

Black-box Membership Inference Attacks against Fine-tuned Diffusion Models

Yan Pang, Tianhao Wang

University of Virginia

Diffusion Models are Everywhere

Diffusion Models can generate **photographic-quality** images.



Diffusion Models with its generation ability can be trained as a **game engine**....



Real-time recordings of people playing the game DOOM simulated entirely by the GameNGen neural model

Gamengen.github.io/



...achieving better performance, the current SOTA diffusion model trains on **1 billion images** and fine-tunes with **30 million**.



stable-diffusion-3

stabilityai.com

Diffusion Models are Everywhere



Diffusion Models can generate **photographic-quality** images.



Diffusion Models with its generation ability can be trained as a **game engine**.... ...achieving better performance, the current SOTA diffusion model trains on **1 billion images** and fine-tunes with **30 million**.



Gamengen.github.io/

stable-diffusion-3

MIA against Image Generator









Existing White-box Attacks





Existing White-box Attacks





Existing Black-box Attacks





Existing Gray-box Attacks





Existing Gray-box Attacks





Existing Black-box Attacks



Monte-Carlo Attack:



Membership Identification Function:

$$\hat{f}_{MC-\epsilon}(x) = \frac{1}{k} \sum_{i=1}^{k} \mathbf{1}_{x'_i \in U_r(x)} \\ U_r(x) = \{x' \mid d(x, x') \le r\}$$

[Hilprecht et al. PETS 19].

Existing Black-box Attacks

Monte-Carlo Attack:



Membership Identification Function:

$$\hat{f}_{MC-\epsilon}(x) = \frac{1}{k} \sum_{i=1}^{k} \mathbf{1}_{x'_i \in U_r(x)} \\ U_r(x) = \{x' \mid d(x, x') \le r\}$$

[Hilprecht et al. PETS 19].



GAN-Leaks:



Membership Identification Function: $P(b_i = 1 | x_i, \theta_v) \propto \frac{1}{k} \sum_{i=1}^k \exp(-L(x, \mathcal{G}(z_i)))$ $z_i \sim P_z$

[Chen et al. CCS 2020]. 12



Membership Identification Function:

$$\hat{f}_{MC-\epsilon}(x) = \frac{1}{k} \sum_{i=1}^{k} \mathbf{1}_{x'_i \in U_r(x)} \\ U_r(x) = \{x' \mid d(x, x') \le r\}$$

[Hilprecht et al. PETS 19].

Membership Identification Function: $P(b_i = 1 | x_i, \theta_v) \propto \frac{1}{k} \sum_{i=1}^k \exp(-L(x, \mathcal{G}(z_i)))$ $z_i \sim P_z$

 $\mathcal{R}(x_2|\mathcal{G}_v)$

 $x_1 \in D_{train}$

 $\mathcal{R}(x_1|\mathcal{G}_v)$

13 [Chen et al. CCS 2020].

Can We Improve the Current Attack?



	Target Model Training Set	Shadow Models
GAN Leaks [1]	Need Access	Not utilized
Int. attack [2]	Need Access	Not utilized
Pixel attack [3]	Need Access	Not utilized
Dist. Attack [4]	Partial Access	Not utilized

Need Access: Use target model's training set as the member set.



[1] Matsumoto et al., IEEE DLSP' 23, [2] Wu et al., Arxiv 2022,[3] Dubinski et al., WACV 2024, [4] Zhang et al., WACV 2024

Can We Improve the Current Attack?



	Target Model Training Set	Shadow Models
GAN Leaks [1]	Need Access	Not utilized
Int. attack [2]	Need Access	Not utilized
Pixel attack [3]	Need Access	Not utilized
Dist. Attack [4]	Partial Access	Not utilized

Need Access: Use target model's training set as the member set.

Existing solutions target the unconditional generator, but the current SOTA model is Stable Diffusion. The model's properties have changed, so the attack needs to be redesigned.

Due to time consumption and computing resource issues, existing solutions make "interesting" assumptions.

[1] Matsumoto et al., IEEE DLSP' 23, [2] Wu et al., Arxiv 2022,[3] Dubinski et al., WACV 2024, [4] Zhang et al., WACV 2024

Can We Improve the Current Attack?



	Target Model Training Set	Shadow Models
GAN Leaks [1]	Need Access	Not utilized
Int. attack [2]	Need Access	Not utilized
Pixel attack [3]	Need Access	Not utilized
Dist. Attack [4]	Partial Access	Not utilized

Need Access: Use target model's training set as the member set.

Existing solutions target t generator, but the curren Diffusion. The model's properties name changed, so the attack needs to be redesigned.

[1] Matsumoto et al., IEEE DLSP' 23, [2] Wu et al., Arxiv 2022,[3] Dubinski et al., WACV 2024, [4] Zhang et al., WACV 2024



on and computing 3 solutions make ns.

Theoretical Foundation





Our Theorem: Assuming we have a pre-trained diffusion model \hat{x}_{θ} with its training set D_{m} , and use a bit b to represent the membership of query sample x {(1 for member and 0 for non-member)}.

$$\Pr\{b = 1 | x, \theta\} \propto - \|x_0 - \hat{x}_{\theta}(x_t, t)\|_2^2$$

Step 1: Synthesize Replicate Samples



No text? No problem! BLIP comes in to fill the gap for any textless query points.

A sunflower, Vincent Willem van Gogh style.



Diffusion Model

Replicate Images

Query Point

Step 2: Calculate Similarity Scores





Similarity Score Vector

Step 3: Execute the Attack





Image embedding color indicates which embedding generated the similarity score with query image.

Three attack models: thresholdbased, distribution-based, classifier-based... all trained with shadow models.





Model: Stable Diffusion v1-5.

Dataset: MS COCO, CelebA-Dialog, WIT.

Metric: Accuracy, AUC ROC, True Positive Rate at low False Positive Rate.

Baseline: Matsumoto et al.[1], Zhang et al.[2]





[1] Matsumoto et al., IEEE DLSP' 23, [2] Zhang et al., IEEE WACV' 23



Based on Query Sample Component and Auxiliary Dataset :

with text	Attack-I	Attack-IV
without text	Attack-II	Attack-III
	Aux ∩ Tar ≠ Ø	Aux \mathbf{n} Tar = $\mathbf{\phi}$



Based on Query Sample Component and Auxiliary Dataset :

♠

with t	with text Attack-I		<-I		Attack-IV		
withou	t text	Attack	<-11		Attack-III		
		Aux n Ta	ar≠Ø		Aux ∩ Tar = Ø	→	Experimental settings fixed; three replicate samples per query. Our
		Attack type		Celeb	A-Dialog		attack outperforms baseline methods.
	1	Attack type	ASR	AUC	TPR@FPR=1%		
	Mats	sumoto et al.[1]	0.52	0.50	0.01		
	Zł	nang et al.[2]	0.51	0.49	0.01		
		Attack-I	0.85	0.93	0.53		
		Attack-II	0.88	0.93	0.60		
		Attack-III	0.87	0.94	0.54		
		Attack-IV	0.87	0.93	0.57		

[1] Matsumoto et al., IEEE DLSP' 23, [2] Zhang et al., IEEE WACV' 23



Based on Query Sample Component and Auxiliary Dataset :

t

with t	ext	Attac	k-1		Attack-IV			
withou	t text	Attack	<-11		Attack-III			
		Aux n Ta	ar≠Ø		Aux ∩ Tar = Ø	→		Experimental settings fixed; three replicate samples per query. Our
-		Attack type		Celeb	A-Dialog			attack outperforms baseline methods.
	1	Attack type	ASR	AUC	TPR@FPR=1%			
-	Mats	sumoto et al.[1]	0.52	0.50	0.01			
	Zł	nang et al.[2]	0.51	0.49	0.01			
		Attack-I	0.85	0.93	0.53		Als	so did ablations, details in paper
		Attack-II	0.88	0.93	0.60			
		Attack-III	0.87	0.94	0.54			
_		Attack-IV	0.87	0.93	0.57			

[1] Matsumoto et al., IEEE DLSP' 23, [2] Zhang et al., IEEE WACV' 23

Case Study





Civitai: The Home of Open-Source Generative Al

Case Study



One of the Original Painting



Download Checkpoint

Member Sample



Original Model



Non-Member Sample

Civitai: The Home of Open-Source Generative Al

Use our attack as an auditing tool target models on Civitai:

ivitai: Th	e Home c	of Open-S	Source Ge	enerative	A

We observe a clear distinction in similarity scores between members and non-members.

Art Style (Artist)	Member	Non-member	Diff.	ROC-AUC
Vincent van Gogh	0.92	0.44	0.48	0.91
Baishi Qi	0.77	0.40	0.37	0.88
Ukiyo-e	0.70	0.45	0.25	0.87
Uemura Shoen	0.88	0.41	0.47	0.89
Wanostyle	0.80	0.49	0.31	0.87
Ken Kelly	0.80	0.47	0.33	0.89
Shanshui Painting	0.88	0.47	0.41	0.86



Case Study

Use our attack as an auditing tool target models on Civitai:

Art Style (Artist)	Member	Non-member	Diff.	ROC-AUC
Vincent van Gogh	0.92	0.44	0.48	0.91
Baishi Qi	0.77	0.40	0.37	0.88
Ukiyo-e	0.70	0.45	0.25	0.87
Uemura Shoen	0.88	0.41	0.47	0.89
Wanostyle	0.80	0.49	0.31	0.87
Ken Kelly	0.80	0.47	0.33	0.89
Shanshui Painting	0.88	0.47	0.41	0.86

Our method can serve as an effective auditing tool to evaluate the model on Civitai.







- We find that membership inference attacks against GANs and VAEs cannot be directly used to target current diffusion models due to architectural and generative differences.
- We propose an attack pipeline targeting conditional diffusion models across four scenarios, designing three attack models and evaluating key factors on three datasets.

Paper (with links to Github and dataset):

