

Fuzz-Testing Meets LLM-Based Agents: An Automated and Efficient Framework for Jailbreaking Text-To-Image Generation Models

Tianshuo Cong¹ Xinlei He² Yun Shen³ Yang Zhang²

Yingkai Dong*, Xiangtao Meng*, Ning Yu*, Zheng Li*†‡, Shanqing Guo*†‡

* School of Cyber Science and Technology, Shandong University

Li*†‡, Shanqing Guo*†‡

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Technology, Shandong University

* School of Cyber Science and Techno

JailFuzzer - A novel prompt-level jailbreak method



B. Methodology

A-1. Text-To-Image Generative Model



A-2. Safety Filter and Jailbreak Attack



- Token-Level: High Perplexity, Low Stealthiness
- Prompt-Level: Fix Pattern, Low Success Rate

C. Evaluation

C-1. Effectiveness

1																		
KS		Target	Safety Filter				One-time Ja			lbreak Prompt				Re-use Jailbreak Prompt				
'	Agent Brain		Type	Method		Byn	ass Rate ((↑)	FID Score (\downarrow)		Queries (\downarrow)		Bynass Rat	te (↑)	(↑) FID Score (↓)			
	Undate		-77-					a	dv. vs. target a		v. vs. real	mean	std		(1)	adv. vs	s. target a	ndv. vs.
((t.	Text-Imag	e text-im	age-classifie	er	100.00%		113.82		132.55	7.04	9.27	50.45%	Ь	15	8.35	177.5
	Update	SD1.4	Text	ter	xt-match		100.00%		122.33		146.27	2.94	3.11	100.009	%	124	4.16	151.3
				text	-classifier		88.30%		104.76		139.43	15.45	14.10	100.009	%	10	0.96	130.4
) →i			_Image	image-classifier			100.00%	0.00% 11		2.63 153.95		6.89 7.26		54.35%		128.82		175.7
				Shinage-clip-classifier			100.00%	121.89			155.75 8.40		10.87	51.49%	Ь	14	8.08	197.4
lake		i.,		dog/cat-image-classifier		fier	97.30%	172.01 (dog/cat)		cat)	-	10.09 14.96		51.38%	b 194.22 (dog/ca		(dog/cat)	-
L		SDXL	Terret	text-match			100.00%		169.29		228.43	4.19 9.90		100.00%		170.04		224.3
	LLaVA		Iext	text-classifier			87.77%		155.21		217.79	11.09 7.4		100.00%		TABLE		219.7
	and			image-classifier			100.00%		184.23		219.43	2.68 3.51		60.97%		196.15		218.0
	Vicuna		Image	image-clip-classifier		er	100.00%		183.74		232.54	3.56 7.70		67.30%		195.06		231.2
				dog/cat-image-classifier			95.95%	185.11 (dog/cat)		cat)	-	6.14	10.17	52.70%		194.32 (dog/cat)		-
		SD3	Teet	text-match			100.00%		160.11		217.70	5.71	5.71 7.50 100.0		% D	ataş	et	225.1
			Iext	text-classifier			89.89%		158.93 2		219.31	11.85	8.87	100.00%		161.27		201.3
ansfer				image-classifier			100.00%		180.51		199.14	2.75 8.08		55.65%		191.46		218.7
main			Image	image-	clip-classifie	er	100.00%		171.85		192.26	3.20	2.73	62.86%	Ь	189	9.01	228.3
				dog/cat-image-classifier		fier	94.15%		181.90 (dog/cat)		-	6.38 10.11		57.26%		191.35 (dog/cat)		-
		DALL·E 3	-	-			81.93%		294.07		309.08	15.26 18.81		67.65%		IFAR-10		284.5
ren	ce w/ T	TA.	-			-						-						
	Safety M	Safety Mechanisms		Bypass rate FID scortarget 94% 142.65 15		core real	re Queri		J:	Jailbreak Defe			By	pass rate tai		FID sc get	real	Quer
	UC					159.90) 5.54	1		None				100% 113 100% 1 C 92% 142		.82	132.55 Do 40	7.04
	PO	POSI SafeGen			161.76	185.31	12.5	5	s	moot	PPL hI I M-In	Insert				4 Г А	168 86	13.7
	Safe				164.72 188.3		.33 16.7	2	S	SmoothLLM-M				88%	133	.84	162.86	12.5
								_			LIIMD			040	121	20	157.27	11.5

C-2. Compare with Baseline



line studios

POWERED BY NETEL

m (Method								
Target	Safety Filter	JailFuzzer	SneakyPrompt	DACA	Ring-A-Bell					
	text-image-classifier	37.56	859.74	42.36	9181.73					
	text-match	34.55	389.56	44.36	15912.84					
SD1.4	text-classifier	30.81	1147.07	80.41	87553.22					
	image-classifier	36.27	708.86	46.05	14532.66					
	image-clip-classifier	32.80	857.38	38.38	6773.42					
	text-match	30.32	423.01	58.30	16474.96					
SDVI	text-classifier	31.37	1082.89	40.65	68108.00					
SDAL	image-classifier	34.43	569.97	54.94	10220.16					
	image-clip-classifier	34.42	440.99	55.16	10268.39					
	text-match	32.35	439.08	56.16	15066.58					
CD2	text-classifier	27.76	618.72	48.19	4984.82					
SD3	image-classifier	38.59	465.89	66.30	12033.25					
	image-clip-classifier	32.74	337.97	61.48	14013.53					
DALL-E 3	-	30.83	797.06	40.69	-					

C-3. Ablation Study



Email: dongyingkai@mail.sdu.edu.cn

Training Process