

# ADAPTIVE METHODS THROUGH THE LENS OF SDEs: THEORETICAL INSIGHTS ON THE ROLE OF NOISE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Despite the vast empirical evidence supporting the efficacy of adaptive optimization methods in deep learning, their theoretical understanding is far from complete. This work introduces novel SDEs for commonly used adaptive optimizers: SignSGD, RMSprop(W), and Adam(W). These SDEs offer a quantitatively accurate description of these optimizers and help illuminate an intricate relationship between adaptivity, gradient noise, and curvature. Our novel analysis of SignSGD highlights a noteworthy and precise contrast to SGD in terms of convergence speed, stationary distribution, and robustness to heavy-tail noise. We extend this analysis to AdamW and RMSpropW, for which we observe that the role of noise is much more complex. Crucially, we support our theoretical analysis with experimental evidence by verifying our insights: this includes numerically integrating our SDEs using Euler-Maruyama discretization on various neural network architectures such as MLPs, CNNs, ResNets, and Transformers. Our SDEs accurately track the behavior of the respective optimizers, especially when compared to previous SDEs derived for Adam and RMSprop. We believe our approach can provide valuable insights into best training practices and novel scaling rules.

## 1 INTRODUCTION

Adaptive optimizers lay the foundation for effective training of modern deep learning models. These methods are typically employed to optimize an objective function expressed as a sum of losses across  $N$  individual data points:  $\min_{x \in \mathbb{R}^d} [f(x) := \frac{1}{N} \sum_{i=1}^N f_i(x)]$ , where  $f, f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $i = 1, \dots, N$ . Due to the practical difficulties of selecting the learning rate of stochastic gradient descent, adaptive methods have grown in popularity over the past decade. At a high level, these optimizers adjust the learning rate for each parameter based on the historical gradients. Popular optimizers that belong to this family are RMSprop (Tieleman and Hinton, 2012), Adam (Kingma and Ba, 2015), SignSGD (Bernstein et al., 2018), AdamW (Loshchilov and Hutter, 2019), and many other variants. SignSGD is often used for compressing gradients in distributed machine learning (Karimireddy et al., 2019a), but it also has gained popularity due to its connection to RMSprop and Adam (Balles and Hennig, 2018). The latter algorithms have emerged as the standard methods for training modern large language models, partly because of enhancements in signal propagation (Noci et al., 2022).

Although adaptive methods are widely favored in practice, their theoretical foundations remain enigmatic. Recent research has illuminated some of their advantages: Zhang et al. (2020b) demonstrated how gradient clipping addresses heavy-tailed gradient noise, Pan and Li (2022) related the success of Adam over SGD to sharpness, and Yang et al. (2024) showed that adaptive methods are more resilient to poor learning rate tuning than SGD. At the same time, many optimization studies focus on worst-case convergence rates: These rates (e.g., Défossez et al. (2022)) are valuable, yet they provide an incomplete depiction of algorithm behavior, showing no quantifiable advantage over standard SGD. One particular aspect still lacking clarity is the precise role of noise in the algorithm trajectory.

Our investigation aims to study how gradient noise influences the dynamics of adaptive optimizers and how it impacts their asymptotic behaviors in terms of expected loss and stationary distribution. In particular, we want to understand which algorithms are more resilient to high (possibly heavy-tailed) gradient noise levels. To do this, we rely on stochastic differential equations (SDEs) which have become popular in the literature to study the behavior of optimization algorithms (Li et al., 2017; Jastrzebski et al., 2018). These continuous-time models unlock powerful tools from Itô calculus, enabling us to establish convergence bounds, determine stationary distributions, unveil implicit

regularization, and elucidate the intricate interplay between landscape and noise. Notably, SDEs facilitate direct comparisons between optimizers by explicitly illustrating how each hyperparameter and certain landscape features influence their dynamics (Orvieto and Lucchi, 2019; Malladi et al., 2022; Compagnoni et al., 2024).

We begin by analyzing SignSGD, showing how the ratio between the gradient and the level of gradient noise affects its dynamics and elucidating the impact of noise at convergence. After examining the case where the gradient noise has an infinite variance, we extend our analysis to Adam and RMSprop with *decoupled* weight decay (Loshchilov and Hutter, 2019) – i.e. AdamW and RMSpropW: for both, we refine batch size scaling rules and compare the role of noise to SignSGD. Our analysis provides some theoretical grounding for the resilience of these adaptive methods to high noise levels. Importantly, we highlight that Adam and RMSprop are byproducts of our analysis and that our novel SDEs are derived under much weaker and more realistic assumptions than those in the literature (Malladi et al., 2022).

**Contributions** We identify our key contributions as follows:

1. We derive the first SDE for SignSGD under very general assumptions: We show that SignSGD exhibits three different phases of the dynamics and characterize the loss behavior in these phases, including the stationary distribution and asymptotic loss value;
2. We prove that for SignSGD, noise inversely affects the convergence rate of both the loss and the iterates. Differently, it has a linear impact on the asymptotic expected loss and the asymptotic variance of the iterates. This is in contrast to SGD, where noise does not influence the convergence speed, but it has a quadratic effect on the loss and variance of the iterates. Finally, we show that, even if the noise has infinite variance, SignSGD is resilient: its performance is only marginally impacted. In the same conditions, SGD diverges;
3. We derive new, improved, SDEs for AdamW and RMSpropW and use them to (i) show a novel batch size scaling rule and (ii) inspect the stationary distribution and stationary loss value in convex quadratics. In particular, we dive into the properties of weight decay: while for vanilla Adam and RMSprop the effect of noise at convergence mimics SignSGD, something different happens in AdamW and RMSpropW — Due to an intricate interaction between noise, curvature, and regularization, *decoupled* weight decay plays a crucial stabilization role at high noise levels near the minimizer;
4. We empirically verify every theoretical insight we derive. Importantly, we integrate our SDEs with Euler-Maruyama to confirm that our SDEs faithfully track their respective optimizers. We do so on an MLP, a CNN, a ResNet, and a Transformer. For RMSprop and Adam, our SDEs exhibit superior modeling power than the SDEs already in the literature. [We emphasize that while our results rely on certain regularity assumptions for loss functions and gradient noise, their applicability extends beyond these. For example, we validate our novel scaling rule for AdamW on a Pythia-like 160M LLM \(Biderman et al., 2023\) trained on 2.5B tokens from the SlimPajama dataset \(Soboleva et al., 2023\).](#)

## 2 RELATED WORK

**SDE approximations and applications.** (Li et al., 2017) introduced a formal theoretical framework aimed at deriving SDEs that effectively model the inherent stochastic nature of optimizers. Ever since, SDEs have found several applications in the field of machine learning, for instance in connection with *stochastic optimal control* to select the stepsize (Li et al., 2017; 2019) and batch size (Zhao et al., 2022), the derivation of *convergence bounds* and *stationary distributions* (Compagnoni et al., 2023; 2024), *implicit regularization* (Smith et al., 2021), and *scaling rules* (Jastrzebski et al., 2018). Previous work by Malladi et al. (2022) has already made strides in deriving SDE models for RMSprop and Adam, albeit under certain restrictive assumptions. They establish a scaling rule which they assert remains valid throughout the entirety of the dynamics. Unfortunately, their derivation is based on the approach of Jastrzebski et al. (2018) which is problematic in the general case (See Appendix E for a detailed discussion). Additionally, we demonstrate that the SDEs derived in Malladi et al. (2022) are only accurate around minima, indicating that their scaling rule is not *globally* valid. (Zhou et al., 2020a) also claimed to have derived a Lévy SDE for Adam. Unfortunately, the quality of their SDE approximation does not come with theoretical guarantees: Their SDE has random coefficients, an approach which is theoretically sound in very limited settings (Kohatsu-Higa et al., 1997; Bishop

and Del Moral, 2019). Xie et al. (2022) modeled AdamW with an SDE and used it to disentangle the effects of learning rate adaptivity and momentum on saddle-point escaping and flat minima selection: Unfortunately, their SDE is not derived within any formal framework and therefore does not come with formal approximation guarantees. Finally, Zhou et al. (2024) informally presented an SDE for (only) the parameters of AdamW: this is achieved under strong assumptions and various approximations, some of which are hard to motivate formally.

**Influence of noise on convergence.** Several empirical papers demonstrate that adaptive algorithms adjust better to the noise during training. Specifically, (Zhang et al., 2020b) noticed a consistent gap in the performance of SGD and Adam on language models and connected that phenomenon with heavy-tailed noise distributions. (Pascanu et al., 2013) suggests using gradient clipping to deal with heavy tail noise, and consequently several follow-up works analyzed clipped SGD under heavy-tailed noise (Zhang et al., 2020a; Mai and Johansson, 2021; Puchkin et al., 2024). Kunstner et al. (2024) present thorough numerical experiments illustrating that a significant contributor to heavy-tailed noise during language model training is class imbalance, where certain words occur much more frequently than others. They demonstrate that adaptive optimization methods such as Adam and SignSGD can better adapt to such class imbalances. However, the theoretical understanding of the influence of noise in the context of adaptive algorithms is much more limited. The first convergence results on Adam and RMSprop were derived under bounded stochastic gradients assumption (De et al., 2018; Zaheer et al., 2018; Chen et al., 2019; Défossez et al., 2022). Later, this noise model was relaxed to weak growth condition (Zhang et al., 2022; Wang et al., 2022) and its coordinate-wise version (Hong and Lin, 2023; Wang et al., 2024) and sub-gaussian noise (Li et al., 2023a). SignSGD and its momentum version Signum were originally studied as a method for compressed communication (Bernstein et al., 2018) under bounded variance assumption, but with a requirement of large batches. Several works provided counterexamples where SignSGD fails to converge if stochastic and full gradients are not correlated enough (Karimireddy et al., 2019b; Safaryan and Richtarik, 2021). In the case of AdamW, (Zhou et al., 2022; 2024) provided convergence guarantees under restrictive assumptions such as bounded gradient and bounded noise. All aforementioned results only show that SignSGD, Adam, and RMSprop at least do not perform worse than vanilla SGD. None of them studied how noise affects the dynamics of the algorithm: In this work, we attempt to close this gap.

### 3 FORMAL STATEMENTS & INSIGHTS: THE SDES

This section provides the general formulations of the SDEs of SignSGD (Theorem 3.2) and AdamW (Theorem 3.12). Due to the technical nature of the analysis, we refer the reader to the appendix for the complete formal statements and proofs and only provide a sketch of the proof of key results.

**Assumptions and notation.** In this section, we collect most of the notation and assumptions used in the paper. All our analysis take place on a filtered probability space  $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$ . The batches  $\gamma$  are of size  $B \geq 1$  and modeled as i.i.d. random variables uniformly distributed on  $\{1, \dots, N\}$ . We assume that the stochastic gradient  $\nabla f_\gamma(x) := \frac{1}{B} \sum_{i \in \gamma} \nabla(f_i(x))$  can be decomposed as  $\nabla f(x) + Z(x)$ , where  $\nabla f(x)$  is the full gradient and  $Z(x)$  is the batch noise. We assume that  $\mathbb{E}[Z(x)] = 0$  and unless we study the cases where the gradient variance is unbounded, we write  $Cov(Z(x)) = \Sigma(x)$ <sup>1</sup> s.t.  $\sqrt{\Sigma(x)}$  is bounded, Lipschitz, satisfies affine growth, and together with its derivatives, it grows at most polynomially fast (Definition 2.5 in Malladi et al. (2022)). Importantly, we assume that  $Z(x)$  has a bounded and smooth probability density function whose derivatives are all integrable: A common assumption in the literature is for  $Z(x)$  to be Gaussian<sup>2</sup> (Ahn et al., 2012; Chen et al., 2014; Mandt et al., 2016; Stephan et al., 2017; Zhu et al., 2019a; Wu et al., 2020; Xie et al., 2021), while our assumption allows for heavy-tailed distributions such as the Student’s t. Specifically, Li et al. (2017); Mertikopoulos and Staudigl (2018); Raginsky and Bouvrie (2012); Zhu et al. (2019b); Mandt et al. (2016); Ahn et al. (2012); Jastrzebski et al. (2018) use a Gaussian noise with constant covariance matrix to model batch noise. To derive the stationary distribution around an optimum, we approximate the loss function with a quadratic convex function  $f(x) = \frac{1}{2}x^\top Hx$  as commonly done in the literature (Ge et al., 2015; Levy, 2016; Jin et al., 2017; Poggio et al., 2017; Mandt et al., 2017; Compagnoni et al., 2023). Finally,  $\eta > 0$  is the step size, the  $\beta$ s refer to momentum parameters,  $\theta > 0$  is the (decoupled)  $L^2$ -regularization parameter, and  $\epsilon > 0$  is a scalar used for numerical stability. Finally, we use  $W_t$  to indicate a Brownian motion.

<sup>1</sup>We omit the size of the batch  $\gamma$  unless relevant.

<sup>2</sup>See Jastrzebski et al. (2018) for the justification why this might be the case.

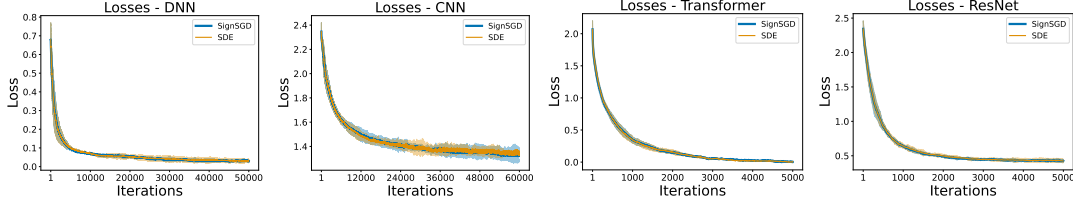


Figure 1: Comparison of SignSGD and its SDE in terms of  $f(x)$ : Our SDE successfully tracks the dynamics of SignSGD on several architectures, datasets, and hyperparameters: DNN on the Breast Cancer dataset (Left); CNN on MNIST (Center-Left); Transformer on MNIST (Center-Right); ResNet on CIFAR-10 (Right).

The following definition formalizes the idea that an SDE can be a “good model” to describe an optimizer. It is drawn from the field of numerical analysis of SDEs (see Mil’shtein (1986)) and it quantifies the disparity between the discrete and the continuous processes.

**Definition 3.1** (Weak Approximation). A continuous-time stochastic process  $\{X_t\}_{t \in [0, T]}$  is an order  $\alpha$  weak approximation (or  $\alpha$ -order SDE) of a discrete stochastic process  $\{x_k\}_{k=0}^{\lfloor T/\eta \rfloor}$  if for every polynomial growth function  $g$ , there exists a positive constant  $C$ , independent of the stepsize  $\eta$ , such that  $\max_{k=0, \dots, \lfloor T/\eta \rfloor} |\mathbb{E}g(x_k) - \mathbb{E}g(X_{k\eta})| \leq C\eta^\alpha$ .

### 3.1 SIGNSGD SDE

In this section, we derive an SDE model for SignSGD, which we believe to be a novel addition to the existing literature. This derivation will reveal the unique manner in which noise influences the dynamics of SignSGD. First, we recall the update equation of SignSGD:

$$x_{k+1} = x_k - \eta \text{sign}(\nabla f_{\gamma_k}(x_k)). \quad (1)$$

The following theorem derives a formal continuous-time model for SignSGD.

**Theorem 3.2** (Informal Statement of Theorem C.16). *Under sufficient regularity conditions, the solution of the following SDE is an order 1 weak approximation of the discrete update of SignSGD:*

$$dX_t = -(1 - 2\mathbb{P}(\nabla f_\gamma(X_t) < 0))dt + \sqrt{\eta} \sqrt{\bar{\Sigma}(X_t)} dW_t, \quad (2)$$

where  $\bar{\Sigma}(x)$  is the noise covariance  $\bar{\Sigma}(x) = \mathbb{E}[\xi_\gamma(x)\xi_\gamma(x)^\top]$ , and  $\xi_\gamma(x) := \text{sign}(\nabla f_\gamma(x)) - 1 + 2\mathbb{P}(\nabla f_\gamma(x) < 0)$  is the noise of  $\text{sign}(\nabla f_\gamma(x))$ .

*Proof idea.* One needs to prove that the first and second moments of the increments of the discretization of the SDE match those of SignSGD up to an error of order  $\mathcal{O}(\eta)$  and  $\mathcal{O}(\eta^2)$ , respectively.  $\square$

For **didactic reasons**, we next present a corollary of Theorem 3.2 that provides a more interpretable SDE. To do so, we model the batch noise with a Gaussian distribution with constant covariance matrix,<sup>3</sup> which is a common approach in the literature (Li et al., 2017; Mertikopoulos and Staudigl, 2018; Raginsky and Boudrie, 2012; Zhu et al., 2019b; Mandt et al., 2016; Ahn et al., 2012; Jastrzebski et al., 2018). Figure 1 shows the empirical validation of this model for various neural network classes: All details are presented in Appendix F.

**Corollary 3.3** (Informal Statement of Corollary C.19). *Under the assumptions of Theorem 3.2, and that the stochastic gradient is  $\nabla f_\gamma(x) = \nabla f(x) + Z$  such that  $Z \sim \mathcal{N}(0, \Sigma)$ ,  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$ , the following SDE provides a 1 weak approximation of the discrete update of SignSGD*

$$dX_t = -\text{Erf}\left(\frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}}\right) dt + \sqrt{\eta} \sqrt{I_d - \text{diag}\left(\text{Erf}\left(\frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}}\right)\right)^2} dW_t, \quad (3)$$

where the error function  $\text{Erf}(x) := \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$  and the square are applied component-wise.

While Eq. 3 may appear intricate at first glance, it becomes apparent upon closer inspection that the properties of the  $\text{Erf}(\cdot)$  function enable a detailed exploration of the dynamics of SignSGD. In particular, we demonstrate that the dynamics of SignSGD can be categorized into three distinct phases. The left of Figure 2 empirically verifies this result on a convex quadratic function.

<sup>3</sup>See Section C.5 for more realistic noise structures.

**Lemma 3.4.** Under the assumptions of Coroll. 3.3,  $Y_t := \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}}$ , and  $|\cdot|$  applied element-wise

1. **Phase 1:** If  $|Y_t| > \frac{3}{2}$ , the SDE coincides with the ODE of SignGD:

$$dX_t = -\text{sign}(\nabla f(X_t))dt; \quad (4)$$

2. **Phase 2:** If  $1 < |Y_t| < \frac{3}{2}$ ,

$$(a) -mY_t - \mathbf{q}^+ \leq \frac{d\mathbb{E}[X_t]}{dt} \leq -mY_t - \mathbf{q}^-;$$

$$(b) \text{ For any } a > 0, \mathbb{P}[\|X_t - \mathbb{E}[X_t]\|_2^2 > a] \leq \frac{\eta}{a} (d - \|mY_t + \mathbf{q}^-\|_2^2);$$

3. **Phase 3:** If  $|Y_t| < 1$ , the SDE is

$$dX_t = -\sqrt{\frac{2}{\pi}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) dt + \sqrt{\eta} \sqrt{I_d - \frac{2}{\pi} \text{diag}(\Sigma^{-\frac{1}{2}} \nabla f(X_t))^2} dW_t. \quad (5)$$

**Remark:** For ease of reading, we will *informally* refer to the gradient  $\nabla f(x)$  as the “signal” and to  $\Sigma$  as the “noise”. Then, Lemma 3.4 tells us that the behavior of SignSGD depends on the size of the “signal-to-noise” ratio. In particular, the SDE itself shows that in Phase 3, the inverse of the scale of the noise  $\Sigma^{-\frac{1}{2}}$  premultiplies  $\nabla f(x)$ , affecting the rate of descent. This is not the case for SGD where  $\Sigma$  only influences the diffusion term.<sup>5</sup> To better understand the role of the noise, we study how it affects the dynamics of the loss on strongly convex functions and compare it with SGD. The dynamics of  $\mathbb{E}[\|\nabla f(X_t)\|_2^2]$  for general non-convex smooth functions is presented in Lemma C.24.

**Lemma 3.5.** Let  $f$  be  $\mu$ -strongly convex,  $\text{Tr}(\nabla^2 f(x)) \leq \mathcal{L}_\tau$ ,  $\sigma_{\max}^2$  be the maximum eigenvalue of  $\Sigma$ , and  $S_t := f(X_t) - f(X_*)$ . Then, during

1. **Phase 1,**  $S_t \leq \frac{1}{4} (\sqrt{\mu t} - 2\sqrt{S_0})^2$ , so SignSGD stays in this phase for at most  $t_* = 2\sqrt{\frac{S_0}{\mu}}$ ;

2. **Phase 2** as  $\Delta := \left( \frac{m}{\sqrt{2}\sigma_{\max}} + \frac{\eta\mu m^2}{4\sigma_{\max}^2} \right)$ :

$$\mathbb{E}[S_t] \leq S_0 e^{-2\mu\Delta t} + \frac{\eta}{2} \frac{(\mathcal{L}_\tau - \mu d \hat{q}^2)}{2\mu\Delta} (1 - e^{-2\mu\Delta t});$$

3. **Phase 3** as  $\Delta := \left( \sqrt{\frac{2}{\pi}} \frac{1}{\sigma_{\max}} + \frac{\eta}{\pi} \frac{\mu}{\sigma_{\max}^2} \right)$ :

$$\mathbb{E}[S_t] \leq S_0 e^{-2\mu\Delta t} + \frac{\eta}{2} \frac{\mathcal{L}_\tau}{2\mu\Delta} (1 - e^{-2\mu\Delta t}).$$

*Proof idea.* For each **Phase**, we use the respective SDE of SignSGD from Lemma 3.4 to derive the SDE of  $S_t$  via Itô’s lemma. Then, we take its expectation to obtain the ODE of  $\mathbb{E}[S_t]$  and leverage the assumptions to establish a bound.  $\square$

As per Eq. 4, during Phase 1 SignSGD behaves like SignGD: Lemma 3.5 shows that, consistently with the analysis of SignGD in (Ma et al., 2022), such a strong decrease in the loss value explains the fast initial convergence of the optimizer as well as of RMSprop and Adam. In this phase, the loss undergoes a steady decrease which ensures the emergence of Phase 2 which in turn triggers that of Phase 3 which is characterized by an exponential decay to an asymptotic loss level: As a practical example, we verify the dynamics of the expected loss around a minimum in the center-left of Fig. 2.

**Lemma 3.6.** For SGD, the expected loss satisfies:  $\mathbb{E}[S_t] \leq S_0 e^{-2\mu t} + \frac{\eta}{2} \frac{\mathcal{L}_\tau \sigma_{\max}^2}{2\mu} (1 - e^{-2\mu t})$ .

<sup>4</sup>Let  $m$  and  $q_1$  are the slope and intercept of the line secant to the graph of  $\text{Erf}(x)$  between the points  $(1, \text{Erf}(1))$  and  $(\frac{3}{2}, \text{Erf}(\frac{3}{2}))$ , while  $q_2$  is the intercept of the line tangent to the graph of  $\text{Erf}(x)$  and slope  $m$ ,  $(\mathbf{q}^+)_i := \begin{cases} q_2 & \text{if } \partial_i f(x) > 0 \\ -q_1 & \text{if } \partial_i f(x) < 0 \end{cases}$ ,  $(\mathbf{q}^-)_i := \begin{cases} q_1 & \text{if } \partial_i f(x) > 0 \\ -q_2 & \text{if } \partial_i f(x) < 0 \end{cases}$ , and  $\hat{q} := \max(q_1, q_2)$ .

<sup>5</sup>The SDE of SGD is  $dX_t = -\nabla f(X_t)dt + \sqrt{\eta}\Sigma^{\frac{1}{2}}dW_t$ .

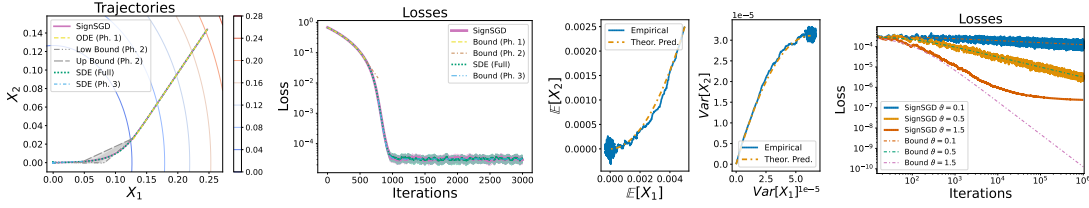


Figure 2: **Phases of SignSGD**: The ODE of Phase 1 and the SDE of Phase 3 overlap with the “Full” SDE as per **Lemma 3.4**. In Phase 2, the dynamics satisfies the prescribed bounds (Left); **Phases of the Loss**: The bounds derived in **Lemma 3.5** for the loss during the different phases correctly track the loss evolution (Center-Left); The **dynamics of the moments of  $X_t$**  predicted in **Lemma 3.7** track the empirical ones (Center-Right); If the **schedulers** satisfy the condition in **Lemma 3.9**, the loss decays to 0 as prescribed. Otherwise, the loss does not converge to 0 (Right). For each figure,  $f(x) = \frac{x^\top H x}{2}$  for  $H = \text{diag}(1, 2)$ ,  $\eta = 0.001$ , and  $\Sigma = \sigma^2 I_2$  where  $\sigma = 0.1$ .

**Remark:** The two key observations are that:

1. Both in Phase 2 and Phase 3, the noise level  $\sigma_{\max}$  inversely affects the exponential convergence speed, while this trend is not observed with SGD;
2. The asymptotic loss of SignSGD is (almost) linear in  $\sigma_{\max}$  while that of SGD is quadratic. Indeed, the asymptotic value of  $\mathbb{E}[S_t]$  in Phase 3 scales with  $\frac{1}{\Delta} = \frac{\pi \sigma_{\max}}{\sqrt{2\pi} + \frac{\eta \mu}{\sigma_{\max}}}$ : when the noise  $\sigma_{\max}$  dominates the learning rate  $\eta$  and/or the minimum eigenvalue  $\mu$  of the Hessian, or in general when  $\frac{\eta \mu}{\sigma_{\max}} \sim 0$ , we can conclude that the scaling is (almost) linear in  $\sigma_{\max}$ .

Additionally, we characterize the stationary distribution of SignSGD around a minimum. To do this, we study the behavior of SignSGD on a quadratic loss function, which is a common approach in the literature (Ge et al., 2015; Levy, 2016; Jin et al., 2017; Poggio et al., 2017; Mandt et al., 2017; Compagnoni et al., 2023). Empirical validation is provided in the center-right of Figure 2.

**Lemma 3.7.** Let  $H = \text{diag}(\lambda_1, \dots, \lambda_d)$  and  $M_t := e^{-2\left(\sqrt{\frac{2}{\pi}}\Sigma^{-\frac{1}{2}}H + \frac{\eta}{\pi}\Sigma^{-1}H^2\right)t}$ . Then,

1.  $\mathbb{E}[X_t] = e^{-\sqrt{\frac{2}{\pi}}\Sigma^{-\frac{1}{2}}Ht} X_0$ ;
2.  $\text{Cov}[X_t] = \left(M_t - e^{-2\sqrt{\frac{2}{\pi}}\Sigma^{-\frac{1}{2}}Ht}\right) X_0^2 + \frac{\eta}{2} \left(\sqrt{\frac{2}{\pi}}I_d + \frac{\eta}{\pi}H\Sigma^{-\frac{1}{2}}\right)^{-1} H^{-1}\Sigma^{\frac{1}{2}}(I_d - M_t)$ .

Therefore, we have that the stationary distribution of SignSGD is:

$$\left(\mathbb{E}[X_\infty], \text{Cov}[X_\infty]\right) = \left(0, \frac{\eta}{2} \left(\sqrt{\frac{2}{\pi}}I_d + \frac{\eta}{\pi}H\Sigma^{-\frac{1}{2}}\right)^{-1} H^{-1}\Sigma^{\frac{1}{2}}\right).$$

*Proof idea.* For the  $\mathbb{E}[X_t]$ , we take the expected value of the SDE of Phase 3 from Lemma 3.4 and integrate the resulting ODE. For  $\text{Cov}[X_t]$ , we derive the SDE of  $X_t X_t^\top$  via Itô’s lemma, take the expectation, and integrate the resulting ODE. Then, we subtract  $\mathbb{E}[X_t]\mathbb{E}[X_t]^\top$ .  $\square$

**Lemma 3.8.** Under the same assumptions as Lemma 3.7, the stationary distribution for SGD is:

$$\mathbb{E}[X_t] = e^{-Ht} X_0 \xrightarrow{t \rightarrow \infty} 0 \quad \text{and} \quad \text{Cov}[X_t] = \frac{\eta}{2} H^{-1} \Sigma (I_d - e^{-2Ht}) \xrightarrow{t \rightarrow \infty} \frac{\eta}{2} H^{-1} \Sigma.$$

As we observed above, the noise inversely affects the convergence rate of the iterates of SignSGD while it does not impact that of SGD. Additionally, while both covariance matrices essentially scale inversely to the Hessian, that of SignSGD scales with  $\Sigma^{\frac{1}{2}}$  while that of SGD scales with  $\Sigma$ .

We conclude this section by presenting a condition on the step size scheduler that ensures the asymptotic convergence of the expected loss to 0 in Phase 3. For general schedulers, we characterize precisely the speed of convergence and the factors influencing it. Empirical validation is provided in the right of Figure 2 for a convex quadratic as we use  $\eta_t^\vartheta = \frac{1}{(t+1)^\vartheta}$  for  $\vartheta \in \{\frac{1}{10}, \frac{1}{2}, \frac{3}{2}\}$ .

**Lemma 3.9.** Under the assumptions of Lemma 3.5, any step size scheduler  $\eta_t$  such that

$$\int_0^\infty \eta_s ds = \infty \text{ and } \lim_{t \rightarrow \infty} \eta_t = 0 \implies \mathbb{E}[f(X_t) - f(X_*)] \xrightarrow{t \rightarrow \infty} \lesssim \frac{\mathcal{L}_\tau \sigma_{\max}}{4\mu} \sqrt{\frac{\pi}{2}} \eta_t \xrightarrow{t \rightarrow \infty} 0. \quad (6)$$

**Remark:** Under the same conditions, SGD satisfies  $\mathbb{E}[f(X_t) - f(X_*)] \xrightarrow{t \rightarrow \infty} \lesssim \frac{\mathcal{L}_\tau \sigma_{\max}^2}{4\mu} \eta_t \xrightarrow{t \rightarrow \infty} 0$ .

**Conclusion:** As noted in Bernstein et al. (2018), the “signal-to-noise” ratio is key in determining the dynamics of SignSGD. Our SDEs help clarify the mechanisms underlying the dynamics of SignSGD: we show that the effect of noise is radically different from SGD: 1) It affects the rate of convergence of the iterates, of the covariance of the iterates, and of the expected loss; 2) The asymptotic loss value and covariance of the iterates scale in  $\Sigma^{\frac{1}{2}}$  while for SGD it does so in  $\Sigma$ . On the one hand, low levels of noise will ensure a faster and steadier loss decrease close to minima for SignSGD than for SGD. On the other, SGD will converge to much lower loss values. A symmetric argument holds for high levels of noise, which suggests that SignSGD is more resilient to high levels of noise.

### 3.1.1 HEAVY-TAILED NOISE

Interestingly, we can replicate the efforts above also in case the noise  $Z(x)$  is heavy-tailed as it is distributed according to a Student’s t distribution. Notably, we derive the SDE for the case where the noise has infinite variance and show how little marginal effect this has on the dynamics of SignSGD.

**Lemma 3.10.** Under the assumptions of Corollary 3.3 but the noise on the gradients  $Z \sim t_\nu(0, I_d)$  where  $\nu \in \mathbb{Z}^+$ : The following SDE is a 1 weak approximation of the discrete update of SignSGD

$$dX_t = -2\Xi \left( \Sigma^{-\frac{1}{2}} \nabla f(X_t) \right) dt + \sqrt{\eta} \sqrt{I_d - 4 \text{diag} \left( \Xi \left( \Sigma^{-\frac{1}{2}} \nabla f(X_t) \right) \right)^2} dW_t, \quad (7)$$

where  $\Xi(x)$  is defined as  $\Xi(x) := x \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\pi\nu}\Gamma(\frac{\nu}{2})} {}_2F_1 \left( \frac{1}{2}, \frac{\nu+1}{2}; \frac{3}{2}; -\frac{x^2}{\nu} \right)$ ,  $\Gamma$  is the gamma function, and  ${}_2F_1$  is the hypergeometric function. Above, the  $\Xi(x)$  and the square are applied component-wise.

We now characterize the dynamics of SignSGD when the noise on the gradient has infinite variance.

**Corollary 3.11.** Under the assumptions of Lemma 3.10 and  $\nu = 2$ , the dynamics in Phase 3 is:

$$dX_t = -\sqrt{\frac{1}{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) dt + \sqrt{\eta} \sqrt{I_d - \frac{1}{2} \text{diag} \left( \Sigma^{-\frac{1}{2}} \nabla f(X_t) \right)^2} dW_t. \quad (8)$$

**Conclusion:** We observe that the dynamics of SignSGD when the noise is Gaussian (Eq. 5) w.r.t. when it is heavy-tailed with unbounded variance (Eq. 8) are very similar: By comparing the constants ( $\sqrt{1/2}$  and  $\sqrt{2/\pi}$ ) in front of the drift terms  $\Sigma^{-\frac{1}{2}} \nabla f(X_t)$ , they are only  $\sim 10\%$  apart, and the diffusion coefficients are comparable. Not only do we once more showcase the resilience of SignSGD to high levels of noise, but in alignment with (Zhang et al., 2020b), we provide theoretical support to the success of Adam in such a scenario where SGD would diverge.

All the results derived above can be extended to this setting.

## 3.2 ADAMW SDE

In the last subsection, we showcased how SDEs can serve as powerful tools to understand the dynamics of the simplest among coordinate-wise adaptive methods: SignSGD. Here, we extend the discussion to Adam with *decoupled* weight decay, i.e. AdamW:

$$\begin{aligned} v_{k+1} &= \beta_2 v_k + (1 - \beta_2) (\nabla f_{\gamma_k}(x_k))^2, & m_{k+1} &= \beta_1 m_k + (1 - \beta_1) \nabla f_{\gamma_k}(x_k), \\ x_{k+1} &= x_k - \eta \frac{\hat{m}_{k+1}}{\sqrt{\hat{v}_{k+1} + \epsilon}} - \eta \theta x_k, & \hat{m}_k &= \frac{m_k}{1 - \beta_1^k}, & \hat{v}_k &= \frac{v_k}{1 - \beta_2^k}, \end{aligned} \quad (9)$$

which, of course, covers Adam, RMSprop, and RMSpropW depending on the values of  $\theta$  and  $\beta_1$ .

The following result proves the SDE of AdamW which we validate in Figure 3 for two simple landscapes and in Figure 4 for a Transformer and a ResNet.

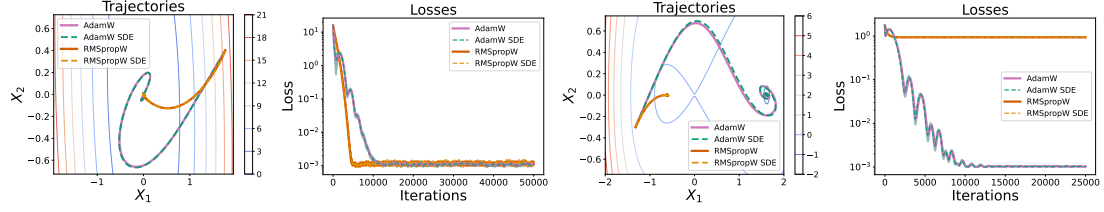


Figure 3: The two images on the left compare the SDEs of AdamW and RMSpropW with the respective optimizers in terms of trajectories and  $f(x)$  for a convex quadratic function while the other two provide a comparison for an embedded saddle. In all cases, we observe good agreements.

**Theorem 3.12** (Informal Statement of Theorem C.53). *Under sufficient regularity conditions,  $\rho_1 = \mathcal{O}(\eta^{-\zeta})$  s.t.  $\zeta \in (0, 1)$ , and  $\rho_2 = \mathcal{O}(1)$ , the order 1 weak approximation of AdamW is:*

$$dX_t = -\frac{\sqrt{\nu_2(t)}}{\nu_1(t)} P_t^{-1} (M_t + \eta \rho_1 (\nabla f(X_t) - M_t)) dt - \theta X_t dt \quad (10)$$

$$dM_t = \rho_1 (\nabla f(X_t) - M_t) dt + \sqrt{\eta} \rho_1 \sqrt{\Sigma(X_t)} dW_t \quad (11)$$

$$dV_t = \rho_2 ((\nabla f(X_t))^2 + \text{diag}(\Sigma(X_t)) - V_t) dt, \quad (12)$$

where  $\beta_i = 1 - \eta \rho_i \sim 1$ ,  $\nu_i(t) = 1 - e^{-\rho_i t}$ ,  $t > t_0$ , and  $P_t = \text{diag} \sqrt{V_t} + \epsilon \sqrt{\nu_2(t)} I_d$ .

$M_t$  and  $V_t$  are the exponential moving averages of the gradient and the squared gradient, respectively.  $P_t^{-1}$  acts as an adaptive preconditioner, scaling the parameter updates  $X_t$  based on the accumulated squared gradients in  $V_t$ . While  $M_t$  is the momentum term and captures the history of gradients to smooth out the updates,  $\eta \rho_1 (\nabla f(X_t) - M_t)$  adjusts  $M_t$  towards the current gradient, ensuring responsiveness to recent changes. Finally  $-\theta X_t dt$  applies regularization by shrinking the parameters.

We highlight that in contrast to Remark 4.3 of Malladi et al. (2022), which suggests that an SDE for Adam is only viable if  $\sigma \gg \|\nabla f(x)\|$  and  $\sigma \sim \frac{1}{\eta}$ , our derivation that does not need these assumptions: See Remark C.46 for a deeper discussion, the implications, and the experimental comparisons.

The following result demonstrates how the asymptotic expected loss of AdamW scales with the noise level. Notably, it introduces the first scaling rule for AdamW, extending the one proposed for Adam in (Malladi et al., 2022) to include weight decay scaling. It is crucial to understand that, unlike the typical approach in the literature (see Jastrzebski et al., 2018; Malladi et al., 2022), our objective in deriving these rules is not to maintain the dynamics of the optimizers or the SDE unchanged. Instead, our goal is to offer a practical strategy for adjusting hyperparameters (e.g., from  $\eta$  to  $\tilde{\eta}$ ) to retain certain performance metrics or optimizer properties as the batch size increases (e.g., from  $B$  to  $\tilde{B}$ ). Therefore, in our upcoming analysis, we aim to derive scaling rules that *preserve* specific relevant aspects of the dynamics, such as the convergence bound on the loss or the speed. See Appendix E for a more detailed discussion motivating our approach.

**Lemma 3.13.** *If  $f$  is  $\mu$ -strongly convex and  $L$ -smooth,  $\text{Tr}(\nabla^2 f(x)) \leq \mathcal{L}_\tau$ ,  $\Sigma(x) = \sigma^2 I_d$ , and  $(\nabla f(x))^2 = \mathcal{O}(\eta)$ ,  $\tilde{\eta} = \kappa \eta$ ,  $\tilde{B} = B \delta$ , and  $\tilde{\rho}_i = \alpha_i \rho_i$ , and  $\tilde{\theta} = \xi \theta$ , AdamW satisfies*

$$\mathbb{E}[f(X_t) - f(X_*)] \stackrel{t \rightarrow \infty}{\leq} \frac{\eta \mathcal{L}_\tau \sigma L}{2} \frac{\kappa}{2\mu \sqrt{B \delta} L + \sigma \xi \theta (L + \mu)}. \quad (13)$$

We derive the novel scaling rule by 1) Preserving the upper bound, which requires that  $\kappa = \sqrt{\delta}$  and  $\xi = \kappa$ ; 2) Preserving the relative speed of  $M_t$ ,  $V_t$  and  $X_t$ , which requires that  $\tilde{\beta}_i = 1 - \kappa^2(1 - \beta_i)$ .

The left of Figure 5 shows the empirical verification of the predicted loss value and scaling rule on a convex quadratic function. Consistently with Lemma 3.13, such a value is bounded w.r.t.  $\sigma$ , meaning that the loss of AdamW does not diverge to infinity even if there is an infinite level of gradient noise: See the right of Figure 6 for an experimental validation. Interestingly, the asymptotic loss value is not influenced by the choice of  $\beta_i$ : We argue that  $\beta_i$  do not impact the asymptotic level of the loss, but rather drive the selection of the basin and speed at which AdamW converges to it — The center-right of Fig. 5 exemplifies this on a simple non-convex landscape. We also extend the verification of our scaling rule to the LLM setting – See Appendix F.8.



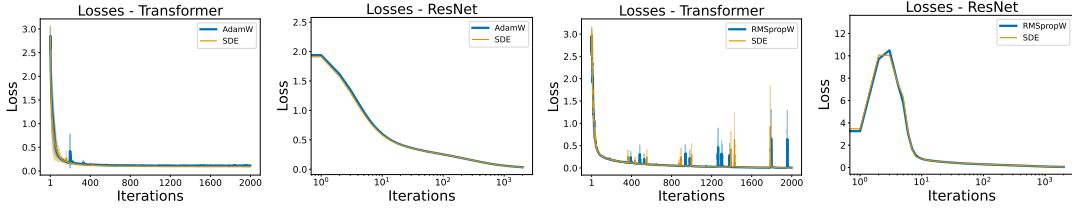


Figure 4: The two images on the left represent the comparison between AdamW and its SDE in terms of  $f(x)$ . The two on the right do the same for RMSpropW. In both cases, the first is a Transformer on MNIST and the second a ResNet on CIFAR-10: Our SDEs match the respective optimizers.

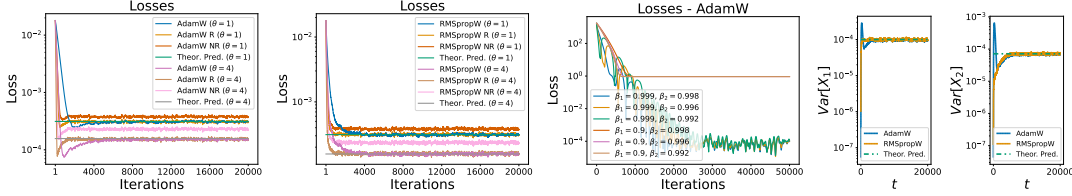


Figure 5: The loss predicted in Lemma 3.13 matches the experimental results on a convex quadratic function. *AdamW* is run with regularization parameter  $\theta = 1$ . *AdamW R* (AdamW Rescaled) is run as we apply the scaling rule with  $\kappa = 2$ . *AdamW NR* (AdamW **Not** Rescaled) is run as we apply the scaling rule with  $\kappa = 2$  on all hyperparameters but  $\theta$ , which is left unchanged: Our scaling rule holds, and failing to rescale  $\theta$  leads the optimizer not to preserve the asymptotic loss level. The same happens for  $\theta = 4$  (Left); The same for RMSpropW (Center-Left); For AdamW,  $\beta_1$  and  $\beta_2$  influence which basin will attract the dynamics and how fast this will converge, but not the asymptotic loss level inside the basin (Center-Right). For both AdamW and RMSpropW, the variance at convergence predicted in Lemma 3.14 matches the experimental results (Right).

Interestingly, the fact that the weight decay is *decoupled* is key to determining the dependency of the asymptotic loss of AdamW w.r.t. the noise level  $\sigma$ . While the asymptotic loss of AdamW is upper-bounded in  $\sigma$ , the same does not hold if we use Adam on the  $L^2$ -regularized loss  $f(x) + \frac{\theta \|x\|_2^2}{2}$ . Under the same assumptions of Lemma 3.13, the dynamics of Adam on  $f(x) + \frac{\theta \|x\|_2^2}{2}$  implies that

$$\mathbb{E}[f(X_t) - f(X_*)] \stackrel{t \rightarrow \infty}{\leq} \frac{\eta \mathcal{L}_\tau \sigma}{2} \frac{L}{2\mu L + \theta(L + \mu)}, \quad (14)$$

meaning that the asymptotic loss level grows linearly in  $\sigma$ : See Figure 14 for empirical validation.

We conclude this section with the stationary distribution of AdamW around a minimum which we empirically validate on the right of Figure 5.

**Lemma 3.14.** *If  $\Sigma(x) = \Sigma$ , the stationary distribution of AdamW is*

$$(\mathbb{E}[X_\infty], Cov[X_\infty]) = \left( 0, \frac{\eta}{2} \left( I_d + \theta H^{-1} \Sigma^{\frac{1}{2}} \right)^{-1} H^{-1} \Sigma^{\frac{1}{2}} \right).$$

**RMSpropW** We derived the analogous results for RMSprop(W) and we reported them in Appendix C.7: importantly, we validate the SDE in Figure 3 for two simple landscapes and in Figure 4 for a Transformer and a ResNet. The results regarding the asymptotic loss level and stationary distributions are validated in the center-left and right of Figure 5 for a convex quadratic function.

**Conclusion:** While for both SignSGD and Adam the asymptotic loss value and the covariance of the iterates scale linearly with  $\Sigma^{\frac{1}{2}}$ , we observe for AdamW this is more intricate: The interaction between curvature, noise, and regularization implies that these two quantities are upper-bounded in  $\Sigma^{\frac{1}{2}}$  and increasing  $\Sigma$  to infinity does not lead to their explosion: *decoupled* weight decay plays a crucial stabilization role at high noise levels near the minimizer — See Figure 6 for a comparison across optimizers. Finally, we argue that  $\beta_i$  play a key role in selecting the basin and the convergence speed to the asymptotic loss value rather than impacting the loss value itself.

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

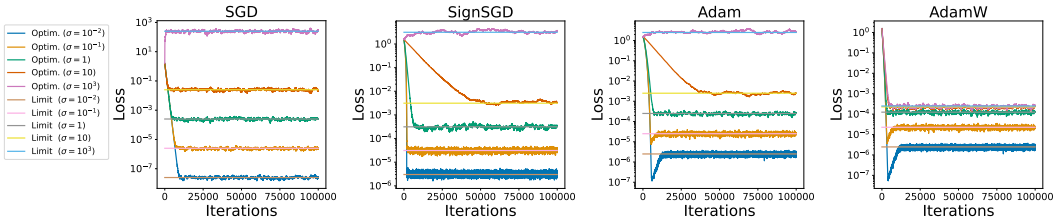


Figure 6: For SGD (Left), SignSGD (Center-Left), Adam (Center-Right), and AdamW: For each optimizer, we plot the loss value on a convex quadratic and compare its asymptotic value with the limits predicted by our theory. As we take  $\Sigma = \sigma^2 I_d$ , we confirm that the loss of SGD scales quadratically in  $\sigma$  (Lemma 3.6), and linearly for SignSGD (Lemma 3.5) and Adam (Lemma 3.13 with  $\theta = 0$ ). For AdamW, the maximum asymptotic loss value is bounded in  $\sigma$  (Lemma 3.13 with  $\theta > 0$ ). In accordance with the experiments, our theory predicts that adaptive methods are more resilient to noise.

#### 4 EXPERIMENTS: SDE VALIDATION

The point of our experiments is to validate the theoretical results derived from the SDEs. Therefore, we first show that our SDEs faithfully represent the dynamics of their respective optimizers. To do so, we integrate the SDEs with Euler-Maruyama (Algorithm 1): This is particularly challenging and expensive as one needs to calculate the full gradients of the DNNs at each iteration.<sup>6</sup> Ours is the first set of validation experiments on various architectures and datasets: An MLP on the Breast Cancer dataset, a CNN and a Transformer on MNIST, and a ResNet on CIFAR-10. Details in Appendix F.

#### 5 CONCLUSION

We derived the first formal SDE for SignSGD, enabling us to demonstrate its dynamics traversing three discernible phases. We characterize how the “signal-to-noise” ratio drives the dynamics of the loss in each of these phases, and we derive the asymptotic value of the loss function, as well as the stationary distribution. Regarding the role of noise, we draw a straightforward comparison with SGD. For SignSGD, the noise level  $\sqrt{\Sigma}$  has an inverse linear effect on the convergence speed of the loss and the iterates. However, it linearly affects the asymptotic expected loss and the asymptotic variance of the iterates. In contrast, for SGD, noise does not influence the convergence speed but has a quadratic impact on the loss level and variance. We also examine the scenario where the noise has infinite variance and demonstrate the resilience of SignSGD, showing that its performance is only marginally affected. Finally, we generalize the analysis to include AdamW and RMSpropW. Specifically, we leverage our novel SDEs to derive the asymptotic value of the loss function, their stationary distribution on a convex quadratic, and a novel scaling rule. The key insight is that, similarly to SignSGD, the loss level and covariance matrix of the iterates of Adam and RMSprop scale linearly in the noise level  $\Sigma^{\frac{1}{2}}$ . For AdamW and RMSpropW, the complex interaction of noise, curvature, and regularization implies that these two quantities are bounded in terms of  $\Sigma^{\frac{1}{2}}$ , showing that *decoupled* weight decay plays a crucial stabilization role at high noise levels near the minimizer. Interestingly, the SDEs for Adam and RMSprop are straightforward corollary of our general results and were derived under much less restrictive and more realistic assumptions than those in the literature. Finally, we thoroughly validate all our theoretical results: We compare the dynamics of the various optimizers with the respective SDEs and find good agreement on simple landscapes and deep neural networks. For Adam and RMSprop, our SDEs track them better than those derived in (Malladi et al., 2022).

**Future work** We believe our results can be extended to other optimizers commonly used in practice such as Signum, AdaGrad, AdaMax, and Nadam. Additionally, inspired by the insights from our SDE analysis, there is potential for designing new optimization algorithms that combine and preserve the strengths of existing methods while mitigating their weaknesses. For example, developing hybrid optimizers that adaptively switch between different strategies based on the training phase or current state of the optimization process could offer superior performance.

<sup>6</sup>Many papers derive SDEs to model optimizers, but most lack validation. Some use toy landscapes, while only Paquette et al. (2021); Compagnoni et al. (2023) validate on simple DNNs. See Appendix A for details.

## REFERENCES

- 540  
541  
542 Ahn, S., Korattikara, A., and Welling, M. (2012). Bayesian posterior sampling via stochastic gradient  
543 fisher scoring. *arXiv preprint arXiv:1206.6380*.
- 544 An, J., Lu, J., and Ying, L. (2020). Stochastic modified equations for the asynchronous stochastic  
545 gradient descent. *Information and Inference: A Journal of the IMA*, 9(4):851–873.
- 546  
547 Ankirchner, S. and Perko, S. (2024). A comparison of continuous-time approximations to stochastic  
548 gradient descent. *Journal of Machine Learning Research*, 25(13):1–55.
- 549 Ayadi, I. and Turinici, G. (2021). Stochastic runge-kutta methods and adaptive sgd-g2 stochastic  
550 gradient descent. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pages  
551 8220–8227. IEEE.
- 552 Balles, L. and Hennig, P. (2018). Dissecting adam: The sign, magnitude and variance of stochastic  
553 gradients. In *International Conference on Machine Learning*, pages 404–413. PMLR.
- 554  
555 Barakat, A. and Bianchi, P. (2021). Convergence and dynamical behavior of the adam algorithm for  
556 nonconvex stochastic optimization. *SIAM Journal on Optimization*, 31(1):244–274.
- 557  
558 Bardi, M. and Kouhkouh, H. (2022). Deep relaxation of controlled stochastic gradient descent via  
559 singular perturbations. *arXiv preprint arXiv:2209.05564*.
- 560 Bercher, A., Gonon, L., Jentzen, A., and Salimova, D. (2020). Weak error analysis for stochastic  
561 gradient descent optimization algorithms. *arXiv preprint arXiv:2007.02723*.
- 562  
563 Bernstein, J., Wang, Y.-X., Azzadenesheli, K., and Anandkumar, A. (2018). signSGD: Compressed  
564 optimisation for non-convex problems. In *Proceedings of the 35th International Conference on*  
565 *Machine Learning*.
- 566 Biderman, S., Schoelkopf, H., Anthony, Q., Bradley, H., O’Brien, K., Hallahan, E., Khan, M. A.,  
567 Purohit, S., Prashanth, U. S., Raff, E., Skowron, A., Sutawika, L., and Van Der Wal, O. (2023).  
568 Pythia: a suite for analyzing large language models across training and scaling. In *Proceedings of*  
569 *the 40th International Conference on Machine Learning, ICML’23*. JMLR.org.
- 570  
571 Bishop, A. N. and Del Moral, P. (2019). Stability properties of systems of linear stochastic differential  
572 equations with random coefficients. *SIAM Journal on Control and Optimization*, 57(2):1023–1042.
- 573  
574 Bock, S. and Weiß, M. G. (2021). Local convergence of adaptive gradient descent optimizers. *arXiv*  
*preprint arXiv:2102.09804*.
- 575  
576 Bradbury, J., Frostig, R., Hawkins, P., Johnson, M. J., Leary, C., Maclaurin, D., Necula, G., Paszke,  
577 A., VanderPlas, J., Wanderman-Milne, S., and Zhang, Q. (2018). JAX: composable transformations  
578 of Python+NumPy programs.
- 579  
580 Chen, P., Lu, J., and Xu, L. (2022). Approximation to stochastic variance reduced gradient langevin dy-  
581 namics by stochastic delay differential equations. *Applied Mathematics & Optimization*, 85(2):15.
- 582  
583 Chen, T., Fox, E., and Guestrin, C. (2014). Stochastic gradient hamiltonian monte carlo. In  
584 *International conference on machine learning*, pages 1683–1691. PMLR.
- 585  
586 Chen, X., Liu, S., Sun, R., and Hong, M. (2019). On the convergence of a class of adam-type  
587 algorithms for non-convex optimization. In *International Conference on Learning Representations*.
- 588  
589 Compagnoni, E. M., Biggio, L., Orvieto, A., Proske, F. N., Kersting, H., and Lucchi, A. (2023). An  
590 sde for modeling sam: Theory and insights. In *International Conference on Machine Learning*,  
591 pages 25209–25253. PMLR.
- 592  
593 Compagnoni, E. M., Orvieto, A., Kersting, H., Proske, F., and Lucchi, A. (2024). Sdes for minimax  
optimization. In *International Conference on Artificial Intelligence and Statistics*, pages 4834–4842.  
PMLR.
- Cui, Z.-X., Fan, Q., and Jia, C. (2020). Momentum methods for stochastic optimization over  
time-varying directed networks. *Signal Processing*, 174:107614.

- 594 Dambrine, M., Dossal, C., Puig, B., and Rondepierre, A. (2024). Stochastic differential equations for  
595 modeling first order optimization methods. *SIAM Journal on Optimization*, 34(2):1402–1426.  
596
- 597 De, S., Mukherjee, A., and Ullah, E. (2018). Convergence guarantees for rmsprop and adam in  
598 non-convex optimization and an empirical comparison to nesterov acceleration. *arXiv preprint*  
599 *arXiv:1807.06766*.
- 600 Défossez, A., Bottou, L., Bach, F., and Usunier, N. (2022). A simple convergence proof of adam and  
601 adagrad. *Transactions on Machine Learning Research*.  
602
- 603 Del Moral, P. and Niclas, A. (2018). A taylor expansion of the square root matrix function. *Journal*  
604 *of Mathematical Analysis and Applications*, 465(1):259–266.
- 605 Deng, L. (2012). The mnist database of handwritten digit images for machine learning research.  
606 *IEEE Signal Processing Magazine*, 29(6):141–142.  
607
- 608 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani,  
609 M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N. (2021). An image is  
610 worth 16x16 words: Transformers for image recognition at scale. In *International Conference on*  
611 *Learning Representations*.
- 612 Dua, D. and Graff, C. (2017). UCI machine learning repository.  
613
- 614 Fontaine, X., De Bortoli, V., and Durmus, A. (2021). Convergence rates and approximation results  
615 for sgd and its continuous-time counterpart. In *Conference on Learning Theory*, pages 1965–2058.  
616 PMLR.
- 617 Ge, R., Huang, F., Jin, C., and Yuan, Y. (2015). Escaping from saddle points—online stochastic  
618 gradient for tensor decomposition. In *Conference on Learning Theory*, pages 797–842.
- 619 Gess, B., Kassing, S., and Konarovskyi, V. (2024). Stochastic modified flows, mean-field limits and  
620 dynamics of stochastic gradient descent. *Journal of Machine Learning Research*, 25(30):1–27.  
621
- 622 Gu, H., Guo, X., and Li, X. (2021). Adversarial training for gradient descent: Analysis through its  
623 continuous-time approximation. *arXiv preprint arXiv:2105.08037*.
- 624 Guo, Z., Xu, Y., Yin, W., Jin, R., and Yang, T. (2021). A novel convergence analysis for algorithms  
625 of the adam family. *arXiv preprint arXiv:2112.03459*.  
626
- 627 Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser,  
628 E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M.,  
629 Haldane, A., del Río, J. F., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K., Reddy,  
630 T., Weckesser, W., Abbasi, H., Gohlke, C., and Oliphant, T. E. (2020). Array programming with  
631 NumPy. *Nature*, 585(7825):357–362.
- 632 Higham, D. J. (2001). An algorithmic introduction to numerical simulation of stochastic differential  
633 equations. *SIAM review*, 43(3):525–546.  
634
- 635 Hong, Y. and Lin, J. (2023). High probability convergence of adam under unbounded gradients and  
636 affine variance noise. *arXiv preprint arXiv:2311.02000*.
- 637 Hu, W., Li, C. J., and Zhou, X. (2019). On the global convergence of continuous-time stochastic  
638 heavy-ball method for nonconvex optimization. In *2019 IEEE International Conference on Big*  
639 *Data (Big Data)*, pages 94–104. IEEE.
- 640 Huang, F., Li, J., and Huang, H. (2021). Super-adam: faster and universal framework of adaptive  
641 gradients. *Advances in Neural Information Processing Systems*, 34:9074–9085.  
642
- 643 Ikeda, N. and Watanabe, S. (2014). *Stochastic differential equations and diffusion processes*. Elsevier.
- 644 Jastrzebski, S., Kenton, Z., Arpit, D., Ballas, N., Fischer, A., Bengio, Y., and Storkey, A. (2018).  
645 Three factors influencing minima in sgd. *ICANN 2018*.  
646
- 647 Jin, C., Ge, R., Netrapalli, P., Kakade, S. M., and Jordan, M. I. (2017). How to escape saddle points  
efficiently. In *International Conference on Machine Learning*, pages 1724–1732. PMLR.

- 648 Karatzas, I. and Shreve, S. (2014). *Brownian motion and stochastic calculus*, volume 113. springer.  
649
- 650 Karimireddy, S. P., Rebjock, Q., Stich, S., and Jaggi, M. (2019a). Error feedback fixes signsgd and  
651 other gradient compression schemes. In *International Conference on Machine Learning*, pages  
652 3252–3261. PMLR.
- 653 Karimireddy, S. P., Rebjock, Q., Stich, S., and Jaggi, M. (2019b). Error feedback fixes SignSGD  
654 and other gradient compression schemes. In *Proceedings of the 36th International Conference on  
655 Machine Learning*.
- 656
- 657 Kingma, D. P. and Ba, J. (2015). Adam: A method for stochastic optimization. In *International  
658 Conference on Learning Representations*.
- 659 Kohatsu-Higa, A., León, J. A., and Nualart, D. (1997). Stochastic differential equations with random  
660 coefficients. *Bernoulli*, pages 233–245.  
661
- 662 Krizhevsky, A., Hinton, G., et al. (2009). Learning multiple layers of features from tiny images.  
663 *Toronto, ON, Canada*.
- 664 Kunin, D., Sagastuy-Brena, J., Gillespie, L., Margalit, E., Tanaka, H., Ganguli, S., and Yamins, D. L.  
665 (2023). The limiting dynamics of sgd: Modified loss, phase-space oscillations, and anomalous  
666 diffusion. *Neural Computation*, 36(1):151–174.  
667
- 668 Kunstner, F., Yadav, R., Milligan, A., Schmidt, M., and Bietti, A. (2024). Heavy-tailed class  
669 imbalance and why adam outperforms gradient descent on language models. *arXiv preprint  
670 arXiv:2402.19449*.
- 671 Lanconelli, A. and Lauria, C. S. (2022). A note on diffusion limits for stochastic gradient descent.  
672 *arXiv preprint arXiv:2210.11257*.  
673
- 674 Levy, K. Y. (2016). The power of normalization: Faster evasion of saddle points. *arXiv preprint  
675 arXiv:1611.04831*.
- 676 Li, H., Rakhlin, A., and Jadbabaie, A. (2023a). Convergence of adam under relaxed assumptions. In  
677 *Thirty-seventh Conference on Neural Information Processing Systems*.
- 678
- 679 Li, L. and Wang, Y. (2022). On uniform-in-time diffusion approximation for stochastic gradient  
680 descent. *arXiv preprint arXiv:2207.04922*.
- 681 Li, Q., Tai, C., and Weinan, E. (2017). Stochastic modified equations and adaptive stochastic gradient  
682 algorithms. In *International Conference on Machine Learning*, pages 2101–2110. PMLR.  
683
- 684 Li, Q., Tai, C., and Weinan, E. (2019). Stochastic modified equations and dynamics of stochastic  
685 gradient algorithms i: Mathematical foundations. *The Journal of Machine Learning Research*,  
686 20(1):1474–1520.
- 687 Li, Z., Malladi, S., and Arora, S. (2021). On the validity of modeling SGD with stochastic differential  
688 equations (SDEs). In Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W., editors,  
689 *Advances in Neural Information Processing Systems*.
- 690
- 691 Li, Z., Wang, Y., and Wang, Z. (2023b). Fast equilibrium of sgd in generic situations. In *The Twelfth  
692 International Conference on Learning Representations*.
- 693
- 694 Liu, T., Chen, Z., Zhou, E., and Zhao, T. (2021). A diffusion approximation theory of momentum  
695 stochastic gradient descent in nonconvex optimization. *Stochastic Systems*.
- 696
- 697 Loshchilov, I. and Hutter, F. (2019). Decoupled weight decay regularization. In *International  
698 Conference on Learning Representations*.
- 699
- 700 Luo, L., Xiong, Y., and Liu, Y. (2019). Adaptive gradient methods with dynamic bound of learning  
701 rate. In *International Conference on Learning Representations*.
- 702
- 703 Ma, C., Wu, L., and Weinan, E. (2022). A qualitative study of the dynamic behavior for adaptive  
704 gradient algorithms. In *Mathematical and Scientific Machine Learning*, pages 671–692. PMLR.

- 702 Mai, V. V. and Johansson, M. (2021). Stability and convergence of stochastic gradient clipping:  
703 Beyond lipschitz continuity and smoothness. In *International Conference on Machine Learning*.  
704
- 705 Malladi, S., Lyu, K., Panigrahi, A., and Arora, S. (2022). On the SDEs and scaling rules for adaptive  
706 gradient algorithms. In *Advances in Neural Information Processing Systems*.  
707
- 708 Mandt, S., Hoffman, M., and Blei, D. (2016). A variational analysis of stochastic gradient algorithms.  
709 In *International conference on machine learning*, pages 354–363. PMLR.
- 710 Mandt, S., Hoffman, M. D., and Blei, D. M. (2017). Stochastic gradient descent as approximate  
711 bayesian inference. *JMLR 2017*.  
712
- 713 Mao, X. (2007). *Stochastic differential equations and applications*. Elsevier.
- 714 Maulen-Soto, R., Fadili, J., Attouch, H., and Ochs, P. (2024). Stochastic inertial dynamics via time  
715 scaling and averaging. *arXiv preprint arXiv:2403.16775*.  
716
- 717 Maulén Soto, R. I. (2021). A continuous-time model of stochastic gradient descent: convergence  
718 rates and complexities under lojasiewicz inequality. *Universidad de Chile*.
- 719 Mertikopoulos, P. and Staudigl, M. (2018). On the convergence of gradient-like flows with noisy  
720 gradient input. *SIAM Journal on Optimization*, 28(1):163–197.  
721
- 722 Milstein, G. N. (2013). *Numerical integration of stochastic differential equations*, volume 313.  
723 Springer Science & Business Media.
- 724 Mil’shtein, G. (1986). Weak approximation of solutions of systems of stochastic differential equations.  
725 *Theory of Probability & Its Applications*, 30(4):750–766.  
726
- 727 Noci, L., Anagnostidis, S., Biggio, L., Orvieto, A., Singh, S. P., and Lucchi, A. (2022). Signal  
728 propagation in transformers: Theoretical perspectives and the role of rank collapse. *Advances in*  
729 *Neural Information Processing Systems*, 35:27198–27211.
- 730 Øksendal, B. (1990). When is a stochastic integral a time change of a diffusion? *Journal of theoretical*  
731 *probability*, 3(2):207–226.  
732
- 733 Orvieto, A. and Lucchi, A. (2019). Continuous-time models for stochastic optimization algorithms.  
734 *Advances in Neural Information Processing Systems*, 32.  
735
- 736 Pan, Y. and Li, Y. (2022). Toward understanding why adam converges faster than SGD for transform-  
737 ers. In *OPT 2022: Optimization for Machine Learning (NeurIPS 2022 Workshop)*.
- 738 Paquette, C., Lee, K., Pedregosa, F., and Paquette, E. (2021). Sgd in the large: Average-case analysis,  
739 asymptotics, and stepsize criticality. In *Conference on Learning Theory*, pages 3548–3626. PMLR.  
740
- 741 Paquette, E., Paquette, C., Xiao, L., and Pennington, J. (2024). 4+ 3 phases of compute-optimal  
742 neural scaling laws. *arXiv preprint arXiv:2405.15074*.
- 743 Pascanu, R., Mikolov, T., and Bengio, Y. (2013). On the difficulty of training recurrent neural  
744 networks. In *International conference on machine learning*.  
745
- 746 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,  
747 Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher,  
748 M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of*  
749 *Machine Learning Research*, 12:2825–2830.
- 750 Poggio, T., Kawaguchi, K., Liao, Q., Miranda, B., Rosasco, L., Boix, X., Hidary, J., and Mhaskar,  
751 H. (2017). Theory of deep learning iii: explaining the non-overfitting puzzle. *arXiv preprint*  
752 *arXiv:1801.00173*.  
753
- 754 Puchkin, N., Gorbunov, E., Kutuzov, N., and Gasnikov, A. (2024). Breaking the heavy-tailed noise  
755 barrier in stochastic optimization problems. In *International Conference on Artificial Intelligence*  
*and Statistics*.

- 756 Raginsky, M. and Bouvrie, J. (2012). Continuous-time stochastic mirror descent on a network:  
757 Variance reduction, consensus, convergence. In *2012 IEEE 51st IEEE Conference on Decision and*  
758 *Control (CDC)*, pages 6793–6800. IEEE.
- 759 Reddi, S. J., Kale, S., and Kumar, S. (2018). On the convergence of adam and beyond. In *International*  
760 *Conference on Learning Representations*.
- 761 Safaryan, M. and Richtarik, P. (2021). Stochastic sign descent methods: New algorithms and better  
762 theory. In *Proceedings of the 38th International Conference on Machine Learning*.
- 763 Shi, N. and Li, D. (2021). Rmsprop converges with proper hyperparameter. In *International*  
764 *conference on learning representation*.
- 765 Smith, S. L., Dherin, B., Barrett, D. G. T., and De, S. (2021). On the origin of implicit regularization  
766 in stochastic gradient descent. *ArXiv*, abs/2101.12176.
- 767 Soboleva, D., Al-Khateeb, F., Myers, R., Steeves, J. R., Hestness, J., and Dey, N. (2023). SlimPajama:  
768 A 627B token cleaned and deduplicated version of RedPajama.
- 769 Soto, R. M., Fadili, J., and Attouch, H. (2022). An sde perspective on stochastic convex optimization.  
770 *arXiv preprint arXiv:2207.02750*.
- 771 Stephan, M., Hoffman, M. D., Blei, D. M., et al. (2017). Stochastic gradient descent as approximate  
772 bayesian inference. *Journal of Machine Learning Research*, 18(134):1–35.
- 773 Su, L. and Lau, V. K. (2023). Accelerated federated learning over wireless fading channels with  
774 adaptive stochastic momentum. *IEEE Internet of Things Journal*.
- 775 Sun, J., Yang, Y., Xun, G., and Zhang, A. (2023). Scheduling hyperparameters to improve generaliza-  
776 tion: From centralized sgd to asynchronous sgd. *ACM Transactions on Knowledge Discovery from*  
777 *Data*, 17(2):1–37.
- 778 Tieleman, T. and Hinton, G. (2012). Lecture 6.5-rmsprop: Divide the gradient by a running average  
779 of its recent magnitude.
- 780 Van Rossum, G. and Drake, F. L. (2009). *Python 3 Reference Manual*. CreateSpace, Scotts Valley,  
781 CA.
- 782 Wang, B., Fu, J., Zhang, H., Zheng, N., and Chen, W. (2024). Closing the gap between the upper  
783 bound and lower bound of adam’s iteration complexity. *Advances in Neural Information Processing*  
784 *Systems*, 36.
- 785 Wang, B., Zhang, Y., Zhang, H., Meng, Q., Ma, Z.-M., Liu, T.-Y., and Chen, W. (2022). Provable  
786 adaptivity in adam. *arXiv preprint arXiv:2208.09900*.
- 787 Wang, Y. and Wu, S. (2020). Asymptotic analysis via stochastic differential equations of gradient  
788 descent algorithms in statistical and computational paradigms. *Journal of machine learning*  
789 *research*, 21(199):1–103.
- 790 Wang, Z. and Mao, Y. (2022). Two facets of sde under an information-theoretic lens: Generalization  
791 of sgd via training trajectories and via terminal states. *arXiv preprint arXiv:2211.10691*.
- 792 Wojtowytsch, S. (2024). Stochastic gradient descent with noise of machine learning type part ii:  
793 Continuous time analysis. *Journal of Nonlinear Science*, 34(1):16.
- 794 Wolkowicz, H. and Styan, G. P. (1980). Bounds for eigenvalues using traces. *Linear algebra and its*  
795 *applications*, 29:471–506.
- 796 Wu, J., Hu, W., Xiong, H., Huan, J., Braverman, V., and Zhu, Z. (2020). On the noisy gradient descent  
797 that generalizes as sgd. In *International Conference on Machine Learning*, pages 10367–10376.  
798 PMLR.
- 799 Wu, L., Wang, M., and Su, W. (2022). The alignment property of sgd noise and how it helps select flat  
800 minima: A stability analysis. *Advances in Neural Information Processing Systems*, 35:4680–4693.

- 810 Xie, Z., Wang, X., Zhang, H., Sato, I., and Sugiyama, M. (2022). Adaptive inertia: Disentangling the  
811 effects of adaptive learning rate and momentum. In *International conference on machine learning*,  
812 pages 24430–24459. PMLR.
- 813
- 814 Xie, Z., Yuan, L., Zhu, Z., and Sugiyama, M. (2021). Positive-negative momentum: Manipulating  
815 stochastic gradient noise to improve generalization. In *Proceedings of the 38th International  
816 Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*,  
817 pages 11448–11458. PMLR.
- 818
- 819 Yang, J., Li, X., Fatkhullin, I., and He, N. (2024). Two sides of one coin: the limits of untuned sgd  
820 and the power of adaptive methods. *Advances in Neural Information Processing Systems*, 36.
- 821
- 822 Zaheer, M., Reddi, S., Sachan, D., Kale, S., and Kumar, S. (2018). Adaptive methods for nonconvex  
823 optimization. *Advances in neural information processing systems*, 31.
- 824
- 825 Zhang, J., He, T., Sra, S., and Jadbabaie, A. (2020a). Why gradient clipping accelerates training: A  
826 theoretical justification for adaptivity. In *International Conference on Learning Representations*.
- 827
- 828 Zhang, J., Karimireddy, S. P., Veit, A., Kim, S., Reddi, S., Kumar, S., and Sra, S. (2020b). Why are  
829 adaptive methods good for attention models? *Advances in Neural Information Processing Systems*.
- 830
- 831 Zhang, Y., Chen, C., Shi, N., Sun, R., and Luo, Z.-Q. (2022). Adam can converge without any  
832 modification on update rules. *Advances in neural information processing systems*.
- 833
- 834 Zhang, Z., Li, Y., Luo, T., and Xu, Z.-Q. J. (2023). Stochastic modified equations and dynamics of  
835 dropout algorithm. *arXiv preprint arXiv:2305.15850*.
- 836
- 837 Zhao, J., Lucchi, A., Prosk, F. N., Orvieto, A., and Kersting, H. (2022). Batch size selection by  
838 stochastic optimal control. In *Has it Trained Yet? NeurIPS 2022 Workshop*.
- 839
- 840 Zhou, P., Feng, J., Ma, C., Xiong, C., Hoi, S. C. H., et al. (2020a). Towards theoretically understanding  
841 why sgd generalizes better than adam in deep learning. *Advances in Neural Information Processing  
842 Systems*, 33:21285–21296.
- 843
- 844 Zhou, P., Xie, X., Lin, Z., and Yan, S. (2024). Towards understanding convergence and generalization  
845 of adamw. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- 846
- 847 Zhou, P., Xie, X., and Shuicheng, Y. (2022). Win: Weight-decay-integrated nesterov accelera-  
848 tion for adaptive gradient algorithms. In *The Eleventh International Conference on Learning  
849 Representations*.
- 850
- 851 Zhou, X., Yuan, H., Li, C. J., and Sun, Q. (2020b). Stochastic modified equations for continuous  
852 limit of stochastic admm. *arXiv preprint arXiv:2003.03532*.
- 853
- 854 Zhu, Y. and Ying, L. (2021). A sharp convergence rate for a model equation of the asynchronous  
855 stochastic gradient descent. *Communications in Mathematical Sciences*.
- 856
- 857 Zhu, Z., Wu, J., Yu, B., Wu, L., and Ma, J. (2019a). The anisotropic noise in stochastic gradient  
858 descent: Its behavior of escaping from sharp minima and regularization effects. *ICML*.
- 859
- 860 Zhu, Z., Wu, J., Yu, B., Wu, L., and Ma, J. (2019b). The anisotropic noise in stochastic gradient  
861 descent: Its behavior of escaping from sharp minima and regularization effects. *ICML 2019*.
- 862
- 863 Ziyin, L., Liu, K., Mori, T., and Ueda, M. (2021). Strength of minibatch noise in sgd. *arXiv preprint  
arXiv:2102.05375*.



## A ADDITIONAL RELATED WORKS

In this section, we list some papers that derived or used SDEs to model optimizers. In particular, we focus on the aspect of empirically verifying the validity of such SDEs in the sense that they indeed track the respective optimizers. We divide these into three categories: Those that did not carry out any type of validation, those that did it on simple landscapes (quadratic functions et similia), and those that did small experiments on neural networks.

None of the following papers carried out any experimental validation of the approximating power of the SDEs they derived. Many of them did not even validate the insights derived from the SDEs: (Liu et al., 2021; Hu et al., 2019; Bercher et al., 2020; Zhu and Ying, 2021; Cui et al., 2020; Maulén Soto, 2021; Wang and Wu, 2020; Lanconelli and Lauria, 2022; Ayadi and Turinici, 2021; Soto et al., 2022; Li and Wang, 2022; Wang and Mao, 2022; Bardi and Kouhkouh, 2022; Chen et al., 2022; Kunin et al., 2023; Zhang et al., 2023; Sun et al., 2023; Li et al., 2023b; Gess et al., 2024; Dambrine et al., 2024; Maulen-Soto et al., 2024).

The following ones carried out validation experiments on artificial landscapes, e.g. quadratic or quartic function, or easy regression tasks: (Li et al., 2017; 2019; Zhou et al., 2020b; An et al., 2020; Fontaine et al., 2021; Gu et al., 2021; Su and Lau, 2023; Ankirchner and Perko, 2024).

The following papers carried out some experiments which include neural networks: (Paquette et al., 2021; Compagnoni et al., 2023). In particular, they both simulate the SDEs with a numerical integrator and compare them with the respective optimizers: The first validates the SDE on a shallow MLP while the second does so on a shallow and a deep MLP. Regarding (Li et al., 2021; Malladi et al., 2022), they do not validate their SDEs: Rather, their approach conceptually proceeds as follows:

1. Derive an SDE for an optimizer which we now dub “*Optimizer A*”;
2. Notice that simulating the SDE is too expensive;
3. Define another discrete-time algorithm called SVAG which also has the same SDE as “*Optimizer A*”. Importantly, SVAG does not numerically integrate the SDE as it does not even require access to it: It does not need access neither to the drift nor to the diffusion term;
4. Simulate SVAG and show that it tracks “*Optimizer A*” successfully;
5. Conclude that the SDE is a good approximation for “*Optimizer A*”.

However, they never validated that the SDE is a good approximation for “*Optimizer A*” or for SVAG either. With the same logic, they could have done the following:

1. Derive an SDE for “*Optimizer A*”;
2. Notice that simulating the SDE is too expensive;
3. Define another discrete-time algorithm called “*Optimizer B*” which coincides with “*Optimizer A*” and thus of course shares the same SDE;
4. Simulate “*Optimizer B*” and show that it tracks “*Optimizer A*” perfectly, as they are the same optimizer by definition;
5. Conclude that the SDE is a good approximation for “*Optimizer A*”.

In particular, the only fact they prove is that SVAG is a discrete-time optimizer that shares the same SDE as “*Optimizer A*” because it describes a discrete trajectory that is a 1st-order approximation of the SDE of “*Optimizer A*”, but NOT that the SDE they derived is THE SDE. One cannot conclude that the SDE derived for “*Optimizer A*” is a good model for “*Optimizer A*” by simply comparing two algorithms “*Optimizer A*” and “*Optimizer B*” that share the same SDE. Otherwise, simply comparing an optimizer “*Optimizer A*” with itself would do the trick. An SDE’s empirical validation can only occur if the SDE is simulated with a numerical integrator that requires access to the drift and diffusion terms (Higham, 2001; Milstein, 2013).

## B STOCHASTIC CALCULUS

In this section, we summarize some important results in the analysis of Stochastic Differential Equations Mao (2007); Øksendal (1990). The notation and the results in this section will be used

extensively in all proofs in this paper. We assume the reader to have some familiarity with Brownian motion and with the definition of stochastic integral (Ch. 1.4 and 1.5 in Mao (2007)).

### B.1 ITÔ'S LEMMA

We start with some notation: Let  $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$  be a filtered probability space. We say that an event  $E \in \mathcal{F}$  holds almost surely (a.s.) in this space if  $\mathbb{P}(E) = 1$ . We call  $\mathcal{L}^p([a, b], \mathbb{R}^d)$ , with  $p > 0$ , the family of  $\mathbb{R}^d$ -valued  $\mathcal{F}_t$ -adapted processes  $\{f_t\}_{a \leq t \leq b}$  such that

$$\int_a^b \|f_t\|^p dt \leq \infty.$$

Moreover, we denote by  $\mathcal{M}^p([a, b], \mathbb{R}^d)$ , with  $p > 0$ , the family of  $\mathbb{R}^d$ -valued processes  $\{f_t\}_{a \leq t \leq b}$  in  $\mathcal{L}([a, b], \mathbb{R}^d)$  such that  $\mathbb{E} \left[ \int_a^b \|f_t\|^p dt \right] \leq \infty$ . We will write  $h \in \mathcal{L}^p(\mathbb{R}_+, \mathbb{R}^d)$ , with  $p > 0$ , if  $h \in \mathcal{L}^p([0, T], \mathbb{R}^d)$  for every  $T > 0$ . Similar definitions hold for matrix-valued functions using the Frobenius norm  $\|A\| := \sqrt{\sum_{ij} |A_{ij}|^2}$ .

Let  $W = \{W_t\}_{t \geq 0}$  be a one-dimensional Brownian motion defined on our probability space and let  $X = \{X_t\}_{t \geq 0}$  be an  $\mathcal{F}_t$ -adapted process taking values on  $\mathbb{R}^d$ .

**Definition B.1.** Let the drift be  $b \in \mathcal{L}^1(\mathbb{R}_+, \mathbb{R}^d)$  and the diffusion term be  $\sigma \in \mathcal{L}^2(\mathbb{R}_+, \mathbb{R}^{d \times m})$ .  $X_t$  is an Itô process if it takes the form

$$X_t = x_0 + \int_0^t b_s ds + \int_0^t \sigma_s dW_s.$$

We shall say that  $X_t$  has the stochastic differential

$$dX_t = b_t dt + \sigma_t dW_t. \quad (15)$$

**Theorem B.2** (Itô's Lemma). *Let  $X_t$  be an Itô process with stochastic differential  $dX_t = b_t dt + \sigma_t dW_t$ . Let  $f(x, t)$  be twice continuously differentiable in  $x$  and continuously differentiable in  $t$ , taking values in  $\mathbb{R}$ . Then  $f(X_t, t)$  is again an Itô process with stochastic differential*

$$df(X_t, t) = \partial_t f(X_t, t) dt + \langle \nabla f(X_t, t), b_t \rangle dt + \frac{1}{2} \text{Tr}(\sigma_t \sigma_t^\top \nabla^2 f(X_t, t)) dt + \langle \nabla f(X_t, t), \sigma_t \rangle dW_t. \quad (16)$$

### B.2 STOCHASTIC DIFFERENTIAL EQUATIONS

Stochastic Differential Equations (SDEs) are equations of the form

$$dX_t = b(X_t, t) dt + \sigma(X_t, t) dW_t.$$

First of all, we need to define what it means for a stochastic process  $X = \{X_t\}_{t \geq 0}$  with values in  $\mathbb{R}^d$  to solve an SDE.

**Definition B.3.** Let  $X_t$  be as above with deterministic initial condition  $X_0 = x_0$ . Assume  $b : \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}^d$  and  $\sigma : \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}^{d \times m}$  are Borel measurable;  $X_t$  is called a solution to the corresponding SDE if

1.  $X_t$  is continuous and  $\mathcal{F}_t$ -adapted;
2.  $b \in \mathcal{L}^1([0, T], \mathbb{R}^d)$ ;
3.  $\sigma \in \mathcal{L}^2([0, T], \mathbb{R}^{d \times m})$ ;
4. For every  $t \in [0, T]$

$$X_t = x_0 + \int_0^t b(X_s, s) ds + \int_0^t \sigma(X_s, s) dW(s) \quad a.s.$$

Moreover, the solution  $X_t$  is said to be unique if any other solution  $X_t^*$  is such that

$$\mathbb{P}\{X_t = X_t^*, \text{ for all } 0 \leq t \leq T\} = 1.$$

Notice that since the solution to an SDE is an Itô process, we can use Itô's lemma. The following theorem gives a sufficient condition on  $b$  and  $\sigma$  for the existence of a solution to the corresponding SDE.

**Theorem B.4.** *Assume that there exist two positive constants  $\bar{K}$  and  $K$  such that*

1. (Global Lipschitz condition) for all  $x, y \in \mathbb{R}^d$  and  $t \in [0, T]$

$$\max\{\|b(x, t) - b(y, t)\|^2, \|\sigma(x, t) - \sigma(y, t)\|^2\} \leq \bar{K}\|x - y\|^2;$$

2. (Linear growth condition) for all  $x \in \mathbb{R}^d$  and  $t \in [0, T]$

$$\max\{\|b(x, t)\|^2, \|\sigma(x, t)\|^2\} \leq K(1 + \|x\|^2).$$

Then, there exists a unique solution  $X_t$  to the corresponding SDE, and  $X_t \in \mathcal{M}^2([0, T], \mathbb{R}^d)$ .

**Numerical approximation.** Often, SDEs are solved numerically. The simplest algorithm to provide a sample path  $(\hat{x}_k)_{k \geq 0}$  for  $X_t$ , so that  $X_{k\Delta t} \approx \hat{x}_k$  for some small  $\Delta t$  and for all  $k\Delta t \leq M$  is called Euler-Maruyama (Algorithm 1). For more details on this integration method and its approximation properties, the reader can check Mao (2007).

---

#### Algorithm 1 Euler-Maruyama Integration Method for SDEs

---

**input** The drift  $b$ , the volatility  $\sigma$ , and the initial condition  $x_0$ .

Fix a stepsize  $\Delta t$ ;

Initialize  $\hat{x}_0 = x_0$ ;

$k = 0$ ;

**while**  $k \leq \lfloor \frac{T}{\Delta t} \rfloor$  **do**

    Sample some  $d$ -dimensional Gaussian noise  $Z_k \sim \mathcal{N}(0, I_d)$ ;

    Compute  $\hat{x}_{k+1} = \hat{x}_k + \Delta t b(\hat{x}_k, k\Delta t) + \sqrt{\Delta t} \sigma(\hat{x}_k, k\Delta t) Z_k$ ;

$k = k + 1$ ;

**end while**

**output** The approximated sample path  $(\hat{x}_k)_{0 \leq k \leq \lfloor \frac{T}{\Delta t} \rfloor}$ .

---

## C THEORETICAL FRAMEWORK - WEAK APPROXIMATION

In this section, we introduce the theoretical framework used in the paper, together with its assumptions and notations.

First of all, many proofs will use Taylor expansions in powers of  $\eta$ . For ease of notation, we introduce the shorthand that whenever we write  $\mathcal{O}(\eta^\alpha)$ , we mean that there exists a function  $K(x) \in G$  such that the error terms are bounded by  $K(x)\eta^\alpha$ . For example, we write

$$b(x + \eta) = b_0(x) + \eta b_1(x) + \mathcal{O}(\eta^2)$$

to mean: there exists  $K \in G$  such that

$$|b(x + \eta) - b_0(x) - \eta b_1(x)| \leq K(x)\eta^2.$$

Additionally, we introduce the following shorthand:

- A multi-index is  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n)$  such that  $\alpha_j \in \{0, 1, 2, \dots\}$ ;
- $|\alpha| := \alpha_1 + \alpha_2 + \dots + \alpha_n$ ;
- $\alpha! := \alpha_1! \alpha_2! \dots \alpha_n!$ ;
- For  $x = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$ , we define  $x^\alpha := x_1^{\alpha_1} x_2^{\alpha_2} \dots x_n^{\alpha_n}$ ;

- For a multi-index  $\beta$ ,  $\partial_\beta^{|\beta|} f(x) := \frac{\partial^{|\beta|}}{\partial x_1^{\beta_1} \partial x_2^{\beta_2} \dots \partial x_n^{\beta_n}} f(x)$ ;
- We also denote the partial derivative with respect to  $x_i$  by  $\partial_{e_i}$ .

**Definition C.1** (*G Set*). Let  $G$  denote the set of continuous functions  $\mathbb{R}^d \rightarrow \mathbb{R}$  of at most polynomial growth, i.e.  $g \in G$  if there exists positive integers  $\nu_1, \nu_2 > 0$  such that  $|g(x)| \leq \nu_1 (1 + |x|^{2\nu_2})$ , for all  $x \in \mathbb{R}^d$ .

**Definition C.2** ( $\mathcal{C}_b^k(\mathbb{R}^n, \mathbb{R})$ ).  $\mathcal{C}_b^k(\mathbb{R}^n, \mathbb{R})$  denotes the space of functions whose  $k$ -th derivatives are bounded.

The next results are inspired by Theorem 1 of Li et al. (2017) and are derived under some regularity assumption on the function  $f$ .

### C.1 ASSUMPTIONS

In general, we assume some regularity in the loss function.

**Assumption C.3.** Assume that the following conditions on  $f, f_i \in \mathcal{C}_b^8(\mathbb{R}^n, \mathbb{R})$ , and their gradients are satisfied:

- $\nabla f, \nabla f_i$  satisfy a Lipschitz condition: there exists  $L > 0$  such that

$$|\nabla f(u) - \nabla f(v)| + \sum_{i=1}^n |\nabla f_i(u) - \nabla f_i(v)| \leq L|u - v|;$$

- $f, f_i$  and its partial derivatives up to order 7 belong to  $G$ ;
- $\nabla f, \nabla f_i$  satisfy a growth condition: there exists  $M > 0$  such that

$$|\nabla f(x)| + \sum_{i=1}^n |\nabla f_i(x)| \leq M(1 + |x|).$$

Regarding the gradient noise, each optimizer has its mild assumptions which are weaker or in line with the literature.

### SignSGD

1. The gradient noise  $Z(x)$  admits a strictly positive density function  $g_x$  for all  $x$  and require that  $g : \mathbb{R}^n \times \mathbb{R}^n \rightarrow [0, \infty)$  s.t.  $(x, y) \mapsto g_x(y)$  is in  $C^8(\mathbb{R}^n \times \mathbb{R}^n)$  such that all partial derivatives of  $g$  up to order 8 are integrable with respect to  $y$  and s.t. their  $L^1$ -norms are uniformly bounded in  $x$ . This assumption covers Gaussian and Student's t, thus being *more general than the literature*. Indeed, the Gaussianity of the noise is commonly assumed: Among others, see Ahn et al. (2012); Chen et al. (2014); Mandt et al. (2016); Stephan et al. (2017); Zhu et al. (2019a); Wu et al. (2020); Xie et al. (2021), while Jastrzebski et al. (2018) offers an intuitive justification as well;

2. For all compact sets  $K$

$$\sup_{x \in K} |g(x, \cdot)| \in L^1(\mathbb{R}^n),$$

which of course covers the Gaussian case, *thus being more general than the literature*.

3. The functions in Eq. 18 to be in  $G$ , which, as we show below, covers Gaussian and Student's t, *thus being more general than the literature*.

### Adam(W) and RMSprop(W)

1. In line with Malladi et al. (2022), we assume that  $\sqrt{\Sigma}(x)$  is: In  $G$  together with its derivatives, Lipschitz, bounded, and satisfy Affine Growth;

2. The term  $(\nabla f(x))^2$  to be Lipschitz and of affine growth, which is a consequence of assuming bounded gradients as often done in the literature on the convergence of RMSprop and Adam: Among many, see (Luo et al., 2019; Défossez et al., 2022; Guo et al., 2021; Huang et al., 2021) together with the discussion in Section 2.1 of Shi and Li (2021).

**Remark** All the assumptions above are *in line with or more general than those commonly found in the literature*. In line with *Remark 11* of the seminal paper Li et al. (2019), we observe that while some of these assumptions might seem strong, loss functions in applications have inward pointing gradients for sufficiently large  $x$ . Therefore, we could simply modify the loss to satisfy the assumptions above.

Regarding the drift and diffusion coefficients, we highlight that many papers in the literature following this framework do not check for their regularity before applying the approximation theorems Hu et al. (2019); An et al. (2020); Zhu and Ying (2021); Cui et al. (2020); Maulén Soto (2021); Wang and Mao (2022); Compagnoni et al. (2023; 2024); Li et al. (2017). At first sight, it would seem that not even the seminal paper Li et al. (2019) checks these conditions carefully. However, a deeper investigation shows that they are restricting their analysis to compact sets to leverage the regularity and convergence properties of mollifiers: The assumption regarding the compactness of the domain is not highlighted nor assumed in any part of the paper. Therefore, we conclude that, willingly or not, most papers implicitly make these assumptions.

## C.2 TECHNICAL RESULTS

In this subsection, we provide some results that will be instrumental in the derivation of the SDEs.

**Lemma C.4.** *Assume the existence of a probability density  $g_x$  of the gradient noise  $Z(x)$  for all  $x$  and require that  $g : \mathbb{R}^n \times \mathbb{R}^n \rightarrow [0, \infty)$ ;  $(x, y) \mapsto g_x(y)$  is in  $C^8(\mathbb{R}^n \times \mathbb{R}^n)$  such that all partial derivatives of  $g$  up to order 8 are integrable with respect to  $y$  and such that their  $L^1$ -norms are uniformly bounded in  $x$ . Further, let  $f \in C^8(\mathbb{R}^n)$  and  $h : \mathbb{R}^n \rightarrow \mathbb{R}$  be a bounded Borel measurable function. Define the function  $k$  by*

$$k(x) = \mathbb{E} [h(\nabla f_\gamma(x))].$$

*Then there exists a version  $\widehat{k}$  of  $k$  with  $\widehat{k} \in C_b^7(\mathbb{R}^n)$ .*

*Proof.* Let  $\varphi$  be smooth and compactly supported. Then for all multi indices  $\beta$  with  $|\beta| \leq 8$ , substitution, Fubini's theorem, and integration by parts imply that

$$\begin{aligned} \int_{\mathbb{R}^n} k(x) \partial_\beta^{|\beta|} \varphi(x) dx &= \int_{\mathbb{R}^n} \mathbb{E} [h(\nabla f_\gamma(x))] \partial_\beta^{|\beta|} \varphi(x) dx \\ &= \int_{\mathbb{R}^n} \int_{\mathbb{R}^n} h(y) g_x(y - \nabla f(x)) dy \partial_\beta^{|\beta|} \varphi(x) dx \\ &= (-1)^{|\beta|} \int_{\mathbb{R}^n} \int_{\mathbb{R}^n} h(y) \partial_\beta^{|\beta|} (g_x(y - \nabla f(x))) dy \varphi(x) dx. \end{aligned}$$

So

$$\int_{\mathbb{R}^n} h(y) \partial_\beta^{|\beta|} (g_x(y - \nabla f(x))) dy$$

is a weak derivative  $\partial_\beta^{|\beta|} k$  of  $k$  on any bounded open set. For compact sets  $K$  we obtain that

$$\begin{aligned} &\int_K \left| \int_{\mathbb{R}^n} h(y) \partial_\beta^{|\beta|} (g_x(y - \nabla f(x))) dy \right|^p dx \\ &\leq \|h\|_\infty^p \lambda^n(K) \left( \sup_{x \in \mathbb{R}^n} \int_{\mathbb{R}^n} \left| \partial_\beta^{|\beta|} (g_x(y - \nabla f(x))) \right| dy \right)^p < \infty \end{aligned}$$

for all  $p \geq 2$  because of our assumptions on  $g$  and  $f$  and substitution ( $\lambda^n$  Lebesgue measure). So it follows from Sobolev embeddings with respect to Hölder spaces that for all bounded and open sets  $\Omega$  there exists a version  $\widehat{k}$  of  $k$  such that  $\widehat{k} \in C^7(\Omega)$ . The latter version can be extended to  $\Omega = \mathbb{R}^n$ , which we also denote by  $\widehat{k}$ . Since  $\partial_\beta^{|\beta|} k$  is bounded for  $|\beta| \leq 8$ , we conclude that  $\widehat{k} \in C_b^7(\mathbb{R}^n)$ .  $\square$

1134 **Lemma C.5.** Assuming that for all compact sets  $K$

$$1135 \sup_{x \in K} |g(x, \cdot)| \in L^1(\mathbb{R}^n),$$

1136 and the positivity of the density functions, we have that for  $m = 1, \dots, 7$  that

$$1137 \left\| \partial_{j_1} \dots \partial_{j_m} A^{1/2}(x) \right\| \leq C l_m(x), \quad (17)$$

1138 where the function  $l_m(x)$  is defined as

$$1139 \begin{aligned} 1140 l_m(x) &:= \sum_{r=0}^{m-1} \left( \frac{1}{m(x) + s(x)(n-1)^{1/2}} \left( 1 + \frac{2s(x)(n-1)^{1/2}}{m(x) - s(x)(n-1)^{-1/2}} \right) \right)^{-(r+1/2)} \\ 1141 &\times \max_{|\beta| \leq m} \left\| \partial_\beta^{|\beta|} A(x) \right\|^{r+1}. \end{aligned} \quad (18)$$

1142 *Proof.* To prove this, we need the fact that the Fréchet derivatives of the square root function  $\varphi$  can be represented as follows (see Theorem 1.1 in Del Moral and Niclas (2018)):

$$1143 \nabla \varphi(A)[H] = \int_0^\infty e^{-t\varphi(A)} H e^{-t\varphi(A)} dt,$$

1144 and higher derivatives of order  $m \geq 2$  are given by

$$1145 \begin{aligned} 1146 \nabla^m \varphi(A)[H, \dots, H] &= -\nabla \varphi(A) \left[ \sum_{p+q=m-2} \frac{m!}{(p+1)!(q+1)!} (\nabla^{p+1} \varphi(A)[H, \dots, H]) \right. \\ 1147 &\quad \left. \times (\nabla^{q+1} \varphi(A)[H, \dots, H]) \right] \end{aligned} \quad (19)$$

1148 for all  $A \in \mathbb{S}$  and symmetric  $n \times n$  matrices  $H$ . Moreover, we have the following estimate for  $m \geq 0$ :

$$1149 \left\| \nabla^{m+1} \varphi(A) \right\| \leq (\sqrt{n})^m (m+1)! C_m 2^{-2(m+1)} \lambda_{\min}(A)^{-(m+1/2)}, \quad (20)$$

1150 where  $\lambda_{\min}(A) > 0$  is the smallest eigenvalue of  $A$  and  $C_m := \frac{1}{m+1} \binom{2m}{m}$ .

1151 We find that  $\partial_l A^{1/2}(x) = \nabla \varphi(A(x))[\partial_l A(x)]$  and

$$1152 \partial_j \partial_l A^{1/2}(x) = \nabla^2 \varphi(A(x))[\partial_j A(x), \partial_l A(x)] + \nabla \varphi(A(x))[\partial_j \partial_l A(x)].$$

1153 Thus, it follows from Eq. (20) that

$$1154 \left\| \partial_l A^{1/2}(x) \right\| \leq C \lambda_{\min}(A(x))^{-1/2} \|\partial_l A(x)\|,$$

1155 and

$$1156 \begin{aligned} 1157 \left\| \partial_j \partial_l A^{1/2}(x) \right\| &\leq C_1 \lambda_{\min}(A(x))^{-(1+1/2)} \|\partial_j A(x)\| \|\partial_l A(x)\| \\ 1158 &\quad + C_2 \lambda_{\min}(A(x))^{-1/2} \|\partial_j \partial_l A(x)\|. \end{aligned}$$

1159 More generally, for  $m = 1, \dots, 7$ ,

$$1160 \begin{aligned} 1161 \left\| \partial_{j_1} \dots \partial_{j_m} A^{1/2}(x) \right\| &\leq C_m \left\{ \sum_{r=0}^{m-1} \lambda_{\min}(A(x))^{-(r+1/2)} \right. \\ 1162 &\quad \left. \times \max_{|\beta| \leq m} \left\| \partial_\beta^{|\beta|} A(x) \right\|^{r+1} \right\}. \end{aligned} \quad (21)$$

1163 Let us now provide a lower bound for  $\lambda_{\min}(A(x))$  in terms of  $\text{tr}(A(x))$  and  $\text{tr}((A(x))^2)$ . Define

$$1164 s^2(x) = n^{-1} \left( \text{tr}((A(x))^2) - \frac{(\text{tr}(A(x)))^2}{n} \right), \quad m(x) = \frac{\text{tr}(A(x))}{n}.$$

Then, from Corollary 2.1, Corollary 2.2, and Theorem 2.1 in Wolkowicz and Styan (1980), we obtain

$$\begin{aligned} \frac{1}{\lambda_{\min}(A(x))} &\leq \frac{1}{\lambda_{\max}(A(x))} \left( 1 + \frac{2s(x)(n-1)^{1/2}}{m(x) - s(x)(n-1)^{-1/2}} \right) \\ &\leq \frac{1}{m(x) + s(x)(n-1)^{1/2}} \left( 1 + \frac{2s(x)(n-1)^{1/2}}{m(x) - s(x)(n-1)^{-1/2}} \right). \end{aligned}$$

Therefore, from Eq. (21), we have for  $m = 1, \dots, 7$  that

$$\left\| \partial_{j_1} \dots \partial_{j_m} A^{1/2}(x) \right\| \leq Cl_m(x), \quad (22)$$

where the function  $l_m(x)$  is defined as

$$\begin{aligned} l_m(x) &:= \sum_{r=0}^{m-1} \left( \frac{1}{m(x) + s(x)(n-1)^{1/2}} \left( 1 + \frac{2s(x)(n-1)^{1/2}}{m(x) - s(x)(n-1)^{-1/2}} \right) \right)^{-(r+1/2)} \\ &\quad \times \max_{|\beta| \leq m} \left\| \partial_{\beta}^{|\beta|} A(x) \right\|^{r+1}. \end{aligned} \quad (23)$$

□

The following results are key to guarantee that an SDE is a weak approximation of an optimizer.

**Lemma C.6** (Lemma 1 Li et al. (2017)). *Let  $0 < \eta < 1$ . Consider a stochastic process  $X_t, t \geq 0$  satisfying the SDE*

$$dX_t = b(X_t) dt + \sqrt{\eta} \sigma(X_t) dW_t$$

*with  $X_0 = x \in \mathbb{R}^d$  and  $b, \sigma$  together with their derivatives belong to  $G$ . Define the one-step difference  $\Delta = X_\eta - x$ , and indicate the  $i$ -th component of  $\Delta$  with  $\Delta_i$ . Then we have*

1.  $\mathbb{E} \Delta_i = b_i \eta + \frac{1}{2} \left[ \sum_{j=1}^d b_j \partial_{e_j} b_i \right] \eta^2 + \mathcal{O}(\eta^3) \quad \forall i = 1, \dots, d;$
2.  $\mathbb{E} \Delta_i \Delta_j = \left[ b_i b_j + \sigma \sigma_{(ij)}^T \right] \eta^2 + \mathcal{O}(\eta^3) \quad \forall i, j = 1, \dots, d;$
3.  $\mathbb{E} \prod_{j=1}^s \Delta_{(i_j)} = \mathcal{O}(\eta^3)$  for all  $s \geq 3, i_j = 1, \dots, d$ .

*All functions above are evaluated at  $x$ .*

**Theorem C.7** (Theorem 2 and Lemma 5, Mil'shtein (1986)). *Let Assumption C.3 hold and let us define  $\bar{\Delta} = x_1 - x$  to be the increment in the discrete-time algorithm, and indicate the  $i$ -th component of  $\bar{\Delta}$  with  $\bar{\Delta}_i$ . If in addition there exists  $K_1, K_2, K_3, K_4 \in G$  so that*

1.  $|\mathbb{E} \Delta_i - \mathbb{E} \bar{\Delta}_i| \leq K_1(x) \eta^2, \quad \forall i = 1, \dots, d;$
2.  $|\mathbb{E} \Delta_i \Delta_j - \mathbb{E} \bar{\Delta}_i \bar{\Delta}_j| \leq K_2(x) \eta^2, \quad \forall i, j = 1, \dots, d;$
3.  $|\mathbb{E} \prod_{j=1}^s \Delta_{i_j} - \mathbb{E} \prod_{j=1}^s \bar{\Delta}_{i_j}| \leq K_3(x) \eta^2, \quad \forall s \geq 3, \quad \forall i_j \in \{1, \dots, d\};$
4.  $\mathbb{E} \prod_{j=1}^3 |\bar{\Delta}_{i_j}| \leq K_4(x) \eta^2, \quad \forall i_j \in \{1, \dots, d\}.$

*Then, there exists a constant  $C$  so that for all  $k = 0, 1, \dots, N$  we have*

$$|\mathbb{E} g(X_{k\eta}) - \mathbb{E} g(x_k)| \leq C\eta.$$

### C.3 LIMITATIONS

Modeling of discrete-time algorithms using SDEs relies on Assumption C.3. As noted by Li et al. (2021), the approximation can fail when the stepsize  $\eta$  is large or if certain conditions on  $\nabla f$  and the noise covariance matrix are not met. Although these issues can be addressed by increasing the order of the weak approximation, we believe that the primary purpose of SDEs is to serve as simplification tools that enhance our intuition: We would not benefit significantly from added complexity.

1242 C.4 FORMAL DERIVATION - SIGNSGD  
1243

1244 In this subsection, we provide the first formal derivation of an SDE model for SignSGD. Let us  
1245 consider the stochastic process  $X_t \in \mathbb{R}^d$  defined as the solution of

$$1246 \quad dX_t = -(1 - 2\mathbb{P}(\nabla f_\gamma(X_t) < 0))dt + \sqrt{\eta} \sqrt{\bar{\Sigma}(X_t)} dW_t, \quad (24)$$

1247 where

$$1248 \quad \bar{\Sigma}(x) = \mathbb{E}[\xi_\gamma(x)\xi_\gamma(x)^\top], \quad (25)$$

1249 and  $\xi_\gamma(x) := \text{sign}(\nabla f_\gamma(x)) - 1 + 2\mathbb{P}(\nabla f_\gamma(x) < 0)$  the noise in the sample  $\text{sign}(\nabla f_\gamma(x))$ . The  
1250 following theorem guarantees that such a process is a 1-order SDE of the discrete-time algorithm of  
1251 SignSGD

$$1252 \quad x_{k+1} = x_k - \eta \text{sign}(f_{\gamma_k}(x_k)), \quad (26)$$

1253 with  $x_0 \in \mathbb{R}^d$ ,  $\eta \in \mathbb{R}^{>0}$  is the step size, the mini-batches  $\{\gamma_k\}$  are modelled as i.i.d. random variables  
1254 uniformly distributed on  $\{1, \dots, N\}$ , and of size  $B \geq 1$ .

1255 Before proceeding, we ensure that the SDE admits a unique solution and that its coefficients are  
1256 sufficiently regular.

1257 **Lemma C.8.** *The drift term  $b(x) := -(1 - 2\mathbb{P}(\nabla f_\gamma(x) < 0))$  is Lipschitz, satisfies affine growth,  
1258 and belongs to the space  $G$  together with its derivatives.*

1259 *Proof.* Since we are assuming that the gradient noise has a smooth and bounded probability density  
1260 function,<sup>7</sup> the drift can be rewritten in terms of the CDF  $F_Z(x)$  of the noise as  $b(x) :=$   
1261  $2F_Z(-\nabla f(x)) - 1$ , whose derivative is  $-2F'_Z(-\nabla f(x))\nabla^2 f(x)$ . Since the density function and  
1262 the Hessian of  $f$  are bounded, we conclude that the derivative is bounded. Therefore, the drift is  
1263 Lipschitz and as regular as  $\nabla f$ , meaning that each entry is in  $G$ , together with its derivatives. Finally,  
1264 since it is bounded, it has affine growth.  $\square$

1265 **Lemma C.9.** *The diffusion coefficient  $\sqrt{\bar{\Sigma}}$  satisfies the affine growth condition.*

1266 *Proof.* Since it is bounded, the result follows immediately.  $\square$

1267 **Lemma C.10.** *Let us assume the same assumptions as Lemma C.4. Additionally, assume that*

$$1268 \quad \sup_{x \in K} |g(x, \cdot)| \in L^1(\mathbb{R}^n)$$

1269 *for all compact sets  $K$ . Then the entries of  $\bar{\Sigma}$  in Eq. 25 are in  $C_b^1(\mathbb{R}^n)$ .*

1270 *Proof.* By the definition of  $\bar{\Sigma}$  in terms of the sign-function and dominated convergence, from the  
1271 additional assumption on  $g$ , it follows that  $\bar{\Sigma}$  is continuous. So Lemma C.4 entails that the entries of  
1272  $\bar{\Sigma}$  are in  $C_b^1(\mathbb{R}^n)$ .  $\square$

1273 **Lemma C.11.** *Under the assumption that*

$$1274 \quad g(x, y) > 0, \quad (27)$$

1275 *the covariance matrix  $\bar{\Sigma}$  is positive definite.*

1276 *Proof.* For  $y = (y_1, \dots, y_n)^T$ , observe that

$$1277 \quad (\bar{\Sigma}(x)y, y) = \sum_{i,j=1}^n y_i \mathbb{E}[\xi_\gamma^i(x)\xi_\gamma^j(x)] y_j = \mathbb{E} \left[ \left( \sum_{i=1}^n \xi_\gamma^i(x) y_i \right)^2 \right].$$

1278 Using the definition of  $\xi_\gamma$  and the positivity of the density  $g$ , we can argue by contradiction and  
1279 see that for  $y \neq 0$ , the right-hand side of the equation must be strictly greater than zero for all  $x$ .  
1280 Therefore,  $\bar{\Sigma}(x) \in \mathbb{S}$  for all  $x$ , where  $\mathbb{S}$  denotes the open set of positive definite matrices in the space  
1281 of symmetric  $n \times n$  matrices.  $\square$

1282 <sup>7</sup>This is commonly assumed in the literature. Among others, Ahn et al. (2012); Chen et al. (2014); Mandt  
1283 et al. (2016); Stephan et al. (2017); Zhu et al. (2019a); Wu et al. (2020); Xie et al. (2021) assume that it is  
1284 Gaussian, while Jastrzebski et al. (2018) offers an intuitive justification.



**Corollary C.12.** Since  $\bar{\Sigma}$  is positive definite and its entries are in  $C_b^7(\mathbb{R}^n)$ ,  $\sqrt{\bar{\Sigma}}$  is Lipschitz.

*Proof.* The function

$$\varphi : \mathbb{S} \rightarrow \mathbb{S}, \quad A \mapsto \sqrt{A}$$

has Fréchet derivatives of any order on  $\mathbb{S}$  (see e.g. Del Moral and Niclas (2018)). Therefore,  $\bar{\Sigma}^{1/2} \in C^7(\mathbb{R}^n)$ , and since  $\bar{\Sigma} \in C_b^7(\mathbb{R}^n)$ ,  $\bar{\Sigma}^{1/2}$  is Lipschitz continuous (see Proposition 6.2 in Ikeda and Watanabe (2014)).  $\square$

**Proposition C.13.** Assume the conditions of Lemma C.5 and assume that the functions  $l_m(x)$  for  $m = 1, \dots, 7$  in Eq. (18) are of polynomial growth. Then  $\bar{\Sigma}^{1/2} \in G$  together with its derivatives.

**Corollary C.14.** If the noise  $Z(x) \sim \mathcal{N}(0, \Sigma)$  or  $Z(x) \sim t_\nu(0, \Sigma)$ , then  $\bar{\Sigma}^{1/2} \in G$  together with its derivatives.

*Proof.* With the definition of  $\Xi(x)$  given in Lemma 3.10, the function  $K(x) := \sqrt{1 - 4\Xi(x)^2}$  is in  $G$  together with its derivative: It is easy to verify that all the derivatives of  $K(x)$  are bounded even in the case  $\nu = 1$ , which is the most pathological one. Therefore,  $\sqrt{\bar{\Sigma}}(x)$  is in  $G$  together with its derivatives.  $\square$

*Remark C.15.* Based on the above results, we have that under mild assumptions on the noise structures (see Sec. C.1) that cover and generalize the well-accepted Gaussianity, e.g. covering Student's  $t$  as well, the SDE of SignSGD admits a unique solution and its coefficients are regular enough to apply Lemma C.6 and Thm. C.7.

**Theorem C.16** (Stochastic modified equations). Let  $0 < \eta < 1, T > 0$  and set  $N = \lfloor T/\eta \rfloor$ . Let  $x_k \in \mathbb{R}^d, 0 \leq k \leq N$  denote a sequence of SignSGD iterations defined by Eq. 26. Consider the stochastic process  $X_t$  defined in Eq. 24 and fix some test function  $g \in G$  and suppose that  $g$  and its partial derivatives up to order 6 belong to  $G$ . Then, under Assumption C.3, there exists a constant  $C > 0$  independent of  $\eta$  such that for all  $k = 0, 1, \dots, N$ , we have

$$|\mathbb{E}g(X_{k\eta}) - \mathbb{E}g(x_k)| \leq C\eta.$$

That is, the SDE 24 is an order 1 weak approximation of the SignSGD iterations 26.

**Lemma C.17.** Under the assumptions of Theorem C.16, let  $0 < \eta < 1$  and consider  $x_k, k \geq 0$  satisfying the SignSGD iterations

$$x_{k+1} = x_k - \eta \text{sign}(\nabla f_{\gamma_k}(x_k))$$

with  $x_0 \in \mathbb{R}^d$ . From the definition the one-step difference  $\bar{\Delta} = x_1 - x$ , then we have

1.  $\mathbb{E}\bar{\Delta}_i = -(1 - 2\mathbb{P}(\partial_i f_\gamma < 0))\eta \quad \forall i = 1, \dots, d;$
2.  $\mathbb{E}\bar{\Delta}_i \bar{\Delta}_j = ((1 - 2\mathbb{P}(\partial_i f_\gamma < 0))(1 - 2\mathbb{P}(\partial_j f_\gamma < 0)) + \bar{\Sigma}_{(ij)})\eta^2 \quad \forall i, j = 1, \dots, d;$
3.  $\mathbb{E}\prod_{j=1}^s \bar{\Delta}_{i_j} = \mathcal{O}(\eta^3) \quad \forall s \geq 3, \quad i_j \in \{1, \dots, d\}.$

All the functions above are evaluated at  $x$ .

*Proof of Lemma C.17.* First of all, we have that by definition

$$\mathbb{E}[x_1^i - x^i] = -\eta \mathbb{E}[\text{sign}(\partial_i f_\gamma(x))], \quad (28)$$

which implies

$$\mathbb{E}\bar{\Delta}_i = -(1 - 2\mathbb{P}(\partial_i f_\gamma(x) < 0))\eta \quad \forall i = 1, \dots, d. \quad (29)$$

Second, we have that by definition

1350

1351

1352

$$\mathbb{E} \left[ (x_1 - x) (x_1 - x)^\top \right] = \mathbb{E} [(x_1 - x)] \mathbb{E} \left[ (x_1 - x)^\top \right] + \quad (30)$$

1353

1354

$$\mathbb{E} \left[ (\text{sign}(\nabla f_\gamma(x)) - 1 + 2\mathbb{P}(\nabla f_\gamma(x) < 0)) \right] \quad (31)$$

1355

1356

$$(\text{sign}(\nabla f_\gamma(x)) - 1 + 2\mathbb{P}(\nabla f_\gamma(x) < 0))^\top \eta^2, \quad (32)$$

1357

which implies that

1358

1359

1360

$$\mathbb{E} \bar{\Delta}_i \bar{\Delta}_j = (1 - 2\mathbb{P}(\partial_i f_\gamma < 0)) (1 - 2\mathbb{P}(\partial_j f_\gamma < 0)) \eta^2 + \bar{\Sigma}_{(ij)} \eta^2 \quad \forall i, j = 1, \dots, d. \quad (33)$$

1361

1362

Finally, by definition

1363

1364

1365

$$\mathbb{E} \prod_{j=1}^s \bar{\Delta}_{i_j} = \mathcal{O}(\eta^3) \quad \forall s \geq 3, \quad i_j \in \{1, \dots, d\}, \quad (34)$$

1366

1367

which concludes our proof.  $\square$

1368

1369

*Proof of Theorem C.16.* To prove this result, all we need to do is check the conditions in Theorem C.7. As we apply Lemma C.6, we make the following choices:

1370

1371

1372

1373

- $b(x) = -(1 - 2\mathbb{P}(\nabla f_\gamma(x) < 0))$ ;

1374

1375

- $\sigma(x) = \sqrt{\bar{\Sigma}(x)}$ .

1376

1377

First of all, we notice that  $\forall i = 1, \dots, d$ , it holds that

1378

1379

1380

- $\mathbb{E} \bar{\Delta}_i \stackrel{1. \text{Lemma C.17}}{=} -(1 - 2\mathbb{P}(\partial_i f_\gamma(x) < 0)) \eta$ ;

1381

1382

- $\mathbb{E} \Delta_i \stackrel{1. \text{Lemma C.6}}{=} -(1 - 2\mathbb{P}(\partial_i f_\gamma(x) < 0)) \eta + \mathcal{O}(\eta^2)$ .

1383

Therefore, we have that for some  $K_1(x) \in G$ ,

1384

1385

$$|\mathbb{E} \Delta_i - \mathbb{E} \bar{\Delta}_i| \leq K_1(x) \eta^2, \quad \forall i = 1, \dots, d. \quad (35)$$

1386

Additionally, we notice that  $\forall i, j = 1, \dots, d$ , it holds that

1387

1388

1389

- $\mathbb{E} \bar{\Delta}_i \bar{\Delta}_j \stackrel{2. \text{Lemma C.17}}{=} (1 - 2\mathbb{P}(\partial_i f_\gamma(x) < 0)) (1 - 2\mathbb{P}(\partial_j f_\gamma(x) < 0)) \eta^2 + \bar{\Sigma}_{(ij)}(x) \eta^2$ ;

1390

1391

1392

- $\mathbb{E} \Delta_i \Delta_j \stackrel{2. \text{Lemma C.6}}{=} ((1 - 2\mathbb{P}(\partial_i f_\gamma(x) < 0)) (1 - 2\mathbb{P}(\partial_j f_\gamma(x) < 0)) + \bar{\Sigma}_{(ij)}(x)) \eta^2 + \mathcal{O}(\eta^3)$ .

1393

1394

1395

Therefore, we have that for some  $K_2(x) \in G$ ,

1396

1397

$$|\mathbb{E} \Delta_i \Delta_j - \mathbb{E} \bar{\Delta}_i \bar{\Delta}_j| \leq K_2(x) \eta^2, \quad \forall i, j = 1, \dots, d. \quad (36)$$

1398

Additionally, we notice that  $\forall s \geq 3, \forall i_j \in \{1, \dots, d\}$ , it holds that

1399

1400

1401

- $\mathbb{E} \prod_{j=1}^s \bar{\Delta}_{i_j} \stackrel{3. \text{Lemma C.17}}{=} \mathcal{O}(\eta^3)$ ;

1402

- $\mathbb{E} \prod_{j=1}^s \Delta_{i_j} \stackrel{3. \text{Lemma C.6}}{=} \mathcal{O}(\eta^3)$ .

1403

Therefore, we have that for some  $K_3(x) \in G$ ,

$$\left| \mathbb{E} \prod_{j=1}^s \Delta_{i_j} - \mathbb{E} \prod_{j=1}^s \bar{\Delta}_{i_j} \right| \leq K_3(x) \eta^2. \quad (37)$$

Additionally, for some  $K_4(x) \in G, \forall i_j \in \{1, \dots, d\}$ ,

$$\mathbb{E} \prod_{j=1}^3 |\bar{\Delta}_{(i_j)}| \stackrel{3. \text{ Lemma C.17}}{\leq} K_4(x) \eta^2. \quad (38)$$

*Remark C.18.* Remembering Remark C.15, and thanks to Eq. 35, Eq. 36, Eq. 37, and Eq. 38, the thesis follows from Lemma C.6 and Thm. C.7.

□

In all the following results, the reader will notice that all the drifts, diffusion terms, and noise assumptions are selected to guarantee that the SDE we derived for SignSGD is indeed a 1 weak approximation for SignSGD even without the mollification argument used in Li et al. (2019) to handle the regularity issues.

**Corollary C.19.** *Let us take the same assumptions of Theorem C.16, and that the stochastic gradient is  $\nabla f_\gamma(x) = \nabla f(x) + Z$  such that  $Z \sim \mathcal{N}(0, \Sigma)$  that does not depend on  $x$ . Then, the following SDE provides a 1 weak approximation of the discrete update of SignSGD*

$$dX_t = -\text{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right) dt + \sqrt{\eta} \sqrt{I_d - \text{diag} \left( \text{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right) \right)^2} dW_t, \quad (39)$$

where the error function  $\text{Erf}(x)$  and the square are applied component-wise, and  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$ .

*Proof of Corollary C.19.* First of all, we observe that

$$1 - 2\mathbb{P}(\nabla f_\gamma(x) < 0) = 1 - 2\mathbb{P}(\nabla f(x) + \Sigma^{\frac{1}{2}} Z < 0) = 1 - 2\Phi(-\Sigma^{-\frac{1}{2}} \nabla f(x)), \quad (40)$$

where  $\Phi$  is the cumulative distribution function of the standardized normal distribution. Remembering that

$$\Phi(x) = \frac{1}{2} \left( 1 + \text{Erf} \left( \frac{x}{\sqrt{2}} \right) \right), \quad (41)$$

we have that

$$1 - 2\mathbb{P}(\nabla f_\gamma(x) < 0) = 1 - 2 \frac{1}{2} \left( 1 + \text{Erf} \left( -\frac{\Sigma^{-\frac{1}{2}} \nabla f(x)}{\sqrt{2}} \right) \right) = \text{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(x)}{\sqrt{2}} \right). \quad (42)$$

Similarly, one can prove that  $\bar{\Sigma}$  defined in 25 becomes

$$\bar{\Sigma} = I_d - \text{diag} \left( \text{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right) \right)^2. \quad (43)$$

□

**Corollary C.20.** *Let us take the same assumptions of Theorem C.16, and that the stochastic gradient is  $\nabla f_\gamma(x) = \nabla f(x) + \sqrt{\Sigma}Z$  such that  $Z \sim t_\nu(0, I_d)$  that does not depend on  $x$  and  $\nu$  is a positive integer number. Then, the following SDE provides a 1 weak approximation of the discrete update of SignSGD*

$$dX_t = -2\Xi \left( \Sigma^{-\frac{1}{2}} \nabla f(X_t) \right) dt + \sqrt{\eta} \sqrt{I_d - 4 \text{diag} \left( \Xi \left( \Sigma^{-\frac{1}{2}} \nabla f(X_t) \right) \right)^2} dW_t, \quad (44)$$

where  $\Xi(x)$  is defined as

$$\Xi(x) := x \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\pi\nu}\Gamma\left(\frac{\nu}{2}\right)} {}_2F_1\left(\frac{1}{2}, \frac{\nu+1}{2}; \frac{3}{2}; -\frac{x^2}{\nu}\right), \quad (45)$$

and  ${}_2F_1(a, b; c; x)$  is the hypergeometric function. Above, function  $\Xi(x)$  and the square are applied component-wise, and  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$ .

*Proof.* First of all, we observe that

$$1 - 2\mathbb{P}(\nabla f_\gamma(x) < 0) = 1 - 2\mathbb{P}\left(\nabla f(x) + \Sigma^{\frac{1}{2}}Z < 0\right) = 1 - 2F_\nu\left(-\Sigma^{-\frac{1}{2}}\nabla f(x)\right), \quad (46)$$

where  $F_\nu(x)$  is the cumulative function of a  $t$  distribution with  $\nu$  degrees of freedom. Remembering that

$$F_\nu(x) = \frac{1}{2} + \Xi_\nu(x), \quad (47)$$

we have that

$$1 - 2\mathbb{P}(\nabla f_\gamma(x) < 0) = 1 - 2\left(\frac{1}{2} + \Xi_\nu(-\Sigma^{-\frac{1}{2}}\nabla f(x))\right) = 2\Xi_\nu(\Sigma^{-\frac{1}{2}}\nabla f(x)). \quad (48)$$

Similarly, one can prove that  $\bar{\Sigma}$  becomes

$$\bar{\Sigma} = I_d - 4 \text{diag} \left( \Xi_\nu \left( \Sigma^{-\frac{1}{2}} \nabla f(X_t) \right) \right)^2. \quad (49)$$

□

**Lemma C.21.** *Under the assumptions of Corollary C.19 and signal-to-noise ratio  $Y_t := \frac{\Sigma^{-\frac{1}{2}}\nabla f(X_t)}{\sqrt{2}}$ ,*

1. **Phase 1:** *If  $|Y_t| > \frac{3}{2}$ , the SDE coincides with the ODE of SignGD:*

$$dX_t = -\text{sign}(\nabla f(X_t))dt; \quad (50)$$

2. **Phase 2:** *If  $1 < |Y_t| < \frac{3}{2}$ :*

- (a)  $mY_t + \mathbf{q}^- \leq \frac{d\mathbb{E}[X_t]}{dt} \leq mY_t + \mathbf{q}^+$ ;
- (b)  $\mathbb{P}\left[\|X_t - \mathbb{E}[X_t]\|_2^2 > a\right] \leq \frac{\eta}{a} (d - \|mY_t + \mathbf{q}^-\|_2^2)$ ;

3. **Phase 3:** *If  $|Y_t| < 1$ , the SDE is*

$$dX_t = -\sqrt{\frac{2}{\pi}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) dt + \sqrt{\eta} \sqrt{I_d - \frac{2}{\pi} \text{diag} \left( \Sigma^{-\frac{1}{2}} \nabla f(X_t) \right)^2} dW_t. \quad (51)$$

*Proof of Lemma C.21.* Exploiting the regularity of the Erf function, we approximate the SDE in Eq. 39 in three different regions:

1. **Phase 1:** *If  $|x| > \frac{3}{2}$ ,  $\text{Erf}(x) \sim \text{sign}(x)$ . Therefore, if  $\left|\frac{\Sigma^{-\frac{1}{2}}\nabla f(X_t)}{\sqrt{2}}\right| > \frac{3}{2}$ ,*

- (a)  $\text{Erf}\left(\frac{\Sigma^{-\frac{1}{2}}\nabla f(X_t)}{\sqrt{2}}\right) \sim \text{sign}\left(\frac{\Sigma^{-\frac{1}{2}}\nabla f(X_t)}{\sqrt{2}}\right) = \text{sign}(\nabla f(X_t))$ ;

$$(b) \operatorname{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right)^2 \sim \operatorname{sign} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right)^2 = (1, \dots, 1).$$

Therefore,

$$\begin{aligned} dX_t &= -\operatorname{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right) dt + \sqrt{\eta} \sqrt{I_d - \operatorname{diag} \left( \operatorname{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right) \right)^2} dW_t \\ &\sim -\operatorname{sign}(\nabla f(X_t)); \end{aligned} \quad (52)$$

2. **Phase 2:** Let  $m$  and  $q_1$  are the slope and intercept of the line secant to the graph of  $\operatorname{Erf}(x)$  between the points  $(1, \operatorname{Erf}(1))$  and  $(\frac{3}{2}, \operatorname{Erf}(\frac{3}{2}))$ , while  $q_2$  is the intercept of the line tangent to the graph of  $\operatorname{Erf}(x)$  and slope  $m$ . If  $1 < x < \frac{3}{2}$ , we have that

$$mx + q_1 < \operatorname{Erf}(x) < mx + q_2. \quad (53)$$

Analogously, if  $-\frac{3}{2} < x < -1$

$$mx - q_2 < \operatorname{Erf}(x) < mx - q_1. \quad (54)$$

Therefore, we have that if  $1 < \left| \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right| < \frac{3}{2}$ , then

(a)

$$\frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^- < \operatorname{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right) < \frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^+, \quad (55)$$

where

$$(\mathbf{q}^+)_i := \begin{cases} q_2 & \text{if } \partial_i f(x) > 0 \\ -q_1 & \text{if } \partial_i f(x) < 0, \end{cases} \quad (56)$$

and

$$(\mathbf{q}^-)_i := \begin{cases} q_1 & \text{if } \partial_i f(x) > 0 \\ -q_2 & \text{if } \partial_i f(x) < 0, \end{cases} \quad (57)$$

Therefore,

$$-\frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) - \mathbf{q}^+ \leq \frac{d\mathbb{E}[X_t]}{dt} \leq -\frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) - \mathbf{q}^-; \quad (58)$$

(b) Similar to the above,

$$\left( \frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^- \right)^2 \leq \operatorname{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right)^2 \leq \left( \frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^+ \right)^2.$$

Therefore,

$$\begin{aligned} \mathbb{P} [\|X_t - \mathbb{E}[X_t]\|_2^2 > a] &\leq \mathbb{P} [\exists i \text{ s.t. } |X_t^i - \mathbb{E}[X_t^i]|^2 > a] \\ &\leq \sum_i \mathbb{P} [|X_t^i - \mathbb{E}[X_t^i]| > \sqrt{a}] \end{aligned} \quad (59)$$

$$\leq \frac{\eta}{a} \sum_i \left( 1 - \operatorname{Erf} \left( \frac{\Sigma_i^{-\frac{1}{2}} \partial_i f(X_t)}{\sqrt{2}} \right)^2 \right) \quad (60)$$

$$< \frac{\eta}{a} \left( d - \left\| \frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^- \right\|_2^2 \right). \quad (61)$$

3. **Phase 3:** If  $|x| < 1$ ,  $\operatorname{Erf}(x) \sim \frac{2}{\sqrt{\pi}}x$ . Therefore, if  $\left| \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right| < 1$ ,

$$(a) \operatorname{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right) \sim \sqrt{\frac{2}{\pi}} \Sigma^{-\frac{1}{2}} \nabla f(X_t);$$

$$(b) \left( \text{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right) \right)^2 \sim \frac{2}{\pi} \left( \Sigma^{-\frac{1}{2}} \nabla f(X_t) \right)^2.$$

Therefore,

$$\begin{aligned} dX_t &= -\text{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right) dt + \sqrt{\eta} \sqrt{I_d - \text{diag} \left( \text{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right) \right)^2} dW_t \\ &\sim -\sqrt{\frac{2}{\pi}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) dt + \sqrt{\eta} \sqrt{I_d - \frac{2}{\pi} \text{diag} \left( \Sigma^{-\frac{1}{2}} \nabla f(X_t) \right)^2} dW_t. \end{aligned} \quad (62)$$

□

**Lemma C.22** (Dynamics of Expected Loss). *Let  $f$  be  $\mu$ -strongly convex,  $\text{Tr}(\nabla^2 f(x)) \leq \mathcal{L}_\tau$ , and  $S_t := f(X_t) - f(X_*)$ . Then, during*

1. Phase 1, the dynamics will stop before  $t_* = 2\sqrt{\frac{S_0}{\mu}}$  because  $S_t \leq \frac{1}{4} (\sqrt{\mu t} - 2\sqrt{S_0})^2$ ;
2. Phase 2 with  $\Delta := \left( \frac{m}{\sqrt{2}\sigma_{\max}} + \frac{\eta\mu m^2}{4\sigma_{\max}^2} \right)$ :  $\mathbb{E}[S_t] \leq S_0 e^{-2\mu\Delta t} + \frac{\eta}{2} \frac{(\mathcal{L}_\tau - \mu d\bar{q}^2)}{2\mu\Delta} (1 - e^{-2\mu\Delta t})$ ;
3. Phase 3 with  $\Delta := \left( \sqrt{\frac{2}{\pi}} \frac{1}{\sigma_{\max}} + \frac{\eta}{\pi} \frac{\mu}{\sigma_{\max}^2} \right)$ :  $\mathbb{E}[S_t] \leq S_0 e^{-2\mu\Delta t} + \frac{\eta}{2} \frac{\mathcal{L}_\tau}{2\mu\Delta} (1 - e^{-2\mu\Delta t})$ .

*Proof of Lemma C.22.* We prove each point by leveraging the shape of the law of  $X_t$  derived in Lemma C.21:

### 1. Phase 1:

$$d(f(X_t) - f(X_*)) = -\nabla f(X_t) \text{sign}(\nabla f(X_t)) dt = -\|\nabla f(X_t)\|_1 dt \leq -\|\nabla f(X_t)\|_2 dt \quad (63)$$

Since  $f$  is  $\mu - PL$ , we have that  $-\|\nabla f(X_t)\|_2^2 < -2\mu(f(X_t) - f(X_*))$ , which implies that

$$f(X_t) - f(X_*) \leq \frac{1}{4} \left( \sqrt{\mu t} - 2\sqrt{f(X_0) - f(X_*)} \right)^2, \quad (64)$$

meaning that the dynamics will stop before  $t_* = 2\sqrt{\frac{f(X_0) - f(X_*)}{\mu}}$ ;

### 2. Phase 2:

By applying the Itô Lemma to  $f(X_t) - f(X_*)$  and that

$$\frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^- < \text{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right) < \frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^+, \quad (65)$$

we have that if  $\hat{q} := \max(q_1, q_2)$ ,

$$d(f(X_t) - f(X_*)) \leq - \left( \frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^- \right)^\top \nabla f(X_t) dt + \mathcal{O}(\text{Noise}) \quad (66)$$

$$+ \frac{\eta}{2} \text{Tr} \left[ \nabla^2 f(X_t) \left( I_d - \text{diag} \left( \frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^- \right)^2 \right) \right] dt \quad (67)$$

$$\leq - \frac{m}{\sqrt{2}} \frac{1}{\sigma_{\max}} \|\nabla f(X_t)\|_2^2 dt - \hat{q} \|\nabla f(X_t)\|_1 dt + \frac{\eta \mathcal{L}_\tau}{2} dt \quad (68)$$

$$- \frac{\eta \mu}{2} \left\| \frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^- \right\|_2^2 dt + \mathcal{O}(\text{Noise}) \quad (69)$$

$$\leq - \frac{m}{\sqrt{2}} \frac{1}{\sigma_{\max}} \|\nabla f(X_t)\|_2^2 dt - \hat{q} \|\nabla f(X_t)\|_1 dt + \frac{\eta \mathcal{L}_\tau}{2} dt \quad (70)$$

$$- \frac{\eta \mu m^2}{4\sigma_{\max}^2} \|\nabla f(X_t)\|_2^2 dt - \frac{\eta \mu d \hat{q}^2}{2} dt - \frac{\sqrt{2} m \hat{q}}{\sigma_{\max}} \|\nabla f(X_t)\|_1 dt \quad (71)$$

$$+ \mathcal{O}(\text{Noise}) \quad (72)$$

$$\leq - 2\mu \left( \frac{m}{\sqrt{2}\sigma_{\max}} + \frac{\eta \mu m^2}{4\sigma_{\max}^2} \right) (f(X_t) - f(X_*)) dt \quad (73)$$

$$+ \frac{\eta}{2} (\mathcal{L}_\tau - \mu d \hat{q}^2) dt + \mathcal{O}(\text{Noise}), \quad (74)$$

which implies that if  $k := 2\mu \left( \frac{m}{\sqrt{2}\sigma_{\max}} + \frac{\eta \mu m^2}{4\sigma_{\max}^2} \right)$ ,

$$\mathbb{E}[f(X_t) - f(X_*)] \leq (f(X_0) - f(X_*)) e^{-kt} + \frac{\eta (\mathcal{L}_\tau - \mu d \hat{q}^2)}{2k} (1 - e^{-kt}). \quad (75)$$

3. **Phase 3:** By applying the Itô Lemma to  $f(X_t) - f(X_*)$ , we have that:

$$d(f(X_t) - f(X_*)) = - \sqrt{\frac{2}{\pi}} \nabla f(X_t)^\top \Sigma^{-\frac{1}{2}} \nabla f(X_t) dt + \mathcal{O}(\text{Noise}) \quad (76)$$

$$+ \frac{\eta}{2} \text{Tr} \left( \left( I_d - \frac{2}{\pi} \text{diag} \left( \Sigma^{-\frac{1}{2}} \nabla f(X_t) \right)^2 \right) \nabla^2 f(X_t) \right) dt \quad (77)$$

$$\leq - \sqrt{\frac{2}{\pi}} \frac{1}{\sigma_{\max}} \|\nabla f(X_t)\|_2^2 dt + \mathcal{O}(\text{Noise}) \quad (78)$$

$$+ \frac{\eta}{2} \text{Tr}(\nabla^2 f(X_t)) dt - \frac{\eta}{\pi} \frac{\mu}{\sigma_{\max}^2} \|\nabla f(X_t)\|_2^2 dt \quad (79)$$

$$\leq - \left( \sqrt{\frac{2}{\pi}} \frac{1}{\sigma_{\max}} + \frac{\eta}{\pi} \frac{\mu}{\sigma_{\max}^2} \right) \|\nabla f(X_t)\|_2^2 dt \quad (80)$$

$$+ \frac{\eta}{2} \text{Tr}(\nabla^2 f(X_t)) dt + \mathcal{O}(\text{Noise}) \quad (81)$$

Since  $f$  is  $\mu$ -Strongly Convex,  $f$  is also  $\mu$ -PL. Therefore, we have

$$d(f(X_t) - f(X_*)) \leq - 2\mu \left( \sqrt{\frac{2}{\pi}} \frac{1}{\sigma_{\max}} + \frac{\eta}{\pi} \frac{\mu}{\sigma_{\max}^2} \right) (f(X_t) - f(X_*)) dt \quad (82)$$

$$+ \frac{\eta}{2} \text{Tr}(\nabla^2 f(X_t)) dt + \mathcal{O}(\text{Noise}). \quad (83)$$

Therefore,

$$d\mathbb{E}[f(X_t) - f(X_*)] \leq - 2\mu \left( \sqrt{\frac{2}{\pi}} \frac{1}{\sigma_{\max}} + \frac{\eta}{\pi} \frac{\mu}{\sigma_{\max}^2} \right) (\mathbb{E}[f(X_t) - f(X_*)]) dt + \frac{\eta}{2} \mathcal{L}_\tau dt, \quad (84)$$

which implies that if  $k := 2\mu \left( \sqrt{\frac{2}{\pi}} \frac{1}{\sigma_{\max}} + \frac{\eta}{\pi} \frac{\mu}{\sigma_{\max}^2} \right)$ ,

$$\mathbb{E}[f(X_t) - f(X_*)] \leq (f(X_0) - f(X_*))e^{-kt} + \frac{\eta \mathcal{L}_\tau}{2k} (1 - e^{-kt}). \quad (85)$$

□

We can weaken the regularity of  $f$  from  $\mu$ -strongly convex to  $\mu$ -PL: This results in less tight bounds as expected.

**Lemma C.23** (Dynamics of Expected Loss). *Let  $f$  be  $\mu$ -PL,  $L$ -smooth, and  $S_t := f(X_t) - f(X_*)$ . Then, during*

1. Phase 1, the dynamics will stop before  $t_* = 2\sqrt{\frac{S_0}{\mu}}$  because  $S_t \leq \frac{1}{4}(\sqrt{\mu t} - 2\sqrt{S_0})^2$ ;
2. Phase 2 with  $\Delta := \frac{m}{\sqrt{2}\sigma_{\max}}$ :  $\mathbb{E}[S_t] \leq S_0 e^{-2\mu\Delta t} + \frac{\eta L d}{4\mu\Delta} (1 - e^{-2\mu\Delta t})$ ;
3. Phase 3 with  $\Delta := \sqrt{\frac{2}{\pi}} \frac{1}{\sigma_{\max}}$ :  $\mathbb{E}[S_t] \leq S_0 e^{-2\mu\Delta t} + \frac{\eta L d}{4\mu\Delta} (1 - e^{-2\mu\Delta t})$ .

*Proof of Lemma C.22.* We prove each point by leveraging the shape of the law of  $X_t$  derived in Lemma C.21:

### 1. Phase 1:

$$d(f(X_t) - f(X_*)) = -\nabla f(X_t) \text{sign}(\nabla f(X_t)) dt = -\|\nabla f(X_t)\|_1 dt \leq -\|\nabla f(X_t)\|_2 dt. \quad (86)$$

Since  $f$  is  $\mu$ -PL, we have that  $-\|\nabla f(X_t)\|_2^2 < -2\mu(f(X_t) - f(X_*))$ , which implies that

$$f(X_t) - f(X_*) \leq \frac{1}{4} \left( \sqrt{\mu t} - 2\sqrt{f(X_0) - f(X_*)} \right)^2, \quad (87)$$

meaning that the dynamics will stop before  $t_* = 2\sqrt{\frac{f(X_0) - f(X_*)}{\mu}}$ ;

### 2. Phase 2:

By applying the Itô Lemma to  $f(X_t) - f(X_*)$  and that

$$\frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^- < \text{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right) < \frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^+, \quad (88)$$

we have that if  $\hat{q} := \max(q_1, q_2)$ ,

$$d(f(X_t) - f(X_*)) \leq - \left( \frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^- \right)^\top \nabla f(X_t) dt + \mathcal{O}(\text{Noise}) \quad (89)$$

$$+ \frac{\eta}{2} \text{Tr} \left[ \nabla^2 f(X_t) \left( I_d - \text{diag} \left( \frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^- \right)^2 \right) \right] dt \quad (90)$$

$$\leq - \frac{m}{\sqrt{2}} \frac{1}{\sigma_{\max}} \|\nabla f(X_t)\|_2^2 dt - \hat{q} \|\nabla f(X_t)\|_1 dt + \frac{\eta L d}{2} dt \quad (91)$$

$$\leq -2\mu \frac{m}{\sqrt{2}\sigma_{\max}} (f(X_t) - f(X_*)) dt + \frac{\eta L d}{2} dt + \mathcal{O}(\text{Noise}), \quad (92)$$

which implies that if  $\Delta := \frac{m}{\sqrt{2}\sigma_{\max}}$ ,

$$\mathbb{E}[f(X_t) - f(X_*)] \leq (f(X_0) - f(X_*))e^{-2\mu\Delta t} + \frac{\eta L d}{4\mu\Delta} (1 - e^{-2\mu\Delta t}). \quad (93)$$



1728 3. **Phase 3:** By applying the Itô Lemma to  $f(X_t) - f(X_*)$ , we have that:

1729  
1730  
1731 
$$d(f(X_t) - f(X_*)) = -\sqrt{\frac{2}{\pi}} \nabla f(X_t)^\top \Sigma^{-\frac{1}{2}} \nabla f(X_t) dt + \mathcal{O}(\text{Noise}) + \frac{\eta L d}{2} dt \quad (94)$$

1732  
1733 
$$\leq -\sqrt{\frac{2}{\pi}} \frac{1}{\sigma_{\max}} \|\nabla f(X_t)\|_2^2 dt + \mathcal{O}(\text{Noise}) + \frac{\eta L d}{2} dt \quad (95)$$

1734 Since  $f$  is  $\mu$ -PL, we have

1735  
1736 
$$d(f(X_t) - f(X_*)) \leq -2\mu \sqrt{\frac{2}{\pi}} \frac{1}{\sigma_{\max}} (f(X_t) - f(X_*)) dt + \frac{\eta L d}{2} dt + \mathcal{O}(\text{Noise}). \quad (96)$$

1737  
1738 Therefore, for  $\Delta := \sqrt{\frac{2}{\pi}} \frac{1}{\sigma_{\max}}$ ,

1739  
1740 
$$\mathbb{E}[f(X_t) - f(X_*)] \leq (f(X_0) - f(X_*)) e^{-2\mu\Delta t} + \frac{\eta L d}{4\mu\Delta} (1 - e^{-2\mu\Delta t}). \quad (97)$$

1741  
1742  
1743  
1744 □

1745 We can weaken the regularity of  $f$  from  $\mu$ -PL to  $L$ -Smooth: Of course, we can only bound the  
1746 expected norm of the gradient.

1747 **Lemma C.24** (Dynamics of Expected Gradient Norm). *Let  $f$  be  $L$ -smooth,  $\eta_t$  be a learning rate  
1748 scheduler such that  $\lim_{t \rightarrow \infty} \frac{\phi_t^2}{\phi_t^1} \xrightarrow{t \rightarrow \infty} 0$  and  $\phi_t^1 \xrightarrow{t \rightarrow \infty} \infty$ , where  $\phi_t^i = \int_0^t (\eta_s)^i ds$ . Then, during*

1749  
1750 1. Phase 1,  $\|\nabla f(X_{\tilde{t}^1})\|_1 \leq \frac{f(X_0) - f(X_*)}{\phi_t^1} \xrightarrow{t \rightarrow \infty} 0$ ;

1751  
1752 2. Phase 2,

1753  
1754 
$$\left( \frac{m}{\sqrt{2}} \mathbb{E} \|\nabla f(X_{\tilde{t}^{(1,2)}})\|_2^2 + \hat{q} \sigma_{\max} \mathbb{E} \|\nabla f(X_{\tilde{t}^{(2,2)}})\|_1 \right) \leq \sigma_{\max} \left( \frac{f(X_0) - f(X_*)}{\phi_t^1} + \frac{\eta L d \phi_t^2}{2 \phi_t^1} \right) \xrightarrow{t \rightarrow \infty} 0;$$

1755  
1756 3. Phase 3,  $\mathbb{E} \|\nabla f(X_{\tilde{t}^3})\|_2^2 \leq \sqrt{\frac{\pi}{2}} \frac{\sigma_{\max} \eta L d \phi_t^2}{2 \phi_t^1} + \sqrt{\frac{\pi}{2}} \sigma_{\max} \frac{f(X_0) - f(X_*)}{\phi_t^1} \xrightarrow{t \rightarrow \infty} 0$ ;

1757  
1758 where  $\tilde{t}^1$ ,  $\tilde{t}^{(1,2)}$ ,  $\tilde{t}^{(2,2)}$ , and  $\tilde{t}^3$  are random times with distribution  $\frac{\eta_t}{\phi_t^1}$ .

1759  
1760  
1761 *Proof of Lemma C.22.* We prove each point by leveraging the shape of the law of  $X_t$  derived in  
1762 Lemma C.21:

1763  
1764 1. **Phase 1:**

1765 
$$d(f(X_t) - f(X_*)) = -\eta_t \nabla f(X_t) \text{sign}(\nabla f(X_t)) dt = -\eta_t \|\nabla f(X_t)\|_1 dt \quad (98)$$

1766  
1767 
$$= -\phi_t^1 \frac{\eta_t \|\nabla f(X_t)\|_1}{\phi_t^1} dt \quad (99)$$

1768  
1769 Therefore, by integrating over time and using the law of the unconscious statistician

1770  
1771 
$$\|\nabla f(X_{\tilde{t}^1})\|_1 \leq \frac{f(X_0) - f(X_*)}{\phi_t^1} \xrightarrow{t \rightarrow \infty} 0; \quad (100)$$

1772  
1773 2. **Phase 2:** By applying the Itô Lemma to  $f(X_t) - f(X_*)$  and that

1774  
1775 
$$\frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^- < \text{Erf} \left( \frac{\Sigma^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2}} \right) < \frac{m}{\sqrt{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) + \mathbf{q}^+, \quad (101)$$

1776  
1777 Similar to what we have shown above, we have that

1778  
1779 
$$d(f(X_t) - f(X_*)) \leq -\frac{m}{\sqrt{2}} \frac{1}{\sigma_{\max}} \eta_t \|\nabla f(X_t)\|_2^2 dt - \eta_t \hat{q} \|\nabla f(X_t)\|_1 dt \quad (102)$$

1780  
1781 
$$+\eta_t^2 \frac{\eta L d}{2} dt + \mathcal{O}(\text{Noise}). \quad (103)$$

Therefore, by integrating over time and using the law of the unconscious statistician we have

$$\frac{m}{\sqrt{2}} \mathbb{E} \|\nabla f(X_{\tilde{t}(1,2)})\|_2^2 + \hat{q} \sigma_{\max} \mathbb{E} \|\nabla f(X_{\tilde{t}(2,2)})\|_1 \leq \frac{\sigma_{\max}}{\phi_t^1} \left( f(X_0) - f(X_*) + \frac{\eta L d \phi_t^2}{2} \right) \xrightarrow{t \rightarrow \infty} 0; \quad (104)$$

3. **Phase 3:** By applying the Itô Lemma to  $f(X_t) - f(X_*)$ , we have that:

$$d(f(X_t) - f(X_*)) \leq -\sqrt{\frac{2}{\pi}} \frac{1}{\sigma_{\max}} \eta_t \|\nabla f(X_t)\|_2^2 dt + \mathcal{O}(\text{Noise}) + \eta_t^2 \frac{\eta L d}{2} dt \quad (105)$$

Therefore, by integrating over time and using the law of the unconscious statistician we have

$$\mathbb{E} \|\nabla f(X_{\tilde{t}^3})\|_2^2 \leq \sqrt{\frac{\pi}{2}} \frac{\sigma_{\max} \eta L d}{2} \frac{\phi_t^2}{\phi_t^1} + \sqrt{\frac{\pi}{2}} \sigma_{\max} \frac{f(X_0) - f(X_*)}{\phi_t^1} \xrightarrow{t \rightarrow \infty} 0. \quad (106)$$

□

**Lemma C.25.** *Under the assumptions of Lemma 3.5, for any step size scheduler  $\eta_t$  such that*

$$\int_0^\infty \eta_s ds = \infty \text{ and } \lim_{t \rightarrow \infty} \eta_t = 0 \implies \mathbb{E}[f(X_t) - f(X_*)] \xrightarrow{t \rightarrow \infty} 0. \quad (107)$$

*Proof of Lemma C.25.* For any scheduler  $\eta_k$  used in

$$x_{k+1} = x_k - \eta \eta_k \text{sign}(f_{\gamma_k}(x_k)), \quad (108)$$

the SDE of Phase 3 is

$$dX_t = -\sqrt{\frac{2}{\pi}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) \eta_t dt + \sqrt{\eta} \eta_t \sqrt{I_d - \frac{2}{\pi} \text{diag}(\Sigma^{-\frac{1}{2}} \nabla f(X_t))^2} dW_t. \quad (109)$$

Therefore, analogously to the calculations in Lemma C.22, we have that

$$\mathbb{E}[f(X_t) - f(X_*)] \leq \frac{f(X_0) - f(X_*) + \frac{\eta \mathcal{L}_\tau}{2} \int_0^t e^{2\mu \int_0^s \left( \sqrt{\frac{2}{\pi}} \frac{1}{\sigma_{\max}} \eta_l + \frac{\eta}{\pi} \frac{\mu}{\sigma_{\max}^2} \eta_l^2 \right) dl} \eta_s^2 ds}{e^{2\mu \int_0^t \left( \sqrt{\frac{2}{\pi}} \frac{1}{\sigma_{\max}} \eta_s + \frac{\eta}{\pi} \frac{\mu}{\sigma_{\max}^2} \eta_s^2 \right) ds}}. \quad (110)$$

Therefore, using l'Hôpital's rule we have that

$$\int_0^\infty \eta_s ds = \infty \text{ and } \lim_{t \rightarrow \infty} \eta_t = 0 \implies \mathbb{E}[f(X_t) - f(X_*)] \xrightarrow{t \rightarrow \infty} 0. \quad (111)$$

□

**Lemma C.26.** *Let  $H = \text{diag}(\lambda_1, \dots, \lambda_d)$  and  $M_t := e^{-2\left(\sqrt{\frac{2}{\pi}} \Sigma^{-\frac{1}{2}} H + \frac{\eta}{\pi} \Sigma^{-1} H^2\right)t}$ . Then,*

1.  $\mathbb{E}[X_t] = e^{-\sqrt{\frac{2}{\pi}} \Sigma^{-\frac{1}{2}} H t} X_0;$
2.  $\text{Var}[X_t] = \left( M_t - e^{-2\sqrt{\frac{2}{\pi}} \Sigma^{-\frac{1}{2}} H t} \right) X_0^2 + \frac{\eta}{2} \left( \sqrt{\frac{2}{\pi}} I_d + \frac{\eta}{\pi} H \Sigma^{-\frac{1}{2}} \right)^{-1} H^{-1} \Sigma^{\frac{1}{2}} (I_d - M_t).$

*Proof of Lemma C.26.* The proof is banal: The expected value derivation leverages the martingale property of the Brownian motion while that of the variance uses the Ito Isomerty. □

**Lemma C.27.** *Let  $H = \text{diag}(\lambda_1, \dots, \lambda_d)$ . Then,  $\mathbb{E} \left[ \frac{X_t^\top H X_t}{2} \right]$  is equal to*

$$\sum_{i=1}^d \frac{\lambda_i (X_0^i)^2}{2} e^{-2\lambda_i \left( \sqrt{\frac{2}{\pi}} \frac{1}{\sigma_i} + \frac{\lambda_i \eta}{\pi \sigma_i^2} \right) t} + \frac{\eta}{4 \left( \sqrt{\frac{2}{\pi}} \frac{1}{\sigma_i} + \frac{\lambda_i \eta}{\pi \sigma_i^2} \right)} \left( 1 - e^{-2\lambda_i \left( \sqrt{\frac{2}{\pi}} \frac{1}{\sigma_i} + \frac{\lambda_i \eta}{\pi \sigma_i^2} \right) t} \right). \quad (112)$$

1836 *Proof of Lemma C.27.* Since the matrix  $H$  is diagonal, we focus on a single component. We apply  
 1837 the Ito Lemma to  $\frac{\lambda_i(X_t^i)^2}{2}$ :

$$1839 \quad d\left(\frac{\lambda_i(X_t^i)^2}{2}\right) = -2\sqrt{\frac{2}{\pi}} \frac{\lambda_i}{\sigma_i} \frac{\lambda_i(X_t^i)^2}{2} dt + \frac{\eta\lambda_i}{2} dt - \frac{2\lambda_i^2\eta}{\pi\sigma_i^2} \frac{\lambda_i(X_t^i)^2}{2} dt + \mathcal{O}(\text{Noise}), \quad (113)$$

1841 which implies that

$$1842 \quad \mathbb{E}\left[\frac{\lambda_i(X_t^i)^2}{2}\right] = \frac{\lambda_i(X_0^i)^2}{2} e^{-2\left(\sqrt{\frac{2}{\pi}} \frac{\lambda_i}{\sigma_i} + \frac{\lambda_i^2\eta}{\pi\sigma_i^2}\right)t} + \frac{\eta}{4\left(\sqrt{\frac{2}{\pi}} \frac{1}{\sigma_i} + \frac{\lambda_i\eta}{\pi\sigma_i^2}\right)} \left(1 - e^{-2\left(\sqrt{\frac{2}{\pi}} \frac{\lambda_i}{\sigma_i} + \frac{\lambda_i^2\eta}{\pi\sigma_i^2}\right)t}\right). \quad (114)$$

1843 Therefore,

$$1844 \quad \mathbb{E}\left[\frac{X_t^\top H X_t}{2}\right] = \sum_{i=1}^d \frac{\lambda_i(X_0^i)^2}{2} e^{-2\lambda_i\left(\sqrt{\frac{2}{\pi}} \frac{1}{\sigma_i} + \frac{\lambda_i\eta}{\pi\sigma_i^2}\right)t} + \frac{\eta}{4\left(\sqrt{\frac{2}{\pi}} \frac{1}{\sigma_i} + \frac{\lambda_i\eta}{\pi\sigma_i^2}\right)} \left(1 - e^{-2\lambda_i\left(\sqrt{\frac{2}{\pi}} \frac{1}{\sigma_i} + \frac{\lambda_i\eta}{\pi\sigma_i^2}\right)t}\right). \quad (115)$$

1845  $\square$

1846 **Lemma C.28.** Under the assumptions of Corollary C.20, where  $\nabla f_\gamma(x) = \nabla f(x) + \sqrt{\Sigma}Z$ , we have  
 1847 that the dynamics of SignSGD in **Phase 3** is:

$$1848 \quad dX_t = -\sqrt{\frac{1}{2}} \Sigma^{-\frac{1}{2}} \nabla f(X_t) dt + \sqrt{\eta} \sqrt{I_d - \frac{1}{2} \text{diag}\left(\Sigma^{-\frac{1}{2}} \nabla f(X_t)\right)^2} dW_t. \quad (116)$$

1849 *Proof of lemma C.28.* We apply Eq. 44 with  $\nu = 2$  and linearly approximate  $\Xi(x)$  as  $|x| < 1$ , where  
 1850  $2\Xi(x) \sim \frac{x}{\sqrt{2}}$ .  $\square$

## 1851 C.5 ALTERNATIVE NOISE ASSUMPTIONS

1852 In this subsection, we report the consequences of assuming different noise structures. We do not  
 1853 provide the proofs as they mimic those of Corollary C.19 and Lemma C.21. We validate our results  
 1854 in Figure 7.

1855 **Assumption from (Ziyin et al., 2021)** As per Eq. (16) in Corollary 2 of (Ziyin et al., 2021), we  
 1856 take  $\Sigma := \sigma^2 f(x_*) \nabla^2 f(x_*)$ , where we added the constant  $\sigma^2$  as a parameter to control the scale  
 1857 of the noise and  $f(x_*) > 0$ . Under this assumption, we have that for  $Y_t := \frac{\nabla^2 f(x_*)^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2\nabla f(x_*)\sigma}}$  and  
 1858  $\mathcal{S}_d(X_t) := \mathbb{E}_\gamma[(\text{sign}(\nabla f_\gamma(X_t)))(\text{sign}(\nabla f_\gamma(X_t)))^\top]$ , Corollary C.19 becomes:

$$1859 \quad dX_t = -\text{Erf}(Y_t) dt + \sqrt{\eta} \sqrt{\mathcal{S}_d(X_t) - \text{Erf}(Y_t) \text{Erf}(Y_t)^\top} dW_t. \quad (117)$$

1860 As a consequence, Lemma C.21 becomes:

1861 **Lemma C.29.** Let  $f$  be  $\mu$ -strongly convex,  $\text{Tr}(\nabla^2 f(x)) \leq \mathcal{L}_\tau$ ,  $\lambda_{\max}$  be the largest eigenvalue of  
 1862  $\nabla^2 f(x_*)$ , and  $S_t := f(X_t) - f(x_*)$ . Then, during

- 1863 1. Phase 1, the loss will reach 0 before  $t_* = 2\sqrt{\frac{S_0}{\mu}}$  because  $S_t \leq \frac{1}{4}(\sqrt{\mu t} - 2\sqrt{S_0})^2$ ;
- 1864 2. Phase 2 with  $\Delta := \left(\frac{m}{\sqrt{2f(x_*)\sigma_{\max}\sqrt{\lambda_{\max}}} + \frac{\eta\mu m^2}{4f(x_*)\sigma_{\max}^2\lambda_{\max}}}\right)$ :  $\mathbb{E}[S_t] \leq S_0 e^{-2\mu\Delta t} +$   
 1865  $\frac{\eta}{2} \frac{(\mathcal{L}_\tau - \mu d\hat{q}^2)}{2\mu\Delta} (1 - e^{-2\mu\Delta t})$ ;
- 1866 3. Phase 3 with  $\Delta := \left(\sqrt{\frac{2}{\pi}} \frac{1}{\sqrt{f(x_*)\sigma_{\max}\sqrt{\lambda_{\max}}} + \frac{\eta}{\pi} \frac{\mu}{f(x_*)\sigma_{\max}^2\lambda_{\max}}}\right)$ :  $\mathbb{E}[S_t] \leq S_0 e^{-2\mu\Delta t} +$   
 1867  $\frac{\eta}{2} \frac{\mathcal{L}_\tau}{2\mu\Delta} (1 - e^{-2\mu\Delta t})$ .

**Assumption from (Wojtowysch, 2024)** (Wojtowysch, 2024) discusses two possible assumptions on  $\Sigma$ :  $\|\Sigma(x)\| \leq Cf(x)$  and  $\|\Sigma(x)\| \leq Cf(x) [1 + |x|^2]$ . As Section 2.4, they ultimately use  $\Sigma = Cf(x)I_d$ . Therefore, we take  $\Sigma := \sigma^2 f(x)I_d$ , where we changed the constant to  $\sigma^2$  to maintain consistency with the rest of our paper. Under this assumption, we have that for  $Y_t := \frac{\nabla f(X_t)}{\sqrt{2f(x)\sigma}}$ , Corollary C.19 becomes:

$$dX_t = -\text{Erf}(Y_t) dt + \sqrt{\eta} \sqrt{I_d - \text{diag}(\text{Erf}(Y_t))^2} dW_t. \quad (118)$$

As a consequence, Lemma C.21 becomes:

**Lemma C.30.** *Let  $f$  be  $\mu$ -strongly convex,  $\text{Tr}(\nabla^2 f(x)) \leq \mathcal{L}_\tau$ , and  $S_t := f(X_t) - f(X_*)$ . Then, during*

1. Phase 1, the loss will reach 0 before  $t_* = 2\sqrt{\frac{S_0}{\mu}}$  because  $S_t \leq \frac{1}{4} (\sqrt{\mu t} - 2\sqrt{S_0})^2$ ;

2. Phase 2 with  $\beta := \frac{\eta}{2} (\mathcal{L}_\tau - \mu d\hat{q}^2 - \frac{m^2 \mu^2}{\sigma^2})$  and  $\alpha := \frac{\sqrt{2m\mu}}{\sigma}$ ,

$$\mathbb{E}[S_t] \leq \frac{\beta^2 \left( \mathcal{W} \left( \frac{(\beta + \sqrt{S_0}\alpha)}{\beta} \exp \left( -\frac{\alpha^2 t - 2\sqrt{S_0}\alpha}{2\beta} - 1 \right) \right) + 1 \right)^2}{\alpha^2} \xrightarrow{t \rightarrow \infty} \frac{\beta^2}{\alpha^2}; \quad (119)$$

3. Phase 3 with  $\beta := \eta \left( \frac{\mathcal{L}_\tau}{2} - \frac{2\mu^2}{\pi\sigma^2} \right)$  and  $\alpha := 2\sqrt{\frac{2}{\pi}} \frac{\mu}{\sigma}$ ,

$$\mathbb{E}[S_t] \leq \frac{\beta^2 \left( \mathcal{W} \left( \frac{(\beta + \sqrt{S_0}\alpha)}{\beta} \exp \left( -\frac{\alpha^2 t - 2\sqrt{S_0}\alpha}{2\beta} - 1 \right) \right) + 1 \right)^2}{\alpha^2} \xrightarrow{t \rightarrow \infty} \frac{\beta^2}{\alpha^2}, \quad (120)$$

where  $\mathcal{W}$  is the Lambert  $\mathcal{W}$  function.

**Assumption from (Wu et al., 2022)** (Wu et al., 2022) proposes a novel structure of  $\Sigma$  as being aligned with the Fisher Information Matrix and proportional to the loss function. Consistently with this, we take  $\Sigma := \sigma^2 f(x) \nabla^2 f(x)$ , where we changed the constants to  $\sigma^2$  to maintain consistency with the rest of our paper. Under this assumption, we have that for  $Y_t := \frac{(\nabla^2 f(X_t))^{-\frac{1}{2}} \nabla f(X_t)}{\sqrt{2\nabla f(x)\sigma}}$  and  $S_d(X_t) := \mathbb{E}_\gamma[(\text{sign}(\nabla f_\gamma(X_t))(\text{sign}(\nabla f_\gamma(X_t)))^\top]$ , Corollary C.19 becomes:

$$dX_t = -\text{Erf}(Y_t) dt + \sqrt{\eta} \sqrt{S_d(X_t) - \text{Erf}(Y_t) \text{Erf}(Y_t)^\top} dW_t. \quad (121)$$

As a consequence, Lemma C.21 becomes:

**Lemma C.31.** *Let  $f$  be  $\mu$ -strongly convex,  $L$ -smooth,  $\text{Tr}(\nabla^2 f(x)) \leq \mathcal{L}_\tau$ , and  $S_t := f(X_t) - f(X_*)$ . Then, during*

1. Phase 1, the loss will reach 0 before  $t_* = 2\sqrt{\frac{S_0}{\mu}}$  because  $S_t \leq \frac{1}{4} (\sqrt{\mu t} - 2\sqrt{S_0})^2$ ;

2. Phase 2 with  $\beta := \frac{\eta}{2} (\mathcal{L}_\tau - \mu d\hat{q}^2 - \frac{m^2 \mu^2}{\sigma^2 L})$  and  $\alpha := \frac{\sqrt{2m\mu}}{\sqrt{L}\sigma}$ ,

$$\mathbb{E}[S_t] \leq \frac{\beta^2 \left( \mathcal{W} \left( \frac{(\beta + \sqrt{S_0}\alpha)}{\beta} \exp \left( -\frac{\alpha^2 t - 2\sqrt{S_0}\alpha}{2\beta} - 1 \right) \right) + 1 \right)^2}{\alpha^2} \xrightarrow{t \rightarrow \infty} \frac{\beta^2}{\alpha^2}; \quad (122)$$

3. Phase 3 with  $\beta := \eta \left( \frac{\mathcal{L}_\tau}{2} - \frac{2\mu^2}{\pi\sigma^2 L} \right)$  and  $\alpha := 2\sqrt{\frac{2}{\pi}} \frac{\mu}{\sqrt{L}\sigma}$ ,

$$\mathbb{E}[S_t] \leq \frac{\beta^2 \left( \mathcal{W} \left( \frac{(\beta + \sqrt{S_0}\alpha)}{\beta} \exp \left( -\frac{\alpha^2 t - 2\sqrt{S_0}\alpha}{2\beta} - 1 \right) \right) + 1 \right)^2}{\alpha^2} \xrightarrow{t \rightarrow \infty} \frac{\beta^2}{\alpha^2}, \quad (123)$$

where  $\mathcal{W}$  is the Lambert  $\mathcal{W}$  function.

1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997

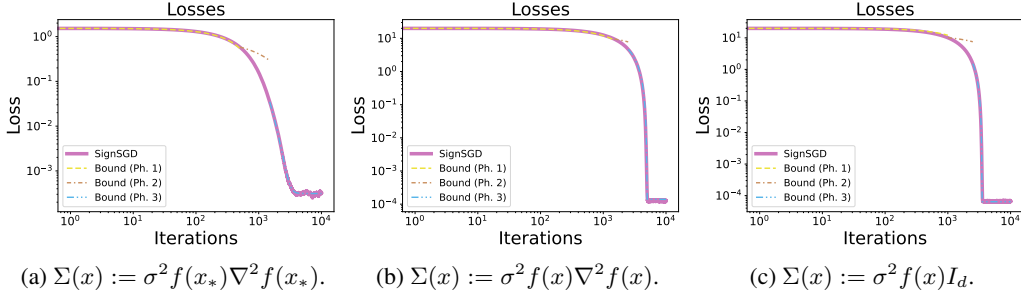


Figure 7: Empirical validation of the bounds for SignSGD derived for the noise structure:  $\Sigma(x) := \sigma^2 f(x_*) \nabla^2 f(x_*)$  (left),  $\Sigma(x) := \sigma^2 f(x) \nabla^2 f(x)$  (center),  $\Sigma(x) := \sigma^2 f(x) I_d$  (right). Each empirical validation has been carried out on strongly convex quadratic loss functions. The experiments are averaged over 500 runs.

**Assumption from (Paquette et al., 2024)** In this section, we adopt the notation of (Paquette et al., 2024). As per Eq. 5 of (Paquette et al., 2024), the loss function can be rewritten as

$$f(\theta) = \frac{1}{2} \langle D(W\theta - b), (W\theta - b) \rangle, \quad \text{where } D = \text{diag}(j^{-2\alpha}) \in \mathbb{R}^{v \times v}. \quad (124)$$

Without loss of generality, we define  $\phi := W\theta - b$ , which implies that

$$f(\theta) = \frac{\phi^\top D \phi}{2}, \quad \text{where } D = \text{diag}(j^{-2\alpha}) \in \mathbb{R}^{v \times v}, \text{ where } 1 \leq j \leq v. \quad (125)$$

The stochastic gradient is unbiased and its covariance is the well-known  $B\Sigma(\phi) = (\phi^\top D \phi)D + D\phi\phi^\top D = 2f(\phi)D + \nabla f(\phi)\nabla f(\phi)^\top$ , where  $B$  is the batch size. Under this assumption, we have that for  $Y_t := \frac{\sqrt{B}(\Sigma(\phi_t))^{-\frac{1}{2}} \nabla f(\phi_t)}{\sqrt{2}}$  and  $\mathcal{S}(\phi_t) = \mathbb{E}[(\text{Sign}(\nabla f_\gamma(\phi_t)))(\text{Sign}(\nabla f_\gamma(\phi_t)))^\top]$ , Corollary C.19 becomes:

$$d\phi_t = -\text{Erf}(Y_t) dt + \sqrt{\eta} \sqrt{\mathcal{S}(\phi_t) - \text{Erf}(Y_t) \text{Erf}(Y_t)^\top} dW_t. \quad (126)$$

As a consequence, Lemma C.21 becomes:

**Lemma C.32.** Let  $f$  be as above,  $f_t := f(\phi_t)$ , and  $\mathcal{L}_\tau := \text{Tr}(D)$ . Let  $\mu$  be the minimum eigenvalue of  $D$ , and  $L$  be its maximum one. Then, during

1. Phase 1, the loss will reach 0 before  $t_* = 2\sqrt{\frac{S_0}{\mu}}$  because  $f_t \leq \frac{1}{4}(\sqrt{\mu t} - 2\sqrt{f_0})^2$ ;

2. Phase 2 with  $\beta := \frac{\eta}{2}\mathcal{L}_\tau$  and  $\alpha := \frac{m\mu\sqrt{B}}{\sqrt{2L}}$ ,

$$\mathbb{E}[S_t] \leq \frac{\beta^2 \left( \mathcal{W} \left( \frac{(\beta + \sqrt{S_0}\alpha)}{\beta} \exp \left( -\frac{\alpha^2 t - 2\sqrt{S_0}\alpha}{2\beta} - 1 \right) \right) + 1 \right)^2}{\alpha^2} \xrightarrow{t \rightarrow \infty} \frac{\beta^2}{\alpha^2}; \quad (127)$$

3. Phase 3 with  $\beta := \frac{\eta}{2}\mathcal{L}_\tau$  and  $\alpha := \sqrt{\frac{2}{\pi}} \frac{\mu\sqrt{B}}{\sqrt{L}}$ ;

$$\mathbb{E}[S_t] \leq \frac{\beta^2 \left( \mathcal{W} \left( \frac{(\beta + \sqrt{S_0}\alpha)}{\beta} \exp \left( -\frac{\alpha^2 t - 2\sqrt{S_0}\alpha}{2\beta} - 1 \right) \right) + 1 \right)^2}{\alpha^2} \xrightarrow{t \rightarrow \infty} \frac{\beta^2}{\alpha^2}, \quad (128)$$

where  $\mathcal{W}$  is the Lambert  $\mathcal{W}$  function.

See Figure 8 for an empirical validation.

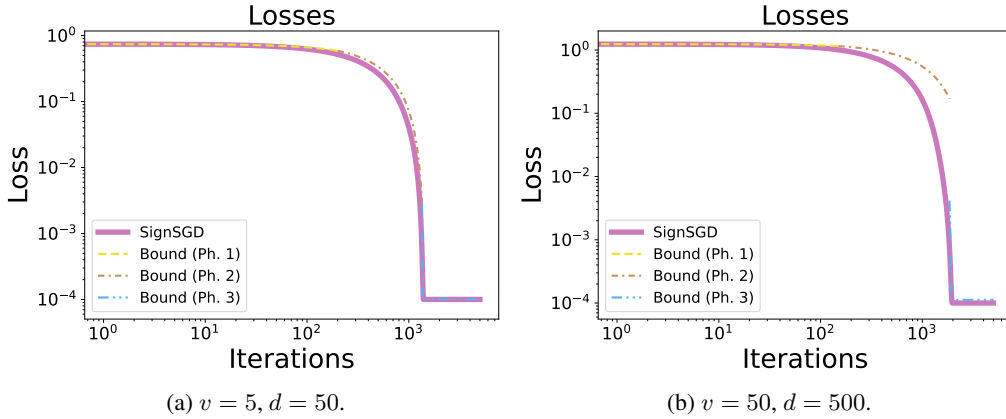


Figure 8: Empirical validation of the bounds for SignSGD derived in Lemma C.32: In both experiments,  $\alpha = 0.25$ ,  $\beta = 2$ ,  $\eta = 0.001$ ,  $B = 256$ ,  $N = 10000$ , and trajectories are averaged over 500 runs.

## C.6 FORMAL DERIVATION - RMSPROP

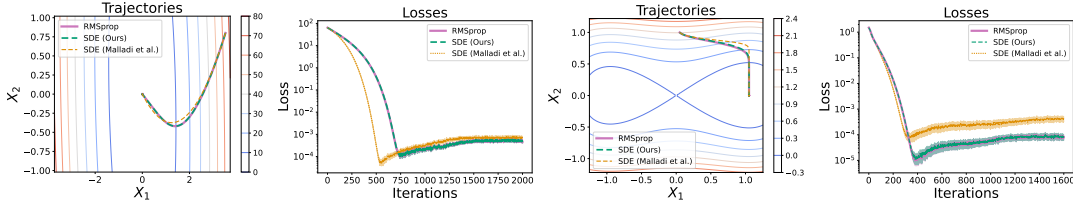


Figure 9: The first two subfigures on the left compare our SDE, that from Malladi et al. (2022), and RMSprop in terms of trajectories and  $f(x)$ , respectively, for a convex quadratic function. The other subfigures do the same for an embedded saddle and one clearly observes that our derived SDE better matches RMSprop.

In this subsection, we provide our formal derivation of an SDE model for RMSprop. Let us consider the stochastic process  $L_t := (X_t, V_t) \in \mathbb{R}^d \times \mathbb{R}^d$  defined as the solution of

$$dX_t = -P_t^{-1}(\nabla f(X_t)dt + \sqrt{\eta}\Sigma(X_t)^{\frac{1}{2}}dW_t) \quad (129)$$

$$dV_t = \rho((\nabla f(X_t))^2 + \text{diag}(\Sigma(X_t)) - V_t)dt, \quad (130)$$

where  $\beta = 1 - \eta\rho$ ,  $\rho = \mathcal{O}(1)$ , and  $P_t := \text{diag}(V_t)^{\frac{1}{2}} + \epsilon I_d$ .

*Remark C.33.* We observe that the term in blue is the only difference w.r.t. the SDE derived in (Malladi et al., 2022) (see Theorem D.2): This is extremely relevant when the gradient size is not negligible. Figure 9 shows the comparison between our SDE, the one derived in (Malladi et al., 2022), and RMSprop itself: It is clear that even on simple landscapes, our SDE matches the algorithm much better. Importantly, one can observe that the SDE derived in (Malladi et al., 2022) is only slightly worse than ours at the end of the dynamics: As we show in Lemma C.37, Theorem D.2 is a corollary of Theorem C.34 when  $\nabla f(x) = \mathcal{O}(\sqrt{\eta})$ : It only describes the dynamics where the gradient is vanishing. In Figure 10, we compare the two SDEs in question with RMSprop on an MLP, a CNN, a ResNet, and a Transformer: Our SDE exhibits a superior description of the dynamics.

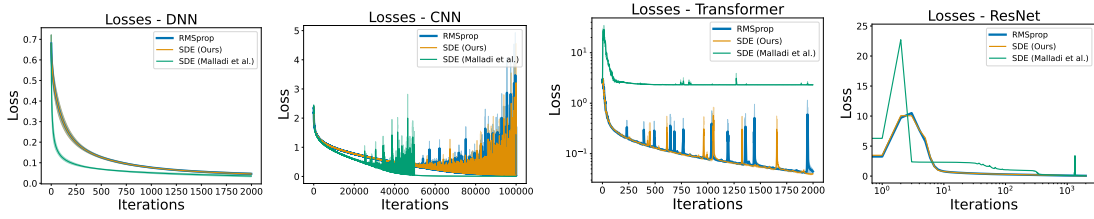


Figure 10: We compare our SDE, that from Malladi et al. (2022), and RMSprop in terms of  $f(x)$ : The first is an MLP on the Breast Cancer dataset, the second a CNN on MNIST, the third a Transformer on MNIST, and the last a ResNet on CIFAR-10: Ours match the algorithms better.

The following theorem guarantees that such a process is a 1-order SDE of the discrete-time algorithm of RMSprop

$$x_{k+1} = x_k - \eta \frac{\nabla f_{\gamma_k}(x_k)}{\sqrt{v_{k+1} + \epsilon I_d}} \quad (131)$$

$$v_{k+1} = \beta v_k + (1 - \beta) (\nabla f_{\gamma_k}(x_k))^2 \quad (132)$$

with  $(x_0, v_0) \in \mathbb{R}^d \times \mathbb{R}^d$ ,  $\eta \in \mathbb{R}^{>0}$  is the step size,  $\beta = 1 - \rho\eta$  for  $\rho = \mathcal{O}(1)$ , the mini-batches  $\{\gamma_k\}$  are modelled as i.i.d. random variables uniformly distributed on  $\{1, \dots, N\}$ , and of size  $B \geq 1$ .

**Theorem C.34** (Stochastic modified equations). *Let  $0 < \eta < 1, T > 0$  and set  $N = \lfloor T/\eta \rfloor$ . Let  $l_k := (x_k, v_k) \in \mathbb{R}^d \times \mathbb{R}^d, 0 \leq k \leq N$  denote a sequence of RMSprop iterations defined by Eq. 131. Consider the stochastic process  $L_t$  defined in Eq. 129 and fix some test function  $g \in G$  and suppose that  $g$  and its partial derivatives up to order 6 belong to  $G$ . Then, under Assumption C.3 and  $\rho = \mathcal{O}(1)$  there exists a constant  $C > 0$  independent of  $\eta$  such that for all  $k = 0, 1, \dots, N$ , we have*

$$|\mathbb{E}g(L_{k\eta}) - \mathbb{E}g(l_k)| \leq C\eta.$$

*That is, the SDE 129 is an order 1 weak approximation of the RMSprop iterations 131.*

*Proof.* The proof is virtually identical to that of Theorem C.16. Therefore, we only report the key steps necessary to conclude the thesis. First of all, we observe that since  $\beta = 1 - \rho\eta$

$$v_{k+1} - v_k = -\eta\rho \left( v_k - (\nabla f_{\gamma_k}(x_k))^2 \right). \quad (133)$$

Then,

$$\frac{1}{\sqrt{v_{k+1}}} = \sqrt{\frac{v_k}{v_{k+1}}} \frac{1}{\sqrt{v_k}} = \sqrt{\frac{v_{k+1} + \mathcal{O}(\eta)}{v_{k+1}}} \frac{1}{\sqrt{v_k}} = \sqrt{1 + \frac{\mathcal{O}(\eta)}{v_{k+1}}} \sqrt{\frac{1}{v_k}} \sim \sqrt{\frac{1}{v_k}} (1 + \mathcal{O}(\eta)). \quad (134)$$

Therefore, we work with the following algorithm as all the approximations below only carry an additional error of order  $\mathcal{O}(\eta^2)$ , which we can ignore. Therefore, we have that

$$x_{k+1} - x_k = -\eta \frac{\nabla f_{\gamma_k}(x_k)}{\sqrt{v_k + \epsilon I_d}} \quad (135)$$

$$v_k - v_{k-1} = -\eta\rho \left( v_{k-1} - (\nabla f_{\gamma_{k-1}}(x_{k-1}))^2 \right). \quad (136)$$

Therefore, if  $\nabla f_{\gamma_j}(x_j) = \nabla f(x_j) + Z_j(x_j)$ ,  $\mathbb{E}[Z_j(x_j)] = 0$ , and  $Cov(Z_j(x_j)) = \Sigma(x_j)$

1.  $\mathbb{E}[x_{k+1} - x_k] = -\eta \text{diag}(v_k + \epsilon I_d)^{-\frac{1}{2}} \nabla f(x_k)$ ;
2.  $\mathbb{E}[v_k - v_{k-1}] = \eta\rho \left[ (\nabla f(x_{k-1}))^2 + \text{diag}(\Sigma(x_k)) - v_{k-1} \right]$ .

Then, we have that if  $\Phi_k := \frac{\nabla f(x_k)}{\sqrt{v_k + \epsilon I_d}} - \frac{\nabla f_{\gamma_k}(x_k)}{\sqrt{v_k + \epsilon I_d}}$

1.

$$\mathbb{E}[(x_{k+1} - x_k)(x_{k+1} - x_k)^\top] = \mathbb{E}[(x_{k+1} - x_k)]\mathbb{E}[(x_{k+1} - x_k)^\top] \quad (137)$$

$$+ \eta^2 \mathbb{E}[(\Phi_k)(\Phi_k)^\top] \quad (138)$$

$$= \mathbb{E}[(x_{k+1} - x_k)]\mathbb{E}[(x_{k+1} - x_k)^\top] \quad (139)$$

$$+ \eta^2 (\text{diag}(v_k) + \epsilon I_d)^{-1} \Sigma(x_k); \quad (140)$$

$$2. \mathbb{E}[(v_k - v_{k-1})(v_k - v_{k-1})^\top] = \mathbb{E}[(v_k - v_{k-1})]\mathbb{E}[(v_k - v_{k-1})^\top] + \mathcal{O}(\rho\eta^2);$$

$$3. \mathbb{E}[(x_{k+1} - x_k)(v_k - v_{k-1})^\top] = \mathbb{E}[(x_{k+1} - x_k)]\mathbb{E}[(v_k - v_{k-1})^\top] + 0.$$

*Remark C.35.* Let us remember that by assumption,  $\nabla f(x)$  and  $\sqrt{\Sigma}(x)$  are Lipschitz, grow at most affinely, and are in  $G$  together with their derivative. Therefore, the drift and diffusion terms of the SDE governing  $X_t$  are the ratio between regular functions and a uniformly lower bounded process. Therefore, they are in turn regular, modulo dividing by  $\sqrt{V_t + \epsilon_V^2} + \epsilon$  s.t.  $\epsilon_V^2 \sim 0$  rather than by  $\sqrt{V_t} + \epsilon$  (See Bock and Weiß (2021) as they experimentally verify that this has no impact on the performance of the optimizer). Regarding the ODE governing  $V_t$ ,  $\Sigma(X_t)$  is Lipschitz because  $\sqrt{\Sigma}(X_t)$  is bounded and Lipschitz. Additionally, it is smooth, and with affine growth. On top of this, we need the term  $(\nabla f(x))^2$  to be Lipschitz and of affine growth, which is a consequence of assuming bounded gradients as often done in the literature on the convergence of RMSprop and Adam: Among many, see (Luo et al., 2019; Défossez et al., 2022; Guo et al., 2021; Huang et al., 2021) together with the discussion in Section 2.1 of Shi and Li (2021). Alternatively, exactly as done in Theorem 9 of Li et al. (2019), one can regularize the drifts and the diffusion terms with mollifiers on a sufficiently large compact Reddi et al. (2018), which automatically implies that drift and diffusion coefficients satisfy all necessary regularity conditions. Importantly, one needs to then send the mollification parameter  $\epsilon$  to 0 to conclude our statement. Therefore, we the SDE of RMSprop for  $P_t := \text{diag}(V_t)^{\frac{1}{2}} + \epsilon I_d$  is

$$dX_t = -P_t^{-1}(\nabla f(X_t)dt + \sqrt{\eta}\Sigma(X_t)^{\frac{1}{2}}dW_t) \quad (141)$$

$$dV_t = \rho((\nabla f(X_t))^2 + \text{diag}(\Sigma(X_t)) - V_t)dt. \quad (142)$$

□

*Remark C.36.* In all the following results, the reader will notice that all the drifts, diffusion terms, and noise assumptions are selected to guarantee that the SDE we derived for RMSprop is indeed a 1 weak approximation for RMSprop even without the mollification argument. Importantly, our analysis of RMSprop focuses on its behavior at convergence, i.e.  $(\nabla f(x))^2 = \mathcal{O}(\eta)$ . Therefore, there is no need to assume bounded gradients or a compact domain.

**Lemma C.37.** *If  $(\nabla f(x))^2 = \mathcal{O}(\eta)$ , Theorem D.2 is a Corollary of Theorem C.34.*

*Proof.* In the proof of Theorem C.34, one drops the term  $\eta(\nabla f(x))^2$  as it is of order  $\eta^2$ . □

**Corollary C.38.** *Under the assumptions of Theorem C.34 with  $\Sigma(x) = \sigma^2 I_d$ ,  $\tilde{\eta} = \kappa\eta$ ,  $\tilde{B} = B\delta$ , and  $\tilde{\rho} = \alpha\rho$ ,*

$$dX_t = \kappa \text{diag}(V_t)^{-\frac{1}{2}} \left( -\nabla f(X_t)dt + \frac{1}{\sqrt{\delta}} \sqrt{\frac{\eta}{B}} \sigma I_d dW_t \right) \quad (143)$$

$$dV_t = \frac{\alpha}{\kappa} \rho \left( (\nabla f(X_t))^2 + \frac{\sigma^2}{B\delta} \mathbf{1} - V_t \right) dt. \quad (144)$$

**Lemma C.39** (Scaling Rule at Convergence). *Under the assumptions of Corollary C.38,  $f$  is  $\mu$ -strongly convex,  $\text{Tr}(\nabla^2 f(x)) \leq \mathcal{L}_\tau$ , and  $(\nabla f(x))^2 = \mathcal{O}(\eta)$ , the asymptotic dynamics of the iterates of RMSprop satisfies the classic scaling rule  $\kappa = \sqrt{\delta}$  because*

$$\mathbb{E}[f(X_t) - f(X_*)] \stackrel{t \rightarrow \infty}{\leq} \frac{\eta\sigma\mathcal{L}_\tau}{4\mu\sqrt{B}} \frac{\kappa}{\sqrt{\delta}}. \quad (145)$$



By enforcing that the speed of  $V_t$  matches that of  $X_t$ , one needs  $\tilde{\rho} = \kappa^2\rho$ , which implies  $\tilde{\beta} = 1 - \kappa^2(1 - \beta)$ .

*Proof of Lemma C.39.* In order to recover the scaling of  $\beta$ , we enforce that the rate at which  $V_t$  converges to its limit matches the speed of  $X_t$ : We need  $\tilde{\rho} = \kappa^2\rho$ , which recovers the classic scaling  $\tilde{\beta} = 1 - \kappa^2(1 - \beta)$ . Additionally, since  $(\nabla f(x))^2 = \mathcal{O}(\eta)$  we have that

$$dX_t = \kappa \text{diag}(V_t)^{-\frac{1}{2}} \left( -\nabla f(X_t)dt + \frac{1}{\sqrt{\delta}} \sqrt{\frac{\eta}{B}} \sigma I_d dW_t \right) \quad (146)$$

$$dV_t = \kappa\rho \left( \frac{\sigma^2}{B\delta} \mathbf{1} - V_t \right) dt. \quad (147)$$

Therefore,  $V_t \xrightarrow{t \rightarrow \infty} \frac{\sigma^2}{B\delta} \mathbf{1}$ , meaning that under these conditions:

$$dX_t = -\frac{\sqrt{B\delta}\kappa}{\sigma} \nabla f(X_t)dt + \kappa\sqrt{\eta} I_d dW_t, \quad (148)$$

which satisfies the following for  $\mu$ -strongly convex functions

$$d\mathbb{E}[f(X_t) - f(X_*)] \leq -2\kappa\mu \frac{\sqrt{B\delta}}{\sigma} \mathbb{E}[f(X_t) - f(X_*)]dt + \frac{\kappa^2\eta\mathcal{L}_\tau}{2} dt, \quad (149)$$

meaning that  $\mathbb{E}[f(X_t) - f(X_*)] \leq \frac{\eta\mathcal{L}_\tau}{4\mu\sqrt{B}} \frac{\kappa}{\sqrt{\delta}}$ .

Since the asymptotic the loss is  $\frac{\eta}{2} \frac{\mathcal{L}_\tau\sigma}{2\mu\sqrt{B}} \frac{\kappa}{\sqrt{\delta}}$  does not depend on  $\kappa$  and  $\delta$  if  $\frac{\kappa}{\sqrt{\delta}} = 1$ , we recover the classic scaling rule.  $\square$

**Remark:** Under the same conditions, SGD satisfies

$$dX_t = -\kappa \nabla f(X_t)dt + \kappa \frac{1}{\sqrt{\delta}} \sqrt{\frac{\eta}{B}} \sigma I_d dW_t \quad (150)$$

and therefore

$$\mathbb{E}[f(X_t) - f(X_*)] \leq (f(X_0) - f(X_*))e^{-2\mu\kappa t} + \frac{\eta}{2} \frac{\mathcal{L}_\tau\sigma^2}{2\mu B} \frac{\kappa}{\delta} (1 - e^{-2\mu\kappa t}), \quad (151)$$

meaning that asymptotically the loss is  $\frac{\eta}{2} \frac{\mathcal{L}_\tau\sigma^2}{2\mu B} \frac{\kappa}{\delta}$  which does not depend on  $\kappa$  and  $\delta$  if  $\frac{\kappa}{\delta} = 1$ .

**Lemma C.40.** For  $f(x) := \frac{x^\top H x}{2}$ , the stationary distribution of RMSprop is  $(\mathbb{E}[X_\infty], \text{Cov}(X_\infty)) = \left(0, \frac{\eta}{2} \Sigma^{\frac{1}{2}} H^{-1}\right)$ .

*Proof.* As  $(\nabla f(x))^2 = \mathcal{O}(\eta)$  and  $t \rightarrow \infty$ , we have

$$dX_t = -\Sigma^{-\frac{1}{2}} H X_t dt + \sqrt{\eta} I_d dW_t \quad (152)$$

which implies that

$$X_t = e^{-\Sigma^{-\frac{1}{2}} H t} \left( X_0 + \sqrt{\eta} \int_0^t e^{\Sigma^{-\frac{1}{2}} H s} dW_s \right). \quad (153)$$

The thesis follows from the martingale property of Brownian motion and the Itô isometry.  $\square$

## C.7 RMSPROPW

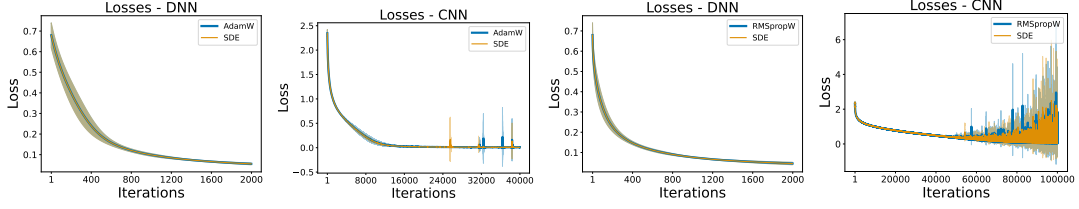


Figure 11: The first two represent the comparison between AdamW and its SDE in terms of  $f(x)$ . The other two do the same for RMSpropW. In both cases, the first is an MLP on the Breast Cancer Dataset and the second a CNN on MNIST: Our SDEs match the respective optimizers.

In this subsection, we derive the SDE of RMSpropW defined as

$$x_{k+1} = x_k - \eta \frac{\nabla f_{\gamma_k}(x_k)}{\sqrt{v_{k+1} + \epsilon I_d}} - \eta \theta x_k \quad (154)$$

$$v_{k+1} = \beta v_k + (1 - \beta) (\nabla f_{\gamma_k}(x_k))^2 \quad (155)$$

with  $(x_0, v_0) \in \mathbb{R}^d \times \mathbb{R}^d$ ,  $\eta \in \mathbb{R}^{>0}$  is the step size,  $\beta = 1 - \rho\eta$  for  $\rho = \mathcal{O}(1)$ ,  $\theta > 0$ , the mini-batches  $\{\gamma_k\}$  are modelled as i.i.d. random variables uniformly distributed on  $\{1, \dots, N\}$ , and of size  $B \geq 1$ .

**Theorem C.41.** *Under the same assumptions as Theorem C.34, the SDE of RMSpropW is*

$$dX_t = -P_t^{-1}(\nabla f(X_t)dt + \sqrt{\eta}\Sigma(X_t)^{\frac{1}{2}}dW_t) - \theta X_t dt \quad (156)$$

$$dV_t = \rho((\nabla f(X_t))^2 + \text{diag}(\Sigma(X_t)) - V_t)dt, \quad (157)$$

where  $\beta = 1 - \eta\rho$ ,  $\rho = \mathcal{O}(1)$ ,  $\theta > 0$ , and  $P_t := \text{diag}(V_t)^{\frac{1}{2}} + \epsilon I_d$ .

*Proof.* The proof is the same as the of Theorem C.34 and the only difference is that  $\eta\theta x_k$  is approximated with  $\theta X_t dt$ .  $\square$

Figure 4 and Figure 11 validate this result on a variety of architectures and datasets.

**Remark C.42.** See Remark C.35 and Remark C.36 for a discussion on the regularity of the SDE derived in Theorem C.41.

**Corollary C.43.** *Under the assumptions of Theorem C.41 with  $\Sigma(x) = \sigma^2 I_d$ ,  $\tilde{\eta} = \kappa\eta$ ,  $\tilde{B} = B\delta$ , and  $\tilde{\rho} = \alpha\rho$ , and  $\tilde{\theta} = \xi\theta$ ,*

$$dX_t = \kappa \text{diag}(V_t)^{-\frac{1}{2}} \left( -\nabla f(X_t)dt + \frac{1}{\sqrt{\delta}} \sqrt{\frac{\eta}{B}} \sigma I_d dW_t \right) - \xi\theta\kappa X_t dt \quad (158)$$

$$dV_t = \frac{\alpha}{\kappa} \rho \left( (\nabla f(X_t))^2 + \frac{\sigma^2}{B\delta} \mathbf{1} - V_t \right) dt. \quad (159)$$

**Lemma C.44** (Scaling Rule at Convergence). *Under the assumptions of Corollary C.43,  $f$  is  $\mu$ -strongly convex and  $L$ -smooth,  $\text{Tr}(\nabla^2 f(x)) \leq \mathcal{L}_\tau$ , and  $(\nabla f(x))^2 = \mathcal{O}(\eta)$ , the asymptotic dynamics of the iterates of RMSpropW satisfies the novel scaling rule if  $\kappa = \sqrt{\delta}$  and  $\xi = \kappa$  because*

$$\mathbb{E}[f(X_t) - f(X_*)] \stackrel{t \rightarrow \infty}{\leq} \frac{\eta \mathcal{L}_\tau \sigma L}{2} \frac{\kappa}{2\mu\sqrt{B\delta}L + \sigma\xi\theta(L + \mu)}. \quad (160)$$

By enforcing that the speed of  $V_t$  matches that of  $X_t$ , one needs  $\tilde{\rho} = \kappa^2\rho$ , which implies  $\tilde{\beta} = 1 - \kappa^2(1 - \beta)$ .

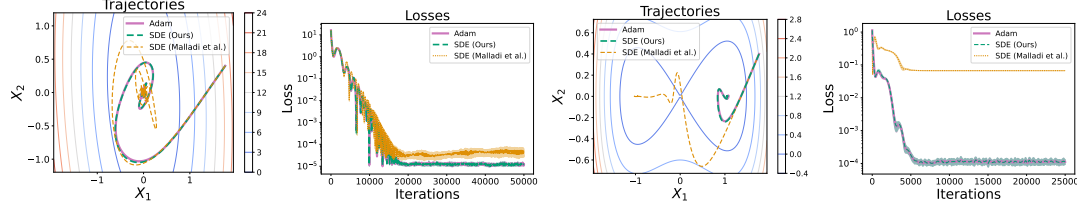


Figure 12: The first two on the left compare our SDE, that from Malladi et al. (2022), and Adam in terms of trajectories and  $f(x)$ , respectively, for a convex quadratic function. The others do the same for an embedded saddle: Ours clearly matches Adam better.

*Proof of Lemma C.44.* In order to recover the scaling of  $\beta$ , we enforce that the rate at which  $V_t$  converges to its limit matches the speed of  $X_t$ : We need  $\tilde{\rho} = \kappa^2 \rho$ , which recovers the classic scaling  $\tilde{\beta} = 1 - \kappa^2(1 - \beta)$ . Additionally, since  $(\nabla f(x))^2 = \mathcal{O}(\eta)$  we have that

$$dX_t = \kappa \text{diag}(V_t)^{-\frac{1}{2}} \left( -\nabla f(X_t) dt + \frac{1}{\sqrt{\delta}} \sqrt{\frac{\eta}{B}} \sigma I_d dW_t \right) - \kappa \xi \theta X_t dt \quad (161)$$

$$dV_t = \kappa \rho \left( \frac{\sigma^2}{B\delta} \mathbf{1} - V_t \right) dt. \quad (162)$$

Therefore,  $V_t \xrightarrow{t \rightarrow \infty} \frac{\sigma^2}{B\delta} \mathbf{1}$ , meaning that under these conditions:

$$dX_t = -\frac{\sqrt{B\delta}\kappa}{\sigma} \nabla f(X_t) dt + \kappa \sqrt{\eta} I_d dW_t - \kappa \xi \theta X_t dt, \quad (163)$$

which satisfies the following for  $\mu$ -strongly convex and  $L$ -smooth functions

$$d\mathbb{E}[f(X_t) - f(X_*)] \leq \kappa \left( 2\mu \frac{\sqrt{B\delta}}{\sigma} + \xi \theta \left( 1 + \frac{\mu}{L} \right) \right) \mathbb{E}[f(X_t) - f(X_*)] dt + \frac{\kappa^2 \eta \mathcal{L}_\tau}{2} dt, \quad (164)$$

meaning that  $\mathbb{E}[f(X_t) - f(X_*)] \leq \frac{\eta \mathcal{L}_\tau \sigma L}{2} \frac{\kappa}{2\mu\sqrt{B\delta}L + \sigma\xi\theta(L+\mu)}$ .

Since the asymptotic the loss  $\frac{\eta \mathcal{L}_\tau \sigma L}{2} \frac{\kappa}{2\mu\sqrt{B\delta}L + \sigma\xi\theta(L+\mu)}$  does not depend on  $\kappa$  and  $\delta$  and  $\xi$  if  $\kappa = \xi = \sqrt{\delta}$ , we recover the novel scaling rule.  $\square$

**Lemma C.45.** For  $f(x) := \frac{x^\top H x}{2}$ , the stationary distribution of RMSpropW is  $(\mathbb{E}[X_\infty], \text{Cov}(X_\infty)) = \left( 0, \frac{\eta}{2} (H \Sigma^{-\frac{1}{2}} + \theta I_d)^{-1} \right)$ .

*Proof.* As  $(\nabla f(x))^2 = \mathcal{O}(\eta)$  and  $t \rightarrow \infty$ , we have

$$dX_t = -\Sigma^{-\frac{1}{2}} H X_t dt + \sqrt{\eta} I_d dW_t - \theta X_t dt \quad (165)$$

which implies that

$$X_t = e^{-(\Sigma^{-\frac{1}{2}} H + \theta I_d)t} \left( X_0 + \sqrt{\eta} \int_0^t e^{(\Sigma^{-\frac{1}{2}} H + \theta I_d)s} dW_s \right). \quad (166)$$

The thesis follows from the martingale property of Brownian motion and the Itô isometry.  $\square$

## C.8 FORMAL DERIVATION - ADAM

In this subsection, we provide our formal derivation of an SDE model for Adam. Let us consider the stochastic process  $L_t := (X_t, M_t, V_t) \in \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}^d$  defined as the solution of

$$dX_t = -\frac{\sqrt{\iota_2(t)}}{\iota_1(t)} P_t^{-1} (M_t + \eta \rho_1 (\nabla f(X_t) - M_t)) dt \quad (167)$$

$$dM_t = \rho_1 (\nabla f(X_t) - M_t) dt + \sqrt{\eta} \rho_1 \Sigma^{1/2}(X_t) dW_t \quad (168)$$

$$dV_t = \rho_2 ((\nabla f(X_t))^2 + \text{diag}(\Sigma(X_t)) - V_t) dt, \quad (169)$$

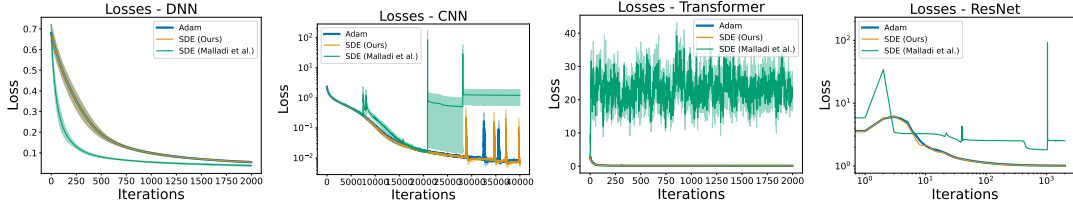


Figure 13: We compare our SDE, that from Malladi et al. (2022), and Adam in terms of  $f(x)$ : The first is an MLP on the Breast Cancer dataset, the second a CNN on MNIST, the third a Transformer on MNIST, and the last a ResNet on CIFAR-10: Ours match the algorithms better.

where  $\beta_i = 1 - \eta\rho_i$ ,  $\iota_i(t) = 1 - e^{-\rho_i t}$ ,  $\rho_1 = \mathcal{O}(\eta^{-\zeta})$  s.t.  $\zeta \in (0, 1)$ ,  $\rho_2 = \mathcal{O}(1)$ ,  $t > t_0$ , and  $P_t = \text{diag} \sqrt{V_t} + \epsilon \sqrt{\iota_2(t)} I_d$ .

*Remark C.46.* The terms in purple and in blue are the two differences w.r.t. that of (Malladi et al., 2022) which is reported in Theorem D.5. The first appears because we assume realistic values of  $\beta_1$  while the second appears because we allow the gradient size to be non-negligible. For two simple landscapes, Figure 12 compares our SDE and that of Malladi et al. (2022) with Adam: In both cases, the first part of the dynamics is perfectly represented only by our SDE. While the discrepancy between the SDE of (Malladi et al., 2022) and Adam is asymptotically negligible in the convex setting, we observe that in the non-convex case, it converges to a different local minimum than ours and of Adam. Finally, Theorem D.5 is a corollary of ours when  $(\nabla f(x))^2 = \mathcal{O}(\eta)$  and  $\rho_1 = \mathcal{O}(1)$ : It only describes the dynamics where the gradient to noise ratio is vanishing and only for unrealistic values of  $\beta_1 = 1 - \eta\rho_1$ . In Figure 13, we compare the dynamics of our SDE, that of Malladi et al. (2022), and Adam on an MLP, a CNN, a ResNet, and a Transformer. One can clearly see that our SDE more accurately captures the dynamics. Details on these experiments are in Appendix F.

The following theorem guarantees that such a process is a 1-order SDE of the discrete-time algorithm of Adam

$$v_{k+1} = \beta_2 v_k + (1 - \beta_2) (\nabla f_{\gamma_k}(x_k))^2 \quad (170)$$

$$m_{k+1} = \beta_1 m_k + (1 - \beta_1) \nabla f_{\gamma_k}(x_k) \quad (171)$$

$$\hat{m}_k = m_k (1 - \beta_1^k)^{-1} \quad (172)$$

$$\hat{v}_k = v_k (1 - \beta_2^k)^{-1} \quad (173)$$

$$x_{k+1} = x_k - \eta \frac{\hat{m}_{k+1}}{\sqrt{\hat{v}_{k+1} + \epsilon I_d}}, \quad (174)$$

with  $(x_0, m_0, v_0) \in \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}^d$ ,  $\eta \in \mathbb{R}^{>0}$  is the step size,  $\beta_i = 1 - \rho_i \eta$  for  $\rho_1 = \mathcal{O}(\eta^{-\zeta})$  s.t.  $\zeta \in (0, 1)$ ,  $\rho_2 = \mathcal{O}(1)$ , the mini-batches  $\{\gamma_k\}$  are modelled as i.i.d. random variables uniformly distributed on  $\{1, \dots, N\}$ , and of size  $B \geq 1$ .

**Theorem C.47** (Stochastic modified equations). *Let  $0 < \eta < 1, T > 0$  and set  $N = \lfloor T/\eta \rfloor$ . Let  $l_k := (x_k, m_k, v_k) \in \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}^d, 0 \leq k \leq N$  denote a sequence of Adam iterations defined by Eq. 170. Consider the stochastic process  $L_t$  defined in Eq. 167 and fix some test function  $g \in G$  and suppose that  $g$  and its partial derivatives up to order 6 belong to  $G$ . Then, under Assumption C.3  $\rho_1 = \mathcal{O}(\eta^{-\zeta})$  s.t.  $\zeta \in (0, 1)$ , while  $\rho_2 = \mathcal{O}(1)$ , there exists a constant  $C > 0$  independent of  $\eta$  such that for all  $k = 0, 1, \dots, N$ , we have*

$$|\mathbb{E}g(L_{k\eta}) - \mathbb{E}g(l_k)| \leq C\eta.$$

*That is, the SDE 167 is an order 1 weak approximation of the Adam iterations 170 for  $t > t_0$ .*

*Proof.* The proof is virtually identical to that of Theorem C.16. Therefore, we only report the key steps necessary to conclude the thesis. First of all, we observe that since  $\beta_1 = 1 - \eta\rho_1$

$$v_{k+1} - v_k = -\eta\rho_1 \left( v_k - (\nabla f_{\gamma_k}(x_k))^2 \right). \quad (175)$$

2376

Then,

2377

2378

2379

$$\frac{1}{\sqrt{v_{k+1}}} = \sqrt{\frac{v_k}{v_{k+1}}} \frac{1}{\sqrt{v_k}} = \sqrt{\frac{v_{k+1} + \mathcal{O}(\eta)}{v_{k+1}}} \frac{1}{\sqrt{v_k}} = \sqrt{1 + \frac{\mathcal{O}(\eta)}{v_{k+1}}} \sqrt{\frac{1}{v_k}} \sim \sqrt{\frac{1}{v_k}} (1 + \mathcal{O}(\eta)). \quad (176)$$

2380

2381

Therefore, we work with the following algorithm as all approximations only carry an additional error of order  $\mathcal{O}(\eta^2)$ , which we can ignore. Therefore, we have that

2382

2383

$$v_k - v_{k-1} = -\eta\rho_2 \left( v_{k-1} - (\nabla f_{\gamma_{k-1}}(x_{k-1}))^2 \right) \quad (177)$$

2384

2385

$$m_{k+1} - m_k = -\eta\rho_1 (m_k - \nabla f_{\gamma_k}(x_k)) \quad (178)$$

2386

2387

$$\hat{m}_k = m_k (1 - \beta_1^k)^{-1} \quad (179)$$

2388

$$\hat{v}_k = v_k (1 - \beta_1^k)^{-1} \quad (180)$$

2389

2390

$$x_{k+1} - x_k = -\frac{\eta}{\sqrt{v_k + \epsilon I_d}} \frac{\sqrt{1 - (1 - \eta\rho_2)^k}}{1 - (1 - \eta\rho_1)^{k+1}} (m_k + \eta\rho_1 (\nabla f_{\gamma_k}(x_k) - m_k)). \quad (181)$$

2391

2392

Therefore, if  $\nabla f_{\gamma_j}(x_j) = \nabla f(x_j) + Z_j(x_j)$  and  $\mathbb{E}[Z_j(x_j)] = 0$ , and  $Cov(Z_j(x_j)) = \Sigma(x_j)$ , we have that

2393

2394

2395

$$1. \mathbb{E}[v_k - v_{k-1}] = \eta\rho_2 \left[ (\nabla f(x_{k-1}))^2 + \text{diag}(\Sigma(x_k)) - v_{k-1} \right];$$

2396

2397

$$2. \mathbb{E}[m_{k+1} - m_k] = \eta\rho_1 [\nabla f(x_k) - m_k];$$

2398

2399

$$3. \mathbb{E}[x_{k+1} - x_k] = -\frac{\eta}{\sqrt{v_k + \epsilon I_d}} \frac{\sqrt{1 - (1 - \eta\rho_2)^k}}{1 - (1 - \eta\rho_1)^{k+1}} (m_k + \eta\rho_1 (\nabla f(x_k) - m_k)).$$

2400

Then, we have

2401

2402

2403

2404

2405

2406

2407

2408

2409

2410

2411

$$1. \mathbb{E}[(x_{k+1} - x_k)(x_{k+1} - x_k)^\top] = \mathbb{E}[(x_{k+1} - x_k)]\mathbb{E}[(x_{k+1} - x_k)]^\top + \mathcal{O}(\eta^4\rho_1^2);$$

$$2. \mathbb{E}[(x_{k+1} - x_k)(m_k - m_{k-1})^\top] = \mathbb{E}[(x_{k+1} - x_k)]\mathbb{E}[(m_k - m_{k-1})]^\top + 0;$$

$$3. \mathbb{E}[(x_{k+1} - x_k)(v_k - v_{k-1})^\top] = \mathbb{E}[(x_{k+1} - x_k)]\mathbb{E}[(v_k - v_{k-1})]^\top + 0;$$

$$4. \mathbb{E}[(v_k - v_{k-1})(v_k - v_{k-1})^\top] = \mathbb{E}[(v_k - v_{k-1})]\mathbb{E}[(v_k - v_{k-1})]^\top + \mathcal{O}(\eta^2\rho_2^2);$$

$$5. \mathbb{E}[(m_k - m_{k-1})(m_k - m_{k-1})^\top] = \mathbb{E}[(m_k - m_{k-1})]\mathbb{E}[(m_k - m_{k-1})]^\top + \eta^2\rho_1^2\Sigma(x_{k-1});$$

$$6. \mathbb{E}[(v_k - v_{k-1})(m_k - m_{k-1})^\top] = \mathbb{E}[(v_k - v_{k-1})]\mathbb{E}[(m_k - m_{k-1})]^\top + \mathcal{O}(\eta^2\rho_1\rho_2).$$

Since in real-world applications,  $\rho_1 = \mathcal{O}(\eta^{-\zeta})$  s.t.  $\zeta \in (0, 1)$ , while  $\rho_2 = \mathcal{O}(1)$ , we have

2412

2413

2414

2415

2416

2417

2418

2419

$$dX_t = -\frac{\sqrt{\iota_2(t)}}{\iota_1(t)} P_t^{-1} (M_t + \eta\rho_1 (\nabla f(X_t) - M_t)) dt \quad (182)$$

$$dM_t = \rho_1 (\nabla f(X_t) - M_t) dt + \sqrt{\eta}\rho_1 \Sigma^{1/2}(X_t) dW_t \quad (183)$$

$$dV_t = \rho_2 \left( (\nabla f(X_t))^2 + \text{diag}(\Sigma(X_t)) - V_t \right) dt. \quad (184)$$

where  $\beta_i = 1 - \eta\rho_i$ ,  $\iota_i(t) = 1 - e^{-\rho_i t}$ ,  $t > t_0$ , and  $P_t = \text{diag} \sqrt{V_t} + \epsilon \sqrt{\iota_2(t)} I_d$ .  $\square$

**Remark C.48.** See Remark C.35 and Remark C.36 for a discussion on the regularity of the SDE derived in Theorem C.47.

2420

2421

2422

2423

**Corollary C.49.** Under the assumptions of Theorem C.47 with  $\Sigma(x) = \sigma^2 I_d$ ,  $\tilde{\eta} = \kappa\eta$ ,  $\tilde{B} = B\delta$ ,  $\tilde{\rho}_1 = \alpha_1\rho_1$ , and  $\tilde{\rho}_2 = \alpha_2\rho_2$

2424

2425

2426

2427

2428

2429

$$dX_t = -\kappa \frac{\sqrt{\iota_2(t)}}{\iota_1(t)} P_t^{-1} (M_t + \eta\alpha_1\rho_1 (\nabla f(X_t) - M_t)) dt \quad (185)$$

$$dM_t = \frac{\alpha_1\rho_1}{\kappa} (\nabla f(X_t) - M_t) dt + \sqrt{\eta} \frac{\alpha_1\rho_1}{\kappa} \frac{\sigma}{\sqrt{B\delta}} I_d dW_t \quad (186)$$

$$dV_t = \frac{\alpha_2\rho_2}{\kappa} \left( (\nabla f(X_t))^2 + \frac{\sigma^2}{B\delta} I_d - V_t \right) dt. \quad (187)$$

**Lemma C.50.** *Under the assumptions of Corollary C.49,  $f$  is  $\mu$ -strongly convex,  $\text{Tr}(\nabla^2 f(x)) \leq \mathcal{L}_\tau$ , and  $(\nabla f(x))^2 = \mathcal{O}(\eta)$ , the asymptotic dynamics of the iterates of Adam satisfies the classic scaling rule  $\kappa = \sqrt{\delta}$  because  $\mathbb{E}[f(X_t)] \stackrel{t \rightarrow \infty}{\leq} \frac{\eta\sigma\mathcal{L}_\tau}{4\sqrt{B}} \frac{\kappa}{\sqrt{\delta}}$ . To enforce that the speed of  $M_t$  and  $V_t$  match that of  $X_t$ , one needs  $\tilde{\rho}_i = \kappa^2 \rho_i$ , which implies  $\tilde{\beta}_i = 1 - \kappa^2(1 - \beta_i)$ .*

*Proof.* First of all, we need to ensure that the relative speeds of  $X_t$ ,  $M_t$ , and  $V_t$  match. Therefore, we select  $\alpha_i = \kappa^2$ , which recovers the scaling rules for  $\tilde{\beta}_i = 1 - \kappa^2(1 - \beta_i)$ . Then, recalling that  $(\nabla f(x))^2 = \mathcal{O}(\eta)$ , we have that as  $t \rightarrow \infty$ ,  $V_t \rightarrow \frac{\sigma^2}{B\delta}$ , and  $M_t \rightarrow \nabla f(X_t)$  with high probability. Therefore,

$$dX_t = -\kappa \frac{\sqrt{B\delta}}{\sigma} \nabla f(X_t) dt \quad (188)$$

$$dM_t = \kappa \sqrt{\eta} \rho_1 \frac{\sigma}{\sqrt{B\delta}} dW_t \quad (189)$$

$$dV_t = 0. \quad (190)$$

Therefore, if  $H(X_t, V_t) := f(X_t) + \frac{\mathcal{L}_\tau \delta B}{\rho_1^2 \sigma^2} \frac{\|M_t\|_2^2}{2}$  and  $\xi \in (0, 1)$  we have that by Itô's lemma,

$$dH(X_t, V_t) = -(\nabla f(X_t))^\top \left( \kappa \frac{\sqrt{B\delta}}{\sigma} \nabla f(X_t) \right) dt + \left( \frac{\mathcal{L}_\tau \delta B}{\rho_1^2 \sigma^2} M_t \right) \kappa \sqrt{\eta} \rho_1 \frac{\sigma}{\sqrt{B\delta}} dW_t \quad (191)$$

$$+ \frac{1}{2} \left( \frac{\mathcal{L}_\tau \delta B}{\rho_1^2 \sigma^2} \right) \kappa^2 \eta \rho_1^2 \frac{\sigma^2}{B\delta} dt \quad (192)$$

$$= - \left( \kappa \frac{\sqrt{B\delta}}{\sigma} \right) \|\nabla f(X_t)\|_2^2 dt + \text{Noise} + \frac{\kappa^2 \eta \mathcal{L}_\tau}{2} dt \quad (193)$$

$$= - \left( \kappa \frac{\sqrt{B\delta}}{\sigma} \right) (\xi \|\nabla f(X_t)\|_2^2 + (1 - \xi) \|\nabla f(X_t)\|_2^2) dt + \text{Noise} + \frac{\kappa^2 \eta \mathcal{L}_\tau}{2} dt \quad (194)$$

$$\leq -2\kappa\mu \frac{\sqrt{B\delta}}{\sigma} \xi \left( f(X_t) + \frac{1 - \xi}{\mu\xi} \frac{\|\nabla f(X_t)\|_2^2}{2} \right) dt + \text{Noise} + \frac{\kappa^2 \eta \mathcal{L}_\tau}{2} dt. \quad (195)$$

Let us now select  $\xi$  such that  $\frac{1-\xi}{\mu\xi} = \frac{\mathcal{L}_\tau \delta B}{\rho_1^2 \sigma^2}$ , this means that  $\xi = \frac{\sigma^2 \rho_1^2}{\sigma^2 \rho_1^2 + \mu \mathcal{L}_\tau \delta B} \in (0, 1)$  and  $\frac{1}{\xi} = 1 + \mu \frac{\mathcal{L}_\tau \delta B}{\rho_1^2 \sigma^2}$ . Since  $M_t \rightarrow \nabla f(X_t)$ , we have that

$$dH(X_t, V_t) \leq -2\kappa\mu \frac{\sqrt{B\delta}}{\sigma} \xi H(X_t, V_t) dt + \frac{\kappa^2 \eta \mathcal{L}_\tau}{2} dt + \text{Noise}. \quad (196)$$

Therefore,

$$\frac{\mathbb{E}[f(X_t)]}{\xi} = \left( 1 + \mu \frac{\mathcal{L}_\tau \delta B}{\rho_1^2 \sigma^2} \right) \mathbb{E}[f(X_t)] \leq \mathbb{E}[H(X_t, V_t)] \stackrel{t \rightarrow \infty}{\leq} \frac{1}{\xi} \frac{\eta\sigma\mathcal{L}_\tau}{4\mu\sqrt{B}} \frac{\kappa}{\sqrt{\delta}}, \quad (197)$$

which implies that

$$\mathbb{E}[f(X_t)] \stackrel{t \rightarrow \infty}{\leq} \frac{\eta\sigma\mathcal{L}_\tau}{4\mu\sqrt{B}} \frac{\kappa}{\sqrt{\delta}}. \quad (198)$$

Analogously,

$$\mathbb{E}[f(X_t) - f(X_*)] \stackrel{t \rightarrow \infty}{\leq} \frac{\eta\sigma\mathcal{L}_\tau}{4\mu\sqrt{B}} \frac{\kappa}{\sqrt{\delta}}. \quad (199)$$

which gives the square root scaling rule.  $\square$

**Lemma C.51.** *Under the assumptions of Corollary C.49,  $f(x) = \frac{x^\top H x}{2}$  s.t.  $H = \text{diag}(\lambda_1, \dots, \lambda_d)$  and  $(\nabla f(x))^2 = \mathcal{O}(\eta)$ , the dynamics of Adam implies that  $f(X_t) \rightarrow \frac{\eta\sigma d}{4\sqrt{B}} \frac{\kappa}{\sqrt{\delta}}$ .*

2484 *Proof.* Recalling that  $(\nabla f(x))^2 = \mathcal{O}(\eta)$ , we have that as  $t \rightarrow \infty$ ,  $V_t \rightarrow \frac{\sigma^2}{B\delta}$ , and  $M_t \rightarrow \lambda X_t$  with  
 2485 high probability. Therefore, in the one-dimensional case  
 2486

$$2487 \quad dX_t = -\kappa \frac{\sqrt{B\delta}}{\sigma} \lambda X_t dt \quad (200)$$

$$2489 \quad dM_t = \kappa \sqrt{\eta} \rho_1 \frac{\sigma}{\sqrt{B\delta}} dW_t \quad (201)$$

$$2491 \quad dV_t = 0. \quad (202)$$

2493 Therefore, if  $H(X_t, V_t) := \frac{\lambda X_t^2}{2} + \frac{\lambda \delta B}{\rho_1^2 \sigma^2} \frac{M_t^2}{2}$ ,<sup>8</sup> we have that by Itô's lemma,  
 2494

$$2496 \quad dH(X_t, V_t) = -(\lambda X_t) \left( \kappa \frac{\sqrt{B\delta}}{\sigma} \lambda X_t \right) dt + \left( \frac{\lambda \delta B}{\rho_1^2 \sigma^2} M_t \right) \kappa \sqrt{\eta} \rho_1 \frac{\sigma}{\sqrt{B\delta}} dW_t \quad (203)$$

$$2499 \quad + \frac{1}{2} \left( \frac{\lambda \delta B}{\rho_1^2 \sigma^2} \right) \kappa^2 \eta \rho_1^2 \frac{\sigma^2}{B\delta} dt \quad (204)$$

$$2502 \quad = -2\kappa \lambda \frac{\sqrt{B\delta}}{\sigma} f(X_t) dt + \frac{\kappa^2 \eta \rho_1^2 \sigma^2}{2B\delta} \frac{\lambda \delta B}{\rho_1^2 \sigma^2} dt + \text{Noise}. \quad (205)$$

$$2504 \quad = -2\kappa \lambda \frac{\sqrt{B\delta}}{\sigma} f(X_t) dt + \frac{\kappa^2 \eta \lambda}{2} dt + \text{Noise}. \quad (206)$$

2506 Once again, since  $M_t \rightarrow \lambda X_t$ , we have that  
 2507

$$2509 \quad H(X_t, V_t) = \frac{\lambda X_t^2}{2} + \frac{\lambda \delta B}{\rho_1^2 \sigma^2} \frac{M_t^2}{2} \rightarrow \frac{\lambda X_t^2}{2} + \lambda \frac{\lambda \delta B}{\rho_1^2 \sigma^2} \frac{\lambda X_t^2}{2} = \left( 1 + \lambda \frac{\lambda \delta B}{\rho_1^2 \sigma^2} \right) \frac{\lambda X_t^2}{2} =: K f(X_t). \quad (207)$$

2512 Therefore,

$$2514 \quad K d\mathbb{E}[f(X_t)] = -2\kappa \lambda \frac{\sqrt{B\delta}}{\sigma} \mathbb{E}[f(X_t)] dt + \frac{\kappa^2 \eta \lambda}{2} dt, \quad (208)$$

2516 which implies that  $\mathbb{E}[f(X_t)] \rightarrow \frac{\eta \sigma}{4\sqrt{B}} \frac{\kappa}{\sqrt{\delta}}$ , which also gives the square root scaling rule. The general-  
 2517 ization to  $d$  dimension is analogous and one needs to sum across all the dimensions.  $\square$   
 2518

2519 **Lemma C.52.** Let  $f(x) := \frac{x^\top H x}{2}$  where  $H = \text{diag}(\lambda_1, \dots, \lambda_d)$ . The stationary distribution of  
 2520 Adam is  $(\mathbb{E}[X_\infty], \text{Cov}(X_\infty)) = \left( 0, \frac{\eta}{2} \Sigma^{\frac{1}{2}} H^{-1} \right)$ .  
 2521

2522 *Proof.* The expected value follows immediately from the fact that  
 2523

$$2525 \quad dX_t = -\Sigma^{-\frac{1}{2}} X_t dt \quad (209)$$

2527 For the covariance, we focus on the one-dimensional case. We define  $H(X_t, V_t) := \frac{X_t^2}{2} + \frac{\lambda^2}{2\sigma^2 \rho_1^2} \frac{M_t^2}{2}$ .  
 2528 With the same arguments as Lemma C.51, we have  
 2529

$$2530 \quad d(X_t)^2 = -\frac{\lambda}{\sigma} X_t^2 dt + \frac{\eta}{2} dt + \text{Noise}, \quad (210)$$

2532 which implies that  
 2533

$$2534 \quad \mathbb{E}[X_t^2] \xrightarrow{t \rightarrow 0} \frac{\eta \sigma}{2 \lambda}. \quad (211)$$

2536 The thesis follows by applying the same logic to multiple dimensions.  $\square$   
 2537

<sup>8</sup>Inspired by (Barakat and Bianchi, 2021)

2538 C.9 ADAMW  
2539

2540 In this subsection, we derive the SDE of AdamW defined as defined as

2541 
$$v_{k+1} = \beta_2 v_k + (1 - \beta_2) (\nabla f_{\gamma_k}(x_k))^2 \quad (212)$$

2542 
$$m_{k+1} = \beta_1 m_k + (1 - \beta_1) \nabla f_{\gamma_k}(x_k) \quad (213)$$

2543 
$$\hat{m}_k = m_k (1 - \beta_1^k)^{-1} \quad (214)$$

2544 
$$\hat{v}_k = v_k (1 - \beta_2^k)^{-1} \quad (215)$$

2545 
$$x_{k+1} = x_k - \eta \frac{\hat{m}_{k+1}}{\sqrt{\hat{v}_{k+1} + \epsilon I_d}} - \eta \theta x_k \quad (216)$$

2546 with  $(x_0, m_0, v_0) \in \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}^d$ ,  $\eta \in \mathbb{R}^{>0}$  is the step size,  $\beta_i = 1 - \rho_i \eta$  for  $\rho_1 = \mathcal{O}(\eta^{-\zeta})$   
2547 s.t.  $\zeta \in (0, 1)$ ,  $\rho_2 = \mathcal{O}(1)$ ,  $\theta > 0$ , the mini-batches  $\{\gamma_k\}$  are modelled as i.i.d. random variables  
2548 uniformly distributed on  $\{1, \dots, N\}$ , and of size  $B \geq 1$ .2549 **Theorem C.53.** *Under the same assumptions as Theorem C.47, the SDE of AdamW is*

2550 
$$dX_t = -\frac{\sqrt{\iota_2(t)}}{\iota_1(t)} P_t^{-1} (M_t + \eta \rho_1 (\nabla f(X_t) - M_t)) dt - \theta X_t dt \quad (217)$$

2551 
$$dM_t = \rho_1 (\nabla f(X_t) - M_t) dt + \sqrt{\eta} \rho_1 \Sigma^{1/2}(X_t) dW_t \quad (218)$$

2552 
$$dV_t = \rho_2 ((\nabla f(X_t))^2 + \text{diag}(\Sigma(X_t)) - V_t) dt. \quad (219)$$

2553 where  $\beta_i = 1 - \eta \rho_i$ ,  $\theta > 0$ ,  $\iota_i(t) = 1 - e^{-\rho_i t}$ ,  $t > t_0$ , and  $P_t = \text{diag}(\sqrt{V_t} + \epsilon \sqrt{\iota_2(t)} I_d)$ .2554 *Proof.* The proof is the same as the of Theorem C.47 and the only difference is that  $\eta \theta x_k$  is  
2555 approximated with  $\theta X_t dt$ .  $\square$ 

2556 Figure 4 and Figure 11 validate this result on a variety of architectures and datasets.

2557 **Remark C.54.** See Remark C.35 and Remark C.36 for a discussion on the regularity of the SDE  
2558 derived in Theorem C.53.2559 **Corollary C.55.** *Under the assumptions of Theorem C.53 with  $\Sigma(x) = \sigma^2 I_d$ ,  $\tilde{\eta} = \kappa \eta$ ,  $\tilde{B} = B \delta$ ,  
2560  $\tilde{\rho}_1 = \alpha_1 \rho_1$ ,  $\tilde{\theta} = \xi \theta$ , and  $\tilde{\rho}_2 = \alpha_2 \rho_2$* 

2561 
$$dX_t = -\kappa \frac{\sqrt{\iota_2(t)}}{\iota_1(t)} P_t^{-1} (M_t + \eta \alpha_1 \rho_1 (\nabla f(X_t) - M_t)) dt - \kappa \xi \theta X_t dt \quad (220)$$

2562 
$$dM_t = \frac{\alpha_1 \rho_1}{\kappa} (\nabla f(X_t) - M_t) dt + \sqrt{\eta} \frac{\alpha_1 \rho_1}{\kappa} \frac{\sigma}{\sqrt{B \delta}} I_d dW_t \quad (221)$$

2563 
$$dV_t = \frac{\alpha_2 \rho_2}{\kappa} \left( (\nabla f(X_t))^2 + \frac{\sigma^2}{B \delta} I_d - V_t \right) dt. \quad (222)$$

2564 **Lemma C.56** (Scaling Rule at Convergence). *Under the assumptions of Corollary C.55,  $f$  is  $\mu$ -  
2565 strongly convex and  $L$ -smooth,  $\text{Tr}(\nabla^2 f(x)) \leq \mathcal{L}_\tau$ , and  $(\nabla f(x))^2 = \mathcal{O}(\eta)$ , the asymptotic dynamics  
2566 of the iterates of AdamW satisfies the novel scaling rule if  $\kappa = \sqrt{\delta}$  and  $\xi = \kappa$  because*

2567 
$$\mathbb{E}[f(X_t) - f(X_*)] \stackrel{t \rightarrow \infty}{\leq} \frac{\eta \mathcal{L}_\tau \sigma L}{2} \frac{\kappa}{2\mu \sqrt{B \delta} L + \sigma \xi \theta (L + \mu)} \quad (223)$$

2568 *By enforcing that the speed of  $V_t$  matches that of  $X_t$ , one needs  $\tilde{\rho} = \kappa^2 \rho$ , which implies  $\tilde{\beta}_i =$   
2569  $1 - \kappa^2 (1 - \beta_i)$ .*2570 *Proof.* The proof is the same as Lemma C.50 where we also use  $L$ -smoothness as in Lemma C.44.  $\square$ 2571 **Lemma C.57.** *For  $f(x) := \frac{x^\top H x}{2}$ , the stationary distribution of AdamW is  $(\mathbb{E}[X_\infty], \text{Cov}(X_\infty)) =$   
2572  $(0, \frac{\eta}{2} (H \Sigma^{-\frac{1}{2}} + \theta I_d)^{-1})$ .*2573 *Proof.* The proof is the same as Lemma C.52.  $\square$



Finally, we prove a generalization of Lemma C.56 to the  $L$ -smooth case.

**Lemma C.58.** *Let  $f$  be  $L$ -smooth,  $\eta_t$  be a learning rate scheduler such that  $\lim_{t \rightarrow \infty} \frac{\phi_t^2}{\phi_t^1} \xrightarrow{t \rightarrow \infty} 0$  and  $\phi_t^1 \xrightarrow{t \rightarrow \infty} \infty$ , where  $\phi_t^i = \int_0^t (\eta_s)^i ds$ . Then*

$$\mathbb{E} \|\nabla f(X_{\tilde{t}})\|_2^2 \leq \left( f(X_0) - f(X_*) + \frac{\mathcal{L}_\tau \delta B}{\rho_1^2 \sigma^2} \frac{\|M_0\|_2^2}{2} + \frac{\phi_{\tilde{t}}^2 \eta \kappa^2 \mathcal{L}_\tau}{2} \right) \frac{\sigma}{\kappa \sqrt{\delta B}} \frac{1}{\phi_{\tilde{t}}^1} \xrightarrow{t \rightarrow \infty} 0, \quad (224)$$

where  $\tilde{t}$  is a random time with distribution  $\frac{\eta_t}{\phi_t^1}$ .

*Proof.* The proof is the same as Lemma C.24.  $\square$

## D SDES FROM THE LITERATURE

**Theorem D.1** (Original Malladi's Statement). *Let  $\sigma_0 := \sigma\eta$ ,  $\epsilon_0 := \epsilon\eta$ , and  $c_2 := \frac{1-\beta}{\eta^2}$ . Define the state of the SDE as  $L_t = (X_t, u_t)$  and the dynamics as*

$$dX_t = -P_t^{-1} \left( \nabla f(X_t) dt + \sigma_0 \Sigma^{1/2}(X_t) dW_t \right) \quad (225)$$

$$du_t = c_2 (\text{diag}(\Sigma(X_t)) - u_t) dt \quad (226)$$

where  $P_t := \sigma_0 \text{diag}(u_t)^{1/2} + \epsilon_0 I_d$ .

**Theorem D.2** (Informal Statement of Theorem C.2 Malladi et al. (2022)). *Under sufficient regularity conditions and  $\nabla f(x) = \mathcal{O}(\sqrt{\eta})$ , the following SDE is an order 1 weak approximation of RMSprop:*

$$dX_t = -P_t^{-1} (\nabla f(X_t) dt + \sqrt{\eta} \Sigma(X_t)^{\frac{1}{2}} dW_t) \quad (227)$$

$$dV_t = \rho (\text{diag}(\Sigma(X_t)) - V_t) dt, \quad (228)$$

where  $\beta = 1 - \eta\rho$ ,  $\rho = \mathcal{O}(1)$ , and  $P_t := \text{diag}(V_t)^{\frac{1}{2}} + \epsilon I_d$ .

**Lemma D.3.** *Theorem D.1 and Theorem D.2 are equivalent.*

*Proof.* It follows applying time rescaling  $t := \eta\xi$  and observing that  $W_t = W_{\eta\xi} = \sqrt{\eta} W_\xi$ .  $\square$

**Theorem D.4** (Original Malladi's Statement). *Let  $c_1 := (1 - \beta_1) / \eta^2$ ,  $c_2 := (1 - \beta_2) / \eta^2$  and define  $\sigma_0, \epsilon_0$  in Theorem D.1. Let  $\iota_1(t) := 1 - \exp(-c_1 t)$  and  $\iota_2(t) := 1 - \exp(-c_2 t)$ . Define the state of the SDE as  $L_t = (X_t, m_t, u_t)$  and the dynamics as*

$$dX_t = -\frac{\sqrt{\iota_2(t)}}{\iota_1(t)} P_t^{-1} m_t dt \quad (229)$$

$$dm_t = c_1 (\nabla f(X_t) - m_t) dt + \sigma_0 c_1 \Sigma^{1/2}(X_t) dW_t, \quad (230)$$

$$du_t = c_2 (\text{diag}(\Sigma(X_t)) - u_t) dt, \quad (231)$$

where  $P_t := \sigma_0 \text{diag}(u_t)^{1/2} + \epsilon_0 \sqrt{\iota_2(t)} I_d$ .

**Theorem D.5** (Informal Statement of Theorem D.2 Malladi et al. (2022)). *Under sufficient regularity conditions and  $\nabla f(x) = \mathcal{O}(\sqrt{\eta})$ , the following SDE is an order 1 weak approximation of Adam:*

$$dX_t = -\frac{\sqrt{\iota_2(t)}}{\iota_1(t)} P_t^{-1} M_t dt \quad (232)$$

$$dM_t = \rho_1 (\nabla f(X_t) - M_t) dt + \sqrt{\eta} \rho_1 \Sigma^{1/2}(X_t) dW_t \quad (233)$$

$$dV_t = \rho_2 (\text{diag}(\Sigma(X_t)) - V_t) dt. \quad (234)$$

where  $\beta_i = 1 - \eta\rho_i$ ,  $\iota_i(t) = 1 - e^{-\rho_i t}$ ,  $\rho_i = \mathcal{O}(1)$ , and  $P_t = \text{diag}(\sqrt{V_t} + \epsilon \sqrt{\iota_2(t)}) I_d$ .

**Lemma D.6.** *Theorem D.4 and Theorem D.5 are equivalent.*

*Proof.* It follows applying time rescaling  $t := \eta\xi$  and observing that  $W_t = W_{\eta\xi} = \sqrt{\eta} W_\xi$ .  $\square$

## E SDE CANNOT BE DERIVED NOR USED NAIVELY

In this section, we provide a gentle introduction to the meaning of deriving an SDE model for an optimizer and discuss how SDEs have been used to derive scaling rules. To aid the intuition of the reader, we informally derive an SDE for SGD with learning rate  $\eta$ , mini-batches  $\gamma_B$  of size  $B$ , and starting point  $x_0 = x$ , which we dub  $\text{SGD}^{(\eta, B)}$ . The iterates are given by:

$$x_{k+1} = x_k - \eta \nabla f_{\gamma_B}(x_k) \quad (235)$$

which for  $U_k := \sqrt{\eta}(\nabla f(x_k) - \nabla f_{\gamma_B}(x_k))$ , we rewrite as

$$x_k - \eta \nabla f(x_k) + \sqrt{\eta} U_k, \quad (236)$$

where  $\mathbb{E}[U_k] = 0$  and  $\text{Cov}(U_k) = \frac{\eta}{B} \Sigma(x_k) = \frac{\eta}{B} \frac{1}{n} \sum_{i=0}^n (\nabla f(x_k) - \nabla f_i(x_k))(\nabla f(x_k) - \nabla f_i(x_k))^\top$ . If we now consider the SDE

$$dX_t = -\nabla f(X_t)dt + \sqrt{\frac{\eta}{B}} \Sigma(X_t)^{\frac{1}{2}} dW_t, \quad (237)$$

its Euler-Maruyama discretization with pace  $\Delta t = \eta$  and  $Z_k \sim \mathcal{N}(0, I_d)$  is

$$X_{k+1} = X_k - \eta \nabla f(X_k) + \sqrt{\eta} \sqrt{\frac{\eta}{B}} \Sigma(X_k)^{\frac{1}{2}} Z_k. \quad (238)$$

Since the Eq. 235 and Eq. 238 share the first two moments, it is reasonable that by identifying  $t = k\eta$ , the SDE in Eq. 237 is a good model to describe the iterates of SGD in Eq. 235.

Informally, we need a “good model”, which is an SDE that is close to the real optimizer. This is formalized in the following definition which comes from the field of numerical analysis of SDEs (see Mil’shtein (1986)) and bounds the disparity between the the discrete and the continuous process.

**Definition E.1** (Weak Approximation). A continuous-time stochastic process  $\{X_t\}_{t \in [0, T]}$  is an order  $\alpha$  weak approximation (or  $\alpha$ -order SDE) of a discrete stochastic process  $\{x_k\}_{k=0}^{\lfloor T/\eta \rfloor}$  if for every polynomial growth function  $g$ , there exists a positive constant  $C$ , independent of the stepsize  $\eta$ , such that  $\max_{k=0, \dots, \lfloor T/\eta \rfloor} |\mathbb{E}g(x_k) - \mathbb{E}g(X_{k\eta})| \leq C\eta^\alpha$ .

To see if an SDE satisfies such a definition, one has to check that for  $\bar{\Delta} = x_1 - x$  and  $\Delta = X_\eta - x$ ,

1.  $|\mathbb{E}\Delta_i - \mathbb{E}\bar{\Delta}_i| = \mathcal{O}(\eta^2), \quad \forall i = 1, \dots, d;$
2.  $|\mathbb{E}\Delta_i \Delta_j - \mathbb{E}\bar{\Delta}_i \bar{\Delta}_j| = \mathcal{O}(\eta^2), \quad \forall i, j = 1, \dots, d.$

**Example:** Let us prove that the SDE in Eq. 237 is a valid approximation of  $\text{SGD}^{(\eta, B)}$ : The first condition is easily verified. Coming to the second condition we have that

1.  $\mathbb{E}\Delta_i \Delta_j = \eta^2 \partial_i f(x) \partial_j f(x) + \frac{\eta^2}{B} \Sigma(x);$
2.  $\mathbb{E}\bar{\Delta}_i \bar{\Delta}_j = \eta^2 \partial_i f(x) \partial_j f(x) + \frac{\eta^2}{B} \Sigma(x) + \mathcal{O}(\eta^3);$

whose difference is of order  $\eta^3$  and thus satisfies the condition. However, we observe that if the scale of the noise is too small w.r.t.  $\eta$ , i.e.  $\Sigma(x) = \mathcal{O}(\eta^\alpha)$  for  $\alpha \geq 0$ , then the **simplest** SDE model describing  $\text{SGD}^{(\eta, B)}$  is the ODE  $dX_t = -\nabla f(X_t)dt$  as in that case

1.  $\mathbb{E}\Delta_i \Delta_j = \eta^2 \partial_i f(x) \partial_j f(x) + \mathcal{O}(\eta^{2+\alpha});$
2.  $\mathbb{E}\bar{\Delta}_i \bar{\Delta}_j = \eta^2 \partial_i f(x) \partial_j f(x) + \mathcal{O}(\eta^2),$

whose difference is also of order  $\eta^2$ . Much differently, if  $\Sigma(x) = \mathcal{O}(\eta^{-\alpha})$  for  $\alpha > 0$ , the simplest model is the SDE in Eq. 237. We highlight that *simplest* does not mean *best*: The SDE is more accurate than the ODE even in a regime with low noise, but this observation serves as a provocation. One has to pay attention when deriving SDEs: Some models are more realistic than others.

Let us dig deeper into this thought as we derive **two** SDEs for SGD with learning rate  $\tilde{\eta} := \kappa\eta$  and batch size  $\tilde{B} := \delta B$  for  $\kappa > 1$  and  $\delta > 1$ , which we dub  $\text{SGD}^{(\tilde{\eta}, \tilde{B})}$ . The first is derived considering that the learning rate is  $\tilde{\eta}$  and carries an error of order  $\mathcal{O}(\tilde{\eta})$  w.r.t.  $\text{SGD}^{(\tilde{\eta}, \tilde{B})}$

$$dX_t = -\nabla f(X_t)dt + \sqrt{\frac{\tilde{\eta}}{\tilde{B}}}\Sigma(X_t)^{\frac{1}{2}}dW_t = -\nabla f(X_t)dt + \sqrt{\frac{\eta\kappa}{B\delta}}\Sigma(X_t)^{\frac{1}{2}}dW_t. \quad (239)$$

The second one instead is derived considering  $\eta$  as the learning rate and  $\kappa$  as a constant “scheduler”. Consistently with (Li et al., 2017), the SDE which carries an error of order  $\mathcal{O}(\eta)$  w.r.t.  $\text{SGD}^{(\tilde{\eta}, \tilde{B})}$  is

$$dX_t = -\kappa\nabla f(X_t)dt + \kappa\sqrt{\frac{\eta}{B\delta}}\Sigma(X_t)^{\frac{1}{2}}dW_t. \quad (240)$$

While they both are valid models, there are three reasons why one should prefer the latter:

1. It fully reflects the fact that a larger learning rate results in a faster and noisier dynamics;
2. It has intrinsically less error than the other;
3. It is consistent with the optimizer in that there is no combination of  $\kappa$  and  $\delta$  that can ever leave the dynamics unchanged.

## E.1 DERIVING SCALING RULES

Jastrzebski et al. (2018) observed that only the ratio between  $\eta$  and  $B$  matters in determining the dynamics of Eq. 238. Therefore, they argue that for  $\kappa = \delta$  the SDE for  $\text{SGD}^{(\kappa\eta, \delta B)}$  coincides with that of  $\text{SGD}^{(\eta, B)}$  and that this implies that the path properties of the optimizers are the same. On the contrary, the path of  $\text{SGD}^{(\eta, B)}$  strongly depends on the hyperparameters: The speed and volatility of the dynamics are driven by  $\eta$ , and no choice of  $B$  can undo this. We remind the reader that the goal of these rules is not to keep the dynamics of the optimizers unaltered, but rather to give a practical way to change a hyperparameter, e.g.  $\eta$ , and have a principled way to adjust the others, e.g.  $B$ , such that the performance of the optimizer is preserved. Therefore, we propose deriving scaling rules as we preserve certain relevant quantities of the dynamics such as the convergence bound on the expected loss. To show this quantitatively, we use this rationale to derive the scaling rule of SGD as we aim at preserving the asymptotic loss level.

**Lemma E.2.** *If  $f$  is a  $\mu$  strongly convex function,  $\text{Tr}(\nabla^2 f(x)) \leq \mathcal{L}_\tau$  and  $\Sigma(x) = \sigma^2 I_d$ , then:*

1. Under the dynamics of Eq. 237 we have:

$$\mathbb{E}[f(X_t) - f(X_*)] \leq (f(X_0) - f(X_*))e^{-2\mu t} + \frac{\eta}{2} \frac{\mathcal{L}_\tau \sigma^2}{2\mu B} (1 - e^{-2\mu t}); \quad (241)$$

2. Under the dynamics of Eq. 239 we have:

$$\mathbb{E}[f(X_t) - f(X_*)] \leq (f(X_0) - f(X_*))e^{-2\mu t} + \frac{\eta}{2} \frac{\mathcal{L}_\tau \sigma^2}{2\mu B} \frac{\kappa}{\delta} (1 - e^{-2\mu t}); \quad (242)$$

3. Under the dynamics of Eq. 240 we have:

$$\mathbb{E}[f(X_t) - f(X_*)] \leq (f(X_0) - f(X_*))e^{-2\mu\kappa t} + \frac{\eta}{2} \frac{\mathcal{L}_\tau \sigma^2}{2\mu B} \frac{\kappa}{\delta} (1 - e^{-2\mu\kappa t}). \quad (243)$$

The first bound implies that the asymptotic limit of the expected loss for  $\text{SGD}^{(\eta, B)}$  is  $\frac{\eta}{2} \frac{\mathcal{L}_\tau \sigma^2}{2\mu B}$ . The last two bounds predict that the asymptotic loss level for  $\text{SGD}^{(\tilde{\eta}, \tilde{B})}$  is  $\frac{\eta}{2} \frac{\mathcal{L}_\tau \sigma^2}{2\mu B} \frac{\kappa}{\delta}$ . Since the objective of the scaling rule is to find  $\kappa$  and  $\delta$  such that  $\text{SGD}^{(\tilde{\eta}, \tilde{B})}$  achieves the same loss level as  $\text{SGD}^{(\eta, B)}$ , we recover the linear scaling rule setting  $\kappa = \delta$ . However, only the last bound can correctly capture the fact that the dynamics of  $\text{SGD}^{(\tilde{\eta}, \tilde{B})}$  is  $\kappa$  times faster than that of  $\text{SGD}^{(\eta, B)}$ .

We conclude the discussion with a simple example of how deriving a scaling rule from the SDE itself inevitably leads to the wrong conclusion. We define the following algorithm which is inspired by AdamW and which we dub SGD<sub>W</sub>:

$$x_{k+1} = x_k - \eta \nabla f_{\gamma_k}(x_k) - \eta \theta x_k. \quad (244)$$

**Lemma E.3.** *The SDE of SGDW is*

$$dX_t = -\nabla f(X_t)dt + \sqrt{\frac{\eta}{B}} \Sigma(X_t)^{\frac{1}{2}} dW_t - \theta X_t dt. \quad (245)$$

Therefore, one would naively deduce that to keep the SDE unchanged, one can simply use the linear scaling rule of SGD and leave  $\theta$  unaltered. However, one can easily derive the upper bound on the expected loss for a convex quadratic function and observe that to preserve that, it is imperative to scale  $\theta$  by  $\kappa$  as well.

We thus conclude that:

1. Eq. 240 is a better model for  $\text{SGD}^{(\bar{\eta}, \bar{B})}$  as it represents the dynamics more accurately;
2. Maintaining the shape of the SDE does not preserve the path properties of the optimizer;
3. Deriving a scaling rule uniquely from the SDE might lead to the wrong conclusions in the general case.

*Remark E.4.* We highlight that Theorem 5.3 of Malladi et al. (2022) claimed to have *formally* derived one for RMSprop: In line with (Jastrzebski et al., 2018), they argue that if they were to find a scaling rule that would leave their SDE unchanged, this would imply that even the dynamics of the iterates of RMSprop itself would be unchanged. First, we remind the reader that an SDE is formally defined as an *equation that drives the dynamics plus an initial condition* (See (Karatzas and Shreve, 2014), Section 5). While their scaling rule does leave the *equation unchanged*, it *alters the initial condition*, thus *changing the SDE* itself: This invalidates their claim and proof. Second, contrary to their claim, the rule is only valid near convergence as their SDE is only valid there. Third, Lemma E.2 offers a shred of concrete evidence that keeping the SDE unchanged does not imply that the path properties of the optimizers are preserved. Fourth, Lemma E.3 is a piece of concrete evidence that deriving scaling rules directly and naively from the SDE might lead to the wrong conclusions.

## F EXPERIMENTS

In this section, we provide the modeling choices and instructions to replicate our experiments. The code is implemented in Python 3 (Van Rossum and Drake, 2009) mainly using Numpy (Harris et al., 2020), scikit-learn (Pedregosa et al., 2011), and JAX (Bradbury et al., 2018).

### F.1 SIGNSGD: SDE VALIDATION (FIGURE 1)

In this subsection, we describe the experiments we run to produce Figure 1: The loss dynamics of SignSGD and that of our SDE match on average.

**DNN on Breast Cancer Dataset (Dua and Graff, 2017)** This paragraph refers to the *left* of Figure 1. The DNN has 10 dense layers with 20 neurons each activated with a ReLU. We minimize the binary cross-entropy loss. We run SignSGD for 50000 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 1$ . The learning rate is  $\eta = 0.001$ . Similarly, we integrate the SignSGD SDE (Eq. 7) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results are averaged over 3 runs and the shaded areas are the average  $\pm$  the standard deviation.

**CNN on MNIST (Deng, 2012)** This paragraph refers to the *center-left* of Figure 1. The CNN has a (3, 3, 32) convolutional layer with stride 1, followed by a ReLU activation, a (2, 2) max pool layer with stride (2, 2), a (3, 3, 32) convolutional layer with stride 1, a ReLU activation, a (2, 2) max pool layer with stride (2, 2). Then the activations are flattened and passed through a dense layer that compresses them into 128 dimensions, a final ReLU activation, and a final dense layer into the output dimension 10. The output finally goes through a softmax as we minimize the cross-entropy loss. We run SignSGD for 60000 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 0.4$ . The learning rate is  $\eta = 0.001$ . Similarly, we integrate the SignSGD SDE (Eq. 7) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results are averaged over 3 run and the shaded areas are the average  $\pm$  the standard deviation.

**Transformer on MNIST** This paragraph refers to the *center-right* of Figure 1. The Architecture is a scaled-down version of (Dosovitskiy et al., 2021), where the hyperparameters are *patch size*=28, *out features*=10, *width*=48, *depth*=3, *num heads*=6, and *dim ffn*=192. We minimize the cross-entropy loss as we run SignSGD for 5000 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 1$ . The learning rate is  $\eta = 0.001$ . Similarly, we integrate the SignSGD SDE (Eq. 7) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results are averaged over 3 runs and the shaded areas are the average  $\pm$  the standard deviation.

**ResNet on CIFAR-10 (Krizhevsky et al., 2009)** This paragraph refers to the *right* of Figure 1. The ResNet has a (3, 3, 32) convolutional layer with stride 1, followed by a ReLu activation, a second (3, 3, 32) convolutional layer with stride 1, followed by a residual connection from the first convolutional layer, then a (2, 2) max pool layer with stride (2, 2). Then the activations are flattened and passed through a dense layer that compresses them into 128 dimensions, a final ReLu activation, and a final dense layer into the output dimension 10. The output finally goes through a softmax as we minimize the cross-entropy loss. We run SignSGD for 5000 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 1$ . The learning rate is  $\eta = 0.001$ . Similarly, we integrate the SignSGD SDE (Eq. 7) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results are averaged over 3 runs and the shaded areas are the average  $\pm$  the standard deviation.

## F.2 SIGNSGD: INSIGHTS VALIDATION (FIGURE 2)

In this subsection, we describe the experiments we run to produce Figure 2: We successfully validate them all.

**Phases: Lemma 3.4 and Lemma 3.5** In this paragraph, we describe how we validated the existence of the phases of SignSGD as predicted in Lemma 3.4 and Lemma 3.5. To produce the *left* of Figure 2), we simulated the *full SDE* (Eq. 24) and the one describing Phase 3 (Eq. 5). The optimized function is  $f(x) = \frac{x^\top H x}{2}$  for  $H = \text{diag}(1, 2)$ ,  $x_0$  drawn (and fixed for all runs) from a normal distribution  $\mathcal{N}(0, 0.01)$ ,  $\eta = 0.001$ , and  $\Sigma = \sigma^2 I_d$  where  $\sigma = 0.1$ . We integrate the SDEs with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$  and for 3000 iterations. Results are averaged over 500 runs and the shaded areas are the average  $\pm$  the standard deviation. Clearly, the two SDEs share the same dynamics.

To produce the *center-left* of Figure 2, we repeat the above as  $x_0$  drawn (and fixed for all runs) from a normal distribution  $\mathcal{N}(0, 1)$ . Then, we plot the average loss values together with the theoretical prediction of Phase 1 and Phase 3: They perfectly overlap.

**Stationary distribution: Lemma 3.7** In this paragraph, we describe how we validated the convergence behavior predicted in Lemma 3.7. To produce the *center-right* of Figure 2), we run SignSGD on  $f(x) = \frac{x^\top H x}{2}$  for  $H = \text{diag}(1, 2)$ ,  $x_0 = (0.001, 0.001)$ ,  $\eta = 0.001$  and  $\Sigma = \sigma^2 I_d$  where  $\sigma = 0.1$ . We run this for 5000 times and report the evolution of the moments. Then, we add lines representing the theoretical predictions derived in Lemma 3.7: They match.

**Schedulers: Lemma 3.9** In this paragraph, we describe how we validated the convergence behavior predicted in Lemma 3.9. To produce the *right* of Figure 2, we run SignSGD on  $f(x) = \frac{x^\top H x}{2}$  for  $H = \text{diag}(1, 2)$ ,  $x_0 = (0.01, 0.01)$ ,  $\eta = 0.01$  and  $\Sigma = \sigma^2 I_d$  where  $\sigma = 0.1$ . We used the scheduler  $\eta_t^\vartheta = \frac{1}{(t+1)^\vartheta}$  for  $\vartheta \in \{0.1, 0.5, 1.5\}$ . For the first two choices of  $\vartheta$ ,  $\eta_t^\vartheta$  satisfies our sufficient condition for the convergence of SignSGD: In the figure, we observe that indeed SignSGD converges to 0 with the same speed as the one predicted in the Lemma. For  $\vartheta = 1.5$ , we observe that SignSGD does not converge following the theoretical curve because it does not satisfy our sufficient condition. Results are averaged over 500 runs.

## F.3 RMSPROP: SDE VALIDATION (FIGURE 9 AND FIGURE 10)

In this subsection, we describe the experiments we run to produce Figure 9 and Figure 10: The dynamics of our SDE matches that of RMSprop better than the SDE derived in (Malladi et al., 2022).

**Quadratic convex function** This paragraph refers to the *left* and *center-left* of Figure 9. We optimize the function  $f(x) = \frac{x^\top Hx}{2}$  where  $H = \text{diag}(10, 2)$ . We run RMSprop for 2000 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 0.1$ . The learning rate is  $\eta = 0.01$ ,  $\beta = 0.99$ . Similarly, we integrate our RMSprop SDE (Eq. 129) and that of Malladi (Eq. 227) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results are averaged over 500 runs and the shaded areas are the average  $\pm$  the standard deviation: Our SDE matches RMSprop much better.

**Embedded saddle** This paragraph refers to the *center-right* and *right* of Figure 9. We optimize the function  $f(x) = \frac{x^\top Hx}{2} + \frac{1}{4}\lambda \sum_{i=1}^2 x_i^4 - \frac{\xi}{3} \sum_{i=1}^2 x_i^3$  where  $H = \text{diag}(-1, 2)$ ,  $\lambda = 1$ , and  $\xi = 0.1$ . We run RMSprop for 1600 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 0.01$ . The learning rate is  $\eta = 0.01$ ,  $\beta = 0.99$ . Similarly, we integrate our RMSprop SDE (Eq. 129) and that of Malladi (Eq. 227) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results are averaged over 500 runs and the shaded areas are the average  $\pm$  the standard deviation: Our SDE matches RMSprop much better.

**DNN on Breast Cancer Dataset** This paragraph refers to the *left* of Figure 10. The architecture and loss are the same as used above for SignSGD. We run RMSprop for 2000 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 10^{-2}$ . The learning rate is  $\eta = 10^{-4}$ ,  $\beta = 0.9995$ . Similarly, we integrate our RMSprop SDE (Eq. 129) and that of Malladi (Eq. 227) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results are averaged over 3 runs and the shaded areas are the average  $\pm$  the standard deviation: Our SDE matches RMSprop much better.

**CNN on MNIST** This paragraph refers to the *center-left* of Figure 10. The architecture and loss are the same as used above for SignSGD. We run RMSprop for 100000 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 10^{-2}$ . The learning rate is  $\eta = 10^{-4}$ ,  $\beta = 0.999$ . Similarly, we integrate our RMSprop SDE (Eq. 129) and that of Malladi (Eq. 227) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results are averaged over 3 run and the shaded areas are the average  $\pm$  the standard deviation: Our SDE matches RMSprop much better.

**Transformer on MNIST** This paragraph refers to the *center-right* of Figure 10. The architecture and loss are the same as used above for SignSGD. We run RMSprop for 2000 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 10^{-2}$ . The learning rate is  $\eta = 10^{-3}$ ,  $\beta = 0.995$ . Similarly, we integrate our RMSprop SDE (Eq. 129) and that of Malladi (Eq. 227) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results are averaged over 3 runs and the shaded areas are the average  $\pm$  the standard deviation: Our SDE matches RMSprop much better.

**ResNet on CIFAR-10** This paragraph refers to the *right* of Figure 10. The architecture and loss are the same as used above for SignSGD. We run RMSprop for 500 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 10^{-4}$ . The learning rate is  $\eta = 10^{-4}$ ,  $\beta = 0.9999$ . Similarly, we integrate our RMSprop SDE (Eq. 129) and that of Malladi (Eq. 227) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results are averaged over 3 runs and the shaded areas are the average  $\pm$  the standard deviation: Our SDE matches RMSprop much better.

#### F.4 ADAM: SDE VALIDATION (FIGURE 12 AND FIGURE 13)

In this subsection, we describe the experiments we run to produce Figure 13 and Figure 12: The dynamics of our SDE matches that of Adam better than that derived in (Malladi et al., 2022).

**Quadratic convex function** This paragraph refers to the *left* and *center-left* of Figure 12. We optimize the function  $f(x) = \frac{x^\top Hx}{2}$  where  $H = \text{diag}(10, 2)$ . We run Adam for 50000 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 0.01$ . The learning rate is  $\eta = 0.001$ ,  $\beta_1 = 0.9$ , and  $\beta_2 = 0.999$ . Similarly, we integrate our Adam SDE (Eq. 167) and that of Malladi (Eq. 232) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results

are averaged over 500 runs and the shaded areas are the average  $\pm$  the standard deviation: Our SDE matches Adam much better.

**Embedded saddle** This paragraph refers to the *center-right* and *right* of Figure 12. We optimize the function  $f(x) = \frac{x^\top Hx}{2} + \frac{1}{4}\lambda \sum_{i=1}^2 x_i^4 - \frac{\xi}{3} \sum_{i=1}^2 x_i^3$  where  $H = \text{diag}(-1, 2)$ ,  $\lambda = 1$ , and  $\xi = 0.1$ . We run Adam as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 0.1$ . The learning rate is  $\eta = 0.001$ ,  $\beta_1 = 0.9$ , and  $\beta_2 = 0.999$ . Similarly, we integrate our Adam SDE (Eq. 167) and that of Malladi (Eq. 232) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results are averaged over 500 runs and the shaded areas are the average  $\pm$  the standard deviation: Our SDE matches Adam much better.

**DNN on Breast Cancer Dataset** This paragraph refers to the *left* of Figure 13. The architecture and loss are the same as used above for SignSGD. We run Adam for 2000 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 10^{-2}$ . The learning rate is  $\eta = 10^{-4}$ ,  $\beta_1 = 0.99$ , and  $\beta_2 = 0.999$ . Similarly, we integrate our Adam SDE (Eq. 167) and that of Malladi (Eq. 232) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results are averaged over 3 runs and the shaded areas are the average  $\pm$  the standard deviation: Our SDE matches Adam much better.

**CNN on MNIST** This paragraph refers to the *center-left* of Figure 13. The architecture and loss are the same as used above for SignSGD. We run Adam for 40000 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 10^{-2}$ . The learning rate is  $\eta = 10^{-3}$ ,  $\beta_1 = 0.99$ , and  $\beta_2 = 0.999$ . Similarly, we integrate our Adam SDE (Eq. 167) and that of Malladi (Eq. 232) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results are averaged over 3 runs and the shaded areas are the average  $\pm$  the standard deviation: Our SDE matches Adam much better.

**Transformer on MNIST** This paragraph refers to the *center-right* of Figure 13. The architecture and loss are the same as used above for SignSGD. We run Adam for 2000 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 10^{-2}$ . The learning rate is  $\eta = 10^{-2}$ ,  $\beta_1 = 0.9$ , and  $\beta_2 = 0.99$ . Similarly, we integrate our Adam SDE (Eq. 167) and that of Malladi (Eq. 232) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results are averaged over 3 runs and the shaded areas are the average  $\pm$  the standard deviation: Our SDE matches Adam much better.

**ResNet on CIFAR-10** This paragraph refers to the *right* of Figure 13. The architecture and loss are the same as used above for SignSGD. We run Adam for 2000 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 10^{-5}$ . The learning rate is  $\eta = 10^{-5}$ ,  $\beta_1 = 0.99$ , and  $\beta_2 = 0.9999$ . Similarly, we integrate our Adam SDE (Eq. 167) and that of Malladi (Eq. 232) with Euler-Maruyama (Algorithm 1) with  $\Delta t = \eta$ . Results are averaged over 3 runs and the shaded areas are the average  $\pm$  the standard deviation: Our SDE matches Adam much better.

## F.5 RMSPROPW & ADAMW: SDE VALIDATION (FIGURE 3, FIGURE 4)

The settings are exactly the same as those for RMSprop and Adam. The regularization parameter used is always  $\theta = 0.01$ . We observe that our SDEs match the respective algorithm with a good agreement.

## F.6 RMSPROPW & ADAMW: INSIGHTS VALIDATION (FIGURE 5)

In this subsection, we describe the experiments we run to produce Figure 5: The theoretically predicted asymptotic loss value and moments of RMSpropW and AdamW match those empirically found.

**Asymptotic loss & scaling rule of AdamW** This paragraph refers to the *left* of Figure 5. We optimize the function  $f(x) = \frac{x^\top Hx}{2}$  where  $H = \text{diag}(1, 3)$ . We run AdamW for 20000 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 1$ . The

learning rate is  $\eta = 0.001$ ,  $\beta_1 = 0.9$ , and  $\beta_2 = 0.999$ . Experiments are run for both  $\theta = 1$  and  $\theta = 4$ . The rescaled versions of the algorithms *AdamWR* follow the novel scaling rule with  $\kappa = 2$ . *AdamWNR* follows the scaling rule but not for  $\theta$  which is left unchanged. We plot the evolution of the loss values with the theoretical predictions of Lemma C.50: Results are averaged over 500 runs.

**Asymptotic loss & scaling rule of RMSpropW** This paragraph refers to the *center-left* of Figure 5: The only difference with the previous paragraph is that we use RMSpropW with  $\beta = 0.999$ .

**AdamW: the role of the  $\beta$ s** This paragraph refers to the *center-right* of Figure 5. We optimize the function  $f(x) = \frac{x^\top Hx}{2} + \frac{1}{4}\lambda \sum_{i=1}^2 x_i^4 - \frac{\xi}{3} \sum_{i=1}^2 x_i^3$  where  $H = \text{diag}(-1, 2)$ ,  $\lambda = 1$ , and  $\xi = 0.1$ . We run AdamW as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 0.1$ . The learning rate is  $\eta = 0.001$ ,  $\theta = 0.1$ ,  $\beta_1 \in \{0.99, 0.999\}$ , and  $\beta_2 \in \{0.992, 0.996, 0.998\}$ : Clearly, three combinations go into a minimum and three go into the other. For each minimum, the three optimizers converge to the same asymptotic loss value independently on the values of  $\beta_1$  and  $\beta_2$ . We argue that  $\beta_1$ , and  $\beta_2$  select the basin and the speed of convergence, not the asymptotic loss value: This is consistent with Lemma 3.13.

**Stationary distribution** This paragraph refers to the *right* of Figure 5. We optimize the function  $f(x) = \frac{x^\top Hx}{2}$  where  $H = \text{diag}(1, 3)$ . We run Adam for 20000 epochs as we calculate the full gradient and inject it with Gaussian noise  $Z \sim \mathcal{N}(0, \sigma^2 I_d)$  where  $\sigma = 0.01$ . The learning rate is  $\eta = 0.001$ ,  $\theta = 4$ ,  $\beta = 0.999$ ,  $\beta_1 = 0.9$ , and  $\beta_2 = 0.999$ . We plot the evolution of the average variances with the theoretical predictions of Lemma C.45 and Lemma 3.14: Results are averaged over 100 runs.

## F.7 EFFECT OF NOISE - VALIDATION (FIGURE 6 AND FIGURE 14)

In this subsection, we describe the experiments run to produce Figure 6 and Figure 14: All bounds on the asymptotic expected loss value for SGD, SignSGD, Adam, and AdamW, and Adam on an  $L^2$ -regularized loss are perfectly verified.

We optimize the loss  $f(x) = \frac{x^\top Hx}{2}$  where  $H = \text{diag}(1, 1)$  as we run each optimizer for 100000 iterations with  $\eta = 0.01$ . We repeat this procedure five times, one for each  $\sigma \in \{0.01, 0.1, 1, 10, 100\}$ . As we train, we inject noise on the gradient as distributed as  $\mathcal{N}(0, \sigma^2 I_d)$ . We plot the average loss together with the respective limits predicted by our Lemmas. For each optimizer and each  $\sigma$ , the average asymptotic loss matches the predicted limit. Therefore, we verify that the loss of SGD scales quadratically in  $\sigma$ , that of Adam on  $f(x)$ , Adam on  $f(x) + \frac{\theta \|x\|_2^2}{2}$ , and SignSGD scales linearly, and that of AdamW is limited in  $\sigma$ .

## F.8 INCREASING WEIGHT DECAY WITH THE BATCH SIZE

The analysis of Malladi et al. (2022) suggests that, when scaling batch size  $B$  by a factor  $\kappa$  one has to scale up ( $\uparrow$ ) the learning rate  $\eta$  by a factor  $\sqrt{\kappa}$  and scale down ( $\downarrow$ )  $\beta_2$  to the value  $1 - \kappa(1 - \beta_2)$ . Our SDE analysis confirms similar rules (Lemma 3.13) but additionally suggests scaling up the decoupled weight decay parameter  $\theta$  by a factor  $\sqrt{\kappa}$ . We test this in the language modeling setting utilizing a Pythia-like 160M parameter transformer architecture (Biderman et al., 2023) trained on 2.5B tokens from the SlimPajama dataset. We run a few experiments and study the match between trajectories at batch-size of 130k or 520k tokens (factor 4). We test this for single runs with (at small batch)  $\beta_2 = 0.99$  and  $\eta \in [0.0005, 0.001, 0.002, 0.004, 0.008]$ ,  $\theta \in [0.005, 0.1, 0.2, 0.4, 0.8]$ . We scale such hyperparameters as prescribed by our law and compare the results with a box plot in Figure 15. Our results suggest that our novel weight decay scaling rule brings trajectories at different batch sizes closer together on average.



3024  
 3025  
 3026  
 3027  
 3028  
 3029  
 3030  
 3031  
 3032  
 3033  
 3034  
 3035  
 3036  
 3037  
 3038  
 3039  
 3040  
 3041  
 3042  
 3043  
 3044  
 3045  
 3046  
 3047  
 3048  
 3049  
 3050  
 3051  
 3052  
 3053  
 3054  
 3055  
 3056  
 3057  
 3058  
 3059  
 3060  
 3061  
 3062  
 3063  
 3064  
 3065  
 3066  
 3067  
 3068  
 3069  
 3070  
 3071  
 3072  
 3073  
 3074  
 3075  
 3076  
 3077

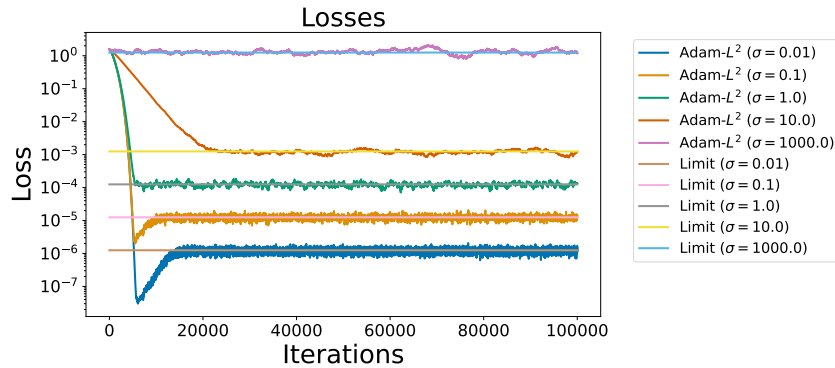


Figure 14: Empirical validation of the bounds for Adam on an  $L^2$ -regularized loss  $f(x) + \frac{\theta \|x\|_2^2}{2}$ . For several levels of noise  $\sigma$ , we find that our theoretical predictions match the experimental results: The loss levels scale linearly in  $\sigma$ .

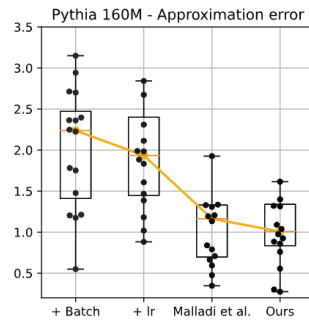


Figure 15: Trajectory match at different batch sizes as a function of the chosen scaling rule. “+batch” refers to scaling the batch size without altering hyperparameters. “+lr” additionally scales the learning rate with a square root rule (i.e. multiplied by 2).