

# Better Diffusion Models Further Improve Adversarial Training

Zekai Wang<sup>\*1</sup>, Tianyu Pang<sup>\*2</sup>, Chao Du<sup>2</sup>, Min Lin<sup>2</sup>,  
Weiwei Liu<sup>1</sup>, Shuicheng Yan<sup>2</sup>

1. School of Computer Science, Wuhan University, China

2. Sea AI Lab, Singapore

## Acknowledgment



Zekai Wang



Tianyu Pang



Chao Du



Min Lin

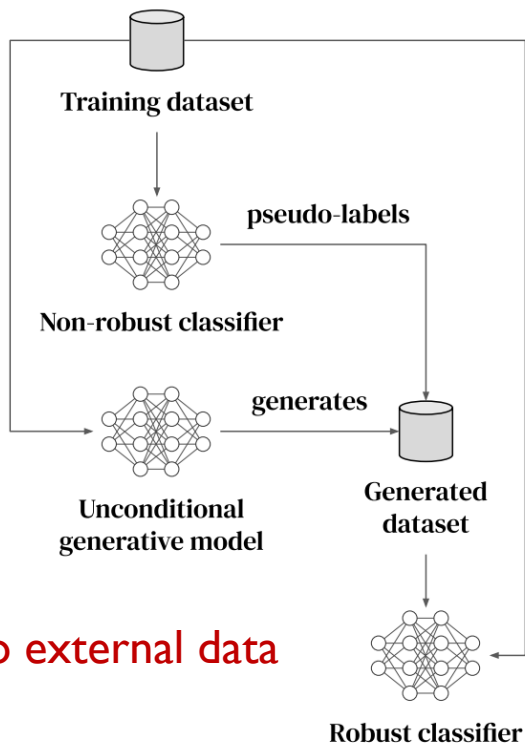


Weiwei Liu

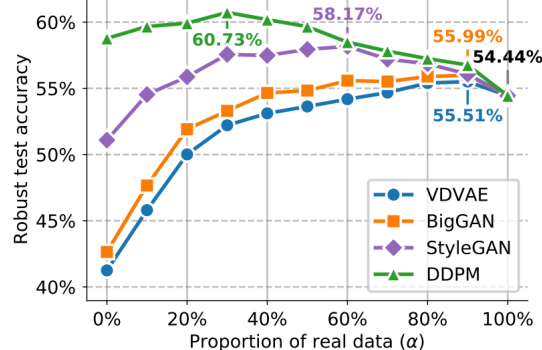
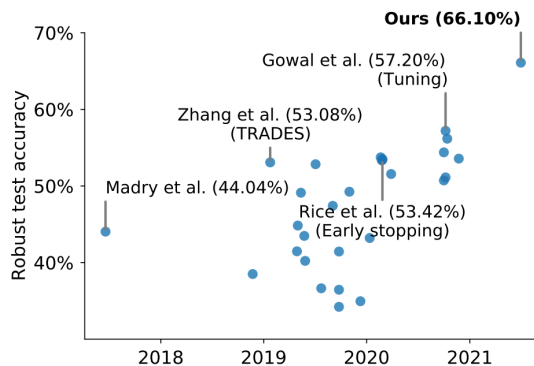


Shuicheng Yan

# Diffusion Models for Trustworthy ML



- **No external data**



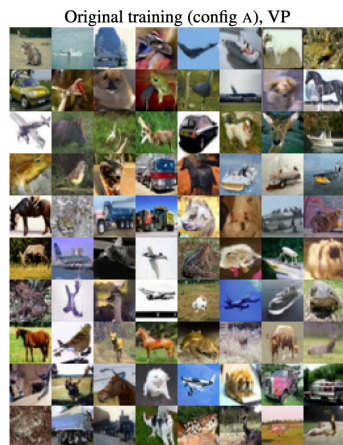
**Dominate**  **ROBUSTBENCH**  
A standardized benchmark for adversarial robustness  
**for two years!**



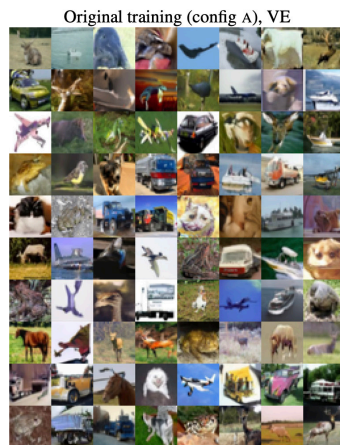
- [1] Rebuffi et al. Fixing Data Augmentation to Improve Adversarial Robustness. NeurIPS 2021
- [2] Gowal et al. Improving Robustness using Generated Data. NeurIPS 2021

# Does Lower FID lead to Better Downstream Performance?

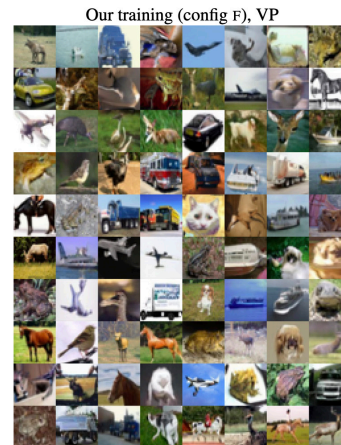
Training configuration	CIFAR-10 [29] at 32×32				FFHQ [27] 64×64		AFHQv2 [7] 64×64	
	Conditional		Unconditional		Unconditional		Unconditional	
	VP	VE	VP	VE	VP	VE	VP	VE
A Baseline [49] (*pre-trained)	2.48	3.11	3.01*	3.77*	3.39	25.95	2.58	18.52
B + Adjust hyperparameters	2.18	2.48	2.51	2.94	3.13	22.53	2.43	23.12
C + Redistribute capacity	2.08	2.52	2.31	2.83	2.78	41.62	2.54	15.04
D + Our preconditioning	2.09	2.64	2.29	3.10	2.94	3.39	2.79	3.81
E + Our loss function	1.88	1.86	2.05	1.99	2.60	2.81	2.29	2.28
F + Non-leaky augmentation	<b>1.79</b>	<b>1.79</b>	<b>1.97</b>	<b>1.98</b>	<b>2.39</b>	<b>2.53</b>	<b>1.96</b>	<b>2.16</b>
NFE	35	35	35	35	79	79	79	79



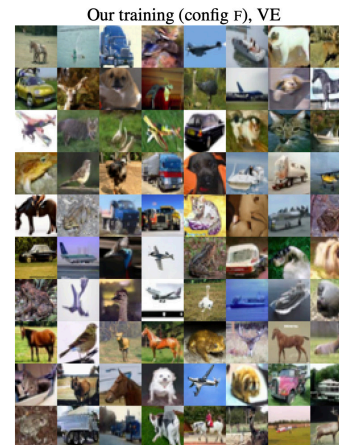
FID 3.01 NFE 35



FID 3.77 NFE 35

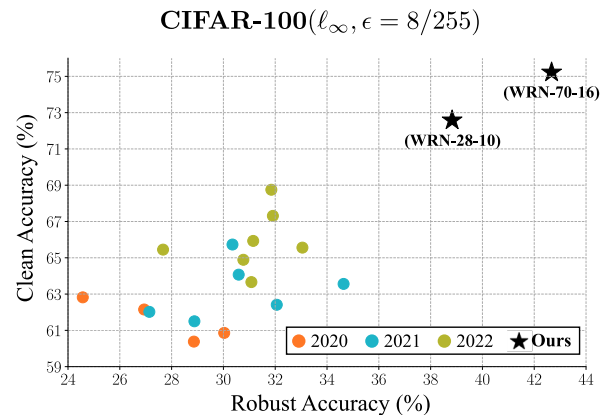
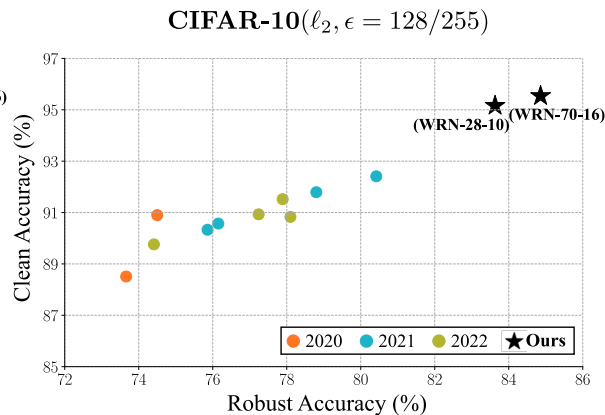
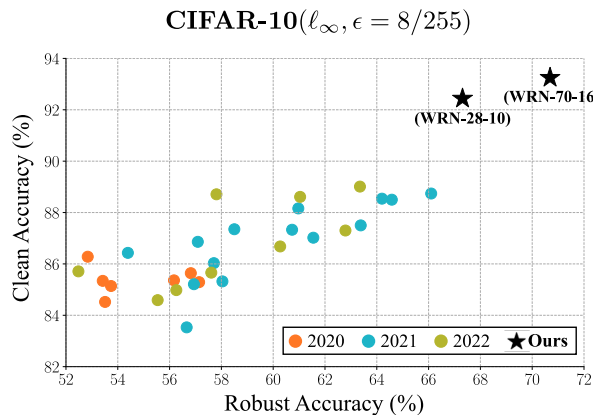


FID 1.97 NFE 35



FID 1.98 NFE 35

# Yes! Better Diffusion Models are Indeed Better



- **New state-of-the-art!**



**ROBUST BENCH**

A standardized benchmark for adversarial robustness

# Yes! Better Diffusion Models are Indeed Better

Table 1. A brief summary comparison of test accuracy (%) between our models and existing Rank #1 models, *with* (✓) and *without* (✗) external datasets, as listed in RobustBench (Croce et al., 2021).

Dataset	Method	External	Clean	AA
<b>CIFAR-10</b> ( $\ell_\infty, \epsilon = 8/255$ )	Rank #1	✗	88.74	66.11
		✓	92.23	66.58
	<b>Ours</b>	✗	<b>93.25</b>	<b>70.69</b>
<b>CIFAR-10</b> ( $\ell_2, \epsilon = 128/255$ )	Rank #1	✗	92.41	80.42
		✓	<b>95.74</b>	82.32
	<b>Ours</b>	✗	95.54	<b>84.86</b>
<b>CIFAR-100</b> ( $\ell_\infty, \epsilon = 8/255$ )	Rank #1	✗	63.56	34.64
		✓	69.15	36.88
	<b>Ours</b>	✗	<b>75.22</b>	<b>42.67</b>

- Even beat previous SOTA that using external datasets
- No extra training time (only extra cost for generating data)

# Yes! Better Diffusion Models are Indeed Better

- Results on SVHN and Tiny-ImageNet

Dataset	Method	Generated	Ratio	Batch	Epoch	Clean	AA
SVHN ( $\ell_\infty, \epsilon = 8/255$ )	Gowal et al. (2021)	$\times$	$\times$	512	400	92.87	56.83
	Gowal et al. (2021)	1M	0.4	1024	800	94.15	60.90
	Rebuffi et al. (2021)	1M	0.4	1024	800	94.39	61.09
	<b>Ours</b>	1M	0.2	1024	800	<b>95.19</b>	<b>61.85</b>
	<b>Ours</b>	50M	0.2	2048	1600	<b>95.56</b>	<b>64.01</b>
TinyImageNet ( $\ell_\infty, \epsilon = 8/255$ )	Gowal et al. (2021)	$\times$	$\times$	512	400	51.56	21.56
	<b>Ours</b>	1M	0.4	512	400	<b>53.62</b>	<b>23.40</b>
	Gowal et al. (2021)*	1M	0.3	1024	800	60.95	26.66
	<b>Ours (ImageNet EDM)</b>	1M	0.2	512	400	<b>65.19</b>	<b>31.30</b>

## Alleviate Overfitting

- The model performs better with a longer training process

Generated	Epoch	Best epoch	Clean			PGD-40			AA		
			Early	Last	Diff	Early	Last	Diff	Early	Last	Diff
x	400	86	84.41	82.18	-2.23	55.23	46.21	-9.02	54.57	44.89	-9.68
	800	88	83.60	82.15	-1.45	53.86	45.75	-8.11	53.13	44.58	-8.55
20M	400	370	91.27	91.45	+0.18	64.65	64.80	+0.15	63.69	63.84	+0.15
	800	755	92.08	92.14	+0.06	66.61	66.72	+0.11	65.66	65.63	+0.03
	1200	1154	92.43	92.32	-0.11	67.45	67.64	+0.19	66.31	66.60	+0.29
	1600	1593	92.51	<b>92.61</b>	+0.10	68.05	67.98	-0.07	67.14	67.10	-0.04
	2000	1978	92.41	92.55	+0.14	68.32	68.30	-0.02	67.22	67.17	-0.05
	2400	2358	<b>92.58</b>	92.54	-0.04	<b>68.43</b>	<b>68.39</b>	-0.04	<b>67.31</b>	<b>67.30</b>	-0.01



# Alleviate Overfitting

- Alleviate overfitting in adversarial training (Amount of Generated Data  $>500K$ )

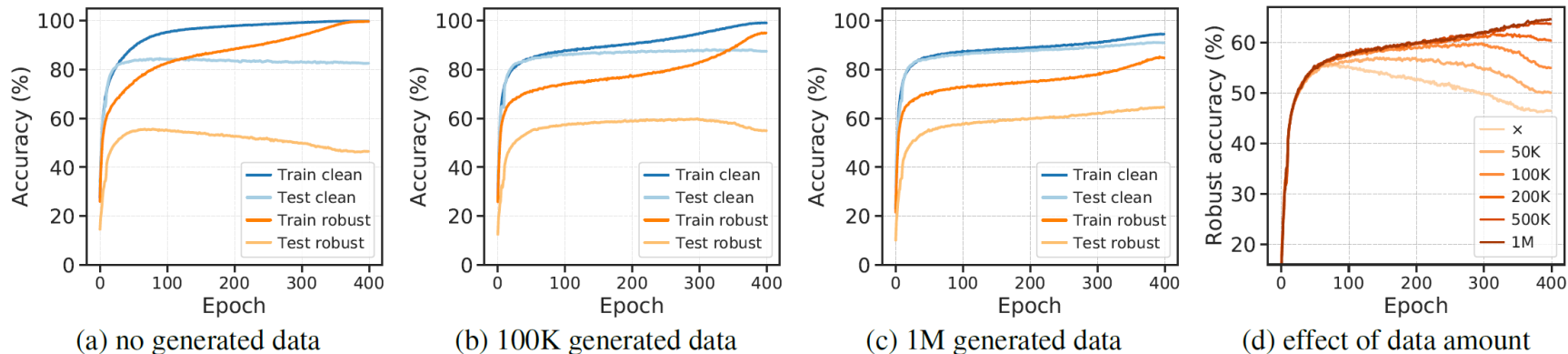


Figure 2. Clean and PGD robust accuracy of AT using (a) no generated data; (b) 100K generated data; (c) 1M generated data. (d) plots the PGD test robust accuracy of AT using different amounts of generated data.

**THANKS**