

# Towards Adversarially Robust Dataset Distillation by Curvature Regularization

Eric Xue<sup>1</sup>, Yijiang Li<sup>2</sup>, Haoyang Liu<sup>3</sup>, Peiran Wang<sup>4</sup>, Yifan Shen<sup>3</sup>, Haohan Wang<sup>3</sup>

<sup>1</sup>UToronto <sup>2</sup>UC San Diego <sup>3</sup>UIUC <sup>4</sup>UCLA

## Introduction

- Background: The **computational demands** for training deep learning models are continuously growing due to the increasing volume of data, prohibitive to entities with limited computational resources.
- Limitation: **Dataset distillation** offers a means to reduce the size of the data while maintaining its utility. However, little attention has been given to the **adversarial robustness** of models trained on distilled datasets, which is important to practical applications.
- Research question: *How can we embed adversarial robustness into the dataset distillation process, thereby generating datasets that lead to more robust models?*

## Preliminaries

- Problem formulation: Robust dataset distillation can be defined as a tri-level optimization problem, where the goal is to find the optimal distilled dataset  $\mathcal{S}$ , such that training a model  $\theta$  on  $\mathcal{S}$  minimizes the *maximum adversarial loss* of  $\theta$  on the original data distribution  $\mathcal{D}_{\mathcal{T}}$ :

$$\mathcal{S}^* = \arg \min_{\mathcal{S}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_{\mathcal{T}}} \left( \max_{\|\mathbf{v}\| \leq \rho} \ell(\mathbf{x} + \mathbf{v}, y; \theta(\mathcal{S})) \right)$$

subject to  $\theta(\mathcal{S}) = \arg \min_{\theta} \mathcal{L}(\mathcal{S}; \theta)$ .

- A natural approach is to integrate **adversarial training** into the distillation process. However, our empirical studies suggest this approach is ineffective. Even mild adversarial training during distillation significantly degrades clean accuracy of the trained model, which in turn limits robust accuracy substantially.

- Possible reason:
  - Cross-over mixture problem: adversarial examples can cross over the decision boundary, providing misleading training signals
  - Create irrelevant artifacts as a kind of data augmentation (DA), consistent with the observation in SRe<sup>2</sup>L that DA can hurt performance

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

Table 1: Performance comparison between our method and SRe<sup>2</sup>L w/wo adversarial training

## Theory

In the paper, we show that the adversarial loss on a distilled image  $\tilde{\ell}_{\rho}^{adv}(\mathbf{x}')$  is upper bounded as below:

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

Through further analysis, we determine that the term  $\frac{1}{2} \rho^2 \mathbb{E}_{\mathbf{x} \sim D_c} \lambda_1(\mathbf{x})$  is the least constrained term and requires optimization, where

- $D_c$  denotes the distribution of data in class  $c$ , and
- $\lambda_1$  is the largest eigenvalue of the Hessian matrix  $H(\ell(\mathbf{x}))$

## Method

We present **GUARD** (Geometric RegUlarization for Adversarially Robust Dataset).

- Since the adversarial loss of distilled datasets is largely upper-bounded by the curvature of the loss function, our goal is to reduce the value of  $\lambda_1$  by incorporating it into the loss function.
- To bypass the expensive calculation of the Hessian matrix, we use an efficient approximate of lambda as below:

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

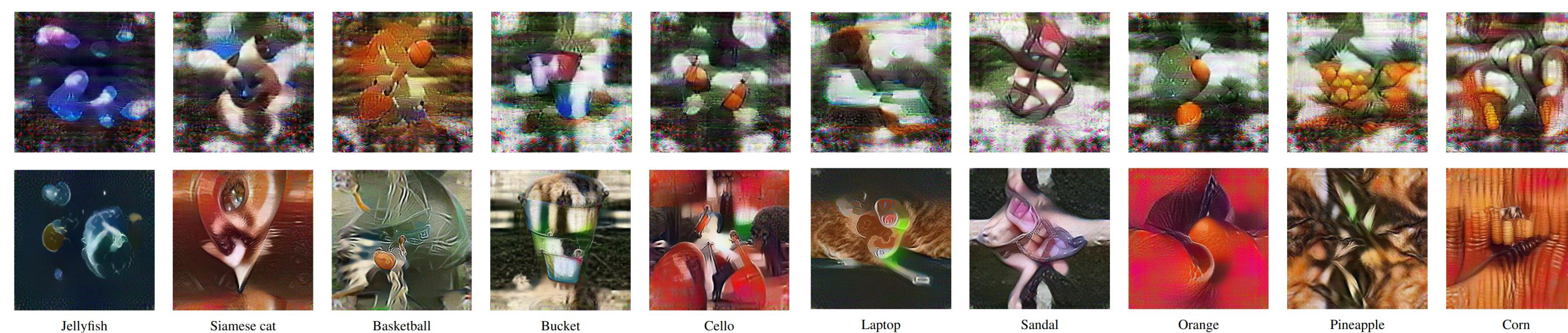
Where  $\mathbf{v}_1$  is the unit eigenvector corresponding to  $\lambda_1$ .

- To make it more efficient, we use the normalized gradient  $\mathbf{z} = \frac{\nabla \ell(\mathbf{x})}{\|\nabla \ell(\mathbf{x})\|}$  to replace  $\mathbf{v}_1$ , inspired by previous works.

The regularized loss function is finally defined as:

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

## Visualization



Upper row – ours; lower row – baseline. Observations:

- Our synthetic data have more distinct object outlines; exhibits high-frequency noise similar to those generated by adversarial attacks

Scan Here  
for Website:



UNIVERSITY OF  
TORONTO

UC San Diego

## Results

Dataset	IPC	Attack	Methods				
			GUARD	SRe <sup>2</sup> L	MTT	TESLA	
TinyImageNet	10	None (Clean)	<b>37.00</b>	33.18	8.14	14.06	
		PGD100	6.39	1.08	4.08	<b>8.40</b>	
		Square	<b>19.53</b>	<u>15.85</u>	2.48	6.31	
		AutoAttack	<u>4.91</u>	0.79	2.44	<b>6.16</b>	
		CW	<b>8.40</b>	3.24	2.50	<u>6.26</u>	
		MIM	<u>6.51</u>	1.10	4.08	<b>8.40</b>	
	50	None (Clean)	<u>55.61</u>	<b>56.42</b>	17.84	28.24	
		PGD100	<b>15.63</b>	0.27	5.62	<u>12.12</u>	
		Square	<b>36.93</b>	<u>15.50</u>	3.84	10.39	
		AutoAttack	<b>13.84</b>	0.16	3.52	<u>10.01</u>	
		CW	<b>20.46</b>	<u>12.12</u>	3.66	10.13	
		MIM	<b>16.09</b>	0.29	5.64	<u>12.12</u>	
	ImageNet-1K	10	None (Clean)	<b>27.25</b>	21.30	-	-
			PGD100	<b>5.25</b>	<u>0.55</u>	-	-
Square			<u>17.88</u>	<b>18.02</b>	-	-	
AutoAttack			<b>3.33</b>	<u>0.34</u>	-	-	
CW			<b>7.68</b>	<u>3.21</u>	-	-	
MIM			<b>5.23</b>	<u>0.51</u>	-	-	
50		None (Clean)	<u>39.89</u>	<b>46.80</b>	-	-	
		PGD100	<b>9.77</b>	0.59	-	-	
		Square	<u>28.39</u>	<b>32.40</b>	-	-	
		AutoAttack	<b>7.03</b>	<u>0.47</u>	-	-	
		CW	<b>14.14</b>	<u>6.31</u>	-	-	
		MIM	<b>9.84</b>	<u>0.64</u>	-	-	

(a)

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

(b)

Table 2: Evaluation results of our GUARD method. (a) Clean accuracy and robustness under various adversarial attacks compared with previous methods; (b) comparison of different DD methods and their GUARD regularized version (marked by †). For more results please see our paper.

- Subfigure (a)(b): GUARD significantly improves robust accuracy across most settings, while maintaining or even improving clean accuracy
- Subfigure (b): GUARD remains effective when applied to other DD methods

## Discussion

- **Robustness Guarantee:** Our theory guarantees the robustness on distilled dataset effectively transfer to the real dataset
- **Computational Overhead:** Only requires an extra forward pass to compute the loss  $\ell(\mathbf{x} + h\mathbf{z})$  within each iteration
  - More than 10x faster than incorporating adversarial training into DD
- **Transferability:** GUARD can theoretically be applied to a large number of dataset distillation methods, beyond the 3 methods in our experiments