# An analysis on "A Pilot Study of Query-Free Adversarial Attack against Stable Diffusion"

Naveen Mysore
nmysore.work@gmail.com

*Abstract*—This analysis examines "A Pilot Study of Query-Free Adversarial Attack against Stable Diffusion," [1] a pioneering investigation into the feasibility and efficacy of query-free adversarial attacks on Text-to-Image (T2I) models, with a focus on the Stable Diffusion [2]. My evaluation centers on the study's conceptual novelty, methodological rigor, clarity of result presentation, and the implications of its findings within the broader context of adversarial machine learning research. By scrutinizing the paper's hypothesis, assumptions, research design, attack methodologies, and the interpretation of results, I aim to assess the contribution of this study to the understanding of T2I model vulnerabilities and its significance in advancing the field's knowledge on adversarial robustness. Additionally, I evaluate the reproducibility of the research, ethical considerations, and the paper's suggestions for future inquiries. This analysis seeks to highlight the strengths and potential areas for improvement, providing a balanced perspective on the study's impact and directions for future research on securing machine learning models against adversarial threats.

*Index Terms*—Text-to-Image models, Diffusion models, Adversarial Attack

## I. INTRODUCTION

The rapid evolution of machine learning models, particularly in the realm of Text-to-Image (T2I) generation, has ushered in groundbreaking applications in art, design, and multimedia content creation. Among these, Stable Diffusion stands out as a prominent example, celebrated for its ability to generate highly detailed and contextually relevant images from natural language prompts. However, the growing capabilities of such models also unveil vulnerabilities to adversarial attacks, posing significant security and integrity challenges. "A Pilot Study of Query-Free Adversarial Attack against Stable Diffusion" delves into this critical issue, proposing a novel approach to generating adversarial attacks without the need for extensive model queries. This paper, therefore, represents a crucial step towards understanding and mitigating the susceptibilities of T2I models to such attacks.

My critical analysis of this pivotal study is structured to offer a comprehensive evaluation of its scientific and practical contributions to the field of adversarial machine learning. I begin by examining the background and significance of the research, placing it within the context of ongoing efforts to enhance the adversarial robustness of T2I models. I then assess the clarity and rigor of the problem statement, the theoretical underpinnings of the proposed query-free attack model, and the methodological framework employed to test the hypotheses. Furthermore, my analysis scrutinizes the results and their interpretation to determine the validity and reliability of the conclusions drawn by the authors. Through this in-depth review, I aim to provide constructive feedback that may inspire further research and development in securing T2I and other machine learning models against adversarial threats, thereby contributing to the advancement of the field.

## II. BACKGROUND AND SIGNIFICANCE

The significance of researching adversarial attacks within machine learning (ML) landscapes is underscored by the escalating dependence on ML models across various sectors, notably those with heightened security implications. Adversarial attacks denote the process of generating inputs aimed at misleading ML models into erroneous outputs. These inputs often appear normal to human observers, highlighting the nuanced challenge these attacks present. The necessity to comprehend these attacks and formulate countermeasures transcends mere technical hurdles, embodying a critical aspect of fortifying the integrity and dependability of ML frameworks. The acknowledgment of neural networks' susceptibility to adversarial manipulations dates back to their inception, with a landmark revelation by Szegedy et al. in 2014. They showcased that minimal, imperceptible alterations to an image could deceive deep neural networks (DNNs) into misclassification [3], catalyzing a burgeoning interest in this domain. Subsequent investigations have introduced various methodologies for executing attacks, such as the Fast Gradient Sign Method (FGSM) by Goodfellow et al., exploiting loss function gradients to craft adversarial samples [4]. Techniques including the Jacobian-based Saliency Map Attack (JSMA) [5], DeepFool [6], and the Carlini Wagner Attacks further exposed the extent of vulnerabilities in ML models [7]. Initially concentrated on image recognition, the scope of adversarial research expanded to encompass domains like natural language processing (NLP), speech recognition, and malware identification, illustrating the universal challenge posed by adversarial vulnerabilities.

In parallel, the exploration of diffusion models within ML signifies the merging of insights from stochastic processes, denoising autoencoders, and concepts like Brownian motion and statistical physics. These models have seen rapid evolution, particularly from the mid-2010s, marking a pivotal phase in generative modeling. Their unique data generation approach, involving iterative noise addition and removal, mirrors the natural diffusion process [8]. This methodology has proven effective in producing high-quality outputs across various modalities, such as images, audio, and text, positioning diffusion models as a critical element in generative

ML. The innovation in these models, exemplified by variants like Stable Diffusion, which enhances efficiency and output quality through latent space manipulation and natural language integration, underscores the dynamic progress in this field. This development not only signifies technical milestones but also marks significant advancements in generative modeling's capabilities. The figure 1 shows the class of various generative models and how duffsion models differes from other models.
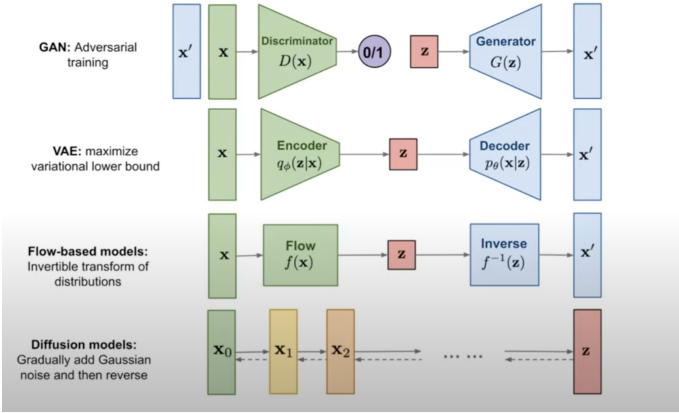


Fig. 1. Class of generative models.

As ML models find applications in vital areas like infrastructure, finance, autonomous navigation, and healthcare, safeguarding against adversarial exploits becomes a paramount concern. The potential repercussions of adversarial attacks accentuate the criticality of this research domain. Delving into adversarial attacks sheds light on ML models' operational limitations and decision-making processes, vital for enhancing model interpretability and transparency. This investigative endeavor is instrumental in evolving ML models to exhibit resilience against adversarial conditions, aligning with the overarching objective of developing AI systems capable of secure and efficacious operation in multifaceted real-world settings. Furthermore, the ethical implications of deploying AI in sensitive domains necessitate a steadfast commitment to ensuring AI systems' resistance to adversarial manipulations, integral to the broader mandate of deploying ethical and trustworthy AI.

## III. ANALYZING PROBLEM STATEMENT

The authors present a well-articulated problem statement, aiming to probe the adversarial robustness of Text-to-Image (T2I) models, with a specific focus on Stable Diffusion. It delves into the intriguing possibility of generating adversarial text prompts devoid of direct, end-to-end queries to the model, an approach termed as "query-free attack generation." This exploration is not only timely but also crucial, considering the expanding deployment of T2I models across various domains and the inherent risks posed by adversarial threats. The spotlight on query-free attacks is particularly noteworthy for its exploration of attack feasibility under constrained access to the

target model—a scenario frequently encountered in practical settings.

The paper's conceptual foundation is robust, strategically zeroing in on the vulnerabilities of T2I models to adversarial text prompts, particularly through the manipulation of text encoders like CLIP [9]. The introduction of a novel methodology to craft both untargeted and targeted query-free attacks by harnessing "steerable key dimensions" within the text embedding space is intellectually stimulating. This approach ingeniously shifts the focus of vulnerability from the image generation phase to the preliminary text encoding stage. Nevertheless, the paper's conceptual scaffold would gain further depth through a more elaborate discourse on the methodologies employed to identify these steerable key dimensions and their inherent susceptibility to adversarial distortions.

Implicit within the paper is a hypothesis suggesting that minor yet strategic perturbations to the text prompt—constrained to alterations of five characters—could profoundly influence the imagery output by Stable Diffusion. Additionally, it posits that such perturbations can be finetuned to impact specific content elements without encroaching upon unrelated areas. While the hypothesis stands on a solid foundational premise, articulating it more explicitly would enhance the structural coherence of the empirical validation process.

The paper anticipates that the delineated query-free attack strategies, including the PGD attack, genetic algorithm, and greedy search, will demonstrably lower the CLIP score of the resultant images. A CLIP score is a metric used to measure the semantic similarity between textual descriptions and corresponding images, based on embeddings generated by the CLIP model. This anticipated reduction in CLIP score is indicative of a weakened semantic linkage between the textual input and the corresponding image output. Moreover, it conjectures that these attack methodologies will not only significantly alter the image output with mere five-character prompt modifications but also enable precise manipulations to exclude or modify particular content segments. These forecasts are well-defined, measurable, and exhibit a strong alignment with the experimental setup delineated within the paper.

Overall, the paper presents a lucidly defined problem statement and an innovative conceptual framework poised to advance our understanding of the adversarial robustness of T2I models. While the hypothesis is inferred rather than explicitly stated, it is underpinned by a logical rationale that seamlessly transitions into tangible, testable predictions. The exploration of query-free attacks introduces a layer of novelty and applicability to the research. Nonetheless, the exposition could be further enriched by an in-depth discussion on the conceptual mechanisms behind steerable key dimensions and a direct articulation of the hypothesis, thereby fortifying the connection between the theoretical underpinnings and the empirical examination.

The authors make the following Assumptions in their attack models:

- Access to CLIP Text Encoder: The assumption of access to the CLIP text encoder is pivotal, given that the text

prompts are processed through this encoder to generate text embeddings that subsequently guide the image synthesis in Stable Diffusion. With the proliferation of open-source models, this assumption holds substantial validity.

- No Direct Access to Stable Diffusion: The premise that the attacker operates without direct access to the Stable Diffusion model, particularly bypassing the diffusion process integral to image generation, underscores the necessity for a query-free approach to the attack. Although recreating diffusion models using theoretical knowledge is feasible, the principal constraint here is attributed to computational costs rather than access.

- Perturbation of Text Prompts: The strategy that attackers can modify text prompts underscores the feasibility and potential effectiveness of the proposed adversarial approach. This foundational assumption is critical for exploring the spectrum of vulnerabilities within the text-to-image synthesis process.

### A. Problem with cosine similarity as the only metric

The authors define the attackers objective by minimizing the cosine similarity between the text embeddings of x and x'. where $\tau_\theta(x)$ denote the text encoder of CLIP with parameters $\theta$

$$\min_{x'} \cos(\tau_\theta(x), \tau_\theta(x'))  \qquad (1)$$

Further in evaluating the "Targeted Attack and Steerable Key Dimensions" section of the paper on query-free adversarial attacks against the Stable Diffusion model, it is clear that the authors present a nuanced approach to crafting adversarial text prompts. This method focuses on precise manipulation of the model's output, specifically targeting the alteration or removal of certain elements within the generated images while leaving the rest of the content largely unaffected. Unlike broader attack strategies that aim for general model disruption, this technique endeavors to achieve deliberate misrepresentations or modifications. The core of this approach lies in the identification and manipulation of "steerable key dimensions" within the text embedding space, indicating a sophisticated strategy for aligning the model's output with the attacker's specific objectives. This evaluation is represented as below. Here $I$ represents the dimension influence by majority vote and that is determined by finding text embedding difference $d_i = \tau_\theta(s) - \tau_\theta(s')$. More details are provided in the authors in their papers but the key analysis here is that even in the end the final evaluation tool cosine similarity metrics.

$$\min_{x'} \cos\left(\tau_\theta(x) \odot I, \tau_\theta(x') \odot I\right)  \qquad (2)$$

the use of cosine similarity metrics as the sole evaluation tool in the context of adversarial attack models, particularly within text-to-image generation processes, will be subjected to scrutiny. While cosine similarity, a metric measuring the similarity between two vectors by assessing their orientation rather than magnitude, has found broad applicability across machine learning, information retrieval, and natural language

processing, its utility in the adversarial domain is not without limitations.

- Magnitude Information Loss: One fundamental limitation of cosine similarity is its focus on vector orientation to the exclusion of magnitude. This characteristic poses a significant drawback in adversarial contexts, where the extent of input alteration (magnitude) is often as critical as the direction of change (orientation) for evaluating an attack's impact or a model's resilience. The neglect of magnitude information could result in an incomplete assessment of adversarial effectiveness or model robustness.

- Challenges in High-Dimensional Spaces: The application of cosine similarity in high-dimensional embedding spaces, typical of text and image data, faces diminishing discriminative power.

- Discrepancy with Human Perceptual Judgments Another critique concerns the metric's alignment with human perceptual judgments. In text-to-image generation tasks, the objective frequently encompasses modifying the perceptual output in human-discernible ways. Cosine similarity, as a purely mathematical measure, may not accurately reflect human perceptions of similarity or difference, leading to discrepancies between computational assessments and human evaluations.

- Dependence on Embedding Quality The reliability of cosine similarity as an evaluation metric is intrinsically linked to the quality of the underlying embeddings. Poorly representative embeddings that fail to capture essential nuances of the data can lead to misleading cosine similarity scores, misrepresenting the true similarity or disparity between adversarial and original inputs.

- Linear Assumptions Limitations The metric's assumption of linear relationships between vectors further constrains its applicability. The input-output transformations of models, especially in generating images from text, are inherently nonlinear. Cosine similarity's linear framework may inadequately capture the complexities and nuances of changes brought about by adversarial perturbations, particularly when subtle input modifications yield disproportionate or nonlinear effects in the output space.

### IV. Evaluation of Methodology

The process begins with generating sentences that encapsulate the attacker's specific targets, followed by the modification of these sentences to remove the targeted elements. This step is crucial for directing the attack but also introduces a significant dependency on the quality and relevance of the generated sentences. The effectiveness of this strategy hinges on the sentences' ability to capture the nuances of the targeted concept adequately. By comparing the embeddings of original $\tau_\theta(s)$ and perturbed sentences $\tau_\theta(s')$, the method identifies key dimensions that influence the generation process. This analytical step is both innovative and fraught with challenges, notably the computational intensity required to process potentially large embedding spaces and the nuanced interpretation

of high-dimensional data. Some of the limitations of these methodology are

- Precision and Efficiency: The method's focus on steerable key dimensions presents a nuanced way of attacking, allowing for precise content manipulation. However, the approach's efficiency is contingent upon the initial quality of sentence generation and the computational resources available for embedding analysis.
- Model-Specific Insights: Insights gained are inherently model-specific, potentially limiting the generalizability of the findings across different models or versions of the text-to-image synthesis process.
- Interpretability Challenges: High-dimensional embedding spaces inherently challenge interpretability, complicating the identification of truly influential dimensions and their semantic impact on the model's output.

The authors could explore Principal Component Analysis (PCA) and eigen vectors to distill high-dimensional embedding spaces into more interpretable, lower-dimensional representations. This approach can simplify the identification of influential directions that impact the generation process, potentially enhancing precision and efficiency by focusing on the dimensions that account for the most variance in the data. However, a drawback of relying on PCA and eigen vectors is that they primarily capture linear relationships, which may not fully represent the complex, non-linear interactions within the embedding space, thus potentially oversimplifying the model's dynamics and overlooking subtle but critical adversarial avenues.

## V. EVALUATION OF EXPERIMENTS

the authors meticulously designed experiments to assess the effectiveness of generating adversarial text prompts without necessitating end-to-end model queries. To achieve this, they employed a variety of attack methods, including the PGD attack, genetic algorithms, and greedy search techniques, each aimed at introducing subtle yet impactful perturbations to the text prompts. These perturbations were designed to manipulate the text-to-image generation process, with the goal of altering or controlling the output images in a targeted manner. The effectiveness of these attacks was evaluated using the CLIP score, which measures the semantic similarity between the text prompts and the generated images, thus providing a quantitative basis for assessing the impact of the adversarial interventions. Through this experimental setup, the authors aimed to shed light on the vulnerabilities inherent in the text-to-image generation process and to explore the potential for query-free adversarial attacks to compromise the integrity of generated images.

In my analysis, the experimental setup, implementation, and evaluation metrics presented in the study on query-free adversarial attacks against the Stable Diffusion model, while pioneering, reveal inherent limitations that warrant critical examination. The specificity of the model version (v1.4) used in the experiments raises concerns about the applicability of the findings to future versions of the model, which may evolve in

architecture, training paradigms, or defense mechanisms. The reliance on a fixed model configuration and a singular dataset for evaluating the attacks potentially overlooks a broader spectrum of adversarial vulnerabilities, thereby constraining the generalizability of the conclusions drawn.

The effectiveness of the adversarial strategies, highly contingent on parameter settings such as the learning rate for PGD and the mutation rate for genetic algorithms, underscores a limitation in the adaptability of the attack methodologies across different inputs and objectives. Moreover, the computational intensity of some attack methods, especially those hinging on evolutionary algorithms, poses a significant barrier to scalability and real-time application, a critical aspect for real-world adversarial scenarios.

Evaluation metrics, particularly the exclusive reliance on cosine similarity and CLIP scores, may not fully encapsulate the multifaceted impact of adversarial perturbations, especially in scenarios where the perceptual or semantic nuances play a crucial role. The lack of perceptual metrics to gauge the human interpretability of the generated images further limits the assessment of attack effectiveness, emphasizing a need for a more holistic approach to evaluating adversarial robustness that transcends singular metric reliance.

Despite these limitations, the experiment setup demonstrates notable strengths, including the simulation of realistic attack scenarios against a widely used text-to-image model and leveraging the advanced semantic understanding inherent in the CLIP model. The comprehensive approach, encompassing a variety of attack methods, enriches the analysis of model vulnerabilities, showcasing the depth of adversarial strategies explored. The inclusion of both quantitative metrics and qualitative assessments facilitates a thorough evaluation of the adversarial attacks' effectiveness, contributing significantly to the understanding of model vulnerabilities and advancing the discourse on the robustness of text-to-image models.

Overall, the experimental framework, while robust in its analysis and insightful in its findings, suggests avenues for further exploration. Future work could extend beyond the current limitations by incorporating a broader array of models, diversified datasets, and multifaceted evaluation metrics, including perceptual assessments, to forge a more comprehensive understanding of adversarial robustness in generative models. Additionally, the development and evaluation of defensive strategies in light of identified limitations emerge as crucial steps forward, ensuring the advancement of the field and the secure deployment of text-to-image models in real-world applications.

## VI. ANALYSIS OF RESULTS

In the analysis of results from the study on query-free adversarial attacks against the Stable Diffusion model, the effectiveness of these attacks in altering the output of text-to-image generation processes is highlighted. The experiments demonstrate a significant ability to disrupt the semantic alignment between input texts and generated images through minimal five-character perturbations, as evidenced by notable

reductions in CLIP scores. A key finding is the genetic algorithm's effectiveness in untargeted attacks, outperforming other methods in achieving pronounced CLIP score drops. This analysis also differentiates between the challenges of untargeted versus targeted attacks, with the latter showing smaller CLIP score decreases, indicating the increased difficulty in generating precise perturbations.

The paper delves into case studies to further elucidate the impact of these perturbations, such as altering an image's thematic content significantly. These instances reveal the intricate mechanism of perturbation, where subtle character adjustments lead to substantial content deviations from the intended output, showcasing the query-free attacks' capability to exploit the model's sensitivities to input variations. The results underscore the robust and versatile nature of these adversarial strategies, underscoring their importance in understanding and potentially mitigating vulnerabilities within text-to-image generative models.

However, the analysis raises critical considerations regarding the generalizability of these findings, the diversity of the dataset, the complexity and computational cost of generating adversarial attacks, the reliance on singular metrics like CLIP scores for evaluating attack success, and the lack of exploration into the model's defenses. These points suggest areas for further research to deepen our understanding of adversarial robustness and to develop more comprehensive defenses against such attacks, thereby enhancing the resilience of text-to-image generative models against adversarial threats.

## VII. CONCLUSION

The paper "A Pilot Study of Query-Free Adversarial Attack against Stable Diffusion" offers an in-depth examination of the novel methodology introduced for conducting query-free adversarial attacks on Text-to-Image (T2I) models, specifically focusing on the vulnerabilities of Stable Diffusion. This analysis critically explores various dimensions of the study, including its innovative approach, methodological rigor, clarity in presenting results, and the broader implications of these findings within the adversarial machine learning landscape.

The paper makes a significant contribution by highlighting the feasibility of executing adversarial attacks without extensive model querying. This approach not only broadens the understanding of adversarial tactics but also signals a shift in addressing model vulnerabilities under constrained access conditions. Such insights are vital for the ongoing efforts to fortify T2I models against potential adversarial threats, thereby enhancing the security and integrity of machine learning applications across diverse domains.

The meticulous evaluation of the proposed query-free attack mechanisms, including the PGD attack, genetic algorithm, and greedy search techniques, underscores the nuanced understanding of model vulnerabilities to textual perturbations. The reliance on the CLIP score as a primary metric for assessing the impact of adversarial interventions, while effective, also invites a deeper reflection on the need for multi-dimensional

evaluation criteria that better capture the complex interplay between textual inputs and image outputs in T2I models.

The paper's exploration of "steerable key dimensions" within text embeddings introduces a novel perspective on manipulating model outputs. However, the methodological and computational challenges associated with this approach, such as the interpretability of high-dimensional embedding spaces and the scalability of attack methodologies, highlight areas for further refinement and exploration.

The findings from this analysis suggest several avenues for future research, including the exploration of alternative metrics for evaluating adversarial effectiveness, the development of more robust defensive mechanisms, and the investigation of the generalizability of these attack methodologies across different models and application contexts. The study's focus on query-free attacks enriches the adversarial machine learning discourse, setting a foundation for subsequent inquiries into the resilience of T2I models.

Furthermore, the ethical considerations surrounding the deployment of T2I models and their susceptibility to adversarial manipulations underscore the imperative for a comprehensive framework that not only addresses technical vulnerabilities but also considers the broader societal impacts of these technologies. The advancement in securing T2I models against adversarial threats is not merely a technical endeavor but a necessary step towards ensuring the ethical and responsible application of AI. This work not only advances scientific knowledge on adversarial robustness but also catalyzes further research aimed at developing secure, reliable, and ethically grounded AI systems.

## REFERENCES

[1] H. Zhuang, Y. Zhang, and S. Liu, "A pilot study of query-free adversarial attack against stable diffusion," 2023.
[2] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," 2022.
[3] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," 2014.
[4] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," 2015.
[5] R. Wiyatno and A. Xu, "Maximal jacobian-based saliency map attack," 2018.
[6] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, "Deepfool: a simple and accurate method to fool deep neural networks," 2016.
[7] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," 2017.
[8] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," 2020.
[9] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," 2021.