Stochastic Resetting Mitigates Latent Gradient Bias of SGD from Label Noise

Youngkyoung Bae^{1*} Yeongwoo Song^{2*} Hawoong Jeong^{2, 3†}

¹Department of Physics and Astronomy, Seoul National University ²Department of Physics, KAIST ³Center for Complex Systems, KAIST tardis_95@snu.ac.kr, ywsong1025@kaist.ac.kr, hjeong@kaist.edu

Abstract

Giving up and starting over may seem wasteful in many situations such as searching for a target or training deep neural networks (DNNs). Our study, though, demonstrates that resetting from a checkpoint can significantly improve generalization performance when training DNNs with noisy labels. In the presence of noisy labels, DNNs initially learn the general patterns of the data but then gradually memorize the corrupted data, leading to overfitting. By deconstructing the dynamics of stochastic gradient descent (SGD), we identify the behavior of a latent gradient bias induced by noisy labels, which harms generalization. To mitigate this negative effect, we apply the stochastic resetting method to SGD, inspired by recent developments in the field of statistical physics achieving efficient target searches. We first theoretically identify the conditions where resetting becomes beneficial, and then we empirically validate our theory, confirming the significant improvements achieved by resetting. We further demonstrate that our method is both easy to implement and compatible with other methods for handling noisy labels. Additionally, this work offers insights into the learning dynamics of DNNs from an interpretability perspective, expanding the potential to analyze training methods through the lens of statistical physics.

1 Introduction

When we explore a search space having complex choices of training schemes or search for appropriate hyperparameters of deep neural networks (DNNs), we often meet circumstances that cause us to give up and train the network all over again. This is akin to our experiences in daily life, where we face various tasks that require solving problems through hit-and-miss. For example, when trying to find a beloved one's face in a crowd, our eyes typically flick back to a certain starting point after scanning the surrounding area. Similarly, when searching for a misplaced wallet after a big night out, we often fail to locate it and restart our search from some original location. These patterns are also frequently observed in animal behavior, such as foraging for food and returning to familiar locations such as nests or dens. In these situations, one might think that revisiting places is a waste of time and resources, potentially diminishing search performance. However, recent developments in statistical physics have proven that resetting to the start or a mid-point can improve the performance of the search process, meaning that this strategy is not so haphazard after all.

This effect of resetting from a particular configuration has been extensively investigated in the field of statistical physics in recent years [1, 2]. These investigations typically involve a blind searcher who evolves their current state stochastically over time without knowledge of the target's location. Surprisingly, it has been found that resetting does not hinder the search process but rather can make

^{*}Equal Contribution.

[†]Correspondence to Hawoong Jeong.



Figure 1: (a) Schematic of stochastic gradient descent (SGD) dynamics with stochastic resetting. The network parameter vector $\boldsymbol{\theta}$ evolves via SGD to find an optimal value $\boldsymbol{\theta}^*$ on the training risk landscape $\mathcal{R}_{\tilde{\mathcal{D}}_{tr}}$ (upper colormap), which differs from the true risk landscape $\mathcal{R}_{\mathcal{D}}$ (lower colormap) due to corrupted data. Here, $\boldsymbol{\theta}$ resets to the checkpoint $\boldsymbol{\theta}_c$ (home icon) with the reset probability r and resets to $\boldsymbol{\theta}_c$. (b) Fraction of correctly predicted data with wrong labels during training with SGD (gray) and SGD with reset (green). The inset shows the validation losses during training.

the searcher more efficient across diverse conditions, including scenarios with high dimensions or the presence of external forces [3–11]. Capitalizing on the success of the resetting strategy, numerous algorithms incorporating this approach have begun to emerge in diverse fields such as molecular dynamics simulations [12, 13] and queuing systems [14].

In parallel, statistical physics has emerged as a valuable framework for understanding the nature of DNNs, offering insights that are both explainable and interpretable [15, 16]. Several techniques rooted in statistical physics, such as spin-glass theories [17, 18] and analytical tools for stochastic systems [19, 20], have been recently applied to advance the understanding of DNNs. Yet it remains unclear how statistical physics, in addition to its ability to analyze the learning process of DNNs, can be effectively utilized to enhance their performance, especially in practical scenarios including low-quality datasets.

In this work, we propose applying the stochastic resetting strategy to supervised learning with noisy labels and show that it can prevent overfitting to corrupted data (also called the memorization effect). During network training, our method resets the model parameters to a checkpoint with a certain probability and restarts the training process [Fig. 1(a)]. By mapping the stochastic gradient descent (SGD) dynamics to the corresponding Langevin dynamics, we explore in-depth to understand the mechanisms and conditions by which resetting can help SGD find the optimal parameters. Our main contributions are summarized as follows.

- We reveal a latent gradient bias in the SGD dynamics induced by noisy labels, which drives the memorization effect in DNNs. Based on this finding, we apply the stochastic resetting method to counteract this effect and explain its core beneficial mechanism (Sec. 3).
- We analyze the key factors for applying the stochastic resetting method. First, we discuss the means of selecting preferable checkpoints to reset to. Then, both theoretically and empirically, we find that the improvements of resetting increase as the stochasticity of the SGD dynamics and the proportion of corrupted training data increase (Sec. 4.1, 4.2).
- We show that the resetting method can be seamlessly integrated into existing approaches and consistently improves generalization performance across several standard benchmark datasets, including those incorporating real-world noise (Sec. 4.4).

2 Related Works

2.1 Search processes and stochastic resetting in statistical physics

Search processes are ubiquitous across various domains, spanning from systems in nature to applications in engineering. For instance, ligands exhibit search processes as they navigate toward target binding sites within proteins [21–23], and similarly, predators employ search strategies to locate their prey in the wild [24, 25]. In engineering, search processes are relevant to finding primary

research studies [26], ranking web pages [27], and determining optimal hyperparameters for training algorithms [28]. Although diverse search strategies are employed depending on the problem at hand, they share a common goal: to identify an efficient search protocol. Efficiency is typically assessed by the time required to reach a target, referred to as the first passage time (FPT) in the context of random walk literature [29]. Numerous search strategies have been investigated to achieve this goal, including the Lévy strategies [30, 31], self-avoiding walks [32, 33], intermittent strategies [34], persistent random walks [35], and more [36]. One recent strategy that has garnered attention is stochastic resetting, with studies showcasing its ability to enhance search performance by revisiting previous places [2, 7, 37, 10]. In particular, these studies have demonstrated that stochastic restarts prevent a random searcher from wandering too far, thereby ensuring a finite mean time to find a target, whereas the mean time is infinite for a diffusive particle without resetting. Drawing from this concept, we introduce a resetting method for training DNNs and illustrate its effectiveness in addressing the noisy label problem.

2.2 Deep learning from noisy labels

While the accessibility of large datasets has propelled remarkable advancements in DNNs, the presence of noisy labels within these datasets often leads to erroneous model prediction [38]. Specifically, DNNs tend to overfit the entire corrupted training dataset by memorizing the wrong labels, which degenerates their generalization performance on a test dataset. Numerous studies have been conducted to address this overfitting phenomenon [39-44], and it has been revealed that DNNs initially learn the clean data (general patterns) during an early learning stage and then gradually memorize the corrupted data (task-specific patterns) [45–47]. The overfitting issue stemming from the memorization effect can be seen in Fig. 1(b), where the model's accuracy in predicting the true labels of the data with noisy labels exhibits an inverted U-shaped curve as the model progressively memorizes the noise. Based on this understanding, the surprising effectiveness of the early-stopping method [48] in alleviating the memorization effect becomes evident; as such, various methods have been proposed to leverage this insight, including Co-teaching [49], SELFIE [50], early learning regularization (ELR) [51], and robust early learning [52]. Our proposed method also capitalizes on this insight by enabling the DNN to reset to a checkpoint, i.e., previously visited parameters during early learning stages, and implicitly serves as a regularization mechanism by indirectly affecting the SGD dynamics. Additionally, our theoretical analysis explores how label noise affects the performance of DNNs from the perspective of optimization strategies, such as those discussed in Refs. [48, 47]. We adopt a different approach based on statistical physics, assuming more practical settings, which offers novel insights and leads to the development of our method. In Sec. 3, we provide a detailed analysis identifying a latent gradient bias of SGD due to label noise that causes the memorization effect, and how stochastic resetting can mitigate such negative effects.

3 Methodology

In this section, we first investigate the SGD dynamics in the presence of label noise and identify the latent gradient bias that leads to memorizing corrupted labels. We then introduce how stochastic resetting can be incorporated into SGD and demonstrate how this improves generalization performance by approximating SGD dynamics into Langevin dynamics.

Problem setup. Consider a *c*-class classification problem, which is a supervised learning task aimed at training a function to map input features to labels through a DNN. Let $\mathcal{X} \subset \mathbb{R}^p$ be the feature space, $\mathcal{Y} = \{0,1\}^c$ be the label space in one-hot vector form, and $f_{\theta} : \mathcal{X} \to \mathcal{Y}$ be a DNN model where $\theta \in \mathbb{R}^d$ encompasses all trainable parameters in the DNN. The goal is to find an optimal θ^* such that f_{θ^*} accurately assigns labels to corresponding input features, given an unknown joint probability distribution $P_{\mathcal{D}}$ over $\mathcal{X} \times \mathcal{Y}$ [Fig. 1(a)]. To obtain this, a training algorithm is applied to minimize the risk $\mathcal{R}_{\mathcal{D}}(\theta) \equiv \langle \mathcal{L}(x, y; \theta) \rangle_{\mathcal{D}}$ during training, where \mathcal{L} denotes a loss function (e.g., cross-entropy loss). Here, $\mathcal{L}(x, y; \theta)$ denotes the loss for a sample (x, y) from $P_{\mathcal{D}}$ with a given model f_{θ} , and $\langle \cdot \rangle_{\mathcal{D}}$ denotes the average over $P_{\mathcal{D}}$. In a typical classification problem, the DNN is trained by minimizing the risk on the training dataset \mathcal{D}_{tr} via SGD and θ^* is selected at the minimum risk on the validation dataset \mathcal{D}_{val} to mitigate overfitting on \mathcal{D}_{tr} , where $\mathcal{D}_{tr(val)} = \{(x_i, y_i)\}_{i=1}^{N_{tr(val)}}$ and each (x_i, y_i) is sampled from $P_{\mathcal{D}}$. Empirically, the risk on the training (validation) dataset is computed as $\mathcal{R}_{\mathcal{D}_{tr(val)}}(\theta) = (1/N_{tr(val)}) \sum_{i=1}^{N_{tr(val)}} \mathcal{L}(x_i, y_i; \theta)$. In the presence of noisy labels,



Figure 2: (a) Schematic of $-\nabla_{\theta} \mathcal{R}_{\tilde{D}_{tr}}(\theta)$, decomposed by two orthogonal terms $-\nabla_{\theta} \mathcal{R}_{\tilde{D}_{tr}^{c}}(\theta)$ and $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{D}_{tr}^{w}}(\theta)$. (b) Cosine similarity between $-\nabla_{\theta} \mathcal{R}_{\tilde{D}_{tr}}(\theta)$ and $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{D}_{tr}^{w}}(\theta)$ (cos ϕ_{tw} ; red), and between $-\nabla_{\theta} \mathcal{R}_{\tilde{D}_{tr}^{c}}(\theta)$ and $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{D}_{tr}^{w}}(\theta)$ (grey), throughout all training iterations for varying noise rate τ . (c) Magnitude difference between the two vectors $\|\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{D}_{tr}^{w}}(\theta)\| - \|\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{D}_{tr}^{c}}(\theta)\|$ throughout all training iterations for varying τ . Here, we set the batch size to B = 8 in Setting 1 described in Sec. 4. Darker colors represent larger values of τ in (b, c).

suppose we have a corrupted training dataset $\tilde{\mathcal{D}}_{tr} = \{(\boldsymbol{x}_i, \tilde{\boldsymbol{y}}_i)\}_{i=1}^{N_{tr}}$, where $\tilde{\boldsymbol{y}}$ is a noisy label that may be corrupted from a ground truth label \boldsymbol{y}_i , and $(\boldsymbol{x}_i, \tilde{\boldsymbol{y}}_i)$ is sampled from the corrupted distribution $P_{\tilde{\mathcal{D}}}$. This corrupted dataset can be partitioned into two subsets, i.e., $\tilde{\mathcal{D}}_{tr} \equiv [\tilde{\mathcal{D}}_{tr}^c, \tilde{\mathcal{D}}_{tr}^w]$, where $\tilde{\mathcal{D}}_{tr}^c$ $(\tilde{\mathcal{D}}_{tr}^w)$ consists of N_{tr}^c (N_{tr}^w) samples with correct (wrong) labels. Note that $N_{tr}^c = (1 - \tau)N_{tr}$ and $N_{tr}^w = \tau N_{tr}$ for an unknown noise rate $\tau \in [0, 1]$.

3.1 Latent gradient bias in SGD by label noise

When we apply the minibatch SGD to minimize the empirical risk $\mathcal{R}_{\mathcal{D}_{tr}}(\theta)$ with respect to θ , the update rules of θ at each training iteration t can be represented by

$$\Delta \boldsymbol{\theta}_{t} = -\frac{\eta}{B} \sum_{(\boldsymbol{x}_{i}, \boldsymbol{y}_{i}) \in \mathcal{B}_{t}} \nabla_{\boldsymbol{\theta}} \mathcal{L}_{i}(\boldsymbol{\theta}_{t})$$

$$= -\frac{\eta}{N_{\text{tr}}} \sum_{(\boldsymbol{x}_{i}, \boldsymbol{y}_{i}) \in \mathcal{D}_{\text{tr}}} \nabla_{\boldsymbol{\theta}} \mathcal{L}_{i}(\boldsymbol{\theta}_{t}) + \left(\frac{\eta}{N_{\text{tr}}} \sum_{(\boldsymbol{x}_{i}, \boldsymbol{y}_{i}) \in \mathcal{D}_{\text{tr}}} \nabla_{\boldsymbol{\theta}} \mathcal{L}_{i}(\boldsymbol{\theta}_{t}) - \frac{\eta}{B} \sum_{(\boldsymbol{x}_{i}, \boldsymbol{y}_{i}) \in \mathcal{B}_{t}} \nabla_{\boldsymbol{\theta}} \mathcal{L}_{i}(\boldsymbol{\theta}_{t})\right),$$
(1)

where θ_t is θ at the *t*-th iteration, $\Delta \theta_t \equiv \theta_{t+1} - \theta_t$, and $\mathcal{L}_i(\theta_t) \equiv \mathcal{L}(\boldsymbol{x}_i, \boldsymbol{y}_i; \theta_t)$ for simplicity. Here, $\eta > 0$ is the learning rate and \mathcal{B}_t is the minibatch of size *B* consisting of independent and identically distributed (i.i.d.) samples from \mathcal{D}_{tr} . While the first term on the right-hand-side (RHS) is deterministic for a given \mathcal{D}_{tr} , the second term on the RHS is stochastic due to the randomly sampled batch at each iteration. Thus, Eq. (1) can be rewritten as

$$\Delta \boldsymbol{\theta}_t = -\boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{D}_{\mathrm{tr}}}(\boldsymbol{\theta}_t) \boldsymbol{\eta} + \boldsymbol{\xi}_t \sqrt{\boldsymbol{\eta}},\tag{2}$$

where a random noise vector $\boldsymbol{\xi}_t \equiv \sqrt{\eta} (\boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{D}_{tr}}(\boldsymbol{\theta}_t) - \boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{B}_t}(\boldsymbol{\theta}_t)) \in \mathbb{R}^d$ satisfies $\langle \boldsymbol{\xi}_t \rangle_{\mathcal{D}_{tr}} = \mathbf{0}$ and $\langle \boldsymbol{\xi}_t \boldsymbol{\xi}_s^{\mathrm{T}} \rangle_{\mathcal{D}_{tr}} = 2\mathrm{D}(\boldsymbol{\theta}_t)\delta_{ts}$ with $\mathrm{D}(\boldsymbol{\theta}_t) \equiv \eta \Sigma(\boldsymbol{\theta}_t)/(2B)$, where δ_{ij} denotes the Kronecker delta (see details in the Supplementary Materials (SM) and Refs. [53–55]). In terms of Langevin dynamics, $\mathcal{R}_{\mathcal{D}_{tr}}(\boldsymbol{\theta})$ and $\mathrm{D}(\boldsymbol{\theta})$ correspond to the potential and diffusion matrix, respectively, where the former generates the deterministic long-term trend called drift and the latter determines the level of stochasticity of the system [56]. As a result, the SGD dynamics of $\boldsymbol{\theta}_t$ can be understood by the Langevin dynamics of a *d*-dimensional particle diffusing with drift $-\boldsymbol{\nabla}_{\boldsymbol{\theta}}\mathcal{R}_{\mathcal{D}_{tr}}(\boldsymbol{\theta}_t)$ and diffusion matrix $\mathrm{D}(\boldsymbol{\theta}_t)$.

For a corrupted dataset $\hat{\mathcal{D}}_{tr}$, the equation of SGD dynamics remains analogous to Eq. (2) when we substitute \mathcal{D}_{tr} with $\tilde{\mathcal{D}}_{tr}$. Then we can divide the drift vector $-\nabla_{\theta}\mathcal{R}_{\tilde{\mathcal{D}}_{tr}}(\theta_t)$ into two components, one originating from $\tilde{\mathcal{D}}_{tr}^c$ and the other from $\tilde{\mathcal{D}}_{tr}^w$ as follows [Fig. 2(a)]:

$$\Delta \boldsymbol{\theta}_{t} = -\left[\boldsymbol{\nabla}_{\boldsymbol{\theta}} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{tr}^{c}}(\boldsymbol{\theta}_{t}) + \boldsymbol{\nabla}_{\boldsymbol{\theta}} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{tr}^{w}}(\boldsymbol{\theta}_{t})\right] \boldsymbol{\eta} + \boldsymbol{\xi}_{t} \sqrt{\boldsymbol{\eta}}, \tag{3}$$

with the gradients from the correct part $\nabla_{\theta} \hat{\mathcal{R}}_{\hat{\mathcal{D}}_{tr}^{c}}(\theta_{t}) \equiv (1 - \tau) \nabla_{\theta} \mathcal{R}_{\hat{\mathcal{D}}_{tr}^{c}}(\theta_{t})$ and the gradients from the wrong part $\nabla_{\theta} \hat{\mathcal{R}}_{\hat{\mathcal{D}}_{tr}^{w}}(\theta_{t}) \equiv \tau \nabla_{\theta} \mathcal{R}_{\hat{\mathcal{D}}_{tr}^{w}}(\theta_{t})$. We refer to the drift from the wrong part $-\nabla_{\theta} \hat{\mathcal{R}}_{\hat{\mathcal{D}}_{tr}^{w}}(\theta_{t})$ as the *latent gradient bias* by label noise. Note that $\langle \nabla_{\theta} \mathcal{R}_{\hat{\mathcal{D}}_{tr}^{c}}(\theta_{t}) \rangle_{\mathcal{D}} = \langle \nabla_{\theta} \mathcal{R}_{\mathcal{D}_{tr}^{c}}(\theta_{t}) \rangle_{\mathcal{D}}$, implying that $-\nabla_{\theta} \hat{\mathcal{R}}_{\hat{\mathcal{D}}_{tr}^{c}}(\theta_{t})$ reflects the gradients toward the true optimum, while $-\nabla_{\theta} \hat{\mathcal{R}}_{\hat{\mathcal{D}}_{tr}^{w}}(\theta_{t})$ reflects the gradients toward the false optimum by memorizing the noisy labels. Additionally, we observe that $-\nabla_{\theta} \hat{\mathcal{R}}_{\hat{\mathcal{D}}_{tr}^{c}}(\theta_{t})$ and $-\nabla_{\theta} \hat{\mathcal{R}}_{\hat{\mathcal{D}}_{tr}^{w}}(\theta_{t})$ are orthogonal to each other [Fig. 2(b) and Fig. S.1 in the SM], leading to $-\nabla_{\theta} \mathcal{R}_{\hat{\mathcal{D}}_{tr}}(\theta_{t})$ being represented by the sum of two orthogonal vectors. Thus, $-\nabla_{\theta} \mathcal{R}_{\hat{\mathcal{D}}_{tr}}(\theta_{t})$ becomes more correlated with $-\nabla_{\theta} \hat{\mathcal{R}}_{\hat{\mathcal{D}}_{tr}^{w}}(\theta_{t})$ as the noise rate τ increases. Figure 2(b) and (c) illustrate that $-\nabla_{\theta} \hat{\mathcal{R}}_{\hat{\mathcal{D}}_{tr}^{w}}(\theta_{t})$ becomes increasingly dominant so that the drift gradually tilts toward wrong directions as τ increases, where $\cos \phi_{tw}$ denotes the cosine similarity between $-\nabla_{\theta} \hat{\mathcal{R}}_{\hat{\mathcal{D}}_{tr}^{w}}(\theta_{t})$ and $-\nabla_{\theta} \hat{\mathcal{R}}_{\hat{\mathcal{D}}_{tr}^{w}}(\theta_{t})$, and $\|\cdot\|$ denotes the Euclidean norm of a vector. This gradually tilting trend toward a wrong direction can also be observed with respect to iteration t, implying that $-\nabla_{\theta} \hat{\mathcal{R}}_{\hat{\mathcal{D}}_{tr}^{w}}(\theta_{t})$ becomes increasingly dominant as the learning process progresses beyond an early learning phase [51]. Therefore, in the presence of noisy labels, we can see that the latent gradient bias emerges and hinders the search for the optimal parameters θ^* .

We note that a similar analysis of the SGD dynamics using statistical physics was previously performed in Ref. [47]. While both studies observed an increasing trend in the effect of latent gradient bias during the learning process, their findings on the cosine similarity between drift components, $\cos \phi_{cw}$, differ from ours. To address this discrepancy, we conducted additional experiments and identified the contributing factors, as detailed in Sec. D.2 of the SM.

3.2 Stochastic resetting method

We now describe how the stochastic resetting method can be integrated into SGD. Based on this, we establish the specific premises of this work. Let θ_c be a reset checkpoint and r be the reset probability at each iteration t, where a checkpoint refers to previously visited model parameters during training. By incorporating the resetting method, Eq. (2) for $\tilde{\mathcal{D}}_{tr}$ can be expressed as

$$\boldsymbol{\theta}_{t+1} = \begin{cases} \boldsymbol{\theta}_c, & \text{with probability } r, \\ \boldsymbol{\theta}_t - \boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{R}_{\tilde{\mathcal{D}}_{tr}}(\boldsymbol{\theta}_t) \eta + \boldsymbol{\xi}_t \sqrt{\eta}, & \text{otherwise, Eq. (3).} \end{cases}$$
(4)

Below we provide the pseudo-code for SGD with stochastic resetting, Algorithm 1. Our implementation is publicly available at https://github.com/qodudrud/stochastic-resetting.

The SGD dynamics with stochastic resetting involves two processes: resetting from a checkpoint θ_c with probability r [top of Eq. (4)], and maintaining the SGD dynamics with probability 1 - r [bottom of Eq. (4)]. Note that Eq. (4) shares the same form as the (overdamped) Langevin equation with Poissonian reset [3], and also that training DNNs to find optimal parameters can be likened to a search process for an unknown target. These parallels imply that similar advantages of stochastic resetting may arise in the training process of DNNs as in the random search process of Langevin dynamics.

Algorithm 1 Stochastic resetting

Require: Corrupted training set \mathcal{D}_{tr} , validation set \mathcal{D}_{val} , reset probability r, threshold \mathcal{T} . 1: Initialize θ_0 and set t = 0, $\theta_c =$ None, and $\theta_{\text{best}} =$ None 2: **for** t = 0 to *T* **do** Update $\boldsymbol{\theta}_t \leftarrow \boldsymbol{\theta}_t - \frac{\eta}{B} \sum_{(\boldsymbol{x}_i, \boldsymbol{y}_i) \in \mathcal{B}_t} \boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{L}_i(\boldsymbol{\theta}_t)$ where \mathcal{B}_t is a randomly sampled batch from $\tilde{\mathcal{D}}_{tr}$ 3: if $\theta_c \neq$ None and rand(0,1) < r then 4: 5: Restart $\boldsymbol{\theta}_t \leftarrow \boldsymbol{\theta}_c$ end if 6: $\theta_{\text{best}} \leftarrow \text{Valid}(\theta_t, \mathcal{D}_{\text{val}})$ where $\text{Valid}(\theta_t, \mathcal{D}_{\text{val}})$ checks whether $\mathcal{R}_{\mathcal{D}_{\text{val}}}(\theta_t)$ is the minimum. 7: if θ_{best} remains unchanged for \mathcal{T} iterations or $\theta_t = \theta_{\text{best}}$ then 8: 9: Set the checkpoint $\theta_c \leftarrow \theta_{\text{best}}$ 10: end if 11: end for



Figure 3: The mean first passage time (MFPT) $\langle T(\gamma) \rangle$ from Eq. (5) with varying (a) diffusion coefficient D and (b) drift v with respect to the reset rate γ . Markers represent the minimum MFPT, $\langle T(\gamma^*) \rangle$, at the optimal reset rate γ^* . We set v = 1 in (a), D = 1 in (b), and L = 1 in both.

To examine this hypothesis, we first explore the beneficial mechanism of stochastic resetting in the search process with a simplified case.

The search efficiency of a random search process is typically quantified by the mean first passage time (MFPT) [29], which represents the average time to find a target. It has been well-established that incorporating stochastic resetting can significantly reduce the MFPT in various complex scenarios, including high-dimensional spaces [1], various confining potentials [57, 7], and a searcher with momentum [6, 58], among others [2]. For a simplified case, let us consider a random searcher in one dimension with diffusion coefficient D, drift v, and reset rate $\gamma \equiv r/\Delta t$ with a time interval Δt between steps), and assume that the searcher starts at the origin (reset point). Then the MFPT for a target located at L (> 0) can be expressed by

$$\langle T(\gamma) \rangle = \frac{1}{\gamma} \left[e^{\frac{L}{2D} \left(\sqrt{v^2 + 4D\gamma} - v \right)} - 1 \right], \tag{5}$$

where $\langle T(\gamma) \rangle$ denotes the MFPT with the reset rate γ (see the derivation in Sec. B of the SM). Examining Eq. (5) provides several insights into the effects of stochastic resetting. When the random searcher either normally diffuses or drifts away from the target ($v \leq 0$), it is straightforward that $\langle T(\gamma) \rangle$ diverges as $\gamma \to 0$ but becomes finite for any $\gamma > 0$. Conversely, when the searcher drifts toward the target (v > 0), while $\langle T(\gamma) \rangle$ is finite as L/v without resetting, introducing stochastic resetting can significantly reduce $\langle T(\gamma) \rangle$ within a certain range of γ , provided the following condition is met:

$$Pe \equiv \frac{Lv}{2D} \le 1.$$
(6)

Here, Pe is known as the Péclet number, which quantifies the ratio between drift and diffusive transport rates, and the beneficial condition (Pe ≤ 1) can be identified by verifying where $[d\langle T(\gamma)\rangle/d\gamma]|_{\gamma\to 0} < 0$. Figure 3 illustrates how the behavior of $\langle T(\gamma)\rangle$ evolves with varying Pe: as Pe decreases, the initially monotonically increasing curve gradually transforms into a U-shaped curve, achieving a minimum $\langle T(\gamma)\rangle$ at the optimal reset rate $\gamma^* > 0$. These findings indicate that resetting is advantageous for a target search when the stochasticity (D) is sufficiently larger than the drift toward a target $\langle v \rangle$. Furthermore, we note that both the optimal reset rate γ^* and the improvement ratio $\langle T(0)\rangle/\langle T(\gamma^*)\rangle$ increase as Pe decreases (Fig. 3), indicating that the benefits of resetting grow as D increases and v decreases.

According to the above observations, the beneficial properties of resetting in a random search process can be summarized as follows:

When the stochasticity is sufficiently larger than the drift toward a target, resetting can be beneficial for random searches and there exists an optimal reset probability.

The advantageous mechanism of resetting in a random search process is to suppress trajectories that move away from the target, thus increasing the chances of finding it. Drawing a parallel between the random search process and the training procedure of DNNs via SGD, we hypothesize this mechanism is also applicable to the noisy label problem in DNNs. Specifically, in supervised learning with noisy labels, the stochasticity of the SGD dynamics increases as batch size *B* decreases, and the drift component toward a target $-\nabla_{\theta} \hat{\mathcal{R}}_{\hat{\mathcal{D}}_{r}^{c}}(\theta_t)$ weakens (i.e., the latent gradient bias strengthens) as the noise rate τ increases. Therefore, the boxed statement suggests that the resetting strategy would be beneficial for searching for optimal parameters in cases with a small batch size and large noise rate in the presence of noisy labels. In Sec. 4, we perform several experiments and empirically show that stochastic resetting enhances generalization performance.

4 Experiments

This section presents the experimental results that support our theory. We perform image classification tasks with noisy labels in the following settings.

Setting 1 (Sec. 4.1, 4.2, and 4.3). To examine the impact of stochastic resetting on the noisy label problem, we first utilize a small dataset called ciFAIR-10 [59], a variant of CIFAR-10 [60]. We employ a vanilla convolutional neural network (VCNN, see Sec. C.1 in the SM) to facilitate straightforward testing of our claims. Training is performed using cross-entropy loss, an SGD optimizer with a fixed learning rate of 10^{-2} , and threshold T = 1000 iterations for the stochastic resetting method. A clean validation set, D_{val} , is used to select the best model and monitor the validation loss during training.

Setting 2 (Sec. 4.4). We assess the generalization performance of our method on two benchmark datasets, CIFAR-10 and CIFAR-100 [60], as well as its compatibility with various existing methods. The model architecture used is ResNet-34 [61], trained with SGD using a momentum of 0.9 and threshold $\mathcal{T} = 5000$ iterations for the stochastic resetting method as default. Additional details for the choice of hyperparameters are provided in Sec. C.2 of the SM. To demonstrate the efficacy of our method, we compare test accuracy with and without resetting. Note that the optimizer and learning rate scheduler do not restart throughout this experiment. To consider practical situations, a corrupted validation set, $\tilde{\mathcal{D}}_{val}$, is used for model selection and validation loss monitoring during training.

Setting 3 (Sec. 4.5). Under the same parameter conditions as Setting 2, we evaluate the performance of the stochastic resetting method on real-world noisy datasets, beyond the synthetic noise scenarios in Settings 1 and 2. Specifically, we use CIFAR-10N/100N, which are controllable, easy-to-use, and moderately sized real-world noisy datasets designed to enable fair comparisons across different benchmarks within accessible computational resources [62]. These datasets contain real-world human annotation errors obtained from Amazon Mechanical Turk. Additionally, we test on ANIMAL-10N, a dataset created by web-crawling image pairs of visually similar animals (e.g., cat and lynx, jaguar and cheetah) [50], to further examine its effectiveness across diverse real-world datasets.

In Settings 1 and 2, we apply symmetric noise with a noise rate τ , where each label in *c* classes is randomly flipped to an incorrect label in other classes with equal probability $\tau/(c-1)$. Note that all results are obtained from the model at the optimal iteration based on minimum validation loss as default, and also that the resulting test accuracy is evaluated on the clean validation set, i.e., the test dataset D_{te} is set to D_{val} .

We use the relative difference in validation loss (RDVLoss) and the relative difference in test accuracy (RDTAcc.) as metrics to indicate the relative improvement compared to the baseline. This metric enables effective comparison of the performance difference between stochastic resetting and original training. These metrics are calculated by $[v(r) - v_{\text{base}}]/v_{\text{base}}$, where v(r) is the resulting value with the reset probability r and v_{base} is the baseline value obtained from the original training $(v_{\text{base}} = v(0))$. The unnormalized results can be found in Sec. D.4 of the SM. We repeated our experiments five times across all settings to report the average and standard error values.

4.1 Which checkpoint would be preferable to reset to?

To introduce the resetting strategy in DNN training, we first explore which checkpoint is suitable to reset to in order to find optimal parameters. A straightforward choice is to select the parameters at the overfitting iteration t_m . Here, the overfitting iteration t_m refers to an iteration where the validation loss ceases to decrease and begins to increase due to the memorization effect [inset of Fig. 1(b)]. The checkpoint at t_m , denoted by θ_m , has the minimum validation loss during training when the double descent phenomenon does not occur [63], and is typically employed as an early-stopping point. Instead of early stopping and considering θ_m as the final model, we utilize θ_m as the checkpoint θ_c to reset to (Algorithm 1), leading to significantly improved results [Fig. 4(a)]. Here, θ_c is initially set to θ_m and adaptively changes to the parameter at a newly found minimum validation loss during training. As can be seen in Fig. 4(a), resetting suppresses the trajectory of θ to be near θ_c , which



Figure 4: (a) Test accuracies of the SGD (gray) and the SGD with our resetting method (green) during training. The inset shows the validation losses. (b,c) Relative difference of validation loss (RDVLoss) with varying the checkpoint to reset to with respect to the reset probability r. In (b), based on the checkpoint at the overfitting iteration t_m , RDVLoss is obtained in earlier iterations (left) and later iterations than t_m (right). $t_m + \delta t$ denotes the iteration where the checkpoint is selected. In (c), RDVLoss is plotted with the perturbed checkpoint parameters $\theta_{c,\epsilon} \equiv \theta_c + \epsilon \hat{n}$, where θ_c denotes the checkpoint and \hat{n} denotes a random unit vector. The shaded areas denote the standard error.

successfully prevents memorizing the noisy labels and increases the chance to find more appropriate parameters. It is important to note that the reset probability r controls the degree of suppression and the results at r = 0 and r = 1 are almost the same due to the overfitting phenomenon (the case of r = 1 corresponds to the early-stopping method). Therefore, the resulting RDVLoss curve for r should be U-shaped, indicating that an optimal r exists to optimize the performance.

While we simply select θ_m as the initial θ_c and adaptively update it, one may ask what effect the choice of θ_c has. To check this, we experiment with a fixed checkpoint both earlier and later than θ_m . For earlier checkpoints [left panel in Fig. 4(b)], the improvement over resetting from θ_m slightly decreases and the value of the optimal r gets smaller as the checkpoint gets earlier. In contrast, for later checkpoints [right panel in Fig. 4(b)], the improvement over resetting from θ_m significantly decreases as the checkpoint gets later. These results support our understanding of the beneficial mechanism of resetting in increasing the chance of finding better parameters, because the chance would decrease as the model memorizes more noise. Therefore, we can conclude that resetting in early learning stages is a good choice: the more memorization occurs, the smaller the improvement.

We additionally experiment to verify how the effect of resetting changes with the distance between the (adaptive) checkpoint θ_c at the minimum validation loss and a perturbed checkpoint $\theta_{c,\epsilon}$. Here, we set the perturbed checkpoint by adding the perturbation $\epsilon \hat{n}$ into θ_c with varying the perturbation magnitude ϵ , where $\hat{n} \equiv n/||n|| \in \mathbb{R}^d$ is a random unit vector with a standard normal random vector n. As shown in Fig. 4(c), it is observed that the benefits of resetting decrease as the distance between the checkpoint to reset to and θ_m increases. This result also supports that while resetting can improve the generalization performance, the choice of checkpoint to reset to can affect the performance, and that the minimum validation loss point is a good choice.

4.2 Impact of stochasticity and drift on stochastic resetting

As the statement highlighted in Sec. 3.2 clarifies, it has been proven in the statistical physics field that the resetting strategy can improve search efficiency as the stochasticity becomes larger than the drift component toward a target. In this section, we verify whether this statement is also valid in the training of DNNs and show under what circumstances resetting is more effective than not resetting.

It is first important to note that the stochasticity and the drift toward a target, i.e., $-\nabla_{\theta} \mathcal{R}_{\tilde{\mathcal{D}}_{tr}^c}$, can be controlled by the batch size B and the noise rate τ , respectively, as illustrated in Sec. 3.1. Particularly, $D(\theta_t) \propto 1/B$ and $\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{tr}^c}(\theta_t) \propto 1 - \tau$, meaning that the stochasticity increases and the drift toward a target decreases as B decreases and τ increases, respectively. We quantitatively examine the RDVLoss and the RDTAcc. values with varying B and τ with respect to the reset probability r. Remarkably, the improvements of RDVLoss and RDTAcc. with resetting become more significant as B decreases [Fig. 5(a)] and τ increases [Fig. 5(b)]. These observations strongly support our claim



Figure 5: Relative difference of validation loss (RDVLoss, left) and relative difference of test accuracy (RDTAcc., right) results with (a) varying the batch size B, and (b) varying the noise rate τ with respect to the reset probability r. We set $\tau = 0.4$ in (a) and B = 16 in (b). The shaded areas denote the standard error.

that stochastic resetting offers more benefits as the stochasticity increases and the drift toward a target decreases.

Moreover, we expect that the optimal reset probability r^* decreases as B decreases and τ increases, but this can only be verified qualitatively because the fluctuation of the results makes it difficult to identify r^* .

4.3 Ablation study on partial resetting

Until now, we have leveraged the memorization effect in our algorithm by utilizing the parameters at the minimum validation loss as the checkpoint for resetting the entire network, a process referred to as full resetting. However, several studies have highlighted that different layers within a DNN exhibit varied learning behaviors, leading to distinct levels of overfitting across these layers [64, 65]. A prevailing explanation for this phenomenon suggests that gradients tend to weaken as they propagate from the latter layers (closer to the output layer) to the former layers (closer to the input layer).

Here, we experiment on which layers, former or latter, play a more dominant role in improving performance with the resetting method. For this, we introduce partial resetting, which involves resetting only one section of the network layers rather than the entire network, while the remaining section of layers continues to follow the standard SGD update rule without resetting. We divide the VCNN structure into former and latter sections, comprising convolutional and linear layers, respectively, and apply partial resetting to one section. Interestingly, our experiments reveal that partial resetting of the latter section can further enhance generalization performance compared to full resetting, whereas partial resetting of the former section does not yield improvements over the case with no resetting (r = 0) (Fig. 6).

Moreover, even when we freeze the latter section by setting r = 1, partial resetting of the latter section still achieves significant improvement. We attribute these findings to a well-established observation: the former layers of CNNs tend to learn general features, while the latter layers tend to specialize in capturing specific features [66–69, 64]. In other words, the latter section composed of linear layers exhibits strong memorization of the corrupted data, leading to improved performances even when we freeze the latter section (r = 1), whereas the former section composed of convolutional layers focuses on learning general features, leading to no improvements even with resetting.

It is essential to note that although our ablation study suggests the effectiveness of partial reset-



Figure 6: Relative difference of validation loss (RDVLoss) and relative difference of test accuracy (RDTAcc.) results with varying one section of the network to reset with respect to the reset probability r. We set $\tau = 0.4$ and B = 16. The shaded areas denote the standard error.

Dataset	Mathad	Noise rate $\tau = 0.2$		Noise r	ate $\tau = 0.4$	Noise rate $\tau = 0.6$	
	Methou	No	Reset	No	Reset	No	Reset
	CE	84.8 ± 0.4	$90.0 \pm 0.5^{***}$	80.8 ± 0.7	$86.3 \pm 0.7^{***}$	72.9 ± 0.7	$79.5 \pm 0.9^{***}$
	Part		$88.5 \pm 0.4^{***}$	_	$83.7 \pm 0.9^{***}$		$75.1 \pm 0.4^{***}$
	MAE	91.0 ± 0.2	91.0 ± 0.2	85.6 ± 3.3	85.9 ± 3.3	67.1 ± 6.8	67.3 ± 6.7
CIFAR-10	GCE	90.6 ± 0.1	$91.0\pm0.2^{*}$	85.3 ± 0.4	$87.1 \pm 0.3^{***}$	76.2 ± 0.5	$79.1 \pm 0.8^{***}$
	SL	91.3 ± 0.3	91.4 ± 0.1	86.6 ± 0.2	$87.9 \pm 0.2^{***}$	80.1 ± 0.5	$81.6 \pm 0.6^{**}$
	ELR*	91.2 ± 0.1	91.3 ± 0.2	88.8 ± 0.2	88.9 ± 0.2	84.6 ± 0.5	84.6 ± 0.5
	SOP+*	94.1 ± 0.2	94.1 ± 0.2	89.9 ± 0.3	89.9 ± 0.3	85.0 ± 0.3	85.1 ± 0.3
	CE	62.7 ± 5.6	64.5 ± 1.5	45.6 ± 2.3	$56.9 \pm 3.0^{***}$	32.9 ± 1.6	$44.1 \pm 3.1^{***}$
	Part	_	64.0 ± 0.9	_	$55.5 \pm 2.1^{***}$	_	$43.5 \pm 1.5^{***}$
	MAE	19.9 ± 2.8	19.9 ± 2.8	11.0 ± 3.8	10.8 ± 3.7	6.7 ± 1.2	6.8 ± 1.3
CIFAR-100	GCE	68.3 ± 0.4	$69.3 \pm 0.2^{**}$	61.2 ± 0.6	$63.2 \pm 0.3^{***}$	50.1 ± 0.6	$53.1 \pm 0.8^{***}$
	SL	66.8 ± 0.5	$68.6 \pm 0.6^{***}$	60.4 ± 0.5	$63.1 \pm 0.7^{***}$	50.4 ± 0.8	$54.1 \pm 0.8^{***}$
	ELR*	67.4 ± 0.3	$70.4 \pm 0.3^{***}$	55.1 ± 0.6	$64.1 \pm 0.6^{***}$	45.9 ± 0.7	$54.0 \pm 1.0^{***}$
	SOP+*	72.1 ± 0.3	72.1 ± 0.3	59.4 ± 0.5	$65.7 \pm 0.9^{***}$	48.4 ± 1.9	$old 54.8 \pm 2.3^{**}$

Table 1: Test accuracies (%) on test datasets with different methods. We compare the performance without resetting (No) and with resetting (Reset) at r = 0.001. Results are presented as the average and the standard deviation. The best results are indicated in **bold** with statistical significance.

ting, our findings do not imply that partial resetting of the latter section always enhances generalization performance compared to full resetting. The extent to which each layer overfits the corrupted data depends on multiple factors, such as the network structure and the choice of loss function. Thus, determining the most effective section of the network to reset also hinges on the specific context. Future research investigating these points would be intriguing and valuable.

4.4 Results on corrupted benchmark datasets

In Setting 2, we investigate the impact of the stochastic resetting strategy on the performance of benchmark datasets, CIFAR-10 and CIFAR-100, under symmetric noise corruption using various methods. Table 1 compares the best results without resetting (No) and with resetting (Reset) at the reset probability r = 0.001. In the table, CE denotes cross-entropy loss, PartRestart denotes cross-entropy loss with partial resetting of only the last linear layer and the last two blocks of ResNet-34 (Sec. 4.3), MAE denotes robust mean absolute error [70], GCE denotes generalized cross-entropy [71], SL denotes symmetric cross-entropy loss [72], ELR denotes early-learning regularization [51], and SOP+ denotes sparse over-parameterization with consistency regularization and class-balance regularization [73]. ELR and SOP+ are representative methods to robustly train DNNs with an additional regularization term to prevent overfitting to corrupted data. However, these methods require additional hyperparameters to be fine-tuned depending on the loss landscape (e.g., dataset, model architecture), and in practical scenarios, it is often costly to find the optimal settings of these methods. In this section, we consider both optimal and non-optimal hyperparameter settings to take such practical situations into account, where we use an asterisk (*) to denote a non-optimal hyperparameter setting (i.e., ELR*, SOP+*). We provide additional details about the hyperparameters for each method in Sec. C.2 of the SM.

Remarkably, in all cases examined, our resetting method consistently achieves either at least equivalent or higher test accuracies compared to the baseline approach involving no resetting (Table 1). Results show that the extent of improvement becomes more pronounced as the noise rate increases, which supports our claim that resetting becomes more advantageous with higher noise rates. Furthermore, while the PartRestart method also obtains improved performance, it does not surpass the benefits of full resetting. We conjecture that the network structure may influence the extent of the additional improvements obtained, as the memorization effect in different layers can vary depending on the network architecture [74]. Minimal improvements are observed in the MAE results for both datasets; this is primarily because the MAE convergence is too slow to identify a suitable checkpoint for resetting, consequently resulting in few instances of resetting in many trials. For ELR and SOP+, no significant improvements are found in the performance when stochastic resetting is used with the known optimal hyperparameter settings, as shown in Table S.4 in the SM. However, for ELR* and SOP+*, incorporating the stochastic resetting method leads to either at least equivalent performance or significant improvements compared to no resetting. This can be explained by the design of ELR and SOP+ that provide regularization terms to prevent memorization, akin to reducing the effect of the

Table 2: Test accuracies (%) on test datasets with real-world datasets CIFAR-10N/100N. We compare the performance without resetting (No) and with resetting (Reset) at r = 0.001. Results are presented as the average and the standard deviation. The best results are indicated in **bold** with statistical significance.

	CIF	AR-10N	CIFAR-10N		CIF	AR-10N	CIFAR-100N	
Method	Random	$1 (\tau \sim 0.09)$	Aggregate ($\tau \sim 0.17$)		Worst ($\tau \sim 0.4$)		Noisy ($\tau \sim 0.4$)	
	No	Reset	No	Reset	No	Reset	No	Reset
CE	82.8 ± 0.5	$87.8 \pm 0.7^{***}$	86.1 ± 0.6	$90.2 \pm 0.3^{***}$	73.8 ± 1.9	$79.7 \pm 0.3^{***}$	46.7 ± 1.6	$56.0 \pm 1.3^{***}$
MAE	90.9 ± 0.3	90.9 ± 0.2	90.0 ± 0.5	90.0 ± 0.5	60.2 ± 3.6	60.2 ± 3.6	4.0 ± 1.3	4.0 ± 1.3
SL	91.7 ± 0.2	91.8 ± 0.2	89.6 ± 0.2	89.9 ± 0.4	80.6 ± 0.5	81.3 ± 0.8	52.3 ± 0.5	$56.5 \pm 0.6^{***}$
GCE	91.7 ± 0.2	91.9 ± 0.2	89.4 ± 0.4	89.8 ± 0.7	78.1 ± 0.9	$82.3 \pm 0.5^{***}$	55.8 ± 0.5	$58.0 \pm 1.2^{***}$
ELR*	89.5 ± 0.3	89.5 ± 0.3	91.3 ± 0.2	91.3 ± 0.2	81.9 ± 0.4	81.9 ± 0.5	56.9 ± 0.4	$60.5 \pm 0.5^{***}$
SOP+*	91.2 ± 0.2	91.2 ± 0.2	93.2 ± 0.1	93.0 ± 0.1	82.4 ± 0.5	82.3 ± 0.7	57.0 ± 0.4	$61.2 \pm 0.9^{***}$
$ViT-T_{rd}$	63.2 ± 0.9	$65.7 \pm 0.7^{***}$	61.0 ± 0.9	$63.6 \pm 0.3^{***}$	55.4 ± 0.2	$57.2 \pm 0.2^{**}$	28.6 ± 0.5	$32.8 \pm 0.5^{***}$
$ViT-T_{pt}$	92.6 ± 0.3	$93.6\pm0.6^*$	91.8 ± 0.5	$93.0 \pm 0.8^{**}$	85.9 ± 1.4	$87.4\pm0.7^*$	63.8 ± 1.1	$68.6 \pm 0.5^{**}$

latent gradient bias. These results indicate that stochastic resetting is compatible with well-established methods, and can also provide at least equal performance compared to the methods without resetting.

We likewise demonstrate the effectiveness of the stochastic resetting method in asymmetric (i.e., class-dependent) noise scenarios (Table S.5 in the SM), and we also verify the improvement of the validation loss when using our stochastic resetting method (Table S.6 in the SM).

4.5 Results on real-world noisy datasets

Finally, we test our stochastic resetting method on real-world noisy datasets, namely CIFAR-10N/100N (Setting 3). There are five different noise types for CIFAR-10N, namely Rand1, Rand2, Rand3, Aggregate, and Worst with a noise rate of 9.03%, 17.23%, 18.12%, 17.64%, and 40.21%, respectively, and a single noise type for CIFAR-100N with a noise rate of 40.20%.

Table 3: Test accuracies (%) on test datasets with real-world datasets ANIMAL-10N. We compare the performance without resetting (No) and with resetting (Reset) at r = 0.001.

Dataset	No	Reset
ANIMAL-10N	80.6 ± 0.6	$85.1\pm0.4^{***}$

For CIFAR-10N, we selected Rand1, Aggregate, and Worst noise types to account for various realworld noise rates. We further evaluate our method on ANIMAL-10N with the default setting of Setting 3 using cross-entropy loss. For the ANIMAL-10N dataset, there is only a single noise type while the ground-truth labels remain unknown. The noise rate τ was estimated as $\tau \sim 0.08$ using cross-validation with a grid search and $\tau \sim 0.06$ based on human inspection, respectively.

Table 2 presents the performance results without resetting (No) and with resetting (Reset) at the reset probability r = 0.001 for CIFAR-10N/100N. Similarly, the results for ANIMAL-10N are shown in Table 3. We provide additional details about the hyperparameters in Sec. C.2 of the SM. For the case of CIFAR-10N, the stochastic resetting method consistently achieves either at least equivalent or higher test accuracies compared to the baseline without resetting, similar to the results in Sec. 4.4 (Table 1). While the experiments on CIFAR-10N show some cases with minimal improvements from the stochastic resetting method, the results on CIFAR-10ON and ANIMAL-10N provide strong evidence of the advantages of applying stochastic resetting in practical scenarios. These results demonstrate that the stochastic resetting method can provide significant improvements while acting as a safeguard to maintain baseline performance.

Furthermore, to assess the applicability of our method beyond CNN-based approaches, we evaluate it on the Vision Transformer (ViT) [75]. As shown in Table 2, the resetting method consistently improves performance on ViT, denoted by ViT- T_{rd} and ViT- T_{pt} , which correspond to ViT-Tiny models [76] trained from randomly initialized and pre-trained weights, respectively. Previous studies have reported that ViTs without pre-training are more prone to overfitting, particularly in the presence of label noise [77, 78], which we also observe in Table 1. Interestingly, regardless of initialization, the resetting method results in significant performance improvements, highlighting its potential for a broader range of recent architectures and scenarios.

5 Discussion

In this work, we identified a latent gradient bias that hinders SGD from generalizing in the presence of label noise. To address this, we developed a stochastic resetting method, motivated by the success of the resetting strategy in statistical physics. By analyzing SGD dynamics through the lens of Langevin dynamics, we theoretically identified factors that influence the effectiveness of resetting, i.e., batch size and noise rate, and then experimentally confirmed the impact of these factors. Experiments showed that the resetting method consistently yields equivalent or improved performance on benchmark datasets compared to existing methods.

As our method can be implemented with minimal code changes and without additional computational costs, flexibility and ease of integration with other approaches are ensured. This simplicity also facilitates extensions to other variants, such as non-Poissonian resetting [79, 80] and state-dependent reset probability [3, 81], etc. [2]. In particular, as illustrated by the U-shaped curves in Figs. 3-5, resetting can excessively constrain the target search of the SGD dynamics for sufficiently large values of r, which can degrade training efficiency. To address this, introducing occasional larger, more dynamic jumps (e.g., those inspired by Lèvy strategies and intermittent search strategies [34]) or adjusting the reset probability based on recent performance (e.g., decreasing the reset probability when no improvement is observed over time) can enhance the performance of our method while mitigating the risk of getting stuck. Moreover, while we evaluated DNN models up to ResNet-34, we confirmed that the computational overhead, including I/O operations associated with resetting, is negligible compared to the overall training time. This makes our method scalable effectively to large-scale models, as further discussed in Sec. D.8 of the SM.

The main limitation of this work is that the proposed method may not be as effective when the double descent phenomenon occurs or the convergence of validation loss is too late, such as the MAE case in Table 1. Moreover, it is challenging to identify the optimal reset probability from a limited number of experiments. In fact, it has been observed that the coefficient of variation of the FPT is unity at the optimal reset probability in a random search problem [82, 83]. Future work investigating whether a similar relationship exists in DNN training may help to identify the optimal reset probability in practice.

We note that the resetting method shares a similar spirit with forgetting, which refers to the loss of previously acquired knowledge [84]. Similar to resetting, forgetting was initially viewed as a catastrophic phenomenon that needed to be addressed [85, 86]; however, recent studies have highlighted its benefits, leading to its use in improving network performance [87, 88]. From this perspective, resetting can be viewed as a form of forgetting memorized task-specific patterns, but unlike general forgetting that primarily erases early experiences, resetting targets the erasure of later experiences. It will be interesting to further explore the connections between resetting and forgetting in future discussions.

We anticipate that our work can influence two different research directions. First, it opens up the possibility of analyzing existing training methods from a statistical physics perspective, and second, it can pave the way to applying various new search strategies, beyond resetting, into neural network training.

Acknowledgments and Disclosure of Funding

This study was supported by the Basic Science Research Program through the National Research Foundation of Korea (Y.S. NRF Grant No. 2022R1A2B5B02001752, H.J. RS-2025-00514776). Y.B. was supported by an NRF grant funded by the Korean government (MSIT) (No. RS-2023-00278985).

References

- [1] M. R. Evans and S. N. Majumdar, "Diffusion with stochastic resetting," <u>Phys. Rev. Lett.</u>, vol. 106, p. 160601, 2011.
- M. R. Evans, S. N. Majumdar, and G. Schehr, "Stochastic resetting and applications," J. Phys. A: Math. Theor., vol. 53, no. 19, p. 193001, 2020.
- [3] M. R. Evans and S. N. Majumdar, "Diffusion with optimal resetting," <u>J. Phys. A: Math. Theor.</u>, vol. 44, no. 43, p. 435001, 2011.
- [4] M. R. Evans and S. N. Majumdar, "Diffusion with resetting in arbitrary spatial dimension," <u>J.</u> Phys. A: Math. Theor., vol. 47, no. 28, p. 285001, 2014.
- [5] S. Ray, D. Mondal, and S. Reuveni, "Péclet number governs transition to acceleratory restart in drift-diffusion," J. Phys. A: Math. Theor., vol. 52, no. 25, p. 255002, 2019.
- [6] D. Gupta, "Stochastic resetting in underdamped brownian motion," J. Stat. Mech., vol. 2019, no. 3, p. 033212, 2019.
- [7] A. Pal, L. Kuśmierz, and S. Reuveni, "Search with home returns provides advantage under high uncertainty," Phys. Rev. Res., vol. 2, p. 043174, 2020.
- [8] O. Tal-Friedman, A. Pal, A. Sekhon, S. Reuveni, and Y. Roichman, "Experimental realization of diffusion with stochastic resetting," <u>J. Phys. Chem. Lett.</u>, vol. 11, no. 17, pp. 7350–7355, 2020.
- [9] S. Ray and S. Reuveni, "Diffusion with resetting in a logarithmic potential," J. Chem. Phys., vol. 152, no. 23, p. 234110, 2020.
- [10] B. De Bruyne, S. N. Majumdar, and G. Schehr, "Optimal resetting brownian bridges via enhanced fluctuations," Phys. Rev. Lett., vol. 128, p. 200603, 2022.
- [11] A. Nagar and S. Gupta, "Stochastic resetting in interacting particle systems: a review," J. Phys. A: Math. Theor., vol. 56, no. 28, p. 283001, 2023.
- [12] O. Blumer, S. Reuveni, and B. Hirshberg, "Stochastic resetting for enhanced sampling," J. Phys. Chem. Lett., vol. 13, no. 48, pp. 11230–11236, 2022.
- [13] O. Blumer, S. Reuveni, and B. Hirshberg, "Combining stochastic resetting with metadynamics to speed-up molecular dynamics simulations," Nat. Commun., vol. 15, no. 1, p. 240, 2024.
- [14] O. L. Bonomo, A. Pal, and S. Reuveni, "Mitigating long queues and waiting times with service resetting," PNAS Nexus, vol. 1, no. 3, p. pgac070, 2022.
- [15] G. Carleo, I. Cirac, K. Cranmer, L. Daudet, M. Schuld, N. Tishby, L. Vogt-Maranto, and L. Zdeborová, "Machine learning and the physical sciences," <u>Rev. Mod. Phys.</u>, vol. 91, p. 045002, Dec 2019.
- [16] Y. Bahri, J. Kadmon, J. Pennington, S. S. Schoenholz, J. Sohl-Dickstein, and S. Ganguli, "Statistical mechanics of deep learning," <u>Ann. Rev. Condens. Matter Phys.</u>, vol. 11, no. Volume 11, 2020, pp. 501–528, 2020.
- [17] M. Baity-Jesi, L. Sagun, M. Geiger, S. Spigler, G. B. Arous, C. Cammarota, Y. LeCun, M. Wyart, and G. Biroli, "Comparing dynamics: Deep neural networks versus glassy systems," in International Conference on Machine Learning, pp. 314–323, PMLR, 2018.
- [18] D. Ghio, Y. Dandi, F. Krzakala, and L. Zdeborová, "Sampling with flows, diffusion, and autoregressive neural networks from a spin-glass perspective," <u>Proc. Natl. Acad. Sci. U.S.A</u>, vol. 121, no. 27, p. e2311810121, 2024.
- [19] N. Yang, C. Tang, and Y. Tu, "Stochastic gradient descent introduces an effective landscapedependent regularization favoring flat solutions," Phys. Rev. Lett., vol. 130, p. 237101, 2023.
- [20] G. Biroli, T. Bonnaire, V. De Bortoli, and M. Mézard, "Dynamical regimes of diffusion models," Nat. Commun., vol. 15, no. 1, p. 9957, 2024.

- [21] O. G. Berg, R. B. Winter, and P. H. Von Hippel, "Diffusion-driven mechanisms of protein translocation on nucleic acids. 1. models and theory," <u>Biochemistry</u>, vol. 20, no. 24, pp. 6929– 6948, 1981. PMID: 7317363.
- [22] M. Coppey, O. Bénichou, R. Voituriez, and M. Moreau, "Kinetics of target site localization of a protein on dna: A stochastic approach," Biophys. J., vol. 87, no. 3, pp. 1640–1649, 2004.
- [23] S. Ghosh, B. Mishra, A. B. Kolomeisky, and D. Chowdhury, "First-passage processes on a filamentous track in a dense traffic: optimizing diffusive search for a target in crowding conditions," J. Stat. Mech., vol. 2018, no. 12, p. 123209, 2018.
- [24] F. Bartumeus and J. Catalan, "Optimal search behavior and classic foraging theory," J. Phys. A: Math. Theor., vol. 42, no. 43, p. 434002, 2009.
- [25] G. M. Viswanathan, M. G. E. da Luz, E. P. Raposo, and H. E. Stanley, <u>The Physics of Foraging:</u> <u>An Introduction to Random Searches and Biological Encounters</u>. Cambridge University Press, 2011.
- [26] O. Dieste, A. Grimán, and N. Juristo, "Developing search strategies for detecting relevant experiments," Empir. Softw. Eng., vol. 14, no. 5, pp. 513–539, 2009.
- [27] S. Brin and L. Page, "The anatomy of a large-scale hypertextual web search engine," <u>Comput.</u> <u>Netw. ISDN Syst.</u>, vol. 30, no. 1, pp. 107–117, 1998. Proceedings of the Seventh International World Wide Web Conference.
- [28] T. Yu and H. Zhu, "Hyper-parameter optimization: A review of algorithms and applications," arXiv preprint arXiv:2003.05689, 2020.
- [29] S. Redner, A Guide to First-Passage Processes. Cambridge University Press, England, 2001.
- [30] G. M. Viswanathan, S. V. Buldyrev, S. Havlin, M. G. E. da Luz, E. P. Raposo, and H. E. Stanley, "Optimizing the success of random searches," Nature, vol. 401, no. 6756, pp. 911–914, 1999.
- [31] M. A. Lomholt, K. Tal, R. Metzler, and K. Joseph, "Lévy strategies in intermittent search processes are advantageous," <u>Proc. Natl. Acad. Sci. U.S.A</u>, vol. 105, no. 32, pp. 11055–11059, 2008.
- [32] D. J. Amit, G. Parisi, and L. Peliti, "Asymptotic behavior of the "true" self-avoiding walk," Phys. Rev. B, vol. 27, pp. 1635–1645, 1983.
- [33] P. Grassberger, "Self-trapping self-repelling random walks," <u>Phys. Rev. Lett.</u>, vol. 119, p. 140601, 2017.
- [34] O. Bénichou, C. Loverdo, M. Moreau, and R. Voituriez, "Intermittent search strategies," <u>Rev.</u> <u>Mod. Phys.</u>, vol. 83, pp. 81–129, 2011.
- [35] V. Tejedor, R. Voituriez, and O. Bénichou, "Optimizing persistent random searches," <u>Phys.</u> Rev. Lett., vol. 108, p. 088103, 2012.
- [36] H. Meyer and H. Rieger, "Optimal non-markovian search strategies with *n*-step memory," Phys. Rev. Lett., vol. 127, p. 070601, 2021.
- [37] Y. Bae, G. Son, and H. Jeong, "Unexpected advantages of exploitation for target searches in complex networks," Chaos, vol. 32, no. 8, p. 083118, 2022.
- [38] T. Xiao, T. Xia, Y. Yang, C. Huang, and X. Wang, "Learning from massive noisy labeled data for image classification," in <u>Proceedings of the IEEE/CVF Conference on Computer Vision</u> and Pattern Recognition (CVPR), 2015.
- [39] B. Han, Q. Yao, T. Liu, G. Niu, I. W. Tsang, J. T. Kwok, and M. Sugiyama, "A survey of labelnoise representation learning: Past, present and future," <u>arXiv preprint arXiv:2011.04406</u>, 2020.

- [40] W. Hu, Z. Li, and D. Yu, "Simple and effective regularization methods for training on noisily labeled data with generalization guarantee," in <u>International Conference on Learning</u> Representations, 2020.
- [41] S. Li, X. Xia, S. Ge, and T. Liu, "Selective-supervised contrastive learning with noisy labels," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 316–325, 2022.
- [42] H. Song, M. Kim, D. Park, Y. Shin, and J.-G. Lee, "Learning from noisy labels with deep neural networks: A survey," <u>IEEE Transactions on Neural Networks and Learning Systems</u>, vol. 34, no. 11, pp. 8135–8153, 2023.
- [43] X. Xia, B. Han, Y. Zhan, J. Yu, M. Gong, C. Gong, and T. Liu, "Combating noisy labels with sample selection by mining high-discrepancy examples," in <u>Proceedings of the IEEE/CVF</u> International Conference on Computer Vision (ICCV), pp. 1833–1843, 2023.
- [44] Z. Huang, J. Zhang, and H. Shan, "Twin contrastive learning with noisy labels," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 11661–11670, 2023.
- [45] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, "Understanding deep learning requires rethinking generalization," in <u>International Conference on Learning Representations</u>, 2017.
- [46] D. Arpit, S. Jastrzębski, N. Ballas, D. Krueger, E. Bengio, M. S. Kanwal, T. Maharaj, A. Fischer, A. Courville, Y. Bengio, and S. Lacoste-Julien, "A closer look at memorization in deep networks," in <u>Proceedings of the 34th International Conference on Machine Learning</u>, vol. 70 of Proceedings of Machine Learning Research, pp. 233–242, PMLR, 2017.
- [47] Y. Feng and Y. Tu, "Phases of learning dynamics in artificial neural networks in the absence or presence of mislabeled data," Mach. Learn.: Sci. Technol., vol. 2, no. 4, p. 043001, 2021.
- [48] M. Li, M. Soltanolkotabi, and S. Oymak, "Gradient descent with early stopping is provably robust to label noise for overparameterized neural networks," in <u>Proceedings of the Twenty Third</u> <u>International Conference on Artificial Intelligence and Statistics</u>, vol. 108 of <u>Proceedings of</u> <u>Machine Learning Research</u>, pp. 4313–4324, PMLR, 2020.
- [49] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. Tsang, and M. Sugiyama, "Co-teaching: Robust training of deep neural networks with extremely noisy labels," in <u>Advances in Neural</u> <u>Information Processing Systems</u>, vol. 31, Curran Associates, Inc., 2018.
- [50] H. Song, M. Kim, and J.-G. Lee, "SELFIE: Refurbishing unclean samples for robust deep learning," in <u>Proceedings of the 36th International Conference on Machine Learning</u>, vol. 97 of Proceedings of Machine Learning Research, pp. 5907–5915, PMLR, 2019.
- [51] S. Liu, J. Niles-Weed, N. Razavian, and C. Fernandez-Granda, "Early-learning regularization prevents memorization of noisy labels," in <u>Advances in Neural Information Processing</u> <u>Systems</u>, vol. 33, pp. 20331–20342, Curran Associates, Inc., 2020.
- [52] X. Xia, T. Liu, B. Han, C. Gong, N. Wang, Z. Ge, and Y. Chang, "Robust early-learning: Hindering the memorization of noisy labels," in <u>International Conference on Learning</u> Representations, 2021.
- [53] Q. Li, C. Tai, and W. E, "Stochastic modified equations and adaptive stochastic gradient algorithms," in <u>Proceedings of the 34th International Conference on Machine Learning</u>, vol. 70 of Proceedings of Machine Learning Research, pp. 2101–2110, PMLR, 2017.
- [54] S. L. Smith and Q. V. Le, "A bayesian perspective on generalization and stochastic gradient descent," in 6th International Conference on Learning Representations ICLR, 2018.
- [55] L. Ziyin, K. Liu, T. Mori, and M. Ueda, "Strength of minibatch noise in SGD," in <u>International</u> Conference on Learning Representations, 2022.
- [56] C. W. Gardiner, Handbook of Stochastic Methods. Springer Berlin, Heidelberg, 2004.

- [57] A. Pal, "Diffusion in a potential landscape with stochastic resetting," <u>Phys. Rev. E</u>, vol. 91, p. 012113, Jan 2015.
- [58] P. Singh, "Random acceleration process under stochastic resetting," J. Phys. A: Math. Theor., vol. 53, no. 40, p. 405005, 2020.
- [59] B. Barz and J. Denzler, "Do we train on test data? purging cifar of near-duplicates," <u>J. Imag.</u>, vol. 6, no. 6, 2020.
- [60] A. Krizhevsky, G. Hinton, et al., "Learning multiple layers of features from tiny images," 2009.
- [61] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016.
- [62] J. Wei, Z. Zhu, H. Cheng, T. Liu, G. Niu, and Y. Liu, "Learning with noisy labels revisited: A study using real-world human annotations," in <u>10th International Conference on Learning</u> Representations, ICLR 2022, Virtual, Apr 25-29, 2022, OpenReview.net, 2022.
- [63] P. Nakkiran, G. Kaplun, Y. Bansal, T. Yang, B. Barak, and I. Sutskever, "Deep double descent: Where bigger models and more data hurt," in <u>International Conference on Learning</u> <u>Representations</u>, 2020.
- [64] Y. Bai, E. Yang, B. Han, Y. Yang, J. Li, Y. Mao, G. Niu, and T. Liu, "Understanding and improving early stopping for learning with noisy labels," in <u>Advances in Neural Information</u> Processing Systems, vol. 34, pp. 24392–24403, Curran Associates, Inc., 2021.
- [65] Y. Chen, A. Yuille, and Z. Zhou, "Which layer is learning faster? a systematic exploration of layer-wise convergence rate for deep neural networks," in <u>The Eleventh International</u> Conference on Learning Representations, 2023.
- [66] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson, "How transferable are features in deep neural networks?," in <u>Advances in Neural Information Processing Systems</u>, vol. 27, Curran Associates, Inc., 2014.
- [67] G. Cohen, G. Sapiro, and R. Giryes, "Dnn or k-nn: That is the generalize vs. memorize question," arXiv preprint arXiv:1805.06822, 2018.
- [68] A. Ansuini, A. Laio, J. H. Macke, and D. Zoccolan, "Intrinsic dimension of data representations in deep neural networks," in <u>Advances in Neural Information Processing Systems</u>, vol. 32, Curran Associates, Inc., 2019.
- [69] H. Maennel, I. M. Alabdulmohsin, I. O. Tolstikhin, R. Baldock, O. Bousquet, S. Gelly, and D. Keysers, "What do neural networks learn when trained with random labels?," in <u>Advances</u> in <u>Neural Information Processing Systems</u>, vol. 33, pp. 19693–19704, Curran Associates, Inc., 2020.
- [70] A. Ghosh, H. Kumar, and P. S. Sastry, "Robust loss functions under label noise for deep neural networks," in <u>Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence</u>, AAAI'17, p. 1919–1925, AAAI Press, 2017.
- [71] Z. Zhang and M. R. Sabuncu, "Generalized cross entropy loss for training deep neural networks with noisy labels," in Proceedings of the 32nd International Conference on Neural Information <u>Processing Systems</u>, NIPS'18, (Red Hook, NY, USA), p. 8792–8802, Curran Associates Inc., 2018.
- [72] Y. Wang, X. Ma, Z. Chen, Y. Luo, J. Yi, and J. Bailey, "Symmetric cross entropy for robust learning with noisy labels," in <u>Proceedings of the IEEE/CVF International Conference on</u> Computer Vision (ICCV), 2019.
- [73] S. Liu, Z. Zhu, Q. Qu, and C. You, "Robust training under label noise by over-parameterization," in <u>Proceedings of the 39th International Conference on Machine Learning</u>, vol. 162 of Proceedings of Machine Learning Research, pp. 14153–14172, PMLR, 2022.

- [74] J. Li, M. Zhang, K. Xu, J. Dickerson, and J. Ba, "How does a neural network's architecture impact its robustness to noisy labels?," in <u>Advances in Neural Information Processing Systems</u>, vol. 34, pp. 9788–9803, Curran Associates, Inc., 2021.
- [75] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in <u>9th International Conference</u> on Learning Representations, ICLR 2021, Virtual, May 3-7, 2021, OpenReview.net, 2021.
- [76] R. Wightman, "Pytorch image models." https://github.com/huggingface/ pytorch-image-models, 2019.
- [77] X. Chen, S. Xie, and K. He, "An empirical study of training self-supervised vision transformers," in <u>Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)</u>, pp. 9620–9629, 2021.
- [78] B. Khanal, P. Shrestha, S. Amgain, B. Khanal, B. Bhattarai, and C. A. Linte, "Investigating the robustness of vision transformers against label noise in medical image classification," <u>arXiv</u> preprint arXiv:2402.16734, 2024.
- [79] A. Pal, A. Kundu, and M. R. Evans, "Diffusion under time-dependent resetting," J. Phys. A: Math. Theor., vol. 49, no. 22, p. 225001, 2016.
- [80] A. Nagar and S. Gupta, "Diffusion with stochastic resetting at power-law times," <u>Phys. Rev. E</u>, vol. 93, p. 060102, 2016.
- [81] E. Roldán and S. Gupta, "Path-integral formalism for stochastic resetting: Exactly solved examples and shortcuts to confinement," Phys. Rev. E, vol. 96, p. 022130, 2017.
- [82] S. Reuveni, "Optimal stochastic restart renders fluctuations in first passage times universal," Phys. Rev. Lett., vol. 116, p. 170601, 2016.
- [83] A. Pal and S. Reuveni, "First passage under restart," <u>Phys. Rev. Lett.</u>, vol. 118, p. 030603, 2017.
- [84] Z. Wang, E. Yang, L. Shen, and H. Huang, "A comprehensive survey of forgetting in deep learning beyond continual learning," arXiv preprint arXiv:2307.09218, 2023.
- [85] M. McCloskey and N. J. Cohen, "Catastrophic interference in connectionist networks: The sequential learning problem," vol. 24 of <u>Psychology of Learning and Motivation</u>, pp. 109–165, Academic Press, 1989.
- [86] G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, and S. Wermter, "Continual lifelong learning with neural networks: A review," Neural Netw., vol. 113, pp. 54–71, 2019.
- [87] H. Zhou, A. Vani, H. Larochelle, and A. Courville, "Fortuitous forgetting in connectionist networks," in International Conference on Learning Representations, 2022.
- [88] E. Nikishin, M. Schwarzer, P. D'Oro, P.-L. Bacon, and A. Courville, "The primacy bias in deep reinforcement learning," in <u>Proceedings of the 39th International Conference on Machine Learning</u>, vol. 162 of <u>Proceedings of Machine Learning Research</u>, pp. 16828–16847, PMLR, 2022.
- [89] P. Chaudhari and S. Soatto, "Stochastic gradient descent performs variational inference, converges to limit cycles for deep networks," in <u>2018 Information Theory and Applications</u> Workshop (ITA), pp. 1–10, 2018.
- [90] M. Gurbuzbalaban, U. Simsekli, and L. Zhu, "The heavy-tail phenomenon in SGD," in Proceedings of the 38th International Conference on Machine Learning (M. Meila and T. Zhang, eds.), vol. 139 of Proceedings of Machine Learning Research, pp. 3964–3975, PMLR, 18–24 Jul 2021.
- [91] S. Jastrzębski, Z. Kenton, D. Arpit, N. Ballas, A. Fischer, Y. Bengio, and A. Storkey, "Three factors influencing minima in SGD," <u>arXiv preprint arXiv:1711.04623</u>, 2017.

- [92] X. Li, Q. Gu, Y. Zhou, T. Chen, and A. Banerjee, <u>Hessian based analysis of SGD for Deep</u> Nets: Dynamics and Generalization, pp. 190–198.
- [93] Z. Xie, I. Sato, and M. Sugiyama, "A diffusion theory for deep learning dynamics: Stochastic gradient descent exponentially favors flat minima," in <u>International Conference on Learning</u> Representations, 2021.
- [94] R. Metzler, J.-H. Jeon, A. G. Cherstvy, and E. Barkai, "Anomalous diffusion models and their properties: non-stationarity, non-ergodicity, and ageing at the centenary of single particle tracking," Phys. Chem. Chem. Phys., vol. 16, pp. 24128–24164, 2014.
- [95] M. Villén-Altamirano, J. Villén-Altamirano, et al., "Restart: a method for accelerating rare event simulations," in Queueing, Performance and Control in ATM (ITC-13), pp. 71–76, 1991.
- [96] M. Luby, A. Sinclair, and D. Zuckerman, "Optimal speedup of las vegas algorithms," <u>Inf.</u> Process. Lett., vol. 47, no. 4, pp. 173–180, 1993.
- [97] H. Tong, C. Faloutsos, and J.-Y. Pan, "Random walk with restart: fast solutions and applications," Knowl. Inf. Syst., vol. 14, pp. 327–346, Mar 2008.
- [98] I. Loshchilov and F. Hutter, "SGDR: Stochastic gradient descent with warm restarts," in International Conference on Learning Representations, 2017.
- [99] S. Ioffe and C. Szegedy, "Batch normalization: accelerating deep network training by reducing internal covariate shift," in <u>Proceedings of the 32nd International Conference on International</u> Conference on Machine Learning - Volume 37, ICML'15, p. 448–456, JMLR.org, 2015.
- [100] H. Daneshmand, A. Joudaki, and F. Bach, "Batch normalization orthogonalizes representations in deep random networks," in <u>Advances in Neural Information Processing Systems</u>, vol. 34, pp. 4896–4906, Curran Associates, Inc., 2021.
- [101] G. Yang, J. Pennington, V. Rao, J. Sohl-Dickstein, and S. S. Schoenholz, "A mean field theory of batch normalization," in International Conference on Learning Representations, 2019.
- [102] A. Joudaki, H. Daneshmand, and F. Bach, "On bridging the gap between mean field and finite width deep random multilayer perceptron with batch normalization," in <u>Proceedings of</u> the 40th International Conference on Machine Learning, vol. 202 of <u>Proceedings of Machine</u> Learning Research, pp. 15388–15400, PMLR, 2023.
- [103] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in <u>Advances in Neural Information Processing</u> <u>Systems</u> (H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, eds.), vol. 33, pp. 1877–1901, Curran Associates, Inc., 2020.
- [104] D. Narayanan, M. Shoeybi, J. Casper, P. LeGresley, M. Patwary, V. Korthikanti, D. Vainbrand, P. Kashinkunti, J. Bernauer, B. Catanzaro, A. Phanishayee, and M. Zaharia, "Efficient large-scale language model training on gpu clusters using megatron-lm," in <u>Proceedings</u> of the International Conference for High Performance Computing, Networking, Storage and <u>Analysis</u>, SC '21, (New York, NY, USA), Association for Computing Machinery, 2021.

A SGD dynamics and Langevin dynamics

A.1 Derivation of Eq. (2)

We explain how the dynamics of SGD can be converted to the discretized Langevin equation. We follow similar procedures as in Refs. [53–55]. As written in the main text, the network parameters θ are updated by

$$\begin{aligned} \Delta \boldsymbol{\theta}_{t} &= -\frac{\eta}{B} \sum_{(\boldsymbol{x}_{i}, \boldsymbol{y}_{i}) \in \mathcal{B}_{t}} \boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{L}_{i}(\boldsymbol{\theta}_{t}) \\ &= -\frac{\eta}{N_{\text{tr}}} \sum_{(\boldsymbol{x}_{i}, \boldsymbol{y}_{i}) \in \mathcal{D}_{\text{tr}}} \boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{L}_{i}(\boldsymbol{\theta}_{t}) + \left(\frac{\eta}{N_{\text{tr}}} \sum_{(\boldsymbol{x}_{i}, \boldsymbol{y}_{i}) \in \mathcal{D}_{\text{tr}}} \boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{L}_{i}(\boldsymbol{\theta}_{t}) - \frac{\eta}{B} \sum_{(\boldsymbol{x}_{i}, \boldsymbol{y}_{i}) \in \mathcal{B}_{t}} \boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{L}_{i}(\boldsymbol{\theta}_{t}) \right) \\ &= -\boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{D}_{\text{tr}}}(\boldsymbol{\theta}_{t}) \eta + \boldsymbol{\xi}_{t} \sqrt{\eta}, \end{aligned}$$

where the random noise vector is $\boldsymbol{\xi}_t \equiv \sqrt{\eta} \left(\nabla_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{D}_{tr}}(\boldsymbol{\theta}_t) - \nabla_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{B}_t}(\boldsymbol{\theta}_t) \right) \in \mathbb{R}^d$, $\nabla_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{D}_{tr}}(\boldsymbol{\theta}_t)$ is the gradient of risk on \mathcal{D}_{tr} , and $\nabla_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{B}_t}(\boldsymbol{\theta}_t)$ is the gradient of risk on a minibatch \mathcal{B}_t . Note that $\langle \nabla_{\boldsymbol{\theta}} \mathcal{L}_i(\boldsymbol{\theta}_t) \rangle_{\mathcal{D}_{tr}} = \nabla_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{D}_{tr}}(\boldsymbol{\theta}_t)$ and $\langle \mathcal{R}_{\mathcal{B}_t}(\boldsymbol{\theta}_t) \rangle_{\mathcal{D}_{tr}} = \mathcal{R}_{\mathcal{D}_{tr}}(\boldsymbol{\theta}_t)$. Using these facts, we can obtain that $\langle \boldsymbol{\xi}_t \rangle_{\mathcal{D}_{tr}} = \mathbf{0}$ and

$$\langle \boldsymbol{\xi}_{t} \boldsymbol{\xi}_{s}^{\mathrm{T}} \rangle_{\mathcal{D}_{\mathrm{tr}}} = \eta \left\langle \left(\boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{D}_{\mathrm{tr}}}(\boldsymbol{\theta}_{t}) - \boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{B}_{t}}(\boldsymbol{\theta}_{t}) \right) \left(\boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{D}_{\mathrm{tr}}}(\boldsymbol{\theta}_{s}) - \boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{B}_{s}}(\boldsymbol{\theta}_{s}) \right)^{\mathrm{T}} \right\rangle_{\mathcal{D}_{\mathrm{tr}}}$$

$$= \eta \left(\left\langle \boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{B}_{t}}(\boldsymbol{\theta}_{t}) \boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{B}_{s}}(\boldsymbol{\theta}_{s})^{\mathrm{T}} \right\rangle_{\mathcal{D}_{\mathrm{tr}}} - \boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{D}_{\mathrm{tr}}}(\boldsymbol{\theta}_{t}) \boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{D}_{\mathrm{tr}}}(\boldsymbol{\theta}_{s})^{\mathrm{T}} \right).$$

$$(S.2)$$

For the $s \neq t$ case, $\langle \nabla_{\theta} \mathcal{R}_{\mathcal{B}_{t}}(\theta_{t}) \nabla_{\theta} \mathcal{R}_{\mathcal{B}_{s}}(\theta_{s})^{\mathrm{T}} \rangle_{\mathcal{D}_{\mathrm{tr}}} = \nabla_{\theta} \mathcal{R}_{\mathcal{D}_{\mathrm{tr}}}(\theta_{t}) \nabla_{\theta} \mathcal{R}_{\mathcal{D}_{\mathrm{tr}}}(\theta_{s})^{\mathrm{T}}$ because each minibatch is independently sampled from $\mathcal{D}_{\mathrm{tr}}$. Applying $\langle \nabla_{\theta} \mathcal{L}_{i}(\theta_{t}) \nabla_{\theta} \mathcal{L}_{j}(\theta_{t}) \rangle_{\mathcal{D}_{\mathrm{tr}}} = \|\nabla_{\theta} \mathcal{R}_{\mathcal{D}_{\mathrm{tr}}}(\theta_{t})\|^{2} + \Sigma(\theta_{t})\delta_{ij}$ with the covariance matrix $\Sigma(\theta_{t})$, where $\|\cdot\|$ denotes the Euclidean norm of a vector, we have

$$\left\langle \|\boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{B}_t}(\boldsymbol{\theta}_t)\|^2 \right\rangle_{\mathcal{D}_{\mathrm{tr}}} = \|\boldsymbol{\nabla}_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{D}_{\mathrm{tr}}}(\boldsymbol{\theta}_t)\|^2 + \frac{1}{B} \Sigma(\boldsymbol{\theta}_t).$$
 (S.3)

Therefore, we obtain $\langle \boldsymbol{\xi}_t \boldsymbol{\xi}_s^T \rangle_{\mathcal{D}_{tr}} = 2\mathsf{D}(\boldsymbol{\theta}_t)\delta_{ts}$ with $\mathsf{D}(\boldsymbol{\theta}_t) = \eta \Sigma(\boldsymbol{\theta}_t)/(2B)$. We assumed $N_{tr} \gg B$ in the above derivation, but the decreasing trend of $\mathsf{D}(\boldsymbol{\theta}_t)$ with *B* is still valid for $N_{tr} \geq B$ [55]. The noise covariance matrix D is typically state-dependent and highly anisotropic, which contributes to the heavy-tailed and nonequilibrium stationary distribution of the SGD dynamics [89, 90]. In addition, it has been shown that D is positively correlated with the Hessian matrix of the loss [91, 92], indicating that the noise strength is greater in directions where the loss landscape is sharper. This relationship suggests that SGD prefers flat minima over sharp minima [93, 19].

Let us consider the overdamped Langevin equation, a first-order stochastic differential equation describing the evolution of a particle where friction dominates over inertia. Applying the Euler method, the overdamped Langevin equation can be approximated with time interval Δt by [56]

$$\Delta \boldsymbol{x}_t = -\boldsymbol{\nabla}_{\boldsymbol{x}} V(\boldsymbol{x}_t) \Delta t + \sqrt{2\mathsf{B}(\boldsymbol{x}_t)} \Delta \boldsymbol{W}_t.$$
(S.4)

Here, x_t is the position of the particle at time step t, $V(x_t)$ is the underlying potential, $B(x_t)$ is the strength of the fluctuations, called the diffusion matrix, and ΔW_t is the random noise vector that satisfies $\langle \Delta W_t \rangle = 0$ and $\langle \Delta W_t \Delta W_s^T \rangle = \Delta t \delta_{ts}$, where $\langle \cdot \rangle$ denotes the ensemble average. When we simulate Eq. (S.4), we randomly sample the random real number ζ_t from a certain probability distribution with zero-mean and unit-variance at each iteration t and represent the noise vector as $\Delta W_t \equiv \zeta_t \sqrt{\Delta t}$. Note that ΔW_t is commonly assumed as an increment of the Wiener process based on the central limit theorem, so that ζ_t is generally sampled from the standard normal distribution. This assumption is often violated in various situations [94], and ζ_t can be sampled from other distributions depending on the system. Comparing Eq. (S.1) with Eq. (S.4), we can easily see that the dynamics of SGD follows the overdamped Langevin equation, with $\mathcal{R}_{D_{tr}}$ and D serving as the potential and the diffusion matrix, respectively. Based on this correspondence, we analyze the SGD dynamics in the language of the Langevin dynamics and apply the stochastic resetting strategy in the main text.

A.2 SGD dynamics for noisy small dataset

In Sec. 3, we derived that latent gradient bias arises caused by label noise hinders the SGD dynamics from reaching the optimal parameters. While the main text primarily focuses on the impact of label noise, other factors can similarly degrade the training of DNNs. One such factor is the limited size of the training dataset. To investigate how dataset size affects DNN training, we recall the equation governing SGD dynamics with label noise [Eq. (3) in the main text]:

$$\Delta \boldsymbol{\theta}_{t} = -\left[\boldsymbol{\nabla}_{\boldsymbol{\theta}} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{tr}^{c}}(\boldsymbol{\theta}_{t}) + \boldsymbol{\nabla}_{\boldsymbol{\theta}} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{tr}^{w}}(\boldsymbol{\theta}_{t})\right] \boldsymbol{\eta} + \boldsymbol{\xi}_{t} \sqrt{\boldsymbol{\eta}}, \tag{S.5}$$

with the label noise rate $0 \le \tau \le 1$ is defined, the gradients from the correct part $\nabla_{\boldsymbol{\theta}} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{tr}^{c}}(\boldsymbol{\theta}_{t}) \equiv (1-\tau) \nabla_{\boldsymbol{\theta}} \mathcal{R}_{\tilde{\mathcal{D}}_{tr}^{c}}(\boldsymbol{\theta}_{t})$, and the gradients from the wrong part $\nabla_{\boldsymbol{\theta}} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{tr}^{w}}(\boldsymbol{\theta}_{t}) \equiv \tau \nabla_{\boldsymbol{\theta}} \mathcal{R}_{\tilde{\mathcal{D}}_{tr}^{w}}(\boldsymbol{\theta}_{t})$. For a sufficiently large training dataset \mathcal{D}_{tr} , $\nabla_{\boldsymbol{\theta}} \mathcal{R}_{\tilde{\mathcal{D}}_{tr}^{c}}(\boldsymbol{\theta}_{t})$ closely approximates the gradients of the true risk, denoted by $\nabla_{\boldsymbol{\theta}} \mathcal{R}_{\mathcal{D}}(\boldsymbol{\theta}_{t})$. However, when \mathcal{D}_{tr} is small, a discrepancy between them emerges due to the small dataset size, which can be expressed by

$$\nabla_{\boldsymbol{\theta}} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{sm}^{c}}^{c}(\boldsymbol{\theta}_{t}) = \nabla_{\boldsymbol{\theta}} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{tr}^{c}}^{c}(\boldsymbol{\theta}_{t}) - \nabla_{\boldsymbol{\theta}} \hat{\mathcal{R}}_{\mathcal{D}}(\boldsymbol{\theta}_{t}), \qquad (S.6)$$

with $\nabla_{\theta} \hat{\mathcal{R}}_{\mathcal{D}}(\theta_t) \equiv (1 - \tau) \nabla_{\theta} \mathcal{R}_{\mathcal{D}}(\theta_t)$. Substituting Eq. (S.6) into Eq. (S.5) yields:

$$\Delta \boldsymbol{\theta}_{t} = -\boldsymbol{\nabla}_{\boldsymbol{\theta}} \hat{\mathcal{R}}_{\mathcal{D}}(\boldsymbol{\theta}_{t}) \eta - \boldsymbol{\nabla}_{\boldsymbol{\theta}} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{sm}^{c}}(\boldsymbol{\theta}_{t}) \eta - \boldsymbol{\nabla}_{\boldsymbol{\theta}} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{tr}^{w}}(\boldsymbol{\theta}_{t}) \eta + \boldsymbol{\xi}_{t} \sqrt{\eta}.$$
(S.7)

Here, the first term in the right-hand-side of Eq. (S.7) represents the drift toward the true optimum, the second term represents the bias caused by the small dataset size, and the third term represents the bias introduced by label noise. Thus, we can identify that latent gradient bias arises from either a small dataset size or label noise, both of which hinder the generalization ability of DNN. Notably, for a large \mathcal{D}_{tr} , the second term $\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{sm}^c}(\theta_t)$ diminishes, leading to $\nabla_{\theta} \hat{\mathcal{R}}_{\mathcal{D}}(\theta_t) \simeq \nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{tr}^c}(\theta_t)$. On the other hand, for $\tau = 0$ case, the third term $\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{tr}^w}(\theta_t)$ disappears, and the first and second terms dominate. Conversely, for $\tau = 1$, the third term persists, and the first and second terms vanish. These results suggest that the label noise rate influences the relative magnitudes of the drift terms, whereas the dataset size affects the direction of the drift by determining how effectively the correct labels within the dataset guide the model toward the true optimum. Our analysis mainly centers on the effects of label noise, but exploring the individual and combined impacts of $\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{sm}^c}(\theta_t)$ on DNN training would be an intriguing direction for future work. Furthermore, the substantial improvements achieved by restarting with a small noisy dataset, as demonstrated in Secs. 4.1–3, highlight the efficacy of our method in this challenging regime.

B Stochastic resetting in statistical physics

We briefly introduce what stochastic resetting is and the conditions under which this strategy can help random searches (see Ref. [2] for a more detailed review). In fact, the resetting method has already been exploited in some stochastic algorithms [95–97], but much attention has been attracted by the theoretical success of stochastic resetting [1]. Several approaches have been made to deal with search processes involving resets, such as calculating the survival probability [1, 5]. Here, we present an easy but general one, and clarify that we follow the same procedure as in Refs. [82, 83, 7].

Let us consider a generic searcher that starts from a resetting point at time zero in a *d*-dimensional space, with the assumption that the searcher resets at a rate of γ if the target is not found. In other words, the search process is completed when the searcher finds the target before resetting; otherwise, the searcher returns to the resetting point and repeats this procedure until the target is found. This procedure can be understood in terms of two random variables, T and R, which represent the time to find a target and the time to reset, respectively. If we draw T and R from their respective distributions, we check whether T > R. If T > R, the searcher resets before finding the target, and the search process begins anew from the resetting point. Conversely, if T < R, the searcher finds the target before resetting, and the search process is completed. Applying this scheme, the time to find a target, known as the first passage time (FPT) and denoted by $T(\gamma)$, can be expressed by the following renewal equation:

$$T(\gamma) = \begin{cases} T, & \text{if } T < R, \\ R + T(\gamma)', & \text{if } R \le T, \end{cases}$$
(S.8)

or equivalently,

$$T(\gamma) = \min(T, R) + I(R \le T)T(\gamma)', \tag{S.9}$$

where $T(\gamma)'$ denotes an independent and identically distributed copy of $T(\gamma)$, and $I(R \le T)$ denotes an indicator function which equals one if $R \le T$ and zero otherwise. Note that $T(\gamma) = T$ if there is no reset ($\gamma = 0$). Taking expectations in Eq. (S.9), we obtain the mean first passage time (MFPT):

$$\langle T(\gamma) \rangle = \frac{\langle \min(T, R) \rangle}{\Pr(T < R)},$$
 (S.10)

with the relations $\langle I(R \leq T) \rangle = \Pr(R \leq T)$ and $\langle T(\gamma) \rangle = \langle T(\gamma)' \rangle$. It is noteworthy that we do not assume the dynamics of the search process and the function of the reset rate. Thus, Eq. (S.10) can be applied to a general search process regardless of the distributions of T and R.

As a simple reset method, we assume a constant reset rate, i.e., the reset probability within a time interval dt is γdt [1]. Then, the distribution of R is exponential with $\gamma e^{-\gamma t}$ at time t and $\langle \min(T, R) \rangle$ can be calculated by

$$\langle \min(T,R) \rangle = \int_0^\infty dt \, [1 - \Pr(\min(T,R) \le t)]$$

= $\int_0^\infty dt \, \Pr(R > t) \Pr(T > t)$
= $\int_0^\infty dt \, e^{-rt} \int_t^\infty dt' f_T(t') = \frac{1}{r} - \frac{1}{r} \int_0^\infty dt \, e^{-rt} f_T(t).$ (S.11)

In addition, $\Pr(T < R) = \int_0^\infty dt \ f_T(t) \Pr(R > t) = \int_0^\infty dt \ e^{-rt} f_T(t)$. Substituting these equations into Eq. (S.10), we obtain

$$\langle T(\gamma) \rangle = \frac{1 - T(\gamma)}{\gamma \tilde{T}(\gamma)}$$
 (S.12)

where $\tilde{T}(\gamma) \equiv \int_0^\infty dt e^{-\gamma t} f_T(t)$ denotes the Laplace transform of T evaluated at γ . Note that $f_T(t)$ is determined by the underlying dynamics of the searcher, which may include factors such as stochasticity or external drift. Therefore, once the dynamics of a searcher are determined and $f_T(t)$ is known, we can calculate the MFPT at γ and identify whether resetting is beneficial.

To obtain the MFPT with resetting, let us specify the dynamics of a searcher. Consider a searcher diffusing in one dimension with a diffusion constant D and a constant drift v and assume that the searcher starts at the origin (reset point) and that the target we want to find is located at L (> 0). Here, D represents the stochasticity and v represents the drift toward a target. The position x(t) of the searcher at time t evolves during a small time interval Δt through the Langevin equation given by

$$x(t + \Delta t) = \begin{cases} x_r, & \text{with probability } \gamma \Delta t, \\ x(t) + v \Delta t + \xi(t) \sqrt{\Delta t}, & \text{otherwise,} \end{cases}$$
(S.13)

where x_r denotes the reset point set as the origin and $\xi(t)$ denotes a stochastic force, typically modeled Gaussian white noise satisfying $\langle \xi(t) \rangle = 0$ and $\langle \xi(t)\xi(s) \rangle = 2D\delta(t-s)$. In this search process, the probability distribution to find the searcher at position x at time t is known to be given by [56]

$$G_0(x,t) = \frac{1}{\sqrt{4\pi Dt}} \left[e^{-\frac{(x-vt)^2}{4Dt}} - e^{\frac{Lv}{D}} e^{-\frac{(x-2L-vt)^2}{4Dt}} \right].$$
 (S.14)

The FPT distribution $f_T(t)$ is then derived from the probability that the searcher has not yet reached the target by time t: $f_T(t) = \frac{d}{dt} \Pr(T \ge t) = \frac{d}{dt} \int_{-\infty}^{L} dx G_0(x,t)$. Using this equation and the definition of $\tilde{T}(\gamma)$, we find $\tilde{T}(\gamma) = 1 - \gamma \int_{-\infty}^{L} dx \tilde{G}_0(x,\gamma)$ where $\tilde{G}_0(x,\gamma) \equiv \int_0^{\infty} dt e^{-\gamma t} G_0(x,t)$ is the Laplace transform of $G_0(x,t)$ evaluated at γ . Thus, upon calculation with these equations and Eq. (S.14), we obtain

$$\tilde{T}(\gamma) = e^{\frac{L}{2D}\left(v - \sqrt{v^2 + 4D\gamma}\right)}.$$
(S.15)

Substituting Eq. (S.15) into Eq. (S.12), we finally have

$$\langle T(\gamma) \rangle = \frac{1}{\gamma} \left[e^{\frac{L}{2D} \left(\sqrt{v^2 + 4D\gamma} - v \right)} - 1 \right].$$
 (S.16)

This final expression gives the MFPT for a searcher with drift and diffusion, including the effect of resetting at a rate γ .

In our paper, we exploit the correspondence between SGD dynamics and Langevin dynamics, as described in Sec. A, and theoretically identify where stochasticity is strengthened and drift toward a target is weakened in DNN training. Although we demonstrated a simple one-dimensional case in this section, numerous studies have explored more complex scenarios, including various confining potentials [57, 7], and related contexts. [2]. Importantly, it has been proven that a diffusive particle with resetting converges to a nonequilibrium steady state even in arbitrary spatial dimensions, and it has also been shown to be beneficial for target search [3, 4]. We believe these findings support the hypothesis that resetting may provide advantages in the search process involved in SGD dynamics for DNN. However, we also acknowledge that the advantages of resetting in terms of MFPT are not directly connected to performance improvements in DNN training. We conjecture and empirically validate that similar beneficial mechanisms successfully operate in DNN training. Nonetheless, it would be necessary and intriguing to find an alternative metric to represent the DNN performance and to theoretically investigate the effect of resetting.

C Description of the experiments

This section provides details on the experiments not included in the main text. For all of the experiments, we perform 5 independent runs to achieve the average and the standard error values. All runs were made independently on a single NVIDIA TITAN V GPU. All results are obtained from the model at the optimal iteration based on minimum validation loss as default, except ELR and SOP+ methods. Additionally, the resulting test accuracy is evaluated on the clean validation set. The objective throughout our experiments is to compare the performance with and without resetting, not to achieve state-of-the-art performances; therefore, we did not heavily tune the hyperparameters for each of the settings. Here we provide the choice of hyperparameters for our experiments.

C.1 Network architecture

We employed a vanilla CNN (VCNN) as mentioned in Secs. 4.1, 4.2, and 4.3 to expedite a straightforward testing of our claims. It consists of simple layers, as outlined in Table S.1. In the table, the inclusion of batch normalization before the activation function and after a layer is indicated by the term "Use BatchNorm". The output dimension of a convolutional layer is represented as (C, W, H), where C denotes the number of channels, and W and H represent the width and height, respectively. Additionally, c represents the number of classes in the dataset.

Table S.1: Network architecture of the VCNN: Layer name, output dimension of the layer, parameters of the convolutional layer (K, P, S), and activation function. Here, K, P, and S represent the size of the filter, padding, and stride, respectively.

Layer name	Output dim.	(K, P, S)	Use BatchNorm	Activation function
Input image	(3, 32, 32)	None	Х	None
Conv2d	(32, 32, 32)	(3, 1, 1)	0	ReLU
Conv2d	(64, 32, 32)	(3, 1, 1)	0	ReLU
MaxPool2d	(64, 16, 16)	(2, 0, 2)	0	None
Conv2d	(128, 16, 16)	(3, 1, 1)	0	ReLU
Conv2d	(128, 16, 16)	(3, 1, 1)	0	ReLU
MaxPool2d	(128, 8, 8)	(2, 0, 2)	Х	None
Conv2d	(256, 8, 8)	(3, 1, 1)	0	ReLU
Conv2d	(256, 8, 8)	(3, 1, 1)	0	ReLU
MaxPool2d	(256, 4, 4)	(2, 0, 2)	Х	None
Dropout $(p = 0.2)$	$256 \times 4 \times 4$	None	Х	None
Linear	1024	None	0	ReLU
Linear	512	None	0	ReLU
Linear	c	None	Х	Softmax

C.2 Experimental setting for results on benchmark datasets in Secs. 4.4 and 4.5

We set the hyperparameters required for each baseline method in Secs. 4.4 and 4.5 following Table S.2, as specified in their original papers.

Mathad	Dataset						
Method	CIFAR-10/10N	CIFAR-100/100N					
GCE [71]	(q,k) = (0.7,0)	(q,k) = (0.7,0)					
SL [72]	$(\alpha,\beta) = (0.1,1.0)$	$(\alpha,\beta) = (6.0,0.1)$					
ELR [51]	$\lambda = 3$	$\lambda = 7$					
ELR*	$\lambda = 0.5$	$\lambda = 1$					
SOP+[73]	$(\alpha_u, \alpha_v) = (10, 10)$	$(\alpha_u, \alpha_v) = (1, 10)$					
SOP+*	$(\alpha_u, \alpha_v) = (1, 1)$	$(\alpha_u, \alpha_v) = (0.1, 1)$					

Table S.2: Hyperparameters for baseline methods used in Secs. 4.4 and 4.5.

We confirm that all results for methods without our proposed approach are consistent with the performance reported in their original papers. We set the batch size to 256, momentum to 0.9, and weight decay to 5×10^{-4} for the cross-entropy loss, MAE, GCE, and SL methods (for SL on CIFAR-100, the batch size is set to 128). For the ELR and SOP methods, we set the batch size to 128, momentum to 0.9, and weight decay to 10^{-3} , following Ref. [51, 73]. We employ a cosine annealing scheduler [98], setting the maximum number of iterations to the total iteration 5×10^4 with an initial learning rate of 0.1 for the cross-entropy loss and 0.01 for MAE [70], GCE, and SL methods. For the ELR and SOP+ methods, we set the initial learning rate as 0.02, and reduce it by 1/100 and 1/10, respectively, after 40 (80) and 80 (120) with a total epoch 150 on CIFAR-10 (CIFAR-100) as indicated in Ref. [51, 73]. Additional regularizer parameters for SOP+ and SOP+* are (λ_C, λ_B) = (0.9, 0.1). Moreover, for the ELR and SOP+ methods, we set the reset checkpoint and select the best model based on validation accuracy instead of validation loss due to the discrepancy in the resulting test accuracies between these accuracy-based and loss-based approaches in the baseline. The threshold \mathcal{T} is set to 30 epochs on ELR and SOP+. Especially in SOP+, our method starts after 80 epochs on CIFAR-100 (Table 1).

When testing on ViT [75], we use the ViT-Tiny (ViT-T) model from the Hugging Face timm library [76]. The model uses a patch size of 16, an embedding dimension of 192, a depth of 12, and 3 attention heads, with an input image size of 224. For training ViT-T, we use cross-entropy loss and SGD with a momentum of 0.9 and no weight decay. The batch size is set to 256, and the initial learning rate is 0.01. We employ a cosine annealing scheduler with the maximum number of iterations to 5×10^4 . For ANIMAL-10N, we followed the default setting of cross-entropy loss, modifying only the batch size to 128.

D Additional results

D.1 Cosine similarities between drifts in Sec. 3

We provide the cosine similarities between drift components, i.e., $-\nabla_{\theta} \mathcal{R}_{\tilde{D}_{tr}}(\theta)$ (total drift), $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{D}_{tr}^c}(\theta)$ (drift from the correct part), and $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{D}_{tr}^w}(\theta)$ (drift from the wrong part, i.e., latent gradient bias by label noise) by varying the noise rate in Fig. S.1. The shaded gray area represents the iteration region before the stochastic resetting takes place, and the white area represents the iteration region after the resetting commences. Note that the stochastic resetting does not directly change the gradient itself, as can be seen by comparing the top row (no reset) in Fig. S.1(a, b, c) and the bottom row (with reset) in Fig. S.1(d, e, f). Instead, stochastic resetting suppresses trajectories of DNNs that drift away from the optimum, thereby increasing the chances of finding it despite the presence of latent gradient bias. Here, we can see that the cosine similarity between $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{D}_{tr}^c}(\theta_t)$ and $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{D}_{tr}^w}(\theta_t)$, denoted by $\cos \phi_{cw}$, is almost zero, indicating that $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{D}_{tr}^c}(\theta_t)$ and $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{D}_{tr}^w}(\theta_t)$ are most likely to be orthogonal to each other (see the further discussion in Sec. D.2). In addition, the cosine similarity between $-\nabla_{\theta} \mathcal{R}_{\tilde{D}_{tr}}(\theta_t)$ and $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{D}_{tr}^c}(\theta_t)$, denoted by $\cos \phi_{tc}$, decreases with increasing τ , implying that $-\nabla_{\theta} \mathcal{R}_{\tilde{D}_{tr}}(\theta_t)$ becomes less correlated with $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{D}_{tr}^c}(\theta_t)$. These results support our claims in Sec. 3 and the schematic in Fig. 2(a).



Figure S.1: Cosine similarities between $-\nabla_{\theta} \mathcal{R}_{\tilde{\mathcal{D}}_{tr}}(\theta)$ (total drift), $-\nabla_{\theta} \mathcal{R}_{\tilde{\mathcal{D}}_{tr}^c}(\theta)$ (drift from the correct label, i.e., correct drift), and $-\nabla_{\theta} \mathcal{R}_{\tilde{\mathcal{D}}_{tr}^w}(\theta)$ (drift from the wrong label, i.e., wrong drift). The first row (a, b, c) shows the cosine similarities when there is no reset, r = 0, and the second row (d, e, f) shows the cosine similarities when resetting is applied, r = 0.01, with (a, d) a noise rate of $\tau = 0.2$, (b, e) $\tau = 0.4$, and (c, f) $\tau = 0.6$. Here, $\cos \phi_{tc}$, $\cos \phi_{tw}$, and $\cos \phi_{cw}$ denote the cosine similarity between total and correct drifts, total and wrong drifts, and correct and wrong drifts, respectively. The gray region represents the iterations before the resetting commences. We set B = 8 in Setting 1 and performed 5 independent runs, where the shaded areas denote the standard error.

D.2 Orthogonality between the drift components

We further explore the orthogonality of $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{t}^{c}}(\theta_{t})$ and $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{t}^{w}}(\theta_{t})$. In contrast to our results, Ref. [47] demonstrated that $\cos \phi_{cw}$ varies significantly during training and can be used as an order parameter to divide the learning phases. To investigate the factors contributing to this discrepancy, we perform additional experiments with DNNs of comparable sizes and settings to those used in Ref. [47], as detailed in Table S.3. Note that for the fully connected network (FCN) structures, the case with $d_l = 50$ number of hidden units matches the configuration used in Ref. [47]. A key difference in our setup is the use of batch normalization [99], which we hypothesize may play a significant role in the observed effects. Interestingly, as shown in Fig. S.2, we observe that batch normalization enforces orthogonality between two drift components during training in both CNN and FCN structures. Moreover, this orthogonality in the presence of batch normalization becomes more pronounced as the number of hidden units d_l increases, implying that the dimensionality of network parameters also affects the orthogonality [Fig. S.2(b,c)]. These observations are in a similar line as findings from Ref. [100], which reported that batch normalization in DNNs induces orthogonalization of hidden representations of samples across layers. As batch normalization is one of the most widely adopted techniques for improving training in DNNs, our findings suggest that the orthogonality between the drift components is likely to generalize across diverse DNN structures and experimental settings. Exploring the mechanisms by which batch normalization enforces this orthogonality would be an intriguing direction for future research, complementing theoretical studies on batch normalization, such as those in Refs. [101, 100, 102].



Figure S.2: Cosine similarities between $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{tr}^c}(\theta)$ and $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{tr}^w}(\theta)$, denoted by $\cos \phi_{cw}$, during training with (w/) and without (w/o) batch normalization (BN) for (a) CNN and (b) FCN structures. (c) $\cos \phi_{cw}$ by varying the number of hidden units d_l in the FCN structure with batch normalization. We set B = 8 in Setting 1 and performed 5 independent runs, where the shaded areas denote the standard error. Additionally, we perform a moving average with a time window of 50 log-iterations to smooth the recorded values of $\cos \phi_{cw}$, which are measured every 20 iterations.

Table S.3: Network architectures of (top) CNN and (bottom) FCN used in Sec. D.2. Here, d_l represents the number of hidden units of the FCN structure. Batch normalization, when applied, is placed before the activation function, except for the output layer.

CNN									
Layer name	Output dim.		K, P, S)	Activation	function				
Input image	(3, 32, 32)	2)	None	Non	e				
Conv2d	(8, 32, 32)	2) ((3, 1, 1)	ReL	U				
Conv2d	(8, 32, 32)	2) ((3, 1, 1)	ReL	U				
MaxPool2d	(8, 16, 16)		(2, 0, 2)	2,0,2) Non					
Dropout $(p = 0.2)$	$8 \times 16 \times$	16	None	None					
Linear	c		None	Softm	nax				
FCN									
Layer nan	e Outpu	t dim.	Activation function						
Input imag	ge 3×32	2×32	None						
Linear	d	ı	ReLU						
Linear	d	ı	ReLU						
Linear	0	;	Softmax						

D.3 Direct impact of label noise on minibatch gradients in SGD dynamics

In the previous sections, we demonstrated the impact of label noise on the drift in SGD dynamics. Here, we extend our analysis to examine the direct impact of label noise on the minibatch gradient. Using a similar decomposition as in Eq. (S.5), the minibatch gradient can be expressed as:

$$\Delta \boldsymbol{\theta}_{t} = -\frac{\eta}{B} \sum_{(\boldsymbol{x}_{i}, \boldsymbol{y}_{i}) \in \tilde{\mathcal{B}}_{t}} \nabla_{\boldsymbol{\theta}} \mathcal{L}_{i}(\boldsymbol{\theta}_{t})$$

$$= -\nabla_{\boldsymbol{\theta}} \mathcal{R}_{\tilde{\mathcal{B}}_{t}}(\boldsymbol{\theta}_{t}) \eta = -\nabla_{\boldsymbol{\theta}} \hat{\mathcal{R}}_{\tilde{\mathcal{B}}_{t}^{c}}(\boldsymbol{\theta}_{t}) \eta - \nabla_{\boldsymbol{\theta}} \hat{\mathcal{R}}_{\tilde{\mathcal{B}}_{t}^{w}}(\boldsymbol{\theta}_{t}) \eta,$$
(S.17)

where the corrupted minibatch $\tilde{\mathcal{B}}_t$ at the *t*-th iteration, sampled from $\tilde{\mathcal{D}}_{tr}$, consists of B_t^c correct labels and B_t^w wrong labels (i.e., $B_t^c + B_t^w = B$). Here, $\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{B}}_t^c}(\theta_t) \equiv (B_t^c/B) \nabla_{\theta} \mathcal{R}_{\tilde{\mathcal{B}}_t^c}(\theta_t)$ denotes the gradients from the correct part on minibatch, and $\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{B}}_t^w}(\theta_t) \equiv (B_t^w/B) \nabla_{\theta} \mathcal{R}_{\tilde{\mathcal{B}}_t^w}(\theta_t)$ denotes the gradients from the wrong part on minibatch. It is important to note that, unlike the decomposed drifts in Eq. (S.5) are determined for a given θ and training dataset $\tilde{\mathcal{D}}_{tr}$, the decomposed minibatch gradients in Eq. (S.17) are inherently stochastic due to the randomness of the minibatch sampling process. This stochasticity introduces greater fluctuations in the minibatch gradients compared to the drifts. However, the fact that the average of minibatch gradients aligns with the corresponding drifts [e.g., $\langle \nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{B}}_t^w}(\theta_t) \rangle = \nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{D}}_{tr}^w}(\theta_t)$] ensures consistent overall trends, i.e., the orthogonal relationship between $\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{B}}_t^c}$ and $\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{B}}_t^w}$ and the increasing contribution of $\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{B}}_t^w}$ with larger τ , as shown in Fig. S.3.



Figure S.3: (a) Cosine similarities between $-\nabla_{\theta} \mathcal{R}_{\tilde{\mathcal{B}}_{t}}(\theta)$ and $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{B}}_{t}^{w}}(\theta_{t})$ (cos ϕ_{tw} ; red), and between $-\nabla_{\theta} \mathcal{R}_{\tilde{\mathcal{B}}_{t}^{c}}(\theta)$ and $-\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{B}}_{t}^{w}}(\theta_{t})$ (cos ϕ_{cw} ; grey), respectively, throughout all training iterations for varying noise rate τ . (c) Magnitude difference between the two vectors $\|\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{B}}_{t}^{w}}(\theta)\| - \|\nabla_{\theta} \hat{\mathcal{R}}_{\tilde{\mathcal{B}}_{t}^{c}}(\theta)\|$ throughout all training iterations for varying τ . Here, we set the batch size to B = 8in setting 1 described in Sec. 4. Darker colors represent larger values of τ .

D.4 Complementary plots for Secs. 4.1, 4.2, 4.3, and 4.4

In order to facilitate a comparison between the resetting method and the original SGD, we normalized the validation loss and test accuracy values in the main text by calculating the relative difference in the metrics (RDVLoss and RDVAcc.). Here, we provide the unnormalized (i.e., original) validation loss and test accuracy values: Figs. S.4, S.5, and S.6 in the Supplementary Materials are the unnormalized results of Figs. 4, 5, and 6 in the main text, respectively. We additionally present test accuracies with and without resetting for varying noise rates in Fig. S.7 with the default setting of cross-entropy loss in Setting 2, and observe that resetting consistently improves performance even in the high-noise regime.



Figure S.4: (a) Validation loss and (b) test accuracy results with varying the checkpoint to reset to with respect to reset probability r. Based on the checkpoint at the overfitting iteration t_m , the results are obtained in earlier iterations (left) and later iterations than t_m (right). Here, $t_m + \delta t$ denotes the iteration where the checkpoint is selected. (c) Validation loss (left) and test accuracy (right) with the perturbed checkpoint parameters $\theta_{c,\epsilon}$. Here, $\theta_{c,\epsilon} \equiv \theta_c + \epsilon \hat{n}$ where θ_c denotes the checkpoint and \hat{n} denotes a random unit vector. The shaded areas denote the standard error.



Figure S.5: Validation loss and test accuracy results with (a) varying the batch size B, and (b) varying the noise rate τ with respect to the reset probability r. We set $\tau = 0.4$ in (a) and B = 16 in (b). The shaded areas denote the standard error.



Figure S.6: Validation loss and test accuracy results with varying one section of the network to reset to with respect to reset probability r. Here, we set $\tau = 0.4$ and B = 16. The shaded areas denote the standard error.



Figure S.7: Test accuracy results for (a) CIFAR-10 and (b) CIFAR-100 dataset by varying the noise rate τ . Inset shows the relative difference in test accuracy (RDTAcc.) to represent the relative improvement compared to r = 0.

D.5 Ablation study for ELR and SOP+

Since the hyperparameters λ and (α_u, α_v) for ELR and SOP+ baselines depend on the loss landscape (e.g., the dataset used, selection of model architecture, etc.), it is often costly to find optimal parameters in practical situations. To ensure that the stochastic resetting method can help users achieve high performance while not depending on the hyperparameter selection, we provide the test accuracy results for both optimal hyperparameter and non-optimal hyperparameter situations using the publicly available codes in Refs. [51, 73]. Table S.4 shows the test accuracies for different settings of ELR and SOP+ baselines. Here, ELR and SOP+ denote the original methods with the hyperparameter sets used in Refs. [51, 73], while ELR* and SOP+* denote the methods with different hyperparameters. The hyperparameters for these methods are listed in Table S.2.

The results corresponding to the use of optimal hyperparameters for each method (first and third rows in Table S.4, ELR and SOP+) show no significant improvement when stochastic resetting is used. On the other hand, under the non-optimal hyperparameter settings (the second and fourth rows, ELR* and SOP+*), the stochastic resetting method improves performance, especially for the CIFAR-100 dataset, which is more complex than CIFAR-10. These results suggest that stochastic resetting is compatible with existing powerful methods and can improve performance while also acting as a safeguard, regardless of the choice of hyperparameters for each method.

D.6 Test accuracies for datasets with asymmetric label noise

Table S.5 shows the test accuracy for the asymmetric noise case, while the rest of the settings remain the same as in Table 1. Similar to the results in Table 1, we find again that the stochastic resetting method makes robust improvement compared to no stochastic resetting. Similarly, results also show a similar trend in improvement as the noise rate increases, and more improvement is obtained in the CIFAR-100 dataset.

Table S.4: Test accuracies (%) on test datasets with ELR and SOP+ methods. We compare the performance without resetting (No) and with resetting (Reset) at r = 0.001. Results are presented as the average and the standard deviation. The best results are indicated in **bold** with statistical significance.

Dataset	Matha d	Noise rate $\tau = 0.2$		Noise r	ate $\tau = 0.4$	Noise rate $\tau = 0.6$		
	Method	No	Reset	No	Reset	No	Reset	
	ELR	91.9 ± 0.2	91.9 ± 0.2	90.2 ± 0.2	90.2 ± 0.2	87.3 ± 0.3	87.3 ± 0.3	
CIFAR-10	ELR*	91.2 ± 0.1	91.3 ± 0.2	88.8 ± 0.2	88.9 ± 0.2	84.6 ± 0.5	84.6 ± 0.5	
	SOP+	95.5 ± 0.1	95.5 ± 0.1	94.8 ± 0.2	94.8 ± 0.2	93.6 ± 0.2	93.5 ± 0.2	
	SOP+*	94.1 ± 0.2	94.1 ± 0.2	89.9 ± 0.3	89.9 ± 0.3	85.0 ± 0.3	85.1 ± 0.3	
	ELR	73.2 ± 0.2	73.0 ± 0.4	69.6 ± 0.4	69.6 ± 0.4	62.8 ± 0.4	62.8 ± 0.4	
CIFAR-100	ELR*	67.4 ± 0.3	$70.4 \pm 0.3^{***}$	55.1 ± 0.6	$64.1 \pm 0.6^{***}$	45.9 ± 0.7	$54.0 \pm 1.0^{***}$	
	SOP+	78.4 ± 0.2	78.5 ± 0.1	76.8 ± 0.3	76.8 ± 0.3	73.6 ± 0.6	73.6 ± 0.6	
	SOP+*	72.1 ± 0.3	72.1 ± 0.3	59.4 ± 0.5	$65.7 \pm 0.9^{***}$	48.4 ± 1.9	$54.8 \pm 2.3^{**}$	

Table S.5: Test accuracies (%) on test datasets with asymmetric noise by different methods. We compare the performance without resetting (No) and with resetting (Reset) at r = 0.001. Results are presented as the average and the standard deviation. The best results are indicated in **bold** with statistical significance.

Deterret	Made a d	Noise rate 0.1		Noise rate 0.2		Noise	e rate 0.3	Noise rate 0.4	
Dataset	Method	No	Reset	No	Reset	No	Reset	No	Reset
	CE	89.8 ± 0.3	$93.3\pm0.3^{***}$	88.5 ± 0.6	$91.6\pm0.6^{***}$	87.3 ± 1.3	$90.0\pm0.8^{***}$	83.3 ± 2.0	$86.7\pm0.8^{***}$
	Part	_	$92.6\pm0.3^{***}$	—	$90.9\pm0.4^{***}$	_	$89.6 \pm 0.6^{**}$	_	83.3 ± 2.6
	MAE	89.1 ± 4.2	89.1 ± 4.2	75.5 ± 2.9	75.6 ± 3.0	57.0 ± 0.1	57.1 ± 0.2	56.9 ± 0.1	56.9 ± 0.1
CIFAR-10	GCE	92.2 ± 0.1	92.4 ± 0.3	89.2 ± 0.5	$89.9\pm0.2^{*}$	85.3 ± 0.7	$86.1\pm0.8^{***}$	78.0 ± 1.6	$80.0 \pm 1.0^{**}$
	SL	92.4 ± 0.1	92.5 ± 0.2	90.1 ± 0.1	$90.7\pm0.3^{*}$	87.6 ± 0.6	$87.9\pm0.3^{*}$	81.1 ± 0.9	$81.7\pm0.8^{***}$
	ELR*	93.4 ± 0.2	93.5 ± 0.2	92.7 ± 0.2	92.7 ± 0.2	91.9 ± 0.2	91.9 ± 0.3	90.4 ± 0.2	90.4 ± 0.3
	SOP+*	94.5 ± 0.1	94.6 ± 0.2	94.0 ± 0.1	94.1 ± 0.2	93.0 ± 0.4	93.0 ± 0.4	91.3 ± 0.6	91.5 ± 0.2
	CE	$72.0\pm0.5^{**}$	69.3 ± 0.9	55.8 ± 0.8	$65.5 \pm 1.5^{***}$	49.7 ± 1.5	$58.6 \pm 1.7^{ au au}$	41.6 ± 1.5	$49.0 \pm 1.6^{***}$
	Part	—	$69.2\pm0.9^{***}$	_	$65.2\pm1.1^{***}$	—	$58.8\pm0.5^{***}$	_	$48.7\pm1.3^{***}$
	MAE	21.7 ± 2.4	21.7 ± 2.4	17.5 ± 1.7	17.6 ± 1.8	16.2 ± 1.3	15.9 ± 1.3	14.2 ± 2.1	14.2 ± 2.1
CIFAR-100	GCE	69.3 ± 0.5	70.0 ± 0.2	62.0 ± 0.8	$63.9 \pm 1.2^{**}$	53.1 ± 0.9	$55.2 \pm 1.5^{**}$	41.7 ± 0.5	$43.6 \pm 1.0^{*}$
	SL	59.1 ± 1.2	$66.0 \pm 1.2^{***}$	54.0 ± 3.2	60.1 ± 2.7	50.0 ± 1.1	$55.4 \pm 1.5^{\circ}$	41.2 ± 1.5	$46.1 \pm 1.0^{***}$
	ELR*	74.9 ± 0.4	74.9 ± 0.4	71.8 ± 0.1	71.8 ± 0.1	66.2 ± 0.5	$68.4 \pm 1.6^{*}$	57.2 ± 0.6	$61.3\pm1.8^{**}$
	SOP+*	74.3 ± 0.4	74.5 ± 0.5	67.4 ± 0.4	$70.9\pm0.5^{**}$	59.6 ± 0.5	$63.7 \pm 1.7^{**}$	50.1 ± 0.6	$52.9 \pm 1.9^{*}$

D.7 Validation loss results in Sec. 4.4

Table S.6 presents the validation loss results from the corresponding models used in Table 1 in the main text. Similar to the test accuracy results in Table 1, Table S.6 shows that our resetting method consistently achieves either at least equivalent or lower validation losses compared to the baseline approach (i.e., no resetting).

Table S.6: Validation losses with different methods. We compare the performance without resetting (No) and with resetting (Reset) at r = 0.001. Results are presented as the average and the standard deviation. The best results are indicated in **bold** with statistical significance.

Dotocet	Mathod	Noise	e rate 0.2	Noise	rate 0.4	Noise rate 0.6		
Dataset	Method	No	Reset	No	Reset	No	Reset	
	CE	1.231 ± 0.018	$1.164 \pm 0.015^{***}$	1.774 ± 0.010	$1.734 \pm 0.018^{**}$	2.124 ± 0.009	$2.102 \pm 0.012^{**}$	
	PartRestart	—	$1.184 \pm 0.012^{**}$	—	$1.748 \pm 0.014^{*}$	—	2.114 ± 0.011	
	MAE	0.545 ± 0.009	0.545 ± 0.009	0.968 ± 0.039	0.963 ± 0.039	1.422 ± 0.042	1.420 ± 0.042	
CIFAR-10	GCE	0.388 ± 0.008	0.383 ± 0.009	0.686 ± 0.007	$0.673 \pm 0.007^{*}$	0.957 ± 0.009	0.945 ± 0.011	
	SL	2.720 ± 0.047	2.707 ± 0.041	4.759 ± 0.045	4.702 ± 0.055	6.617 ± 0.059	6.558 ± 0.069	
	ELR*	1.572 ± 0.081	1.572 ± 0.081	1.718 ± 0.021	1.723 ± 0.022	2.012 ± 0.006	2.016 ± 0.006	
	SOP+*	2.067 ± 0.035	2.029 ± 0.061	1.611 ± 0.007	1.632 ± 0.021	2.004 ± 0.01	2.011 ± 0.012	
	CE	2.759 ± 0.053	$2.583 \pm 0.051^{***}$	3.675 ± 0.064	$3.540 \pm 0.074^{*}$	4.283 ± 0.023	$4.201 \pm 0.026^{***}$	
	PartRestart	—	$2.564 \pm 0.047^{***}$	—	$3.549 \pm 0.067^*$	—	$4.206 \pm 0.019^{***}$	
	MAE	1.690 ± 0.039	1.690 ± 0.039	1.864 ± 0.044	1.864 ± 0.045	1.932 ± 0.010	1.931 ± 0.010	
CIFAR-100	GCE	0.644 ± 0.007	0.636 ± 0.008	0.896 ± 0.007	$0.883 \pm 0.006^{*}$	1.129 ± 0.004	$1.116 \pm 0.006^{**}$	
	SL	4.837 ± 0.066	$4.727 \pm 0.066^{*}$	6.578 ± 0.036	$6.427 \pm 0.059^{**}$	8.108 ± 0.043	$7.996 \pm 0.038^{**}$	
	ELR*	3.393 ± 0.072	$2.678 \pm 0.099^{***}$	3.985 ± 0.477	3.556 ± 0.073	4.242 ± 0.039	4.253 ± 0.068	
	SOP+*	2.600 ± 0.061	2.605 ± 0.047	3.616 ± 0.338	3.43 ± 0.082	4.185 ± 0.074	4.180 ± 0.088	

D.8 Training time comparison & scalability discussion for large-scale models

In our experiments, we used DNN models up to ResNet-34, which has approximately 21M parameters and requires about 85MB in FP32 precision. For these models, the memory required to store the checkpoint did not significantly burden GPU memory, allowing us to maintain both the current model and the checkpoint in GPU memory during training. Consequently, this approach did not result in any noticeable increase in training time. To validate this, we measure the training time ratios between runs with and without resetting in our experiments on the small dataset ciFAIR-10 (Setting 1 in Sec. 4.2) and the real-world noisy datasets CIFAR-10N/100N (Setting 3 in Sec. 4.5). As shown in Fig. S.8, the training time ratios consistently fall within the error bars near 1 regardless of the reset probability r. Note that the training time results for CIFAR-10N/100N in Setting 3 are almost the same as the results for CIFAR-10/100 in Setting 2. These results demonstrate that our method does not involve significant time overhead when GPU memory is sufficient to maintain two DNN models simultaneously during training.

However, for much larger models, it may not be feasible to keep both the current model and the checkpoint in GPU memory. In such cases, the checkpoint would need to be stored on a storage device and reloaded at each reset operation, potentially incurring additional I/O overhead. To estimate this cost for large-scale models, let us consider GPT-3 [103], which has 175B parameters (350GB in FP16), and single storage PCIe Gen4 x4 NVMe SSD with a transfer speed of 7GB/s. Then, loading or saving the model would take approximately 50 seconds per operation. GPT-3 was trained on 300B tokens with a global batch size of 3.2M tokens per iteration, resulting in approximately 94,000 iterations. Assuming 1,024 V100 GPUs [104], the total training time for GPT-3 is calculated as 3.55 days. Applying our method with $r = 10^{-3}$ to GPT-3 would only add 47,000 seconds (0.05 days) to the total training time, which is negligible in comparison to the overall training duration. Furthermore, while this calculation assumes the use of a single storage device, distributed storage systems are commonly employed for large-scale model training. These systems can significantly reduce loading times, minimizing the I/O overhead associated with resetting.

These findings demonstrate both the scalability and applicability of our method to large-scale models without significant time overhead. The performance improvements achieved on ViTs in Table 2 further support our claims. Additionally, the partial resetting method described in Sec. 4.3 offers an effective solution to mitigate I/O costs by resetting only portions of the DNN rather than the entire model, making it even more practical for large-scale implementations.



Figure S.8: Training time ratio of runs with resetting compared to runs with no reset (a) on ciFAIR-10 with varying the noise rate τ and the reset probability r (Sec. 4.2), and (b) on real-world noisy datasets, CIFAR-10N/100N (Sec. 4.5). Error bars indicate the standard deviation.