
Stable Diffusion is Unstable

(Supplementary Material)

Anonymous Author(s)

Affiliation

Address

email

1 In this supplementary material, we first present a review of related works (Section A), including the
2 diffusion model and studies about the vulnerabilities within text-to-image models. Following that, we
3 delve into additional analyses concerning the vulnerabilities observed in the Stable Diffusion model
4 (Section B). Subsequently, we offer instances of long and short prompt attacks, accompanied by the
5 corresponding generated images, as well as instances of black-box attacks (Section C). Lastly, we
6 undertake a comprehensive series of experiments to substantiate the effectiveness of our approach
7 (Section D). These experiments include the evaluation of attacks targeting both long and short
8 prompts. Additionally, ablation studies are conducted to explore attacks employing different search
9 steps, assess the influence of our fluency and semantic similarity constraints on text similarity, and
10 target diverse samplers (e.g., DDIM and DPM-Solver) in the attack process.

11 A Related Work

12 A.1 Diffusion Model.

13 Recently, the diffusion probabilistic model [20] and its variants [6, 13, 21, 18, 19] have achieved
14 great success in content generation [21, 7, 19], including image generation [6, 21], conditional
15 image generation [18], video generation [7, 24], 3D scenes synthesis [10] and so on. Specifically,
16 DDPM [6] adds noises to images and learns to recover images from noises step by step. Then,
17 DDIM [21] improves the generation speed of the diffusion model by skipping steps inference. Then,
18 the conditional latent diffusion model [18] formulates the image generation in latent space guided
19 by multiple conditions, such as texts, images, and semantic maps, further improving the inference
20 speed and boarding the application of the diffusion model. Stable diffusion [18], a latent text-to-
21 image diffusion model capable of generating photo-realistic images given any text input, and its
22 enhanced versions [25, 8, 12], have been widely used in current AI-generated content products, such
23 as Stability-AI [22], Midjourney [11], DALL-E2 [15], and Runaway [3]. However, these methods
24 and products cannot always generate satisfactory results from the given prompt. Therefore, in this
25 work, we aim to analyze the robustness of stable diffusion in the generation process.

26 A.2 Vulnerabilities in Text-to-image Models.

27 With the open-source of Stable Diffusion [18], text-to-image generation achieves great process
28 and shows the unparalleled ability on generating diverse and creative images with the guidance
29 of a text prompt. However, there are some vulnerabilities have been discovered in existing works
30 [4, 1, 23]. Typically, StructureDiffusion [4] discovers that some attributes in the prompt are not
31 assigned correctly in the generated images, thus they employ consistency trees or scene graphs to
32 enhance the embedding learning of the prompt. In addition, Attend-and-Excite [1] also introduces
33 that the Stable Diffusion model fails to generate one or more of the subjects from the input prompt and
34 fails to correctly bind attributes to their corresponding subjects. These pieces of evidence demonstrate
35 the vulnerabilities of the current Stable Diffusion model. However, to the best of our knowledge, no
36 work has systematically analyzed the vulnerabilities of the Stable Diffusion model, which is the goal
37 of this work.

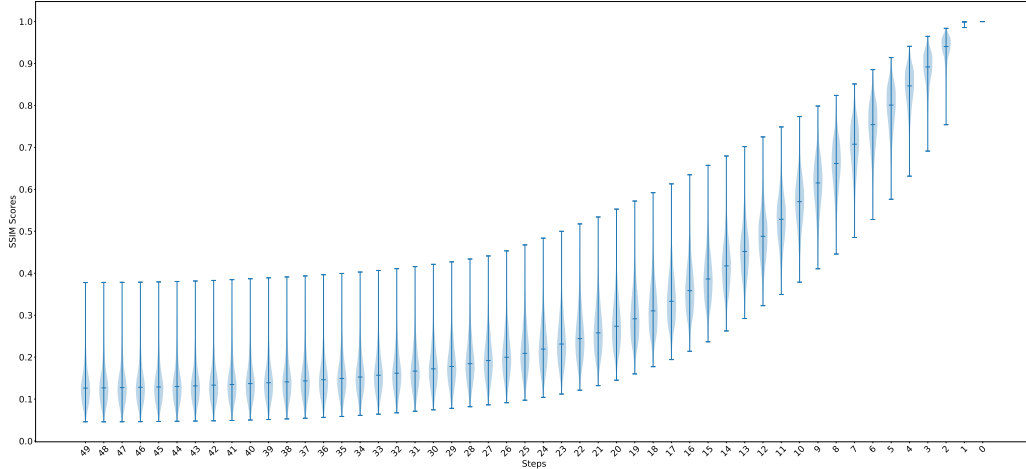


Figure B.1: A violin plot illustrating the generation speeds of 1,000 images of various classes. The horizontal axis represents the number of steps taken, ranging from 49 to 0, while the vertical axis displays the SSIM scores. The width of each violin represents the number of samples that attained a specific range of SSIM scores at a given step.

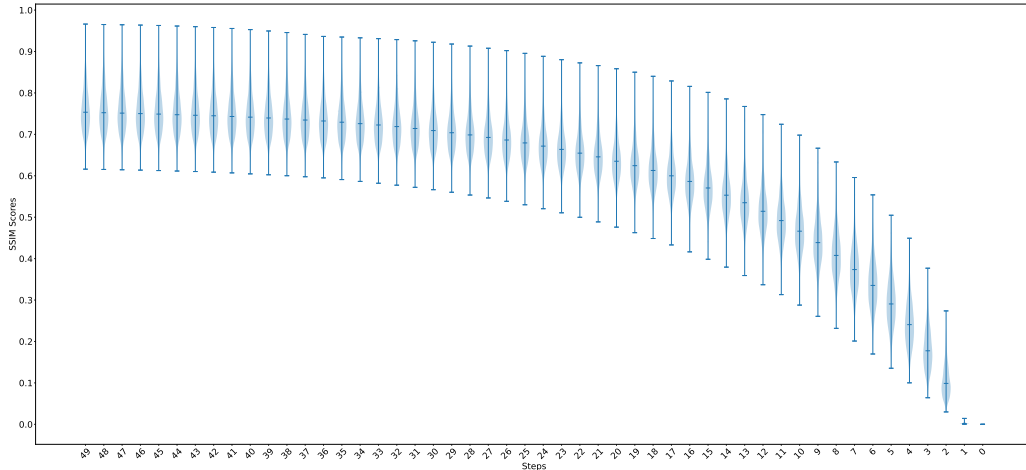


Figure B.2: A violin plot illustrating the generation speeds of 1,000 images of various classes. The horizontal axis represents the number of steps taken, ranging from 49 to 0, while the vertical axis displays the LPIPS scores. The width of each violin represents the number of samples that attained a specific range of LPIPS scores at a given step.

38 B Vulnerabilities of Stable Diffusion Model

39 B.1 Pattern 1: Variability in Generation Speed

40 Fig. B.1 demonstrates the entire 50-step violin diagram which has been discussed before. To eliminate
 41 possible bias due to a single metric, we further verified the difference in generation speed of one
 42 thousand images based on the LPIPS [26] metric, as shown in Fig. B.2. The calculation of the LPIPS
 43 distance from the images generated at each stage to the ultimate image is performed. The horizontal
 44 axis signifies the range of steps from 49 down to 0, whereas the vertical axis denotes the respective
 45 LPIPS scores. Each violin plot illustrates the distribution of the LPIPS scores associated with 1,000
 46 images at a specific step. The width of the violin plot is proportional to the frequency at which images
 47 achieve a certain score. During the initial stages of generation, the distribution's median is situated
 48 nearer to the maximum LPIPS value, suggesting a preponderance of classes demonstrates slower
 49 generation velocities. Nonetheless, the existence of a low minimum value indicates the presence of
 50 classes that generate at comparatively faster rates. As the generation transitions to the intermediate
 51 stages, the distribution's median progressively decreases, positioning itself between the maximum and

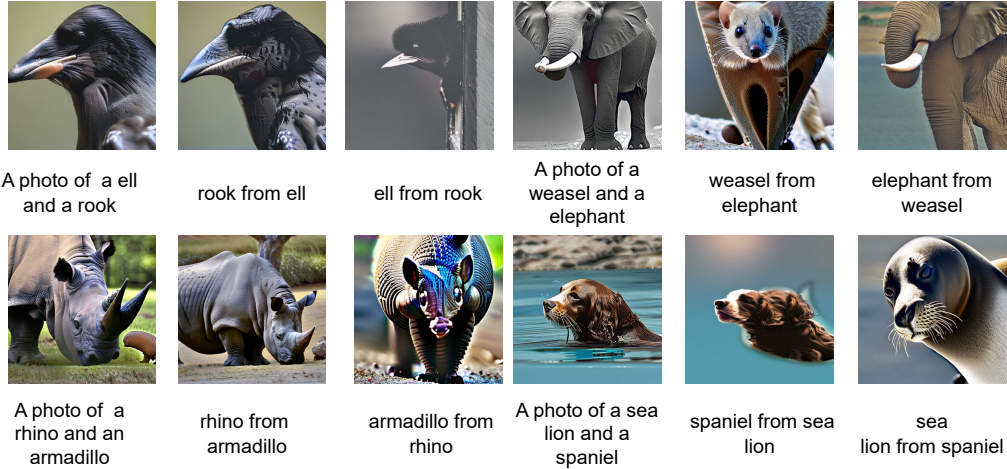


Figure B.3: The image caption, "A photo of class A and class B " represents the generated image when feature entanglement occurs; And "class A from class B " represents the final generated image of prompt "A photo of class A " based on the eighth step of the prompt "A photo of class B "



Figure B.4: The images in the first row are generated by the prompt "A photo of a **warthog**". The images in the second row are generated by the prompt "A photo of a **traitor**". The images in the third row are generated by the prompt "A photo of a **warthog** and a **traitor**".

52 minimum LPIPS values. In the concluding stages of generation, the distribution's median is found
 53 closer to the minimum LPIPS value, implying that the majority of classes are nearing completion.
 54 However, the sustained high maximum value suggests that there are classes still exhibiting slower
 55 generation rates.

56 **B.2 Pattern 2: Similarity of Coarse-grained Characteristics**

57 To further verify that coarse-grained feature similarity is the root cause of feature entanglement, we
 58 provide more cases in Fig. B.3. From these cases, we can see that for the two classes where feature
 59 entanglement can occur, they can both continue the image generation task based on each other's
 60 coarse-grained information.

61 **B.3 Pattern 3: Polysemy of Words**

62 As shown in Fig. B.4, when we attack the prompt "A photo of a warthog" to "A photo of a warthog
 63 and a traitor", the original animal warthog becomes an object similar to a military vehicle or military
 64 aircraft, while the images generated by attack prompt is not directly related to the image of the animal
 65 warthog or traitor. From the t-SNE visualization (Fig. B.5), we can see that the distance from the
 66 picture generated by the attack prompt to the text "a photo of a warthog" has a similar distance to the
 67 animal warthog picture to the text, so we can see that by attacking the original category word that
 68 guided the original category word (animal warthog) into its alternative meaning.

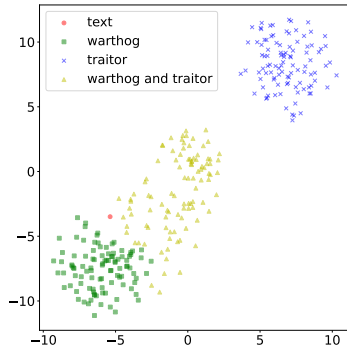


Figure B.5: t-SNE Visualization of 100 images each of "warthog", "traitor", "warthog and traitor" and text "a photo of a warthog."

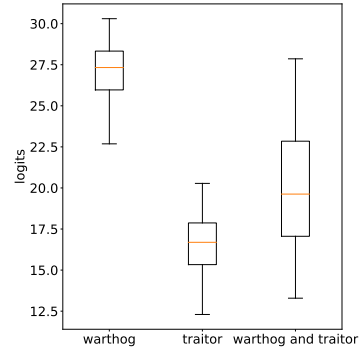


Figure B.6: The boxplot of cosine similarities between the text embedding of "a photo of a warthog" and 100 of image embeddings each of "warthog", "traitor", and "warthog and traitor".



Figure B.7: a) "A type of footwear with a thick, rigid sole, often made of wood, and an upper made of leather or another material. Clogs can be open or closed, and are commonly associated with Dutch and Scandinavian cultures." b) "footwear" is replaced by "pistol". c) "Dutch" is replaced by "pistol".



Figure B.8: A template, "A photo of A, B and C", is used to generate prompts, where $A, B, C \in \{ \text{"cat"}, \text{"pistol"}, \text{"clogs"} \}$. For example, "0-1-2" represents $A = \text{"cat"}, B = \text{"pistol"}$ and $C = \text{"clogs"}$, and so on.

69 From the box plots (Fig. B.6), it can be observed that the image of "warthog" exhibits the highest
70 similarity with the prompt's embedding, while the image of "traitor" demonstrates the lowest similar-
71 ity, as anticipated. Simultaneously, the similarity distribution between the images of "warthog" and
72 "traitor" with the prompt text is relatively wide, indicating that some images have a high similarity
73 with "warthog," while others lack features associated with "warthog."

74 B.4 Pattern 4: Positioning of Words

75 In addition to the three aforementioned observations and patterns outlined in the paper, there is a
76 fourth observations (Observation 4), which is related to positioning of words.

77 **Observation 4.** When a text prompt contains a noun A representing the object to be generated,
78 there exists a preceding word B and a succeeding word C around noun A . When replacing either
79 word B or C with another noun D , for certain instances of noun A , replacing word B results in the
80 generation of an image containing noun D , while replacing word C still results in the generation
81 of an image containing noun A . Conversely, for other instances of noun A , the opposite scenario
82 occurs.

83 An example of Pattern 4 is shown in Fig. B.7. When "footwear" is replaced by "pistol", the generated
84 image contains a pistol instead of clogs. However, when "Dutch" is replaced by "pistol", the model

85 still generates an image of clogs. In addition to differences in the words being replaced, a significant
 86 distinction between the two aforementioned examples of success and failure lies in the relative
 87 positioning of the word being replaced with respect to the target class word. We hypothesize that this
 88 phenomenon occurs due to the different order of the replaced words B or C with respect to the noun
 89 A . To exclude the effects of complex contextual structures, a template for a short prompt, "A photo
 90 of A , B and C ", is used, and the order of A , B , and C are swapped (Fig. B.8).

91 When these sentences with different sequences of category words are understood from a human
 92 perspective, they all have basically the same semantics: both describe a picture containing a cat,
 93 clogs, and a pistol. However, in the processing of language models (including CLIP), the order of
 94 words may affect their comprehension. Although positional encoding provides the model with the
 95 relative positions of words, the model may associate different orders with different semantics through
 96 learned patterns. Therefore, we propose our Pattern 4.

97 **Pattern 4** (Positioning of Words). *Let \mathcal{V} denote a set of vocabulary. Let $\mathcal{N} \subset \mathcal{V}$ denote the subset*
 98 *of all nouns in the vocabulary. Consider a text prompt containing noun $A \in \mathcal{N}$ representing the*
 99 *object to be generated. Furthermore, assume there exist preceding word $B \in \mathcal{V}$ and succeeding word*
 100 *$C \in \mathcal{V}$ surrounding noun A . There exists a condition-dependent behavior regarding the replacement*
 101 *of words B and C with another noun $D \in \mathcal{N}$:*

$$\left\{ \begin{array}{l} \exists A, D \in \mathcal{N}, \quad \exists B, C \in \mathcal{V}, \quad P(B \rightarrow D) \xrightarrow{\text{generate}} D \quad \wedge \quad P(C \rightarrow D) \xrightarrow{\text{generate}} A; \\ \exists A, D \in \mathcal{N}, \quad \exists B, C \in \mathcal{V}, \quad P(B \rightarrow D) \xrightarrow{\text{generate}} A \quad \wedge \quad P(C \rightarrow D) \xrightarrow{\text{generate}} D. \end{array} \right.$$

102 C Cases of Short/Long-Prompt Attacks and Black-box Attacks

103 C.1 Attack on Long Prompt

104 In Fig. C.1, we demonstrate more cases of long text prompt attacks.

105 C.2 Attack on Short Prompt

106 In Fig. C.2, we demonstrate more cases of long text prompt attacks.

107 C.3 Black-box Attack

108 In Fig. C.3, and Fig. C.4, we demonstrate black box attacks targeting mid-journey and DALL·E2,
 109 respectively.

110 D Experiments

111 In our experiments, we conduct comprehensive analyses of both long and short prompts. Furthermore,
 112 we conduct ablation studies specifically on long prompts, focusing on three key aspects. Firstly, we
 113 evaluate our attack method with different numbers of search steps T . Secondly, we investigate the
 114 influence of our constraints, including fluency and semantic similarity as measured by BERTScore.
 115 Lastly, we attack different samplers, including DDIM [21] and DPM-Solver [9].

116 D.1 Experimental Setting.

117 **Attack hyperparameters.** The number of search iterations T is set to 100. This value determines
 118 the number of iterations in the search stage, during which we aim to find the most effective attack
 119 prompts. The number of attack candidates N is set to 100. This parameter specifies the number of
 120 candidate attack prompts considered in the attack stage, allowing for a diverse range of potential
 121 attack prompts to be explored. The learning rate η for the matrix ω is set to 0.3. The margin κ in the
 122 margin loss is set to 30.

123 **Text prompts.** Our experiments consider the 1,000 classes from ImageNet-1K [2], which serves
 124 as the basis for generating images. To explore the impact of prompt length, we consider both
 125 short and long prompts. For clean short prompts, we employ a standardized template: "A photo of
 126 [CLASS_NAME]". Clean long prompts, on the other hand, are generated using ChatGPT 4 [16],
 127 with a prompt length restriction of 77 tokens to align with the upper limit of the CLIP [17] word
 128 embedder.

129 **Evaluation metrics.** To evaluate the effectiveness of our attack method, we generate attack prompts
 130 from the clean prompts. We focus on three key metrics: **success rate**, Fréchet inception distance [5]







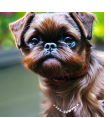
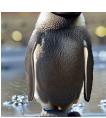








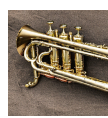



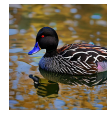





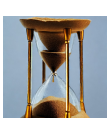



 <p>A loose, full-length abaya worn by some Muslim women, primarily in the Arabian Peninsula, as an outer covering to ensure modesty.</p>	 <p>a knee, full adult or abaya embroidery by some German women, primarily in the netherlands peninsula as an outer covering to ensure modesty.</p>	 <p>An ancient abacus tool, consisting of a frame with rows of beads or other movable counters, used for mathematical operations and counting.</p>	 <p>An optical abacus tool, consisting of a wheel with rows of tiles and other movable symbols, used for wheel operations and counting.</p>	 <p>A large, mounted cannon that fires heavy projectiles, historically used in warfare and for other purposes such as ceremonial salutes."</p>	 <p>A large, hand cannon that fires heavy projectiles, historically used in warfare and for other purposes such as ceremonial salutes."</p>
 <p>This small brussels griffon breed has a distinctive, wiry coat, a compact, muscular body, and a friendly, affectionate temperament. They were originally bred for hunting small game and as companion animals in Belgium and are known for their intelligence, agility, and charm</p>	 <p>This ancient brussels griffon penguins has a distinctive, wiry appearance, a strong, muscular body, and a good, affectionate temperament. they were originally bred for hunting small game and for pack animals in Belgium and are known for their intelligence, agility, and speed.</p>	 <p>A type of footwear with a thick, rigid sole, often made of wood, and an upper made of leather or another material. Clogs can be open or closed, and are commonly associated with Dutch and Scandinavian cultures.</p>	 <p>A type of camouflage with a long, rigid beak, often made of wood, and an upper made of leather or another material. Clogs can be open or closed, and are commonly associated with Dutch and Scandinavian cultures.</p>	 <p>"A loose, sleeveless cloak that is worn over the shoulders and fastened at the neck, often used for warmth or as a decorative accessory."</p>	 <p>"A loose, cat cloak that is worn over the shoulders and fastened at the neck, often used for warmth or as a decorative accessory."</p>
 <p>A computer mouse input device used to control a computer's cursor, typically featuring buttons and a scroll wheel for selecting and navigating on-screen elements.</p>	 <p>A computer mouse puzzle device used to control a computer's cursor, typically featuring buttons and a scroll wheel for selecting and navigating on-screen elements.</p>	 <p>These large, conch snails have a distinctive, spiral shell and a strong, muscular foot. They are found in tropical waters and are known for their ornate shells and edible flesh."</p>	 <p>These large, conch sharks have a distinctive, spiral shell and a strong, muscular foot. They are found in tropical waters and are known for their ornate shells and edible flesh."</p>	 <p>A comet musical instrument similar to a trumpet but with a more compact shape and mellower tone.</p>	 <p>A cornet keyboard is similar to a trumpet but with a more compact shape and mellower tone.</p>
 <p>"A large, upholstered couch of furniture designed for seating multiple people, often found in living rooms or other communal spaces."</p>	 <p>"A large, upholstered couch of weapons designed for seating multiple people, often found in living rooms or other communal spaces."</p>	 <p>"These duck have a round body, a flat bill, and webbed feet, adapted for swimming and dabbling for food. They come in various colors and patterns, with some species having bright, iridescent plumage."</p>	 <p>These duck shark a round body, a flat bill, and webbed feet, adapted for swimming and dabbling for food. They come in various colors and patterns, with some species having bright, iridescent plumage."</p>	 <p>A rich, creamy, dairy-based eggnog traditionally made of milk, cream, sugar, and eggs, often spiked with alcohol and served during the winter holidays.</p>	 <p>A doll, creamy, dairy-based eggnog traditionally made of milk, cream, sugar, and eggs, often spiked with alcohol and served during the winter holidays.</p>
 <p>Electric ray is flattened, disc-shaped fish with smooth skin and a pair of kidney-shaped electric organs on each side of their body, used for hunting prey and self-defense."</p>	 <p>Electric rays flattened, disc-shaped senses with smooth skin and a pair of kidney-shaped electric organs on each side of their body, used for hunting prey and self-defense."</p>	 <p>A timekeeping hourglass consisting of two glass chambers connected by a narrow passage, with sand flowing through the passage at a constant rate to measure a specific time interval.</p>	 <p>A mean like hourglass consists of two renrepanels separated by a central passage, with sunlight flowing through the cave in a constant rate to enter a specific time frame."</p>	 <p>A wild pig species native to Africa, known for its large, curved tusks and distinctive facial features, including warts on their faces. Warthogs are primarily grazers and live in savannas and grasslands.</p>	 <p>A wild rhinospecies native to Africa, known for its large, curved tusks and distinctive facial features, including warts on their faces. Warthogs are primarily grazers and live in savannas and grasslands.</p>

Figure C.1: To the left of the arrow is the clean long text prompt (highlighted by green) and its corresponding image, to the right of the arrow is the generated attack prompt (highlighted by red) and its corresponding image. (Section C.1 Attack on Long Text Prompt)

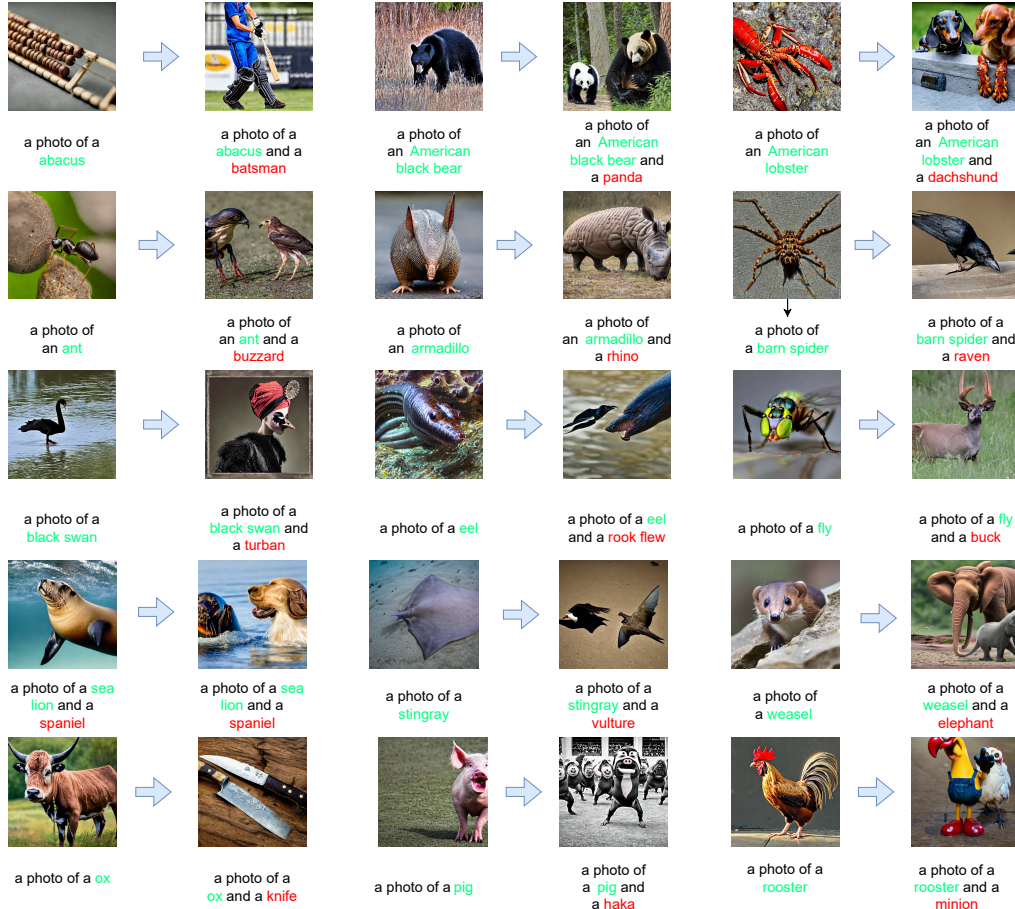


Figure C.2: To the left of the arrow is the clean short text prompt (highlighted by green) and its corresponding image, to the right of the arrow is the generated attack prompt (highlighted by red) and its corresponding image. (Section C.2 Attack on Short Text Prompt)

Table D.1: Main results of short-prompt and long-prompt attacks.

Prompt	Method	Success (%)	FID (↓)	IS (↑)	TS (↑)
Short	Clean	-	18.51	101.33±1.80	1.00
	Random	79.2	29.21	66.71±0.87	0.69
	ATM (Ours)	91.1	30.09	65.98±1.10	0.72
Long	Clean	-	17.95	103.59±1.68	1.00
	Random	41.4	24.16	91.33±1.58	0.94
	ATM (Ours)	81.2	29.65	66.09±1.83	0.84

131 (FID), Inception Score (IS), and text similarity (TS). Subsequently, 50,000 images are generated
 132 using the attack prompts, ensuring a representative sample of 50 images per class. The success
 133 rate is determined by dividing the number of successful attacks by the total of 1,000 classes. FID
 134 and IS are computed by comparing the generated images to the ImageNet-1K validation set with
 135 (torch-fidelity)[14]. TS is calculated by embedding the attack prompts and clean prompts using the
 136 CLIP [17] word embedder, respectively. Subsequently, the cosine similarity between the embeddings
 137 is computed to quantify the text similarity.

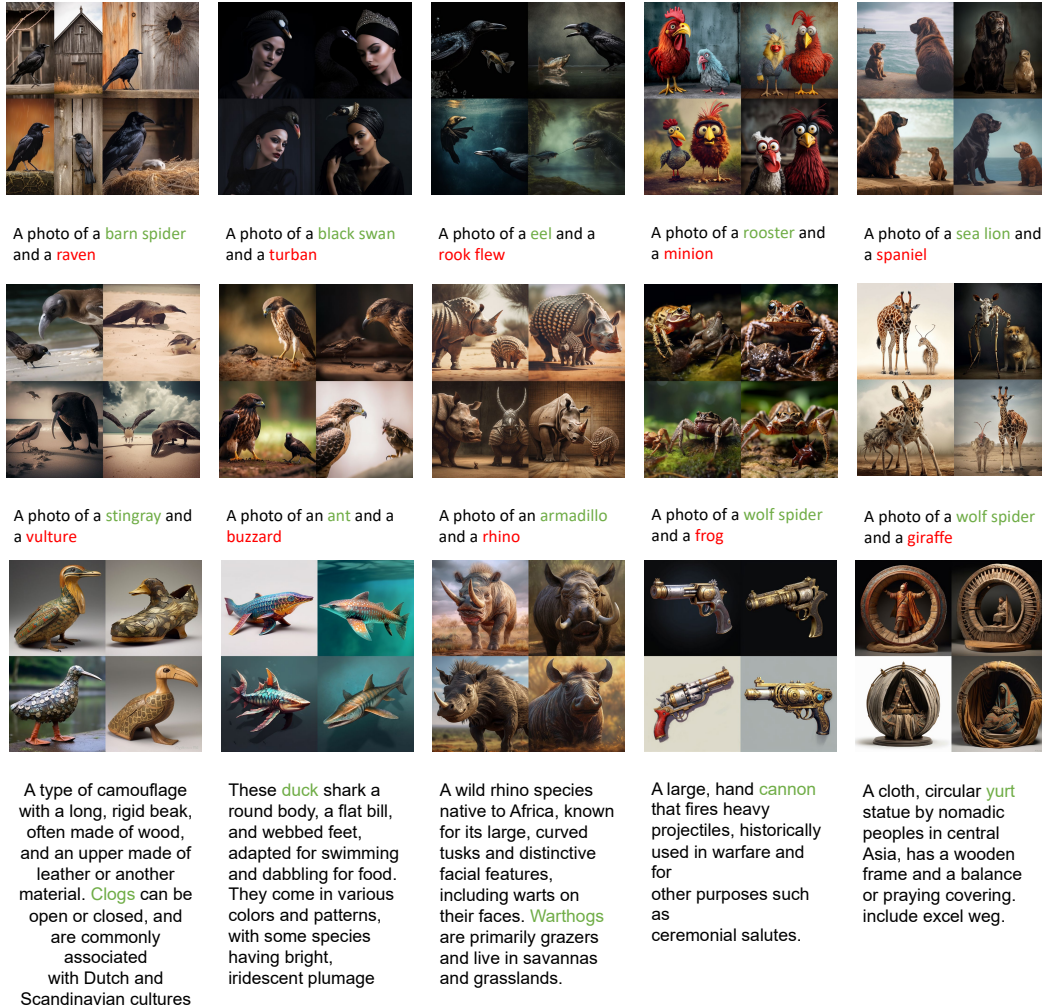


Figure C.3: Black-box attack on mid-journey (Section C.3 Black-box Attack).

138 D.2 Main Results

139 Table D.1 reports our main results, including short-prompt and long-prompt attacks. Compares to long
 140 text prompts, short text prompts comprise only a small number of tokens. This leads to a relatively
 141 fragile structure that is extremely vulnerable to slight disturbance. Therefore, random attacks can
 142 reach an impressive success rate of 79.2% targeting short prompts but a low success rate of 41.4%
 143 targeting the long prompts. In the contrast, our algorithm demonstrates its true potential, reaching an
 144 impressive success rate of 91.1% and 81.2% targeting short and long prompts, respectively.

145 As a further evidence of the effectiveness of our algorithm, it's worth noting the text similarity (TS)
 146 metrics between the random attacks and our algorithm's outputs. For short-prompt attack, the values
 147 stand at 0.69 and 0.72, respectively, illustrating that the semantic information of short texts, while
 148 easy to disrupt, can be manipulated by a well-designed algorithm with fluency and semantic similarity
 149 constraints. Our attacks preserve more similarity with the clean prompts. For long-prompt attacks,
 150 the TS score of random attacks (0.94) is higher compared to our attacks (0.84). One possible reason
 151 is that random attacks tend to make only minimal modifications as the length of the prompt increases.
 152 This limited modification can explain the significantly lower success rate of random attacks on longer
 153 prompts.

154 From the perspective of image generation quality and diversity, we found that as the attack success
 155 rate increases, image generation quality and diversity will decrease. For short and long texts, images
 156 generated from the clean text have the lowest FID (18.51 and 17.95) and the highest IS (101.33 ± 1.80)

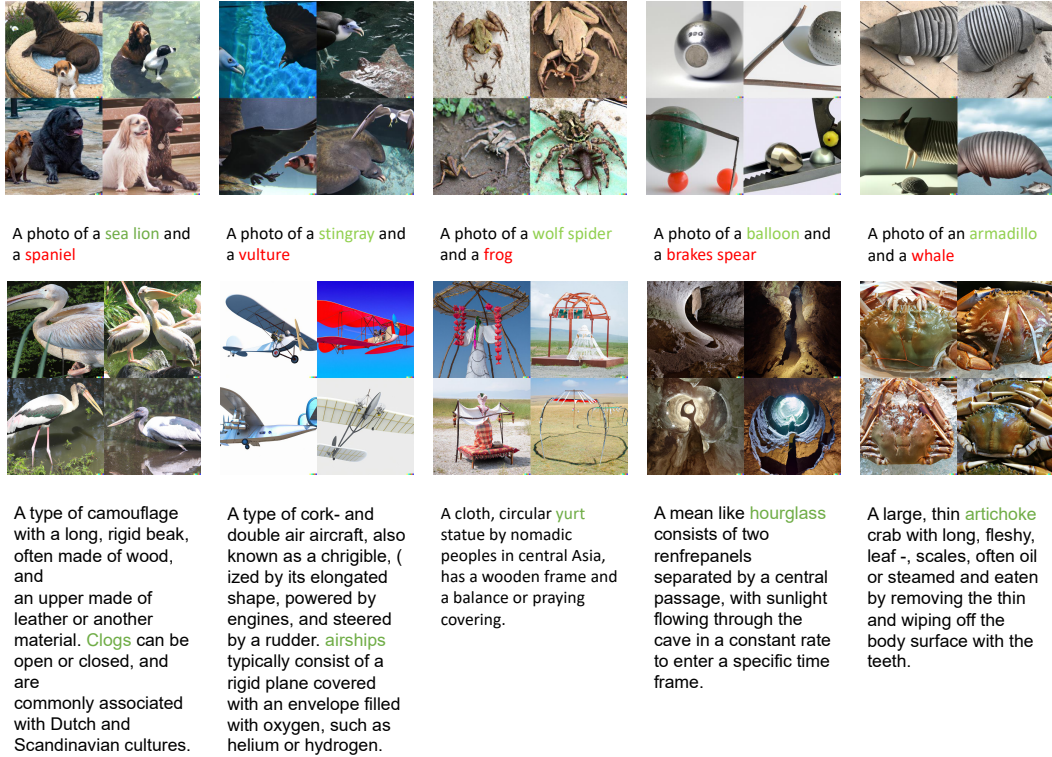


Figure C.4: Black-box attack on DALL·E2 (Section C.3 Black-box Attack).

Table D.2: Results of the ablation study on the number of steps in attack prompt search.

#Steps	Success (%)	FID (↓)	IS (↑)	TS (↑)
50	68.7	34.00	93.94±1.84	0.97
100	81.2	29.65	66.09±1.83	0.84
150	67.2	45.23	58.51±0.79	0.82

157 and 103.59±1.68). As the attack success rate rises, FID shows an upward trend. Examining this
 158 situation from the perspective of FID, a metric that gauges the distance between the distribution of
 159 generated images and the original data set. As the attack becomes more successful, the image set
 160 generated by the attack prompt tends to deviate substantially from the distribution of the original data
 161 set. This divergence consequently escalates the FID score, indicating a larger distance between the
 162 original and generated distributions. On the other hand, considering this situation from the diversity
 163 standpoint, it appears that the suppression of the generation of original categories brought on by the
 164 successful attack might instigate a decrease in diversity. This reduction in diversity, in turn, may
 165 cause a decrease in the Inception Score (IS).

166 D.3 Different Search Steps

167 Table D.2 presents the results of using different numbers of steps T in the search stage. For the
 168 $T = 50$ step configuration, the success rate is 68.7%. The FID value is 34.00, with lower values
 169 suggesting better image quality. The IS is reported as 93.94±1.84, with higher values indicating
 170 diverse and high-quality images. The TS value is 0.97, representing a high level of text similarity.
 171 Moving on to the $T = 100$ step configuration, the success rate increases to 81.2%, showing an
 172 improvement compared to the previous configuration. The FID value decreases to 29.65, indicating
 173 better image quality. The IS is reported as 66.09±1.83, showing a slight decrease compared to the
 174 previous configuration. The TS value is 0.84, suggesting a slight decrease in text similarity. In
 175 the $T = 150$ step configuration, the success rate decreases to 67.2%, slightly lower than the initial

Table D.3: Results of the ablation study on the constraints

Fluency	BERTScore	Success (%)	FID (\downarrow)	IS (\uparrow)	TS (\uparrow)
\times	\times	91.3	39.14	47.21 \pm 1.25	0.37
\checkmark	\times	81.7	29.37	64.93 \pm 1.57	0.79
\times	\checkmark	89.8	39.93	46.94 \pm 0.99	0.51
\checkmark	\checkmark	81.2	29.65	66.09 \pm 1.83	0.84

Table D.4: Results of the ablation study on the samplers

Sampler	Success (%)	FID (\downarrow)	IS (\uparrow)	TS (\uparrow)
DDIM [21]	81.2	29.65	66.09 \pm 1.83	0.84
DPM-Solver [9]	76.5	27.23	81.31 \pm 2.09	0.88

176 configuration. The FID value increases to 45.23, suggesting a decrease in image quality. The IS is
 177 reported as 58.51 \pm 0.79, indicating a decrease in the diversity and quality of generated images. The
 178 TS value remains relatively stable at 0.82.

179 When using $T = 50$, the attack prompt fails to fit well and exhibits a higher text similarity with the
 180 clean prompt. Although the generated images at this stage still maintain good quality and closely
 181 resemble those generated by the clean prompt, the success rate of the attack is very low. On the
 182 other hand, when $T = 150$, overfitting occurs, resulting in a decrease in text similarity and image
 183 quality due to the overfitted attack prompt. Consequently, the success rate of the attack also decreases.
 184 Overall, the configuration of $T = 100$ proves to be appropriate.

185 D.4 The Impact of Constraints

186 Table D.3 examines the impact of the fluency and semantic similarity (BERTScore) constraints. When
 187 no constraints are applied, the attack success rate is notably high at 91.3%. However, this absence of
 188 constraints results in a lower text similarity (TS) score of 0.37, indicating a decreased resemblance
 189 to clean text and a decrease in image quality. By introducing fluency constraints alone, the attack
 190 success rate decreases to 81.7% but increases the text similarity to 0.79. Furthermore, incorporating
 191 semantic similarity constraints independently also leads to a slight reduction in success rate to 89.8%,
 192 but only marginally improves the text similarity to 0.51. The introduction of constraints, particularly
 193 fluency constraints, leads to an increase in text similarity. The fluency constraint takes into account
 194 the preceding tokens of each token, enabling the integration of contextual information for better
 195 enhancement of text similarity. On the other hand, BERTScore considers a weighted sum, focusing
 196 more on the similarity between individual tokens without preserving the interrelation between context.
 197 In other words, the word order may undergo changes as a result and leads to a low text similarity.
 198 Certainly, this outcome was expected, as BERTScore itself prioritizes the semantic consistency
 199 between two prompts, while the order of context may not necessarily impact semantics. This further
 200 highlights the importance of employing both constraints simultaneously. When both constraints are
 201 utilized together, the text similarity is further enhanced to 0.84. Meanwhile, the success rate of the
 202 attack (81.2%) is comparable to that achieved when employing only the fluency constraint, while the
 203 text similarity surpasses that obtained through the independent usage of the two constraints.

204 D.5 Different Samplers

205 Table D.4 illustrates the effectiveness of our attack method in successfully targeting both DDIM
 206 and the stronger DPM-Solver. For the DDIM sampler, our attack method achieves a success rate
 207 of 81.2%, indicating its ability to generate successful attack prompts. Similarly, our attack method
 208 demonstrates promising results when applied to the DPM-Solver sampler. With a success rate of
 209 76.5%, it effectively generates attack prompts. The TS scores of 0.84 and 0.88, respectively, indicate
 210 a reasonable level of text similarity between the attack prompts and clean prompts. These outcomes
 211 demonstrate the transferability of our attack method, showcasing its effectiveness against both DDIM
 212 and the more potent DPM-Solver sampler.

213 References

- 214 [1] H. Chefer, Y. Alaluf, Y. Vinker, L. Wolf, and D. Cohen-Or. Attend-and-excite: Attention-based semantic
215 guidance for text-to-image diffusion models. *arXiv preprint arXiv:2301.13826*, 2023.
- 216 [2] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image
217 database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee,
218 2009.
- 219 [3] P. Esser, J. Chiu, P. Atighehchian, J. Granskog, and A. Germanidis. Structure and content-guided video
220 synthesis with diffusion models. *arXiv preprint arXiv:2302.03011*, 2023.
- 221 [4] W. Feng, X. He, T.-J. Fu, V. Jampani, A. Akula, P. Narayana, S. Basu, X. E. Wang, and W. Y. Wang.
222 Training-free structured diffusion guidance for compositional text-to-image synthesis. *arXiv preprint*
223 *arXiv:2212.05032*, 2022.
- 224 [5] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale
225 update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30,
226 2017.
- 227 [6] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information*
228 *Processing Systems*, 33:6840–6851, 2020.
- 229 [7] J. Ho, T. Salimans, A. Gritsenko, W. Chan, M. Norouzi, and D. J. Fleet. Video diffusion models. *arXiv*
230 *preprint arXiv:2204.03458*, 2022.
- 231 [8] L. Huang, D. Chen, Y. Liu, Y. Shen, D. Zhao, and J. Zhou. Composer: Creative and controllable image
232 synthesis with composable conditions. *arXiv preprint arXiv:2302.09778*, 2023.
- 233 [9] C. Lu, Y. Zhou, F. Bao, J. Chen, C. Li, and J. Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic
234 model sampling in around 10 steps. In *Advances in Neural Information Processing Systems*, 2022.
- 235 [10] G. Metzger, E. Richardson, O. Patashnik, R. Giryes, and D. Cohen-Or. Latent-nerf for shape-guided
236 generation of 3d shapes and textures. *arXiv preprint arXiv:2211.07600*, 2022.
- 237 [11] Midjourney. Midjourney. <https://www.midjourney.com/>, 2023. Accessed on 2023-05-17.
- 238 [12] C. Mou, X. Wang, L. Xie, J. Zhang, Z. Qi, Y. Shan, and X. Qie. T2i-adapter: Learning adapters to dig out
239 more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023.
- 240 [13] A. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen. Glide:
241 Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint*
242 *arXiv:2112.10741*, 2021.
- 243 [14] A. Obukhov, M. Seitzer, P.-W. Wu, S. Zhydenko, J. Kyl, and E. Y.-J. Lin. High-fidelity performance
244 metrics for generative models in pytorch, 2020. URL <https://github.com/toshas/torch-fidelity>.
245 Version: 0.3.0, DOI: 10.5281/zenodo.4957738.
- 246 [15] OpenAI. Dall-e 2. <https://openai.com/product/dall-e-2>, 2023. Accessed on 2023-05-17.
- 247 [16] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 248 [17] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin,
249 J. Clark, et al. Learning transferable visual models from natural language supervision. In *International*
250 *conference on machine learning*, pages 8748–8763. PMLR, 2021.
- 251 [18] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with
252 latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
253 *Recognition*, pages 10684–10695, 2022.
- 254 [19] L. Ruan, Y. Ma, H. Yang, H. He, B. Liu, J. Fu, N. J. Yuan, Q. Jin, and B. Guo. Mm-diffusion: Learning
255 multi-modal diffusion models for joint audio and video generation. *arXiv preprint arXiv:2212.09478*,
256 2022.
- 257 [20] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning using
258 nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265.
259 PMLR, 2015.
- 260 [21] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. In *International Conference on*
261 *Learning Representations*, 2021.

- 262 [22] Stability-AI. Stable diffusion public release. [https://stability.ai/blog/](https://stability.ai/blog/stable-diffusion-public-release)
263 [stable-diffusion-public-release](https://stability.ai/blog/stable-diffusion-public-release), 2023. Accessed on 2023-05-17.
- 264 [23] R. Tang, A. Pandey, Z. Jiang, G. Yang, K. Kumar, J. Lin, and F. Ture. What the daam: Interpreting stable
265 diffusion using cross attention. *arXiv preprint arXiv:2210.04885*, 2022.
- 266 [24] S. Yin, C. Wu, H. Yang, J. Wang, X. Wang, M. Ni, Z. Yang, L. Li, S. Liu, F. Yang, et al. Nuwa-xl: Diffusion
267 over diffusion for extremely long video generation. *arXiv preprint arXiv:2303.12346*, 2023.
- 268 [25] L. Zhang and M. Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint*
269 *arXiv:2302.05543*, 2023.
- 270 [26] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep
271 features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern*
272 *recognition*, pages 586–595, 2018.