiently small  $\delta$ , we have

$$\mathbb{P}\left((A(\lceil(1/\eta)^\beta\rceil, \eta, \frac{\epsilon}{4\exp(2M\hat{t}(\epsilon))}, \delta))^c\right) = o(\eta^N). \quad (\text{B.45})$$

Combining this result with (B.43)(B.45) and Lemma B.7, we get

$$\sup_{x \in \Omega: |x-s_-| \wedge |x-s_+| > \epsilon, x \geq -L+\bar{\epsilon}} \mathbb{P}_x\left(T_{\text{return}}(\eta, \epsilon) < \sigma(x, \eta), T_{\text{return}}(\eta, \epsilon) \leq 2\hat{t}(\epsilon)/\eta\right) \geq 1 - \eta^N - \eta^{\alpha-1-\Delta_\alpha} \quad (\text{B.46})$$

for any sufficiently small  $\eta$ . To conclude the proof, we only to combine strong Markov property (at  $\tilde{T}_{\text{escape}}^\eta$ ) with bounds in (B.44)(B.46).  $\square$

In the next result, we show that, once entering a  $\epsilon$ -small neighborhood of the local minimum, the SGD iterates will most likely stay there until the next large jump.

**Lemma B.12.** *Given  $N_0 > 0$ , the following claim holds for any  $\epsilon, \delta > 0$  that are sufficiently small:*

$$\sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}\left(\exists n < T_1^\eta(\delta) \text{ s.t. } |X_n^\eta(x)| > 3\epsilon\right) = o(\eta^{N_0})$$

as  $\eta \downarrow 0$ .

*Proof.* Fix  $\epsilon$  small enough such that  $3\epsilon < \epsilon_0$  (see Assumption 1 for the constant  $\epsilon_0$ ). Also, fix some  $\Delta_\alpha \in (0, 1), \beta \in (1, \alpha), N > \alpha + \Delta_\alpha - \beta + N_0$ . Due to Lemma B.10, for any  $\delta$  sufficiently small, we will have

$$\mathbb{P}\left(\max_{j=1,2,\dots,\lceil(1/\eta)^\beta\rceil} \eta|Z_1^{\leq,\delta} + \dots + Z_j^{\leq,\delta}| > \frac{\epsilon}{\exp(2M)}\right) = o(\eta^N). \quad (\text{B.47})$$

Fix such  $\delta > 0$ . We now show that the desired claim is true for the chosen  $\epsilon, \delta$ .

First of all, from Lemma B.1, we know the existence of some  $\theta > 0$  such that

$$\mathbb{P}(T_1^\eta(\delta) > 1/\eta^{\alpha+\Delta\alpha}) = o(\exp(-1/\eta^\theta)). \quad (\text{B.48})$$

Next, let us zoom in on the first  $\lceil (1/\eta)^\beta \rceil$  SGD iterates. For any  $\eta$  small enough, we will have  $\mathbf{y}_n^\eta(x) \in [-2\epsilon, 2\epsilon]$  for any  $n \geq 1$  and  $\mathbf{y}_{\lceil (1/\eta)^\beta \rceil}^\eta(x) \in [-\epsilon, \epsilon]$  given  $x \in [-2\epsilon, 2\epsilon]$ . From now on we only consider such  $\eta$ . Due to Lemma B.7, we know that on event  $\left\{ \max_{j=1,2,\dots,\lceil (1/\eta)^\beta \rceil} \eta |Z_1^{\leq,\delta} + \dots + Z_j^{\leq,\delta}| > \frac{\epsilon}{\exp(2M)} \right\}$ , we have

$$|X_n^\eta(x)| \leq 3\epsilon \quad \forall n \leq \lceil 1/\eta^\beta \rceil \wedge (T_1^\eta(\delta) - 1)$$

and on event  $\left\{ \max_{j=1,2,\dots,\lceil (1/\eta)^\beta \rceil} \eta |Z_1^{\leq,\delta} + \dots + Z_j^{\leq,\delta}| > \frac{\epsilon}{\exp(2M)} \right\} \cap \{T_1^\eta(\delta) > \lceil (1/\eta)^\beta \rceil\}$ , we have  $X_{T_1^\eta(\delta)}^\eta(x) \in [-2\epsilon, 2\epsilon]$ . Now by repeating the same argument inductively for  $\lceil 1/\eta^{\alpha+\Delta\alpha-\beta} \rceil$  times, we can show that on event

$$\left\{ \forall i = 1, 2, \dots, \lceil 1/\eta^{\alpha+\Delta\alpha-\beta} \rceil, \max_{j=1,2,\dots,\lceil (1/\eta)^\beta \rceil} \eta |Z_{i\lceil (1/\eta)^\beta \rceil+1}^{\leq,\delta} + \dots + Z_{i\lceil (1/\eta)^\beta \rceil+j}^{\leq,\delta}| > \frac{\epsilon}{\exp(2M)} \right\}, \quad (\text{B.49})$$

we have  $|X_n^\eta(x)| \leq 3\epsilon \quad \forall n \leq 1/\eta^{\alpha+\Delta\alpha} \wedge (T_1^\eta(\delta) - 1)$ . To conclude the proof, we only need to combine this fact with (B.47).  $\square$

The following result is an immediate product from the geometric tail bound in (B.48) and the inductive argument in (B.49).

**Corollary B.13.** *Given  $N > 0, \tilde{\epsilon} > 0$ , the following claim holds for all sufficiently small  $\delta > 0$ :*

$$\mathbb{P}\left(\max_{j=1,2,\dots,T_1^\eta(\delta)-1} \eta |Z_1 + \dots + Z_j| > \tilde{\epsilon}\right) = o(\eta^N). \quad (\text{B.50})$$

We introduce a few concepts that will be crucial in the analysis below. Recall the definition of perturbed ODE  $\tilde{\mathbf{x}}^\eta$  in (B.18)-(B.20) (note that we will drop the notational dependency on learning rate  $\eta$  when we choose  $\eta = 1$ ). Consider the definition of the following two mappings from where  $\mathbf{w} = (w_1, \dots, w_{l^*})$  is a sequence of real numbers and  $\mathbf{t} = (t_1, t_2, \dots, t_{l^*})$  with  $0 = t_1 < t_2 < t_3 < \dots$  as

$$h(\mathbf{w}, \mathbf{t}) = \tilde{\mathbf{x}}(t_{l^*}, 0; \mathbf{t}, \mathbf{w}).$$

Next, define sets (for any  $\epsilon \in (-\bar{\epsilon}, \bar{\epsilon})$ )

$$E(\epsilon) = \{(\mathbf{w}, \mathbf{t}) \subseteq \mathbb{R}^{l^*} \times \left(\mathbb{R}_+\right)^{l^*-1} : h(\mathbf{w}, \mathbf{t}) \notin [s_- - \epsilon, s_+ + \epsilon]\}. \quad (\text{B.51})$$

We add a few remarks about the two types of sets defined above.

- Intuitively speaking,  $E(\epsilon)$  contains all the perturbations (with times and sizes) that can send the ODE out of the current attraction field (allowing for some error with size  $\epsilon$ );
- From the definition of  $\bar{t}, \bar{\delta}$  in (B.21)(B.22) and Corollary B.5, one can easily see that for a fixed  $\epsilon \in (-\bar{\epsilon}, \bar{\epsilon})$ ,

$$(\mathbf{w}, \mathbf{t}) \in E(\epsilon) \Rightarrow |w_j| > \bar{\delta}, t_j - t_{j-1} \leq \bar{t} \quad \forall j;$$

- Lastly,  $E(\epsilon)$  are open sets due to  $f \in C^2$ .

Use  $\mathbf{Leb}_+$  to denote the Lebesgue measure restricted on  $[0, \infty)$ , and define (Borel) measure  $\nu_\alpha$  with density on  $\mathbb{R} \setminus \{0\}$ :

$$\nu_\alpha(dx) = \mathbb{1}\{x > 0\} \frac{\alpha p_+}{x^{\alpha+1}} + \mathbb{1}\{x < 0\} \frac{\alpha p_-}{|x|^{\alpha+1}}$$

where  $\alpha > 1$  is the regular variation index for the distribution of  $Z_1$  and  $p_-, p_+ \in (0, 1)$  are constants in Assumption 2. Now we can define a Borel measure  $\mu$  on  $\mathbb{R}^{l^*} \times \left(\mathbb{R}_+\right)^{l^*-1}$  as product measure

$$\mu = (\nu_\alpha)^{l^*} \times (\mathbf{Leb}_+)^{l^*-1}. \quad (\text{B.52})$$

Due to remarks above, one can see that for  $\epsilon \in (-\bar{\epsilon}, \bar{\epsilon})$ , we have  $\mu(E(\epsilon)) < \infty$ . We are now ready to analyze a specific type of noise  $Z_n$ .

**Definition B.1.** For any  $n \geq 1$  and any  $\epsilon \in (-\bar{\epsilon}, \bar{\epsilon}), \delta \in (0, b \wedge \bar{\delta}), \eta > 0$ , we say that the jump  $Z_n$  has  $(\epsilon, \delta, \eta)$ -**overflow** if

- $\eta|Z_n| > \delta$ ;
- In the set  $\{n+1, \dots, n+2\lceil l^* \bar{t}/\eta \rceil\}$ , there are at least  $(l^* - 1)$  elements (ordered as  $n < t_2 < t_3 < \dots < t_{l^*}$ ) such that  $\eta|Z_{t_i}| > \delta$  for any  $i = 2, \dots, l^*$ ;
- Let  $t_1 = n$  and  $t'_i = t_i - t_{i-1}$  for any  $i = 2, \dots, l^*$ ,  $w_i = \eta Z_{t_i}$  for any  $i = 1, \dots, l^*$ , for real sequence  $\mathbf{w} = (w_1, w_2, \dots, w_{l^*})$  and a sequence of positive number  $\mathbf{t} = (\eta(t_i - n))_{i=2}^{l^*}$ , we have

$$(\mathbf{w}, \mathbf{t}) \in E(\epsilon).$$

Moreover, if  $Z_n$  has  $(\epsilon, \delta, \eta)$ -overflow, then we call  $h(\mathbf{w}, \mathbf{t})$  as its  $(\epsilon, \delta, \eta)$ -**overflow endpoint**.

Due to the iid nature of  $(Z_j)_{j \geq 1}$ , let us consider an iid sequence  $(V_j)_{j \geq 0}$  where the sequence has the same law of  $Z_1$ . Note that for any fixed  $n \geq 1$ , the probability that  $Z_n$  has  $(\epsilon, \delta, \eta)$ -overflow is equal to the probability that  $V_0$  has  $(\epsilon, \delta, \eta)$ -overflow. More specifically, we know that  $\mathbb{P}(\eta|V_0| > \delta) = H(\delta/\eta)$ , and now we focus on conditional probability admitting the following form:

$$p(\epsilon, \delta, \eta) = \mathbb{P}\left(V_0 \text{ has } (\epsilon, \delta, \eta)\text{-overflow} \mid \eta|V_0| > \delta\right). \quad (\text{B.53})$$

For any open interval  $A = (a_1, a_2)$  such that  $A \cap [s_- + \bar{\epsilon}, s_+ - \bar{\epsilon}] = \emptyset$ , we also define

$$p(\epsilon, \delta, \eta; A) = \mathbb{P}\left(V_0 \text{ has } (\epsilon, \delta, \eta)\text{-overflow and the endpoint is in } A \mid \eta|V_0| > \delta\right). \quad (\text{B.54})$$

**Lemma B.14.** For any  $\epsilon \in (-\bar{\epsilon}, \bar{\epsilon}), \delta \in (0, b \wedge \bar{\delta})$ , and any open interval  $A = (a_1, a_2)$  such that  $|a_1| \wedge |a_2| > r - \bar{\epsilon}$  and  $|a_1| \neq L, |a_2| \neq L$ , we have

$$\lim_{\eta \downarrow 0} \frac{p(\epsilon, \delta, \eta; A)}{\delta^\alpha \left(\frac{H(1/\eta)}{\eta}\right)^{l^*-1}} = \mu(E(\epsilon) \cap h^{-1}(A))$$

where  $\mu$  is the measure defined in (B.52), and  $p(\cdot, \cdot, \cdot; A)$  is the conditional probability defined in (B.54).

*Proof.* Let us start by fixing some notations. Let  $T_1 = 0$ , and define stopping times  $T_j = \min\{n > T_{j-1} : \eta|V_n| > \delta\}$  and inter-arrival times  $T'_j = T_j - T_{j-1}$  for any  $j \geq 1$ , and large jump  $W_j = V_{T_j}$  for any  $j \geq 0$ . Note that: first of all, the pair  $(T'_i, W_i)$  is independent of  $(T'_j, W_j)$  whenever  $i \neq j$ ; besides,  $W_j$  and  $T'_j$  are independent for all  $j \geq 1$ .

Define the following sequence (of random elements)  $\mathbf{w} = (w_1, \dots, w_{l^*})$  and  $\mathbf{t} = (t_1, \dots, t_{l^*})$  by

$$w_j = \eta W_j, \quad t_j = \eta T_j.$$

If  $V_0$  has  $(\epsilon, \delta, \eta)$ -overflow, then the following two events must occur:

- $T'_j \leq 2\bar{t}/\eta$  for any  $j = 2, \dots, l^*$ ;
- $\eta|W_j| > \bar{\delta}$  for any  $j = 1, 2, \dots, l^*$ ;
- $(\mathbf{w}, \mathbf{t}) \in E(\epsilon)$

Therefore, for sufficiently small  $\eta$ , we now have

$$\begin{aligned}
& p(\epsilon, \delta, \eta) \\
&= \left( \mathbb{P}(T'_1 \leq 2\bar{t}/\eta) \right)^{l^*-1} \cdot \int \mathbb{1}\left\{ (\mathbf{w}, \mathbf{t}) \in E(\epsilon) \right\} \\
&\quad \cdot \mathbb{P}(\eta W_1 = dw_1 | \eta|W_1| > \delta) \cdots \mathbb{P}(\eta W_{l^*} = dw_{l^*} | \eta|W_{l^*}| > \delta) \\
&\quad \cdot \mathbb{P}(\eta T'_2 = dt'_2 | \eta T'_2 \leq 2\bar{t}) \cdots \mathbb{P}(\eta T'_{l^*} = dt'_{l^*} | \eta T'_{l^*} \leq 2\bar{t}) \\
&= \left( \mathbb{P}(T'_1 \leq 2\bar{t}/\eta) \right)^{l^*-1} \cdot \mathbb{Q}_{\eta, \delta}(E(\epsilon) \cap h^{-1}(A))
\end{aligned} \tag{B.55}$$

where  $\mathbb{Q}_{\eta, \delta}$  is the Borel-measurable probability measure on  $\mathbb{R}^{l^*} \times \left(\mathbb{R}_+\right)^{l^*-1}$  induced by a sequence of independent random variables  $(W_1^\uparrow(\eta, \delta), \dots, W_{l^*}^\uparrow(\eta, \delta), T_2^\uparrow(\eta, \delta), \dots, T_{l^*}^\uparrow(\eta, \delta))$  such that

- For any  $i = 1, \dots, l^*$ , the distribution of  $W_i^\uparrow(\eta, \delta)$  follows from  $\mathbb{P}\left(\eta W_1 \in \cdot \mid \eta|W_{l^*}| > \delta\right)$ ;
- For any  $i = 2, \dots, l^*$ , the distribution of  $T_i^\uparrow(\eta, \delta)$  follows from  $\mathbb{P}\left(\eta T_1 \in \cdot \mid \eta T_1 \leq 2\bar{t}\right)$ ;
- $\mathbb{Q}_{\eta, \delta}(\cdot) = \mathbb{P}\left((\eta W_1^\uparrow(\eta, \delta), \dots, \eta W_{l^*}^\uparrow(\eta, \delta), \eta T_2^\uparrow(\eta, \delta), \dots, \eta \sum_{j=2}^{l^*} T_j^\uparrow(\eta, \delta)) \in \cdot\right)$ .

Now we study the weak convergence of  $W_1^\uparrow, T_1^\uparrow$ :

- Due to the regularly varying nature of distribution of  $Z_1$  (hence for  $W_1$ ), we know that: for any  $x > \delta$ ,

$$\lim_{\eta \downarrow 0} \mathbb{P}\left(\eta W_1 > x \mid \eta|W_{l^*}| > \delta\right) = p_+ \frac{\delta^\alpha}{x^\alpha}, \quad \lim_{\eta \downarrow 0} \mathbb{P}\left(\eta W_1 < -x \mid \eta|W_{l^*}| > \delta\right) = p_- \frac{\delta^\alpha}{x^\alpha};$$

therefore,  $W_1^\uparrow(\eta, \delta)$  weakly converges to a (randomly signed) Pareto RV that admits the density

$$\nu_{\alpha, \delta}(dx) = \mathbb{1}\{x > 0\} p_+ \frac{\alpha \delta^\alpha}{x^{\alpha+1}} + \mathbb{1}\{x < 0\} p_- \frac{\alpha \delta^\alpha}{|x|^{\alpha+1}}$$

as  $\eta \downarrow 0$ ;

- For any  $x \in [0, 2\bar{t}]$ , since  $\lim_{\eta \downarrow 0} \lfloor x/\eta \rfloor H(\delta/\eta) = 0$ , it is easy to show that

$$\lim_{\eta \downarrow 0} \frac{1 - (1 - H(\delta/\eta))^{\lfloor x/\eta \rfloor}}{\lfloor x/\eta \rfloor H(\delta/\eta)} = 1;$$

therefore, we have (for any  $x \in (0, 2\bar{t}]$ )

$$\mathbb{P}(\eta T_1 \leq x \mid \eta T_1 \leq 2\bar{t}) = \frac{1 - (1 - H(\delta/\eta))^{\lfloor x/\eta \rfloor}}{1 - (1 - H(\delta/\eta))^{\lfloor 2\bar{t}/\eta \rfloor}} \rightarrow \frac{x}{2\bar{t}}$$

as  $\eta \downarrow 0$ , which implies that  $T_1^\uparrow$  converges weakly to a uniform RV on  $[0, 2\bar{t}]$ .



Let us denote the weak limit of measure  $\mathbb{Q}_{\eta,\delta}$  as  $\mu_{\delta,2\bar{t}}$ . In the discussion before the Lemma we have shown that, for any  $(\mathbf{w}, \mathbf{t}) \in E(\epsilon)$  (with  $\delta \in (0, \bar{\delta})$ ), we have  $|w_i| \geq \bar{\delta}$  and  $|t'_i| \leq 2\bar{t}$ ; since we require  $\delta < \bar{\delta}$ , by definition of measures  $\mu$  and  $\mu_{\delta,2\bar{t}}$  we have

$$\mu_{\delta,2\bar{t}}(E(\epsilon) \cap h^{-1}(A)) = \frac{\delta^{\alpha l^*}}{(2\bar{t})^{l^*-1}} \cdot \mu(E(\epsilon) \cap h^{-1}(A)).$$

For simplicity of notations, we let  $E(\epsilon, A) \triangleq E(\epsilon) \cap h^{-1}(A)$ . By definition of the set  $E(\epsilon)$ , we have (recall that  $A$  is an open interval  $(a_1, a_2)$  that does not overlap with  $[s_- + \bar{\epsilon}, s_+ - \bar{\epsilon}]$ )

$$\begin{aligned} E(\epsilon, A) &= h^{-1}((-\infty, s_- - \epsilon) \cup (s_+ + \epsilon, \infty)) \cap h^{-1}((a_1, a_2)) \\ &= h^{-1}\left(\left((-\infty, s_- - \epsilon) \cup (s_+ + \epsilon, \infty)\right) \cap (a_1, a_2)\right) \\ &= h^{-1}(F(\epsilon, a_1, a_2)) \end{aligned}$$

where  $F(\epsilon, a_1, a_2) \triangleq ((-\infty, s_- - \epsilon) \cup (s_+ + \epsilon, \infty)) \cap (a_1, a_2)$ . Meanwhile, it is easy to see that  $h$  is a continuous mapping, hence

$$(\mathbf{w}, \mathbf{t}) \in \partial E(\epsilon, A) \Rightarrow h(\mathbf{w}, \mathbf{t}) \in \{s_- + \epsilon, s_+ - \epsilon, a_1, a_2\}.$$

Fix some  $s$  with  $s \neq \pm L, |s| > (l^* - 1)b + \bar{\epsilon}$ . For any fixed real numbers  $t_2, \dots, t_{l^*-1}, w_1, \dots, w_{l^*}$ , if  $h(w_1, \dots, w_{l^*}, t_2, \dots, t_{l^*-1}, t) = s$ , then since  $\tilde{\mathbf{x}}(t_{l^*} - 1, 0; w_1, \dots, w_{l^*}, t_2, \dots, t_{l^*-1}, t) \in [s - b, s + b]$ , due to Assumption 1 (in particular, there is no point  $x$  on this interval with  $|f'(x)| \leq c_0$ ), there exists at most one possible  $t$  that makes  $h(w_1, \dots, w_{l^*}, t_2, \dots, t_{l^*-1}, t) = s$ . Therefore, let  $W_j^*$  be iid RVs from law  $\nu_{\alpha,\delta}$  defined above, and  $(T_j^{*,'})_{j \geq 2}$  be iid RVs from  $\text{Unif}[0, 2\bar{t}]$ ,  $T_0^* = 0, T_k^* = \sum_{j=2}^k T_j^{*,'}$ . By conditioning on all  $W_j^*$  and all  $T_2^{*,'}, \dots, T_{l^*-1}^{*,'}$ , we must have

$$\mathbb{P}\left(h(W_1^*, \dots, W_j^*, T_2^*, \dots, T_{l^*}^*) = s \mid W_1^* = dw_1, \dots, W_{l^*}^* = dw_{l^*}, T_2^{*,'} = dt_2, \dots, T_{l^*-1}^{*,'} = dt_{l^*-1}\right) = 0 \quad (\text{B.56})$$

which implies

$$\mathbb{P}\left(h(W_1^*, \dots, W_j^*, T_2^*, \dots, T_{l^*}^*) = s\right) = 0$$

hence

$$\mu(\partial E(\epsilon, A)) = 0.$$

By Portmanteau theorem (see Theorem 2.1 of [1]) we have

$$\lim_{\eta \downarrow 0} \mathbb{Q}_{\eta,\delta}(E(\epsilon, A)) = \mu_{\delta,2\bar{t}}(E(\epsilon, A)).$$

Collecting the results we have and using (B.55), we can see that

$$\begin{aligned} \limsup_{\eta \downarrow 0} \frac{p(\epsilon, \delta, \eta; A)}{\left(\frac{H(1/\eta)}{\eta}\right)^{l^*-1} \delta^\alpha} &= \limsup_{\eta \downarrow 0} \frac{(2\bar{t})^{l^*-1} \cdot p(\epsilon, \delta, \eta; A)}{\delta^{\alpha l^*} \cdot \left(\mathbb{P}(T_1' \leq 2\bar{t}/\eta)\right)^{l^*-1}} \cdot \left(\frac{\delta^\alpha}{2\bar{t}} \cdot \frac{\mathbb{P}(T_1' \leq 2\bar{t}/\eta)}{H(1/\eta)/\eta}\right)^{l^*-1} \\ &\leq \limsup_{\eta \downarrow 0} \frac{(2\bar{t})^{l^*-1} \cdot p(\epsilon, \delta, \eta; A)}{\delta^{\alpha l^*} \cdot \left(\mathbb{P}(T_1' \leq 2\bar{t}/\eta)\right)^{l^*-1}} \cdot \limsup_{\eta \downarrow 0} \left(\frac{\delta^\alpha}{2\bar{t}} \cdot \frac{\mathbb{P}(T_1' \leq 2\bar{t}/\eta)}{H(1/\eta)/\eta}\right)^{l^*-1} \\ &\leq \mu(E(\epsilon, A)) \cdot \limsup_{\eta \downarrow 0} \left(\frac{\delta^\alpha}{2\bar{t}} \cdot \frac{\mathbb{P}(T_1' \leq 2\bar{t}/\eta)}{H(1/\eta)/\eta}\right)^{l^*-1}. \end{aligned}$$

Fix some  $\kappa > 1$ . From Lemma B.2 and the regularly varying nature of function  $H$ , we get

$$\limsup_{\eta \downarrow 0} \left( \frac{\delta^\alpha}{2\bar{t}} \cdot \frac{\mathbb{P}(T'_1 \leq 2\bar{t}/\eta)}{H(1/\eta)/\eta} \right)^{l^*-1} \leq \kappa^{l^*-1} \limsup_{\eta \downarrow 0} \left( \frac{\delta^\alpha}{2\bar{t}} \cdot \frac{2\bar{t}H(\delta/\eta)/\eta}{H(1/\eta)/\eta} \right)^{l^*-1} = \kappa^{l^*-1}.$$

Due to the arbitrariness of  $\kappa > 1$ , we have established that

$$\limsup_{\eta \downarrow 0} \frac{p(\epsilon, \delta, \eta; A)}{\left( \frac{H(1/\eta)}{\eta} \right)^{l^*-1} \delta^\alpha} \leq \mu(E(\epsilon)).$$

The lower bound can be shown by an argument symmetric to the one for upper bound.  $\square$

The following result is an immediate corollary of Lemma B.14.

**Corollary B.15.** *For any  $\epsilon \in (-\bar{\epsilon}, \bar{\epsilon})$ ,  $\delta \in (0, b \wedge \bar{\delta})$ , we have*

$$\lim_{\eta \downarrow 0} \frac{p(\epsilon, \delta, \eta)}{\delta^\alpha \left( \frac{H(1/\eta)}{\eta} \right)^{l^*-1}} = \mu(E(\epsilon))$$

where  $\mu$  is the measure defined in (B.52), and  $p(\cdot, \cdot, \cdot)$  is the conditional probability defined in (B.53).

Define the following stopping times:

$$\sigma(\eta) = \min\{n \geq 0 : X_n^\eta \notin \Omega\}; \quad (\text{B.57})$$

$$R(\epsilon, \delta, \eta) = \min\{n \geq T_1^\eta(\delta) : X_n^\eta \in [-2\epsilon, 2\epsilon]\}. \quad (\text{B.58})$$

$\sigma$  indicate the time that the iterates escape the current attraction field, while  $R$  denotes the time the SGD iterates return to a small neighborhood of the local minimum after first exit from this small neighborhood. In the next few results, we study the probability of several atypical scenarios when SGD iterates make attempts to escape  $\Omega$  or return to local minimum after the attempt fails. First, we show that, when starting from the local minimum, it is very unlikely to escape with less than  $l^*$  big jumps.

**Lemma B.16.** *Given  $\epsilon \in (0, \bar{\epsilon})$ ,  $N > 0$ , the following claim holds for any sufficiently small  $\delta > 0$ :*

$$\sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x \left( \sigma(\eta) < R(\epsilon, \eta), \sigma(\eta) < T_{l^*}^\eta(\delta) \right) = o(\eta^N)$$

as  $\eta \downarrow 0$ .

*Proof.* Based on the given  $\epsilon > 0$ , fix some  $\tilde{\epsilon} = \frac{\epsilon}{4 \exp(2M\tilde{t}(\epsilon))}$ . Recall the definition of  $\hat{t}(\epsilon)$  in (B.24).

First, using Lemma B.12, we know that for sufficiently small  $\delta$ , we have

$$\sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P} \left( A_1^\times(\epsilon, \delta, \eta) \right) = o(\eta^N) \quad (\text{B.59})$$

where

$$A_1^\times(\epsilon, \delta, \eta) = \left\{ \exists n < T_1^\eta(\delta) \text{ s.t. } |X_n^\eta(x)| > 3\epsilon \right\}.$$

Define event

$$A_2^\times(\tilde{\epsilon}, \delta, \eta) \triangleq \left\{ \exists j = 2, \dots, l^* \text{ s.t. } \max_{k=1,2,\dots, T_j^\eta(\delta) - T_{j-1}^\eta(\delta) - 1} \eta |Z_{T_{j-1}^\eta(\delta)+1} + \dots + Z_{T_{j-1}^\eta(\delta)+j}| > \tilde{\epsilon} \right\}.$$

From Corollary B.13, we know that for sufficiently small  $\delta > 0$ ,

$$\mathbb{P}\left(A_2^\times(\tilde{\epsilon}, \delta, \eta)\right) = o(\eta^N). \quad (\text{B.60})$$

From now on, we only consider such  $\delta$  that (B.59)(B.60) hold.

On event  $\left(A_1^\times \cup A_2^\times\right)^c \cap \{\sigma(\eta) < R(\epsilon, \eta)\} \cap \{\sigma(\eta) > T_1^\eta(\delta)\}$ , we must have  $\sigma(\eta) > T_1^\eta(\delta)$  and

$$T_2^\eta(\delta) \wedge \sigma(\eta) - T_1^\eta(\delta) < 2\hat{t}(\epsilon)/\eta.$$

Otherwise, due to Lemma B.4 and B.7, we know that at step  $\tilde{t} = T_1^\eta(\delta) + \lfloor \hat{t}(\epsilon)/\eta \rfloor$ , we have

$$|X_{\tilde{t}}^\eta| < 2\epsilon, \text{ and } |X_n^\eta| \leq \bar{\epsilon} \quad \forall n \leq \tilde{t}$$

for any sufficiently small  $\eta$ . By repeating this argument inductively, we obtain the following result: define

$$J = \min\{j = 1, 2, \dots : \sigma(\eta) \in [T_j^\eta(\delta), T_{j+1}^\eta(\delta))\},$$

then on event  $\left(A_1^\times \cup A_2^\times\right)^c \cap \{\sigma(\eta) < R(\epsilon, \eta), \sigma(\eta) < T_{l^*}^\eta(\delta)\}$ , we must have

$$T_j^\eta(\delta) \wedge \sigma(\eta) - T_{j-1}^\eta(\delta) \wedge \sigma(\eta) < 2\hat{t}(\epsilon)/\eta \quad \forall j = 2, 3, \dots, J. \quad (\text{B.61})$$

Furthermore, using this bound and Lemma B.7, we know that on event  $\left(A_1^\times \cup A_2^\times\right)^c \cap \{\sigma(\eta) < R(\epsilon, \eta), \sigma(\eta) < T_{l^*}^\eta(\delta)\}$ ,

- $|X_{T_j^\eta(\delta)}^\eta| \leq |X_{T_{j-1}^\eta(\delta)}^\eta| + b + \epsilon + \bar{\epsilon}$  for all  $j = 2, 3, J-1$ ,
- $|X_{\sigma(\eta)}^\eta| \leq |X_{T_{J-1}^\eta(\delta)}^\eta| + \epsilon + \bar{\epsilon}$

However, this implies

$$|X_{\sigma(\eta)}^\eta| \leq l^*(\bar{\epsilon} + \epsilon) + (l^* - 1)b < r$$

and contradicts the definition of  $\sigma(\eta)$ . In summary,

$$\sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x\left(\sigma(\eta) < R(\epsilon, \eta), \sigma(\eta) < T_{l^*}^\eta(\delta)\right) \leq \mathbb{P}\left(A_1^\times(\epsilon, \delta, \eta) \cup A_2^\times(\tilde{\epsilon}, \delta, \eta)\right) = o(\eta^N).$$

□

The following two results follow immediately from the proof above, especially the inductive argument leading to bound (B.61), and we state them without repeating the details of the proof.

**Corollary B.17.** *Given  $\epsilon \in (0, \bar{\epsilon})$ ,  $N > 0$ , the following claim holds for any sufficiently small  $\delta > 0$ :*

$$\sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x\left(T_{l^*}^\eta(\delta) \leq \sigma(\eta) \wedge R(\epsilon, \eta), \text{ and } \exists j = 2, 3, \dots, l^* \text{ s.t. } T_j^\eta(\delta) - T_{j-1}^\eta(\delta) > 2\hat{t}(\epsilon)/\eta\right) = o(\eta^N)$$

as  $\eta \downarrow 0$ .

**Corollary B.18.** *Given  $\epsilon \in (0, \bar{\epsilon})$ ,  $N > 0$ , the following claim holds for any sufficiently small  $\delta > 0$ :*

$$\sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x\left(R(\epsilon, \eta) < T_{l^*}^\eta(\delta) \wedge \sigma(\eta), R(\epsilon, \eta) - T_1^\eta(\delta) > 2l^*\hat{t}(\epsilon)/\eta\right) = o(\eta^N)$$

as  $\eta \downarrow 0$ .

In the next result, we show that, if the inter-arrival time between some large jumps are too long, or some large jumps are still not *large enough*, then it is very unlikely that the SGD iterates could escape at the time of  $l^*$ -th large jump (or even get close enough to the boundary of the attraction field).

**Lemma B.19.** *Given  $\epsilon \in (0, \bar{\epsilon})$  and any  $N > 0$ , the following holds for all  $\delta > 0$  that are sufficiently small:*

$$\sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x(B_2^\times(\epsilon, \delta, \eta)) = o(\eta^N).$$

where

$$B_2^\times(\epsilon, \delta, \eta) = \{T_{l^*}^\eta(\delta) \leq \sigma(\eta) \wedge R(\epsilon, \eta)\} \cap \left\{ \exists j = 2, 3, \dots, l^* \text{ s.t. } T_j^\eta(\delta) - T_{j-1}^\eta(\delta) > \bar{t}/\eta \right. \\ \left. \text{or } \exists j = 1, 2, \dots, l^* \text{ s.t. } \eta|W_j^\eta(\delta)| \leq \bar{\delta} \right\} \cap \{|X_{T_{l^*}^\eta}^\eta| \geq r - \bar{\epsilon}\}.$$

*Proof.* Let  $A_1^\times, A_2^\times$  be the events defined in the proof of Lemma B.16. Based on the given  $\epsilon > 0$ , fix some  $\tilde{\epsilon} = \frac{\epsilon}{4 \exp(2M\bar{t}(\epsilon))}$ .

Let  $J = \min\{j = 2, 3, \dots : T_j^\eta(\delta) - T_{j-1}^\eta(\delta) > \bar{t}/\eta\}$ . On event  $\left(A_1^\times(\epsilon, \delta, \eta) \cup A_2^\times(\tilde{\epsilon}, \delta, \eta)\right)^c \cap \{J \leq l^*\}$ , from Lemma B.4 and B.7 and the definition of constant  $\bar{t}$ , we know that

- $|X_{T_j^\eta(\delta)}^\eta| \leq |X_{T_{j-1}^\eta(\delta)}^\eta| + b + \epsilon + \bar{\epsilon}$  for all  $j = 2, 3, J-1$ ;
- $|X_{T_J^\eta(\delta)}^\eta| \leq 2\bar{\epsilon}$
- $|X_n^\eta| < s - \bar{\epsilon} \quad \forall n \leq T_J^\eta(\delta)$

Now starting from step  $T_J^\eta(\delta)$ , by using Lemma B.4 and B.7 again one can see that

- $|X_{T_j^\eta(\delta)}^\eta| \leq |X_{T_{j-1}^\eta(\delta)}^\eta| + b + \epsilon + \bar{\epsilon}$  for all  $j = J+1, \dots, l^*$ .

Combining these results, we have that  $|X_{T_{l^*}^\eta}^\eta| < r - \bar{\epsilon}$  on event  $\left(A_1^\times(\epsilon, \delta, \eta) \cup A_2^\times(\tilde{\epsilon}, \delta, \eta)\right)^c \cap \{J \leq l^*\}$ .

Next, define  $J' = \min\{j = 1, 2, \dots : \eta|W_j^\eta(\delta)| \leq \bar{\delta}\}$ . Similarly, on event  $\left(A_1^\times(\epsilon, \delta, \eta) \cup A_2^\times(\tilde{\epsilon}, \delta, \eta)\right)^c \cap \{J > l^*\} \cap \{J' \leq l^*\}$ , using Lemma B.4 and B.7 again one can see that

- $|X_{T_j^\eta(\delta)}^\eta| \leq |X_{T_{j-1}^\eta(\delta)}^\eta| + b + \epsilon + \bar{\epsilon}$  for all  $j = 1, 2, \dots, l^*, j \neq J'$ ;
- $|X_{T_{J'}^\eta(\delta)}^\eta| \leq |X_{T_{J'-1}^\eta(\delta)}^\eta| + \bar{\delta} + \epsilon + \bar{\epsilon}$  for all  $j = 1, 2, \dots, l^*, j \neq J'$ .

Since  $\bar{\delta} \in (0, \bar{\epsilon})$ , we have  $|X_{T_{l^*}^\eta}^\eta| < r - \bar{\epsilon}$  on this event.

In summary, the following bound

$$\sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x(B_2^\times) \leq \mathbb{P}(A_1^\times(\epsilon, \delta, \eta) \cup A_2^\times(\tilde{\epsilon}, \delta, \eta)) = o(\eta^N)$$

holds for any  $\delta$  that is sufficiently small, which is established in Lemma B.16. This conclude the proof.  $\square$

In the next lemma, we show that, starting from the local minimum, it is unlikely that the SGD iterates will be right at the boundary of the attraction field after  $l^*$  large jumps. Recall that there are  $n_{\min}$  attraction fields on  $f$ , and excluding  $s_0 = -\infty, s_{n_{\min}} = \infty$  the remaining points  $s_1, \dots, s_{n_{\min}-1}$  are the boundaries of the attraction fields.

**Lemma B.20.** *There exists a function  $\Psi(\cdot) : \mathbb{R}^+ \mapsto \mathbb{R}^+$  satisfying  $\lim_{\epsilon \downarrow 0} \Psi(\epsilon) = 0$  such that the following claim folds. Given  $\epsilon \in (0, \bar{\epsilon}/(3\bar{\rho} + 3\tilde{\rho} + 9))$ , we have*

$$\limsup_{\eta \downarrow 0} \frac{\sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x(B_3^\times(\epsilon, \delta, \eta))}{\left(H(1/\eta)/\eta\right)^{l^*-1}} \leq \delta^\alpha \Psi(\epsilon)$$

for all  $\delta$  sufficiently small, where  $\bar{\rho}$  and  $\tilde{\rho}$  are the constants defined in Corollary B.5 and B.8, and the event is defined as

$$B_3^\times(\epsilon, \delta, \eta) = \left\{ T_{l^*}^\eta(\delta) \leq \sigma(\eta) \wedge R(\epsilon, \eta) \right\} \cap \left\{ \exists k \in [n_{\min} - 1] \text{ such that } X_{T_{l^*}^\eta(\delta)}^\eta \in [s_k - \epsilon, s_k + \epsilon] \right\}.$$

*Proof.* Let  $A_1^\times, A_2^\times$  be the events defined in the proof of Lemma B.16. Based on the given  $\epsilon > 0$ , fix some  $\tilde{\epsilon} = \frac{\epsilon}{4 \exp(2M\hat{t}(\epsilon))}$ . Fix some  $N > \alpha l^*$ .

Choose  $\delta$  small enough so that claim in Lemma B.19 holds for the  $\epsilon$  prescribed. Using the same arguments in Lemma B.19, we have the following inclusion of events:

$$\begin{aligned} & B_3^\times(\epsilon, \delta, \eta) \cap \left( A_1^\times(\epsilon, \delta, \eta) \cup A_2^\times(\tilde{\epsilon}, \delta, \eta) \right)^c \\ & \subseteq \left\{ \forall j = 2, 3, \dots, l^*, T_j^\eta(\delta) - T_{j-1}^\eta(\delta) \leq \bar{t}/\eta \right\} \cap \left\{ \forall j = 1, 2, 3, \dots, l^*, \eta |W_1^\eta(\delta)| > \bar{\delta} \right\}. \end{aligned}$$

Therefore, on event  $B_3^\times(\epsilon, \delta, \eta) \cap \left( A_1^\times(\epsilon, \delta, \eta) \cup A_2^\times(\tilde{\epsilon}, \delta, \eta) \right)^c$ , we can apply Corollary B.5 and B.8 and conclude that  $Z_{T_1^\eta(\delta)}$  has  $(-\bar{\epsilon}, \delta, \eta)$ -overflow, and its  $(-\bar{\epsilon}, \delta, \eta)$ -overflow endpoint lies

$$(s_k - 3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, s_k + 3(\bar{\rho} + \tilde{\rho} + 3)\epsilon)$$

for some  $k \in [n_{\min} - 1]$ . Using Lemma B.14 and Corollary B.15, we have that (for any sufficiently small  $\eta$ )

$$\begin{aligned} & \mathbb{P}\left(B_3^\times(\epsilon, \delta, \eta) \cap \left( A_1^\times(\epsilon, \delta, \eta) \cup A_2^\times(\tilde{\epsilon}, \delta, \eta) \right)^c\right) \\ & \leq \delta^\alpha \left( \frac{H(1/\eta)}{\eta} \right)^{l^*-1} \cdot \sum_{k=1}^{n_{\min}-1} \mu(E(-\bar{\epsilon}) \cap h^{-1}((s_k - \hat{\epsilon}, s_k + \hat{\epsilon}))) \end{aligned}$$

where  $\hat{\epsilon} = 3(\bar{\rho} + \tilde{\rho} + 3)\epsilon$ . Besides, as established in the proof of Lemma B.16, we have

$$\mathbb{P}\left(A_1^\times(\epsilon, \delta, \eta) \cup A_2^\times(\tilde{\epsilon}, \delta, \eta)\right) = o(\eta^N)$$

for all sufficiently small  $\delta$ . In conclusion, we only need to choose

$$\Psi(\epsilon) = \sum_{k=1}^{n_{\min}-1} \mu\left(E(-\bar{\epsilon}) \cap h^{-1}((s_k - 3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, s_k + 3(\bar{\rho} + \tilde{\rho} + 3)\epsilon))\right).$$

To conclude the proof, just note that by combining the continuity of measure with the conditional probability argument leading to (B.56), we can show that  $\lim_{\epsilon \downarrow 0} \Psi(\epsilon) = 0$ .  $\square$

Lastly, we establish the lower bound for the probability of the *most likely* way for SGD iterates to exit the current attraction field: making  $l^*$  large jumps in a relatively short period of time. Recall that  $\bar{\epsilon}$  is the fixed constant in (B.13)-(B.16).

**Lemma B.21.** *Given  $\epsilon \in (0, \bar{\epsilon}/3)$ , it holds for any sufficiently small  $\delta > 0$  such that*

$$\liminf_{\eta \downarrow 0} \frac{\inf_{|x| \leq 2\epsilon} \mathbb{P}_x(A^\circ(\epsilon, \delta, \eta))}{(H(1/\eta)/\eta)^{l^*-1}} \geq c_* \delta^\alpha$$

where the event is defined as

$$\begin{aligned} A^\circ(\epsilon, \delta, \eta) \triangleq & \left\{ \sigma(\eta) < R(\epsilon, \eta), \sigma(\eta) = T_{l^*}^\eta(\delta), X_{T_{l^*}^\eta}^\eta \notin [s_- - \epsilon, s_+ + \epsilon] \right\} \\ & \cap \left\{ T_j^\eta(\delta) - T_{j-1}^\eta(\delta) \leq \frac{\bar{\epsilon}}{2M} \lceil 1/\eta \rceil \quad \forall j = 2, 3, \dots, l^* \right\} \end{aligned}$$

and the constant

$$c_* = \frac{1}{2} \left( \frac{1}{2b} \right)^{l^* \alpha} \left( \frac{\bar{\epsilon}}{4M} \right)^{l^* - 1}$$

is strictly positive and does not vary with  $\epsilon, \delta$ .

*Proof.* Let  $A_1^\times, A_2^\times$  be the events defined in the proof of Lemma B.16. Fix some  $N$  such that  $N > \alpha l^*$ . Based on the given  $\epsilon > 0$ , fix some  $\tilde{\epsilon} = \frac{\epsilon}{4 \exp(2Mt(\epsilon))}$ . We only consider  $\delta < M$ . Furthermore, choose  $\delta$  small enough so that (B.59) and (B.60) hold for the chosen  $N$  and  $\epsilon$ . Also, we only consider  $\eta$  small enough so that  $\eta M < b \wedge \bar{\epsilon}$ .

Due to (B.13)-(B.16), we can, without loss of generality, assume that  $r = s_+$ , and in this case we will have

$$l^* b - 100l^* \bar{\epsilon} > s_+ + 100l^* \bar{\epsilon}.$$

Under this assumption, we will now focus on providing a lower bound for the following event that describes the exit from the right side of  $\Omega$  (in other words, by crossing  $s_+$ )

$$\begin{aligned} A_{\rightarrow}^\circ(\epsilon, \delta, \eta) \triangleq & \left\{ \sigma(\eta) < R(\epsilon, \eta), \sigma(\eta) = T_{l^*}^\eta(\delta), X_{T_{l^*}^\eta}^\eta > s_+ + \epsilon \right\} \\ & \cap \left\{ T_j^\eta(\delta) - T_{j-1}^\eta(\delta) \leq \frac{\bar{\epsilon}}{2M} \lceil 1/\eta \rceil \quad \forall j = 2, 3, \dots, l^* \right\}. \end{aligned}$$

First, define event

$$A_3^\circ(\delta, \eta) = \left\{ W_j^\eta(\delta) \geq 2b \quad \forall j = 1, \dots, l^*, T_j^\eta(\delta) - T_{j-1}^\eta(\delta) \leq \frac{\bar{\epsilon}}{2M} \lceil 1/\eta \rceil \quad \forall j = 2, \dots, l^* \right\},$$

and observe some facts on event  $A_3^\circ(\delta, \eta) \cap \left( A_1^\times(\epsilon, \delta, \eta) \cup A_2^\times(\tilde{\epsilon}, \delta, \eta) \right)^c$ .

- $|X_k^\eta| \leq 3\epsilon \forall n < T_1^\eta(\delta)$ ; (due to  $A_1^\times$  not occurring)
- $X_{T_1^\eta(\delta)}^\eta \in [b - 3\epsilon, b + 3\epsilon]$ ; (due to  $W_1^\eta \geq 2b$  and the effect of gradient clipping at step  $T_1^\eta$ , as well as the fact that  $X_{T_1^\eta-1}^\eta \in [-3\epsilon, 3\epsilon]$  from the previous bullet point)
- Due to  $|f'(\cdot)| \leq M$  and  $\delta < M$ , one can see that (for any  $n \geq 1$ )

$$\sup_{x \in [-L, L]} |\eta f'(x)| + |\eta Z_n^{\leq \delta, \eta}| \leq 2\eta M;$$

this provides an upper bound for the change in SGD iterates at each step, and gives us

$$X_n^\eta \in [b - 3\epsilon - \bar{\epsilon}, b + 3\epsilon + \bar{\epsilon}] \quad \forall T_1^\eta(\delta) \leq n < T_2^\eta(\delta)$$

where we also used  $T_2^\eta(\delta) - T_1^\eta(\delta) \leq \frac{\bar{\epsilon}}{2M} \lceil 1/\eta \rceil$

- Therefore, at the arrival time of the second large jump, we must have  $X_{T_2^\eta(\delta)}^\eta \geq 2b - 3\epsilon - \bar{\epsilon}$ ;
- By repeating the argument above inductively, we can show that (for all  $j = 1, 2, \dots, l^*$ )

$$\begin{aligned} X_n^\eta &\in [(j-1)b - 3\epsilon - (j-1)\bar{\epsilon}, (j-1)b + 3\epsilon + (j-1)\bar{\epsilon}] \quad \forall T_{j-1}^\eta \leq n < T_j^\eta \\ X_{T_j^\eta}^\eta &\in [jb - 3\epsilon - (j-1)\bar{\epsilon}, jb + 3\epsilon + (j-1)\bar{\epsilon}]; \end{aligned}$$

In particular, we know that  $X_n^\eta \in \Omega$  for any  $n < T_{l^*}^\eta$  (so the exit does not occur before  $T_{l^*}^\eta$ ), and at the arrival of the  $l^*$ -th large jump, we have (using  $3\epsilon < \bar{\epsilon}$ )

$$X_{T_{l^*}^\eta(\delta)}^\eta \geq l^*b - l^*\bar{\epsilon} > s_+ + \epsilon.$$

In summary, we have shown that

$$A_3^\circ(\delta, \eta) \cap \left( A_1^\times(\epsilon, \delta, \eta) \cup A_2^\times(\bar{\epsilon}, \delta, \eta) \right)^c \subseteq A_{\rightarrow}^\circ(\epsilon, \delta, \eta).$$

To conclude the proof, just notice that (for sufficiently small  $\eta$ )

$$\begin{aligned} &\mathbb{P}\left(A_3^\circ(\delta, \eta) \cap \left( A_1^\times(\epsilon, \delta, \eta) \cup A_2^\times(\bar{\epsilon}, \delta, \eta) \right)^c\right) \\ &\geq \mathbb{P}(A_3^\circ(\delta, \eta)) - \mathbb{P}(A_1^\times(\epsilon, \delta, \eta)) - \mathbb{P}(A_2^\times(\bar{\epsilon}, \delta, \eta)) \\ &\geq \mathbb{P}(A_3^\circ(\delta, \eta)) - \eta^N \quad \text{due to (B.59) and (B.60)} \\ &\geq \left( \frac{H(2b/\eta)}{H(\delta/\eta)} \right)^{l^*} \left( \frac{\bar{\epsilon}}{4M} H(\delta/\eta)/\eta \right)^{l^*-1} - \eta^N \quad \text{due to Lemma B.2} \\ &\geq 2c_*\delta^\alpha (H(1/\eta)/\eta)^{l^*-1} - \eta^N \quad \text{for all } \eta \text{ sufficiently small, due to } H \in \mathcal{RV}_{-\alpha} \\ &\geq c_*\delta^\alpha (H(1/\eta)/\eta)^{l^*-1}. \end{aligned}$$

□

In order to present the main result of this section, we need to take into account the loss landscape outside of the current attraction field  $\Omega$ . Recall that there are  $n_{\min}$  attraction fields on  $f$ . For all the attraction fields different from  $\Omega$ , we call them  $(\tilde{\Omega}_k)_{k=1}^{n_{\min}-1}$  where, for each  $k \in [n_{\min} - 1]$ , the attraction field  $\tilde{\Omega}_k = (s_k^-, s_k^+)$  with the corresponding local minimum located at  $\tilde{m}_k$ . Also, recall that  $\sigma(\eta)$  is the first time  $X_n^\eta$  exits from  $\Omega$ . Building upon these concepts, we can define a stopping time

$$\tau(\eta, \epsilon) \triangleq \min\{n \geq 0 : X_n^\eta \in \bigcup_{k=1}^{n_{\min}-1} [\tilde{m}_k - 2\epsilon, \tilde{m}_k + 2\epsilon]\} \quad (\text{B.62})$$

as the first time the SGD iterates visit a minimizer in an attraction field that is different from  $\Omega$ . Besides, let index  $J_\sigma(\eta)$  be such that

$$J_\sigma(\eta) = j \iff X_{\sigma(\eta)}^\eta \in \tilde{\Omega}_j \quad \forall j \in [n_{\min} - 1]. \quad (\text{B.63})$$

In other words, it is the label of the attraction field that  $X_n^\eta$  escapes to. Lastly, define

$$\lambda(\eta) \triangleq H(1/\eta) \left( H(1/\eta)/\eta \right)^{l^*-1}, \quad (\text{B.64})$$

$$\nu^\Omega \triangleq \mu(E(0)), \quad (\text{B.65})$$

$$\nu_k^\Omega \triangleq \mu(E(0) \cap h^{-1}(\tilde{\Omega}_k)) \quad \forall k \in [n_{\min} - 1]. \quad (\text{B.66})$$

For definitions of the measure  $\mu$ , set  $E$ , and mapping  $h$ , see (B.51) and (B.52).

Now we are ready to state Proposition B.22, the most important technical tool in this section. In (B.67) and (B.68), we provide upper and lower bounds for the joint distribution of first exit time  $\sigma$  and the label  $J_\sigma$  indexing the attraction field we escape to; it is worth noticing that the claims hold uniformly for all  $u > C$ . In (B.69) and (B.70), we provide upper and lower bounds for the joint distribution of when we first visit a different local minimum (which is equal to  $\tau$ ) and which one we visit (indicated by  $X_\tau^\eta$ ). The similarity between (B.67) (B.68) and (B.69)(B.70) suggests a strong correlation between the behavior of the SGD iterates at time  $\sigma(\eta)$  and that of time  $\tau(\eta, \epsilon)$ , and this is corroborated by (B.71): we show that it is almost always the case that  $\tau$  is very close to  $\sigma$ , and on the short time interval  $[\sigma(\eta), \tau(\eta, \epsilon)]$  the SGD iterates stay within the same attraction field.

**Proposition B.22.** *Given  $C > 0$  and some  $k' \in [n_{\min} - 1]$ , the following claims hold for all  $\epsilon > 0$  that is sufficiently small:*

$$\limsup_{\eta \downarrow 0} \sup_{u \in (C, \infty)} \sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x \left( \nu^\Omega \lambda(\eta) \sigma(\eta) > u, J_\sigma(\eta) = k' \right) \leq 2C + \exp(-(1-C)^3 u) \frac{\nu_{k'}^\Omega + C}{\nu^\Omega}, \quad (\text{B.67})$$

$$\liminf_{\eta \downarrow 0} \inf_{u \in (C, \infty)} \inf_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x \left( \nu^\Omega \lambda(\eta) \sigma(\eta) > u, J_\sigma(\eta) = k' \right) \geq -2C + \exp(-(1+C)^3 u) \frac{\nu_{k'}^\Omega - C}{\nu^\Omega}, \quad (\text{B.68})$$

$$\limsup_{\eta \downarrow 0} \sup_{u \in (C, \infty)} \sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x \left( \nu^\Omega \lambda(\eta) \tau(\eta, \epsilon) > u, X_{\tau(\eta, \epsilon)}^\eta \in B(\tilde{m}_{k'}, 2\epsilon) \right) \leq 4C + \exp(-(1-C)^3 u) \frac{\nu_{k'}^\Omega + C}{\nu^\Omega}, \quad (\text{B.69})$$

$$\liminf_{\eta \downarrow 0} \inf_{u \in (C, \infty)} \inf_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x \left( \nu^\Omega \lambda(\eta) \tau(\eta, \epsilon) > u, X_{\tau(\eta, \epsilon)}^\eta \in B(\tilde{m}_{k'}, 2\epsilon) \right) \geq -4C + \exp(-(1+C)^3 u) \frac{\nu_{k'}^\Omega - C}{\nu^\Omega}, \quad (\text{B.70})$$

$$\liminf_{\eta \downarrow 0} \inf_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x \left( \lambda(\eta) (\tau(\eta, \epsilon) - \sigma(\eta)) < C, X_n^\eta \in \tilde{\Omega}_{J_\sigma(\eta)} \quad \forall n \in [\sigma(\eta), \tau(\eta, \epsilon)] \right) \geq 1 - C. \quad (\text{B.71})$$

Before presenting the proof to Proposition B.22, we make some preparations. First, we introduce stopping times (for all  $k \geq 1$ )

$$\begin{aligned} \tau_k(\epsilon, \delta, \eta) &= \min\{n > \tilde{\tau}_{k-1}(\epsilon, \delta, \eta) : |\eta| Z_n| > \delta\} \\ \tilde{\tau}_k(\epsilon, \delta, \eta) &= \min\{n \geq \tau_k(\epsilon, \delta, \eta) : |X_n^\eta| \leq 2\epsilon\} \end{aligned}$$

with the convention that  $\tau_0(\epsilon, \delta, \eta) = \tilde{\tau}_0(\epsilon, \delta, \eta) = 0$ . The intuitive interpretation is as follows. For the fixed  $\epsilon$  we treat  $[-2\epsilon, 2\epsilon]$  as a small neighborhood of the local minimum of the attraction field  $\Omega$ . All the  $\tilde{\tau}_k$  partitioned the entire timeline into different *attempts* of escaping  $\Omega$ . The interval  $[\tilde{\tau}_{k-1}, \tilde{\tau}_k]$  can be viewed as the  $k$ -th *attempt*. If for  $\sigma(\eta)$ , the first exit time defined in (B.57), we have  $\sigma(\eta) > \tilde{\tau}_k$ , then we consider the  $k$ -th attempt of escape as a *failure* because the SGD iterates returned to this small neighborhood of the local minimum again without exiting the attraction field. On the other hand, the stopping times  $\tau_{k-1}$  indicate the arrival time of the first large jump during the  $k$ -th *attempt*. The proviso that  $\tilde{\tau}_k \geq \tau_{k-1}$  can be interpreted, intuitively, as that an *attempt* is considered failed only if, after some significant efforts to exit (for instance, a large jump) has been observed, the SGD iterates still returned to the small neighborhood  $[-2\epsilon, 2\epsilon]$ . Regarding the notations, we add a remark that when there is no ambiguity we will drop the dependency on  $\epsilon, \delta, \eta$  and simply write  $\tau_k, \tilde{\tau}_k$ .

To facilitate the characterization of events during each *attempt*, we introduce the following definitions. First, for all  $k \geq 1$ , let

$$\mathbf{j}_k \triangleq \#\{n = \tau_{k-1}(\epsilon, \delta, \eta), \tau_{k-1}(\epsilon, \delta, \eta) + 1, \dots, \tilde{\tau}_k(\epsilon, \delta, \eta) \wedge \sigma(\eta) : |\eta| Z_n| > \delta\}$$

be the number of large jumps during the  $k$ -th *attempt*. Two implications of this definition:



- First, for any  $k$  with  $\sigma(\eta) < \tilde{\tau}_k$ , we have  $\mathbf{j}_k = 0$ . Note that this proposition concerns the dynamics of SGD up until  $\sigma(\eta)$ , the first time the SGD iterates escaped from  $\Omega$ , so there is no need to consider an attempt that is after  $\sigma(\eta)$ , and we will not do so in the analysis below;
- Besides, the random variable  $\mathbf{j}_k$  is measurable w.r.t.  $\mathcal{F}_{\tilde{\tau}_k \wedge \sigma(\eta)}$ , the stopped  $\sigma$ -algebra generated by the stopping time  $\tilde{\tau}_k \wedge \sigma(\eta)$ .

Furthermore, for each  $k = 1, 2, \dots$ , let

$$\begin{aligned} T_{k,1}(\epsilon, \delta, \eta) &= \tau_{k-1}(\epsilon, \delta, \eta) \wedge \sigma(\eta), \\ T_{k,j}(\epsilon, \delta, \eta) &= \min\{n > T_{k,j-1}(\epsilon, \delta, \eta) : \eta|Z_n| > \delta\} \wedge \sigma(\eta) \wedge \tilde{\tau}_k \quad \forall j \geq 2, \\ W_{k,j}(\epsilon, \delta, \eta) &= Z_{T_{k,j}(\epsilon, \delta, \eta)} \quad \forall j \geq 1 \end{aligned}$$

with the convention  $T_{k,0}(\epsilon, \delta, \eta) = \tilde{\tau}_{k-1}(\epsilon, \delta, \eta)$ . Note that for any  $k \geq 1, j \geq 1$ ,  $T_{k,j}$  is a stopping time. Besides, from the definition of  $\mathbf{j}_k$ , one can see that

$$\tilde{\tau}_{k-1} + 1 \leq T_{k,j} \leq \tilde{\tau}_k \wedge \sigma(\eta) \quad \forall j \in [\mathbf{j}_k], \quad (\text{B.72})$$

and the sequences  $(T_{k,j})_{j=1}^{\mathbf{j}_k}$  and  $(W_{k,j})_{j=1}^{\mathbf{j}_k}$  are the arrival times and sizes of *large* jumps during the  $k$ -th attempt, respectively. Again, when there is no ambiguity we will drop the dependency on  $\epsilon, \delta, \eta$  and simply write  $T_{k,j}$  and  $W_{k,j}$ .

In order to prove Proposition B.22, we analyze the most likely scenario that the exit from  $\Omega$  would happen. Specifically, we will introduce a series of events with superscript  $\times$  or  $\circ$ , where  $\times$  indicates that the event is *atypical* or unlikely to happen and  $\circ$  means that it is a *typical* event and is likely to be observed before the first exit from the attraction field  $\Omega$ . Besides, the subscript  $k$  indicates that the event in discussion concerns the dynamics of SGD during the  $k$ -th attempt. Our goal is to show that for some event  $\mathbf{A}^\times(\epsilon, \delta, \eta)$  its probability becomes sufficiently small as learning rate  $\eta$  tends to 0, so the escape from  $\Omega$  almost always occurs in the manner described by  $(\mathbf{A}^\times(\epsilon, \delta, \eta))^c$ . In particular, the definition of this *atypical* scenario  $\mathbf{A}^\times$  involves the union of some *atypical* events  $\mathbf{A}_k^\times, \mathbf{B}_k^\times$  that occur in the  $k$ -th attempt. In other words, the intuition of  $\mathbf{A}^\times$  is that something *abnormal* happened during one of the attempts before the final exit.

Here is one more comment for the general naming convention of these events. Events with label  $\mathbf{A}$  often describe the “efforts” made in an attempt to get out of  $\Omega$  (such as large noises), while those with label  $\mathbf{B}$  concern how the SGD iterates return to  $[-2\epsilon, 2\epsilon]$  (and how this attempt fails). For instance,  $\mathbf{A}_k^\times$  discusses the unlikely scenario before  $T_{k,l^*}$ , the arrival of the  $l^*$ -th large jump in this attempt, while  $\mathbf{B}_k^\times$  in general discusses the abnormal cases after  $T_{k,l^*}$  and before the return to  $[-2\epsilon, 2\epsilon]$ . On the other hand,  $\mathbf{A}_k^\circ$  describes a successful escape during  $k$ -th attempt, while  $\mathbf{B}_k^\circ$  means that during this attempt the iterates return to without spending too much time.

Now we proceed and provide a formal definition and analysis of the aforementioned series of events. As building blocks, we inspect the process  $(X_n^\eta)_{n \geq 1}$  at a even finer granularity, and bound the probability of some events  $(\mathbf{A}_{k,i}^\times)_{i \geq 0}, (\mathbf{B}_{k,i}^\times)_{i \geq 1}$  detailing several cases that are *unlikely* to occur during the escape from or return to local minimum in the  $k$ -th attempt. First, for each  $k \geq 1$ , define the event

$$\mathbf{A}_{k,0}^\times(\epsilon, \delta, \eta) \triangleq \left\{ \exists i = 0, 1, \dots, l^* \wedge \mathbf{j}_k \text{ s.t. } \max_{j=T_{k,i}+1, \dots, (T_{k,i+1}-1) \wedge \tilde{\tau}_k \wedge \sigma(\eta)} \eta|Z_{T_{k,i}+1} + \dots + Z_j| > \tilde{\epsilon} \right\}. \quad (\text{B.73})$$

Intuitively speaking, the event characterizes the atypical scenario where, during the  $k$ -th attempt, there is some large fluctuations (compared to  $\tilde{\epsilon}$ ) between any of the first  $l^*$  large jumps (or the first  $\mathbf{j}_k$  large jumps in case that  $\mathbf{j}_k < l^*$ ). Similarly, consider event (for all  $k \geq 1$ )

$$\mathbf{A}_{k,1}^\times(\epsilon, \delta, \eta) \triangleq \left\{ \sigma(\eta) < \tilde{\tau}_k, \mathbf{j}_k < l^* \right\} \quad (\text{B.74})$$

that describes the atypical case where the exit occurs during the  $k$ -th attempt with less than  $l^*$  large jumps. Next, for all  $k \geq 1$  we have another atypical event (note that from (B.72) we can see that, for any  $j \geq 1$ ,  $\mathbf{j}_k \geq j$  implies  $T_{k,j} \leq \sigma(\eta) \wedge \tilde{\tau}_k$ )

$$\mathbf{A}_{k,2}^\times \triangleq \left\{ \mathbf{j}_k \geq l^*, \exists j = 2, 3, \dots, l^* \text{ s.t. } T_{k,j} - T_{k,j-1} > 2\hat{t}(\epsilon)/\eta \right\}. \quad (\text{B.75})$$

representing the case where we have at least  $l^*$  large noises during the  $k$ -th attempt, but for some of the large noise (from the 2nd to the  $l^*$ -th), the inter-arrival time is unusually long. Moving on, we consider the following events (defined for all  $k \geq 1$ )

$$\mathbf{A}_{k,3}^\times \triangleq \left\{ \mathbf{j}_k < l^*, \tilde{\tau}_k < \sigma(\eta), \tilde{\tau}_k - T_{k,1} > 2l^*\hat{t}(\epsilon)/\eta \right\} \quad (\text{B.76})$$

that describes the atypical case where the  $k$ -th attempt failed but the return to the small neighborhood  $[-2\epsilon, 2\epsilon]$  took unusually long time.

The following event also concerns the scenario where there are at least  $l^*$  large noises during the  $k$ -th attempt:

$$\mathbf{A}_{k,4}^\times \triangleq \left\{ \mathbf{j}_k \geq l^*, |X_{T_{k,l^*}}^\eta| \geq r - \bar{\epsilon}, \exists j = 1, 2, \dots, l^* \text{ s.t. } \eta|W_{k,j}| \leq \bar{\delta} \right\}; \quad (\text{B.77})$$

specifically, it describes the atypical case where, during this attempt, right after the  $l^*$ -large noise the SGD iterate is far enough from the local minimum yet some of the large noises are not that *large*. Lastly, by defining events

$$\mathbf{A}_{k,5}^\times \triangleq \left\{ \mathbf{j}_k \geq l^*, T_{k,l^*} \leq \sigma(\eta) \wedge \tilde{\tau}_k, X_{T_{k,l^*}}^\eta \in \bigcup_{j \in [n_{\min}-1]} [s_j - \epsilon, s_j + \epsilon] \right\}, \quad (\text{B.78})$$

we analyze an atypical case where the SGD iterates arrive at somewhere too close to the boundaries of  $\Omega$  at the arrival time of the  $l^*$  large noise during this attempt. As an amalgamation of these atypical scenarios, we let

$$\mathbf{A}_k^\times(\epsilon, \delta, \eta) \triangleq \bigcup_{i=0}^5 \mathbf{A}_{k,i}^\times(\epsilon, \delta, \eta). \quad (\text{B.79})$$

Also, we analyze the probability of some events  $(\mathbf{B}_k^\times)_{k \geq 1}$  that concern the SGD dynamics after the  $l^*$ -th large noise during the  $k$ -th attempt. Let us define

$$\begin{aligned} \mathbf{B}_{k,1}^\times(\epsilon, \delta, \eta) &\triangleq \left\{ \mathbf{j}_k \geq l^*, X_{T_{k,l^*}}^\eta \in [s_- + \epsilon, s_+ - \epsilon], T_{k,j} - T_{k,j-1} \leq 2\frac{\hat{t}(\epsilon)}{\eta} \forall j = 2, 3, \dots, l^* \right\} \\ \mathbf{B}_{k,2}^\times(\epsilon, \delta, \eta) &\triangleq \{\tilde{\tau}_k - T_{k,l^*} > \rho(\epsilon)/\eta\} \cup \{\sigma(\eta) < \tilde{\tau}_k\} \\ \mathbf{B}_k^\times(\epsilon, \delta, \eta) &\triangleq \mathbf{B}_{k,1}^\times \cap \mathbf{B}_{k,2}^\times \end{aligned} \quad (\text{B.80})$$

where  $\rho(\cdot)$  is the function in Lemma B.11. From the definition of  $\mathbf{B}_k^\times$ , in particular the inclusion of  $\mathbf{B}_{k,2}^\times$ , one can see that the intuitive interpretation of event  $\mathbf{B}_k^\times$  is that the SGD iterates *did not return* to local minimum efficiently (or simply escaped from the attraction field) after the  $l^*$ -th large noise during the  $k$ -th attempt. In comparison, the following events will characterize what would typically happen during each attempt:

$$\mathbf{A}_k^\circ(\epsilon, \delta, \eta) \triangleq \left\{ \mathbf{j}_k \geq l^*, \sigma(\eta) = T_{k,l^*}, X_{T_{k,l^*}}^\eta \notin [s_- - \epsilon, s_+ + \epsilon], T_{k,j} - T_{k,j-1} \leq \frac{2\hat{t}(\epsilon)}{\eta} \forall j = 2, 3, \dots, l^* \right\}, \quad (\text{B.81})$$

$$\mathbf{B}_k^\circ(\epsilon, \delta, \eta) \triangleq \left\{ \sigma(\eta) > \tilde{\tau}_k, \tilde{\tau}_k - T_{k,1} \leq \frac{2l^*\hat{t}(\epsilon) + \rho(\epsilon)}{\eta} \right\}. \quad (\text{B.82})$$

Intuitively speaking,  $\mathbf{A}_k^\circ$  tells us that the exit happened right at  $T_{k,l^*}$ , the arrival time of the  $l^*$ -th large noise during the  $k$ -th attempt, and  $\mathbf{B}_k^\circ$  tells us that the first exit from  $\Omega$  did not occur during the  $k$ -th attempt, and the SGD iterates returned to local minimum rather efficiently. All the preparations above allow use to define

$$\mathbf{A}^\times(\epsilon, \delta, \eta) \triangleq \bigcup_{k \geq 1} \left( \bigcap_{i=1}^{k-1} (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times \cup \mathbf{A}_i^\circ)^c \right) \cap (\mathbf{A}_k^\times \cup \mathbf{B}_k^\times). \quad (\text{B.83})$$

We need the next lemma in the proof of Proposition B.22. As mentioned earlier, the takeaway is that  $\mathbf{A}^\times$  is indeed *atypical* in the sense that we will almost always observe  $(\mathbf{A}^\times)^c$ .

**Lemma B.23.** *Given any  $C > 0$ , the following claim holds for all  $\epsilon > 0, \delta > 0$  sufficiently small:*

$$\limsup_{\eta \downarrow 0} \sup_{|x| \leq 2\epsilon} \mathbb{P}_x(\mathbf{A}^\times(\epsilon, \delta, \eta)) < C.$$

*Proof.* We fix some parameters for the proof. First, with out loss of generality we only consider  $C \in (0, 1)$ , and we fix some  $N > \alpha l^*$ . Next we discuss the valid range of  $\epsilon$  for the claim to hold. We only consider  $\epsilon > 0$  such that

$$\epsilon < \frac{\bar{\epsilon}}{6(\bar{\rho} + \tilde{\rho} + 3)} \wedge \frac{\epsilon_0}{3}$$

where  $\bar{\rho}$  and  $\tilde{\rho}$  are the constants in Corollary B.5 and Corollary B.8 respectively, and  $\epsilon_0$  is the constant in (B.1). Moreover, recall function  $\Psi$  in Lemma B.20 and the constant  $c_* > 0$  in Lemma B.21. Due to  $\lim_{\epsilon \downarrow 0} \Psi(\epsilon) = 0$ , it holds for all  $\epsilon$  small enough such that

$$\frac{3\Psi(\epsilon)}{c_*} < C \quad (\text{B.84})$$

In our proof we only consider  $\epsilon$  small enough so the inequality above holds, and the claim in Lemma B.11 holds. Now we specify the valid range of parameter  $\delta$  that will be used below:

- For all sufficiently small  $\delta > 0$ , the claim in Lemma B.12 will hold for the prescribed  $\epsilon$  and with  $N_0 = N$ ;
- For all sufficiently small  $\delta > 0$ , the claims in Lemma B.16, Corollary B.17, Corollary B.18 and Lemma B.19 will hold with the prescribed  $\epsilon$  and  $N$ ;
- For all sufficiently small  $\delta > 0$ , the inequalities in Lemma B.20 and B.21 will hold for the  $\epsilon$  we fixed at the beginning;
- Lastly, let

$$\begin{aligned} \tilde{\epsilon} &= \min \left\{ \frac{\epsilon}{2 \exp(2\bar{t}M)}, \frac{\epsilon}{4 \exp(2M\hat{t}(\epsilon))} \right\} \\ \hat{\epsilon} &= 3(\bar{\rho} + \tilde{\rho} + 3)\epsilon. \end{aligned}$$

where the function  $\hat{t}(\cdot)$  is defined in (B.24); due to Corollary B.13 we know that for all sufficiently small  $\delta > 0$  the claim (B.50) holds with the chosen  $\tilde{\epsilon}$  and  $N$ .

We show that the claim holds for any  $\epsilon, \delta$  small enough to satisfy the conditions above.

First, recall that

$$\mathbf{A}_{k,0}^\times(\epsilon, \delta, \eta) \triangleq \left\{ \exists i = 0, 1, \dots, l^* \wedge \mathbf{j}_k \text{ s.t. } \max_{j=T_{k,i+1}, \dots, (T_{k,i+1}-1) \wedge \tilde{\tau}_k \wedge \sigma(\eta)} \eta |Z_{T_{k,i+1}} + \dots + Z_j| > \tilde{\epsilon} \right\}.$$

Due to our choice of  $\delta$  stated earlier and Corollary B.13, there exists some  $\eta_0 > 0$  such that for all  $\eta \in (0, \eta_0)$ ,

$$\mathbb{P}(\mathbf{A}_{k,0}^\times(\epsilon, \delta, \eta)) \leq \eta^N \quad \forall k \geq 1. \quad (\text{B.85})$$

Similarly, recall that  $\mathbf{A}_{k,1}^\times(\epsilon, \delta, \eta) \triangleq \{\sigma(\eta) < \tilde{\tau}_k, \mathbf{j}_k < l^*\}$ . Let us temporarily focus on the first attempt (namely the case  $k = 1$ ). From Lemma B.16 and our choice of  $\epsilon$  and  $\delta$ , we know the existence of some  $\eta_1 > 0$  such that

$$\sup_{|x| \leq 2\epsilon} \mathbb{P}_x(\mathbf{A}_{1,1}^\times(\epsilon, \delta, \eta)) \leq \eta^N \quad \forall \eta \in (0, \eta_1). \quad (\text{B.86})$$

Next, for  $\mathbf{A}_{k,2}^\times \triangleq \{\mathbf{j}_k \geq l^*, \exists j = 2, 3, \dots, l^* \text{ s.t. } T_{k,j} - T_{k,j-1} > 2\hat{t}(\epsilon)/\eta\}$ , from Corollary B.17 and our choice of  $\delta$  at the beginning, we have the existence of some  $\eta_2 > 0$  such that

$$\sup_{|x| \leq 2\epsilon} \mathbb{P}_x(\mathbf{A}_{1,2}^\times(\epsilon, \delta, \eta)) \leq \eta^N \quad \forall \eta \in (0, \eta_2). \quad (\text{B.87})$$

Moving on, for  $\mathbf{A}_{k,3}^\times \triangleq \{\mathbf{j}_k < l^*, \tilde{\tau}_k < \sigma(\eta), \tilde{\tau}_k - T_{k,1} > 2l^*\hat{t}(\epsilon)/\eta\}$ , due to Corollary B.18 and our choice of  $\epsilon, \delta$ , we have the existence of some  $\eta_3 > 0$  such that

$$\sup_{|x| \leq 2\epsilon} \mathbb{P}_x(\mathbf{A}_{1,3}^\times(\epsilon, \delta, \eta)) \leq \eta^N \quad \forall \eta \in (0, \eta_3). \quad (\text{B.88})$$

As for  $\mathbf{A}_{k,4}^\times \triangleq \{\mathbf{j}_k \geq l^*, |X_{T_{k,l^*}}^\eta| \geq r - \bar{\epsilon}, \exists j = 1, 2, \dots, l^* \text{ s.t. } \eta|W_{k,j}| \leq \bar{\delta}\}$ , from Lemma B.19, one can see the existence of  $\eta_4 > 0$  such that

$$\sup_{|x| \leq 2\epsilon} \mathbb{P}_x(\mathbf{A}_{1,4}^\times(\epsilon, \delta, \eta)) \leq \eta^N \quad \forall \eta \in (0, \eta_4). \quad (\text{B.89})$$

Lastly, for  $\mathbf{A}_{k,5}^\times \triangleq \{\mathbf{j}_k \geq l^*, T_{k,l^*} \leq \sigma(\eta) \wedge \tilde{\tau}_k, X_{T_{k,l^*}}^\eta \in \bigcup_{j \in [n_{\min}-1]} [s_j - \epsilon, s_j + \epsilon]\}$ , from Lemma B.20 we see the existence of  $\eta_5 > 0$  such that

$$\sup_{|x| \leq 2\epsilon} \mathbb{P}_x(\mathbf{A}_{1,5}^\times(\epsilon, \delta, \eta)) \leq 2\delta^\alpha \Psi(\epsilon) \left( H(1/\eta)/\eta \right)^{l^*-1} \quad \forall \eta \in (0, \eta_5). \quad (\text{B.90})$$

Recall that  $\mathbf{A}_k^\times(\epsilon, \delta, \eta) = \bigcup_{i=0}^5 \mathbf{A}_{k,i}^\times(\epsilon, \delta, \eta)$ . Also, for definitions of  $\mathbf{B}_k^\times, \mathbf{A}_k^\circ, \mathbf{B}_k^\circ$ , see (B.80), (B.81), (B.82) respectively. Our next goal is to establish bounds regarding the probabilities of these events. First, if we consider the event  $\bigcap_{j=1}^k (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ$ , then the inclusion of the  $(\mathbf{B}_j^\circ)_{j=1}^k$  implies that during the first  $k$  attempts the SGD iterates have never left the attraction field, so

$$\bigcap_{j=1}^k (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ = \left( \bigcap_{j=1}^k (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ \right) \cap \{\sigma(\eta) > \tilde{\tau}_k\}.$$

Next, note that

$$\begin{aligned} & \mathbb{P}_x(\mathbf{B}_k^\times \mid \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \\ &= \mathbb{P}_x(\mathbf{B}_{k,1}^\times \mid \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \mathbb{P}_x(\mathbf{B}_{k,2}^\times \mid (\bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \cap \mathbf{B}_{k,1}^\times) \end{aligned}$$

$$\leq \mathbb{P}_x(\mathbf{j}_k \geq l^*, T_{k,j} - T_{k,j-1} \leq \frac{2\hat{t}(\epsilon)}{\eta} \forall j = 2, 3, \dots, l^* \mid \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \\ \cdot \mathbb{P}_x(\mathbf{B}_{k,2}^\times \mid (\bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \cap \mathbf{B}_{k,1}^\times).$$

From the definition of the events  $\mathbf{A}_j^\times, \mathbf{B}_j^\times, \mathbf{B}_j^\circ$ , one can see that  $\bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ \in \mathcal{F}_{\tilde{\tau}_{k-1} \wedge \sigma(\eta)}$ , and on this event  $\bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ$  we have  $\sigma(\eta) > \tilde{\tau}_{k-1}$ . So by applying strong Markov property at stopping time  $\tilde{\tau}_{k-1} \wedge \sigma(\eta)$ , we have

$$\mathbb{P}_x(\mathbf{B}_k^\times \mid \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \\ \leq \mathbb{P}\left(T_j^\eta(\delta) - T_{j-1}^\eta(\delta) \leq 2\hat{t}(\epsilon)/\eta \forall j \in [l^* - 1]\right) \cdot \mathbb{P}_x(\mathbf{B}_{k,2}^\times \mid (\bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \cap \mathbf{B}_{k,1}^\times) \\ \leq 2\left(H(\delta/\eta)\hat{t}(\epsilon)/\eta\right)^{l^*-1} \cdot \mathbb{P}_x(\mathbf{B}_{k,2}^\times \mid (\bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \cap \mathbf{B}_{k,1}^\times) \quad \text{for all } \eta \text{ sufficiently small due to Lemma B.2,} \\ \leq 4\left(\frac{\hat{t}(\epsilon)}{\delta^\alpha}\right)^{l^*-1} \left(\frac{H(1/\eta)}{\eta}\right)^{l^*-1} \cdot \mathbb{P}_x(\mathbf{B}_{k,2}^\times \mid (\bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \cap \mathbf{B}_{k,1}^\times)$$

for all sufficiently small  $\eta$ , due to  $H \in \mathcal{RV}_{-\alpha}(\eta)$ . Meanwhile, note that

- $(\bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \cap \mathbf{B}_{k,1}^\times \in \mathcal{F}_{T_{k,l^*}}$ ;
- on this event  $(\bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \cap \mathbf{B}_{k,1}^\times$  we have  $\sigma(\eta) \wedge \tilde{\tau}_k > T_{k,l^*}$  and  $X_{T_{k,l^*}}^\eta \in [s_- + \epsilon, s_+ - \epsilon]$ .

Therefore, using Lemma B.11 and strong Markov property again (at stopping time  $T_{k,l^*}$ ), we know the following inequality holds for all  $\eta$  sufficiently small:

$$\sup_{k \geq 1} \sup_{|x| \leq 2\epsilon} \mathbb{P}_x(\mathbf{B}_{k,2}^\times \mid (\bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \cap \mathbf{B}_{k,1}^\times) \\ = \sup_{k \geq 1} \sup_{|x| \leq 2\epsilon} \mathbb{P}_x(\{\sigma(\eta) > \tilde{\tau}_k, \tilde{\tau}_k - T_{k,l^*} \leq \rho(\epsilon)/\eta\}^c \mid (\bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \cap \mathbf{B}_{k,1}^\times) \\ \leq \Psi(\epsilon) \frac{\delta^\alpha}{4(\hat{t}(\epsilon)/\delta^\alpha)^{l^*-1}}.$$

Therefore, we know the existence of some  $\eta_6 > 0$  such that

$$\sup_{|x| \leq 2\epsilon} \mathbb{P}_x(\mathbf{B}_k^\times \mid \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \leq \Psi(\epsilon) \delta^\alpha \left(\frac{H(1/\eta)}{\eta}\right)^{l^*-1} \quad \forall \eta \in (0, \eta_6), \forall k \geq 1. \quad (\text{B.91})$$

Similarly, we can bound conditional probabilities of the form  $\mathbb{P}_x(\mathbf{A}_k^\times \mid \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ)$ . To be specific, recall that  $\mathbf{A}_k^\times = \bigcup_{i=0}^5 \mathbf{A}_{k,i}^\times$ . By combining (B.85)-(B.90) with Markov property, we know the existence of some  $\eta_7 > 0$  such that

$$\sup_{|x| \leq 2\epsilon} \mathbb{P}_x(\mathbf{A}_k^\times \mid \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \leq 5\eta^N + 2\Psi(\epsilon) \delta^\alpha \left(\frac{H(1/\eta)}{\eta}\right)^{l^*-1} \quad \forall \eta \in (0, \eta_7), \forall k \geq 1. \quad (\text{B.92})$$

On the other hand, a lower bound can be established for conditional probability involving  $\mathbf{A}_k^\circ$ , the event defined in (B.81) describing the exit from  $\Omega$  during an attempt with exactly  $l^*$  large noises. Using Lemma B.21 and Markov property of  $(X_n^\eta)_{n \geq 1}$ , one can see the existence of some  $\eta_8 > 0$  such that

$$\inf_{|x| \leq 2\epsilon} \mathbb{P}_x(\mathbf{A}_k^\circ \mid \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ) \geq c_* \delta^\alpha \left( \frac{H(1/\eta)}{\eta} \right)^{l^*-1} \quad \forall \eta \in (0, \eta_8). \quad (\text{B.93})$$

In order to apply the bounds (B.91)-(B.93), we make use of the following inclusion relationship:

$$\left( \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ \right) \cap (\mathbf{A}_k^\times \cup \mathbf{B}_k^\times)^c \subseteq \mathbf{A}_k^\circ \cup \mathbf{B}_k^\circ. \quad (\text{B.94})$$

To see why this is true, let us consider a decomposition of the event on the L.H.S. of (B.94). As mentioned above, on event  $\bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ$  we know that  $\sigma(\eta) > \tilde{\tau}_{k-1}$ , so the  $k$ -th attempt occurred and there are only three possibilities on this event:

- $\mathbf{j}_k < l^*$ ;
- $\mathbf{j}_k \geq l^*$ ,  $X_{T_{k,l^*}}^\eta \notin \Omega$ ;
- $\mathbf{j}_k \geq l^*$ ,  $X_{T_{k,l^*}}^\eta \in \Omega$ .

Let us partition the said event accordingly and analyze them one by one.

- On  $\left( \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ \right) \cap (\mathbf{A}_k^\times \cup \mathbf{B}_k^\times)^c \cap \{\mathbf{j}_k < l^*\}$ , due to the exclusion of  $\mathbf{A}_k^\times$  (especially  $\mathbf{A}_{k,1}^\times$  and  $\mathbf{A}_{k,3}^\times$ ), we can see that if  $\mathbf{j}_k < l^*$ , then we must have  $\sigma(\eta) > \tilde{\tau}_k$  and  $\tilde{\tau}_k - T_{k,1} \leq 2l^* \hat{t}(\epsilon)/\eta$ . Therefore,

$$\left( \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ \right) \cap (\mathbf{A}_k^\times \cup \mathbf{B}_k^\times)^c \cap \{\mathbf{j}_k < l^*\} \subseteq \mathbf{B}_k^\circ.$$

- On  $\left( \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ \right) \cap (\mathbf{A}_k^\times \cup \mathbf{B}_k^\times)^c \cap \{\mathbf{j}_k \geq l^*, X_{T_{k,l^*}}^\eta \notin \Omega\}$ , then the exclusion of  $\mathbf{A}_{k,2}^\times$  implies that  $T_{k,j} - T_{k,j-1} \leq 2\hat{t}(\epsilon)/\eta$  for all  $j = 2, \dots, l^*$ , and the exclusion of  $\mathbf{A}_{k,5}^\times$  tells us that if  $X_{T_{k,l^*}}^\eta \notin \Omega$ , then we have  $X_{T_{k,l^*}}^\eta \notin [s_- - \epsilon, s_+ + \epsilon]$ . In summary,

$$\left( \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ \right) \cap (\mathbf{A}_k^\times \cup \mathbf{B}_k^\times)^c \cap \{\mathbf{j}_k \geq l^*, X_{T_{k,l^*}}^\eta \notin \Omega\} \subseteq \mathbf{A}_k^\circ.$$

- On  $\left( \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ \right) \cap (\mathbf{A}_k^\times \cup \mathbf{B}_k^\times)^c \cap \{\mathbf{j}_k \geq l^*, X_{T_{k,l^*}}^\eta \in \Omega\}$ , the exclusion of  $\mathbf{A}_{k,2}^\times$  again implies that  $T_{k,j} - T_{k,j-1} \leq 2\hat{t}(\epsilon)/\eta$  for all  $j = 2, \dots, l^*$ , hence  $T_{k,l^*} - T_{k,1} \leq 2l^* \hat{t}(\epsilon)/\eta$ . Similarly, the exclusion of  $\mathbf{A}_{k,5}^\times$  tells us that if  $X_{T_{k,l^*}}^\eta \in \Omega$ , then we have  $X_{T_{k,l^*}}^\eta \in [s_- + \epsilon, s_+ - \epsilon]$ . Now since  $\mathbf{B}_k^\times$  did not occur (see the definition in (B.80)), we must have  $\sigma(\eta) > \tilde{\tau}_k$  and  $\tilde{\tau}_k - T_{k,l^*} \leq \rho(\epsilon)/\eta$ , hence  $\tilde{\tau}_k - T_{k,1} \leq \frac{2l^* \hat{t}(\epsilon) + \rho(\epsilon)}{\eta}$ . Therefore,

$$\left( \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ \right) \cap (\mathbf{A}_k^\times \cup \mathbf{B}_k^\times)^c \cap \{\mathbf{j}_k \geq l^*, X_{T_{k,l^*}}^\eta \in \Omega\} \subseteq \mathbf{B}_k^\circ.$$

Collecting results above, we have (B.94). Now we discuss some of its implications. First, from (B.94) we can immediately get that

$$\left( \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ \right) \cap (\mathbf{A}_k^\times \cup \mathbf{B}_k^\times)^c = \left( \bigcap_{j=1}^{k-1} (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times)^c \cap \mathbf{B}_j^\circ \right) \cap (\mathbf{A}_k^\times \cup \mathbf{B}_k^\times)^c \cap (\mathbf{A}_k^\circ \cup \mathbf{B}_k^\circ). \quad (\text{B.95})$$

Next, recall the definitions of  $\mathbf{A}_k^\circ$  in (B.81) and  $\mathbf{B}_k^\circ$  in (B.82), and one can see that  $\mathbf{A}_k^\circ$  and  $\mathbf{B}_k^\circ$  are mutually exclusive, since the former implies that the first exit occurs during the  $k$ -th attempt while the latter implies that this attempt fails. This fact and (B.95) allow us to conclude that

$$\bigcap_{i=1}^k (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times \cup \mathbf{A}_i^\circ)^c = \bigcap_{i=1}^k (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times)^c \cap \mathbf{B}_i^\circ = \left( \bigcap_{i=1}^{k-1} (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times)^c \cap \mathbf{B}_i^\circ \right) \cap (\mathbf{A}_k^\times \cup \mathbf{B}_k^\times \cup \mathbf{A}_k^\circ)^c. \quad (\text{B.96})$$

Now we use the results obtained so far to bound the probability of

$$\mathbf{A}^\times(\epsilon, \delta, \eta) \triangleq \bigcup_{k \geq 1} \left( \bigcap_{i=1}^{k-1} (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times \cup \mathbf{A}_i^\circ)^c \right) \cap (\mathbf{A}_k^\times \cup \mathbf{B}_k^\times).$$

Using (B.96), we can see that (for any  $x \in [-2\epsilon, 2\epsilon]$ )

$$\begin{aligned} & \mathbb{P}_x(\mathbf{A}^\times(\epsilon, \delta, \eta)) \\ &= \sum_{k \geq 1} \mathbb{P}_x \left( \left( \bigcap_{i=1}^{k-1} (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times \cup \mathbf{A}_i^\circ)^c \right) \cap (\mathbf{A}_k^\times \cup \mathbf{B}_k^\times) \right) \\ &= \sum_{k \geq 1} \mathbb{P}_x \left( \left( \bigcap_{i=1}^{k-1} (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times)^c \cap \mathbf{B}_i^\circ \right) \cap (\mathbf{A}_k^\times \cup \mathbf{B}_k^\times) \right) \\ &= \sum_{k \geq 1} \mathbb{P}_x \left( \mathbf{A}_k^\times \cup \mathbf{B}_k^\times \mid \bigcap_{i=1}^{k-1} (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times)^c \cap \mathbf{B}_i^\circ \right) \cdot \prod_{j=1}^{k-1} \mathbb{P}_x \left( \bigcap_{i=1}^j (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times)^c \cap \mathbf{B}_i^\circ \mid \bigcap_{i=1}^{j-1} (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times)^c \cap \mathbf{B}_i^\circ \right) \\ &= \sum_{k \geq 1} \mathbb{P}_x \left( \mathbf{A}_k^\times \cup \mathbf{B}_k^\times \mid \bigcap_{i=1}^{k-1} (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times)^c \cap \mathbf{B}_i^\circ \right) \\ & \quad \cdot \prod_{j=1}^{k-1} \mathbb{P}_x \left( \left( \bigcap_{i=1}^{j-1} (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times)^c \cap \mathbf{B}_i^\circ \right) \cap (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times \cup \mathbf{A}_j^\circ)^c \mid \bigcap_{i=1}^{j-1} (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times)^c \cap \mathbf{B}_i^\circ \right) \\ &= \sum_{k \geq 1} \mathbb{P}_x \left( \mathbf{A}_k^\times \cup \mathbf{B}_k^\times \mid \bigcap_{i=1}^{k-1} (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times)^c \cap \mathbf{B}_i^\circ \right) \cdot \prod_{j=1}^{k-1} \mathbb{P}_x \left( (\mathbf{A}_j^\times \cup \mathbf{B}_j^\times \cup \mathbf{A}_j^\circ)^c \mid \bigcap_{i=1}^{j-1} (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times)^c \cap \mathbf{B}_i^\circ \right) \\ &\leq \sum_{k \geq 1} \mathbb{P}_x \left( \mathbf{A}_k^\times \cup \mathbf{B}_k^\times \mid \bigcap_{i=1}^{k-1} (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times)^c \cap \mathbf{B}_i^\circ \right) \cdot \prod_{j=1}^{k-1} \left( 1 - \mathbb{P}_x \left( \mathbf{A}_j^\circ \mid \bigcap_{i=1}^{j-1} (\mathbf{A}_i^\times \cup \mathbf{B}_i^\times)^c \cap \mathbf{B}_i^\circ \right) \right). \end{aligned}$$

This allows us to apply (B.91)-(B.93) and conclude that (here we only consider  $\eta < \min\{\eta_i : i \in [8]\}$ ),

$$\begin{aligned} & \sup_{|x| \leq 2\epsilon} \mathbb{P}_x(\mathbf{A}^\times(\epsilon, \delta, \eta)) \\ &\leq \sum_{k \geq 1} \left( 5\eta^N + 2\Psi(\epsilon)\delta^\alpha \left( \frac{H(1/\eta)}{\eta} \right)^{l^*-1} \right) \cdot \left( 1 - c_*\delta^\alpha \left( \frac{H(1/\eta)}{\eta} \right)^{l^*-1} \right)^{k-1} \end{aligned}$$

$$\begin{aligned}
&= \frac{5\eta^N + 2\Psi(\epsilon)\delta^\alpha \left(\frac{H(1/\eta)}{\eta}\right)^{l^*-1}}{c_*\delta^\alpha \left(\frac{H(1/\eta)}{\eta}\right)^{l^*-1}} \\
&\leq \frac{2\Psi(\epsilon) + 5\eta^\alpha}{c_*} \quad \text{for sufficiently small } \eta, \text{ due to } H \in \mathcal{RV}_{-\alpha}(\eta) \text{ and our choice of } N > \alpha l^* \\
&\leq \frac{3\Psi(\epsilon)}{c_*} < C \quad \text{for all } \eta \text{ small enough such that } 5\eta^\alpha < \Psi(\epsilon).
\end{aligned}$$

The last inequality follows from our choice of  $\epsilon$  in (B.84). This concludes the proof.  $\square$

Having established Lemma B.23, we return to Proposition B.22 and give a proof. Recall that, aside from the attraction field  $\Omega = (s_-, s_+)$ , there are  $n_{\min} - 1$  other attraction fields  $\tilde{\Omega}_k = (s_k^-, s_k^+)$  (for each  $k \in [n_{\min} - 1]$ ). Besides, the function  $\lambda(\cdot)$  and constants  $\nu^\Omega, \nu_k^\Omega$  are defined in (B.64)-(B.66).

*Proof of Proposition B.22.* We fix some parameters for the proof. First, with out loss of generality we only need to consider  $C \in (0, 1)$ . Next we discuss the valid range of  $\epsilon$  for the claim to hold. We only consider  $\epsilon > 0$  such that

$$\epsilon < \frac{\bar{\epsilon}}{6(\bar{\rho} + \tilde{\rho} + 3)} \wedge \frac{\epsilon_0}{3}$$

where  $\bar{\rho}$  and  $\tilde{\rho}$  are the constants in Corollary B.5 and Corollary B.8 respectively, and  $\epsilon_0$  is the constant in (B.1). Due to continuity of measure  $\mu$ , it holds for all  $\epsilon$  small enough such that (let  $\hat{\epsilon} = 3(\bar{\rho} + \tilde{\rho} + 3)\epsilon$ )

$$\frac{\mu(E(0))}{\mu(E(\hat{\epsilon}))} < 1/(1 - C), \quad (\text{B.97})$$

$$\frac{\mu(E(0))}{\mu(E(-\hat{\epsilon}))} > 1/(1 + C), \quad (\text{B.98})$$

$$\frac{\mu\left(h^{-1}((s_- - 2\hat{\epsilon}, s_- + 2\hat{\epsilon}) \cup (s_+ - 2\hat{\epsilon}, s_+ + 2\hat{\epsilon}))\right)}{\mu(E(-\hat{\epsilon}))} \leq C \quad (\text{B.99})$$

$$\frac{\mu\left(E(\hat{\epsilon}) \cap (s_{k'}^- - \hat{\epsilon}, s_{k'}^+ + \hat{\epsilon})\right)}{\mu(E(\hat{\epsilon}))} \leq \frac{\nu_{k'}^\Omega + C}{\nu^\Omega} \quad (\text{B.100})$$

$$\frac{\mu\left(E(-\hat{\epsilon}) \cap (s_{k'}^- + 2\hat{\epsilon}, s_{k'}^+ - 2\hat{\epsilon})\right)}{\mu(E(-\hat{\epsilon}))} \geq \frac{\nu_{k'}^\Omega - C}{\nu^\Omega} \quad (\text{B.101})$$

In our proof we only consider  $\epsilon$  small enough so the inequality above holds, and the claims in Lemma B.11 hold. Moreover, we only consider  $\epsilon$  and  $\delta$  small enough so that Lemma B.23 hold and we have

$$\lim_{\eta \downarrow 0} \sup_{|x| \leq 2\epsilon} \mathbb{P}_x(\mathbf{A}^\times(\epsilon, \delta, \eta)) < C. \quad (\text{B.102})$$

We show that the desired claims hold for all  $\epsilon, \delta$  sufficiently small that satisfy conditions above.

First, in order to show (B.71), we define event

$$\begin{aligned}
&\tilde{\mathbf{A}}^\times(\epsilon, \delta, \eta) \\
&\triangleq (\mathbf{A}^\times(\epsilon, \delta, \eta))^c \cap \left\{ \lambda(\eta)(\tau(\eta, \epsilon) - \sigma(\eta)) \geq C \text{ or } \exists n = \sigma(\eta) + 1, \dots, \tau(\eta, \epsilon) \text{ such that } X_n^\eta \notin \tilde{\Omega}_{J_\sigma(\eta)} \right\}
\end{aligned}$$

Since  $\lambda \in \mathcal{RV}_{-1-l^*(\alpha-1)}$  and  $\alpha > 1$ , for the  $\epsilon$  we fixed at the beginning of this proof,  $\rho(\epsilon)$  is a fixed constant as well (the function  $\rho$  is defined in Lemma B.11) and we have  $\lim_{\eta \downarrow 0} \lambda(\eta)\rho(\epsilon)/\eta = 0$ . Next, the occurrence of  $(\mathbf{A}^\times(\epsilon, \delta, \eta))^c$  (in particular, the exclusion of all the  $\mathbf{A}_{k,5}^\times$  defined in (B.78)),



we know that  $X_{\sigma(\eta)}^\eta \notin [s_{J_\sigma}^- - \epsilon, s_{J_\sigma}^- + \epsilon] \cup [s_{J_\sigma}^+ - \epsilon, s_{J_\sigma}^+ + \epsilon]$  (recall that for any  $k \in [n_{\min} - 1]$ , we have  $\tilde{\Omega}_j = (s_j^-, s_j^+)$ ; for definition of  $J_\sigma$  see (B.63)). Meanwhile, for all  $\eta$  sufficiently small, we have  $\epsilon/\lambda(\eta) > \rho(\epsilon)/\eta$ . Therefore, using Lemma B.11 we can see that (for all  $\eta$  sufficiently small)

$$\sup_{|x| \leq 2\epsilon} \mathbb{P}_x \left( \tilde{\mathbf{A}}^\times(\epsilon, \delta, \eta) \mid (\mathbf{A}^\times(\epsilon, \delta, \eta))^c \right) \leq C. \quad (\text{B.103})$$

Lastly, observe that

$$\begin{aligned} & \mathbb{P} \left( \left\{ \lambda(\eta)(\tau(\eta, \epsilon) - \sigma(\eta)) \geq C \text{ or } \exists n = \sigma(\eta) + 1, \dots, \tau(\eta, \epsilon) \text{ such that } X_n^\eta \notin \tilde{\Omega}_{J_\sigma(\eta)} \right\} \right) \\ & \leq \mathbb{P}_x \left( (\mathbf{A}^\times)^c \cap \left\{ \lambda(\eta)(\tau(\eta, \epsilon) - \sigma(\eta)) \geq C \text{ or } \exists n = \sigma(\eta) + 1, \dots, \tau(\eta, \epsilon) \text{ such that } X_n^\eta \notin \tilde{\Omega}_{J_\sigma(\eta)} \right\} \right) \\ & + \mathbb{P}_x(\mathbf{A}^\times) \end{aligned}$$

so by combining (B.102) with (B.103), we can obtain (B.71).

Moving on, we discuss the upper bounds (B.67) and (B.69). Recall that the fixed constant  $k' \in [n_{\min} - 1]$  is prescribed in the description of this proposition. Let us observe some facts on event  $(\mathbf{A}^\times(\epsilon, \eta, \delta))^c \cap \{J_\sigma(\eta) = k'\}$ : If we let  $J(\epsilon, \delta, \eta) \triangleq \sup\{k \geq 0 : \tilde{\tau}_k < \sigma(\eta)\}$  be the number of attempts it took to escape, and

$$J^\uparrow(\epsilon, \delta, \eta) \triangleq \min\{k \geq 1 : T_{k,1} \text{ has } (3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, \delta, \eta)\text{-overflow}\},$$

then for all  $\eta$  sufficiently small, we must have  $J \leq J^\uparrow$  on event  $(\mathbf{A}^\times(\epsilon, \eta, \delta))^c \cap \{J_\sigma(\eta) = k'\}$ . To see this via a proof of contradiction, let us assume that, for some arbitrary positive integer  $j$ , there exists some sample path on  $(\mathbf{A}^\times)^c \cap \{J_\sigma(\eta) = k'\}$  such that  $J^\uparrow = j < J$ . Then from the definition of  $(\mathbf{A}^\times)^c$ , in particular the exclusion of event  $\mathbf{A}_{j,0}^\times$  (see the definition in (B.73)), for all sufficiently small  $\eta$ , we are able to apply Corollary B.8 and B.5 and conclude that  $X_{T_{j,l^*}}^\eta \notin \Omega$ : indeed, using Corollary B.8 and B.5 we can show that the distance between  $X_{T_{j,l^*}}^\eta$  and the perturbed ODE

$$\tilde{\mathbf{x}}^\eta \left( T_{j,l^*} - T_{j,1}, 0; (0, T_{j,2} - T_{j,1}, \dots, T_{j,l^*} - T_{j,1}), (\eta W_{j,1}, \dots, \eta W_{j,l^*}) \right)$$

is strictly less than  $3(\bar{\rho} + \tilde{\rho} + 3)\epsilon$ ; on the other hand, the definition of  $(3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, \delta, \eta)\text{-overflow}$  implies that

$$\begin{aligned} & \tilde{\mathbf{x}}^\eta \left( T_{j,l^*} - T_{j,1}, 0; (0, T_{j,2} - T_{j,1}, \dots, T_{j,l^*} - T_{j,1}), (\eta W_{j,1}, \dots, \eta W_{j,l^*}) \right) \\ & \notin [s_- - 3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, s_+ + 3(\bar{\rho} + \tilde{\rho} + 3)\epsilon]. \end{aligned}$$

Therefore, we must have  $X_{T_{j,l^*}}^\eta \notin \Omega$ , which contradicts our assumption  $j = J^\uparrow < J$ . In summary, we have shown that, on  $(\mathbf{A}^\times)^c \cap \{J_\sigma = k'\}$ , we have  $J^\uparrow(\epsilon, \delta, \eta) \geq J(\epsilon, \delta, \eta)$ . Similarly, if we consider

$$J^\downarrow(\epsilon, \delta, \eta) \triangleq \min\{k \geq 1 : T_{k,1} \text{ has } (-3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, \delta, \eta)\text{-overflow}\},$$

then by the same argument above we can show that  $J^\downarrow(\epsilon, \delta, \eta) \leq J(\epsilon, \delta, \eta)$ . Now consider the following decomposition of events.

- On  $\{J^\downarrow < J^\uparrow\}$ , we know that for the first  $k$  such that  $T_{k,1}$  has  $(-3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, \delta, \eta)\text{-overflow}$ , it does not have  $(3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, \delta, \eta)\text{-overflow}$ . Now we analyze the probability that  $Z_0$  does not have  $(\hat{\epsilon}, \delta, \eta)\text{-overflow}$  conditioning on that it does have  $(-\hat{\epsilon}, \delta, \eta)\text{-overflow}$  (recall that we let  $\hat{\epsilon} = 3(\bar{\rho} + \tilde{\rho} + 3)$ ). Using Lemma B.14 and the bound (B.99), we know that for all  $\eta$  sufficiently small,

$$\sup_{|x| \leq 2\epsilon} \mathbb{P}_x \left( (\mathbf{A}^\times)^c \cap \{J^\downarrow < J^\uparrow\} \right) \leq \frac{\mu \left( h^{-1}((s_- - 2\hat{\epsilon}, s_- + 2\hat{\epsilon}) \cup (s_+ - 2\hat{\epsilon}, s_+ + 2\hat{\epsilon})) \right)}{\mu(E(-\hat{\epsilon}))} \leq C. \quad (\text{B.104})$$

- On  $(\mathbf{A}^\times)^c \cap \{J_\sigma = k'\} \cap \{J^\uparrow = J^\downarrow\}$ , due to  $J^\uparrow = J^\downarrow = J$  we know that  $T_{J(\epsilon, \delta, \eta), 1}$  is the first among all  $T_{k, 1}$  to have  $(\hat{\epsilon}, \delta, \eta)$ -overflow. Moreover, due to  $\{J_\sigma = k'\}$  and using Corollary B.8 and B.5 again as we did above, we know that the overflow endpoint of  $T_{J(\epsilon, \delta, \eta), 1}$  is in  $(s_{k'}^- - \hat{\epsilon}, s_{k'}^+ + \hat{\epsilon})$  (recall that  $\tilde{\Omega}_{k'} = (s_{k'}^-, s_{k'}^+)$ ). In summary, for any  $n \geq 0$

$$\begin{aligned} & (\mathbf{A}^\times)^c \cap \{J_\sigma = k'\} \cap \{J^\uparrow = J^\downarrow > n\} \\ & \subseteq (\mathbf{A}^\times)^c \cap \{J^\uparrow > n\} \cap \left\{ T_{J^\uparrow, 1} \text{ has overflow endpoint in } (s_{k'}^- - \hat{\epsilon}, s_{k'}^+ + \hat{\epsilon}) \right\} \end{aligned}$$

so using Lemma B.14, we obtain that (for all  $\eta$  sufficiently small)

$$\begin{aligned} & \sup_{|x| \leq 2\epsilon} \mathbb{P}_x \left( (\mathbf{A}^\times)^c \cap \{J_\sigma = k'\} \cap \{J^\uparrow = J^\downarrow > n\} \right) \\ & \leq \sup_{|x| \leq 2\epsilon} \mathbb{P}_x \left( (\mathbf{A}^\times)^c \cap \{J^\uparrow > n\} \right) \cdot \frac{p(\hat{\epsilon}, \delta, \eta; (s_{k'}^- - \hat{\epsilon}, s_{k'}^+ + \hat{\epsilon}))}{p(\hat{\epsilon}, \delta, \eta)} \\ & \leq \sup_{|x| \leq 2\epsilon} \mathbb{P}_x \left( (\mathbf{A}^\times)^c \cap \{J^\uparrow > n\} \right) \cdot \frac{\nu_{k'}^\Omega + C}{\nu^\Omega}. \end{aligned} \tag{B.105}$$

uniformly for any  $n = 0, 1, 2, \dots$  due to (B.100).

- On the other hand, on  $(\mathbf{A}^\times)^c$ , if  $T_{J^\downarrow, 1}$  has overflow endpoint in  $(s_{k'}^- + 2\hat{\epsilon}, s_{k'}^+ - 2\hat{\epsilon})$ , then from Definition B.1 we know that  $T_{J^\downarrow, 1}$  also has  $(\hat{\epsilon}, \delta, \eta)$ -overflow, hence  $J^\downarrow = J^\uparrow = J$ . Moreover, using Corollary B.8 and B.5 again, we know that  $X_{T_{J^\downarrow, 1}^*}^\eta \in (s_{k'}^-, s_{k'}^+)$  so  $J_\sigma = k'$ . In summary, for any  $n \geq 0$ ,

$$\begin{aligned} & (\mathbf{A}^\times)^c \cap \{J_\sigma = k'\} \cap \{J^\uparrow = J^\downarrow > n\} \\ & \supseteq (\mathbf{A}^\times)^c \cap \{J^\downarrow > n\} \cap \left\{ T_{J^\downarrow, 1} \text{ has overflow endpoint in } (s_{k'}^- + 2\hat{\epsilon}, s_{k'}^+ - 2\hat{\epsilon}) \right\} \end{aligned}$$

so using Lemma B.14, we obtain that (for all  $\eta$  sufficiently small)

$$\begin{aligned} & \inf_{|x| \leq 2\epsilon} \mathbb{P}_x \left( (\mathbf{A}^\times)^c \cap \{J_\sigma = k'\} \cap \{J^\uparrow = J^\downarrow > n\} \right) \\ & \geq \inf_{|x| \leq 2\epsilon} \mathbb{P}_x \left( (\mathbf{A}^\times)^c \cap \{J^\downarrow > n\} \right) \cdot \frac{p(-\hat{\epsilon}, \delta, \eta; (s_{k'}^- + 2\hat{\epsilon}, s_{k'}^+ - 2\hat{\epsilon}))}{p(-\hat{\epsilon}, \delta, \eta)} \\ & \geq \inf_{|x| \leq 2\epsilon} \mathbb{P}_x \left( (\mathbf{A}^\times)^c \cap \{J^\downarrow > n\} \right) \cdot \frac{\nu_{k'}^\Omega - C}{\nu^\Omega}. \end{aligned} \tag{B.106}$$

uniformly for any  $n = 0, 1, 2, \dots$  due to (B.101).

Besides, the following claim holds on event  $(\mathbf{A}^\times)^c$ .

- From (B.96), the definition of  $\mathbf{B}_k^\circ$  as well as the definition of event  $\mathbf{A}_k^\circ$  (see (B.81)), one can see that for any  $j = 1, 2, \dots, J$ , we have

$$\tilde{\tau}_j \wedge \sigma(\eta) - T_{j, 1} \leq \frac{2l^* \hat{t}(\epsilon) + \rho(\epsilon)}{\eta}.$$

- Now if we turn to the interval  $(\tilde{\tau}_{j-1}, T_{j, 1}]$  (the time between the start of the  $j$ -th attempt and the arrival of the first large noise during this attempt) for each  $j = 1, 2, \dots, J$ , and the following sequence constructed by concatenating these intervals

$$\mathbf{S}(\epsilon, \delta, \eta)$$

$$\triangleq (1, 2, \dots, T_{1,1}, \tilde{\tau}_1 + 1, \tilde{\tau}_1 + 2, \dots, T_{2,1}, \dots, \tilde{\tau}_k + 1, \tilde{\tau}_k + 2, \dots, T_{k+1,1}, \tilde{\tau}_{k+1} + 1, \tilde{\tau}_{k+2} + 1, \dots),$$

then the discussion above have shown that, for

$$\min\{n \in \mathbf{S}(\epsilon, \delta, \eta) : Z_n \text{ has } (3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, \delta, \eta)\text{-overflow}\} \geq T_{J,1}.$$

Meanwhile, from the definition of overflow we know that the probability that  $Z_1$  has  $(3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, \delta, \eta)\text{-overflow}$  is equal to

$$H(\delta/\eta)p(3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, \delta, \eta).$$

- Therefore, if, within the duration of each attempt, we split the attempt into two parts at the arrival time of the first large jump  $(T_{k,1})_{k \geq 1}$  at each attempt, and define (here the subscript *before* or *after* indicates that we are counting the steps before or after the first large jump in an attempt)

$$\begin{aligned} \mathbf{S}_{\text{before}}(\epsilon, \delta, \eta) &\triangleq \{n \in \mathbf{S}(\epsilon, \delta, \eta) : n \leq \sigma(\eta)\}, \quad I_{\text{before}}(\epsilon, \delta, \eta) \triangleq \#\mathbf{S}_{\text{before}}(\epsilon, \delta, \eta), \\ \mathbf{S}_{\text{after}}(\epsilon, \delta, \eta) &\triangleq \{n \notin \mathbf{S}(\epsilon, \delta, \eta) : n \leq \sigma(\eta)\}, \quad I_{\text{after}}(\epsilon, \delta, \eta) \triangleq \#\mathbf{S}_{\text{after}}(\epsilon, \delta, \eta), \end{aligned}$$

then we have  $\sigma(\eta) = I_{\text{before}} + I_{\text{after}}$ . Moreover, the discussion above implies that

$$\begin{aligned} I_{\text{after}} &\leq J(2l^*\hat{t}(\epsilon) + \rho(\epsilon))/\eta \\ I_{\text{before}} &\leq \min\{n \in \mathbf{S}(\epsilon, \delta, \eta) : Z_n \text{ has } (3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, \delta, \eta)\text{-overflow}\} \end{aligned}$$

and on event  $(\mathbf{A}^\times)^c$ .

Define geometric random variables with the following success rates

$$\begin{aligned} U_1(\epsilon, \delta, \eta) &\sim \text{Geom}\left(p(3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, \delta, \eta)\right) \\ U_2(\epsilon, \delta, \eta) &\sim \text{Geom}\left(H(\delta/\eta)p(3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, \delta, \eta)\right). \end{aligned}$$

Using results above to bound  $I_{\text{before}}$  and  $I_{\text{after}}$  separately on event  $(\mathbf{A}^\times)^c$ , we can show that (for all  $\eta$  sufficiently small and any  $u > 0$ )

$$\begin{aligned} &\sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x\left(v^\Omega \lambda(\eta) \sigma(\eta) > u, J_\sigma(\eta) = k'\right) \\ &\leq \sup_{|x| \leq 2\epsilon} \mathbb{P}_x(\mathbf{A}^\times(\epsilon, \delta, \eta)) + \sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x\left(\{v^\Omega \lambda(\eta) \sigma(\eta) > u, J_\sigma(\eta) = k'\} \cap (\mathbf{A}^\times(\epsilon, \delta, \eta))^c\right) \\ &\leq C + \sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x\left(\{v^\Omega \lambda(\eta) \sigma(\eta) > u, J_\sigma(\eta) = k'\} \cap (\mathbf{A}^\times(\epsilon, \delta, \eta))^c\right) \quad \text{due to (B.102)} \\ &\leq C + \sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x\left(\{v^\Omega \lambda(\eta) I_{\text{before}}(\epsilon, \delta, \eta) > (1 - C)u, J_\sigma(\eta) = k'\} \cap (\mathbf{A}^\times(\epsilon, \delta, \eta))^c\right) \\ &\quad + \sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x\left(\{v^\Omega \lambda(\eta) I_{\text{after}}(\epsilon, \delta, \eta) > Cu\} \cap (\mathbf{A}^\times(\epsilon, \delta, \eta))^c\right) \\ &\leq C + \sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x\left(\{v^\Omega \lambda(\eta) I_{\text{before}}(\epsilon, \delta, \eta) > (1 - C)u, J_\sigma(\eta) = k'\} \cap (\mathbf{A}^\times(\epsilon, \delta, \eta))^c\right) \\ &\quad + \mathbb{P}\left(v^\Omega \lambda(\eta) \frac{\rho(\epsilon) + 2l^*\hat{t}(\epsilon)}{\eta} \cdot U_1(\epsilon, \delta, \eta) > Cu\right) \end{aligned}$$

$$\begin{aligned}
&\leq C + \sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x \left( \{v^\Omega \lambda(\eta) I_{\text{before}}(\epsilon, \delta, \eta) > (1-C)u, J_\sigma(\eta) = k'\} \cap (\mathbf{A}^\times(\epsilon, \delta, \eta))^c \cap \{J^\downarrow = J^\uparrow\} \right) \\
&+ \sup_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x \left( (\mathbf{A}^\times)^c \cap \{J^\downarrow < J^\uparrow\} \right) + \mathbb{P} \left( v^\Omega \lambda(\eta) \frac{\rho(\epsilon) + 2l^* \hat{t}(\epsilon)}{\eta} \cdot U_1(\epsilon, \delta, \eta) > Cu \right) \\
&\leq 2C + \mathbb{P} \left( v^\Omega \lambda(\eta) U_2(\epsilon, \delta, \eta) > (1-C)u \right) \frac{\nu_{k'}^\Omega + C}{\nu^\Omega} + \mathbb{P} \left( v^\Omega \lambda(\eta) \frac{\rho(\epsilon) + 2l^* \hat{t}(\epsilon)}{\eta} \cdot U_1(\epsilon, \delta, \eta) > Cu \right)
\end{aligned} \tag{B.107}$$

where the last inequality follows from (B.104) and (B.105). Now let us analyze the probability terms on the last row of the display above. For the first term, let  $a(\eta) = H(\delta/\eta)p(3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, \delta, \eta)$ . Due to Lemma B.14, we have (recall that  $\nu^\Omega = \mu(E(0))$ )

$$\lim_{\eta \downarrow 0} \frac{a(\eta)}{\lambda(\eta)\mu(E(3(\bar{\rho} + \tilde{\rho} + 3)\epsilon))} = 1.$$

Combining this with (B.97), one can see that for all  $\eta$  sufficiently small,

$$\mathbb{P} \left( v^\Omega \lambda(\eta) U_2(\epsilon, \delta, \eta) > (1-C)u \right) \leq \mathbb{P} \left( a(\eta) \text{Geom}(a(\eta)) > (1-C)^2 u \right) \quad \forall u > 0.$$

Next, let  $b(\eta, u) = \mathbb{P} \left( a(\eta) \text{Geom}(a(\eta)) > (1-C)^2 u \right) = \mathbb{P} \left( \text{Geom}(a(\eta)) > \frac{(1-C)^2 u}{a(\eta)} \right)$ . For  $g(y) = \log(1-y)$ , we know the existence of some  $y_0 > 0$  such that for all  $y \in (0, y_0)$ , we have  $\log(1-y) \leq -(1-C)y$ . So one can see that for all  $\eta$  sufficiently small,

$$\begin{aligned}
\log b(u, \eta) &\leq \frac{(1-C)^2 u}{a(\eta)} \log(1-a(\eta)) \leq -(1-C)^3 u \\
&\Rightarrow b(u, \eta) \leq \exp(-(1-C)^3 u)
\end{aligned} \tag{B.108}$$

uniformly for all  $u > 0$ .

For the second probability term, if we only consider  $u \geq C$ , then

$$\mathbb{P} \left( v^\Omega \lambda(\eta) \frac{\rho(\epsilon) + 2l^* \hat{t}(\epsilon)}{\eta} \cdot U_1(\epsilon, \delta, \eta) > Cu \right) \leq \mathbb{P} \left( v^\Omega \lambda(\eta) \frac{\rho(\epsilon) + 2l^* \hat{t}(\epsilon)}{\eta} \cdot U_1(\epsilon, \delta, \eta) > C^2 \right).$$

Using  $H \in \mathcal{RV}_{-\alpha}(\eta)$  with  $\alpha > 1$ , we get

$$p(3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, \delta, \eta) U_1(\epsilon, \delta, \eta) \xrightarrow{d} \text{Exp}(1) \quad \text{as } \eta \downarrow 0$$

due to the nature of the Geometric random variable  $U_1$ . Besides, due to  $H \in \mathcal{RV}_{-\alpha}(\eta)$  with  $\alpha > 1$  and Lemma B.14, it is easy to show that

$$\lim_{\eta \downarrow 0} \frac{\lambda(\eta) \frac{\rho(\epsilon) + 2l^* \hat{t}(\epsilon)}{\eta}}{p(3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, \delta, \eta)} = 0.$$

Combining these results with Slutsky's theorem, we now obtain

$$\mu(E(0))\lambda(\eta) \frac{\rho(\epsilon) + 2l^* \hat{t}(\epsilon)}{\eta} \cdot U_1(\epsilon, \delta, \eta) \xrightarrow{d} 0 \quad \text{as } \eta \downarrow 0.$$

Therefore,

$$\limsup_{\eta \downarrow 0} \sup_{u \geq C} \mathbb{P} \left( \mu(E(0))\lambda(\eta) \frac{\rho(\epsilon) + 2l^* \hat{t}(\epsilon)}{\eta} \cdot U_1(\delta, \eta) > Cu \right) = 0. \tag{B.109}$$

Plugging (B.108) and (B.109) back into (B.107), we can establish the upper bound in (B.67). To show (B.69), note that for event

$$E(\epsilon, \eta) = \{\nu^\Omega \lambda(\eta) \tau(\eta, \epsilon) > u, X_{\tau(\eta, \epsilon)}^\eta \in B(\tilde{m}_{k'}, 2\epsilon)\},$$

we have (for definitions of  $\tau$ , see (B.62))

$$\begin{aligned} E(\epsilon, \eta) &\supseteq \{v^\Omega \lambda(\eta) \sigma(\eta) > u, J_\sigma(\eta) = k'\} \cap \{X_n^\eta \in \tilde{\Omega}_{J_\sigma(\eta)} \quad \forall n \in [\sigma(\eta), \tau(\eta, \epsilon)]\}, \\ E(\epsilon, \eta) \cap \{v^\Omega \lambda(\eta) \sigma(\eta) > u, J_\sigma(\eta) = j\} &\cap \{X_n^\eta \in \tilde{\Omega}_{J_\sigma(\eta)} \quad \forall n \in [\sigma(\eta), \tau(\eta, \epsilon)]\} = \emptyset \quad \forall j \neq k'. \end{aligned}$$

Therefore, for all  $\eta$  sufficiently small,

$$\begin{aligned} &\sup_{|x| \leq 2\epsilon} \mathbb{P}_x(E(\epsilon, \eta)) \\ &\leq \sup_{|x| \leq 2\epsilon} \mathbb{P}_x(\mathbf{A}^\times) + \sup_{|x| \leq 2\epsilon} \mathbb{P}_x((\mathbf{A}^\times)^c \cap \{X_n^\eta \notin \tilde{\Omega}_{J_\sigma(\eta)} \text{ for some } n \in [\sigma(\eta), \tau(\eta, \epsilon)]\}) \\ &\quad + \sup_{|x| \leq 2\epsilon} \mathbb{P}_x\left((\mathbf{A}^\times)^c \cap \{v^\Omega \lambda(\eta) \sigma(\eta) > u, J_\sigma(\eta) = k'\} \cap \{X_n^\eta \in \tilde{\Omega}_{J_\sigma(\eta)} \quad \forall n \in [\sigma(\eta), \tau(\eta, \epsilon)]\}\right) \\ &\leq 4C + \exp(-(1-C)^3 u) \frac{\nu_{k'}^\Omega + C}{\nu^\Omega} \end{aligned}$$

uniformly for all  $u \geq C$ , due to (B.102), (B.71) and (B.107).

The lower bound can be shown by an almost identical approach. In particular, analogous to (B.107), we can show that (for any  $u > 0$ )

$$\begin{aligned} &\inf_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x\left(v^\Omega \lambda(\eta) \sigma(\eta) > u, J_\sigma(\eta) = k'\right) \\ &\geq \inf_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x\left(v^\Omega \lambda(\eta) \sigma(\eta) > u, J_\sigma(\eta) = k'\} \cap (\mathbf{A}^\times(\epsilon, \delta, \eta))^c\right) \\ &\geq \inf_{x \in [-2\epsilon, 2\epsilon]} \mathbb{P}_x\left(\{v^\Omega \lambda(\eta) I_{\text{before}}(\epsilon, \delta, \eta) > u, J_\sigma(\eta) = k'\} \cap (\mathbf{A}^\times(\epsilon, \delta, \eta))^c\right) \\ &\geq \mathbb{P}\left(v^\Omega \lambda(\eta) U'_2(\epsilon, \delta, \eta) > (1-C)u\right) \frac{\nu_{k'}^\Omega + C}{\nu^\Omega} - 2C \text{ due to } \mathbb{P}(E \setminus F) \geq \mathbb{P}(E) - \mathbb{P}(F) \text{ and (B.102)(B.104)(B.106)} \end{aligned}$$

where

$$U'_2(\epsilon, \delta, \eta) \sim \text{Geom}\left(H(\delta/\eta)p(-3(\bar{\rho} + \tilde{\rho} + 3)\epsilon, \delta, \eta)\right).$$

Using the similar argument leading to (B.108), we are able to show (B.68), (B.70) and conclude the proof.  $\square$

Recall that  $\sigma_i(\eta) = \min\{n \geq 0 : X_n \notin \Omega_i\}$  and that value of constants  $q_i, q_{i,j}$  are specified via (5)-(10). Define

$$\tau_i^{\min}(\eta, \epsilon) \triangleq \min\{n \geq \sigma_i(\eta) : X_n^\eta \in \bigcup_j [m_j - 2\epsilon, m_j + 2\epsilon]\}, \quad (\text{B.110})$$

$$J_i(\eta) = j \iff X_{\sigma_i(\eta)}^\eta \in \Omega_j \quad \forall j \in [n_{\min}]. \quad (\text{B.111})$$

The following result is simply a restatement of Proposition B.22 under the new system of notations. Despite the reiteration, we still state it here because this is the version that will be used to prove Lemma 8, which is the key tool for establishing Theorem 1, as well as many other results in Appendix C.

**Proposition B.24.** *Given  $C > 0$  and  $i, j \in [n_{\min}]$  such that  $i \neq j$ , the following claims hold for all  $\epsilon > 0$  that are sufficiently small:*

$$\begin{aligned}
& \limsup_{\eta \downarrow 0} \sup_{u \in (C, \infty)} \sup_{x \in (m_i - 2\epsilon, m_i + 2\epsilon)} \mathbb{P}_x \left( q_i \lambda_i(\eta) \sigma_i(\eta) > u, X_{\sigma_i(\eta)}^\eta \in \Omega_j \right) \leq C + \exp \left( - (1 - C)u \right) \frac{q_{i,j} + C}{q_i}, \\
& \liminf_{\eta \downarrow 0} \inf_{u \in (C, \infty)} \inf_{x \in (m_i - 2\epsilon, m_i + 2\epsilon)} \mathbb{P}_x \left( q_i \lambda_i(\eta) \sigma_i(\eta) > u, X_{\sigma_i(\eta)}^\eta \in \Omega_j \right) \geq -C + \exp \left( - (1 + C)u \right) \frac{q_{i,j} - C}{q_i}, \\
& \limsup_{\eta \downarrow 0} \sup_{u \in (C, \infty)} \sup_{x \in (m_i - 2\epsilon, m_i + 2\epsilon)} \mathbb{P}_x \left( q_i \lambda_i(\eta) \tau_i^{\min}(\eta, \epsilon) > u, X_{\tau_i^{\min}(\eta, \epsilon)}^\eta \in \Omega_j \right) \leq C + \exp \left( - (1 - C)u \right) \frac{q_{i,j} + C}{q_i}, \\
& \liminf_{\eta \downarrow 0} \inf_{u \in (C, \infty)} \inf_{x \in (m_i - 2\epsilon, m_i + 2\epsilon)} \mathbb{P}_x \left( q_i \lambda_i(\eta) \tau_i^{\min}(\eta, \epsilon) > u, X_{\tau_i^{\min}(\eta, \epsilon)}^\eta \in \Omega_j \right) \geq -C + \exp \left( - (1 + C)u \right) \frac{q_{i,j} - C}{q_i}, \\
& \liminf_{\eta \downarrow 0} \inf_{x \in (m_i - 2\epsilon, m_i + 2\epsilon)} \mathbb{P}_x \left( q_i \lambda_i(\eta) (\tau_i^{\min}(\eta, \epsilon) - \sigma_i(\eta)) < C, X_n^\eta \in \Omega_{J_i(\eta)} \quad \forall n \in [\sigma_i(\eta), \tau_i^{\min}(\eta, \epsilon)] \right) \geq 1 - C.
\end{aligned}$$

Concluding this section, we apply Proposition B.24 and prove Lemma 8.

*Proof of Lemma 8.* Fix some  $C \in (0, 1)$ ,  $u > 0$ , and some  $k, l \in [n_{\min}]$  with  $k \neq l$ . Let  $q_i, q_{i,j}$  be the constants defined in (10).

Fix some  $C_0 \in (0, \frac{C}{n_{\min}} \wedge \frac{q_k}{n_{\min}} C)$ . Using Proposition B.24, we know that for all  $\epsilon$  sufficiently small, we have

$$\limsup_{\eta \downarrow 0} \sup_{x \in (m_k - 2\epsilon, m_k + 2\epsilon)} \mathbb{P}_x \left( q_k \lambda_k(\eta) \sigma_k(\eta) > u, X_{\sigma_k(\eta)}^\eta \in \Omega_j \right) \leq C_0 + \exp \left( - (1 - C)u \right) \frac{q_{k,j} + C_0}{q_k} \quad \forall j \in [n_{\min}].$$

Summing up the inequality above over all  $j \in [n_{\min}]$ , we can obtain (17). The lower bound (18) can be established using an identical approach.

In order to show (20), note that we can find  $C_1 \in (0, u)$  sufficiently small so that

$$-C_1 + \exp \left( - (1 + C_1) \cdot 2C_1 \right) \frac{q_{k,l} - C_1}{q_k} \geq \frac{q_{k,l} - C}{q_k}.$$

Fix such  $C_1$ . From Proposition B.24, we also know that for all  $\epsilon$  small enough, we have

$$\liminf_{\eta \downarrow 0} \inf_{x \in (m_k - 2\epsilon, m_k + 2\epsilon)} \mathbb{P}_x \left( q_k \lambda_k(\eta) \sigma_k(\eta) > u, X_{\sigma_k(\eta)}^\eta \in \Omega_l \right) \geq -C_1 + \exp \left( - (1 + C_1) \cdot 2C_1 \right) \frac{q_{k,l} - C_1}{q_k}.$$

Then using  $\mathbb{P}_x \left( X_{\sigma_k(\eta)}^\eta \in \Omega_l \right) \geq \mathbb{P}_x \left( q_k \lambda_k(\eta) \sigma_k(\eta) > u, X_{\sigma_k(\eta)}^\eta \in \Omega_l \right)$  we conclude the proof for (20).

Moving on, we show (19) in the following way. Note that we can find  $C_2 \in (0, u)$  small enough so that

$$2C_2 + \frac{q_{k,l} + C_2}{q_k} < \frac{q_{k,l} + C}{q_k}. \quad (\text{B.112})$$

Fix such  $C_2$ . Since (18) has been established already, we can find some  $u_2 > 0$  such that for all  $\epsilon$  small enough,

$$\limsup_{\eta \downarrow 0} \sup_{x \in (m_k - 2\epsilon, m_k + 2\epsilon)} \mathbb{P}_x \left( q_k \lambda_k(\eta) \sigma_k(\eta) \leq u_2 \right) < C_2 \quad (\text{B.113})$$

Fix such  $u_2$ . Meanwhile, fix some  $C_3 \in (0, C_2 \wedge u_2)$ . From Proposition B.24 we know that for all  $\epsilon$  sufficiently small,

$$\limsup_{\eta \downarrow 0} \sup_{x \in (m_k - 2\epsilon, m_k + 2\epsilon)} \mathbb{P}_x \left( q_k \lambda_k(\eta) \sigma_k(\eta) > u_2, X_{\sigma_k(\eta)}^\eta \in \Omega_l \right) \leq C_3 + \exp \left( - (1 - C_3)u_2 \right) \frac{q_{k,l} + C_3}{q_k}$$

$$\leq C_2 + \frac{q_{k,l} + C_2}{q_k}. \quad (\text{B.114})$$

Lastly, observe the following decomposition of events (for any  $x \in \Omega_k$ )

$$\mathbb{P}_x(X_{\sigma_k(\eta)}^\eta \in \Omega_l) \leq \mathbb{P}_x(q_k \lambda_k(\eta) \sigma_k(\eta) \leq u_2) + \mathbb{P}_x(q_k \lambda_k(\eta) \sigma_k(\eta) > u_2, X_{\sigma_k(\eta)}^\eta \in \Omega_l).$$

Combining this bound with (B.112)-(B.114), we complete the proof.  $\square$

## C Proofs for Technical Lemmas in Section 2.6

### C.1 Proof of Lemma 9

First, we introduce another dichotomy for *small* and *large* noises. For any  $\tilde{\gamma} > 0$  and any learning rate  $\eta > 0$ , we say that a noise  $Z_n$  is small if

$$\eta|Z_n| > \eta^{\tilde{\gamma}}$$

and we say  $Z_n$  is large otherwise. For this new classification of small and large noises, we introduce the following notations and definitions:

$$Z_n^{\leq, \tilde{\gamma}, \eta} = Z_n \mathbb{1}\{\eta|Z_n| \leq \eta^{\tilde{\gamma}}\}, \quad (\text{C.1})$$

$$Z_n^{>, \tilde{\gamma}, \eta} = Z_n \mathbb{1}\{\eta|Z_n| > \eta^{\tilde{\gamma}}\}, \quad (\text{C.2})$$

$$\tilde{T}_1^\eta(\tilde{\gamma}) \triangleq \min\{n \geq 1 : \eta|Z_n| > \eta^{\tilde{\gamma}}\}. \quad (\text{C.3})$$

Similar to Lemma B.10, the following result is a direct application of Lemma B.9, and shows that it is rather unlikely to observe large perturbation that are caused only by *small* noises. Specifically, since  $\alpha > 1$  we can always find

$$\begin{aligned} \tilde{\gamma} &\in (0, (1 - \frac{1}{\alpha} \wedge \frac{1}{2})) \\ \beta &\in (1, (2 - 2\tilde{\gamma}) \wedge \alpha(1 - \tilde{\gamma})). \end{aligned}$$

Now in Lemma B.9, if we let  $\Delta = \tilde{\gamma}$ ,  $\tilde{\Delta} = \Delta/2$  and  $\epsilon = \delta = 1$  (in other words,  $u(\eta) = 1/\eta^{1-\tilde{\gamma}}$ ,  $v(\eta) = \eta^{\tilde{\gamma}/2}$ ), then for any positive integer  $j$  the condition (B.38) is satisfied, allowing us to draw the following conclusion immediately as a corollary from Lemma B.9.

**Lemma C.1.** *Given  $N > 0$  and*

$$\tilde{\gamma} \in (0, (1 - \frac{1}{\alpha}) \wedge \frac{1}{2}), \quad \beta \in (1, (2 - 2\tilde{\gamma}) \wedge (\alpha - \alpha\tilde{\gamma})),$$

*we have (as  $\eta \downarrow 0$ )*

$$\mathbb{P}\left(\max_{j=1,2,\dots,\lceil(1/\eta)^\beta\rceil} \eta|Z_1^{\leq, \tilde{\gamma}, \eta} + \dots + Z_j^{\leq, \tilde{\gamma}, \eta}| > \eta^{\tilde{\gamma}/2}\right) = o(\eta^N).$$

The flavor of the next lemma is similar to that of Lemma B.11. Specifically, we show that, with high probability, the SGD iterates would quickly return to the local minimum as long as they start from somewhere that are not too close the boundary of an attraction field (namely, the points  $s_1, s_2, \dots, s_{n_{\min}}$ ). To this end, we consider a refinement of function  $\hat{t}(\cdot)$  defined in (B.24). For any  $i = 1, 2, \dots, n_{\min}$ , any  $x \in \Omega_i$  and any  $\eta > 0, \gamma \in (0, 1)$ , we can define the return time to  $\eta^\gamma$ -neighborhood for the ODE  $\mathbf{x}^\eta$  as

$$\hat{t}_\gamma^{(i)}(x, \eta) \triangleq \min\{t \geq 0 : |\mathbf{x}^\eta(t, x) - m_i| \leq \eta^\gamma\}.$$

Given the bound in (B.23) (which is stated for a specific attraction field) and the fact that there only exists finitely many attraction fields, we know the existence of some  $c_2 < \infty$  such that for any  $i = 1, 2, \dots, n_{\min}$ , any  $\eta > 0$ , any  $\gamma \in (0, 1)$  and any  $x \in \Omega_i$  such that  $|x - s_i| \vee |x - s_{i-1}| > \eta^\gamma$ , we have

$$\hat{t}_\gamma^{(i)}(x, \eta) \leq c_2 \gamma \log(1/\eta)/\eta$$

and define function  $t^\dagger$  as

$$t^\dagger(\eta, \gamma) \triangleq c_2 \gamma \log(1/\eta).$$

Lastly, define the following stopping time for any  $i = 1, 2, \dots, n_{\min}$ , any  $x \in \Omega_i$  and any  $\Delta > 0$

$$T_{\text{return}}^{(i)}(\eta, \Delta) \triangleq \min\{n \geq 0 : X_n^\eta(x) \in B(m_i, 2\Delta)\}$$

where we adopt the notation  $B(u, v) \triangleq [u - v, u + v]$  for the  $v$ -neighborhood around point  $u$ .

**Lemma C.2.** *Given*

$$\tilde{\gamma} \in (0, (1 - \frac{1}{\alpha}) \wedge (\frac{1}{2})), \quad \gamma \in (0, \frac{\tilde{\gamma}}{16Mc_2} \wedge \frac{\tilde{\gamma}}{4}),$$

*and any  $i = 1, 2, \dots, n_{\min}$ , any  $\Delta > 0$ , we have*

$$\liminf_{\eta \downarrow 0} \inf_{x \in \Omega_i : |x - s_{i-1}| \vee |x - s_i| \geq 2\eta^\gamma} \mathbb{P}_x \left( T_{\text{return}}^{(i)}(\eta, \Delta) \leq 2c_2 \gamma \log(1/\eta)/\eta, \quad X_n^\eta \in \Omega_i \quad \forall n \leq T_{\text{return}}^{(i)}(\eta, \Delta) \right) = 1.$$

*Proof.* Throughout this proof, we only consider  $\eta$  small enough such that

$$2c_2 \gamma \log(1/\eta)/\eta < \lceil (1/\eta)^\beta \rceil, \quad \eta M \leq \eta^{2\tilde{\gamma}}, \quad 2\eta^\gamma < \Delta/2, \quad 2\eta^{\tilde{\gamma}/4} < \eta^\gamma. \quad (\text{C.4})$$

The condition above holds for all  $\eta > 0$  sufficiently small because  $\beta > 1$ ,  $2\tilde{\gamma} < 1$ , and  $\gamma < \tilde{\gamma}/4$ . Also, fix some  $\beta \in (1, (2 - 2\tilde{\gamma}) \wedge (\alpha - \alpha\tilde{\gamma}))$

Define the following events

$$\begin{aligned} A_1^\times(\eta) &\triangleq \left\{ \max_{j=1,2,\dots,\lceil (1/\eta)^\beta \rceil} \eta |Z_1^{\leq, \tilde{\gamma}, \eta} + \dots + Z_j^{\leq, \tilde{\gamma}, \eta}| > \eta^{\tilde{\gamma}/2} \right\} \\ A_2^\times(\eta) &\triangleq \{ \tilde{T}_1^\eta(\tilde{\gamma}) \leq \lceil (1/\eta)^\beta \rceil \} \quad (\text{see (C.3) for definition of the stopping time involved}) \end{aligned}$$

and fix some  $N > 0$ . From Lemma C.1, we see that (for all sufficiently small  $\eta$ )

$$\mathbb{P}(A_1^\times(\eta)) \leq \eta^N.$$

Besides, using Lemma B.2 together with the fact that  $\beta < \alpha(1 - \tilde{\gamma})$ , we know the existence of some constant  $\theta > 0$  such that

$$\mathbb{P}(A_2^\times(\eta)) \leq \eta^\theta$$

for all sufficiently small  $\eta$ .

Now we focus on the behavior of the SGD iterates on event  $\left( A_1^\times(\eta) \cap A_2^\times(\eta) \right)^c$ . Let us arbitrarily choose some  $x \in \Omega_i$  such that  $|x - s_i| \vee |x - s_{i-1}| > 2\eta^\gamma$ . First, from Lemma B.4 and (C.4), we know that

$$|\mathbf{x}_t^\eta(x) - \mathbf{y}_{\lfloor t \rfloor}^\eta(x)| \leq 2\eta M \exp(2Mc_2 \gamma \log(1/\eta)) \leq 2\eta^{2\tilde{\gamma} - 2Mc_2 \gamma} \leq 2\eta^{\tilde{\gamma}} \leq \eta^\gamma \quad \forall t \leq 2c_2 \gamma \log(1/\eta)/\eta \quad (\text{C.5})$$



Next, from the definition of the function  $t^\dagger(\cdot)$  and (C.5), we know that for

$$T_{\text{GD},\text{return}}(x; \eta, \Delta) \triangleq \min\{n \geq 0 : \mathbf{y}_n^\eta(x) \in B(m_i, \frac{\Delta}{2} + \eta^\gamma)\}, \quad (\text{C.6})$$

we have

$$T_{\text{GD},\text{return}}(x; \eta, \Delta) \leq 2c_2\gamma \log(1/\eta)/\eta \quad (\text{C.7})$$

$$\mathbf{y}_n^\eta > s_{i-1} + \eta^\gamma, \quad \mathbf{y}_n^\eta < s_i - \eta^\gamma \quad \forall n \leq 2c_2\gamma \log(1/\eta)/\eta. \quad (\text{C.8})$$

Furthermore, on event  $\left(A_1^\times(\eta) \cap A_2^\times(\eta)\right)^c$ , due to Lemma B.3 and (C.4), we have that

$$|X_n^\eta(x) - \mathbf{y}_n^\eta(x)| \leq \eta^{\tilde{\gamma}/2} \exp(2Mc_2 \log(1/\eta)) = \eta^{\frac{\tilde{\gamma}}{2} - 2Mc_2\gamma} \leq \eta^{\tilde{\gamma}/4} < \eta^\gamma \quad \forall n \leq 2c_2\gamma \log(1/\eta)/\eta.$$

Combining this with (C.7), (C.8), we can conclude that (recall that due to (C.4) we have  $2\eta^\gamma < \Delta/2$ )

$$\begin{aligned} T_{\text{return}}^{(i)}(\eta, \Delta) &\leq 2c_2\gamma \log(1/\eta)/\eta \\ X_n^\eta &\in \Omega_i \quad \forall n \leq 2c_2\gamma \log(1/\eta)/\eta \end{aligned}$$

on event  $\left(A_1^\times(\eta) \cap A_2^\times(\eta)\right)^c$ . Therefore,

$$\begin{aligned} &\liminf_{\eta \downarrow 0} \inf_{x \in \Omega_i : |x - s_{i-1}| \vee |x - s_i| \geq 2\eta^\gamma} \mathbb{P}_x \left( T_{\text{return}}^{(i)}(\eta, \Delta) \leq 2c_2\gamma \log(1/\eta)/\eta, \quad X_n^\eta \in \Omega_i \quad \forall n \leq T_{\text{return}}^{(i)}(\eta, \Delta) \right) \\ &\geq \liminf_{\eta \downarrow 0} \mathbb{P} \left( \left( A_1^\times(\eta) \cap A_2^\times(\eta) \right)^c \right) \geq \liminf_{\eta \downarrow 0} 1 - \eta^N - \eta^\theta = 1. \end{aligned}$$

This concludes the proof.  $\square$

The takeaway of the next lemma is that, almost always, the SGD iterates will quickly escape from the neighborhood of any  $s_i$ , the boundaries of each attraction fields.

**Lemma C.3.** *Given any  $\gamma \in (0, 1), t > 0$ , we have*

$$\liminf_{\eta \downarrow 0} \inf_{x \in [-L, L]} \mathbb{P}_x \left( \min \{n \geq 0 : X_n^\eta \notin \bigcup_i B(s_i, 2\eta^\gamma)\} \leq \frac{t}{H(1/\eta)} \right) = 1.$$

*Proof.* We only consider  $\eta$  small enough so that

$$\begin{aligned} \min_{i=2,3,\dots,n_{\min}-1} |s_i - s_{i-1}| &> 3\eta^{\frac{1+\gamma}{2}}, \\ \eta M &< \eta^\gamma. \end{aligned}$$

Also, the claim is trivial if  $x \notin \cup_j B(s_j, 2\eta^\gamma)$ , so without loss of generality we only consider the case where there is some  $j \in [n_{\min}]$  and  $x \in [-L, L], x \in B(s_j, 2\eta^\gamma)$ . Let us define stopping times

$$T^\gamma \triangleq \min\{n \geq 1 : \eta|Z_n| > 5\eta^\gamma\}; \quad (\text{C.9})$$

$$T_{\text{escape}}^\gamma \triangleq \min\{n \geq 0 : X_n^\eta \notin \cup_j B(s_j, 2\eta^\gamma)\}, \quad (\text{C.10})$$

and the following two events

$$\begin{aligned} A_1^\times(\eta) &\triangleq \{T^\gamma > \frac{t}{H(1/\eta)}\}, \\ A_2^\times(\eta) &\triangleq \{\eta|Z_{T^\gamma}| > \eta^{\frac{1+\gamma}{2}}\}. \end{aligned}$$

First, using Lemma B.1 and the regularly varying nature of  $H(\cdot) = \mathbb{P}(|Z_1| > \cdot)$ , we know the existence of some  $\theta > 0$  such that

$$\mathbb{P}(A_1^\times(\eta)) \leq \exp(-1/\eta^\theta)$$

for all  $\eta > 0$  sufficiently small. Next, by definition of  $T^\gamma$ , one can see that (for any  $\eta \in (0, 1)$ )

$$\mathbb{P}(A_2^\times(\eta)) = \frac{H(1/\eta^{\frac{1-\gamma}{2}})}{H(5/\eta^{1-\gamma})}.$$

Again, due to  $H \in \mathcal{RV}_{-\alpha}$  and  $1 - \gamma > 0$ , we know the existence of some  $\theta_1 > 0$  such that

$$\mathbb{P}(A_2^\times(\eta)) < \eta^{\theta_1}$$

for all  $\eta > 0$  sufficiently small. To conclude the proof, we only need to note the following fact on event  $(A_1^\times(\eta) \cup A_2^\times(\eta))^c$ . There are only two possibilities on this event:  $T_{\text{escape}}^\gamma \leq T^\gamma - 1$ , or  $T_{\text{escape}}^\gamma \geq T^\gamma$ . Now we analyze the two cases respectively.

- On  $(A_1^\times(\eta) \cup A_2^\times(\eta))^c \cap \{T_{\text{escape}}^\gamma \leq T^\gamma - 1\}$ , we must have  $T_{\text{escape}}^\gamma < T^\gamma \leq t/H(1/\eta)$ .
- On  $(A_1^\times(\eta) \cup A_2^\times(\eta))^c \cap \{T_{\text{escape}}^\gamma \geq T^\gamma\}$ , we know that at  $n = T^\gamma - 1$ , there exists an integer  $j \in \{1, 2, \dots, n_{\min} - 1\}$  such that  $X_n^\eta \in B(s_j, 2\eta^\gamma)$ . Now since  $\eta M < \eta^\gamma$  and  $\eta|Z_{T^\gamma}| > 5\eta^\gamma$ , we must have

$$|X_{T^\gamma}^\eta - X_{T^\gamma-1}^\eta| > 4\eta^\gamma \Rightarrow X_{T^\gamma}^\eta \notin B(s_j, 2\eta^\gamma).$$

On the other hand, the exclusion of event  $A_2^\times(\eta)$  tells us that  $|X_{T^\gamma}^\eta - X_{T^\gamma-1}^\eta| < 2\eta^{\frac{1+\gamma}{2}}$ . Due to (C.9), we then have  $X_n^\eta \notin \cup_i B(s_i, 2\eta^\gamma)$ .

In summary,  $(A_1^\times(\eta) \cup A_2^\times(\eta))^c \subseteq \{T_{\text{return}}^\gamma \leq t/H(1/\eta)\}$  and this conclude the proof.  $\square$

In the next lemma, we analyze the number of transitions needed to visit a certain local minimizer in the loss landscape. In general, we focus on a communication class  $G$  and, for now, assume it is absorbing. Next, we introduce the following concepts to record the transitions between different local minimum. To be specific, for any  $\eta > 0$  and any  $\Delta > 0$  small enough so that  $B(m_j, \Delta) \cap \Omega_j^c = \emptyset$  for all  $j$ , define

$$T_0(\eta, \Delta) = \min\{n \geq 0 : X_n^\eta \in \cup_j B(m_j, 2\Delta)\}; \quad (\text{C.11})$$

$$I_0(\eta, \Delta) = j \text{ iff } X_{T_0(\eta, \Delta)}^\eta \in B(m_j, 2\Delta); \quad (\text{C.12})$$

$$T_k(\eta, \Delta) = \min\{n > T_{k-1}(\eta, \Delta) : X_n^\eta \in \cup_{j \neq I_{k-1}(\eta, \Delta)} B(m_j, 2\Delta)\} \quad \forall k \geq 1 \quad (\text{C.13})$$

$$I_k(\eta, \Delta) = j \text{ iff } X_{T_k(\eta, \Delta)}^\eta \in B(m_j, 2\Delta) \quad \forall k \geq 1. \quad (\text{C.14})$$

As mentioned earlier, the next goal is to analyze the transitions between attraction fields it takes to visit  $m_j$  when starting from  $m_i$  when  $m_i, m_j \in G$ . Define

$$K_i(\eta, \Delta) \triangleq \min\{k \geq 0 : I_k(\eta, \Delta) = i\}.$$

**Lemma C.4.** *Assume that  $G$  is an absorbing communication class on the graph  $\mathcal{G}$ . Then there exists some constant  $p > 0$  such that for any  $i$  with  $m_i \in G$ , any  $\epsilon > 0$ , and any  $\Delta > 0$ ,*

$$\begin{aligned} \sup_{j: m_j \in G; x \in B(m_j, 2\Delta)} \mathbb{P}_x(K_i(\eta, \Delta) > u \cdot n_{\min}) &\leq \mathbb{P}(\text{Geom}(p) \geq u) + \epsilon \quad \forall u = 1, 2, \dots, \\ \sup_{j: m_j \in G; x \in B(m_j, 2\Delta)} \mathbb{P}_x(\exists k \in [K_i(\eta, \Delta)] \text{ s.t. } m_{I_k(\eta, \Delta)} \notin G) &\leq \epsilon \end{aligned}$$

hold for all  $\eta > 0$  sufficiently small.

*Proof.* The claim is trivial if, for the initial condition, we have  $x \in B(m_i, 2\Delta)$ . Next, let us observe the following facts.

- Define (recall the definitions of measure  $\mu_i$  and sets  $E_i, E_{i,j}$  in (4)(8)(9))

$$J(j) \triangleq \arg \min_{\tilde{j}: \mu_i(E_{j,\tilde{j}}) > 0} |i - \tilde{j}| \quad \forall j \neq i$$

$$p^* \triangleq \min_{j: j \neq i, m_j \in G} \frac{\mu_j(E_{j,J(j)})}{\mu_j(E_j)}.$$

- From the definition of  $J(j)$  and the fact that there are only finitely many attraction fields we can see that  $p^* > 0$ . Moreover,  $G$  being a communication class implies that

$$|J(j) - i| < |j - i| \quad \forall j \neq i, m_j \in G.$$

Indeed, if  $i < j$ , then since  $G$  is a communication class and there are some  $m_i \in G$  with  $i < j$ , we will at least have  $\mu_j(E_{j,j-1}) > 0$ , so  $|J(j) - i| \leq |i - j| - 1$ ; the case that  $i > j$  can be approached analogously.

- Now from the definition of  $J(j)$  and Proposition B.24, together with the previous bullet point, we know that for all  $\eta$  sufficiently small,

$$\inf_{x \in [-L, L]} \mathbb{P}_x(|I_{k+1} - i| \leq |I_k - i| - 1, m_{I_{k+1}} \in G \mid K_i(\eta, \Delta) > k, m_{I_k} \in G) \geq p^*/2$$

uniformly for all  $k \geq 0$ .

- Meanwhile, since  $p^* > 0$ , we are able to fix some  $\delta > 0$  small enough such that

$$\frac{n_{\min} \delta}{(p^*/2)^{n_{\min}}} < \epsilon.$$

- On the other hand, for any  $\tilde{j}$  with  $m_{\tilde{j}} \notin G$ , by definition of the typical transition graph we must have  $\mu_j(E_{j,\tilde{j}}) = 0$  for any  $j$  with  $m_j \in G$ . Then due to Proposition B.24 again, one can see that for all  $\eta > 0$  that is sufficiently small,

$$\sup_{x \in [-L, L]} \mathbb{P}_x(m_{I_{k+1}} \notin G \mid K_i(\eta, \Delta) > k, m_{I_k} \in G) < \delta$$

uniformly for all  $k \geq 0$ .

- Repeat this argument for  $n_{\min}$  times, and we can see that for all  $\eta$  sufficiently small

$$\inf_{x \in [-L, L]} \mathbb{P}_x(K_i(\eta, \Delta) \leq k + n_{\min} \mid K_i(\eta, \Delta) > k, m_{I_k} \in G) \geq \left(\frac{p^*}{2}\right)^{n_{\min}}$$

$$\sup_{x \in [-L, L]} \mathbb{P}_x(\exists l \in [n_{\min}] \text{ s.t. } m_{I_{k+l}} \notin G \mid K_i(\eta, \Delta) > k, m_{I_k} \in G) \leq n_{\min} \delta$$

uniformly for all  $k \geq 1$ .

- Lastly, to apply the bounds established above, we will make use of the following expression of several probabilities. For any  $j$  with  $m_j \in G$  and  $x \in B(m_j, 2\Delta)$  and any  $u = 1, 2, \dots$ ,

$$\mathbb{P}_x(\exists l \in [un_{\min}] \text{ s.t. } m_{I_l} \notin G, K_i(\eta, \Delta) > u \cdot n_{\min})$$

$$= \sum_{v=0}^{u-1} \mathbb{P}_x\left(\exists k \in [n_{\min}] \text{ s.t. } m_{I_{k+vn_{\min}}} \notin G \mid K_i(\eta, \Delta) > v \cdot n_{\min}, m_{I_l} \in G \forall l \leq vn_{\min}\right)$$

$$\begin{aligned}
& \cdot \prod_{w=0}^{v-1} \mathbb{P}_x \left( K_i(\eta, \Delta) > (w+1)n_{\min}, m_{I_{k+wn_{\min}}} \in G \ \forall k \in [n_{\min}] \mid K_i(\eta, \Delta) > w \cdot n_{\min}, m_{I_l} \in G \ \forall l \leq wn_{\min} \right) \\
& \mathbb{P}_x(m_{I_l} \in G \ \forall l \in [un_{\min}], K_i(\eta, \Delta) > u \cdot n_{\min}) \\
& = \prod_{v=0}^{u-1} \mathbb{P} \left( K_i(\eta, \Delta) > (v+1)n_{\min}, m_{I_{k+vn_{\min}}} \in G \ \forall k \in [n_{\min}] \mid K_i(\eta, \Delta) > vn_{\min}, m_{I_k} \in G \ \forall k \in [vn_{\min}] \right)
\end{aligned}$$

In summary, now we can see that (for sufficiently small  $\eta$ )

$$\begin{aligned}
& \sup_{j: m_j \in G; x \in B(m_j, 2\Delta)} \mathbb{P}_x \left( \exists k \in [K_i(\eta, \Delta)] \text{ s.t. } m_{I_k(\eta, \Delta)} \notin G \right) \\
& \leq \sum_{u=0}^{\infty} \sup_{j: m_j \in G; x \in B(m_j, 2\Delta)} \mathbb{P}_x \left( \exists v \in [n_{\min}] \text{ such that } m_{I_{v+un_{\min}}} \notin G, K_i(\eta, \Delta) > un_{\min}, m_{I_k} \in G \ \forall k \in [un_{\min}] \right) \\
& \leq \sum_{u=0}^{\infty} \sup_{j: m_j \in G; x \in B(m_j, 2\Delta)} \mathbb{P}_x \left( \exists v \in [n_{\min}] \text{ such that } m_{I_{v+un_{\min}}} \notin G \mid K_i(\eta, \Delta) > un_{\min}, m_{I_k} \in G \ \forall k \in [un_{\min}] \right) \\
& \cdot \prod_{v=0}^{u-1} \sup_{j: m_j \in G; x \in B(m_j, 2\Delta)} \mathbb{P}_x \left( K_i(\eta, \Delta) > (v+1)n_{\min}, m_{I_{k+vn_{\min}}} \in G \ \forall k \in [n_{\min}] \mid K_i(\eta, \Delta) > vn_{\min}, m_{I_k} \in G \ \forall k \in [vn_{\min}] \right) \\
& \leq \sum_{u \geq 0} n_{\min} \delta \left( 1 - \left( \frac{p^*}{2} \right)^{n_{\min}} \right)^{u-1} = \frac{n_{\min} \delta}{(p^*/2)^{n_{\min}}} \leq \epsilon
\end{aligned}$$

and

$$\begin{aligned}
& \sup_{j: m_j \in G, x \in B(m_j, 2\Delta)} \mathbb{P}_x(K_i(\eta, \Delta) > u \cdot n_{\min}) \\
& \leq \sup_{j: m_j \in G, x \in B(m_j, 2\Delta)} \mathbb{P}_x(\exists l \in [un_{\min}] \text{ s.t. } m_{I_l} \notin G, K_i(\eta, \Delta) > u \cdot n_{\min}) \\
& + \sup_{j: m_j \in G, x \in B(m_j, 2\Delta)} \mathbb{P}_x(m_{I_l} \in G \ \forall l \in [un_{\min}], K_i(\eta, \Delta) > u \cdot n_{\min}) \\
& \leq \sum_{v=1}^u n_{\min} \delta \left( 1 - \left( \frac{p^*}{2} \right)^{n_{\min}} \right)^{v-1} + \left( 1 - \left( \frac{p^*}{2} \right)^{n_{\min}} \right)^u \\
& \leq \frac{n_{\min} \delta}{(p^*/2)^{n_{\min}}} + \left( 1 - \left( \frac{p^*}{2} \right)^{n_{\min}} \right)^u \\
& \leq \epsilon + \left( 1 - \left( \frac{p^*}{2} \right)^{n_{\min}} \right)^u
\end{aligned}$$

uniformly for all  $u = 1, 2, \dots$ . To conclude the proof, it suffices to set  $p = (\frac{p^*}{2})^{n_{\min}}$ .  $\square$

The proof above can be easily adapted to the case when the communication class  $G$  is transient. Define

$$K_i^G(\eta, \Delta) \triangleq \min\{k \geq 0 : I_k(\eta, \Delta) = i \text{ or } m_{I_k(\eta, \Delta)} \notin G\}.$$

**Lemma C.5.** *Assume that  $G$  is a transient communication class on the graph  $\mathcal{G}$ . Then there exists some constant  $p > 0$  such that for any  $i$  with  $m_i \in G$  and any  $\Delta \in (0, \bar{\epsilon}/3)$ ,*

$$\sup_{j: m_j \in G; x \in B(m_j, 2\Delta)} \mathbb{P}_x(K_i^G(\eta, \Delta) > u \cdot n_{\min}) \leq \mathbb{P}(\text{Geom}(p) \geq u) \quad \forall u = 1, 2, \dots \quad (\text{C.15})$$

hold for all  $\eta > 0$  sufficiently small.

*Proof.* The structure of this proof is analogous to that of Lemma C.4. Again, the claim is trivial if, for the initial condition, we have  $x \in B(m_i, 2\Delta)$ . Next, let us observe the following facts.

- Define (recall the definitions of measure  $\mu_i$  and sets  $E_i, E_{i,j}$  in (4)(8)(9))

$$J(j) \triangleq \arg \min_{\tilde{j}: \mu_i(E_{j,\tilde{j}}) > 0} |i - \tilde{j}| \quad \forall j \neq i$$

$$p^* \triangleq \min_{j: j \neq i, m_j \in G} \frac{\mu_j(E_{j,J(j)})}{\mu_j(E_j)}.$$

- From the definition of  $J(j)$  and the fact that there are only finitely many attraction fields we can see that  $p^* > 0$ . Moreover,  $G$  being a communication class implies that

$$|J(j) - i| < |j - i| \quad \forall j \neq i, m_j \in G.$$

Indeed, if  $i < j$ , then since  $G$  is a communication class and there are some  $m_i \in G$  with  $i < j$ , we will at least have  $\mu_j(E_{j,j-1}) > 0$ , so  $|J(j) - i| \leq |i - j| - 1$ ; the case that  $i > j$  can be approached analogously.

- Now from the definition of  $J(j)$  and Proposition B.24, together with the previous bullet point, we know that for all  $\eta$  sufficiently small,

$$\inf_{x \in [-L, L]} \mathbb{P}_x(|I_{k+1} - i| \leq |I_k - i| - 1, m_{I_{k+1}} \in G \mid K_i^G(\eta, \Delta) > k \geq p^*/2)$$

uniformly for all  $k \geq 0$ .

- Repeat this argument for  $n_{\min}$  times, and we can see that for all  $\eta$  sufficiently small

$$\inf_{x \in [-L, L]} \mathbb{P}_x(K_i^G(\eta, \Delta) \leq k + n_{\min} \mid K_i^G(\eta, \Delta) > k) \geq \left(\frac{p^*}{2}\right)^{n_{\min}}$$

uniformly for all  $k \geq 1$ .

- Lastly, for any  $j \neq i$  with  $m_j \in G$  and  $x \in B(m_j, 2\Delta)$  and any  $u = 1, 2, \dots$ ,

$$\begin{aligned} & \mathbb{P}_x(K_i^G(\eta, \Delta) > u \cdot n_{\min}) \\ &= \prod_{v=0}^{u-1} \mathbb{P}_x\left(K_i^G(\eta, \Delta) > (v+1)n_{\min} \mid K_i^G(\eta, \Delta) > v \cdot n_{\min}\right) \\ &= \prod_{v=0}^{u-1} \left(1 - \mathbb{P}_x\left(K_i^G(\eta, \Delta) \leq (v+1)n_{\min} \mid K_i^G(\eta, \Delta) > v \cdot n_{\min}\right)\right) \end{aligned}$$

In summary, now we can see that (for sufficiently small  $\eta$ )

$$\sup_{j: m_j \in G, x \in B(m_j, 2\Delta)} \mathbb{P}_x(K_i^G(\eta, \Delta) \geq u \cdot n_{\min}) \leq \left(1 - \left(\frac{p^*}{2}\right)^{n_{\min}}\right)^u$$

uniformly for all  $u = 1, 2, \dots$ . To conclude the proof, it suffices to set  $p = \left(\frac{p^*}{2}\right)^{n_{\min}}$ .  $\square$

We are now ready to prove Lemma 9, which, as demonstrated earlier, is the key tool in proof of Theorem 2.

*Proof of Lemma 9.* The claim is trivial if  $l^{\text{large}} = 1$ , so we focus on the case where  $l^{\text{large}} \geq 2$ . Fix some

$$\tilde{\gamma} \in (0, (1 - \frac{1}{\alpha}) \wedge (\frac{1}{2})), \quad \beta \in (1, (2 - 2\tilde{\gamma}) \wedge (\alpha - \alpha\tilde{\gamma})), \quad \gamma \in (0, \frac{\tilde{\gamma}}{16Mc_2} \wedge \frac{\tilde{\gamma}}{4}).$$

Let  $q^* = \max_j \mu_j(E_j(0))$ . We show that for any  $t \in (0, \frac{\delta}{4q^*})$  the claim is true.

Now we only consider  $\Delta \in (0, \bar{\epsilon}/3)$  and  $\eta$  small enough so that  $\eta M \leq \eta^\gamma$  and  $\eta^\gamma < \Delta$ . Consider the following stopping times

$$\begin{aligned} T_{\text{escape}}^\gamma &\triangleq \min\{n \geq 0 : X_n^\eta \notin \cup_j B(s_j, 2\eta^\gamma)\}; \\ T_{\text{return}}^\gamma &\triangleq \min\{n \geq 0 : X_n^\eta \in \cup_j B(m_j, 2\eta^\gamma)\}. \end{aligned}$$

First, from Lemma C.3, we know that

$$\sup_{x \in [-L, L]} \mathbb{P}_x(T_{\text{escape}}^\gamma > 1/H(1/\eta)) < \delta/2$$

for all  $\eta$  sufficiently small. Besides, by combining Lemma C.2 with Markov property (applied at  $T_{\text{escape}}^\gamma$ ), we have

$$\sup_{x \in [-L, L]} \mathbb{P}_x\left(T_{\text{return}}^\gamma - T_{\text{escape}}^\gamma > 2c_2\gamma \log(1/\eta)/\eta \mid T_{\text{escape}}^\gamma \leq \frac{1}{H(1/\eta)}\right) < \delta/2$$

for all  $\eta$  sufficiently small. Therefore, for all  $\eta$  sufficiently small,

$$\sup_{x \in [-L, L]} \mathbb{P}\left(T_{\text{return}}^\gamma > \frac{1}{H(1/\eta)} + 2c_2\gamma \frac{\log(1/\eta)}{\eta}\right) < \delta. \quad (\text{C.16})$$

Let  $J$  be the unique index such that  $X_{T_{\text{return}}^\gamma}^\eta \in \Omega_J$ . Our next goal is to show that, almost always, the SGD iterates will visit the local minimum at some *large* attraction fields. Therefore, without loss of generality, we can assume that  $m_J \notin M^{\text{large}}$ , and define

$$T_{\text{large}}^\gamma \triangleq \min\{n \geq T_{\text{return}}^\gamma : X_n^\eta \in \bigcup_{i: m_i \in M^{\text{large}}} B(m_i, 2\Delta)\}$$

and introduce the following definitions:

$$\begin{aligned} \tau_0 &\triangleq T_{\text{return}}^\gamma, \quad J_0 \triangleq J \\ \tau_k &\triangleq \min\{n > \tau_{k-1} : X_n^\eta \in \bigcup_{j \neq J_{k-1}} B(m_j, 2\Delta)\} \\ J_k &= j \Leftrightarrow X_{\tau_k}^\eta \in \Omega_j \quad \forall k \geq 1 \\ K &\triangleq \min\{k \geq 0 : m_{J_k} \in M^{\text{large}}\}. \end{aligned}$$

In other words, the sequence of stopping times  $(\tau_k)_{k \geq 1}$  is the time that, starting from  $T_{\text{return}}^\gamma$ , the SGD iterates visited a local minimum that is different from the one visited at  $\tau_{k-1}$ , and  $(J_k)_{k \geq 0}$  records the label of the visited local minima. The random variable  $K$  is the number of transitions required to visit a minimizer in a *large* attraction field. From Lemma C.4, we know the existence of some  $p^* > 0$  such that (for all  $\eta$  sufficiently small)

$$\sup_{x \in [-L, L]} \mathbb{P}_x(K \geq u \cdot n_{\min}) \leq \mathbb{P}\left(\text{Geom}(p^*) \geq u\right) + \frac{\delta}{2} \quad \forall u = 1, 2, 3, \dots$$

where  $\text{Geom}(a)$  is a Geometric random variable with success rate  $a \in (0, 1)$ . Therefore, one can find integer  $N(\delta)$  such that (for all sufficiently small  $\eta$ )

$$\sup_{x \in [-L, L]} \mathbb{P}(K \geq N(\delta)) \leq \delta. \quad (\text{C.17})$$

Next, given results in Proposition B.24 and the fact that there are only finitely many attraction fields, one can find a real number  $u(\delta)$  such that (for all sufficiently small  $\eta$ )

$$\sup_{x \in [-L, L]} \mathbb{P}_x(\tau_k - \tau_{k-1} \leq \frac{u(\delta)}{\lambda_{J_{k-1}}(\eta)}) \leq \delta/N(\delta) \quad (\text{C.18})$$

uniformly for all  $k = 1, 2, \dots, N(\delta)$ . From (C.16), (C.17), (C.18), we now have

$$\sup_{x \in [-L, L]} \mathbb{P}_x \left( X_n^\eta \notin \bigcup_{j: m_j \in M^{\text{large}}} B(m_j, 2\Delta) \ \forall n \leq N(\delta)u(\delta) \frac{H(1/\eta)/\eta}{\lambda^{\text{large}}(\eta)} + \frac{1}{H(1/\eta)} + 2c_2\gamma \frac{\log(1/\eta)}{\eta} \right) \leq 3\delta \quad (\text{C.19})$$

for any sufficiently small  $\eta$ . To conclude the proof we just observe the following facts. First, due to  $H \in \mathcal{RV}_{-\alpha}$  and  $l^{\text{large}} \geq 2$ , we have

$$\lim_{\eta \downarrow 0} H(1/\eta)/\eta = 0, \quad \lim_{\eta \downarrow 0} \frac{\lambda^{\text{large}}(\eta)}{H(1/\eta)} = 0, \quad \lim_{\eta \downarrow 0} \frac{\log(1/\eta)}{\eta} \lambda^{\text{large}}(\eta) = 0.$$

Therefore, for sufficiently small  $\eta$ , we will have (note that  $\epsilon, \delta$  are fixed constants in this proof, so  $N(\delta), u(\delta)$  are also fixed)

$$\frac{N(\delta)u(\delta) \frac{H(1/\eta)/\eta}{\lambda^{\text{large}}(\eta)} + \frac{1}{H(1/\eta)} + 2c_2\gamma \frac{\log(1/\eta)}{\eta}}{\lfloor t/\lambda^{\text{large}}(\eta) \rfloor} \leq \epsilon. \quad (\text{C.20})$$

Second, recall that we fixed some  $t \in (0, \frac{\delta}{4q^*})$  where  $q^* = \max_j \mu_j(E_j)$ . Also, choose some  $C > 0$  small enough so that

$$C < \delta/2, \quad 2(1+C)^2 < 4.$$

From Proposition B.22 and the fact that there are only finitely many attraction fields, there exists some  $\bar{\eta}_0 > 0$  such that for any  $\eta \in (0, \bar{\eta}_0)$  and any  $\Delta > 0$  sufficiently small,

$$\begin{aligned} & \sup_{i: m_i \in M^{\text{large}}} \sup_{x \in [m_i - 2\Delta, m_i + 2\Delta]} \mathbb{P}_x \left( \sigma_i(\eta) \leq \frac{t}{\lambda^{\text{large}}(\eta)} \right) \\ & \leq \sup_{i: m_i \in M^{\text{large}}} \sup_{x \in [m_i - 2\Delta, m_i + 2\Delta]} \mathbb{P}_x \left( \mu_i(E_i) \lambda^{\text{large}}(\eta) \sigma_i(\eta) \leq q^* t \right) \\ & \leq C + 2(1+C)^2 q^* t \leq 2\delta. \end{aligned}$$

Combine this bound with Markov property (applied at  $\tau_K$ ), and we obtain that

$$\sup_{x \in [-L, L]} \mathbb{P}_x \left( \exists n \in [\lfloor t/\lambda^{\text{large}}(\eta) \rfloor] \text{ s.t. } X_{n+\tau_K}^\eta \notin \bigcup_{i: m_i \in M^{\text{large}}} \Omega_i \right) \leq 2\delta$$

for all  $\eta$  sufficiently small. Together with (C.19)(C.20), we have shown that

$$\sup_{x \in [-L, L]} \mathbb{P}_x \left( V^{\text{small}}(\eta, \epsilon, t) > \epsilon \right) \leq 5\delta$$

holds for all  $\eta$  sufficiently small.  $\square$

## C.2 Proof of Lemma 10, 11

We shall return to the discussion about the dynamics of SGD iterates on a communication class  $G$ . Recall that

$$G^{\text{large}} = \{m_1^{\text{large}}, \dots, m_{i_G}^{\text{large}}\}, \quad G^{\text{small}} = \{m_1^{\text{small}}, \dots, m_{i'_G}^{\text{small}}\}.$$

If  $X_n^\eta$  is initialized at some sharp minimum on  $G$ , then we are interested in the behavior of  $X_n^\eta$  at the first visit to some large attraction fields on  $G$ . Define

$$T_G(\eta, \Delta) \triangleq \min\{n \geq 0 : X_n^\eta \in \bigcup_{i: m_i \in G^{\text{large}}} B(m_i, 2\Delta) \text{ or } X_n^\eta \notin \cup_{i: m_i \in G} \Omega_i\}. \quad (\text{C.21})$$

Not only is this definition of  $T_G$  analogous to the one for  $T_G^{DTMC}$  in (34), but, as illustrated in the next lemma,  $T_G$  also behaves similarly as  $T_G$  on a communication class  $G$  in the following sense: the probabilities  $p_{i,j}$  defined in (35) govern the dynamics regarding which large attraction field on  $G$  is the first one to be visited. Besides,  $T_G$  is usually rather small, meaning that the SGD iterates would efficiently arrive at a large attraction field on  $G$  or simply escape from  $G$ .

**Lemma C.6.** *Given any  $\theta \in (0, (\alpha-1)/2)$ ,  $\epsilon \in (0, 1)$ ,  $i, j \in [n_{\min}]$  such that  $m_i \in G^{\text{small}}$ ,  $m_j \notin G^{\text{large}}$ , the following claims hold for all  $\Delta > 0$  that is sufficiently small:*

$$\begin{aligned} \limsup_{\eta \downarrow 0} \sup_{x \in B(m_i, 2\Delta)} \mathbb{P}_x \left( T_G(\eta, \Delta) \leq \frac{\eta^\theta}{\lambda_G(\eta)}, \quad X_{T_G}^\eta \in B(m_j, 2\Delta) \right) &\leq p_{i,j} + 5\epsilon, \\ \liminf_{\eta \downarrow 0} \inf_{x \in B(m_i, 2\Delta)} \mathbb{P}_x \left( T_G(\eta, \Delta) \leq \frac{\eta^\theta}{\lambda_G(\eta)}, \quad X_{T_G}^\eta \in B(m_j, 2\Delta) \right) &\geq p_{i,j} - 5\epsilon, \\ \limsup_{\eta \downarrow 0} \sup_{x \in B(m_i, 2\Delta)} \mathbb{P}_x \left( T_G(\eta, \Delta) > \frac{\eta^\theta}{\lambda_G(\eta)} \right) &\leq 2\epsilon. \end{aligned}$$

*Proof.* For  $G^{\text{small}} \neq \emptyset$  to hold (and the discussion to be meaningful), we must have  $l_G^* \geq 2$ . Throughout the proof, we assume this is the case. Besides, we require that  $\Delta \in (0, \bar{\epsilon}/3)$  so we have

$$B(m_i, 3\Delta) \cap \Omega_i^c = \emptyset \quad \forall i \in [n_{\min}]$$

and the  $3\Delta$ -neighborhood of each local minimum will not intersect with each other. In this proof we will only consider  $\Delta$  in this range.

From Lemma C.4 (if  $G$  is absorbing) or Lemma C.5 (if  $G$  is transient), we know the existence of some integer  $N(\epsilon)$  such that for (see the definition of  $I_k$  in (C.11)-(C.14))

$$N_G(\eta, \Delta) \triangleq \min\{k \geq 0 : m_{I_k(\eta, \Delta)} \in G^{\text{large}} \text{ or } m_{I_k(\eta, \Delta)} \notin G\},$$

we have

$$\sup_{x \in B(m_i, 2\Delta)} \mathbb{P}_x \left( N_G(\eta, \Delta) > N(\epsilon) \right) < \epsilon$$

for all  $\eta$  sufficiently small. Fix such  $N(\epsilon)$ . Next, from Proposition B.24, we can find  $u(\epsilon) \in (0, \infty)$  and  $\bar{\Delta} \in (0, \bar{\epsilon}/3)$  such that for all  $\Delta \in (0, \bar{\Delta})$ , we have

$$\sup_{x \in B(m_i, 2\Delta)} \mathbb{P}_x \left( T_k(\eta, \Delta) - T_{k-1}(\eta, \Delta) > u(\epsilon)/\Lambda(I_{k-1}(\eta, \Delta), \eta) \right) \leq \epsilon/N(\epsilon) \quad \forall k \in [N(\epsilon)]$$

for all  $\eta$  sufficiently small. Fix such  $u(\epsilon)$  and  $\bar{\Delta}$ . Now note that on the event

$$A \triangleq \left\{ N_G \leq N(\epsilon) \right\} \cap \left\{ T_k(\eta, \Delta) - T_{k-1}(\eta, \Delta) \leq u(\epsilon)/\Lambda(I_{k-1}(\eta, \Delta), \eta) \quad \forall k \in [N(\epsilon)] \right\},$$



due to the choice of  $\theta \in (0, (\alpha - 1)/2)$  and  $H \in \mathcal{RV}_{-\alpha}$ , we have (when  $\eta \in (0, 1)$ )

$$\begin{aligned} T_k(\eta, \Delta) - T_{k-1}(\eta, \Delta) &\leq \frac{\eta^{2\theta}}{\lambda_G(\eta)} \quad \forall k < N_G(\eta, \Delta) \\ \Rightarrow T_G(\eta, \Delta) &= T_{N_G(\eta, \Delta)}(\eta, \Delta) \leq N(\epsilon)u(\epsilon) \frac{\eta^{2\theta}}{\lambda_G(\eta)}. \end{aligned}$$

For any  $\eta$  sufficiently small, we will have  $N(\epsilon)u(\epsilon) \frac{\eta^{2\theta}}{\lambda_G(\eta)} < \frac{\eta^\theta}{\lambda_G(\eta)}$ . In summary, we have established that for all  $\Delta \in (0, \bar{\Delta})$ ,

$$\limsup_{\eta \downarrow 0} \sup_{x \in B(m_i, 2\Delta)} \mathbb{P}_x \left( T_G > \frac{\eta^\theta}{\lambda_G(\eta)} \right) \leq \limsup_{\eta \downarrow 0} \sup_{x \in B(m_i, 2\Delta)} \mathbb{P}_x(A^c) < 2\epsilon. \quad (\text{C.22})$$

Next, let

$$\mathbf{S}(\epsilon) \triangleq \left\{ (m'_1, \dots, m'_{N(\epsilon)}) \in \{m_1, \dots, m_{n_{\min}}\}^{N(\epsilon)} : \exists k \in [N(\epsilon)] \text{ s.t. } m'_k = m_j \right\}.$$

We can see that  $\mathbf{S}(\epsilon)$  contains all the possible transition path for  $Y^{DTMC}$  where the state  $m_j$  is visited within the first  $N(\epsilon)$  steps. Obviously,  $|\mathbf{S}(\epsilon)| < \infty$ . Let  $\epsilon_1 = \epsilon/|\mathbf{S}(\epsilon)|$ . If we are able to show the existence of some  $\bar{\Delta}_1 > 0$  such that for all  $\Delta \in (0, \bar{\Delta}_1)$ , the following claim holds for any  $(m'_k)_{k=1}^{N(\epsilon)} \in \mathbf{S}(\epsilon)$ :

$$\limsup_{\eta \downarrow 0} \sup_{x \in B(m_i, 2\Delta)} \left| \mathbb{P}_x \left( m_{I_k} = m'_k \quad \forall k \in [N(\epsilon)] \right) - \mathbb{P} \left( Y_k^{DTMC}(m_i) = m'_k \quad \forall k \in [N(\epsilon)] \right) \right| < \epsilon_1, \quad (\text{C.23})$$

then we must have (for all  $\Delta \in (0, \bar{\Delta} \wedge \bar{\Delta}_1)$ )

$$\begin{aligned} &\limsup_{\eta \downarrow 0} \sup_{x \in B(m_i, 2\Delta)} \left| \mathbb{P}_x \left( X_{T_G}^\eta \in B(m_j, 2\Delta) \right) - p_{i,j} \right| \\ &= \limsup_{\eta \downarrow 0} \sup_{x \in B(m_i, 2\Delta)} \left| \mathbb{P}_x \left( X_{T_G}^\eta \in B(m_j, 2\Delta), T_G \leq N(\epsilon) \right) + \mathbb{P}_x \left( X_{T_G}^\eta \in B(m_j, 2\Delta), T_G > N(\epsilon) \right) \right. \\ &\quad \left. - \mathbb{P} \left( Y_{T_G}^{DTMC}(m_i) = m_j, T_G^{DTMC} \leq N(\epsilon) \right) - \mathbb{P} \left( Y_{T_G}^{DTMC}(m_i) = m_j, T_G^{DTMC} > N(\epsilon) \right) \right| \\ &\leq \limsup_{\eta \downarrow 0} \sup_{x \in B(m_i, 2\Delta)} \left| \mathbb{P}_x \left( X_{T_G}^\eta \in B(m_j, 2\Delta), T_G \leq N(\epsilon) \right) - \mathbb{P} \left( Y_{T_G}^{DTMC}(m_i) = m_j, T_G^{DTMC} \leq N(\epsilon) \right) \right| \\ &\quad + \limsup_{\eta \downarrow 0} \sup_{x \in B(m_i, 2\Delta)} \mathbb{P}_x(T_G > N(\epsilon)) + \mathbb{P}(T_G^{DTMC}(m_i) > N(\epsilon)) \\ &\leq |\mathbf{S}(\epsilon)|\epsilon_1 + \limsup_{\eta \downarrow 0} \sup_{x \in B(m_i, 2\Delta)} \mathbb{P}_x(T_G > N(\epsilon)) + \mathbb{P}(T_G^{DTMC}(m_i) > N(\epsilon)) \\ &\leq 3\epsilon. \end{aligned}$$

To show that (C.23) is true, we fix some  $(m'_k)_{k=1}^{N(\epsilon)} \in \mathbf{S}(\epsilon)$  and let  $(\mathbf{k}'(k))_{k=1}^{N(\epsilon)}$  be the sequence with  $m_{\mathbf{k}'(k)} = m'_k$  for each  $k \in [N(\epsilon)]$ . From the definition of  $Y^{DTMC}$  we have (let  $\mathbf{k}'(0) = i$ )

$$\mathbb{P} \left( Y_k^{DTMC}(m_i) = m'_k \quad \forall k \in [N(\epsilon)] \right) = \prod_{k=0}^{N(\epsilon)-1} \frac{\mu_{\mathbf{k}'(k)}(E_{\mathbf{k}'(k), \mathbf{k}'(k+1)})}{\mu_{\mathbf{k}'(k)}(E_{\mathbf{k}'(k)})}.$$

On the other hand, using Proposition B.24, we know that for any arbitrarily chosen  $\epsilon' \in (0, 1)$ , we have

$$\limsup_{\eta \downarrow 0} \sup_{x \in B(m_i, 2\Delta)} \mathbb{P}_x \left( m_{I_k(\eta, \Delta)} = m'_k \quad \forall k \in [N(\epsilon)] \right) \leq \prod_{k=0}^{N(\epsilon)-1} \frac{\mu_{\mathbf{k}'(k)}(E_{\mathbf{k}'(k), \mathbf{k}'(k+1)})}{\mu_{\mathbf{k}'(k)}(E_{\mathbf{k}'(k)})} \cdot (1 + \epsilon'),$$

$$\liminf_{\eta \downarrow 0} \inf_{x \in B(m_i, 2\Delta)} \mathbb{P}_x \left( m_{I_k(\eta, \Delta)} = m'_k \ \forall k \in [N(\epsilon)] \right) \geq \prod_{k=0}^{N(\epsilon)-1} \frac{\mu_{\mathbf{k}'(k)}(E_{\mathbf{k}'(k), \mathbf{k}'(k+1)})}{\mu_{\mathbf{k}'(k)}(E_{\mathbf{k}'(k)})} \cdot (1 - \epsilon'),$$

for all  $\Delta > 0$  sufficiently small. The arbitrariness of  $\epsilon' > 0$ , together with  $|\mathbf{S}(\epsilon)| < \infty$ , allows us to see the existence of some  $\bar{\Delta}_1 > 0$  such that with  $\Delta \in (0, \bar{\Delta}_1)$ , (C.23) holds for any  $(m'_k)_{k=1}^{N(\epsilon)} \in \mathbf{S}(\epsilon)$ . To conclude the proof, observe that

$$\begin{aligned} & \limsup_{\eta \downarrow 0} \sup_{x \in B(m_i, 2\Delta)} \left| \mathbb{P}_x \left( X_{T_G}^\eta \in B(m_j, 2\Delta) \right) - \mathbb{P}_x \left( T_G(\eta, \Delta) \leq \frac{\eta^\theta}{\lambda_G(\eta)}, X_{T_G}^\eta \in B(m_j, 2\Delta) \right) \right| \\ & \leq \limsup_{\eta \downarrow 0} \sup_{x \in B(m_i, 2\Delta)} \mathbb{P}_x \left( T_G > \frac{\eta^\theta}{\lambda_G(\eta)} \right) < \epsilon \end{aligned}$$

due to (C.22).  $\square$

Recall that continuous-time process  $X^{*,\eta}$  is the *scaled* version of  $X^\eta$  defined in (31), and the mapping  $\mathbf{T}^*(n, \eta) \triangleq n\lambda_G(\eta)$  returns the timestamp  $t$  for  $X_t^{*,\eta}$  corresponding to the unscaled step  $n$  on the time horizon of  $X_n^\eta$ . As an *inverse* mapping of  $\mathbf{T}^*$ , we define the mapping  $\mathbf{N}^*(t, \eta) = \lfloor t/\lambda_G(\eta) \rfloor$  that maps the scaled timestamp  $t$  back to the step number  $n$  for the unscaled process  $X^\eta$ .

In the next lemma, we show that, provided that  $X^{*,\eta}$  stays on a communication class  $G$  before some time  $t$ , the scaled process  $X_t^{*,\eta}$  is almost always in the largest attraction fields of a communication class  $G$ .

**Lemma C.7.** *Let  $G$  be a communication class on the graph  $\mathcal{G}$ . Given any  $\epsilon_1 > 0$ ,  $t > 0$  and any  $x \in \Omega_i$  with  $m_i \in G$ , the following claim holds for all  $\Delta > 0$  small enough:*

$$\limsup_{\eta \downarrow 0} \mathbb{P}_x \left( \left\{ X_t^{*,\eta} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, 3\Delta) \right\} \cap \left\{ X_s^{*,\eta} \in \bigcup_{k: m_k \in G} \Omega_k \ \forall s \in [0, t] \right\} \right) \leq 2\epsilon_1.$$

*Proof.* Let  $\Delta \in (0, \bar{\epsilon}/3)$  for the constant  $\bar{\epsilon}$  in (B.15)(B.16), so we are certain that each  $B(m_i, 2\Delta)$  lies entirely in  $\Omega_i$  and would not intersect with each other since

$$B(m_i, 3\Delta) \cap \Omega_i^c = \emptyset \ \forall i \in [n_{\min}].$$

Besides, with  $\epsilon = \Delta/3$ , we know the existence of some  $\delta > 0$  such that claims in Lemma B.12 would hold of the chosen  $\epsilon, \delta$ . Fix such  $\delta$  for the entirety of this proof. Lastly, fix some

$$\tilde{\gamma} \in (0, (1 - \frac{1}{\alpha}) \wedge (\frac{1}{2})), \ \beta \in (1, (2 - 2\tilde{\gamma}) \wedge (\alpha - \alpha\tilde{\gamma})), \ \gamma \in (0, \frac{\tilde{\gamma}}{16Mc_2} \wedge \frac{\tilde{\gamma}}{4}).$$

The blueprint of this proof is as follows. We will define a sequence of stopping times  $(N_j)_{j=1}^6$  such that the corresponding scaled timestamps  $\mathbf{T}_j^* = \mathbf{T}^*(N_j, \eta)$  gradually approach  $t$ . By analyzing the behavior of  $X^{*,\eta}$  on a time interval  $[t - \Delta_t, t]$  that is very close to  $t$  (in particular, on the aforementioned stopping times  $\mathbf{T}_j^*$ ), we are able to establish the properties of a series of events  $A_1 \supseteq A_2 \supseteq A_3$ . Moreover, we will show that  $A_3 \subseteq \{X_t^{*,\eta, \Delta} \in \bigcup_{i: m_i \in G^{\text{large}}} B(m_i, 3\Delta)\}$ , so the properties about events  $A, A_2, A_3$  can be used to bound the probability of the target event.

Arbitrarily choose some  $\Delta_t \in (0, t)$ . To proceed, let  $N_0 \triangleq \mathbf{N}^*(t - \Delta_t, \eta)$  be the stopping time corresponding to timestamp  $t - \Delta_t$  for the scaled process. Using Lemma C.3, we know that for stopping time  $N_1 \triangleq \min\{n \geq N_0 : X_n^\eta \notin \cup_j B(s_j, 2\eta^\gamma)\}$ , we have

$$\liminf_{\eta \downarrow 0} \inf_{x \in [-L, L]} \mathbb{P}_x(N_1 - N_0 < \frac{\Delta_t/4}{H(1/\eta)}) = 1.$$

Next, let  $N_2 \triangleq \min\{n \geq N_1 : X_n^\eta \in \cup_j B(m_j, 2\Delta)\}$ . From Lemma C.2 and  $H \in \mathcal{RV}_{-\alpha}$  (so that  $\log(1/\eta)/\eta = o(H(1/\eta))$ ), we have

$$\liminf_{\eta \downarrow 0} \inf_{x \in [-L, L]} \mathbb{P}_x(N_2 - N_1 < \frac{\Delta_t/4}{H(1/\eta)}) = 1.$$

Collecting results above, we have

$$\liminf_{\eta \downarrow 0} \inf_{x \in [-L, L]} \mathbb{P}_x(N_2 - N_0 < \frac{\Delta_t/2}{H(1/\eta)}) = 1. \quad (\text{C.24})$$

Now note the following fact on the event  $\{N_2 - N_0 < \frac{\Delta_t/2}{H(1/\eta)}\}$ . The definition of the mapping  $\mathbf{T}^*$  implies that, for any pair of positive integers  $n_1 \leq n_2$ , we have  $\mathbf{T}^*(n_2, \eta) - \mathbf{T}^*(n_1, \eta) = (n_2 - n_1)\lambda_G(\eta) \leq (n_2 - n_1) \cdot H(1/\eta)$ . Therefore, on  $\{N_2 - N_0 \leq \frac{\Delta_t/2}{H(1/\eta)}\}$  we have

$$\mathbf{T}^*(N_2, \eta) - \mathbf{T}^*(N_0, \eta) < \Delta_t/2 \Rightarrow \mathbf{T}^*(N_2, \eta) < t - \frac{\Delta_t}{2}.$$

Besides, let  $\mathbf{T}_2^* = \mathbf{T}^*(N_2, \eta)$ . Now we can see that for event

$$A_0 \triangleq \left\{ X_s^{*,\eta} \in \bigcup_{k: m_k \in G} \Omega_k \ \forall s \in [0, t] \right\} \cap \left\{ N_2 - N_0 < \frac{\Delta_t/2}{H(1/\eta)} \right\},$$

we have

$$A_0 \subseteq A_1 \triangleq \left\{ \mathbf{T}_2^* < t - \frac{\Delta_t}{2} \right\} \cap \left\{ X_s^{*,\eta} \in \bigcup_{k: m_k \in G} \Omega_k \ \forall s \in [0, \mathbf{T}_2^*] \right\}.$$

Meanwhile, from (C.24) we obtain that

$$\limsup_{\eta \downarrow 0} \sup_{x \in [-L, L]} \mathbb{P}_x \left( A_1^c \cap \left\{ X_s^{*,\eta} \in \bigcup_{k: m_k \in G} \Omega_k \ \forall s \in [0, t] \right\} \right) = 0. \quad (\text{C.25})$$

Moving on, we consider the following stopping times

$$N_3 \triangleq \min\{n \geq N_2 : X_n^\eta \in \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, 2\Delta) \text{ or } X_n^\eta \notin \bigcup_{j: m_j \in G} \Omega_j\},$$

$$\mathbf{T}_3^* \triangleq \mathbf{T}^*(N_3, \eta).$$

Using Lemma C.6, we have

$$\limsup_{\eta \downarrow 0} \mathbb{P}_x \left( N_3 - N_2 > \frac{\Delta_t/4}{\lambda_G(\eta)} \mid A_1 \right) \leq \epsilon_1. \quad (\text{C.26})$$

Meanwhile, on event  $A_1 \cap \left\{ N_3 - N_2 \leq \frac{\Delta_t/4}{\lambda_G(\eta)} \right\} \cap \left\{ X_s^{*,\eta} \in \bigcup_{k: m_k \in G} \Omega_k \ \forall s \in [0, t] \right\}$ , we have  $\mathbf{T}_3^* - \mathbf{T}_2^* \leq \Delta_t/4$ , hence  $\mathbf{T}_3^* \in [t - \Delta_t, t - \Delta_t/4]$ . In summary,

$$\begin{aligned} & A_1 \cap \left\{ N_3 - N_2 \leq \frac{\Delta_t/4}{\lambda_G(\eta)} \right\} \cap \left\{ X_s^{*,\eta} \in \bigcup_{k: m_k \in G} \Omega_k \ \forall s \in [0, t] \right\} \\ & \subseteq \left\{ \mathbf{T}_3^* \in [t - \Delta_t, t - \Delta_t/4] \right\} \cap \left\{ X_s^{*,\eta} \in \bigcup_{k: m_k \in G} \Omega_k \ \forall s \in [0, \mathbf{T}_3^*] \right\} \end{aligned}$$

Moreover, on event  $\left\{X_s^{*,\eta} \in \bigcup_{k: m_k \in G} \Omega_k \ \forall s \in [0, t]\right\}$ , if we let  $J_3$  be label of the local minimum visited at  $\mathbf{T}_3^*$  such that  $J_3 = j \iff X_{T_3^*}^{*,\eta} \in B(m_j, 2\Delta)$ , then we must have  $m_{J_3} \in G^{\text{large}}$ . Meanwhile, consider the following stopping times

$$\mathbf{T}^\sigma \triangleq \min\{s > \mathbf{T}_3^* : X_s^{*,\eta} \notin \Omega_{J_3}\}.$$

From Proposition B.24, we know that

$$\begin{aligned} & \limsup_{\eta \downarrow 0} \mathbb{P}_x \left( \mathbf{T}^\sigma - \mathbf{T}_3^* \leq \Delta_t \mid \left\{ \mathbf{T}_3^* \in [t - \Delta_t, t - \Delta_t/4] \right\} \cap \left\{ X_s^{*,\eta} \in \bigcup_{k: m_k \in G} \Omega_k \ \forall s \in [0, \mathbf{T}_3^*] \right\} \right) \\ & \leq \epsilon_1 + 1 - \exp \left( - (1 + \epsilon_1) q^* \Delta_t \right) \end{aligned} \quad (\text{C.27})$$

where  $q^* = \max_j \mu_j(E_j)$ . Now we define the event

$$A_2 \triangleq A_1 \cap \left\{ N_3 - N_2 \leq \frac{\Delta_t/4}{\lambda_G(\eta)} \right\} \cap \left\{ \mathbf{T}^\sigma - \mathbf{T}_3^* > \Delta_t \right\} \cap \left\{ X_s^{*,\eta} \in \bigcup_{k: m_k \in G} \Omega_k \ \forall s \in [0, \mathbf{T}_3^*] \right\}.$$

Using (C.25)-(C.27), we get

$$\limsup_{\eta \downarrow 0} \sup_{x \in [-L, L]} \mathbb{P}_x \left( A_2^c \cap \left\{ X_s^{*,\eta} \in \bigcup_{k: m_k \in G} \Omega_k \ \forall s \in [0, t] \right\} \right) \leq 2\epsilon_1 + 1 - \exp \left( - (1 + \epsilon_1) q^* \Delta_t \right). \quad (\text{C.28})$$

Furthermore, on event  $A_2$ , due to  $\mathbf{T}_3^* \in [t - \Delta_t, t - \Delta_t/4]$  as established above, we must have

$$X_s^{*,\eta} \in \Omega_{J_3} \ \forall s \in [\mathbf{T}_3^*, t].$$

Now let us focus on a timestamp  $\mathbf{T}_4^* = t - \frac{\Delta_t/8}{H(1/\eta)} \lambda_G(\eta)$  and  $N_4 = \mathbf{N}^*(\mathbf{T}_4^*, \eta)$ . Obviously,  $\mathbf{T}_4^* > \mathbf{T}_3^*$  on event  $A_2$ . Next, define

$$\begin{aligned} N_5 & \triangleq \min\{n \geq N_4 : X_n^\eta \in \bigcup_j B(m_j, 2\Delta)\} \\ \mathbf{T}_5^* & \triangleq \mathbf{T}^*(N_5, \eta). \end{aligned}$$

Using Lemma C.2 and C.3 again as we did above when obtaining (C.24), we can show that

$$\limsup_{\eta \downarrow 0} \mathbb{P}_x \left( N_5 - N_4 > \frac{\Delta_t/16}{H(1/\eta)} \right) = 0. \quad (\text{C.29})$$

On the other hand, on event  $A_2 \cap \left\{ N_5 - N_4 \leq \frac{\Delta_t/16}{H(1/\eta)} \right\}$  we must have

- $\mathbf{T}_5^* - \mathbf{T}_4^* \leq \frac{\Delta_t/16}{H(1/\eta)} \lambda_G(\eta)$ , so  $\mathbf{T}_5^* \in [t - \frac{\Delta_t/8}{H(1/\eta)} \lambda_G(\eta), t - \frac{\Delta_t/16}{H(1/\eta)} \lambda_G(\eta)]$ ;
- $X_{\mathbf{T}_5^*}^{*,\eta} \in \Omega_{J_3}$ , due to  $\mathbf{T}^\sigma - \mathbf{T}_3^* > \Delta_t$ .

This implies that for event

$$\tilde{A} \triangleq \left\{ \mathbf{T}_5^* \in [t - \frac{\Delta_t/8}{H(1/\eta)} \lambda_G(\eta), t - \frac{\Delta_t/16}{H(1/\eta)} \lambda_G(\eta)], X_{\mathbf{T}_5^*}^{*,\eta} \in \Omega_{J_3} \right\} \cap \left\{ X_s^{*,\eta} \in \bigcup_{k: m_k \in G} \Omega_k \ \forall s \in [0, \mathbf{T}_5^*] \right\},$$

we have  $A_2 \cap \left\{ N_5 - N_4 \leq \frac{\Delta_t/16}{H(1/\eta)} \right\} \subseteq \tilde{A}$ . Lastly, observe that

- From Lemma B.2, we know that for  $N_6(\delta) \triangleq \min\{n > N_5 : \eta|Z_n| > \delta\}$  we have

$$\limsup_{\eta \downarrow 0} \mathbb{P}(N_6(\delta) - N_5 \leq \Delta_t/H(1/\eta)) \leq \Delta_t/\delta^\alpha;$$

- As stated at the beginning of the proof, our choice of  $\delta$  allows us to apply Lemma B.12 and show that

$$\limsup_{\eta \downarrow 0} \sup_{x \in [-L, L]} \mathbb{P}_x(\exists n = N_5, \dots, N_6 - 1 \text{ s.t. } X_n^\eta \notin B(m_{J_3}, 3\Delta) \mid \tilde{A}) = 0;$$

- Combining the two bullet points above, we get

$$\limsup_{\eta \downarrow 0} \mathbb{P}_x(\exists s \in [\mathbf{T}_5^*, t] \text{ such that } X_s^{*,\eta} \notin B(m_{J_3}, 3\Delta) \mid \tilde{A}) \leq \Delta_t / \delta^\alpha. \quad (\text{C.30})$$

On the other hand,

$$\tilde{A} \cap \{X_s^{*,\eta} \in B(m_{J_3}, 3\Delta) \mid \forall s \in [\mathbf{T}_5^*, t]\} \subseteq \{X_t^{*,\eta} \in \bigcup_{k: m_k \in G^{\text{large}}} B(m_k, 3\Delta)\}.$$

In summary, for event

$$A_3 \triangleq A_2 \cap \left\{ N_5 - N_4 \leq \frac{\Delta_t / 16}{H(1/\eta)} \right\} \cap \left\{ X_s^{*,\eta} \in B(m_{J_3}, 3\Delta) \mid \forall s \in [\mathbf{T}_5^*, t] \right\},$$

we have  $A_3 \subseteq \{X_t^{*,\eta} \in \bigcup_{k: m_k \in G^{\text{large}}} B(m_k, 3\Delta)\}$ . Besides, due to (C.28)(C.29)(C.30), we get

$$\begin{aligned} & \limsup_{\eta \downarrow 0} \sup_{x \in [-L, L]} \mathbb{P}_x \left( A_3^c \cap \{X_s^{*,\eta} \in \bigcup_{k: m_k \in G} \Omega_k \mid \forall s \in [0, t]\} \right) \\ & \leq 2\epsilon_1 + 1 - \exp(-(1 + \epsilon_1)q^*\Delta_t) + \frac{\Delta_t}{\delta^\alpha}. \end{aligned}$$

Remember that  $\delta, \epsilon_1, q^*$  are fixed constants while  $\Delta_t$  can be made arbitrarily small, so by driving  $\Delta_t$  to 0 we can conclude the proof.  $\square$

Recall the definition of jump processes in Definition 2. Central to the proof of Lemma 10, the next result provides a set of sufficient conditions for the convergence of a sequence of such jump processes in the sense of finite dimensional distributions.

**Lemma C.8.** *For a sequence of processes  $(Y^n)_{n \geq 1}$  that, for each  $n \geq 1$ ,  $Y^n$  is a  $((U_j^n)_{j \geq 0}, (V_j^n)_{j \geq 0})$  jump process, and a  $((U_j)_{j \geq 0}, (V_j)_{j \geq 0})$  jump process  $Y$ , if*

- $U_0 \equiv 0$ ;
- $(U_0^n, V_0^n, U_1^n, V_1^n, U_2^n, V_2^n, \dots)$  converges in distribution to  $(0, V_0, U_1, V_1, U_2, V_2, \dots)$  as  $n \rightarrow \infty$ ;
- For any  $x > 0$  and any  $n \geq 1$ ,

$$\mathbb{P}(U_1 + \dots + U_n = x) = 0;$$

- For any  $x > 0$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}(U_1 + U_2 + \dots + U_n > x) = 1,$$

then the finite dimensional distribution of  $Y^n$  converges to that of  $Y$  in the following sense: for any  $k \in \mathbb{N}$  and any  $0 < t_1 < t_2 < \dots < t_k < \infty$ , the random element  $(Y_{t_1}^n, \dots, Y_{t_k}^n)$  converges in distribution to  $(Y_{t_1}, \dots, Y_{t_k})$  as  $n \rightarrow \infty$ .

*Proof.* Fix some  $k \in \mathbb{N}$  and  $0 < t_1 < t_2 < \dots < t_k < \infty$ . For notational simplicity, let  $t = t_k$ . Let  $(\mathbb{D}, \mathbf{d})$  be the metric space where  $\mathbb{D} = \mathbb{D}_{[0,t]}$ , the space of all càdlàg functions in  $\mathbb{R}$  on the time interval  $[0, t]$ , and  $\mathbf{d}$  is the Skorokhod metric defined as

$$\mathbf{d}(\zeta_1, \zeta_2) \triangleq \inf_{\lambda \in \Lambda} \|\zeta_1 - \zeta_2 \circ \lambda\| \vee \|\lambda - I\|$$

where  $\Lambda$  is the set of all nondecreasing homeomorphism from  $[0, t]$  onto itself, and  $I(s) = s$  is the identity mapping. Also, we arbitrarily choose some  $\epsilon \in (0, 1)$  and some open set  $A \subseteq \mathbb{R}^k$ .

From the assumption, we can find integer  $J(\epsilon)$  such that  $\mathbb{P}\left(\sum_{j=1}^{J(\epsilon)} U_j \leq t\right) < \epsilon$ , as well as an integer  $N(\epsilon)$  such that, for all  $n \geq N(\epsilon)$ , we have  $\mathbb{P}\left(\sum_{j=1}^{J(\epsilon)} U_j^n \leq t\right) + \mathbb{P}\left(U_0^n \geq t_1\right) < \epsilon$ . We fix such  $J(\epsilon), N(\epsilon)$  (we may abuse the notation slightly and simply write  $J, N$  when there is no ambiguity).

Using Skorokhod's representation theorem, we can construct a probability space  $(\Omega, \mathcal{F}, \mathbb{Q})$  that supports random variables  $(\tilde{U}_0^n, \tilde{V}_0^n, \dots, \tilde{U}_J^n, \tilde{V}_J^n)_{n \geq 1}$  and  $(\tilde{U}_0, \tilde{V}_0, \dots, \tilde{U}_J, \tilde{V}_J)$  and satisfies the following conditions:

- $\mathcal{L}(U_0^n, V_0^n, \dots, U_J^n, V_J^n) = \mathcal{L}(\tilde{U}_0^n, \tilde{V}_0^n, \dots, \tilde{U}_J^n, \tilde{V}_J^n)$  for all  $n \geq 1$ ;
- $\mathcal{L}(U_0, V_0, \dots, U_J, V_J) = \mathcal{L}(\tilde{U}_0, \tilde{V}_0, \dots, \tilde{U}_J, \tilde{V}_J)$ ;
- $U_j^n \xrightarrow{a.s.} U_j$  and  $V_j^n \xrightarrow{a.s.} V_j$  as  $n \rightarrow \infty$  for all  $j \in [J]$ .

Therefore, on  $(\Omega, \mathcal{F}, \mathbb{Q})$  we can define the following random elements (taking values in the space of càdlàg functions):

$$Y_s^{n, \downarrow J} = \begin{cases} \tilde{V}_0^n & \text{if } s < \tilde{U}_0^n \\ \sum_{j=0}^J \tilde{V}_j^n \mathbb{1}_{[\tilde{U}_0^n + \tilde{U}_1^n + \dots + \tilde{U}_j^n, \tilde{U}_0^n + \tilde{U}_1^n + \dots + \tilde{U}_{j+1}^n)}(s) & \text{otherwise} \end{cases},$$

$$Y_s^{\downarrow J} = \sum_{j=0}^J \tilde{V}_j \mathbb{1}_{[\tilde{U}_1 + \dots + \tilde{U}_j, \tilde{U}_1 + \dots + \tilde{U}_{j+1})}(s) \quad \forall s \geq 0.$$

Note that (1) for the first jump time of  $Y^{\downarrow J}$  we have  $\tilde{U}_0 \equiv 0$ , hence  $Y_0^{\downarrow J} = \tilde{V}_0$ ; (2) when defining  $Y^{n, \downarrow J}$  we set its value on  $[0, \tilde{U}_0^n)$  to be  $\tilde{V}_0^n$  instead of 0.

Since  $U_j^n \xrightarrow{a.s.} U_j$  and  $V_j^n \xrightarrow{a.s.} V_j$  as  $n \rightarrow \infty$  for all  $j \in [J]$ , we must have

$$\lim_n \mathbf{d}(Y_s^{n, \downarrow J}, Y_s^{\downarrow J}) = 0$$

almost surely, which further implies that  $Y_s^{n, \downarrow J} \Rightarrow Y_s^{\downarrow J}$  as  $n \rightarrow \infty$  on  $(\mathbb{D}, \mathbf{d})$ . Now from our assumption that, for the jump times  $U_1 + \dots + U_j$ , we have  $\mathbb{P}(U_1 + \dots + U_j = x) = 0 \quad \forall x > 0, j \geq 1$ , as well as (13.3) in [1], we then obtain

$$(Y_{t_1}^{n, \downarrow J}, \dots, Y_{t_k}^{n, \downarrow J}) \Rightarrow (Y_{t_1}^{\downarrow J}, \dots, Y_{t_k}^{\downarrow J}) \quad (\text{C.31})$$

as  $n \rightarrow \infty$ . Recall that  $A$  is the open set we arbitrarily chose at the beginning of the proof, and  $\epsilon > 0$  is also chosen arbitrarily. Now we observe the following facts.

- Using (C.31), we can see that

$$\liminf_n \mathbb{Q}((Y_{t_1}^{n, \downarrow J}, \dots, Y_{t_k}^{n, \downarrow J}) \in A) \geq \mathbb{Q}(Y_{t_1}^{\downarrow J}, \dots, Y_{t_k}^{\downarrow J}) \in A).$$

- The choice of  $N(\epsilon)$  and  $J(\epsilon)$  above implies that

$$|\mathbb{Q}((Y_{t_1}^{\downarrow J(\epsilon)}, \dots, Y_{t_k}^{\downarrow J(\epsilon)}) \in A) - \mathbb{P}((Y_{t_1}, \dots, Y_{t_k}) \in A)| \leq \mathbb{P}\left(\sum_{j=1}^{J(\epsilon)} U_j \leq t\right) < \epsilon$$

$$|\mathbb{Q}((Y_{t_1}^{n, \downarrow J(\epsilon)}, \dots, Y_{t_k}^{n, \downarrow J(\epsilon)}) \in A) - \mathbb{P}((Y_{t_1}^n, \dots, Y_{t_k}^n) \in A)| \leq \mathbb{P}\left(\sum_{j=1}^{J(\epsilon)} U_j^n \leq t\right) + \mathbb{P}(U_0^n \geq t_1) < \epsilon \quad \forall n \geq N(\epsilon).$$

Collecting the two results above, we have established that

$$\liminf_n \mathbb{P}((Y_{t_1}^n, \dots, Y_{t_k}^n) \in A) \geq \mathbb{P}((Y_{t_1}, \dots, Y_{t_k}) \in A) - 2\epsilon.$$

From Portmanteau theorem, together with arbitrariness of  $\epsilon > 0$  and open set  $A$ , we can now conclude that  $(Y_{t_1}^n, \dots, Y_{t_k}^n)$  converges in distribution to  $(Y_{t_1}, \dots, Y_{t_k})$ .  $\square$

The following lemma concerns the scaled version of the marker process  $\hat{X}^{*, \eta, \Delta}$  defined in (31)-(32). Obviously, it is a jump process that complies with Definition 2. When there is no ambiguity about the sequences  $(\eta_n)_{n \geq 1}$  and  $(\Delta_n)_{n \geq 1}$ , let  $\hat{X}_t^{(n)} \triangleq \hat{X}_t^{*, \eta_n, \Delta_n}$ . From (24)-(29) and (33), we know that for any  $n \geq 1$ ,  $\hat{X}^{(n)}$  is a  $\left((\tau_k^*(\eta_n, \Delta_n) - \tau_{k-1}^*(\eta_n, \Delta_n))_{k \geq 0}, (m_{I_k(\eta_n, \Delta_n)})_{k \geq 0}\right)$ -jump process (with the convention that  $\tau_{-1}^* = 0$ ). Also, for clarity of the exposition, we let (for all  $n \geq 1, k \geq 0$ )

$$\begin{aligned} \tilde{S}_k^{(n)} &= \sigma_k^*(\eta_n, \Delta_n) - \tau_{k-1}^*(\eta_n, \Delta_n), \\ S_k^{(n)} &= \tau_k^*(\eta_n, \Delta_n) - \tau_{k-1}^*(\eta_n, \Delta_n), \\ \widetilde{W}_k^{(n)} &= m_{\tilde{I}_k^G(\eta_n, \Delta_n)}, \\ W_k^{(n)} &= m_{I_k^G(\eta_n, \Delta_n)}. \end{aligned}$$

Lastly, remember that  $Y$  is the continuous-time Markov chain defined in (36)-(39) and  $\pi_G(\cdot)$  is the random mapping defined in (40) that is used to initialize  $Y$ . Besides,  $Y$  is a  $((S_k)_{k \geq 0}, (W_k)_{k \geq 0})$  jump process under Definition 2, with  $S_0 = 0$  and  $W_0 = \pi_G(m_i)$  (here  $x \in \Omega_i$  and  $X_0^\eta = x$ , so  $i$  is the index of the attraction field where the SGD iterate is initialized). The following result states that, given a sequence of learning rates  $(\eta_n)_{n \geq 1}$  that tend to 0, we are able to find a sequence of  $(\Delta_n)_{n \geq 1}$  to parametrize  $\hat{X}^{(n)} = \hat{X}^{*, \eta_n, \Delta_n}$ ,  $X^{(n)} = X^{*, \eta_n, \Delta_n}$  so that they have several useful properties, one of which is that the jump times and locations of  $\hat{X}^{(n)}$  converges in distribution to those of  $Y(\pi_G(m_i))$ .

**Lemma C.9.** *Assume the communication class  $G$  is absorbing. Given any  $m_i \in G$ ,  $x \in \Omega_i$ , finitely many real numbers  $(t_l)_{l=1}^{k'}$  such that  $0 < t_1 < t_2 < \dots < t_{k'}$ , and a sequence of strictly positive real numbers  $(\eta_n)_{n \geq 1}$  with  $\lim_{n \rightarrow \infty} \eta_n = 0$ , there exists a sequence of strictly positive real numbers  $(\Delta_n)_{n \geq 1}$  with  $\lim_n \Delta_n = 0$  such that*

- Under  $\mathbb{P}_x$  (so  $X_0^\eta = x$ ), as  $n$  tends to  $\infty$ ,

$$(S_0^{(n)}, W_0^{(n)}, S_1^{(n)}, W_1^{(n)}, S_2^{(n)}, W_2^{(n)}, \dots) \Rightarrow (S_0, W_0, S_1, W_1, S_2, W_2, \dots) \quad (\text{C.32})$$

- (Recall the definition of  $T_k, I_k$  in (C.11)-(C.14)) Given any  $\epsilon > 0$ , the following claim holds for all  $n$  sufficiently large:

$$\sup_{k \geq 0} \mathbb{P}_x \left( \exists j \in [T_k(\eta_n, \Delta_n), T_{k+1}(\eta_n, \Delta_n)] \text{ s.t. } X_j^\eta \notin \bigcup_{l: m_l \in G} \Omega_l \mid m_{I_k(\eta_n, \Delta_n)} \in G \right) < \epsilon; \quad (\text{C.33})$$

- Given any  $\epsilon > 0$ , the following claim holds for all  $n$  sufficiently large,

$$\sup_{k \geq 0} \mathbb{P}_x \left( m_{I_k(\eta_n, \Delta_n) + v} \notin G^{\text{large}} \mid m_{I_k(\eta_n, \Delta_n)} \in G \right) \leq \mathbb{P}(\text{Geom}(p^*) \geq u) + \epsilon \quad \forall u = 1, 2, \dots; \quad (\text{C.34})$$

- For any  $l \in [k']$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}_x \left( X_{t_l}^{*, \eta_n} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, \Delta_n), X_s^{*, \eta_n} \in \bigcup_{j: m_j \in G} \Omega_j \ \forall s \in [0, t_{k'}] \right) = 0 \quad (\text{C.35})$$

where  $p^* > 0$  is a constant that does not vary with our choices of  $\eta_n$  or  $\Delta_n$ .

*Proof.* Let

$$\begin{aligned} \nu_j &\triangleq q_j = \mu_j(E_j) \\ \nu_{j,k} &\triangleq \mu_j(E_{j,k}) \end{aligned}$$

so from the definition of  $q_{j,k}$  we have  $q_{j,k} = \mathbb{1}\{j \neq k\} \nu_{j,k} + \sum_{l: m_l \in G^{\text{small}}} \nu_{j,l} p_{l,k}$ .

In order to specify our choice of  $(\Delta_n)_{n \geq 1}$ , we consider a construction of sequences  $(\bar{\Delta}(j))_{j \geq 0}, (\bar{\eta}(j))_{j \geq 0}$  as follows. Fix some  $\theta \in (0, \alpha - 1)/2$ . Let  $\bar{\Delta}(0) = \bar{\eta}(0) = 1$ . One can see the existence of some  $(\bar{\Delta}(j))_{j \geq 1}, (\bar{\eta}(j))_{j \geq 1}$  such that

- $\bar{\Delta}(j) \in (0, \bar{\Delta}(j-1)/2]$ ,  $\bar{\eta}(j) \in (0, \bar{\eta}(j-1)/2]$  for all  $j \geq 1$ ;
- (Due to Lemma C.2) for any  $j \geq 1$ ,  $\eta \in (0, \bar{\eta}(j)]$ , (remember that  $x$  and  $i$  are the fixed constants prescribed in the description of the lemma)

$$\mathbb{P}_x \left( \sigma_0^*(\eta, \bar{\Delta}(j)) < \eta^\theta, \tilde{I}_0^G(\eta, \bar{\Delta}(j)) = i \right) > 1 - \frac{1}{2^j}.$$

For definitions of  $\sigma_k^G, \tau_k^G, I_k^G, \tilde{I}_k^G$ , see (24)-(29).

- (Due to Lemma C.6) for any  $j \geq 1$ ,  $\eta \in (0, \bar{\eta}(j)]$ ,

$$\left| \mathbb{P}_x \left( \tau_k^*(\eta, \bar{\Delta}(j)) - \sigma_k^*(\eta, \bar{\Delta}(j)) < \eta^\theta, I_k^G(\eta, \bar{\Delta}(j)) = i_2 \mid \tilde{I}_k^G(\eta, \bar{\Delta}(j)) = i_1 \right) - p_{i_1, i_2} \right| < 1/2^j$$

uniformly for all  $k \geq 0$  and all  $m_{i_1} \in G^{\text{small}}, m_{i_2} \in G^{\text{large}}$ . Also, by definition of  $\sigma^*$  and  $\tau^*$ , we must have

$$\mathbb{P}_x \left( \tau_k^*(\eta, \bar{\Delta}(j)) - \sigma_k^*(\eta, \bar{\Delta}(j)) = 0, I_k^G(\eta, \bar{\Delta}(j)) = i_1 \mid \tilde{I}_k^G(\eta, \bar{\Delta}(j)) = i_1 \right) = 1$$

for all  $k \geq 0$  and  $m_{i_1} \in G^{\text{large}}$ .

- (Due to Proposition B.24) for any  $j \geq 1$ ,  $\eta \in (0, \bar{\eta}(j)]$ ,

$$\begin{aligned} & -\frac{1}{2^j} + \exp \left( - \left( 1 + \frac{1}{2^j} \right) q_{i_1} u \right) \frac{\nu_{i_1, i_2} - \frac{1}{2^j}}{q_{i_1}} \\ & \leq \mathbb{P}_x \left( \sigma_{k+1}^*(\eta, \bar{\Delta}(j)) - \tau_k^*(\eta, \bar{\Delta}(j)) > u, \tilde{I}_{k+1}^G = i_2 \mid I_k^G(\eta, \bar{\Delta}(j)) = i_1 \right) \\ & \leq \frac{1}{2^j} + \exp \left( - \left( 1 - \frac{1}{2^j} \right) q_{i_1} u \right) \frac{\nu_{i_1, i_2} + \frac{1}{2^j}}{q_{i_1}} \end{aligned}$$

uniformly for all  $k \geq 1$ , all  $u > 1/2^j$ , and all  $m_{i_1} \in G^{\text{large}}, m_{i_2} \in G$ .

- (Due to  $G$  being absorbing and, again, Proposition B.24) for any  $j \geq 1$ , for any  $j \geq 1$ ,  $\eta \in (0, \bar{\eta}(j)]$ , (Recall the definition of  $T_k, I_k$  in (C.11)-(C.14))

$$\mathbb{P}_x \left( I_{k+1}(\eta, \bar{\Delta}(j)) = i_2 \mid I_k(\eta, \bar{\Delta}(j)) = i_1 \right) < \frac{1}{2^j} \quad (\text{C.36})$$

uniformly for all  $k \geq 0$ ,  $m_{i_1} \in G, m_{i_2} \notin G$ .



- (Due to Lemma C.4) There exists some  $p^* > 0$  such that for any  $j \geq 1$ , for any  $j \geq 1$ ,  $\eta \in (0, \bar{\eta}(j)]$ ,

$$\mathbb{P}_x \left( m_{v+I_k(\eta, \bar{\Delta}(j))} \notin G^{\text{large}} \mid \forall v \in [un_{\min}] \mid I_k(\eta, \bar{\Delta}(j)) = i_1 \right) \leq \mathbb{P}(\text{Geom}(p^*) \geq u) + 1/2^j \quad (\text{C.37})$$

uniformly for all  $k \geq 0, u \geq 1$  and  $m_{i_1} \in G$ .

- (Due to Lemma C.7) for any  $j \geq 1$ , for any  $j \geq 1$ ,  $\eta \in (0, \bar{\eta}(j)]$ ,

$$\mathbb{P}_x \left( X_{t_k}^{*, \eta_n} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, \bar{\Delta}(j)), X_s^{*, \eta_n} \in \bigcup_{j: m_j \in G} \Omega_j \mid \forall s \in [0, t_{k'}] \right) < 1/2^j \quad (\text{C.38})$$

uniformly for all  $k \in [k']$ .

Fix such  $(\bar{\Delta}(j))_{j \geq 0}, (\bar{\eta}(j))_{j \geq 0}$ . Define a function  $\mathbf{J}(\cdot) : \mathbb{N} \mapsto \mathbb{N}$  as

$$\mathbf{J}(n) = 0 \vee \max\{j \geq 0 : \bar{\eta}(j) \geq \eta_n\}$$

with the convention that  $\max \emptyset = -\infty$ . Lastly, let

$$\Delta_n = \bar{\Delta}(\mathbf{J}(n)) \quad \forall n \geq 1.$$

Note that, due to  $\lim_n \eta_n = 0$ , we have  $\lim_n \mathbf{J}(n) = \infty$ , hence  $\lim_n \Delta_n = 0$ . Besides, the definition of  $\mathbf{J}(\cdot)$  tells us that in case that  $\mathbf{J}(n) \geq 1$  (which will hold for all  $n$  sufficiently large), the claims above holds with  $\eta = \eta_n$  and  $j = \mathbf{J}(n)$ . In particular, by combining  $\lim_n \mathbf{J}(n) = \infty$  with (C.36)(C.37)(C.38) respectively, we have (C.33)(C.34)(C.35).

Now it remains to prove (C.32). To this end, it suffices to show that, for any positive integer  $K$ , we have  $(S_0^{(n)}, W_0^{(n)}, \dots, S_K^{(n)}, W_K^{(n)})$  converges in distribution  $(S_0, W_0, \dots, S_K, W_K)$  as  $n$  tends to infinity. In particular, note that  $S_0 = 0, W_0 = \pi_G(m_i)$ , so  $W_0 = m_j$  with probability  $p_{i,j}$  if  $m_i \in G^{\text{small}}$ , and  $W_0 \equiv m_i$  if  $m_i \in G^{\text{large}}$ .

For clarity of the exposition, we restate some important claims above under the new notational system with  $\tilde{S}_k^{(n)}, \tilde{W}_k^{(n)}, S_k^{(n)}, W_k^{(n)}$  we introduced right above this lemma. Given any  $\epsilon > 0$ , the following claims hold for all  $n$  sufficiently large:

- First of all,

$$\mathbb{P}_x \left( \tilde{S}_0^{(n)} < \eta_n^\theta, \tilde{W}_0^{(n)} = m_i \right) > 1 - \epsilon. \quad (\text{C.39})$$

- For all  $k \geq 0$  and all  $m_{i_1} \in G^{\text{small}}, m_{i_2} \in G^{\text{large}}$ ,

$$\left| \mathbb{P}_x \left( S_k^{(n)} - \tilde{S}_k^{(n)} < \eta_n^\theta, W_k^{(n)} = m_{i_2} \mid \tilde{W}_k^{(n)} = m_{i_1} \right) - p_{i_1, i_2} \right| < \epsilon. \quad (\text{C.40})$$

- For all  $k \geq 0$  and all  $m_{i_1} \in G^{\text{large}}$ ,

$$\mathbb{P}_x \left( S_k^{(n)} - \tilde{S}_k^{(n)} = 0, W_k^{(n)} = m_{i_1} \mid \tilde{W}_k^{(n)} = m_{i_1} \right) = 1. \quad (\text{C.41})$$

- For all  $k \geq 0$ , all  $m_{i_1} \in G^{\text{large}}, m_{i_2} \in G$  and all  $u > \epsilon$ ,

$$\begin{aligned} & -\epsilon + \exp \left( -(1 + \epsilon)q_{i_1}u \right) \frac{\nu_{i_1, i_2} - \epsilon}{q_{i_1}} \\ & \leq \mathbb{P}_x \left( \tilde{S}_{k+1}^{(n)} - S_k^{(n)} > u, \tilde{W}_{k+1}^{(n)} = m_{i_2} \mid W_k^{(n)} = m_{i_1} \right) \\ & \leq \mathbb{P}_x \left( \tilde{S}_{k+1}^{(n)} - S_k^{(n)} > u - \eta_n^\theta, \tilde{W}_{k+1}^{(n)} = m_{i_2} \mid W_k^{(n)} = m_{i_1} \right) \\ & \leq \epsilon + \exp \left( -(1 - \epsilon)q_{i_1}u \right) \frac{\nu_{i_1, i_2} + \epsilon}{q_{i_1}} \end{aligned} \quad (\text{C.42})$$

- Here is one implication of (C.40). Since  $|G| \leq n_{\min}$ , we have

$$\mathbb{P}_x \left( S_k^{(n)} - \tilde{S}_k^{(n)} \geq \eta_n^\theta \mid \widetilde{W}_k^{(n)} = m_{i_1} \right) < n_{\min} \cdot \epsilon \quad (\text{C.43})$$

for all  $k \geq 0$  and  $m_{i_1} \in G^{\text{small}}$ .

- Note that for any  $m_{i_1}, m_{i_2} \in G^{\text{large}}$  and any  $k \geq 0$

$$\begin{aligned} & \mathbb{P}_x \left( S_{k+1}^{(n)} - S_k^{(n)} > u, W_{k+1}^{(n)} = m_{i_2} \mid W_k^{(n)} = m_{i_1} \right) \\ &= \mathbb{1}\{i_2 \neq i_1\} \mathbb{P}_x \left( \tilde{S}_{k+1}^{(n)} - S_k^{(n)} > u, \widetilde{W}_{k+1}^{(n)} = m_{i_2} \mid W_k^{(n)} = m_{i_1} \right) \\ &+ \sum_{i_3: m_{i_3} \in G^{\text{small}}} \int_{s>0} \mathbb{P}_x \left( S_{k+1}^{(n)} - \tilde{S}_{k+1}^{(n)} \geq (u-s) \vee 0, \widetilde{W}_{k+1}^{(n)} = m_{i_2} \mid \widetilde{W}_k^{(n)} = m_{i_3} \right) \\ &\quad \cdot \mathbb{P}_x \left( \tilde{S}_{k+1}^{(n)} - S_k^{(n)} = ds, \widetilde{W}_{k+1}^{(n)} = m_{i_3} \mid W_k^{(n)} = m_{i_1} \right). \end{aligned}$$

Fix some  $i_3$  with  $m_{i_3} \in G^{\text{small}}$ . On the one hand, due to (C.43),

$$\begin{aligned} & \int_{s \in (0, u - \eta_n^\theta]} \mathbb{P}_x \left( S_{k+1}^{(n)} - \tilde{S}_{k+1}^{(n)} \geq u - s, \widetilde{W}_{k+1}^{(n)} = m_{i_2} \mid \widetilde{W}_k^{(n)} = m_{i_3} \right) \\ &\quad \cdot \mathbb{P}_x \left( \tilde{S}_{k+1}^{(n)} - S_k^{(n)} = ds, \widetilde{W}_{k+1}^{(n)} = m_{i_3} \mid W_k^{(n)} = m_{i_1} \right) \\ &\leq n_{\min} \epsilon. \end{aligned}$$

On the other hand, by considering the integral on  $(u - \eta_n^\theta, \infty)$ , we get

$$\begin{aligned} & \int_{s \in (u - \eta_n^\theta, \infty)} \mathbb{P}_x \left( S_{k+1}^{(n)} - \tilde{S}_{k+1}^{(n)} \geq (u-s) \vee 0, \widetilde{W}_{k+1}^{(n)} = m_{i_2} \mid \widetilde{W}_k^{(n)} = m_{i_3} \right) \\ &\quad \cdot \mathbb{P}_x \left( \tilde{S}_{k+1}^{(n)} - S_k^{(n)} = ds, \widetilde{W}_{k+1}^{(n)} = m_{i_3} \mid W_k^{(n)} = m_{i_1} \right) \\ &\geq \int_{s \in (u, \infty)} \mathbb{P}_x \left( \widetilde{W}_{k+1}^{(n)} = m_{i_2} \mid \widetilde{W}_k^{(n)} = m_{i_3} \right) \cdot \mathbb{P}_x \left( \tilde{S}_{k+1}^{(n)} - S_k^{(n)} = ds, \widetilde{W}_{k+1}^{(n)} = m_{i_3} \mid W_k^{(n)} = m_{i_1} \right) \\ &\geq (p_{i_3, i_2} - \epsilon) \left( -\epsilon + \exp \left( -(1 + \epsilon) q_{i_1} u \right) \frac{\nu_{i_1, i_3} - \epsilon}{q_{i_1}} \right) \end{aligned}$$

due to (C.40) and (C.42). Meanwhile,

$$\begin{aligned} & \int_{s \in (u - \eta_n^\theta, \infty)} \mathbb{P}_x \left( S_{k+1}^{(n)} - \tilde{S}_{k+1}^{(n)} \geq (u-s) \vee 0, \widetilde{W}_{k+1}^{(n)} = m_{i_2} \mid \widetilde{W}_k^{(n)} = m_{i_3} \right) \\ &\quad \cdot \mathbb{P}_x \left( \tilde{S}_{k+1}^{(n)} - S_k^{(n)} = ds, \widetilde{W}_{k+1}^{(n)} = m_{i_3} \mid W_k^{(n)} = m_{i_1} \right) \\ &\leq (n_{\min} \epsilon + p_{i_3, i_2} + \epsilon) \left( \epsilon + \exp \left( -(1 - \epsilon) q_{i_1} u \right) \frac{\nu_{i_1, i_3} + \epsilon}{q_{i_1}} \right) \end{aligned}$$

due to (C.40), (C.42) and (C.43).

- Therefore, for any  $m_{i_1}, m_{i_2} \in G^{\text{large}}$  and any  $k \geq 0$ ,

$$\begin{aligned} & \mathbb{P}_x \left( S_{k+1}^{(n)} - S_k^{(n)} > u, W_{k+1}^{(n)} = m_{i_2} \mid W_k^{(n)} = m_{i_1} \right) \\ &\leq g(\epsilon) + \exp \left( -(1 - \epsilon) q_{i_1} u \right) \frac{\mathbb{1}\{i_2 \neq i_1\} \nu_{i_1, i_2} + \sum_{i_3: m_{i_3} \in G^{\text{small}}} \nu_{i_1, i_3} p_{i_3, i_2}}{q_{i_1}} \end{aligned}$$

$$\leq g(\epsilon) + \exp\left(-(1-\epsilon)q_{i_1}u\right) \frac{q_{i_1,i_2}}{q_{i_1}} \quad (\text{C.44})$$

and

$$\begin{aligned} & \mathbb{P}_x\left(S_{k+1}^{(n)} - S_k^{(n)} > u, W_{k+1}^{(n)} = m_{i_2} \mid W_k^{(n)} = m_{i_1}\right) \\ & \geq -g(\epsilon) + \exp\left(-(1+\epsilon)q_{i_1}u\right) \frac{\mathbb{1}\{i_2 \neq i_1\}\nu_{i_1,i_2} + \sum_{i_3: m_{i_3} \in G^{\text{small}}} \nu_{i_1,i_3}p_{i_3,i_2}}{q_{i_1}} \\ & \geq -g(\epsilon) + \exp\left(-(1+\epsilon)q_{i_1}u\right) \frac{q_{i_1,i_2}}{q_{i_1}} \end{aligned} \quad (\text{C.45})$$

where  $q^* = \max_i q_i$  and

$$g(\epsilon) \triangleq \epsilon + \frac{\epsilon}{q^*} + n_{\min}(1+\epsilon)\epsilon + \epsilon \frac{1+\epsilon}{q^*} n_{\min} + n_{\min}\left(\epsilon + \frac{\epsilon}{q^*}\right) + n_{\min}(n_{\min}\epsilon + \epsilon + 1)\left(1 + \frac{1}{q^*}\right)\epsilon.$$

Note that  $\lim_{\epsilon \downarrow 0} g(\epsilon) = 0$ .

Now we apply the bounds in (C.39)(C.44)(C.45) to establish the weak convergence claim regarding  $(S_0^{(n)}, W_0^{(n)}, \dots, S_K^{(n)}, W_K^{(n)})$ . Fix some positive integer  $K$ , some strictly positive real numbers  $(s_k)_{k=0}^K$ , a sequence  $(w_k)_{k=0}^K \in (G^{\text{large}})^{K+1}$  with  $w_k = m_{i_k}$  for each  $k$ , and some  $\epsilon > 0$  such that  $\epsilon < \min_{k=0,1,\dots,K} \mathbb{1}\{s_k\}$ . On the one hand, the definition of the CTMC  $Y$  implies that

$$\begin{aligned} & \mathbb{P}\left(S_0 < t_0, W_0 = w_0; S_k > s_k \text{ and } W_k = w_k \ \forall k \in [K]\right) \\ & = \mathbb{P}(\pi_G(m_i) = w_0) \prod_{k=1}^K \left(S_k > s_k \text{ and } W_k = w_k \mid W_{k-1} = w_{k-1}\right) \\ & = \left(\mathbb{1}\{m_i \in G^{\text{large}}, i_0 = i\} + \mathbb{1}\{m_i \in G^{\text{small}}\}p_{i_0,i_1}\right) \cdot \prod_{k=1}^K \exp(-q_{i_{k-1}}s_k) \frac{q_{i_{k-1},i_k}}{q_{i_{k-1}}}. \end{aligned}$$

On the other hand, using (C.39)(C.44)(C.45), we know that for all  $n$  sufficiently large,

$$\begin{aligned} & \mathbb{P}_x\left(S_0^{(n)} < s_0, W_0^{(n)} = w_0; S_k^{(n)} > s_k \text{ and } W_k^{(n)} = w_k \ \forall k \in [K]\right) \\ & \geq (1-\epsilon) \left(\mathbb{1}\{m_i \in G^{\text{large}}, i_0 = i\} + \mathbb{1}\{m_i \in G^{\text{small}}\}(p_{i_0,i_1} - \epsilon)\right) \cdot \prod_{k=1}^K \left(-g(\epsilon) + \exp(-(1+\epsilon)q_{i_{k-1}}s_k) \frac{q_{i_{k-1},i_k}}{q_{i_{k-1}}}\right) \end{aligned}$$

and

$$\begin{aligned} & \mathbb{P}_x\left(S_0^{(n)} < s_0, W_0^{(n)} = w_0; S_k^{(n)} > s_k \text{ and } W_k^{(n)} = w_k \ \forall k \in [K]\right) \\ & \leq \left(\mathbb{1}\{m_i \in G^{\text{large}}, i_0 = i\} + \mathbb{1}\{m_i \in G^{\text{small}}\}(p_{i_0,i_1} + \epsilon)\right) \cdot \prod_{k=1}^K \left(g(\epsilon) + \exp(-(1-\epsilon)q_{i_{k-1}}s_k) \frac{q_{i_{k-1},i_k}}{q_{i_{k-1}}}\right). \end{aligned}$$

Since  $\epsilon > 0$  can be arbitrarily small, we now obtain

$$\begin{aligned} & \lim_{n \rightarrow \infty} \mathbb{P}_x\left(S_0^{(n)} < t_0, W_0^{(n)} = w_0; S_k^{(n)} > s_k \text{ and } W_k^{(n)} = m_k \ \forall k \in [K]\right) \\ & = \mathbb{P}\left(S_0 < s_0, W_0 = w_0; S_k > s_k \text{ and } W_k = w_k \ \forall k \in [K]\right), \end{aligned}$$

and the arbitrariness of the integer  $K$ , the strictly positive real numbers  $(s_k)_{k=0}^K$ , and the sequence  $(w_k)_{k=0}^K \in (G^{\text{large}})^{K+1}$  allows us to conclude the proof.  $\square$

To extend the result above to the case where the communication class  $G$  is transient, we revisit the definition of the  $Y^\dagger$  in (44). Let  $\bar{G} = G^{\text{large}} \cup \{\dagger\}$  and let  $m_0 = \dagger$  (remember that all the local minimizers of  $f$  on  $[-L, L]$  are  $m_1, \dots, m_{\min}$ ). Meanwhile, using  $q_i$  and  $q_{i,j}$  in (38)(39), we can define

$$q_{i,j}^\dagger = \begin{cases} q_{i,j} & \text{if } i \geq 1, j \geq 1 \\ \mathbb{1}\{j = 0\} & \text{if } i = 0 \\ \sum_{j \in [n_{\min}], m_j \notin G} q_{i,j} & \text{if } i \geq 1, j = 0. \end{cases}$$

and  $q_0^\dagger = 1$ ,  $q_i^\dagger = q_i \forall i \geq 1$ . Next, fix some  $i$  with  $m_i \in G$  and  $x \in \Omega_i$ . Define a sequence of random variables  $(S_k^\dagger)_{k \geq 0}, (W_k^\dagger)_{k \geq 0}$  such that  $S_k^\dagger = 0$  and  $W_0^\dagger = \dagger \mathbb{1}\{W_0 \notin G^{\text{large}}\} + W_0 \mathbb{1}\{W_0 \in G^{\text{large}}\}$  (see the definition of random mapping  $\pi_G$  in (40)) and (for all  $k \geq 0$  and  $i, j$  with  $m_j, m_l \in \bar{G}$ )

$$\mathbb{P}\left(W_{k+1}^\dagger = m_l, S_{k+1}^\dagger > t \mid W_k^\dagger = m_j, (W_l^\dagger)_{l=0}^{k-1}, (S_l^\dagger)_{l=0}^k\right) \quad (\text{C.46})$$

$$= \mathbb{P}\left(W_{k+1}^\dagger = m_l, S_{k+1}^\dagger > t \mid W_k^\dagger = m_j\right) = \exp(-q_j^\dagger t) \frac{q_{j,l}^\dagger}{q_j^\dagger} \forall t > 0 \quad (\text{C.47})$$

Then it is easy to see that  $Y^\dagger(\pi_G(m_i))$  defined in (44) is a  $((S_k^\dagger)_{k \geq 0}, (W_k^\dagger)_{k \geq 0})$  jump process. In particular, from at any state that is not  $\dagger$  (namely, any  $m_j$  with  $m_j \in G^{\text{large}}$ ), the probability that  $Y^\dagger$  moves to  $\dagger$  in the next transition is equal to the chance that, starting from the same state,  $Y$  moves to a state that is not in  $G$ . Once entering  $m_0 = \dagger$ , the process  $Y^\dagger$  will only make dummy jumps (with interarrival times being iid  $\text{Exp}(1)$ ): indeed, we have  $q_0^\dagger = q_{0,0}^\dagger = 1$  and  $q_{0,j}^\dagger = 0$  for any  $j \geq 1$ , implying that, given  $W_k^\dagger = m_0 = \dagger$ , we must have  $W_{k+1}^\dagger = m_0 = \dagger$ . These dummy jumps ensure that  $Y^\dagger$  is stuck at the cemetery state  $\dagger$  after visiting it.

Similarly, we can characterize the jump times and locations of the jump process  $\hat{X}^{\dagger,*,\eta,\Delta}$  (for the definition, see (43)). When there is no ambiguity about the sequences  $(\eta_n)_{n \geq 1}, (\Delta_n)_{n \geq 1}$ , let  $\hat{X}^{\dagger,(n)} = \hat{X}^{\dagger,*,\eta_n,\Delta_n}$  and  $X^{\dagger,(n)} = X^{*,\eta_n,\Delta_n}$ . Also, recall that  $\tau_G$  defined in (13) is the step  $n$  when  $X_n^\eta$  exits the communication class  $G$ . Now let  $(E_k)_{k \geq 0}$  be a sequence of iid  $\text{Exp}(1)$  random variables that is also independent of the noises  $(Z_k)_{k \geq 1}$  (so they are independent from the SGD iterates  $X_n^\eta$ ). For all  $n \geq 1, k \geq 0$ , define (see (24)-(29) and (33) for definitions of the quantities involved)

$$\begin{aligned} \tilde{S}_k^{\dagger,(n)} &= \begin{cases} \sigma_k^*(\eta_n, \Delta_n) \wedge \mathbf{T}^*(\tau_G(\eta_n), \eta_n) - \tau_{k-1}^*(\eta_n, \Delta_n) & \text{if } \tau_{k-1}^*(\eta_n, \Delta_n) < \mathbf{T}^*(\tau_G(\eta_n), \eta_n) \\ 0 & \text{otherwise} \end{cases} \\ S_k^{\dagger,(n)} &= \begin{cases} \tau_k^*(\eta_n, \Delta_n) \wedge \mathbf{T}^*(\tau_G(\eta_n), \eta_n) - \tau_{k-1}^*(\eta_n, \Delta_n) & \text{if } \tau_{k-1}^*(\eta_n, \Delta_n) < \mathbf{T}^*(\tau_G(\eta_n), \eta_n) \\ E_k & \text{otherwise} \end{cases} \\ \widetilde{W}_k^{\dagger,(n)} &= \begin{cases} m_{\tilde{I}_k^G(\eta_n, \Delta_n)} & \text{if } \sigma_k^*(\eta_n, \Delta_n) < \mathbf{T}^*(\tau_G(\eta_n), \eta_n) \\ \dagger & \text{otherwise} \end{cases} \\ W_k^{\dagger,(n)} &= \begin{cases} m_{I_k^G(\eta_n, \Delta_n)} & \text{if } \tau_k^*(\eta_n, \Delta_n) < \mathbf{T}^*(\tau_G(\eta_n), \eta_n) \\ \dagger & \text{otherwise} \end{cases} \end{aligned}$$

with the convention that  $\tau_{-1}^* = 0$ . Note that  $\mathbf{T}^*(\tau_G(\eta_n), \eta_n)$  is the scaled timestamp for  $X^{(n)} = X^{*,\eta}$  corresponding to  $\tau_G(\eta_n)$ , hence  $\mathbf{T}^*(\tau_G(\eta_n), \eta_n) = \min\{t \geq 0 : X^{(n)} \notin \bigcup_{j: m_j \in G} \Omega_j\}$ . One can see that  $\hat{X}^{\dagger,(n)}$  is a  $((S_k^{\dagger,(n)})_{k \geq 0}, (W_k^{\dagger,(n)})_{k \geq 0})$  jump process. The next lemma is similar to Lemma C.9 and discusses the convergence of the jump times and locations of  $\hat{X}^{\dagger,(n)}$  on a communication class  $G$  in the transient case.

**Lemma C.10.** *Assume that the communication class  $G$  is transient. Given any  $m_i \in G$ ,  $x \in \Omega_i$ , finitely many real numbers  $(t_l)_{l=1}^{k'}$  such that  $0 < t_1 < t_2 < \dots < t_{k'}$ , and a sequence of strictly positive*

real numbers  $(\eta_n)_{n \geq 1}$  with  $\lim_{n \rightarrow 0} \eta_n = 0$ , there exists a sequence of strictly positive real numbers  $(\Delta_n)_{n \geq 1}$  with  $\lim_n \Delta_n = 0$  such that

- Under  $\mathbb{P}_x$  (so  $X_0^\eta = x$ ), as  $n$  tends to  $\infty$ ,

$$(S_0^{\dagger,(n)}, W_0^{\dagger,(n)}, S_1^{\dagger,(n)}, W_1^{\dagger,(n)}, S_2^{\dagger,(n)}, W_2^{\dagger,(n)}, \dots) \Rightarrow (S_0^\dagger, W_0^\dagger, S_1^\dagger, W_1^\dagger, S_2^\dagger, W_2^\dagger, \dots) \quad (\text{C.48})$$

- For any  $l \in [k']$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}_x \left( X_{t_l}^{\dagger,(n)} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, \Delta_n), X_s^{\dagger,(n)} \in \bigcup_{j: m_j \in G} \Omega_j \ \forall s \in [0, t_l] \right) = 0 \quad (\text{C.49})$$

*Proof.* Let

$$\begin{aligned} \nu_{j,k} &\triangleq \mu_j(E_{j,k}) \ \forall j, k \geq 1, \ j \neq k \\ p_{j,\dagger} &\triangleq \sum_{\tilde{j}: m_{\tilde{j}} \notin G} p_{j,\tilde{j}} \ \forall m_j \in G^{\text{small}} \\ q_{j,\dagger} &\triangleq \sum_{k: m_k \notin G} \nu_{j,k} + \sum_{k: m_k \in G^{\text{small}}} \nu_{j,k} p_{k,\dagger} \ \forall m_j \in G^{\text{large}}. \end{aligned}$$

In order to specify our choice of  $(\Delta_n)_{n \geq 1}$ , we consider a construction of sequences  $(\bar{\Delta}(j))_{j \geq 0}, (\bar{\eta}(j))_{j \geq 0}$  as follows. Fix some  $\theta \in (0, \alpha - 1)/2$ . Let  $\bar{\Delta}(0) = \bar{\eta}(0) = 1$ . One can see the existence of some  $(\bar{\Delta}(j))_{j \geq 1}, (\bar{\eta}(j))_{j \geq 1}$  such that

- $\bar{\Delta}(j) \in (0, \bar{\Delta}(j-1)/2]$ ,  $\bar{\eta}(j) \in (0, \bar{\eta}(j-1)/2]$  for all  $j \geq 1$ ;
- (Due to Lemma C.2) for any  $j \geq 1$ ,  $\eta \in (0, \bar{\eta}(j)]$ , (remember that  $x$  and  $i$  are the fixed constants prescribed in the description of the lemma)

$$\mathbb{P}_x \left( \sigma_0^*(\eta, \bar{\Delta}(j)) < \eta^\theta, \tilde{I}_0^G(\eta, \bar{\Delta}(j)) = i \right) > 1 - \frac{1}{2^j}.$$

For definitions of  $\sigma_k^G, \tau_k^G, I_k^G, \tilde{I}_k^G$ , see (24)-(29).

- (Due to Lemma C.6) for any  $j \geq 1$ ,  $\eta \in (0, \bar{\eta}(j)]$ ,

$$\left| \mathbb{P}_x \left( \tau_k^*(\eta, \bar{\Delta}(j)) - \sigma_k^*(\eta, \bar{\Delta}(j)) < \eta^\theta, I_k^G(\eta, \bar{\Delta}(j)) = i_2 \mid \tilde{I}_k^G(\eta, \bar{\Delta}(j)) = i_1 \right) - p_{i_1, i_2} \right| < 1/2^j$$

uniformly for all  $k \geq 0$  and all  $m_{i_1} \in G^{\text{small}}, m_{i_2} \notin G^{\text{small}}$ . Also, by definition of  $\sigma^*$  and  $\tau^*$ , we must have

$$\mathbb{P}_x \left( \tau_k^*(\eta, \bar{\Delta}(j)) - \sigma_k^*(\eta, \bar{\Delta}(j)) = 0, I_k^G(\eta, \bar{\Delta}(j)) = i_1 \mid \tilde{I}_k^G(\eta, \bar{\Delta}(j)) = i_1 \right) = 1$$

for all  $k \geq 0$  and  $m_{i_1} \in G^{\text{large}}$ .

- (Due to Proposition B.24) for any  $j \geq 1$ ,  $\eta \in (0, \bar{\eta}(j)]$ ,

$$\begin{aligned} & - \frac{1}{2^j} + \exp \left( - \left( 1 + \frac{1}{2^j} \right) q_{i_1} u \right) \frac{\nu_{i_1, i_2} - \frac{1}{2^j}}{q_{i_1}} \\ & \leq \mathbb{P}_x \left( \sigma_{k+1}^*(\eta, \bar{\Delta}(j)) - \tau_k^*(\eta, \bar{\Delta}(j)) > u, \tilde{I}_{k+1}^G = i_2 \mid I_k^G(\eta, \bar{\Delta}(j)) = i_1 \right) \\ & \leq \frac{1}{2^j} + \exp \left( - \left( 1 - \frac{1}{2^j} \right) q_{i_1} u \right) \frac{\nu_{i_1, i_2} + \frac{1}{2^j}}{q_{i_1}} \end{aligned}$$

uniformly for all  $k \geq 1$ , all  $u > 1/2^j$ , and all  $m_{i_1} \in G^{\text{large}}, m_{i_2} \in G$ .

- (Due to Lemma C.7) for any  $j \geq 1$ , for any  $j \geq 1$ ,  $\eta \in (0, \bar{\eta}(j)]$ ,

$$\mathbb{P}_x \left( X_{t_k}^{(n)} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, \bar{\Delta}(j)), X_s^{(n)} \in \bigcup_{j: m_j \in G} \Omega_j \ \forall s \in [0, t_k] \right) < 1/2^j \quad (\text{C.50})$$

uniformly for all  $k \in [k']$ .

Fix such  $(\bar{\Delta}(j))_{j \geq 0}, (\bar{\eta}(j))_{j \geq 0}$ . Define a function  $\mathbf{J}(\cdot) : \mathbb{N} \mapsto \mathbb{N}$  as

$$\mathbf{J}(n) = 0 \vee \max\{j \geq 0 : \bar{\eta}(j) \geq \eta_n\}$$

with the convention that  $\max \emptyset = -\infty$ . Lastly, let

$$\Delta_n = \bar{\Delta}(\mathbf{J}(n)) \ \forall n \geq 1.$$

Note that, due to  $\lim_n \eta_n = 0$ , we have  $\lim_n \mathbf{J}(n) = \infty$ , hence  $\lim_n \Delta_n = 0$ . Besides, since  $X_t^{\dagger, (n)} = X_t^{(n)}$  given  $X_s^{(n)} \in \bigcup_{j: m_j \in G} \Omega_j$  for all  $s \in [0, t]$ , by combining  $\lim_n \mathbf{J}(n) = \infty$  with (C.50) we obtain (C.49).

Now it remains to prove (C.48). To this end, it suffices to show that, for any positive integer  $K$ , we have  $(S_0^{\dagger, (n)}, W_0^{\dagger, (n)}, \dots, S_K^{\dagger, (n)}, W_K^{\dagger, (n)})$  converges in distribution  $(S_0^{\dagger}, W_0^{\dagger}, \dots, S_K^{\dagger}, W_K^{\dagger})$  as  $n$  tends to infinity. In particular, due to introduction of the dummy jumps, we know that for any  $k$  with  $\tau_k^*(\eta_n, \Delta_n) \geq \tau_G(\eta_n)$  (in other words,  $\hat{X}^{\dagger, (n)}$  has reached state  $\dagger$  within the first  $k$  jumps) we have  $S_{k+1}^{\dagger, (n)} \sim \text{Exp}(1)$  and  $W_{k+1}^{\dagger, (n)} \equiv \dagger$ . Similarly, for any  $k$  with  $S_0^{\dagger} + \dots + S_k^{\dagger} \leq \tau_G^Y$ , we have  $S_{k+1}^{\dagger} \sim \text{Exp}(1)$  and  $W_{k+1}^{\dagger} \equiv \dagger$ . Therefore, it suffices to show that, for any positive integer  $K$ , any series of strictly positive real numbers  $(s_k)_{k=0}^K$ , any sequence  $(w_k)_{k=0}^K \in (\bar{G})^{K+1}$  such that  $w_j \neq \dagger$  for any  $j < K$ , indices  $i_k$  such that  $w_k = m_{i_k}$  for each  $k$ , we have

$$\begin{aligned} & \lim_{n \rightarrow \infty} \mathbb{P}_x \left( S_0^{\dagger, (n)} < t_0, W_0^{\dagger, (n)} = w_0; S_k^{\dagger, (n)} > s_k \text{ and } W_k^{\dagger, (n)} = w_k \ \forall k \in [K] \right) \\ &= \mathbb{P} \left( S_0^{\dagger} < s_0, W_0^{\dagger} = w_0; S_k^{\dagger} > s_k \text{ and } W_k^{\dagger} = w_k \ \forall k \in [K] \right) \end{aligned} \quad (\text{C.51})$$

Fix some  $(s_k)_{k=0}^K, (w_k)_{k=0}^K \in (\bar{G})^{K+1}$ , and indices  $(i_k)_{k=1}^K$  satisfying the conditions above. Besides, arbitrarily choose some  $\epsilon > 0$  so that  $\epsilon < \min_{k=0, \dots, K} s_k$ . To proceed, we start by translating the inequalities established above under the new system of notations.

- First of all, for all  $n$  sufficiently large, (remember that  $x$  and  $i$  are prescribed constants in the description of this lemma)

$$\mathbb{P}_x \left( \tilde{S}_0^{\dagger, (n)} < \eta_n^{\theta}, \tilde{W}_0^{\dagger, (n)} = m_i \right) > 1 - \epsilon. \quad (\text{C.52})$$

- For all  $k \geq 0$  and all  $m_{i_1} \in G^{\text{small}}, m_{i_2} \in G^{\text{large}}$ , it holds for all  $n$  sufficiently large that

$$\left| \mathbb{P}_x \left( S_k^{\dagger, (n)} - \tilde{S}_k^{\dagger, (n)} < \eta_n^{\theta}, W_k^{\dagger, (n)} = m_{i_2} \mid \tilde{W}_k^{\dagger, (n)} = m_{i_1} \right) - p_{i_1, i_2} \right| < \epsilon. \quad (\text{C.53})$$

- On the other hand, for all  $k \geq 0$  and all  $m_{i_1} \in G^{\text{small}}$ , it holds for all  $n$  sufficiently large that

$$\begin{aligned} & \left| \mathbb{P}_x \left( S_k^{\dagger, (n)} - \tilde{S}_k^{\dagger, (n)} < \eta_n^{\theta}, W_k^{\dagger, (n)} = \dagger \mid \tilde{W}_k^{\dagger, (n)} = m_{i_1} \right) - \sum_{i_2: m_{i_2} \notin G} p_{i_1, i_2} \right| \\ &= \left| \mathbb{P}_x \left( S_k^{\dagger, (n)} - \tilde{S}_k^{\dagger, (n)} < \eta_n^{\theta}, W_k^{\dagger, (n)} = \dagger \mid \tilde{W}_k^{\dagger, (n)} = m_{i_1} \right) - p_{i_1, \dagger} \right| < \epsilon. \end{aligned} \quad (\text{C.54})$$

- For all  $k \geq 0$  and all  $m_{i_1} \in G^{\text{large}}$ , it holds for all  $n$  that

$$\mathbb{P}_x \left( S_k^{\dagger, (n)} - \tilde{S}_k^{\dagger, (n)} = 0, W_k^{\dagger, (n)} = m_{i_1} \mid \widetilde{W}_k^{\dagger, (n)} = m_{i_1} \right) = 1. \quad (\text{C.55})$$

- For all  $k \geq 0$ , all  $m_{i_1} \in G^{\text{large}}$ ,  $m_{i_2} \in G$  and all  $u > \epsilon$ , the following claim holds for all  $n$  sufficiently large:

$$\begin{aligned} & -\epsilon + \exp \left( -(1+\epsilon)q_{i_1}u \right) \frac{\nu_{i_1, i_2} - \epsilon}{q_{i_1}} \\ & \leq \mathbb{P}_x \left( \tilde{S}_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} > u, \widetilde{W}_{k+1}^{\dagger, (n)} = m_{i_2} \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\ & \leq \mathbb{P}_x \left( \tilde{S}_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} > u - \eta_n^\theta, \widetilde{W}_{k+1}^{\dagger, (n)} = m_{i_2} \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\ & \leq \epsilon + \exp \left( -(1-\epsilon)q_{i_1}u \right) \frac{\nu_{i_1, i_2} + \epsilon}{q_{i_1}} \end{aligned} \quad (\text{C.56})$$

- On the other hand, for all  $k \geq 0$ , all  $m_{i_1} \in G^{\text{large}}$ , the following claim holds for all  $n$  sufficiently large:

$$\begin{aligned} & -\epsilon + \exp \left( -(1+\epsilon)q_{i_1}u \right) \frac{-\epsilon + \sum_{i_2: m_{i_2} \notin G} \nu_{i_1, i_2}}{q_{i_1}} \\ & \leq \mathbb{P}_x \left( \tilde{S}_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} > u, \widetilde{W}_{k+1}^{\dagger, (n)} = \dagger \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\ & \leq \mathbb{P}_x \left( \tilde{S}_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} > u - \eta_n^\theta, \widetilde{W}_{k+1}^{\dagger, (n)} = \dagger \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\ & \leq \epsilon + \exp \left( -(1-\epsilon)q_{i_1}u \right) \frac{\epsilon + \sum_{i_2: m_{i_2} \notin G} \nu_{i_1, i_2}}{q_{i_1}} \end{aligned} \quad (\text{C.57})$$

- Here is one implication of (C.53)(C.54). Since  $|G| \leq n_{\min}$ , we have (when  $n$  is sufficiently large)

$$\mathbb{P}_x \left( S_k^{(n)} - \tilde{S}_k^{(n)} \geq \eta_n^\theta \mid \widetilde{W}_k^{(n)} = m_{i_1} \right) < n_{\min} \cdot \epsilon \quad (\text{C.58})$$

for all  $k \geq 0$  and  $m_{i_1} \in G^{\text{small}}$ .

- Note that for any  $m_{i_1}, m_{i_2} \in G^{\text{large}}$  and any  $k \geq 0$

$$\begin{aligned} & \mathbb{P}_x \left( S_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} > u, W_{k+1}^{\dagger, (n)} = m_{i_2} \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\ & = \mathbb{1}_{\{i_2 \neq i_1\}} \mathbb{P}_x \left( \tilde{S}_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} > u, \widetilde{W}_{k+1}^{\dagger, (n)} = m_{i_2} \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\ & + \sum_{i_3: m_{i_3} \in G^{\text{small}}} \int_{s>0} \mathbb{P}_x \left( S_{k+1}^{\dagger, (n)} - \tilde{S}_{k+1}^{\dagger, (n)} \geq (u-s) \vee 0, \widetilde{W}_{k+1}^{\dagger, (n)} = m_{i_2} \mid \widetilde{W}_k^{\dagger, (n)} = m_{i_3} \right) \\ & \cdot \mathbb{P}_x \left( \tilde{S}_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} = ds, \widetilde{W}_{k+1}^{\dagger, (n)} = m_{i_3} \mid W_k^{\dagger, (n)} = m_{i_1} \right). \end{aligned}$$

Fix some  $i_3$  with  $m_{i_3} \in G^{\text{small}}$ . Due to (C.58)

$$\begin{aligned} & \int_{s \in (0, u - \eta_n^\theta]} \mathbb{P}_x \left( S_{k+1}^{\dagger, (n)} - \tilde{S}_{k+1}^{\dagger, (n)} \geq u - s, \widetilde{W}_{k+1}^{\dagger, (n)} = m_{i_2} \mid \widetilde{W}_k^{\dagger, (n)} = m_{i_3} \right) \\ & \cdot \mathbb{P}_x \left( \tilde{S}_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} = ds, \widetilde{W}_{k+1}^{\dagger, (n)} = m_{i_3} \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\ & \leq n_{\min} \epsilon. \end{aligned}$$

Meanwhile, by considering the integral on  $(u - \eta_n^\theta, \infty)$ , we get

$$\begin{aligned}
& \int_{s \in (u - \eta_n^\theta, \infty)} \mathbb{P}_x \left( S_{k+1}^{\dagger, (n)} - \tilde{S}_{k+1}^{\dagger, (n)} \geq (u - s) \vee 0, \widetilde{W}_{k+1}^{\dagger, (n)} = m_{i_2} \mid \widetilde{W}_k^{\dagger, (n)} = m_{i_3} \right) \\
& \quad \cdot \mathbb{P}_x \left( \tilde{S}_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} = ds, \widetilde{W}_{k+1}^{\dagger, (n)} = m_{i_3} \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\
& \geq \int_{s \in (u, \infty)} \mathbb{P}_x \left( \widetilde{W}_{k+1}^{\dagger, (n)} = m_{i_2} \mid \widetilde{W}_k^{\dagger, (n)} = m_{i_3} \right) \cdot \mathbb{P}_x \left( \tilde{S}_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} = ds, \widetilde{W}_{k+1}^{\dagger, (n)} = m_{i_3} \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\
& \geq (p_{i_3, i_2} - \epsilon) \left( -\epsilon + \exp \left( -(1 + \epsilon) q_{i_1} u \right) \frac{\nu_{i_1, i_3} - \epsilon}{q_{i_1}} \right)
\end{aligned}$$

due to (C.53) and (C.56). As for the upper bound,

$$\begin{aligned}
& \int_{s \in (u - \eta_n^\theta, \infty)} \mathbb{P}_x \left( S_{k+1}^{\dagger, (n)} - \tilde{S}_{k+1}^{\dagger, (n)} \geq (u - s) \vee 0, \widetilde{W}_{k+1}^{\dagger, (n)} = m_{i_2} \mid \widetilde{W}_k^{\dagger, (n)} = m_{i_3} \right) \\
& \quad \cdot \mathbb{P}_x \left( \tilde{S}_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} = ds, \widetilde{W}_{k+1}^{\dagger, (n)} = m_{i_3} \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\
& \leq (n_{\min} \epsilon + p_{i_3, i_2} + \epsilon) \left( \epsilon + \exp \left( -(1 - \epsilon) q_{i_1} u \right) \frac{\nu_{i_1, i_3} + \epsilon}{q_{i_1}} \right)
\end{aligned}$$

due to (C.53), (C.56) and (C.58).

- Therefore, for any  $m_{i_1}, m_{i_2} \in G^{\text{large}}$  and any  $k \geq 0$ ,

$$\begin{aligned}
& \mathbb{P}_x \left( S_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} > u, W_{k+1}^{\dagger, (n)} = m_{i_2} \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\
& \leq g(\epsilon) + \exp \left( -(1 - \epsilon) q_{i_1} u \right) \frac{\mathbb{I}\{i_2 \neq i_1\} \nu_{i_1, i_2} + \sum_{i_3: m_{i_3} \in G^{\text{small}}} \nu_{i_1, i_3} p_{i_3, i_2}}{q_{i_1}} \\
& \leq g(\epsilon) + \exp \left( -(1 - \epsilon) q_{i_1} u \right) \frac{q_{i_1, i_2}}{q_{i_1}}
\end{aligned} \tag{C.59}$$

and

$$\begin{aligned}
& \mathbb{P}_x \left( S_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} > u, W_{k+1}^{\dagger, (n)} = m_{i_2} \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\
& \geq -g(\epsilon) + \exp \left( -(1 + \epsilon) q_{i_1} u \right) \frac{\mathbb{I}\{i_2 \neq i_1\} \nu_{i_1, i_2} + \sum_{i_3: m_{i_3} \in G^{\text{small}}} \nu_{i_1, i_3} p_{i_3, i_2}}{q_{i_1}} \\
& \geq -g(\epsilon) + \exp \left( -(1 + \epsilon) q_{i_1} u \right) \frac{q_{i_1, i_2}}{q_{i_1}}
\end{aligned} \tag{C.60}$$

where  $q^* = \max_i q_i$  and

$$g(\epsilon) \triangleq 2\epsilon + \frac{\epsilon}{q^*} + n_{\min}(1 + \epsilon)\epsilon + \epsilon \frac{1 + \epsilon}{q^*} n_{\min} + n_{\min} \left( \epsilon + \frac{\epsilon}{q^*} \right) + n_{\min}(n_{\min} \epsilon + \epsilon + 1) \left( 1 + \frac{1}{q^*} \right) \epsilon.$$

Note that  $\lim_{\epsilon \downarrow 0} g(\epsilon) = 0$ .

- On the other hand, for the case where the marker process  $\hat{X}^{\dagger, (n)}$  jumps to the cemetery state  $\dagger$  from some  $m_{i_1} \in G^{\text{large}}$ , note that

$$\begin{aligned}
& \mathbb{P}_x \left( S_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} > u, W_{k+1}^{\dagger, (n)} = \dagger \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\
& = \mathbb{P}_x \left( \tilde{S}_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} > u, \widetilde{W}_{k+1}^{\dagger, (n)} = \dagger \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\
& + \sum_{i_3: m_{i_3} \in G^{\text{small}}} \int_{s > 0} \mathbb{P}_x \left( S_{k+1}^{\dagger, (n)} - \tilde{S}_{k+1}^{\dagger, (n)} \geq (u - s) \vee 0, \widetilde{W}_{k+1}^{\dagger, (n)} = \dagger \mid \widetilde{W}_k^{\dagger, (n)} = m_{i_3} \right)
\end{aligned}$$



$$\cdot \mathbb{P}_x \left( \widetilde{S}_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} = ds, \widetilde{W}_{k+1}^{\dagger, (n)} = m_{i_3} \mid W_k^{\dagger, (n)} = m_{i_1} \right).$$

Arguing similarly as we did above by considering the integral on  $[0, u - \eta_n^\theta]$  and  $(u - \eta_n^\theta, \infty)$  separately, and using (C.54) and (C.57), we then get (for all  $n$  sufficiently large)

$$\begin{aligned} & \mathbb{P}_x \left( S_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} > u, W_{k+1}^{\dagger, (n)} = \dagger \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\ & \leq g(\epsilon) + \exp \left( - (1 - \epsilon) q_{i_1} u \right) \frac{\sum_{i_2: m_{i_2} \notin G} \nu_{i_1, i_2} + \sum_{i_2: m_{i_2} \in G^{\text{small}}} \nu_{i_1, i_2} p_{i_2, \dagger}}{q_{i_1}} \\ & = g(\epsilon) + \exp \left( - (1 - \epsilon) q_{i_1} u \right) \frac{q_{i_1, \dagger}}{q_{i_1}} \end{aligned} \quad (\text{C.61})$$

and

$$\begin{aligned} & \mathbb{P}_x \left( S_{k+1}^{\dagger, (n)} - S_k^{\dagger, (n)} > u, W_{k+1}^{\dagger, (n)} = \dagger \mid W_k^{\dagger, (n)} = m_{i_1} \right) \\ & \geq -g(\epsilon) + \exp \left( - (1 + \epsilon) q_{i_1} u \right) \frac{\sum_{i_2: m_{i_2} \notin G} \nu_{i_1, i_2} + \sum_{i_2: m_{i_2} \in G^{\text{small}}} \nu_{i_1, i_2} p_{i_2, \dagger}}{q_{i_1}} \\ & = -g(\epsilon) + \exp \left( - (1 + \epsilon) q_{i_1} u \right) \frac{q_{i_1, \dagger}}{q_{i_1}} \end{aligned} \quad (\text{C.62})$$

For simplicity of presentation, we also let  $q_{j,0} = q_{j,\dagger}$  and  $p_{j,0} = p_{j,\dagger}$ . First of all, remember that we have fixed some series of strictly positive real numbers  $(s_k)_{k=0}^K$ , some sequence  $(w_k)_{k=0}^K \in (\bar{G})^{K+1}$  such that  $w_j \neq \dagger$  for any  $j < K$ , and indices  $i_k$  such that  $w_k = m_{i_k}$  for each  $k$ . The definition of the continuous-time Markov chain  $Y^\dagger$  implies that

$$\begin{aligned} & \mathbb{P} \left( S_0 < t_0, W_0 = w_0; S_k > s_k \text{ and } W_k = w_k \ \forall k \in [K] \right) \\ & = \mathbb{P}(\pi_G(m_i) = w_0) \prod_{k=1}^K \left( S_k > s_k \text{ and } W_k = m_k \mid W_{k-1} = w_{k-1} \right) \\ & = \left( \mathbb{1}\{m_i \in G^{\text{large}}, i_0 = i\} + \mathbb{1}\{m_i \in G^{\text{small}}\} p_{i_0, i_1} \right) \cdot \prod_{k=1}^K \exp(-q_{i_{k-1}} s_k) \frac{q_{i_{k-1}, i_k}}{q_{i_{k-1}}}. \end{aligned}$$

On the other hand, using (C.59)-(C.62), we know that for all  $n$  sufficiently large,

$$\begin{aligned} & \mathbb{P}_x \left( S_0^{(n)} < s_0, W_0^{(n)} = w_0; S_k^{(n)} > s_k \text{ and } W_k^{(n)} = w_k \ \forall k \in [K] \right) \\ & \geq (1 - \epsilon) \left( \mathbb{1}\{m_i \in G^{\text{large}}, i_0 = i\} + \mathbb{1}\{m_i \in G^{\text{small}}\} (p_{i_0, i_1} - \epsilon) \right) \cdot \prod_{k=1}^K \left( -g(\epsilon) + \exp(-(1 + \epsilon) q_{i_{k-1}} s_k) \frac{q_{i_{k-1}, i_k}}{q_{i_{k-1}}} \right) \end{aligned}$$

and

$$\begin{aligned} & \mathbb{P}_x \left( S_0^{(n)} < s_0, W_0^{(n)} = w_0; S_k^{(n)} > s_k \text{ and } W_k^{(n)} = m_k \ \forall k \in [K] \right) \\ & \leq \left( \mathbb{1}\{m_i \in G^{\text{large}}, i_0 = i\} + \mathbb{1}\{m_i \in G^{\text{small}}\} (p_{i_0, i_1} + \epsilon) \right) \cdot \prod_{k=1}^K \left( g(\epsilon) + \exp(-(1 - \epsilon) q_{i_{k-1}} s_k) \frac{q_{i_{k-1}, i_k}}{q_{i_{k-1}}} \right). \end{aligned}$$

The arbitrariness of  $\epsilon > 0$  then allows us to establish (C.51) and conclude the proof.  $\square$

Now we are ready to prove Lemma 10 and Lemma 11.

*Proof of Lemma 10.* From Lemma C.9, one can see the existence of some  $(\Delta_n)_{n \geq 1}$  with  $\lim_n \Delta_n = 0$  such that (C.32)-(C.35) hold. For simplicity of notations, we let  $\hat{X}^{(n)} = \hat{X}^{*, \eta_n, \Delta_n}$ ,  $X^{(n)} = X^{*, \eta_n}$ , and let  $\bar{t} = t_{k'}$ .

Combine (C.32) with Lemma C.8, and we immediately get (41). In order to prove (42), it suffices to show that for any  $\epsilon > 0$ ,

$$\limsup_n \mathbb{P}_x \left( X_{t_k}^{(n)} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, \Delta_n) \right) \leq 4\epsilon \quad \forall k \in [k'].$$

Fix  $\epsilon > 0$ , and observe following bound by decomposing the events

$$\begin{aligned} & \mathbb{P}_x \left( X_{t_k}^{(n)} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, \Delta_n) \right) \\ & \leq \mathbb{P}_x \left( X_{t_k}^{(n)} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, \Delta_n), X_t^{(n)} \in \bigcup_{j: m_j \in G} \Omega_j \quad \forall t \in [0, \bar{t}] \right) \\ & + \mathbb{P}_x \left( \exists t \in [0, \bar{t}] \text{ such that } X_t^{(n)} \notin \bigcup_{j: m_j \in G} \Omega_j \right) \end{aligned}$$

Therefore, given (C.35), it suffices to prove

$$\limsup_n \mathbb{P}_x \left( \exists t \in [0, \bar{t}] \text{ such that } X_t^{(n)} \notin \bigcup_{j: m_j \in G} \Omega_j \right) \leq 3\epsilon. \quad (\text{C.63})$$

Let

$$\begin{aligned} T_0^{(n)} & \triangleq \min\{t \geq 0 : X_t^{(n)} \in \bigcup_{j: m_j \in G} B(m_j, 2\Delta_n)\} \\ I_0^{(n)} = j & \iff X_{T_0^{(n)}}^{(n)} \in B(m_j, 2\Delta_n) \\ T_k^{(n)} & \triangleq \min\{t > T_{k-1}^{(n)} : X_t^{(n)} \in \bigcup_{j: m_j \in G, j \neq I_{k-1}^{(n)}} B(m_j, 2\Delta_n)\} \\ I_k^{(n)} = j & \iff X_{T_k^{(n)}}^{(n)} \in B(m_j, 2\Delta_n). \end{aligned}$$

Building upon this definition, we define the following stopping times and marks that only records the hitting time to minimizer in *large* attraction fields in  $G$  (with convention  $\mathbf{k}^{(n), \text{large}}(-1) = -1$ ,  $T_{-1}^{(n), \text{large}} = 0$ ,  $T_{-1}^{(n)} = 0$ )

$$\begin{aligned} \mathbf{k}^{(n), \text{large}}(k) & \triangleq \min\{l > \mathbf{k}^{(n), \text{large}}(k-1) : m_{I_l^{(n)}} \in G^{\text{large}}\} \\ T_k^{(n), \text{large}} & \triangleq T_{\mathbf{k}^{(n), \text{large}}(k)}^{(n)}, \quad I_k^{(n), \text{large}} \triangleq I_{\mathbf{k}^{(n), \text{large}}(k)}^{(n)}. \end{aligned}$$

Now by defining

$$\begin{aligned} J^{(n)}(t) & \triangleq \#\{k \geq 0 : T_k^{(n)} \leq t\}, \\ J_{\text{large}}^{(n)}(t) & \triangleq \#\{k \geq 0 : T_k^{(n), \text{large}} \leq t\}, \\ J^{(n)}(s, t) & \triangleq \#\{k \geq 0 : T_k^{(n)} \in [s, t]\}, \end{aligned}$$

we use  $J^{(n)}(t)$  to count the numbers of visits to local minima on  $G$ , and  $J_{\text{large}}^{(n)}(t)$  for the number of visits to minimizers in the *large* attraction fields on  $G$ .  $J^{(n)}(s, t)$  counts the indices  $k$  such that at  $T_k^{(n)}$  a minimizer on  $G$  is visited and regarding the hitting time we have  $T_k^{(n)} \in [s, t]$ .

First of all, the weak convergence result in (C.32) implies the existence of some positive integer  $N(\epsilon)$  such that

$$\limsup_n \mathbb{P}_x(J_{\text{large}}^{(n)}(\bar{t}) > N(\epsilon)) < \epsilon.$$

Fix such  $N(\epsilon)$ . Next, from (C.34), we know the existence of some integer  $K(\epsilon)$  such that

$$\limsup_n \sup_{k \geq 0} \mathbb{P}_x \left( J^{(n)}(T_{k-1}^{(n), \text{large}}, T_k^{(n), \text{large}}) > K(\epsilon) \right) \leq \epsilon/N(\epsilon).$$

Fix such  $K(\epsilon)$  as well. From the results above, we know that for event

$$A_1(n) \triangleq \{J_{\text{large}}^{(n)}(\bar{t}) \leq N(\epsilon)\} \cap \left\{ J^{(n)}(T_{k-1}^{(n), \text{large}}, T_k^{(n), \text{large}}) \leq K(\epsilon) \ \forall k \in [N(\epsilon)] \right\},$$

we have  $\limsup_n \mathbb{P}_x \left( (A_1(n))^c \right) \leq 2\epsilon$ . On the other hand, on event  $A_1(n)$ , we must have

$$J^{(n)}(\bar{t}) \leq N(\epsilon)K(\epsilon).$$

Meanwhile, it follows immediately from (C.33) that

$$\limsup_n \sup_{k \geq 0} \mathbb{P} \left( \exists t \in [T_{k-1}^{(n)}, T_k^{(n)}] \text{ such that } X_t^{(n)} \notin \bigcup_{j: m_j \in G} \Omega_j \right) < \frac{\epsilon}{N(\epsilon)K(\epsilon)},$$

hence for event

$$A_2(n) \triangleq \left\{ X_t^{(n)} \in \bigcup_{j: m_j \in G} \Omega_j \ \forall t \in [0, T_{N(\epsilon)K(\epsilon)}^{(n)}] \right\},$$

we must have  $\limsup_n \mathbb{P}_x \left( (A_2(n))^c \right) \leq \epsilon$ . To conclude the proof, note that

$$\begin{aligned} A_1(n) \cap A_2(n) &\subseteq \{J^{(n)}(\bar{t}) \leq N(\epsilon)K(\epsilon)\} \cap \left\{ X_t^{(n)} \in \bigcup_{j: m_j \in G} \Omega_j \ \forall t \in [0, T_{N(\epsilon)K(\epsilon)}^{(n)}] \right\} \\ &= \{T_{N(\epsilon)K(\epsilon)}^{(n)} \geq \bar{t}\} \cap \left\{ X_t^{(n)} \in \bigcup_{j: m_j \in G} \Omega_j \ \forall t \in [0, T_{N(\epsilon)K(\epsilon)}^{(n)}] \right\} \\ &\subseteq \left\{ X_t^{(n)} \in \bigcup_{j: m_j \in G} \Omega_j \ \forall t \in [0, \bar{t}] \right\} \end{aligned}$$

so we have established (C.63).  $\square$

*Proof of Lemma 11.* From Lemma C.10, one can see the existence of some  $(\Delta_n)_{n \geq 1}$  with  $\lim_n \Delta_n = 0$  such that (C.48) and (C.49) hold. For simplicity of notations, we let  $\hat{X}^{\dagger, (n)} = \hat{X}^{\dagger, *, \eta_n, \Delta_n}$ ,  $X^{\dagger, (n)} = X^{\dagger, *, \eta_n}$ , and let  $\bar{t} = t_k$ .

Combine (C.48) with Lemma C.8, and we immediately get (45). In order to prove (46), note that

$$\begin{aligned} &\left\{ X_{t_k}^{\dagger, (n)} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, \Delta_n) \text{ and } X_{t_k}^{\dagger, (n)} \neq \dagger \right\} \\ &= \left\{ X_{t_k}^{\dagger, (n)} \notin \bigcup_{j: m_j \in G^{\text{large}}} B(m_j, \Delta_n) \text{ and } X_s^{\dagger, (n)} \in \bigcup_{j: m_j \in G} \Omega_j \ \forall s \in [0, t_k] \right\} \end{aligned}$$

so the conclusion of the proof follows directly from (C.49).  $\square$

## D Proofs of Lemma B.1, B.2, B.3

*Proof of Lemma B.1.* For any  $\epsilon > 0$ ,

$$\mathbb{P}\left(U(\epsilon) > \frac{1}{b(\epsilon)}\right) = \left(1 - a(\epsilon)\right)^{\lfloor 1/b(\epsilon) \rfloor}.$$

By taking logarithm on both sides, we have

$$\begin{aligned} \ln \mathbb{P}\left(U(\epsilon) > \frac{1}{b(\epsilon)}\right) &= \lfloor 1/b(\epsilon) \rfloor \ln(1 - a(\epsilon)) \\ &= \frac{\lfloor 1/b(\epsilon) \rfloor \ln(1 - a(\epsilon))}{-a(\epsilon)} \frac{-a(\epsilon)}{b(\epsilon)}. \end{aligned}$$

Since  $\lim_{x \rightarrow 0} \frac{\ln(1+x)}{x} = 1$ , we know that for  $\epsilon$  sufficiently small, we will have

$$-c \frac{a(\epsilon)}{b(\epsilon)} \leq \ln \mathbb{P}\left(U(\epsilon) > \frac{1}{b(\epsilon)}\right) \leq -\frac{a(\epsilon)}{c \cdot b(\epsilon)}. \quad (\text{D.1})$$

By taking exponential on both sides, we conclude the proof.  $\square$

*Proof of Lemma B.2.* To begin with, for any  $\epsilon > 0$  we have

$$\mathbb{P}\left(U(\epsilon) \leq \frac{1}{b(\epsilon)}\right) = 1 - \mathbb{P}\left(U(\epsilon) > \frac{1}{b(\epsilon)}\right).$$

Using bound (D.1), we know that for  $\epsilon$  sufficiently small,  $\mathbb{P}(U(\epsilon) > 1/b(\epsilon)) \geq \exp(-c \cdot a(\epsilon)/b(\epsilon))$ . The upper bound follows from the generic bound  $1 - \exp(-x) \leq x$ ,  $\forall x \in \mathbb{R}$  with  $x = c \cdot a(\epsilon)/b(\epsilon)$ .

Now we move onto the lower bound. Again, from bound (D.1), we know that for sufficiently small  $\epsilon$ , we will have

$$\mathbb{P}\left(U(\epsilon) \leq \frac{1}{b(\epsilon)}\right) \geq 1 - \exp\left(-\frac{1}{\sqrt{c}} \cdot \frac{a(\epsilon)}{b(\epsilon)}\right).$$

Due to the assumption that  $\lim_{\epsilon \downarrow 0} a(\epsilon)/b(\epsilon) = 0$  and the fact that  $1 - \exp(-x) \geq \frac{x}{\sqrt{c}}$  for  $x > 0$  sufficiently close to 0, we will have (for  $\epsilon$  small enough)  $\mathbb{P}\left(U(\epsilon) \leq \frac{1}{b(\epsilon)}\right) \geq \frac{1}{c} \cdot \frac{a(\epsilon)}{b(\epsilon)}$ .  $\square$

*Proof of Lemma B.3.* Let  $a_k \triangleq x_k - \tilde{x}_k$ . From the fact that  $g \in \mathcal{C}^2$  and Taylor expansion of  $g'$ , one can easily see that

$$\begin{aligned} a_k &= \eta \sum_{j=1}^k \left( g'(\tilde{x}_{j-1}) - g'(x_{j-1}) \right) + \eta(z_1 + \cdots + z_k) + x - \tilde{x}; \\ \Rightarrow |a_k| &\leq \eta C(|a_0| + \cdots + |a_{k-1}|) + \tilde{c}. \end{aligned}$$

The desired bound then follows immediately from Gronwall's inequality (see Theorem 68, Chapter V of [24], where we let function  $\alpha(t)$  be  $\alpha(t) = |a_{\lfloor t \rfloor}|$ ).  $\square$

## E Notations

Table E.1 lists the notations used in Appendix B.

Table E.1: Summary of notations frequently used in Appendix B

$[k]$	$\{1, 2, \dots, k\}$	
$\eta$	Learning rate (gradient descent step size)	
$b$	Truncation threshold of stochastic gradient	
$\epsilon$	An accuracy parameter; typically used to denote an $\epsilon$ -neighborhood of $s_i, m_i$	
$\delta$	A threshold parameter used to define <i>large</i> noises	
$\bar{\epsilon}$	A constant defined for (B.15)-(B.16). Since $\bar{\epsilon} < \epsilon_0$ , in (B.1) the claim holds for $ x - y  < \bar{\epsilon}$ . Note that the value of the constant $\bar{\epsilon}$ does not vary with our choice of $\eta, \epsilon, \delta$ .	
$M$	Upper bound of $ f' $ and $ f'' $	(B.3)
$L$	Radius of training domain	(B.3)
$\Omega$	The open interval $(s_-, s_+)$ ; a simplified notation for $\Omega_i$	
$\varphi, \varphi_c$	$\varphi_c(w) \triangleq \varphi(w, c) \triangleq (w \wedge c) \vee (-c)$	truncation operator at level $c > 0$
$Z_n^{\leq \delta, \eta}$	$Z_n \mathbb{1}\{\eta Z_n  \leq \delta\}$	“small” noise (B.7)
$Z_n^{> \delta, \eta}$	$Z_n \mathbb{1}\{\eta Z_n  > \delta\}$	“large” noise (B.8)
$T_j^\eta(\delta)$	$\min\{n > T_{j-1}^\eta(\delta) : \eta Z_n  > \delta\}$	arrival time of $j$ -th large noise (B.9)
$W_j^\eta(\delta)$	$Z_{T_j^\eta(\delta)}$	size of $j$ -th large noise (B.10)
$X_n^\eta(x)$	$X_{n+1}^\eta(x) = \varphi_L\left(X_n^\eta(x) - \varphi_b(\eta(f'(X_n^\eta(x)) - Z_{n+1}))\right), \quad X_0^\eta(x) = x$	SGD
$\mathbf{y}_n^\eta(x)$	$\mathbf{y}_n^\eta(x) = \mathbf{y}_{n-1}^\eta(x) - \eta f'(\mathbf{y}_{n-1}^\eta(x)), \quad \mathbf{y}_0^\eta(x) = x$	GD
$Y_n^\eta(x)$	$\mathbf{y}_n^\eta(x)$ perturbed by large noises $(\mathbf{T}^\eta(\delta), \mathbf{W}^\eta(\delta))$	GD + large jump
$\tilde{\mathbf{y}}_n^\eta(x; \mathbf{t}, \mathbf{w})$	$\mathbf{y}_n^\eta(x)$ perturbed by noise vector $(\mathbf{t}, \mathbf{w})$	perturbed GD
$\mathbf{x}^\eta(t, x)$	$d\mathbf{x}^\eta(t; x) = -\eta f'(\mathbf{x}^\eta(t; x))dt, \quad \mathbf{x}^\eta(0; x) = x$	ODE
$\mathbf{x}(t, x)$	$\mathbf{x}^1(t, x)$	
$\tilde{\mathbf{x}}^\eta(t, x; \mathbf{t}, \mathbf{w})$	$\mathbf{x}^\eta(t, x)$ perturbed by noise vector $(\mathbf{t}, \mathbf{w})$	perturbed ODE
$A(n, \eta, \epsilon, \delta)$	$\left\{ \max_{k \in [n \wedge (T_1^\eta(\delta) - 1)]} \eta Z_1 + \dots + Z_k  \leq \epsilon \right\}$	(B.28)
$r$	$r \triangleq \min\{-s_-, s_+\}$ . Effective radius of the attraction field $\Omega$ .	
$l^*$	$l^* \triangleq \lceil r/b \rceil$ . The minimum number of jumps required to escape $\Omega$ when starting from its local minimum $m = 0$ .	
$h(\mathbf{w}, \mathbf{t})$	A mapping defined as $h(\mathbf{w}, \mathbf{t}) = \tilde{\mathbf{x}}(t_{l^*}, 0; \mathbf{t}, \mathbf{w})$ .	(B.21)-(B.22)
$\bar{t}, \bar{\delta}$	Necessary conditions for $h(\mathbf{w}, \mathbf{t})$ to be outside of $\Omega$	(B.21)-(B.22)
$\hat{t}(\epsilon)$	$\hat{t}(\epsilon) \triangleq c_1 \log(1/\epsilon)$ . The quantity $\hat{t}(\epsilon)/\eta$ provides an upper bound for the time it takes $\mathbf{x}^\eta$ to return to $2\epsilon$ -neighborhood of local minimum $m = 0$ when starting from somewhere $\epsilon$ -away from $s_-, s_+$ . See (B.23).	
$E(\epsilon)$	$\{(\mathbf{w}, \mathbf{t}) \subseteq \mathbb{R}^{l^*} \times \mathbb{R}_+^{l^*-1} : h(\mathbf{w}, \mathbf{t}) \notin [(s_- - \epsilon) \vee (-L), (s_+ + \epsilon) \wedge L]\}$	
$p(\epsilon, \delta, \eta)$	The probability that, for $\mathbf{t} = (T_j^\eta(\delta) - 1)_{j=1}^{l^*}$ and $\mathbf{w} = (\eta W_j^\eta(\delta))_{j=1}^{l^*}$ , we have $(\mathbf{w}, \mathbf{t}) \in E(\epsilon)$ conditioning on $\{T_1^\eta(\delta) = 1\}$ . Intuitively speaking, it characterizes the probability that the first $l^*$ <i>large</i> noises alone can drive the ODE out of the attraction field. Defined in (B.53).	

$\nu_\alpha$  The Borel measure on  $\mathbb{R}$  with density

$$\nu_\alpha(dx) = \mathbb{1}\{x > 0\} \frac{\alpha p_+}{x^{\alpha+1}} + \mathbb{1}\{x < 0\} \frac{\alpha p_-}{|x|^{\alpha+1}}$$

where  $p_-, p_+$  are constants in Assumption 2 in the main paper.

$\mu$  The product measure  $\mu = (\nu_\alpha)^{l^*} \times (\mathbf{Leb}_+)^{l^*-1}$ .

$\sigma(\eta)$   $\min\{n \geq 0 : X_n^\eta \notin \Omega\}$ . first exit time

$H(x)$   $\mathbb{P}(|Z_1| > x) = x^{-\alpha} L(x)$

$T_{\text{return}}(\epsilon, \eta)$   $\min\{n \geq 0 : X_n^\eta(x) \in [-2\epsilon, 2\epsilon]\}$