

Towards Adversarially Robust Dataset Distillation by Curvature Regularization

Eric Xue¹, Yijiang Li², Haoyang Liu³, Peiran Wang³, Yifan Shen⁴, Haohan Wang³

¹Department of Computer Science, University of Toronto

²Electrical and Computer Engineering, University of California San Diego

³School of Information Sciences, University of Illinois Urbana-Champaign

⁴Siebel School of Computing and Data Science, University of Illinois Urbana-Champaign

e.xue@mail.utoronto.ca, yijiangli@ucsd.edu, whilebug@gmail.com, {hl57, yifan26, haohanw}@illinois.edu,

Abstract

Dataset distillation (DD) allows datasets to be distilled to fractions of their original size while preserving the rich distributional information so that models trained on the distilled datasets can achieve a comparable accuracy while saving significant computational loads. Recent research in this area has been focusing on improving the accuracy of models trained on distilled datasets. In this paper, we aim to explore a new perspective of DD. We study how to embed adversarial robustness in distilled datasets, so that models trained on these datasets maintain the high accuracy and meanwhile acquire better adversarial robustness. We propose a new method that achieves this goal by incorporating curvature regularization into the distillation process with much less computational overhead than standard adversarial training. Extensive empirical experiments suggest that our method not only outperforms standard adversarial training on both accuracy and robustness with less computation overhead but is also capable of generating robust distilled datasets that can withstand various adversarial attacks.

Introduction

In the era of big data, the computational demands for training deep learning models are continuously growing due to the increasing volume of data. This presents substantial challenges, particularly for entities with limited computational resources. To mitigate such issues, concepts like dataset distillation (Wang et al. 2018) and dataset condensation (Zhao, Mopuri, and Bilen 2021; Zhao and Bilen 2021, 2023) have emerged, offering a means to reduce the size of the data while maintaining its utility. A successful implementation of dataset distillation can bring many benefits, such as enabling more cost-effective research on large datasets and models.

Dataset distillation (DD) refers to the task of synthesizing a smaller dataset such that models trained on this smaller set yield high performance when tested against the original, larger dataset. The dataset distillation algorithm takes a large dataset as input and generates a compact, synthetic dataset. The efficacy of the distilled dataset is assessed by evaluating models trained on it against the original dataset.

Conventionally, distilled datasets are evaluated based on their standard test accuracy. Therefore, recent research has

expanded rapidly in the direction of improving the test accuracy following the evaluation procedure (Sachdeva and McAuley 2023). Additionally, many studies focus on improving the efficiency of the distillation process (Sachdeva and McAuley 2023).

Less attention, however, has been given to an equally important aspect of this area of research: the adversarial robustness of models trained on distilled datasets. Adversarial robustness is a key indicator of a model’s resilience against malicious inputs, making it a crucial aspect of trustworthy machine learning. Given the potential of dataset distillation to safeguard the privacy of the original dataset (Geng et al. 2023; Chen et al. 2023), exploring its capability to also enhance model robustness opens a promising avenue for advancing research in trustworthy machine learning (Liu, Chaudhary, and Wang 2023). Thus, our work seeks to bridge this gap and focuses on the following question: **How can we embed adversarial robustness into the dataset distillation process, thereby generating datasets that lead to more robust models?**

Motivated by this question, we explore potential methods to accomplish this goal. As it turns out, it is not as simple as adding adversarial training to the distillation process. To find a more consistent method, we study the theoretical connection between adversarial robustness and dataset distillation. Our theory suggests that we can directly improve the robustness of the distilled dataset by minimizing the curvature of the loss function with respect to the real data. Based on our findings, we propose a novel method, GUARD (Geometric regularization for Adversarial Robust Dataset), which incorporates curvature regularization into the distillation process. We then evaluate GUARD against existing distillation methods on ImageNet, Tiny ImageNet, and ImageNet datasets. In summary, the contributions of this paper are as follows

- Empirical and theoretical exploration of adversarial robustness in distilled datasets
- A theory-motivated method, GUARD, that offers robust dataset distillation with minimal computational overhead
- Detailed evaluation of GUARD to demonstrate its effectiveness across multiple aspects

Related Works

Dataset Distillation

Aiming to address the issue of the increasing amount of data required to train deep learning models, the goal of dataset distillation is to efficiently train neural networks using a small set of synthesized training examples from a larger dataset. Dataset distillation (DD) (Wang et al. 2018) was one of the first such methods developed, and it showed that training on a few synthetic images can achieve similar performance on MNIST and CIFAR10 as training on the original dataset. Later, Cazenavette et al. (2022); Zhao and Bilen (2021); Zhao, Mopuri, and Bilen (2021); Lee et al. (2022) explored different methods of distillation, including gradient and trajectory matching w.r.t. the real and synthetic data, with stronger supervision for the training process. Instead of matching the weights of the neural network, another thread of works (Wang et al. 2022; Zhao and Bilen 2023; Zhang et al. 2024; Liu et al. 2023) focuses on matching feature distributions of the real and synthetic data in the embedding space to better align features or preserve real-feature distribution. Considering the lack of efficiency of the bi-level optimization in previous methods, Nguyen et al. (2021); Zhou, Nezhadarya, and Ba (2022) aim to address the significant amount of meta gradient computation challenges. Nguyen, Chen, and Lee (2020) proposed a kernel-inducing points meta-learning algorithm and they further leverage the connection between the infinitely wide ConvNet and kernel ridge regression for better performance. Furthermore, Sucholutsky and Schonlau (2021) addresses the simultaneous distillation of images and their corresponding soft labels. Later, some works focused on further improving efficiency of the process, such as Yin, Xing, and Shen (2023) that introduced SRe^2L , which optimizes the distillation process by dividing it into three distinct steps for greater efficiency, and Xu et al. (2024), which proposed an approach to enhance both the efficiency and performance by first pruning the original dataset. Finally, Li et al. (2024) further advanced the process by dynamically pruning the original dataset based on the desired compression ratio and extracting information from deeper layers of the network.

Dataset distillation approaches can be broadly classified into four families based on their underlying principles: meta-model matching, gradient matching, distribution matching, and trajectory matching (Sachdeva and McAuley 2023). Regardless of the particular approach, most of the existing methods rely on optimizing the distilled dataset w.r.t. a network trained with real data, such methods include DD (Wang et al. 2018), DC (Zhao, Mopuri, and Bilen 2021), DSA (Zhao and Bilen 2021), MTT (Cazenavette et al. 2022), DCC (Lee et al. 2022), SRe^2L (Yin, Xing, and Shen 2023), ATT (Liu et al. 2024) and many more.

In a related direction, some works also address the robustness of dataset distillation, but specifically focusing on out-of-distribution (OOD) robustness. For instance, Vahidian et al. (2024) employs risk minimization techniques to ensure robustness, while TrustDD (Ma et al. 2024) incorporates outliers during the distillation process to facilitate OOD detection.

Adversarial Attacks

Adversarial attacks are a significant concern in the field of machine learning, as they can cause models to make incorrect predictions even when presented with seemingly similar input. Kurakin, Goodfellow, and Bengio (2017) demonstrates the real-world implications of these attacks. Many different types of adversarial attacks have been proposed in the literature (Goodfellow, Shlens, and Szegedy 2015; Madry et al. 2018). In particular, Projected Gradient Descent (PGD) is a widely used adversarial attack that has been shown to be highly effective against a variety of machine learning models (Madry et al. 2018). The limitations of defensive distillation, a technique initially proposed for increasing the robustness of machine learning models, were explored by Papernot et al. (2017). Moosavi-Dezfooli, Fawzi, and Frossard (2016) introduced DeepFool, an efficient method to compute adversarial perturbations. Other notable works include the study of the transferability of adversarial attacks by Papernot, McDaniel, and Goodfellow (2016), the simple and effective black-box attack by Narodytska and Kasiviswanathan (2016), and the zeroth-order optimization-based attack by Chen et al. (2017). More recently, Athalye, Carlini, and Wagner (2018) investigated the robustness of obfuscated gradients, and Wong, Schmidt, and Kolter (2019) introduced the Wasserstein smoothing as a novel defense against adversarial attacks. Croce and Hein (2020) introduced AutoAttack, which is a suite of adversarial attacks consisting of four diverse and parameter-free attacks that are designed to provide a comprehensive evaluation of a model's robustness to adversarial attacks.

Adversarial Defense

Numerous defenses against adversarial attacks have been proposed. Among these, adversarial training stands out as a widely adopted defense mechanism that entails training machine learning models on both clean and adversarial examples (Goodfellow, Shlens, and Szegedy 2015). Several derivatives of the adversarial training approach have been proposed, such as ensemble adversarial training (Tramer et al. 2018), and randomized smoothing (Cohen, Rosenfeld, and Kolter 2019) — a method that incorporates random noise to obstruct the generation of effective adversarial examples. However, while adversarial training can be effective, it bears the drawback of being computationally expensive and time-consuming.

Some defense mechanisms adopt a geometrical approach to robustness. One such defense mechanism is CURE, a method that seeks to improve model robustness by modifying the curvature loss function used during training (Moosavi-Dezfooli et al. 2019). The primary aim of CURE is to minimize the sensitivity of the model to adversarial perturbations in the input space to make it more difficult for an attacker to find adversarial examples which cross this boundary. Miyato et al. (2015) focuses on improving the smoothness of the output distribution to make models more resistant to adversarial attacks, while Cisse et al. (2017b) introduced Parseval networks, which enforce Lipschitz constant to improve model robustness. Ross and Doshi-Velez

(2018) also presented a method for improving the robustness of deep learning models using input gradient regularization.

Several other types of defense techniques have also been proposed, such as corrupting with additional noise and pre-processing with denoising autoencoders by Gu and Rigazio (2014), the defensive distillation approach by Papernot et al. (2016), and the Houdini adversarial examples by Cisse et al. (2017a).

Preliminary

Dataset Distillation

Before we delve into the theory of robustness in dataset distillation methods, we will formally introduce the formulation of dataset distillation in this section.

Notations Let \mathcal{T} represent the real dataset, drawn from the distribution $\mathcal{D}_{\mathcal{T}}$. The dataset \mathcal{T} comprises n image-label pairs, defined as $\mathcal{T} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$. Similarly, let \mathcal{S} denote the distilled dataset, drawn from the distribution $\mathcal{D}_{\mathcal{S}}$, and consisting of m image-label pairs, defined as $\mathcal{S} = \{(\tilde{\mathbf{x}}_j, \tilde{y}_j)\}_{j=1}^m$, where $m \ll n$. Conventionally, instead of directly expressing the size of the distilled dataset as $|\mathcal{S}|$, it is more common to describe it in terms of "images per class" (IPC). We denote the loss function of a model parameterized by θ on a sample (\mathbf{x}, y) as $\ell(\mathbf{x}, y; \theta)$. The empirical loss on \mathcal{T} is represented as $\mathcal{L}(\mathcal{T}; \theta)$, defined as $\mathcal{L}(\mathcal{T}; \theta) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, y_i; \theta)$.

Given the real training set \mathcal{T} , dataset distillation aims to find the optimal synthetic dataset \mathcal{S}^* by solving the following bi-level optimization problem:

$$\begin{aligned} \mathcal{S}^* = \arg \min_{\mathcal{S}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_{\mathcal{T}}} \ell(\mathbf{x}, y; \theta(\mathcal{S})) \\ \text{subject to } \theta(\mathcal{S}) = \arg \min_{\theta} \mathcal{L}(\mathcal{S}; \theta). \end{aligned} \quad (1)$$

Directly solving this problem requires searching for the optimal parameters in the inner problem and unrolling the gradient descent steps in the computation graph to find the hypergradient with respect to \mathcal{S} , which is computationally expensive. One common alternative approach is to align a model trained on the distilled set with one trained on the real dataset. Conceptually, it can be summarized in the below equation:

$$\begin{aligned} \min_{\mathcal{S}} D(\theta(\mathcal{S}), \theta(\mathcal{T})) \\ \text{subject to } \theta(\mathcal{S}) = \arg \min_{\theta} \mathcal{L}(\mathcal{S}; \theta) \\ \text{and } \theta(\mathcal{T}) = \arg \min_{\theta} \mathcal{L}(\mathcal{T}; \theta) \end{aligned} \quad (2)$$

where D is a manually chosen distance function. Recent works have proliferated along this direction, with methods such as gradient matching (Zhao, Mopuri, and Bilen 2021) and trajectory matching (Cazenavette et al. 2022), each focusing on aligning different aspects of the model's optimization dynamics. Some works have also tried to align the distribution of the distilled data with that of the real data (Zhao and Bilen 2023; Zhang et al. 2024; Liu et al. 2023), or recover a distilled version of the training data from a trained model (Yin, Xing, and Shen 2023; Buzaglo et al. 2023).

These methods do not rely on the computation of second-order gradients, leading to improved efficiency and performance on large-scale datasets.

Despite the wide spectrum of methods for dataset distillation, they were primarily designed for improving the standard test accuracy, and significantly less attention has been paid to the adversarial robustness. In the following, we conduct a preliminary study to show that adversarial robustness cannot be easily incorporated into the distilled data by the common approach of adversarial training, necessitating more refined analysis.

The Limitation of Adversarial Training in Dataset Distillation

IPC	Attack	GUARD	SRe ² L	SRe ² L +Adv
1	None (Clean)	37.49	27.97	11.61
	PGD100	16.22	12.05	10.03
	Square	26.74	18.62	11.18
	AutoAttack	15.81	12.12	10.03
	CW	29.14	20.38	10.31
	MIM	16.32	12.05	10.03
10	None (Clean)	57.93	42.42	12.81
	PGD100	23.87	4.76	9.93
	Square	44.07	22.77	11.46
	AutoAttack	19.69	4.99	9.96
	CW	58.67	22.11	10.90
	MIM	21.80	4.76	9.96

Table 1: Accuracy of ResNet18 on ImageNette trained on distilled datasets from GUARD, SRe²L, and SRe²L with adversarial training

In the supervised learning setting, one of the most commonly used methods to enhance model robustness is adversarial training, which involves training the model on adversarial examples that are algorithmically searched for or crafted (Goodfellow, Shlens, and Szegedy 2015). This can be formulated as

$$\min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left(\max_{\|\mathbf{v}\| \leq \rho} \ell(\mathbf{x} + \mathbf{v}, y; \theta) \right), \quad (3)$$

where \mathbf{v} is some perturbation within the ℓ_p ball with radius ρ , and \mathcal{D} is the data distribution.

Analogously, in the dataset distillation setting, one intuitive way to distill robust datasets would be to synthesize a distilled dataset using a robust model trained with adversarial training. As mentioned in the related works section, many dataset distillation methods utilize a model trained on the original dataset as a comparison target, therefore this technique can be easily integrated to those methods.

While embedding adversarial training directly within the dataset distillation process may seem like an intuitive and straightforward approach, our comprehensive analysis reveals its limitations across various distillation methods. As an example, we show the evaluation of one such implementation using SRe²L, an efficient dataset distillation method in Table 1. The results indicate a significant decline in clean accuracy for models trained on datasets distilled using this

technique, in contrast to those synthesized by the original method. Moreover, the improvements in robustness achieved are very inconsistent. In our experiment, we only employed a weak PGD attack with $\epsilon = 1/255$ to generate adversarial examples for adversarial training, leading to the conclusion that even minimal adversarial training can detrimentally impact model performance when integrated into the dataset distillation process.

Such outcomes are not entirely unexpected. Previous studies, such as those by Zhang et al. (2020), have indicated that adversarial training can significantly alter the semantics of images through perturbations, even when adhering to set norm constraints. This can lead to the cross-over mixture problem, severely degrading the clean accuracy. We hypothesize that these adverse effects might be magnified during the distillation process, where the distilled dataset’s constrained size results in a distribution that is vastly different from that of the original dataset.

Methods

Formulation of the Robust Distillation Problem

Extending the distillation problem to the adversarial robustness setting, robust dataset distillation can be formulated as a tri-level optimization problem as below:

$$\begin{aligned} \mathcal{S}^* = \arg \min_{\mathcal{S}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_{\mathcal{T}}} \left(\max_{\|\mathbf{v}\| \leq \rho} \ell(\mathbf{x} + \mathbf{v}, y; \boldsymbol{\theta}(\mathcal{S})) \right) \\ \text{subject to } \boldsymbol{\theta}(\mathcal{S}) = \arg \min_{\boldsymbol{\theta}} \mathcal{L}(\mathcal{S}; \boldsymbol{\theta}). \end{aligned} \quad (4)$$

If we choose to directly optimize for the robust dataset distillation objective, the tri-level optimization problem will result in a hugely inefficient process. Instead, we will uncover a theoretical relationship between dataset distillation and adversarial robustness to come up with a more efficient method that avoids the tri-level optimization process.

Theoretical Bound of Robustness

Our aim is to create a method that allows us to efficiently and reliably introduce robustness into distilled datasets, thus we will start by exploring the theoretical connections between dataset distillation and adversarial robustness. Conveniently, previous works (Jetley, Lord, and Torr 2018; Fawzi et al. 2018) have studied the adversarial robustness of neural networks via the geometry of the loss landscape. Inspired by Moosavi-Dezfooli, Fawzi, and Frossard (2016), here we find connections between standard training procedures and dataset distillation to provide a theoretical bound for the adversarial loss of models trained with distilled datasets.

Let $\ell(\mathbf{x}, y; \boldsymbol{\theta})$ denote the loss function of the neural network, or $\ell(\mathbf{x})$ for simplicity, and \mathbf{v} denote a perturbation vector. By Taylor’s Theorem,

$$\ell(\mathbf{x} + \mathbf{v}) = \ell(\mathbf{x}) + \nabla \ell(\mathbf{x})^\top \mathbf{v} + \frac{1}{2} \mathbf{v}^\top \mathbf{H} \mathbf{v} + o(\|\mathbf{v}\|^2). \quad (5)$$

We are interested in the property of $\ell(\cdot)$ in the locality of \mathbf{x} , so we focus on the quadratic approximation $\tilde{\ell}(\mathbf{x} + \mathbf{v}) = \ell(\mathbf{x}) + \nabla \ell(\mathbf{x})^\top \mathbf{v} + \frac{1}{2} \mathbf{v}^\top \mathbf{H} \mathbf{v}$. We define the adversarial loss

on real data as $\tilde{\ell}_\rho^{adv}(\mathbf{x}) = \max_{\|\mathbf{v}\| \leq \rho} \tilde{\ell}(\mathbf{x} + \mathbf{v})$. We can expand this and take the expectation over the distribution with class label c , denoted as D_c , to get the following:

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim D_c} \tilde{\ell}_\rho^{adv}(\mathbf{x}) \leq \mathbb{E}_{\mathbf{x} \sim D_c} \ell(\mathbf{x}) + \rho \mathbb{E}_{\mathbf{x} \sim D_c} \|\nabla \ell(\mathbf{x})\| + \\ \frac{1}{2} \rho^2 \mathbb{E}_{\mathbf{x} \sim D_c} \lambda_1(\mathbf{x}), \end{aligned} \quad (6)$$

where λ_1 is the largest eigenvalue of the Hessian matrix $\mathbf{H}(\ell(\mathbf{x}))$. Then, we have the proposition:

Proposition 1. *Let \mathbf{x}' be a distilled datum with the label c and satisfies $\|h(\mathbf{x}') - \mathbb{E}_{\mathbf{x} \sim D_c}[h(\mathbf{x})]\| \leq \sigma$, where $h(\cdot)$ is a feature extractor. Assume $\ell(\cdot)$ is convex in \mathbf{x} and $\tilde{\ell}_\rho^{adv}(\cdot)$ is L -Lipschitz in the feature space, then the below inequality holds:*

$$\begin{aligned} \tilde{\ell}_\rho^{adv}(\mathbf{x}') \leq \mathbb{E}_{\mathbf{x} \sim D_c} \ell(\mathbf{x}) + \rho \mathbb{E}_{\mathbf{x} \sim D_c} \|\nabla \ell(\mathbf{x})\| + \\ \frac{1}{2} \rho^2 \mathbb{E}_{\mathbf{x} \sim D_c} \lambda_1(\mathbf{x}) + L\sigma. \end{aligned} \quad (7)$$

Given the assumption of convexity in the loss function, we can further observe that in a convex landscape the gradient magnitude tends to be lower, particularly near the optimal points. Therefore, in the context of a convex loss function and a well-distilled dataset, the gradient term $\rho \mathbb{E}_{\mathbf{x} \sim D_c} \|\nabla \ell(\mathbf{x})\|$ contribute insignificantly to the overall value of the inequality. This insignificance is amplified by the presence of the curvature term, $\frac{1}{2} \rho^2 \mathbb{E}_{\mathbf{x} \sim D_c} \lambda_1(\mathbf{x})$, which provides a sufficient descriptor of the loss landscape under our assumptions. Hence, it is reasonable to simplify the expression by omitting the gradient term, resulting in a focus on the curvature term, which is more representative of the convexity assumption and the characteristics of a well-distilled dataset. The revised expression would then be:

$$\tilde{\ell}_\rho^{adv}(\mathbf{x}') \leq \mathbb{E}_{\mathbf{x} \sim D_c} \ell(\mathbf{x}) + \frac{1}{2} \rho^2 \mathbb{E}_{\mathbf{x} \sim D_c} \lambda_1(\mathbf{x}) + L\sigma. \quad (8)$$

Dataset distillation methods usually already optimizes for $\ell(\mathbf{x})$, and we can also assume that the σ for a well-distilled dataset is small. Hence, we can conclude that the upper bound of adversarial loss of distilled datasets is largely affected by the curvature of the loss function in the locality of real data samples.

Geometric Regularization for Adversarial Robust Dataset

Based on our theoretical discussion, we propose a method, GUARD (Geometric regUlarization for Adversarial Robust Dataset). Since the theorem suggests that the upper bound of the adversarial loss is mainly determined by the curvature of the loss function, we modify the distillation process so that the trained model has a loss function with a low curvature with respect to real data.

To reduce λ_1 in Eq. 8 requires computing the Hessian matrix and get the largest eigenvalue λ_1 , which is quite computationally expensive. Here we find an efficient approximation of it. Let \mathbf{v}_1 be the unit eigenvector corresponding to

λ_1 , then the Hessian-vector product is

$$\mathbf{H}\mathbf{v}_1 = \lambda_1\mathbf{v}_1 = \lim_{h \rightarrow 0} \frac{\nabla\ell(\mathbf{x} + h\mathbf{v}_1) - \nabla\ell(\mathbf{x})}{h}. \quad (9)$$

We take the differential approximation of the Hessian-vector product, because we are interested in the curvature in a local area of x rather than its asymptotic property. Therefore, for a small h ,

$$\lambda_1 = \|\lambda_1\mathbf{v}_1\| \approx \left\| \frac{\nabla\ell(\mathbf{x} + h\mathbf{v}_1) - \nabla\ell(\mathbf{x})}{h} \right\|. \quad (10)$$

Previous works (Fawzi et al. 2018; Jetley, Lord, and Torr 2018; Moosavi-Dezfooli et al. 2019) have empirically shown that the direction of the gradient has a large cosine similarity with the direction of \mathbf{v}_1 in the input space of neural networks. Instead of calculating \mathbf{v}_1 directly, it is more efficient to take the gradient direction as a surrogate of \mathbf{v}_1 to perturb the input \mathbf{x} . So we replace the \mathbf{v}_1 above with the normalized gradient $\mathbf{z} = \frac{\nabla\ell(\mathbf{x})}{\|\nabla\ell(\mathbf{x})\|}$, and define the regularized loss ℓ_R to encourage linearity in the input space:

$$\ell_R(\mathbf{x}) = \ell(\mathbf{x}) + \lambda \|\nabla\ell(\mathbf{x} + h\mathbf{z}) - \nabla\ell(\mathbf{x})\|^2, \quad (11)$$

where ℓ is the original loss function, h is the discretization step, and the denominator h is merged with the regularization coefficient λ .

Engineering Specification

In order to evaluate the effectiveness of our method, we implemented GUARD using the SRe²L method as a baseline. We incorporated our regularizer into the squeeze step of SRe²L by substituting the standard training loss with the modified loss outlined in Eq. 11. In the case of SRe²L, this helps to synthesize a robust distilled dataset by allowing images to be recovered from a robust model in the subsequent recover step.

Experiments

Experiment Settings

For a systematic evaluation of our method, we investigate the top-1 classification accuracy of models trained on data distilled from three commonly-used datasets in this domain: ImageNette (Howard 2018), Tiny ImageNet (Le and Yang 2015), and ImageNet (Deng et al. 2009). ImageNette is a subset of ImageNet containing 10 easy-to-classify classes. Tiny ImageNet is a scaled-down subset of ImageNet, containing 200 classes and 100,000 downsized 64x64 images. We train networks using the distilled datasets and subsequently verify the network’s performance on the original datasets. For consistency in our experiments across all datasets, we use the standard ResNet18 architecture (He et al. 2016) to synthesize the distilled datasets and evaluate their performance.

During the squeeze step of the distillation process, we train the model on the original dataset over 50 epochs using a learning rate of 0.025. Based on preliminary experiments, we determined that the settings $h = 3$ and $\lambda = 100$ provide an optimal configuration for our regularizer. In the recover step, we perform 2000 iterations to synthesize the images and run 300 epochs to generate the soft labels to obtain

the full distilled dataset. In the evaluation phase, we train a ResNet18 model on the distilled dataset for 300 epochs, before assessing it on the test set of the original dataset.

Comparison with Other Methods

As of now, there is only a small number of dataset distillation methods that can achieve good performance on ImageNet-level datasets, therefore our choices for comparison is small. Here, we first compare our method to the original SRe²L (Yin, Xing, and Shen 2023) to observe the direct effect of our regularizer on the adversarial robustness of the trained model. We also compare with MTT (Cazenavette et al. 2022) and TESLA (Cui et al. 2023) on the same datasets to gain a further understanding on the differences in robustness between our method and other dataset distillation methods. We utilized the exact ConvNet architecture described in the papers of MTT and TESLA for their distillation and evaluation, as their performance on ResNet seems to be significantly lower.

We evaluate all the methods on three distillation scales: 10 IPC, 50 IPC, and 100 IPC. We also employed a range of attacks to evaluate the robustness of the model, including PGD100 (Madry et al. 2017), Square (Andriushchenko et al. 2020), AutoAttack (Croce and Hein 2020), CW (Carlini and Wagner 2017), and MIM (Dong et al. 2017). This assortment includes both white-box and black-box attacks, providing a comprehensive evaluation of GUARD. For all adversarial attacks, with the exception of CW attack, we use the setting $\epsilon = 1/255$. For CW specifically, we set the box constraint c to $1e^{-5}$. Due to computational limits, we were not able to obtain results for MTT and TESLA with the 100 IPC setting on ImageNet, as well as the 100 IPC setting on ImageNet for all methods.

Results

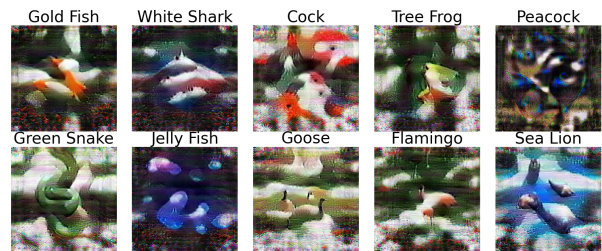


Figure 1: Visualization of distilled images generated using GUARD with 1 IPC setting from ImageNet.

The results are detailed in Table 2. It can be observed that GUARD consistently outperforms both SRe²L and MTT in terms of robustness across various attacks. Interestingly, we observe an increase in clean accuracy upon incorporating GUARD across various settings. While enhancing clean accuracy was not the primary goal of GUARD, this outcome aligns with its function as a regularizer, potentially aiding in model generalization. In the context of dataset distillation, where the goal is to distill essential features of the original dataset into a smaller subset, improving the generalization

is expected to have positive effects on the performance. We also provide a visualization of the distilled images generated by GUARD in Figure 1, utilizing a distillation scale of 1 image per class among selected ImageNet classes. It can be seen that the images exhibit characteristics that resemble a blend of multiple objects within their assigned class, highlighting the method’s capacity to capture essential features.

Ablation Study on Gradient Regularization

Eq. 7 showed that the adversarial loss is upper-bounded by the normal loss, the gradient magnitude, and the curvature term. GUARD regularizes the curvature term while disregarding the gradient magnitude, which could theoretically reduce the upper bound of the loss as well. Here, we investigate the effect of regularizing gradient instead of curvature and present the results in Table 3. The results indicate that GUARD outperforms the gradient regularization alternatives, regardless of the regularization parameter.

Discussion

Robustness Guarantee

Due to the nature of dataset distillation, it is impossible to optimize the robustness of the final model with respect to the real dataset. Therefore, most approaches in this direction, including ours, have to optimize the adversarial loss of the model with respect to the distilled dataset. Unfortunately, there is always a distribution shift between the real and distilled datasets, which raises uncertainties about whether robustness on the distilled dataset will be effectively transferred when evaluated against the real dataset. Nevertheless, our theoretical framework offers assurances regarding this concern. A comparison between Eq. 6 with Eq. 7 reveals that the bounds of adversarial loss for real data and distilled data differ only by $L\sigma$. For a well-distilled dataset, σ should be relatively small. We have thus demonstrated that the disparity between minimizing adversarial loss on the distilled dataset and on the real dataset is confined to this constant. This conclusion of our theory allows future robust dataset distillation methods to exclusively enhance robustness with respect to the distilled dataset, without worrying if the robustness can transfer well to the real dataset.

Computational Overhead

The structure of robust dataset distillation, as outlined in Eq. 4, inherently presents a tri-level optimization challenge. Typically, addressing such a problem could entail employing complex tri-level optimization algorithms, resulting in significant computational demands. One example of this is the integration of adversarial training within the distillation framework, which necessitates an additional optimization loop for generating adversarial examples within each iteration. However, GUARD’s approach, as detailed in Eq. 11, introduces an efficient alternative. GUARD’s regularization loss only requires an extra forward pass to compute the loss $\ell(\mathbf{x} + h\mathbf{z})$ within each iteration. Therefore, integrating GUARD’s regularizer into an existing method does not significantly increase the overall computational complexity, ensuring that the computational overhead remains minimal.

Dataset	IPC	Attack	Methods			
			GUARD	SRe ² L	MTT	TESLA
ImageNette	10	None (Clean)	<u>57.93</u>	42.42	58.43	36.84
		PGD100	23.87	4.76	39.85	<u>28.10</u>
		Square	44.07	22.77	<u>34.79</u>	24.61
		AutoAttack	19.69	4.99	33.72	<u>24.48</u>
		CW	41.47	22.11	<u>34.57</u>	24.61
		MIM	21.80	4.76	39.20	<u>28.15</u>
	50	None (Clean)	80.86	<u>80.15</u>	59.69	36.21
		PGD100	41.42	12.30	<u>41.13</u>	28.72
		Square	72.81	<u>61.50</u>	<u>36.72</u>	25.34
		AutoAttack	42.47	12.91	<u>35.46</u>	27.21
		CW	58.67	<u>53.42</u>	36.54	29.01
		MIM	43.23	12.43	<u>41.69</u>	30.12
	100	None (Clean)	<u>83.39</u>	85.83	64.33	45.04
		PGD100	57.50	31.65	<u>44.89</u>	33.98
		Square	77.68	19.18	<u>40.41</u>	29.27
		AutoAttack	64.84	17.93	<u>39.46</u>	28.99
		CW	69.35	<u>68.20</u>	40.66	29.32
		MIM	65.07	18.98	<u>44.89</u>	33.98
TinyImageNet	10	None (Clean)	37.00	<u>33.18</u>	8.14	14.06
		PGD100	6.39	1.08	4.08	8.40
		Square	<u>19.53</u>	<u>15.85</u>	2.48	6.31
		AutoAttack	<u>4.91</u>	0.79	2.44	6.16
		CW	8.40	3.24	2.50	<u>6.26</u>
		MIM	<u>6.51</u>	1.10	4.08	8.40
	50	None (Clean)	<u>55.61</u>	56.42	17.84	28.24
		PGD100	15.63	0.27	5.62	<u>12.12</u>
		Square	36.93	<u>15.50</u>	3.84	10.39
		AutoAttack	13.84	0.16	3.52	<u>10.01</u>
		CW	20.46	<u>12.12</u>	3.66	10.13
		MIM	16.09	0.29	5.64	<u>12.12</u>
	100	None (Clean)	60.13	<u>59.30</u>	29.16	30.48
		PGD100	<u>13.79</u>	0.25	8.63	14.45
		Square	37.06	<u>17.74</u>	7.29	12.02
		AutoAttack	12.76	0.19	6.75	<u>11.57</u>
		CW	20.05	<u>14.02</u>	6.93	11.57
		MIM	<u>14.35</u>	0.24	8.63	14.45
ImageNet-1K	10	None (Clean)	27.25	21.30	-	-
		PGD100	5.25	<u>0.55</u>	-	-
		Square	<u>17.88</u>	18.02	-	-
		AutoAttack	3.33	<u>0.34</u>	-	-
		CW	7.68	<u>3.21</u>	-	-
		MIM	5.23	<u>0.51</u>	-	-
	50	None (Clean)	<u>39.89</u>	46.80	-	-
		PGD100	9.77	<u>0.59</u>	-	-
		Square	<u>28.39</u>	32.40	-	-
		AutoAttack	7.03	<u>0.47</u>	-	-
		CW	14.14	<u>6.31</u>	-	-
		MIM	9.84	<u>0.64</u>	-	-

Table 2: Evaluation of different dataset distillation methods under adversarial attacks on ImageNette, TinyImageNet, and ImageNet. The best results among all methods are highlighted in bold, second best are underlined.

This efficiency is particularly notable given the computa-

IPC	Attack	Methods						
		SRe ² L	GUARD	10 ⁻⁴	10 ⁻³	10 ⁻²	0.1	1
1	(Clean)	27.97	37.49	13.72	16.15	16.41	17.58	18.75
	PGD100	12.05	16.22	6.39	9.32	10.27	10.39	13.27
	Square	18.62	26.74	9.71	11.69	13.68	12.94	15.97
	AA	12.12	15.81	6.32	9.35	10.17	10.42	13.25
	CW	20.38	29.14	7.62	11.06	11.64	11.31	14.47
	MIM	12.05	16.32	6.39	9.25	10.39	10.39	13.27
10	(Clean)	42.42	57.93	44.82	41.81	42.80	43.34	40.31
	PGD100	4.76	23.87	15.39	14.47	15.41	13.68	16.92
	Square	22.77	44.07	34.06	31.80	34.73	31.85	31.03
	AA	4.99	19.69	15.62	14.52	15.64	13.78	16.87
	CW	22.11	41.47	21.58	20.64	21.17	18.85	21.53
	MIM	4.76	21.80	15.41	14.60	15.34	13.66	17.41
50	(Clean)	80.15	80.86	76.46	74.06	73.12	74.52	-
	PGD100	12.30	41.42	35.54	35.36	33.48	29.17	-
	Square	61.50	72.81	67.08	63.24	63.97	65.53	-
	AA	12.91	42.47	35.59	35.54	33.52	29.32	-
	CW	53.42	58.67	46.06	45.43	44.05	41.02	-
	MIM	12.43	43.23	35.61	35.39	33.66	29.07	-

Table 3: Accuracy on ImageNette of original SRe²L, GUARD, and gradient regularization on SRe²L with regularization parameters ($\lambda_g = 10^{-4}, 10^{-3}, 10^{-2}, 0.1, 1$, where λ_g is omitted in table columns for brevity). AA stands for AutoAttack. The best results among all methods are highlighted in bold.

tionally intensive nature of tri-level optimization in robust dataset distillation. In Table 4, we present a comparison of the time per iteration required for a tri-level optimization algorithm, such as the one used for embedded adversarial training, against the time required for GUARD. The findings show that GUARD is much more computationally efficient and has a lower memory usage as well.

Method	Time (s) / Iter	Peak Mem
GUARD	$0.007 \pm 2.661e^{-4}$	3851MiB
Adv Training	$2.198 \pm 6.735e^{-4}$	4211MiB

Table 4: Computation overhead of GUARD compared with embedded adversarial training. Experiments are performed on one NVIDIA A100 80GB PCIe GPU with batch size 32. We measure 5 times per iteration training time and report the average and standard deviation.

Transferability

Our investigation focuses on studying the effectiveness of the curvature regularizer within the SRe²L framework. Theoretically, this method can be extended to a broad spectrum of dataset distillation methods. GUARD’s application is feasible for any distillation approach that utilizes a model trained on the original dataset as a comparison target during the distillation phase — a strategy commonly seen across many dataset distillation techniques as noted in the related

works section. This criterion is met by the majority of dataset distillation paradigms, with the exception of those following the distribution matching approach, which may not consistently employ a comparison model (Sachdeva and McAuley 2023). This observation suggests GUARD’s potential compatibility with a wide array of dataset distillation strategies. To demonstrate this, we explored two additional implementations of GUARD using DC (Zhao, Mopuri, and Bilen 2021) and CDA (Yin and Shen 2023) as baseline distillation methods. DC represents an earlier, simpler approach that leverages gradient matching for distillation purposes, whereas CDA is a more recent distillation technique, specifically designed for very large datasets. As shown in Table 5, GUARD consistently improves both clean accuracy and robustness across various dataset distillation methods.

IPC	Attack	Methods					
		DC	<i>DC</i>	SRe ² L	<i>SRe²L</i>	CDA	<i>CDA</i>
1	None (Clean)	29.96	30.95	17.13	22.88	14.98	23.18
	PGD100	24.59	46.88	13.56	19.21	12.69	18.70
	Square	24.72	48.56	13.75	19.55	12.84	19.17
	AutoAttack	24.33	14.99	13.43	18.91	12.63	18.42
	CW	24.58	15.19	13.52	18.95	12.62	18.52
	MIM	24.62	15.27	13.57	19.22	12.69	18.71
10	None (Clean)	45.38	46.83	26.58	30.76	20.55	30.65
	PGD100	31.84	32.36	18.24	22.31	14.60	24.33
	Square	33.71	33.54	19.99	24.16	15.93	25.66
	AutoAttack	31.05	31.84	18.11	21.58	14.47	24.04
	CW	31.95	32.35	18.73	21.98	14.83	24.51
	MIM	31.89	32.37	18.25	22.35	14.62	24.34
50	None (Clean)	-	-	43.96	44.05	36.32	43.05
	PGD100	-	-	24.74	33.12	21.58	33.02
	Square	-	-	29.68	35.22	25.76	35.19
	AutoAttack	-	-	24.45	32.24	21.46	31.96
	CW	-	-	26.09	32.67	22.54	32.56
	MIM	-	-	24.81	33.12	21.61	33.03

Table 5: Direct comparison of the original DC, SRe²L, and CDA methods with the addition of GUARD regularizer (denoted by italic font) on CIFAR10. The best results among each pair of compared methods are highlighted in bold.

Conclusions

Our work focuses on a novel perspective on dataset distillation by emphasizing its adversarial robustness characteristics. Upon reaching the theoretical conclusion that the adversarial loss of distilled datasets is bounded by the curvature, we proposed GUARD, a method that can be integrated into many dataset distillation methods to provide robustness against diverse types of attacks and potentially improve clean accuracy. Our theory also provided the insight that the optimization of robustness with respect to distilled and real datasets is differentiated only by a constant term, which may open up various potentials for subsequent research in the field.

References

- Andriushchenko, M.; Croce, F.; Flammarion, N.; and Hein, M. 2020. Square attack: a query-efficient black-box adversarial attack via random search. In *Proceedings of the European Conference on Computer Vision*.
- Athalye, A.; Carlini, N.; and Wagner, D. 2018. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. *arXiv preprint arXiv:1802.00420*.
- Buzaglo, G.; Haim, N.; Yehudai, G.; Vardi, G.; Oz, Y.; Nikankin, Y.; and Irani, M. 2023. Deconstructing Data Reconstruction: Multiclass, Weight Decay and General Losses. *arXiv preprint arXiv:2307.01827*.
- Carlini, N.; and Wagner, D. 2017. Towards evaluating the robustness of neural networks. In *IEEE Symposium on Security and Privacy*, 39–57.
- Cazenavette, G.; He, K.; Torralba, A.; Efros, A. A.; and Zhu, J.-Y. 2022. Dataset distillation by matching training trajectories. In *IEEE Conference on Computer Vision and Pattern Recognition*.
- Chen, P.-Y.; Zhang, H.; Sharma, Y.; Yi, J.; and Hsieh, C.-J. 2017. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, 15–26. ACM.
- Chen, Z.; Geng, J.; Zhu, D.; Woisetschlaeger, H.; Li, Q.; Schimmler, S.; Mayer, R.; and Rong, C. 2023. A Comprehensive Study on Dataset Distillation: Performance, Privacy, Robustness and Fairness. *arXiv preprint arXiv:2305.03355*.
- Cisse, M.; Adi, Y.; Neverova, N.; and Keshet, J. 2017a. Houdini: Fooling deep structured prediction models. *arXiv preprint arXiv:1707.05373*.
- Cisse, M.; Bojanowski, P.; Grave, E.; Dauphin, Y.; and Usunier, N. 2017b. Parseval networks: Improving robustness to adversarial examples. In *Proceedings of the 34th International Conference on Machine Learning*, 854–863.
- Cohen, J.; Rosenfield, E.; and Kolter, Z. 2019. Certified adversarial robustness via randomized smoothing. In *International Conference on Machine Learning*.
- Croce, F.; and Hein, M. 2020. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *International Conference on Machine Learning*.
- Cui, J.; Wang, R.; Si, S.; and Hsieh, C.-J. 2023. Scaling Up Dataset Distillation to ImageNet-1K with Constant Memory. In *International Conference on Machine Learning*.
- Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, 248–255. Ieee.
- Dong, Y.; Pang, T.; Su, H.; Zhu, J.; Hu, X.; and Li, J. 2017. Boosting Adversarial Attacks with Momentum. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Fawzi, A.; Moosavi-Dezfooli, S.-M.; Frossard, P.; and Soatto, S. 2018. Empirical Study of the Topology and Geometry of Deep Networks. In *IEEE Conference on Computer Vision and Pattern Recognition*.
- Geng, J.; Chen, Z.; Wang, Y.; Woisetschlaeger, H.; Li, Q.; Schimmler, S.; Mayer, R.; Zhao, Z.; and Rong, C. 2023. A Survey on Dataset Distillation: Approaches, Applications, and Future Directions. *arXiv preprint arXiv:2305.01975*.
- Goodfellow, I. J.; Shlens, J.; and Szegedy, C. 2015. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*.
- Gu, S.; and Rigazio, L. 2014. Towards Deep Neural Network Architectures Robust to Adversarial Examples. *arXiv:1412.5068*.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep Residual Learning for Image Recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, 770–778.
- Howard, J. 2018. Imagenette.
- Jetley, S.; Lord, N.; and Torr, P. 2018. With Friends Like These, Who Needs Adversaries? In *Advances in neural information processing systems*.
- Kurakin, A.; Goodfellow, I.; and Bengio, S. 2017. Adversarial examples in the physical world. *International Conference on Learning Representations Workshops*.
- Le, Y.; and Yang, X. 2015. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7): 3.
- Lee, S.; Chun, S.; Jung, S.; Yun, S.; and Yoon, S. 2022. Dataset Condensation with Contrastive Signals. In *International Conference on Machine Learning*.
- Li, Z.; Guo, Z.; Zhao, W.; Zhang, T.; Cheng, Z.-Q.; Khaki, S.; Zhang, K.; Sajedi, A.; Plataniotis, K. N.; Wang, K.; and You, Y. 2024. Prioritize Alignment in Dataset Distillation. *arXiv:2408.03360*.
- Liu, D.; Gu, J.; Cao, H.; Trinitis, C.; and Schulz, M. 2024. Dataset Distillation by Automatic Training Trajectories. *arXiv:2407.14245*.
- Liu, H.; Chaudhary, M.; and Wang, H. 2023. Towards Trustworthy and Aligned Machine Learning: A Data-centric Survey with Causality Perspectives. *arXiv:2307.16851*.
- Liu, H.; Xing, T.; Li, L.; Dalal, V.; He, J.; and Wang, H. 2023. Dataset Distillation via the Wasserstein Metric. *arXiv preprint arXiv:2311.18531*.
- Ma, S.; Zhu, F.; Cheng, Z.; and Zhang, X.-Y. 2024. Towards trustworthy dataset distillation. *Pattern Recognition*, 110875.
- Madry, A.; Makelov, A.; Schmidt, L.; Tsipras, D.; and Vladu, A. 2017. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*.
- Madry, A.; Makelov, A.; Schmidt, L.; Tsipras, D.; and Vladu, A. 2018. Towards Deep Learning Models Resistant to Adversarial Attacks. In *Proceedings of the International Conference on Learning Representations*.
- Miyato, T.; Ichi Maeda, S.; Koyama, M.; Nakae, K.; and Ishii, S. 2015. Distributional Smoothing with Virtual Adversarial Training. *arXiv:1507.00677*.
- Moosavi-Dezfooli, S.-M.; Fawzi, A.; and Frossard, P. 2016. Deepfool: a simple and accurate method to fool deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2574–2582.

- Moosavi-Dezfooli, S.-M.; Fawzi, A.; Uesato, J.; and Frossard, P. 2019. Robustness via curvature regularization, and vice versa. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 9078–9086.
- Narodytska, N.; and Kasiviswanathan, S. P. 2016. Simple Black-Box Adversarial Perturbations for Deep Networks. *arXiv:1612.06299*.
- Nguyen, T.; Chen, Z.; and Lee, J. 2020. Dataset meta-learning from kernel ridge-regression. *arXiv preprint arXiv:2011.00050*.
- Nguyen, T.; Novak, R.; Xiao, L.; and Lee, J. 2021. Dataset distillation with infinitely wide convolutional networks. *Advances in Neural Information Processing Systems*, 34: 5186–5198.
- Papernot, N.; McDaniel, P.; and Goodfellow, I. 2016. Transferability in Machine Learning: from Phenomena to Black-Box Attacks using Adversarial Samples. *arXiv:1605.07277*.
- Papernot, N.; McDaniel, P.; Jha, S.; Fredrikson, M.; Celik, Z. B.; and Swami, A. 2017. The limitations of deep learning in adversarial settings. In *2016 IEEE European Symposium on Security and Privacy*.
- Papernot, N.; McDaniel, P.; Wu, X.; Jha, S.; and Swami, A. 2016. Distillation as a defense to adversarial perturbations against deep neural networks. In *2016 IEEE Symposium on Security and Privacy (SP)*, 582–597. IEEE.
- Ross, A. S.; and Doshi-Velez, F. 2018. Improving the adversarial robustness and interpretability of deep neural networks by regularizing their input gradients. In *AAAI Conference on Artificial Intelligence*.
- Sachdeva, N.; and McAuley, J. 2023. Data distillation: a survey. *arXiv preprint arXiv:2301.04272*.
- Sucholutsky, I.; and Schonlau, M. 2021. Soft-label dataset distillation and text dataset distillation. In *2021 International Joint Conference on Neural Networks (IJCNN)*, 1–8. IEEE.
- Tramèr, F.; Kurakin, A.; Papernot, N.; Goodfellow, I.; Boneh, D.; and McDaniel, P. 2018. Ensemble adversarial training: attacks and defenses. In *International Conference on Learning Representations*.
- Vahidian, S.; Wang, M.; Gu, J.; Kungurtsev, V.; Jiang, W.; and Chen, Y. 2024. Group Distributionally Robust Dataset Distillation with Risk Minimization. *arXiv:2402.04676*.
- Wang, K.; Zhao, B.; Peng, X.; Zhu, Z.; Yang, S.; Wang, S.; Huang, G.; Bilén, H.; Wang, X.; and You, Y. 2022. CAFE: Learning to Condense Dataset by Aligning Features. In *IEEE Conference on Computer Vision and Pattern Recognition*.
- Wang, T.; Zhu, J.-Y.; Torralba, A.; and Efros, A. A. 2018. Dataset distillation. *arXiv preprint arXiv:1811.10959*.
- Wong, E.; Schmidt, F. R.; and Kolter, J. Z. 2019. Wasserstein Adversarial Examples via Projected Sinkhorn Iterations. *arXiv preprint arXiv:1902.07906*.
- Xu, Y.; Li, Y.-L.; Cui, K.; Wang, Z.; Lu, C.; Tai, Y.-W.; and Tang, C.-K. 2024. Distill Gold from Massive Ores: Bi-level Data Pruning towards Efficient Dataset Distillation. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- Yin, Z.; and Shen, Z. 2023. Dataset Distillation in Large Data Era. *arXiv preprint arXiv:2311.18838*.
- Yin, Z.; Xing, E.; and Shen, Z. 2023. Squeeze, Recover and Relabel: Dataset Condensation at ImageNet Scale From A New Perspective. In *Advances in Neural Information Processing Systems*.
- Zhang, H.; Li, S.; Wang, P.; and Zeng, S., Dan Ge. 2024. M3D: Dataset Condensation by Minimizing Maximum Mean Discrepancy. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*.
- Zhang, J.; Xu, X.; Han, B.; Niu, G.; Cui, L.; Sugiyama, M.; and Kankanhalli, M. 2020. Attacks which do not kill training make adversarial learning stronger. In *International Conference on Machine Learning*, 11278–11287. PMLR.
- Zhao, B.; and Bilén, H. 2021. Dataset condensation with differentiable siamese augmentation. In *International Conference on Machine Learning*.
- Zhao, B.; and Bilén, H. 2023. Dataset condensation with distribution matching. In *IEEE Winter Conference on Applications of Computer Vision*.
- Zhao, B.; Mopuri, K. R.; and Bilén, H. 2021. Dataset condensation with gradient matching. In *International Conference on Learning Representations*.
- Zhou, Y.; Nezhadarya, E.; and Ba, J. 2022. Dataset distillation using neural feature regression. In *Advances in Neural Information Processing Systems*.