

# Additional Experiments for ICML Submission 4722: Trustworthy Machine Learning through Data-Specific Indistinguishability

## Experiment 1: Copyright/Contribution Control in Finetuning Diffusion Models with DSI Framework

	Original Artwork	Before Tuning	Regular Tuning ( $\epsilon = +\infty$ )	$\epsilon = 1$	$\epsilon = 2$	$\epsilon = 4$	$\epsilon = 6$	$\epsilon = 8$
Qi.								
Xu.								
Remb.								

Table 1: Finetuning a Stable Diffusion (v1-4) on 425 paintings collected from 10 artists with/out DSI noise on **per-sample contribution** with epoch number selected to be **10**. The original artwork, and the generated images before tuning, after tuning without noise ( $\epsilon = \infty$ ), and after provable-trust tuning with DSI noises and various indistinguishability budget in  $\epsilon$ , ( $\delta = 0.002$ ) of three selected artists, Baishi Qi (Qi.), Beihong Xu (Xu.) and Rembrandt Harmenszoon van Rijn (Remb.) are presented.

	Original Artwork	Stable Diffusion	$\epsilon = +\infty$	$\epsilon = 1$	$\epsilon = 2$	$\epsilon = 4$	$\epsilon = 6$	$\epsilon = 8$
Qi.								
Xu.								
Rem.								

Table 2: Finetuning a Stable Diffusion (v1-4) on 425 paintings collected from 10 artists with/out DSI noise on **per-sample contribution** with epoch number selected to be **20**. The original artwork, and the generated images before tuning, after tuning without noise ( $\epsilon = \infty$ ), and after provable-trust tuning with DSI noises and various indistinguishability budget in  $\epsilon$ , ( $\delta = 0.002$ ) of three selected artists, Baishi Qi (Qi.), Beihong Xu (Xu.) and Rembrandt Harmenszoon van Rijn (Remb.) are presented.

## Experiment 2: Comparison between DSI Local SGD with Indistinguishability Control Methods from DP-SGD

Model	Method	$\epsilon$	1	2	3	4	5	6	7	8
		$\infty$	55.3	63.1	67.6	72.4	73.7	74.3	75.8	76.0
ResNet-20	[XXWD23]	91.7	/	<b>59.7</b>	/	/	<b>70.1</b>	/	/	<b>74.9</b>
	[YZCL21]		57.4	71.2	73.8	78.7	79.8	83.1	83.6	83.9
	DSI-Local-SGD		56.8	64.9	69.2	71.9	74.1	77.0	78.8	79.5
WideResNet-16	[DBH <sup>+</sup> 22]	94.6	<b>57.2</b>	<b>64.6</b>	/	<b>70.5</b>	/	/	/	<b>79.8</b>
	[BPBB23]		63.7	76.1	80.5	82.4	84.4	86.6	86.8	87.1
	DSI-Local-SGD									

Table 3: **Test Accuracy (%)** Comparison between standard distinguishability control through per-sample gradient clipping and isotropic noise in DP-SGD [XXWD23, DBH<sup>+</sup>22] and the augmented versions with additional public data – [YZCL21] projects per-sample gradients into a 2000-rank subspace estimated by public gradients and [BPBB23] conducts mixup between every datapoint and synthetic data – and DSI-Local-SGD with  $\mathcal{O}$  being 20-local-GD-iteration with  $R_i$  as a leaving-one subset of CIFAR-10 training data  $U$  from scratch across different  $\epsilon$  selections with  $\delta = 10^{-5}$ .

### Experiment 3: Comparison between DSI Noise and Isotropic Noise in Defending Backdoor Attacks

Table 4: Comparison on indistinguishability control and defense efficiency in Adversarial Success Rate (ASR) against **Low-Frequency Attacks** [ZPMJ21] between **DSI Noise** and **Isotropic Noise** with **fixed** test accuracy on clean data.

(a): $m = 10$ Sources			(b): $m = 20$ Sources		
Test ACC (%)	Ind. Guarantee	ASR (%)	Test ACC (%)	Ind. Guarantee	ASR (%)
75.6	$(\epsilon = 8, \delta = 10^{-5})$ (DSI) $(\epsilon = 2019, \delta = 10^{-5})$ (Iso.)	<b>13.9</b> 25.6	79.1	$(\epsilon = 8, \delta = 10^{-5})$ (DSI) $(\epsilon = 1661, \delta = 10^{-5})$ (Iso.)	<b>8.4</b> 16.2
65.5	$(\epsilon = 4, \delta = 10^{-5})$ (DSI) $(\epsilon = 575, \delta = 10^{-5})$ (Iso.)	<b>2.0</b> 20.2	70.9	$(\epsilon = 4, \delta = 10^{-5})$ (DSI) $(\epsilon = 479, \delta = 10^{-5})$ (Iso.)	<b>5.3</b> 11.5
(c): $m = 40$ Sources			(d): $m = 80$ Sources		
Test ACC (%)	Ind. Guarantee	ASR (%)	Test ACC (%)	Ind. Guarantee	ASR (%)
78.4	$(\epsilon = 8, \delta = 10^{-5})$ (DSI) $(\epsilon = 627, \delta = 10^{-5})$ (Iso.)	6.4 <b>7.3</b>	79.5	$(\epsilon = 8, \delta = 10^{-5})$ (DSI) $(\epsilon = 350, \delta = 10^{-5})$ (Iso.)	8.9 <b>8.2</b>
73.6	$(\epsilon = 4, \delta = 10^{-5})$ (DSI) $(\epsilon = 312, \delta = 10^{-5})$ (Iso.)	<b>6.6</b> 7.2	74.1	$(\epsilon = 4, \delta = 10^{-5})$ (DSI) $(\epsilon = 114, \delta = 10^{-5})$ (Iso.)	<b>6.1</b> 7.7

Table 5: Comparison on indistinguishability control and defense efficiency in Adversarial Success Rate (ASR) against **Blended Attacks** [CLL<sup>+</sup>17] between **DSI Noise** and **Isotropic Noise** with **fixed** test accuracy on clean data.

(a): 10 Sources			(b): 20 Sources		
Test ACC (%)	Ind. Guarantee	ASR (%)	Test ACC (%)	Ind. Guarantee	ASR (%)
73.9	$(\epsilon = 8, \delta = 10^{-5})$ (DSI) $(\epsilon = 2019, \delta = 10^{-5})$ (Iso.)	15.7 <b>9.9</b>	77.9	$(\epsilon = 8, \delta = 10^{-5})$ (DSI) $(\epsilon = 1661, \delta = 10^{-5})$ (Iso.)	<b>7.7</b> 8.8
57.5	$(\epsilon = 4, \delta = 10^{-5})$ (DSI) $(\epsilon = 350, \delta = 10^{-5})$ (Iso.)	<b>0.1</b> 9.8	67.1	$(\epsilon = 4, \delta = 10^{-5})$ (DSI) $(\epsilon = 241, \delta = 10^{-5})$ (Iso.)	<b>1.1</b> 2.4
(c): 40 Sources			(d): 80 Sources		
Test ACC (%)	Ind. Guarantee	ASR (%)	Test ACC (%)	Ind. Guarantee	ASR (%)
78.8	$(\epsilon = 8, \delta = 10^{-5})$ (DSI) $(\epsilon = 627, \delta = 10^{-5})$ (Iso.)	7.1 <b>4.1</b>	79.3	$(\epsilon = 8, \delta = 10^{-5})$ (DSI) $(\epsilon = 350, \delta = 10^{-5})$ (Iso.)	<b>3.4</b> 3.9
73.5	$(\epsilon = 4, \delta = 10^{-5})$ (DSI) $(\epsilon = 312, \delta = 10^{-5})$ (Iso.)	12.3 <b>5.1</b>	73.7	$(\epsilon = 4, \delta = 10^{-5})$ (DSI) $(\epsilon = 179, \delta = 10^{-5})$ (Iso.)	4.7 <b>3.6</b>

## References

- [BPBB23] Wenzuan Bao, Francesco Pittaluga, Vijay Kumar BG, and Vincent Bindschaedler. Dp-mix: mixup-based data augmentation for differentially private learning. *Advances in Neural Information Processing Systems*, 36:12154–12170, 2023.
- [CLL<sup>+</sup>17] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017.
- [DBH<sup>+</sup>22] Soham De, Leonard Berrada, Jamie Hayes, Samuel L Smith, and Borja Balle. Unlocking high-accuracy differentially private image classification through scale. *arXiv preprint arXiv:2204.13650*, 2022.
- [XXWD23] Hanshen Xiao, Zihang Xiang, Di Wang, and Srinivas Devadas. A theory to instruct differentially-private learning via clipping bias reduction. In *2023 IEEE Symposium on Security and Privacy (SP)*, pages 2170–2189. IEEE, 2023.
- [YZCL21] Da Yu, Huishuai Zhang, Wei Chen, and Tie-Yan Liu. Do not let privacy overbill utility: Gradient embedding perturbation for private learning. In *International Conference on Learning Representations*, 2021.
- [ZPMJ21] Yi Zeng, Won Park, Z Morley Mao, and Ruoxi Jia. Rethinking the backdoor attacks’ triggers: A frequency perspective. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 16473–16481, 2021.