
Exploring Semantic-constrained Adversarial Example with Instruction Uncertainty Reduction

Jin Hu^{1,2} Jiakai Wang^{2✉} Linna Jing¹ Haolin Li³ Haodong Liu³
Haotong Qin⁴ Aishan Liu¹ Ke Xu^{1,2} Xianglong Liu^{1,2}

¹State Key Laboratory of Complex & Critical Software Environment (CCSE), Beihang University

²Zhongguancun Laboratory ³School of Computer Science and Engineering, Beihang University

⁴Department of Information Technology and Electrical Engineering, ETH Zurich
hujin@buaa.edu.cn wangjk@mail.zgclab.edu.cn

Abstract

Recently, semantically constrained adversarial examples (SemanticAE), which are directly generated from natural language instructions, have become a promising avenue for future research due to their flexible attacking forms, but have not been thoroughly explored yet. To generate SemanticAEs, current methods fall short of satisfactory attacking ability as the key underlying factors of semantic uncertainty in human instructions, such as *referring diversity*, *descriptive incompleteness*, and *boundary ambiguity*, have not been fully investigated. To tackle the issues, this paper develops a multi-dimensional **instruction uncertainty reduction (InsUR)** framework to generate more satisfactory SemanticAE, *i.e.*, transferable, adaptive, and effective. Specifically, in the dimension of the sampling method, we propose the residual-driven attacking direction stabilization to alleviate the unstable adversarial optimization caused by the diversity of language references. By coarsely predicting the language-guided sampling process, the optimization process will be stabilized by the designed ResAdv-DDIM sampler, therefore releasing the transferable and robust adversarial capability of multi-step diffusion models. In task modeling, we propose the context-encoded attacking scenario constraint to supplement the missing knowledge from incomplete human instructions. Guidance masking and renderer integration are proposed to regulate the constraints of 2D/3D SemanticAE, activating stronger scenario-adapted attacks. Moreover, in the dimension of generator evaluation, we propose the semantic-abstacted attacking evaluation enhancement by clarifying the evaluation boundary based on the label taxonomy, facilitating the development of more effective SemanticAE generators. Extensive experiments demonstrate the superiority of the transfer attack performance of InsUR. Besides, it is worth highlighting that we realize the reference-free generation of semantically constrained 3D adversarial examples by utilizing language-guided 3D generation models for the first time.

1 Introduction

Adversarial example (AE), showing that small perturbations can impact the performance of deep learning models, is broadly focused due to its potential to promote model robustness and secure applications in practice. A series of studies has uncovered several forms of AEs, including physical-world AEs [1, 2], transfer AEs [3, 4], and naturalistic AEs [5, 6], as well as the applications in evaluating autonomous driving [7, 8, 9] or LLM systems [10, 11].

While most adversarial example research focuses on finding AEs around existing data, generating AEs from natural language instructions without referenced data has not yet been thoroughly explored, *i.e.*, to find *Semantic-Constrained Adversarial Examples* (SemanticAE). Specifically, given a certain

natural language description, we aim to generate the data that corresponds to its real semantic meaning but is hardly to be correctly recognized by deep learning models trained in related tasks. Recent works have employed techniques related to naturalistic AEs to accomplish a similar objective [12, 13, 14], , but the de facto potential of SemanticAE has still not been fully released in performing transferable, adaptive, and effective attacks, limiting the applicability. In light of the recent advancements in language-driven multimodal intelligence and the increasing demand for alignment [15, 16, 17], it is necessary to take a step further in SemanticAE generation and facilitate more versatile AE generation.

To push the boundary of the current technology, we focus on the key underlying factor limiting the adversarial capability of SemanticAEs: the inherent uncertainty within human instructions that defines semantic constraints. We categorize three major forms of uncertainty in instructions: ① *Referring diversity* introduces a barrier in SemanticAE optimization via the multi-step generative models, since it leads to the inconsistent language-guidance that the adversarial optimization should collaborate with. ② *Descriptive incompleteness*, which conceptualizes the gap between the precise model of the attack scenario and the instructions given by potential users, restricts the application scenarios. ③ *Boundary ambiguity* of the semantic constraint is hard to characterize in task definitions, affecting the evaluation of SemanticAE generators.

We propose a multidimensional **instruction uncertainty reduction (InSUR)** framework to tackle the issues and generate more transferable, adaptive, and effective SemanticAE. Specifically, for referring diversity, we propose residual-driven attacking direction stabilization via the novel ResAdv-DDIM sampler that stabilizes optimization through coarsely predicting the language-guided sampling process, releasing the capability of multistep diffusion models on adversarial transferability and robustness. For descriptive incompleteness, we propose the context-encoded attacking scenario constraint for both 2D and 3D generation problems by scenario knowledge integration, tackling the scenario adaptation problem by addressing the descriptions' incompleteness problem, achieving the first 3D SemanticAE generation. For boundary ambiguity, we propose the semantic-abstacted attacking evaluation enhancement based on label taxonomy. Our contribution can be summarized as:

- We conceptualize the SemanticAE generation problem and propose a multi-dimensional instruction uncertainty reduction framework, InSUR, to address the challenges.
- In the dimension of the sampling method, we propose the residual-driven attacking direction stabilization to achieve better adversarial optimization. In task modeling, we propose the context-encoded attacking scenario constraint to realize scenario-adapted attacks. In generator evaluation, we propose the semantic-abstacted attacking evaluation enhancement to facilitate the development of SemanticAE generators.
- Extensive experiments demonstrate the superiority in the transfer attack performance of generated 2D SemanticAEs, and for the first time, we realize the reference-free generation of 3D SemanticAE by utilizing language-guided 3D generation models.

2 Backgrounds

Adversarial Example Generation Adversarial attack generating algorithms can be categorized as iterative optimization in the data space, e.g. *FGSM* [18], *PGD* [19], *AutoAttack* [20], iterative optimization in the latent space of generative models, e.g., *DiffPGD* [6], *NAP* [21], and *AdvFlow* [22], and training a neural network for adversarial attack generation [23, 24, 25]. Adversarial examples may not be robust in the physical world. For physical attacks, expectation-over-transformation (EoT) [26] and 3D simulation [27, 28] are proposed to bridge the digital-physical gap. Research has also focused on attacking unknown models, *i.e.*, transfer attacks [3, 4]. An extended technical background and related works are provided in Appendix A.

Semantic-constrained Adversarial Example [29] first proposes the *Unrestricted Adversarial Example*(UAE), of which the restriction is defined by the human's cognition instead of l_p -norm on existing data, and proposes a generative learning method for its generation. A line of studies further develops optimization techniques in latent spaces [30, 21, 31], and terms it as *Natural Adversarial Example* (NAE), while another line of study focuses on constructing the perceptual constraints for UAEs. A difference between them is that NAE studies also focus on generating

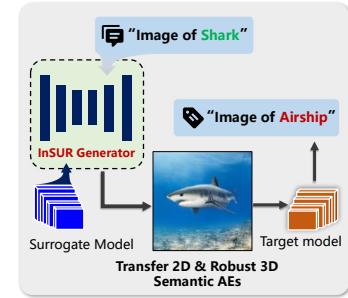


Figure 1: SemanticAEs are generated directly by instructions.

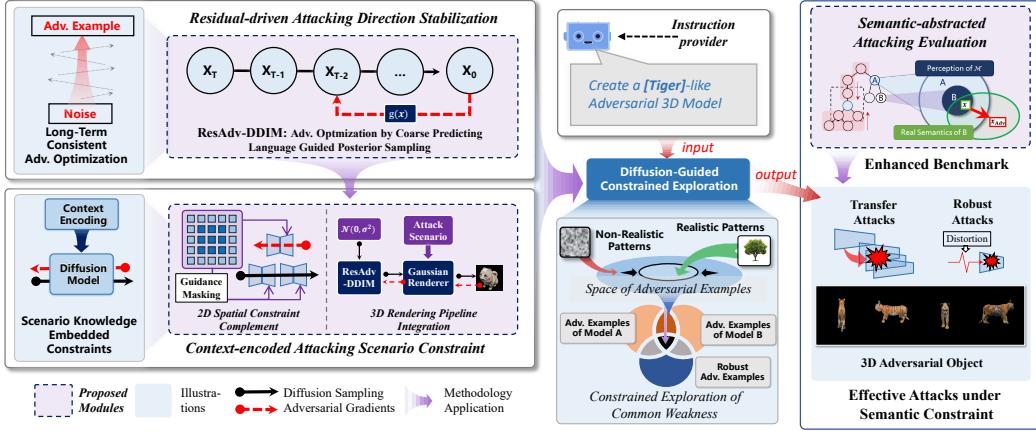


Figure 2: Overview of multi-dimensional instruction uncertainty reduction (**InSUR**) framework.

diverse-distributed adversarial examples without referencing a static image. We formulate the adversarial example generating task constrained by natural language’s semantics as the *semantically-constrained adversarial example generation* problem. Recent works([6, 31, 32, 13, 14]) focus on integrating the pre-trained diffusion model and iterative optimization to constrain the naturalness and improve transferability. Furthermore, generating 3D adversarial examples that are more aligned with the physical world and satisfy the semantic constraints is still an open problem.

3 Methodology

3.1 Problem Formulation and Analysis

Semantic-Constrained Adversarial Example (SemanticAE) Generation Problem We define SemanticAE generation problem as generating an adversarial example x_{adv} that fools the target model and satisfies the semantics constraint defined by the user’s instruction Text. Formally, we formulate the SemanticAE generation problem as follows:

$$\text{find } x_{\text{adv}} \in \mathcal{S}(\text{Text}) \text{ s.t. } \mathcal{M}(x_{\text{adv}}) \in A_{\text{Text}}, \quad (1)$$

where $\mathcal{S}(\text{Text})$ is the set of data with semantic meaning corresponding to Text, \mathcal{M} represents the target model, and A_{Text} defines the types of target model’s output that are conceptually different from Text, representing a successful attack. In the strict black-box setting, which is the focus of this paper, both \mathcal{S} and \mathcal{M} are unknown to the generation algorithm.

The goals of SemanticAE generation are *to build a red-team model \mathcal{G} that automatically finds the alignment problem between the intelligent model \mathcal{M} and the implicit semantics $\mathcal{S}(\text{Text})$ reflecting the social consensus or the physical world*. This goal leads to the following constraints: firstly, to achieve automatic alignment with limited supervision, the instruction Text is not required to characterize semantic constraints precisely. Secondly, from the perspective of data value, the generated SemanticAE x_{adv} should be able to perform transfer attacks.

Challenges in SemanticAE Generation As shown in the middle card of Figure 2, generative or diffusion models can constrain the pattern of generated AEs and facilitate transfer attacks with better in-manifold constraint [33]. We take a step further in SemanticAE, focusing on the inherent challenge related to instruction uncertainty: ① Reference diversity challenges adversarial optimization. The language guidance that is learned from the mapping between Text and $\mathcal{S}(\text{Text})$ is non-linear since $\mathcal{S}(\text{Text})$ is diverse. This makes collaborating with adversarial optimization and the diffusion model for better transfer attacks and robust attacks a non-trivial problem. ② Descriptive incompleteness requires scenario-knowledge integration for scenario-adapted generation. The challenges are identifying the missing contexts in pretrained models and establishing practical knowledge embedding methodologies, thereby further eliminating the reference diversity from the external perspective. ③ Boundary ambiguity makes defining \mathcal{S} and A for evaluating the generator also challenging, which lies in the fact that inappropriate evaluation leads to inaccurate results.

Multi-dimensional Instruction Uncertainty Reduction (InSUR) Framework As shown in Figure 2, for the reference diversity problem, we propose the residual-driven attacking direction

stabilization with the designed ResAdv-DDIM sampler. For the contextual incompleteness, we propose the context-encoded attacking scenario constraint methods for scenario-knowledge integration in representative 2D and 3D SemanticAE generation tasks. Moreover, since the unclear semantic boundary makes evaluating the generator difficult, semantic-abstracted attacking evaluation enhancement is proposed to facilitate further developments of SemanticAE generation.

3.2 Residual-driven Attacking Direction Stabilization with ResAdv-DDIM

Semantic-constrained Optimization Problem Referring to the generative-model-based adversarial examples, we solve SemanticAE generation by maximizing the loss \mathcal{L}_{ATK} under the constraint of the posterior sampling process defined by the natural language guidance Text. However, if the posterior sampling process is more complex, *e.g.*, multi-step diffusion de-noising, tackling this maximization problem is challenging. Recent work utilizes a deterministic sampling process, *e.g.*, DDIM [34], as a constraint defined by Text [6, 12]. For simplicity, we denote the sampling step as $f_{\theta, \Delta T}(x_t) \sim q_\theta(x_{t-\Delta T} | x_t, x_0, \text{Text})$, and such optimization can be formulated as:

$$\max \mathcal{L}_{\text{ATK}}(\underbrace{\mathcal{M}(f_{\theta, \Delta T} \circ f_{\theta, \Delta t} \circ \dots \circ f_{\theta, \Delta t}(x_T))}_{T/\Delta T \text{ times}}), \quad (2)$$

where \circ denotes function composition. However, this may trigger the robust problem of f , since it is hard to determine whether x_0 is an adversarial example of f or \mathcal{M} . This results in the instruction misalignment of SD-NAE shown in section 4.4. Also, this optimization is computationally expensive. Another solution is to tackle the challenges by approximating the gradient [6] or directly altering the sampling process [13], which can be re-formulated as:

$$x'_{t-\Delta T} = f_{\theta, \Delta T}(\arg \max_{x'_t} \mathcal{L}_{\text{ATK}}(\mathcal{M}(x'_t))), \quad \forall t \in [\Delta T, t_s], \quad (3)$$

where t_s is a selected intermediate step, and the \max_{x_t} optimization could be a single iteration. An advantage is that, since the maximization algorithm does not retrieve the information of f , or f has been protected from adversarial attacks, it alleviates the robustness problem of f . However, the optimization used $\nabla_{x_0} \mathcal{L}_{\text{ATK}}$ to approximate $\nabla_{x_t} \mathcal{L}_{\text{ATK}}$, and the inconsistency introduces noise to the optimization. As shown in Figure 3, the optimization direction may vary, or be non-linear, with respect to different x_t . This misalignment makes the adversarial pattern optimized in the initial denoising stage ineffective in the latter stage, thereby limiting the multi-step regularization opportunity of diffusion models for better transfer attacks.

Overall, this technical challenge originates from the conflict between (1) the accurate estimation of $\nabla_{x_t} \mathcal{L}_{\text{ATK}}$, causing the robust problems of the language guidance defined by f and the computational problems, and (2) the approximated estimation of $\nabla_{x_t} \mathcal{L}_{\text{ATK}}$, causing the non-optimality of attack optimization. We solve the problem by improving the approximation with the novel *ResAdv-DDIM* posterior sampler, enabling the discovery of more robust adversarial patterns with better regularization.

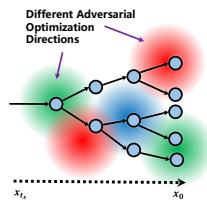


Figure 3:
Inconsistent adv.
direction problem.

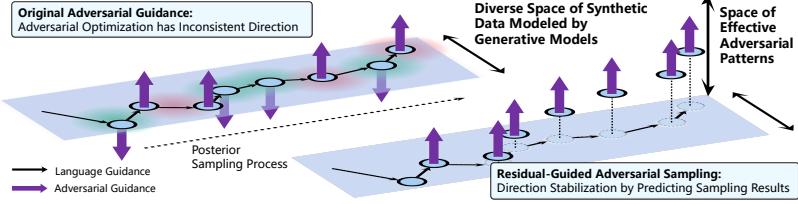


Figure 4: Residual-driven attacking direction stabilization. *ResAdv-DDIM* is designed for the efficient and thorough exploration of new adversarial patterns constrained by multi-step sampling processes.

Residual-Guided Adversarial DDIM Sampler Inspired by *Learning to Optimize* [35], we handle the challenge by **predicting a coarse sketch of the future-step denoising result x_0** for estimating the attack optimization direction with \mathcal{L}_{ATK} . Our key insight is that since multi-modal models acquire general capabilities through training in the task of *predicting diverse human responses* and *complementing missing information across diverse data*, we should fully leverage the model’s intrinsic multi-granularity predictive capabilities to achieve stable generation under semantic uncertainty.

Specifically, we leverage DDIM’s multi-step posterior sampling capabilities to achieve a coarse prediction of x_0 from current x_t , *i.e.*, $\hat{x}_0 = g(x_t)$, which allows for a more accurate estimate of $\nabla_{x_t} \mathcal{L}_{\text{ATK}}$ by using $\nabla_{\hat{x}_0} \mathcal{L}_{\text{ATK}}$ compared to $\nabla_{x_0} \mathcal{L}_{\text{ATK}}$. We formulate the generation process as:

$$g_\theta(x_t) = \underbrace{f_{\theta, \Delta T_1} \circ f_{\theta, \Delta T_2} \circ \cdots \circ f_{\theta, \Delta T_k}}_{k \text{ times}, k \ll T/\Delta T}(x_t), \text{ where } \sum_{i=1}^k \Delta T_i = t$$

$$x_{t-\Delta T} = f_{\theta, \Delta T}(\arg \max_{x_t} \mathcal{L}_{\text{ATK}}(\mathcal{M}(g_\theta(x_t)))), \forall t \in [\Delta T, t_s]. \quad (4)$$

The notation is the same as Eq 2. $g_\theta(x_t)$ is the coarse estimation of x_0 , and k is a small number of iterations that could be selected from $\{1, 2, 3, 4\}$. Since the sampling process takes a residual shortcut to x_0 , we name it as *Residual-Guided Adversarial DDIM Sampler (ResAdv-DDIM)*.

To further establish the concrete adversarial attack algorithm, we further propose the following method. **(1) Constraining the Semantics.** To ensure the generated sample satisfies $x_{\text{adv}} \in \mathcal{S}(\text{Text})$, we constrain the discrepancy of the sample trajectory with the l_2 -norm between DDIM-generated samples and the adversarially optimized samples after determining x_{t_s} :

$$\| \text{Denoise}_{\text{DDIM}}(x_{t_s-\Delta T}) - \text{Denoise}_{\text{Adv}}(x_{t_s-\Delta T}) \|_2 < \epsilon. \quad (5)$$

Such constrained optimization problem could be performed by simultaneously sampling with Denoise_{DDIM} and Denoise_{Adv}, and clipping x'_t in each step. **(2) Adaptive Attack Optimization.** To solve the maximization in Eq. 4, we introduce an early-stop mechanism that terminates optimization:

- (1) after the first iteration, and the estimated probability of unsuccessful attack $< \xi_1$, or,
- (2) at the first iteration, when this probability satisfies a stricter threshold $\xi_2 < \xi_1$.

The probabilities are estimated from $\mathcal{M}(g_\theta(x_t))$. These conditions reduce the expected number and lower bound of optimization steps, while maintaining attack performance alongside the denoising step, which degrades the adversarial capability. Lower threshold could result in better attack performance and slower optimization, and we set $\xi_1 = 0.1$, $\xi_2 = 0.01$ as a feasible setting across all experiments. In addition, we integrate momentum optimization, introduced in [14, 3], to improve the attack transferability. Detailed implementation and analysis are shown in Appendix B.1.

3.3 Context-encoded Attacking Scenario Constraint for 2D and 3D Generation

In the application scenarios, the instruction Text might be ambiguous or incomplete, which requires integrating learned guidance with external knowledge. For effective task adaptation, we provide knowledge embedding strategies on the key data structures that collaborate with the ResAdv-DDIM sampler, achieving better 2D SemanticAE generation and realizing 3D SemanticAE generation.

Spatial Constraint Complement for 2D SemanticAE Generation We focus on the problem of the incompleteness of the spatial constraint given by the instruction Text representing the object label. Specifically, an effective SemanticAE generator shall leverage the optimization space of the image backgrounds and generate patterns that amplify the attack’s effectiveness. However, the diffusion model’s conditional background generation is overly uniform because the attack functionalities were not considered during the original training. To solve the problem, we encode the context of the attack application through the fine-grained control of the denoising guidance. We leverage the application of guidance-masking from edit area control of single guidance [36] to the re-distribution of multiple guidance, and embed the masking into the key guidance function f_θ applied in both posterior sampling and adversarial optimization process in ResAdv-DDIM (as shown in Figure 2). Together with the deterministic DDIM, f_θ is formulated as:

$$f_{\theta, \Delta T}(x_t) = \sqrt{\bar{\alpha}_{t-\Delta T}/\bar{\alpha}_t} (x_t - \sqrt{1-\bar{\alpha}_t} \epsilon_\theta(x_t, t)) + \sqrt{1-\bar{\alpha}_{t-\Delta T}} \cdot \epsilon_\theta(x_t, t),$$

$$\epsilon_\theta(x_t, t) = (1-M) \cdot \epsilon_{\theta, \text{Unconditional}}(x_t, t) + M \cdot \epsilon_{\theta, \text{Conditional}}(x_t, t, \text{Text}), \quad (6)$$

where α defines the noise ratio in the diffusion model, ϵ_θ is the noise estimating network, and M is the guidance masking that regularizes the spatial distribution of semantic guidance Text. Detailed definition of M is in Appendix B.1, and the integration is illustrated in Figure 2.

Differentiable Rendering Pipeline Integration for 3D SemanticAE Generation 3D Data is valuable for world modeling [37, 38]. We focus on the problem of generating 3D SemanticAE $x_{\text{adv}}^{(3D)}$ for target models \mathcal{M} operated under the 2D inputs. To fill the gap between the 3D scenarios 2D target models, additional physical rendering knowledge shall be efficiently encoded. Leveraging the proposed ResAdv-DDIM sampler and the advancements in Gaussian-splatting [38], we fill the gap between the 3D generation and 2D target models. As shown in Figure 5, we optimize latents with the

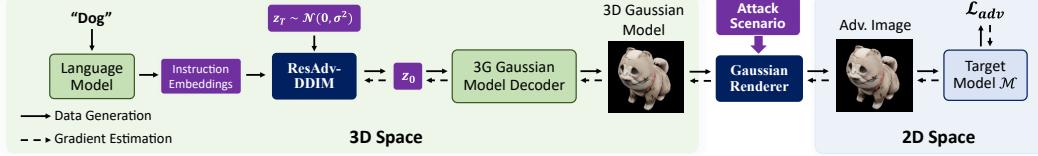


Figure 5: 3D optimization pipeline. Scenario knowledge is encoded with the Gaussian renderer.

gradient back-propagated through the 3D data structure and integrate the scenario knowledge during the differentiable rendering process. The optimization pipeline could be implemented concisely with the *Trellis* [39] framework, which is a recently published diffusion-based 3D access generation framework. For simplicity, we reformulate the *Trellis* sampling and rendering process as:

$$\begin{aligned} \mathbf{pos} &= \text{Coords}(\mathcal{D}_{slat}(z_0^{slat})), \quad \text{Model}_{GS} = \mathcal{D}_{GS}(z_0, \mathbf{pos}), \quad x = \text{Renderer}_{GS}(\text{Model}_{GS}, \text{Camera}), \\ \mathcal{D}_{GS} : \{(z^i, pos^i)\}_{i=1}^L &\rightarrow \{\{(x_i^k, c_i^k, s_i^k, \alpha_i^k, r_i^k)\}_{k=1}^K\}_{i=1}^L \end{aligned} \quad (7)$$

where z_0 and z_0^{slat} are latents sampled by the diffusion model, and are represented by sparse and dense tensors, respectively. \mathcal{D}_{slat} is the coarse structure decoder, Coords transforms the voxel to point positions \mathbf{pos} , \mathcal{D}_{GS} is the refined structure decoder that decodes each vertex into multiple Gaussian points and Renderer_{GS} renders the Gaussian model to 2D images x with the camera parameter. For SemanticAE generation, the refined feature generation process $\{z_T, z_{T-\Delta T}, \dots, z_0\}$ is replaced with ResAdv-DDIM sampling in Eq. 4 with the rendering model embedded in g_θ :

$$\begin{aligned} g_\theta(z_t, \mathbf{pos}, \text{Camera}) &:= \text{Render}_{GS}(\mathcal{D}_{GS}(f_{\theta, \Delta T_1} \circ \dots \circ f_{\theta, \Delta T_k}(z_t, \mathbf{pos}), \mathbf{pos}), \text{Camera}) \\ z_{t-\Delta T} &:= f_{\theta, \Delta T}(\arg \max_{z_t} \mathbb{E}_{\text{Camera} \sim P_{\text{Cam}}} [\mathcal{L}_{\text{ATK}}(\mathcal{M}(g_\theta(z_t, \text{Camera}, \mathbf{pos})))]) \end{aligned} \quad (8)$$

We use the EoT method with gradient accumulation to optimize z_t for unknown camera positioning, *i.e.*, samples Camera from P_{Cam} in each iteration. The scenario knowledge is encoded in P_{Cam} and the rendering background. Due to the stabilized guidance in ResAdv-DDIM, the gradient of the previous steps could be utilized as current-step gradient estimation, and therefore, fewer EoT steps are required. Since texture and localized positioning perturbation are enough for adversarial attacks, z_0^{slat} is not included in the parameter space and serves as a semantic anchor. With the collaborative constraint of 3D diffusion and renderer model, semantically-constrained and multi-view adapted SemanticAEs are efficiently generated. Detailed implementation is shown in Appendix B.2.

3.4 Semantic-abstracted Attacking Evaluation Enhancement with Label Taxonomy

The evaluation of SemanticAE generator requires the benchmark to judge whether $x \in \mathcal{S}(\text{Text})$ and define A_{Text} , which determines the adversarial attack and semantic alignment performance of the generator, and is still a blank in practice. To address the issue, we provide a task construction method for automatic evaluation based on the application goal of the SemanticAE generation task. Note that our method evaluates the SemanticAE generator instead of the adversarial example.

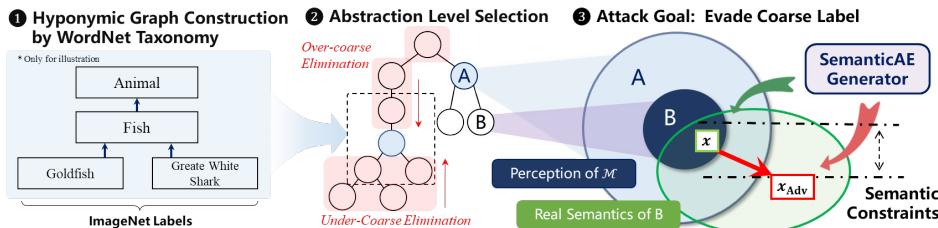


Figure 6: Construction of abstract label evasion evaluation task.

The task is constructed based on semantic abstracting with the label taxonomy. Firstly, the attack targets of existing non-target evaluation methods based on ImageNet labels are often too simple, while the constraint space of SemanticAE is relatively loose, making it easy for the attack generation model \mathcal{G} to achieve successful attacks easily. For example, it is unreasonable to use the ImageNet label “tiger-shark” as the misclassification category A_{Text} for the instruction Text “great-white-shark”, since achieving successful attacks in this task may not show the capability of successful attacks in real scenarios. To clarify the boundary, we re-construct the evaluation label with a better abstraction level by leveraging *WordNet* [40] taxonomy. As shown in Figure 6, we firstly construct the hyponymic

graph based on the hyponymic relation defined by *WordNet*, then select the proper abstraction level, and finally define the attack goal as the evasion attack on the abstracted label. Specifically, under the definition of SemanticAE generation, the evaluation task is formulated as:

$$\begin{aligned} \text{Text} &:= \text{"Realistic image of [AbstractedLabel], specifically, [label]"}, \\ A_{\text{Text}} &:= \{\text{label}_{\text{Adv}} \mid \text{AbstractedLabel} \notin \text{Ancestors}(\text{label}_{\text{Adv}})\}, \\ \text{AbstractedLabel} &\in \{c \in \mathbb{L}' \mid \exists l \in \mathbb{L} \text{ s.t. } c \in \text{Ancestors}(l) \wedge \text{Count}_{\text{Children}}(c) > 0\}, \end{aligned} \quad (9)$$

where \mathbb{L} is the transitive closure of *ImageNet* labels on the hyponymic graph, the construction of AbstractedLabel is equivalent to: (1) Remove overly coarse-grained labels through annotation to obtain the label subset \mathbb{L}' . (2) For each linear path, select the node with the lowest height as a candidate label, constraining the upper bound of the abstracting level. (3) Eliminate descendant labels of the candidate labels to constrain the lower bound of the abstracting level.

Secondly, from the perspective of semantic constraint evaluation, using another deep-learning model for evaluation, *e.g.*, *CLIP*, will limit the benchmark to the robust region of such models. Drawing upon previous discussions and attempts in evaluation enhancement [14], we further conceptualize the sub-task of non-adversarial exemplar generation. As shown in Figure 6, the adversarial generator \mathcal{G} is required to simultaneously generate a nearby sample $x_{\text{exemplar}} \in \mathcal{X}_{\text{exemplar}}$ as a proof that x_{adv} complies with semantic constraints. We further propose the evaluation method based on attack success rate and pair-wise semantic metric as a complement to the single-image assessments:

$$ASR_{\text{Relative}} = \frac{\sum_{i=1}^K \text{Attack Success}(x_{\text{adv}}^{(i)}) \wedge \text{Classification Correct}(x_{\text{exemplar}}^{(i)})}{K \cdot \text{Accuracy}(\mathcal{X}_{\text{exemplar}})} \in [0, 1], \quad (10)$$

$$\text{SemanticDiff}_{\mathcal{S}} = \langle x_{\text{exemplar}}, x_{\text{adv}} \rangle_{\mathcal{S}},$$

where K is the amount of samples, \mathcal{S} is a visual similarity metric, such as *LPIPS* or *MS-SSIM*. Measuring local similarity is easier since high-level feature extraction, which could be attacked, is less required. By assuming the generator \mathcal{G} is not motivated towards finding a *positive adversarial example*, achieving a high score on both metrics can sufficiently show both adversarial capability and instruction compliance of \mathcal{G} . Notably, ASR_{Relative} evaluation metric imposes a **more rigorous** assessment for the masked language guidance in Section 3.3, as it eliminates the confounding variations in benign example generation methods through regularizing with Classification Correct(x_{exemplar})).

4 Experiments

4.1 Experiment Settings

Tasks and Baselines We evaluate different generation methods by generating 6 samples for each label in the ImageNet 1000-class label evasion task and the proposed abstracted label evasion task. The baseline method is constructed by combining **diverse** ① Surrogate models, ② Transfer attack methods (MI-FGSM [3], DeCoWA [41]), and ③ Diffusion-based naturalistic AE generation methods (AdvDiff [13], SD-NAE [12], VENOM [14]). Note that DeCoWa develops on top of MI-FGSM. For fairness, we incorporate the same pretrained diffusion model as SD-NAE and VENOM. For 3D SemanticAE we evaluate the classification models' performance on the generated video showing the rotating object. Detailed implementation and settings are shown in Appendix C.

Evaluation Metrics Referenced image quality assessment metrics, including *LPIPS* [42] and *MSSSIM* [43], are been employed to measure the similarity between x_{exemplar} and x_{adv} under the proposed non-adversarial exemplar evaluation. Also, we supplement the non-reference image quality assessment *CLIP-IQA* [44] for the generated images. We do not use *FID* and *IS* as a primary metric since their adapted vision backbones are simple and might be adversarially attacked, and keeping the feature distribution consistent is not the primary goal of the SemanticAE generation. For attack evaluation, we computed the classification accuracy and ASR_{Relative} , defined in Eq. 10, across diverse targets set $\mathcal{T} = \{\text{ResNet50} [45], \text{ViT-B/16} [46], \text{ConvNeXt-T} [47], \text{ResNet152}, \text{InceptionV3} [48], \text{Swin-Transformer-B} [49]\}$. The first three models are used individually as surrogates in experiments. The main paper presents the average ASR_{relative} and accuracy (ACC) on the target model, while the detailed transfer attack performance on different target models is shown in Appendix D.1. In addition, we fuse the input-transformation-based module *DeCoWa* with ResNet50 as one of the surrogate models to evaluate the collaboration capability of the generation algorithms. 2D / 3D generation times are benchmarked on a single 4090 or A800 GPU, respectively, by generating 100 samples with abstracted labels, and are presented in the results with standard deviation.

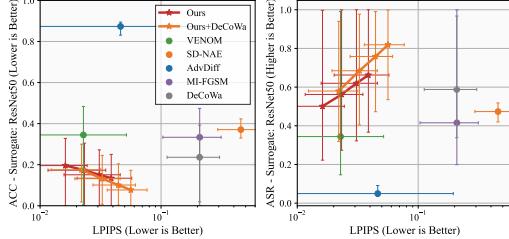


Figure 7: ImageNet label results.

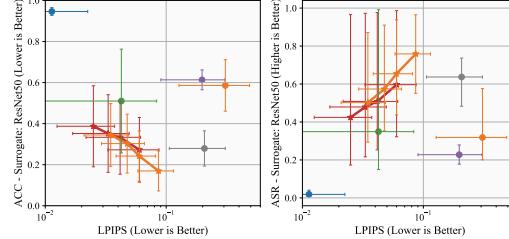


Figure 8: Abstracted label results.

Table 1: Results on more surrogate models. Appendix D shows our Pareto optimality in each setting.

Attacker Settings Surrogates	Method	ImageNet Label SemanticAE Generation Task					Coarse Label Evasion Task (Section 3.4)					Time(s)↓
		Acc.↓	ASR↑	Clip _Q ↑	MSSSIM↑	LPIPS↓	Acc.↓	ASR↑	Clip _Q ↑	MSSSIM↑	LPIPS↓	
ResNet50	MI-FGSM	33.4%	41.5%	0.548	0.880	0.201	61.3%	22.8%	0.551	0.885	0.198	1.43 _{±0.02}
	AdvDiff	87.3%	4.9%	0.634	0.939	0.046	94.6%	1.8%	0.621	0.992	0.011	19.6 _{±0.01}
	SD-NAE	37.1%	47.4%	0.841	0.433	0.457	58.6%	31.8%	0.771	0.599	0.308	24.43 _{±0.14}
	VENOM	34.5%	34.4%	0.795	0.972	0.023	51.0%	34.9%	0.779	0.951	0.043	3.09 _{±0.52}
ViT-B	Ours	15.1%	62.0%	0.815	0.961	0.031	35.2%	47.9%	0.808	0.958	0.033	7.26 _{±2.57}
	MI-FGSM	32.7%	42.4%	0.524	0.855	0.205	58.2%	25.7%	0.521	0.860	0.207	1.46 _{±0.01}
	AdvDiff	65.6%	30.3%	0.638	0.430	0.390	94.1%	2.2%	0.628	0.972	0.026	20.5 _{±0.01}
	SD-NAE	33.7%	51.7%	0.844	0.441	0.459	56.1%	33.6%	0.787	0.609	0.300	24.5 _{±0.10}
	VENOM	30.5%	40.6%	0.796	0.977	0.021	46.3%	40.3%	0.780	0.958	0.040	3.07 _{±0.33}
ConvNeXt	Ours	10.9%	69.7%	0.815	0.956	0.038	28.7%	55.4%	0.814	0.955	0.039	7.23 _{±2.15}
	MI-FGSM	31.9%	44.9%	0.543	0.877	0.204	41.2%	46.4%	0.532	0.88	0.202	1.46 _{±0.01}
	AdvDiff	52.9%	44.4%	0.636	0.312	0.471	93.3%	3.2%	0.627	0.985	0.017	19.9 _{±0.11}
	SD-NAE	22.4%	67.3%	0.848	0.432	0.458	53.6%	36.1%	0.782	0.603	0.308	24.5 _{±0.10}
	VENOM	28.8%	44.6%	0.796	0.978	0.020	42.6%	45.0%	0.785	0.961	0.037	3.14 _{±0.65}
ResNet50 +DeCoWa	Ours	9.1%	75.8%	0.817	0.958	0.036	28.6%	57.4%	0.812	0.957	0.036	6.96 _{±1.96}
	MI-FGSM	23.6%	58.8%	0.535	0.869	0.201	27.9%	63.7%	0.474	0.870	0.207	3.01 _{±0.04}
	AdvDiff	67.9%	27.5%	0.597	0.761	0.293	93.4%	3.2%	0.629	0.934	0.081	20.8 _{±0.18}
	SD-NAE	26.6%	61.7%	0.845	0.421	0.470	64.4%	26.3%	0.782	0.568	0.331	24.3 _{±0.07}
	VENOM	36.1%	31.4%	0.805	0.968	0.027	56.2%	28.1%	0.796	0.944	0.048	4.93 _{±2.67}
	Ours	10.1%	75.8%	0.810	0.947	0.044	30.3%	57.4%	0.808	0.943	0.048	10.7 _{±3.60}

* Attack performance is averaged across targets. Detailed results with the specified target models are shown in Appendix D.

4.2 Overall Performance Evaluation

2D SemanticAE Generation The overall results on 2D SemanticAE generation is shown in Figure 7 and 8 with the strength of semantic constraint of our method set in $\epsilon = \{1.5, 2, 2.5, 3\}$ and $\epsilon = \{2, 2.5, 3, 4\}$ via Eq. 5, respectively. Multiple ϵ are applied since it is difficult to control and align the distortion strength for baselines. The min/max ASR across target models and the standard deviation of LPIPS across generated images are plotted as bars. Table 1 shows the results of our method set as $\epsilon = 2.5$. More results are in the appendix. Overall, in **any** of the 4 surrogate and 2 task settings, **InSUR** is able to achieve **at least 1.19×** average ASR and **1.08×** minimal ASR across **all target models in \mathcal{T}** , and maintains with lower LPIPS (unsuccessful baseline generation with avg. ASR < 5% are not considered), showing the **consistent superiority**. The Pareto improvement shown by the figure is more significant. Moreover, ① ϵ -based semantic constraint in Eq. 5 achieves more consistent LPIPS across images generated in identical settings. ② For $Clip_Q$, our method performs better in the challenging abstracted-label evasion task, while SD-NAE is higher in original tasks.

3D SemanticAE Generation We export the video visualization of the object under MPEG4 encoding, and evaluate the attack performance by reading it. The surrogate and target models are both ResNet50. The results are in Table 2. There is no 3D SemanticAE previously available. It shows that our results show satisfactory attack performance, validating the cross-task scalability of **InSUR**. Note that since 3D-diffusion research is still under development, the clean accuracy on the generated 3D samples is not high, while making **InSUR** a growable research.

4.3 Key Ablation Studies

Residual Approximation We evaluate the effects of residual approximation steps k in the adversarial guidance g in Eq 4, under the 2D abstracted label task and the 3D task. In implementation, k is controlled by the upper-bound parameter K in Eq. 15 in Appendix B.1. The performance improvement is significant and consistent compared to the result without future-step sampling prediction ($K = 0$), which represents the original DDIM sampling with adversarial guidance. The naturalness metrics, including $Clip_Q$, MSSSIM, and LPIPS, are also slightly improved. Moreover, attributed to

Table 2: 3D generation results.

Generator	Acc.	ASR	MSSSIM	LPIPS
Non-Adversarial	21.5%	—	—	—
Ours w/o ResAdv	17.9%	45.1%	0.658	0.261
Ours	2.8%	92.2%	0.665	0.258

the adaptive iteration mechanism, more accurate estimation leads to lower optimization steps, and therefore, the increase in time consumption is sub-linear. Detailed results on more settings are shown in the Appendix D.3. Parameter and time consumption analysis are shown with the implementation details in Appendix B.

The effect of different k on the process of adversarial optimization is shown in Figure 9. The setting of Text is *a dog sitting on the floor*, and the evasion label is *dog*. As k increases, the estimated ASR after adversarially sampling x_t increases earlier. Since the white-box ASR in both settings is nearly 100%, the improvements are from: (1) by eliminating the *referring diversify*, the initial sampling steps receive more accurate guidance. (2) More effective adversarial optimization on earlier diffusion denoising steps provides better on-manifold regularization, leading to better visual quality and better adversarial transferability.

Table 3: Ablation of residual approximation

Surro.	K	Acc.	ASR	Clip_Q	MSSSIM	LPIPS	Time
ViT-B	0	43.1%	43.3%	0.794	0.941	0.055	4.53 ± 0.19
	1	33.5%	54.7%	0.812	0.942	0.049	6.87 ± 1.48
	2	31.3%	56.6%	0.816	0.942	0.049	7.44 ± 1.92
	3	29.2%	57.5%	0.807	0.943	0.048	7.75 ± 2.37
	4	27.7%	60.1%	0.813	0.944	0.047	7.87 ± 2.31
ResNet50 +DeCoWa	0	39.8%	47.1%	0.764	0.946	0.053	5.65 ± 0.38
	1	36.3%	50.4%	0.812	0.954	0.040	9.64 ± 2.64
	2	32.8%	52.6%	0.818	0.955	0.039	10.7 ± 3.27
	3	30.4%	53.9%	0.817	0.955	0.039	11.2 ± 3.67
	4	29.5%	54.2%	0.814	0.955	0.039	11.5 ± 4.00
3D Gen.	0	17.9%	45.1%	–	0.658	0.261	7.49 ± 1.49
	1	3.4%	90.2%	–	0.658	0.262	30.4 ± 35.02
	2	3.4%	91.2%	–	0.659	0.261	32.7 ± 39.17
	3	2.9%	91.2%	–	0.671	0.255	41.1 ± 53.75
	4	2.8%	92.2%	–	0.665	0.258	40.1 ± 59.25

Spatial Masking of Language Guidance We evaluate the effect of guidance masking by setting $M_{\text{edge}}/M_{\text{mid}}$ to different values in Eq. 6, and setting to 1.0 represent removing the masking. The results are shown in Table 4 and Figure 10, indicating that (1) decreasing M_{edge} leads to an increase in unconditional guidance, which enriches the background diversity as shown in the figure, and may lead to the improvements of Clip_Q . (2) When the budget ϵ is small ($\epsilon = 2$) and the optimization space is relatively narrow, diversifying the background is beneficial. Note that ϵ is large, the improvement is marginal since there is no need for expanding the optimization space. Overall, the guidance masking design improves diffusion models’ adaptability to strongly constrained SemanticAE generation.

Table 4: Ablation of guidance masking.

ϵ	$M_{\text{edge}}/M_{\text{mid}}$	Acc.	ASR	Clip_Q	MSSSIM	LPIPS
2	0.0	32.8%	52.3%	0.813	0.958	0.035
	0.1	34.5%	50.2%	0.809	0.957	0.035
	1.0	36.9%	48.3%	0.788	0.958	0.034
4	0.0	16.6%	75.9%	0.804	0.901	0.087
	1.0	17.3%	75.9%	0.775	0.904	0.084

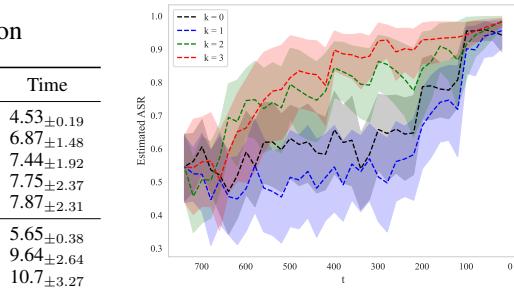


Figure 9: Estimated ASR of \hat{x}_0 after adv. sampling x_t , under different settings of k in Eq. 4. Earlier increase of estimated ASR leads to better naturalness and higher attack transferability.

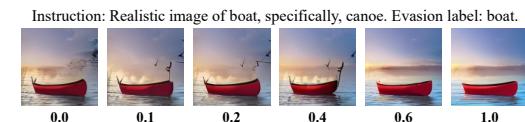


Figure 10: Vis. of different guidance masking. The value under images denotes $M_{\text{edge}}/M_{\text{mid}}$.

4.4 Visualization of Generated SemanticAEs

We selected the 2D / 3D samples with x_{exemplar} correctly classified and visualized in Fig. 11 and Fig. 12. 2D samples are generated with the original ImageNet-label and the surrogate model is DeCoWa+ResNet. 3D samples are generated and visualized as the main experiment. The results are coherent with the main table, showing that MI-FGSM is not natural, SD-NAE disturbs more global semantics, while ours achieves both global semantic preservation and naturalness. Through observation, our method generates local in-manifold patterns to achieve the strong attacks, *e.g.*, adding fog in the castle image or altering the lightning in the jellyfish image. For 3D results, although the MSSIM and LPIPS metrics are not excellent in the main experiment, the generated 3D objects are natural and follow the semantic constraint. More visualizations are shown in Appendix E.

4.5 Discussions

Extending to Attacking Large Vision Language Models The InSUR framework is agnostic to specific target models and tasks, and can be directly applied to adversarial attack evaluation

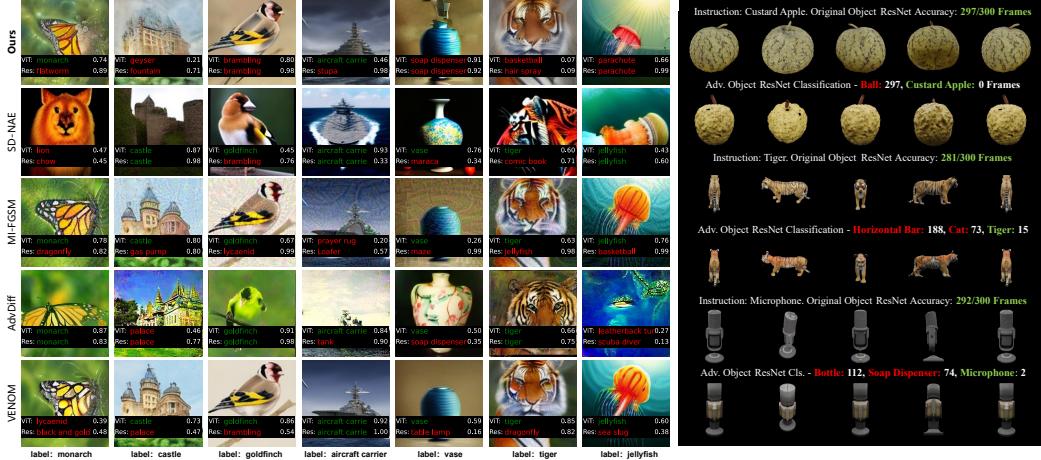


Figure 11: Comparison of 2D SemanticAEs.

Figure 12: 3D SemanticAEs

in Visual Question Answering (VQA) scenarios. To illustrate the capability, we conduct VQA experiments using OpenAI CLIP [50] model *RN50* (with DeCoWA augmentation) as the surrogate, using *stabilityai/stable-diffusion-2-1-base* [51] as the diffusion generator, and generate adversarial examples for the VQA problem. The question and options are “*What is in the picture:*” and {*Police car, Ambulance, Taxi, School Bus*}. We evaluate ***weak-to-strong*** transferability across OpenAI CLIP variants {RN50, RN101, RN50x4, RN50x16, ViT-B/32, ViT-B/16} and larger vision language models (VLM), LLaVA [52] and Qwen-7B [53]. The results are in Table 5, showing that our technical contribution is also effective in the VQA transfer attack task.

Table 5: ASR_{relative} (%) of attacks. The surrogate model is ResNet50 CLIP model.

Residual Approximation	Target Model							
	RN50	RN101	RN50x4	RN50x16	ViT-B/32	ViT-B/16	LLaVa1.5-7B	Qwen2.5VL-7B
✗ ($K = 0$)	77.50	27.50	30.00	11.11	12.50	10.53	11.11	16.22
✓ ($K = 3$)	97.50	57.50	62.50	27.78	40.00	47.37	22.22	35.14

Broader Applications Our approach also has the potential to extend to real-world adversarial attacks under the 3DGS representations [28], providing technical tools for security-related research. We tested our generated 3D SemanticAE in with lightning, material, and camera settings with *Blender* environment, and the ResNet50 ACC is **59.1%, 20.0%**, for the two settings in Figure 13, respectively. Although the attack performance is worse than the original rendering (ACC=7.7%), the relative improvement is significant. Our framework can also be integrated with existing 3D scene generation pipelines [54], enabling more effective incorporation of diffusion-based adversarial optimization into world models. This supports safety-critical scenario generation for applications such as autonomous driving [55, 56] and embodied agents [57]. On another front, recent work [58] demonstrates that diffusion-generated hard samples can achieve higher sample efficiency, especially with adversarial guidance, suggesting that our InSUR framework also holds promise for supporting the generation of adversarial training data.

5 Conclusion

This paper proposes multi-dimensional uncertainty reduction frameworks for SemanticAE generation, pushes the boundary of the 2D generation, and opens the door to 3D generation. The proposed technology consistently improves the attack performance and has the potential to scale to other tasks. Moreover, we believe it could provide valuable insights for the test-time scaling of the red-teaming framework. For limitations, there is scope for improvement in the generation quality, the evaluation on larger models, and the application in the real world, suggesting future research avenues of the SemanticAE generation algorithms based on concrete generative models and the scenario adaptation methods oriented to real-world scenarios.

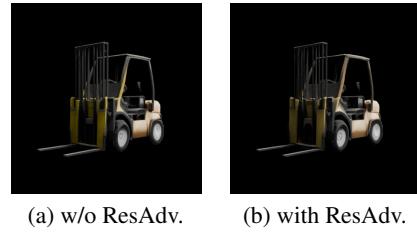


Figure 13: *Blender* re-rendered 3D SemanticAEs with altered lightning.

Acknowledgments and Disclosure of Funding

This work was supported by the National Science Foundation of China (NSFC) (Grant No. 62476018), and was supported by Zhongguancun Laboratory.

References

- [1] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *arXiv preprint arXiv:1607.02533*, 2016.
- [2] Hui Wei, Hao Tang, Xuemei Jia, Zhixiang Wang, Hanxun Yu, Zhubo Li, Shin’ichi Satoh, Luc Van Gool, and Zheng Wang. Physical adversarial attack meets computer vision: A decade survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12):9797–9817, 2024.
- [3] Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boosting adversarial attacks with momentum. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9185–9193, 2018.
- [4] Jiakai Wang, Aishan Liu, Zixin Yin, Shunchang Liu, Shiyu Tang, and Xianglong Liu. Dual attention suppression attack: Generate adversarial camouflage in physical world. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8565–8574, 2021.
- [5] Yu-Chih-Tuan Hu, Bo-Han Kung, Daniel Stanley Tan, Jun-Cheng Chen, Kai-Lung Hua, and Wen-Huang Cheng. Naturalistic physical adversarial patch for object detectors. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7848–7857, 2021.
- [6] Haotian Xue, Alexandre Araujo, Bin Hu, and Yongxin Chen. Diffusion-based adversarial sample generation for improved stealthiness and controllability. In *Proceedings of the Advances in Neural Information Processing Systems*, volume 36, pages 2894–2921, 2023.
- [7] Zelun Kong, Junfeng Guo, Ang Li, and Cong Liu. Physgan: Generating physical-world-resilient adversarial examples for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14254–14263, 2020.
- [8] Wenhao Ding, Chejian Xu, Mansur Arief, Haohong Lin, Bo Li, and Ding Zhao. A survey on safety-critical driving scenario generation—a methodological perspective. *IEEE Transactions on Intelligent Transportation Systems*, 24(7):6971–6988, 2023.
- [9] Chen Ma, Ningfei Wang, Qi Alfred Chen, and Chao Shen. Slowtrack: Increasing the latency of camera-based perception in autonomous driving using adversarial examples. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 4062–4070, 2024.
- [10] Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023.
- [11] Jia-Yu Yao, Kun-Peng Ning, Zhen-Hui Liu, Mu-Nan Ning, Yu-Yang Liu, and Li Yuan. Llm lies: Hallucinations are not bugs, but features as adversarial examples. *arXiv preprint arXiv:2310.01469*, 2023.
- [12] Yueqian Lin, Jingyang Zhang, Yiran Chen, and Hai Li. SD-NAE: Generating natural adversarial examples with stable diffusion. In *The Second Tiny Papers Track at ICLR 2024*, 2024.
- [13] Xuelong Dai, Kaisheng Liang, and Bin Xiao. Advdiff: Generating unrestricted adversarial examples using diffusion models. In *European Conference on Computer Vision*, pages 93–109. Springer, 2024.
- [14] Hui Kuurila-Zhang, Haoyu Chen, and Guoying Zhao. Venom: Text-driven unrestricted adversarial example generation with diffusion models. *arXiv preprint arXiv:2501.07922*, 2025.

- [15] Songtao Li and Hao Tang. Multimodal alignment and fusion: A survey. *arXiv preprint arXiv:2411.17040*, 2024.
- [16] Yixu Wang, Jiaxin Song, Yifeng Gao, Xin Wang, Yang Yao, Yan Teng, Xingjun Ma, Yingchun Wang, and Yu-Gang Jiang. Safevid: Toward safety aligned video large multimodal models. *arXiv preprint arXiv:2505.11926*, 2025.
- [17] Shiji Zhao, Ranjie Duan, Fengxiang Wang, Chi Chen, Caixin Kang, Shouwei Ruan, Jialing Tao, YueFeng Chen, Hui Xue, and Xingxing Wei. Jailbreaking multimodal large language models via shuffle inconsistency. *ICCV*, 2025.
- [18] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015.
- [19] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.
- [20] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *International conference on machine learning*, pages 2206–2216. PMLR, 2020.
- [21] Yu-Chih-Tuan Hu, Bo-Han Kung, Daniel Stanley Tan, Jun-Cheng Chen, Kai-Lung Hua, and Wen-Huang Cheng. Naturalistic physical adversarial patch for object detectors. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7848–7857, October 2021.
- [22] Hadi Mohaghegh Dolatabadi, Sarah Erfani, and Christopher Leckie. Advflow: Inconspicuous black-box adversarial attacks using normalizing flows. In *Neural Information Processing Systems*, volume 33, pages 15871–15884, 2020.
- [23] Shumeet Baluja and Ian Fischer. Learning to attack: Adversarial transformation networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [24] Chaowei Xiao, Bo Li, Jun-Yan Zhu, Warren He, Mingyan Liu, and Dawn Song. Generating adversarial examples with adversarial networks. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pages 3905–3911, 2018.
- [25] Jin Hu, Xianglong Liu, Jiakai Wang, Junkai Zhang, Xianqi Yang, Haotong Qin, Yuqing Ma, and Ke Xu. Dynamiccpae: Generating scene-aware physical adversarial examples in real-time. *arXiv preprint arXiv:2412.08053*, 2024.
- [26] Anish Athalye, Logan Engstrom, Andrew Ilyas, and Kevin Kwok. Synthesizing robust adversarial examples. In *International conference on machine learning*, pages 284–293. PMLR, 2018.
- [27] Yanjie Li, Bin Xie, Songtao Guo, Yuanyuan Yang, and Bin Xiao. A survey of robustness and safety of 2d and 3d deep learning models against adversarial attacks. *ACM Comput. Surv.*, 56(6), January 2024.
- [28] Tianrui Lou, Xiaojun Jia, Siyuan Liang, Jiawei Liang, Ming Zhang, Yanjun Xiao, and Xiaochun Cao. 3d gaussian splatting driven multi-view robust physical adversarial camouflage generation. *International Conference on Computer Vision*, 2025.
- [29] Yang Song, Rui Shu, Nate Kushman, and Stefano Ermon. Constructing unrestricted adversarial examples with generative models. In *Proceedings of the Advances in Neural Information Processing Systems*, volume 31, 2018.
- [30] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15262–15271, 2021.
- [31] Xinquan Chen, Xitong Gao, Juanjuan Zhao, Kejiang Ye, and Cheng-Zhong Xu. Advdifuser: Natural adversarial example synthesis with diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4562–4572, 2023.

- [32] Jianqi Chen, Hao Chen, Keyan Chen, Yilan Zhang, Zhengxia Zou, and Zhenwei Shi. Diffusion models for imperceptible and transferable adversarial attack. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [33] Yanbo Chen and Weiwei Liu. A theory of transfer-based black-box attacks: Explanation and implications. In *Proceedings of the Advances in Neural Information Processing Systems*, volume 36, pages 13887–13907, 2023.
- [34] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021.
- [35] Tianlong Chen, Xiaohan Chen, Wuyang Chen, Howard Heaton, Jialin Liu, Zhangyang Wang, and Wotao Yin. Learning to optimize: A primer and a benchmark. *Journal of Machine Learning Research*, 23(189):1–59, 2022.
- [36] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based semantic image editing with mask guidance. In *The Eleventh International Conference on Learning Representations*, 2023.
- [37] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.
- [38] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):139–1, 2023.
- [39] Jianfeng Xiang, Zelong Lv, Sicheng Xu, Yu Deng, Ruicheng Wang, Bowen Zhang, Dong Chen, Xin Tong, and Jiaolong Yang. Structured 3d latents for scalable and versatile 3d generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 21469–21480, 2025.
- [40] George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995.
- [41] Qinliang Lin, Cheng Luo, Zenghao Niu, Xilin He, Weicheng Xie, Yuanbo Hou, Linlin Shen, and Siyang Song. Boosting adversarial transferability across model genus by deformation-constrained warping. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 3459–3467, 2024.
- [42] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.
- [43] Zhou Wang, Eero P Simoncelli, and Alan C Bovik. Multiscale structural similarity for image quality assessment. In *The Thrity-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*, volume 2, pages 1398–1402. Ieee, 2003.
- [44] Jianyi Wang, Kelvin CK Chan, and Chen Change Loy. Exploring clip for assessing the look and feel of images. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pages 2555–2563, 2023.
- [45] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [46] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- [47] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11976–11986, 2022.

- [48] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016.
- [49] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021.
- [50] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [51] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, June 2022.
- [52] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Proceedings of the Advances in neural information processing systems*, volume 36, pages 34892–34916, 2023.
- [53] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- [54] Beichen Wen, Haozhe Xie, Zhaoxi Chen, Fangzhou Hong, and Ziwei Liu. 3d scene generation: A survey. *arXiv preprint arXiv:2505.05474*, 2025.
- [55] Yuping Wang, Shuo Xing, Cui Can, Renjie Li, Hongyuan Hua, Kexin Tian, Zhaobin Mo, Xiangbo Gao, Keshu Wu, Sulong Zhou, et al. Generative ai for autonomous driving: Frontiers and opportunities. *arXiv preprint arXiv:2505.08854*, 2025.
- [56] Chejian Xu, Aleksandr Petushko, Ding Zhao, and Bo Li. Diffscene: Diffusion-based safety-critical scenario generation for autonomous vehicles. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 8797–8805, 2025.
- [57] Xinjie Wang, Liu Liu, Yu Cao, Ruiqi Wu, Wenkang Qin, Dehui Wang, Wei Sui, and Zhizhong Su. Embodiedgen: Towards a generative 3d world engine for embodied intelligence. *arXiv preprint arXiv:2506.10600*, 2025.
- [58] Reyhane Askari-Hemmat, Mohammad Pezeshki, Elvis Dohmatob, Florian Bordes, Pietro Astolfi, Melissa Hall, Jakob Verbeek, Michal Drozdzal, and Adriana Romero-Soriano. Improving the scaling laws of synthetic data with deliberate practice. In *International Conference on Machine Learning*. PMLR, 2025.
- [59] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.
- [60] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. In *International Conference on Machine Learning*, pages 32211–32252. PMLR, 2023.
- [61] Zhengyang Geng, Mingyang Deng, Xingjian Bai, J Zico Kolter, and Kaiming He. Mean flows for one-step generative modeling. In *Proceedings of the Advances in neural information processing systems*, 2025.
- [62] Viraj Prabhu, Sriram Yenamandra, Prithvijit Chattopadhyay, and Judy Hoffman. Lance: Stress-testing visual models by generating language-guided counterfactual images. In *Proceedings of the Advances in Neural Information Processing Systems*, volume 36, pages 25165–25184, 2023.
- [63] Chenshuang Zhang, Fei Pan, Junmo Kim, In So Kweon, and Chengzhi Mao. Imagenet-d: Benchmarking neural network robustness on diffusion synthetic object. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21752–21762, 2024.

- [64] Sreyan Ghosh, Chandra Kiran Reddy Evuru, Sonal Kumar, Utkarsh Tyagi, S. Sakshi, Sanjoy Chowdhury, and Dinesh Manocha. ASPIRE: language-guided data augmentation for improving robustness against spurious correlations. In Lun-Wei Ku, Andre Martins, and Vivek Srikanth, editors, *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 386–406. Association for Computational Linguistics, 2024.
- [65] Xiaopei Zhu, Peiyang Xu, Guanning Zeng, Yinpeng Dong, and Xiaolin Hu. Natural language induced adversarial images. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 10872–10881, 2024.
- [66] Wenkai Yang, Shiqi Shen, Guangyao Shen, Wei Yao, Yong Liu, Zhi Gong, Yankai Lin, and Ji-Rong Wen. Super(ficial)-alignment: Strong models may deceive weak models in weak-to-strong generalization. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [67] Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples and black-box attacks. In *International Conference on Learning Representations*, 2017.
- [68] Martin Gubri, Maxime Cordy, Mike Papadakis, Yves Le Traon, and Koushik Sen. Lgv: Boosting adversarial example transferability from large geometric vicinity. In *European Conference on Computer Vision*, pages 603–618. Springer, 2022.
- [69] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Representations*, 2021.
- [70] Kunyu Wang, Xuanran He, Wenxuan Wang, and Xiaosen Wang. Boosting adversarial transferability by block shuffle and rotation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24336–24346, 2024.
- [71] Xin Wang, Jie Ren, Shuyun Lin, Xiangming Zhu, Yisen Wang, and Quanshi Zhang. A unified approach to interpreting and boosting adversarial transferability. In *The Ninth International Conference on Learning Representations*, 2021.
- [72] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan. Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4541–4550, 2019.
- [73] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *Proceedings of the European conference on computer vision*, pages 52–67, 2018.
- [74] Chaowei Xiao, Dawei Yang, Bo Li, Jia Deng, and Mingyan Liu. Meshadv: Adversarial meshes for visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6898–6907, 2019.
- [75] Yao Huang, Yinpeng Dong, Shouwei Ruan, Xiao Yang, Hang Su, and Xingxing Wei. Towards transferable targeted 3d adversarial attack in the physical world. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24512–24522, 2024.
- [76] Yanjie Li, Kaisheng Liang, and Bin Xiao. UV-attack: Physical-world adversarial attacks on person detection via dynamic-neRF-based UV mapping. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [77] Sergey Kastrulyin, Jamil Zakirov, Denis Prokopenko, and Dmitry V Dylov. Pytorch image quality: Metrics for image quality assessment. *arXiv preprint arXiv:2208.14818*, 2022.
- [78] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. CLIPScore: A reference-free evaluation metric for image captioning. In *Proceedings of the Empirical Methods in Natural Language Processing*, pages 7514–7528, 2021.

Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The primary claims of the paper are addressed with the provided methodology and the experiment evaluates the effectiveness.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We discuss the limitations in the conclusion section of the generation performance, the evaluation on larger models, and the application in the real world

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: No theoretical results are included in the main paper.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We make it clear how to reproduce that algorithm as a new algorithm. The detailed implementation is provided in Appendix B.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: We will release the code after publication.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer:[Yes]

Justification: We specify the Experimental setting in the Experimental section and will include the details in Appendix C.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

We have selected statistical significance verification methods based on the properties of the metrics, including the standard deviation, extremum difference across different task settings.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We provide the execution time in Table 1 and 3 on a single 4090 or A800 GPU. Detailed experiment settings are shown in Appendix C.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The research conducted in the paper conforms, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: Societal impacts are expounded upon in the discussion section in the appendix.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: No risk of misuse.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We have reflected Licenses for existing assets in our references.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The paper does not release new assets

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.

Contents

1	Introduction	1
2	Backgrounds	2
3	Methodology	3
3.1	Problem Formulation and Analysis	3
3.2	Residual-driven Attacking Direction Stabilization with ResAdv-DDIM	4
3.3	Context-encoded Attacking Scenario Constraint for 2D and 3D Generation	5
3.4	Semantic-abstracted Attacking Evaluation Enhancement with Label Taxonomy	6
4	Experiments	7
4.1	Experiment Settings	7
4.2	Overall Performance Evaluation	8
4.3	Key Ablation Studies	8
4.4	Visualization of Generated SemanticAEs	9
4.5	Discussions	9
5	Conclusion	10
A	Technical Background and Related Works	24
A.1	Diffusion models	24
A.2	Generating Hard Samples from Language	24
A.3	Transfer Attack Methodology	25
A.4	3D Generation and Gaussian Splatting	25
B	Detailed Implementation of InSUR Framework	25
B.1	2D Image Generation	25
B.2	3D Object Generation	28
B.3	Construction and Analysis of Semantic-Abstracted Evaluation Task	30
C	Detailed Experiment Settings	32
C.1	2D Experiment Settings	32
C.2	3D Experiment Settings	33
D	Detailed Results and Discussions	34
D.1	Transfer Attack Analysis	34
D.2	Attack Robustness Analysis	35
D.3	On the Role of Residual-driven Attacking Direction Stabilization in 2D and 3D SemanticAE Applications	35
D.4	On the Potential Adversarial Transferability to Semantic Evaluator	37
E	Results Visualization	38
E.1	2D Visualized Results	38
E.2	3D Visualized Results	39
F	Boarder Impacts	40

A Technical Background and Related Works

A.1 Diffusion models

Diffusion Model This work mainly applies the diffusion model as the language-guided data generator. As a brief review, diffusion models establish the theoretical and technical route of estimating the score function $\nabla_x \log p(x)$ by training UNet models on the dataset $\{x\}$, and sampling the data from the score function with numerical methods. One of the key insight of diffusion models is alleviating the training difficulty by disturbing the data x with noise, which re-constructs the as a forward process from x_0 to x_T with the process $x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon$, $\epsilon \sim \mathcal{N}(0, I)$ scheduled by β . The training is performed by learning the denoising process $\epsilon_\theta(x_t, t)$, and during the inference, the data is sampled from x_T to x_0 step by step with ϵ_θ . DDIM model reinterprets the forward process as $p(x_t|x_{t-1}, x_0)$, which abandons the Markov property, and constructs the sampling method by $q_\theta(x_{t-1}|x_t, x_0)$. It provides the theoretical grounding for the step jumping in the sampling process. The DDIM [34] sampling procedure could be formulated as (deterministic version):

$$x_{t-\Delta T} = \sqrt{\bar{\alpha}_{t-\Delta T}/\bar{\alpha}_t} (x_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_\theta(x_t, t)) + \sqrt{1 - \bar{\alpha}_{t-\Delta T}} \cdot \epsilon_\theta(x_t, t), \quad (11)$$

where $\bar{\alpha} = \prod_{s=1}^t (1 - \beta_s)$, and ΔT is the step interval. We adapted this notation in the main paper. The language-guided generation problem is modeled by adding the conditional term Context that corresponds to the data x and sampling the data by $\epsilon_\theta(x_t, t, \text{Context})$, where Context is the encoding given by the language model. To balance the generation diversity and the instruction following, conditional and unconditional guidance are integrated, *i.e.*, $\epsilon_\theta = -\omega\epsilon_\theta(x_t, t, \text{Unconditional}) + (1 + \omega)\epsilon_\theta(x_t, t, \text{Context})$, in the *classifier-free guidance* [59].

Discussions In the application in content edit [36], masking has been applied to constrain the editing area, *i.e.* $x_{t-1} = \text{Mask} \cdot x_{t,\text{edit}} + (1 - \text{Mask}) \cdot x_{t,\text{original}}$. In our 2D generation task, we take a further step to model the interaction between conditional and unconditional guidance, and achieve the new function of re-distributing the spatial strength of the semantic constraint. In recent years, training techniques that enable single-step generation have been proposed [60, 61]. Although our method is designed with multi-step diffusion models, our method still has practical value since (1) there is a performance gap between multi-step and single-step diffusion models, and (2) as the adversarial optimization is performed with local perturbations, sampling in small step sizes might provide a better regularization for transfer attacks.

A.2 Generating Hard Samples from Language

Language-guided adversarial example generation is still under development. We categorize two major paradigms. ① Perturbing the input text without altering the language-guided data generation model. A line of study has focused on utilizing the language-guided image edition model to evaluate the robustness of visual models [62, 63, 64]. Compared to adversarial attacks, they focus more on enriching the dataset than finding hard examples for victim models, while the adversarial capability is relatively limited. Related to the discussion in Section 3.2 of the main paper, overly perturbing the text guidance may lead to unsatisfactory semantic alignment. To solve this problem, [65] integrates the additional *CLIP* supervision on the generated image. ② Altering the language-guided data generation model by introducing the adversary capability. Although it has difficulties in global optimization, this paradigm has an advantage in the fine-grained control of the generated samples.

We focus on the latter paradigm, since (1) Perturbing the text input only will limit the capability of the generation method within the semantic knowledge learned by both discriminative and generative models. *e.g.*, a *CLIP*-based semantic supervision may not attack the *CLIP* model itself. (2) Optimizing the input text for the entire data generation process with the adversarial feedback is time-consuming. These features are crucial from the perspective of **SuperAlignment** [66], *i.e.*, creating efficient methodologies that utilize small models with limited capability to build the oversight framework for larger models. Also, we believe that the two methods could be implemented into an integrated red-team system, *i.e.*, the language-space reject sampling as the high-level generation method and the latent-space optimization as the mid or low-level generation module. Therefore, our proposed evaluation task focuses on the second paradigm, serving as an evaluation of the key component of the red-team framework.

A.3 Transfer Attack Methodology

Transferable attack denotes finding adversarial examples that can attack other unknown or even stronger black-box models, which is practical for revealing real-world AI safety problems. The general method is to find a surrogate model(s) in the related task and optimize the adversarial example for it with the regularization. From the perspective of the optimization pipeline, regularization on three modules has been investigated to improve transferability: ① The gradient decent. Momentum is the general method to boost the transferability. ② The surrogate model. Model ensemble can effectively enhance transferability by combining gradient features from multiple distinct models [67, 68]. Constructing a proper objective using the Sharpness-Aware Minimization (SAM) [69] or the attention mechanism [4], to model the cross-model vulnerability, is also beneficial for transferability. ③ The input transformation. Transfer improvement could be achieved by inserting a random transformation between the current-step adversarial example and the surrogate model’s input, the [70, 41]. This could be regarded as improving the diversity of the surrogate model. From the theoretical perspective, decreasing the cooperation within the adversarial pattern [71] and in-manifold constraint [33] is beneficial for transferability. In the adversarial example generation pipeline, the generative-model-based method, which this work focuses on, could be regarded as a replacement of gradient descent with a better in-manifold constraint. We construct baselines and evaluation tasks based on this background.

A.4 3D Generation and Gaussian Splatting

3D geometric data employs diverse representations that are crucial in 3D generation. These representations can be categorized into three types: ① explicit representations, such as point clouds [72] and meshes [73], ② implicit representations, such as Neural Radiance Fields (NeRFs) [37], ③ hybrid representations, such as 3D Gaussian [38]. Currently, 3D Gaussian splatting is widely used in 3D generation due to the geometric editability, high-frequency detail capture capability and real-time rendering efficiency. Each 3D Gaussian point can be represented by the following parameters: position x , spherical harmonics coefficients c , opacity α , a rotation matrix r , and a scaling matrix s . Then, 3D Gaussian points can be projected onto the image plane via the viewing transformation W and the Jacobi affine approximation matrix of the projection transformation matrix in geometry. In terms of appearance, we can calculate the color of every pixel in the 2D image by blending N ordered points overlapping the pixel using spherical harmonics coefficients c and opacities α .

With the development of 3D representations, 3D adversarial examples for different 3D structures have also been advanced [74, 75, 76]. Unlike other adversarial methods, which are based on the geometric or texture data of existing objects, for the first time, we realize the reference-free generation of semantically constrained 3D adversarial examples by utilizing language-guided 3D generation models. We implement language-guided 3D adversarial example generation through the ResAdv-DDIM sampler referencing the optimization pipeline of Trellis in latent space and decode into 3D Gaussian formats. Then, we achieve 3D adversarial attacks on 2D target models using 3D Gaussian rendering.

B Detailed Implementation of InSUR Framework

Adversarial attacks often require reengineering existing tools to achieve new functions, and their implementation is not simple, especially when achieving new characteristics. Therefore, we provide the design principles and the key techniques in the main paper, and supplement the detailed implementation in this section as a complement. In the following subsections, we describe the implementation of the SemanticAE generator based on the different scenarios, and then provide more details about the *Semantic-abstacted Attacking Evaluation Enhancement*.

B.1 2D Image Generation

ImageNet Object Generation with Guidance Masking For the classical image generation problem, we start with the original DDIM sampling process:

$$x_{t-\Delta T} = \sqrt{\bar{\alpha}_{t-\Delta T}/\bar{\alpha}_t} (x_t - \sqrt{1-\bar{\alpha}_t} \epsilon_\theta(x_t, t)) + \sqrt{1-\bar{\alpha}_{t-\Delta T}} \cdot \epsilon_\theta(x_t, t), \quad (12)$$

where $\bar{\alpha} = \prod_{s=1}^t (1 - \beta_s)$ is the sampling parameter, and ΔT is the step interval. As described in section 3.3, masked guidance is adapted as the conditional diffusion guidance, which formulates the denoise function $x_{t-\Delta T} := f_{\theta, \Delta T}(x_t)$ as:

$$\begin{aligned} f_{\theta, \Delta T}^{(t, \text{Text})}(x_t) &= \sqrt{\frac{\bar{\alpha}_{t-\Delta T}}{\bar{\alpha}_t}} (x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_{\theta}(x_t, t, \text{Text})) + \sqrt{1 - \bar{\alpha}_{t-\Delta T}} \cdot \epsilon_{\theta}(x_t, t, \text{Text}), \\ \epsilon_{\theta}(x_t, t, \text{Text}) &:= (1 - M) \cdot \epsilon_{\theta, \text{Unconditional}}(x_t, t) + M \cdot \epsilon_{\theta, \text{Conditional}}(x_t, t, \text{Text}). \end{aligned} \quad (13)$$

The guidance mask is defined as a matrix with different values in the border elements:

$$M_{ij} := \begin{cases} M_{\text{mid}} & \frac{h}{16} \leq i < \frac{15h}{16}, \frac{w}{16} \leq j < \frac{15w}{16}, \\ M_{\text{edge}} & \text{otherwise.} \end{cases} \quad (14)$$

In the main paper, we use the simplified notation of f without parameters Text and t explicitly written. Here we use the detailed notations. For the experiments, we adapted $M_{\text{mid}} = 3.0$ and $M_{\text{edge}} = 0.3$ in the main experiment and selected $M_{\text{edge}} = \{0.0, 0.3, 3.0\}$ in the ablation study.

Residual Approximation Recall that ResAdv-DDIM defines a residual estimation function g for the adversarial feedback of the target model. We construct the step size of g as evenly distributed:

$$\begin{aligned} g_{\theta}^{(t, \text{Text})}(x_t) &:= \underbrace{f_{\theta, \Delta T_1}^{(\Delta T_1, \text{Text})} \circ f_{\theta, \Delta T_2}^{(\Delta T_2 + \Delta T_1, \text{Text})} \circ \dots \circ f_{\theta, \Delta T_k}^{(t, \text{Text})}(x_t)}_{k \text{ times}}, \text{ where } \sum_{i=1}^k \Delta T_i = t, \\ k &:= \lceil \frac{t}{[T/K]} \rceil, \quad \Delta T_i(t) = \begin{cases} \lfloor T/K \rfloor, & 1 < i \leq k, \\ t \bmod \lfloor T/K \rfloor, & i = 1 \end{cases}, \end{aligned} \quad (15)$$

where K is the maximal iteration number, corresponding to Iter_{\max} in the ablation study in the main paper. For brevity, we define T as the number of sampling steps of the original DDIM generator and use the sampling interval of the original DDIM generator as the unit of ΔT . We let t_s (default set as 0.75) be the start step of the adversarial optimization, and the sampling process with adversarial optimization is formulated as:

$$x_{t-\Delta T} = \begin{cases} f_{\theta, \Delta T}^{(t, \text{Text})} \left(\arg \max_{x_t} \mathcal{L}_{\text{ATK}}(\mathcal{M} \circ g_{\theta}^{(t, \text{Text})}(x_t)) \right), & t \leq t_s, \\ f_{\theta, \Delta T}^{(t, \text{Text})}(x_t), & t > t_s. \end{cases} \quad (16)$$

Collaborating with Guidance Masking To generate attack-related image backgrounds without disturbing the foreground semantics, as the goal of *Context-encoded Attacking Scenario Constrain*, we let the exemplar generation $x'_t \rightarrow x'_{t-\Delta T}$ communicate with the adversarial generation $x_t \rightarrow x_{t-\Delta T}$. Specifically, at the beginning of the adv. optimization in step t_s , we set the benign sample generated from $x'_{t_s} \leftarrow x_{t_s}$ as the initialization of the constraint anchor. At the end of the adversarial optimization, the optimized background is written back $x'_0 [M = M_{\text{edge}}]_{\text{Select}} \leftarrow \frac{1}{2}(x_0 + x'_0) [M = M_{\text{edge}}]_{\text{Select}}$ after the generation.

Adaptive Optimization Iteration To improve the efficiency in multi-step diffusion-based adversarial optimization, we implement the adaptive iteration mechanism for ResAdv-DDIM by early-stopping the optimization problem in Eq 16, which is formulated as:

$$\begin{aligned} n &:= \begin{cases} M, & t = t_s \vee t < 4, \\ m, & \text{otherwise.} \end{cases} \\ \text{PerformOptimize} &:= (i \leq n \wedge \underbrace{\arg \max_{l \notin A_{\text{Text}}} \text{P}(\mathcal{M}(x_t) = l) > \xi_1}_{\text{If confidence} \leq \xi_1, \text{early stop. Maximal optimize iterations is } n.}) \vee \\ &\quad (i = 1 \wedge \underbrace{\arg \max_{l \notin A_{\text{Text}}} \text{P}(\mathcal{M}(x_t) = l) > \xi_2}_{\text{If confidence} \leq \xi_2, \text{do not optimize at the first step.}}), \end{aligned} \quad (17)$$

where i is the current iteration number, and n stands for the maximal adversarial optimization iteration in the single diffusion step. We set $\xi_1 = 0.1$, $\xi_2 = 0.01$, and set $m = 3$ and $M = 10$ for diversified

Table 6: Analysis of the different setting of ξ

	ξ_1	0.15	0.15	0.15	0.1	0.1	0.1	0.05	0.05	0.05
	ξ_2	0.02	0.01	0.005	0.02	0.01	0.005	0.02	0.01	0.005
ResNet50 ACC↓		0.0224	0.0192	0.0214	0.0246	0.0246	0.0246	0.0224	0.0256	0.0246
ViT-B/16 ACC↓		0.2949	0.2885	0.2885	0.2853	0.2724	0.2714	0.2404	0.2382	0.2511

strategy in different sampling steps, *i.e.*, the initial and final steps are set with a higher maximal iteration number ($n = M$).

In Table 6, we present a parameter analysis of ξ_1 and ξ_2 using ResNet50 as the surrogate model on the Abstract Label Evasion Task. In white-box attacks, all tested threshold values achieve strong performance (accuracy after attack < 3%). In black-box attacks, smaller ξ_1 values yield higher transferability, because lower ξ_1 leads to more aggressive adversarial optimization in early denoising steps—optimization that tends to be more transferable (this is consistent with the design principle of ResAdv-DDIM, which aims to enhance early-step optimization). However, this also increases computational cost. A practical deployment could balance effectiveness and efficiency.

Attack Loss Construction For the classification task, given the incorrect label set A_{Text} , we implement the loss as:

$$\mathcal{L}_{\text{ATK}} := -\text{Logits}[\hat{L}_{\text{Tar}}]. \quad (18)$$

Where Logits is the model’s prediction, and \hat{L}_{Tar} is the estimated highest confidence label in the incorrect label set A_{Text} . We maximize \mathcal{L}_{ATK} to perform the attack. For the current denoising step x_t , the logits are estimated based on the residual approximation function and the surrogate model, *i.e.*, $\text{Logits} = \mathcal{M}(g_{\theta}^{(t, \text{Text})}(x_t))$. To further eliminate guidance fluctuation, we adapt the fixed \hat{L}_{tar} after its initialization in the initial steps of attack optimization (*target label update delay*). Note that our adversarial attack method is not designed for the classification task only, and the loss construction could be substituted regarding the specific adversarial example generation task.

The overall SemanticAE generation algorithm is constructed by further integrating *global perturbation constraints, momentum gradient optimization, target label update delay* (after $t < t_k = 0.4T$) through the bi-level optimization. The pseudo-code is shown in Algorithm 1. The default configuration for other parameters is coherent with related baselines, *i.e.*, $\beta = 0.5$, $s = 0.7$, $T = 100$.

Algorithm 1: ResAdv-DDIM

Require: Input: Text, \mathcal{M} , A_{Text} , f , ϵ , optimization parameters T, K, t_s, t_k, β, s
 Init $x_T \sim \mathcal{N}(0, \mathbf{I})$, $v_x \leftarrow 0$, $x'_T \leftarrow x_T$.

```

for  $t = T, \dots, 1$  do
    if  $t \leq t_s$  then
        for  $i = 1, 2, \dots, n$  do
            if  $\neg \text{PerformOptimize (Eq. 17)}$  then
                break                                /* Adaptive Iteration */
            if  $t = t_s \vee t < t_k$  then
                 $\hat{L}_{\text{Tar}} = \arg \max_{L \in A_{\text{Text}}} \mathcal{M}(g_{\theta}^{(t, \text{Text})}(x_t)) \text{logits}[L]$  /* Target Update */
                 $v_x \leftarrow \beta v_x - (1 - \beta) \nabla_{x_t} \mathcal{L}_{\text{ATK}}$                                 /* Momentum Updates */
                 $x_t \leftarrow x_t + s * v_x$                                 /* Adversarial Optimization */
            if  $t = t_s$  then
                 $x'_{t-1} \leftarrow x_{t-1}$                                 /* Determine Exemplar on Step  $t_s$  */
                 $x_{t-1} \leftarrow f_{\theta, 1}^{(t, \text{Text})}(x_t)(x_t)$                                 /* DDIM Sampling Step for  $x_{\text{adv}}$  */
             $x'_{t-1} \leftarrow f_{\theta, 1}^{(t, \text{Text})}(x_t)(x'_t)$                                 /* DDIM Sampling Step for  $x_{\text{exemplar}}$  */
             $x_{t-1} \leftarrow x_{t-1} + \min\{\epsilon, \|x'_{t-1} - x_{t-1}\|_2\} \cdot \frac{x'_{t-1} - x_{t-1}}{\|x'_{t-1} - x_{t-1}\|_2}$ . /* Semantic Constraint */
    return  $x_0, x'_0$ .

```

By counting, the step number of the forward or backward process of diffusion-UNet is less than $(2mt_s + 8M)K + t_s + T$. The maximal memory cost of the backward process is $\lceil \frac{t_s}{\lfloor T/K \rfloor} \rceil \times$ the parameters in the feature maps of the diffusion-UNet, combined with other modules, including the surrogate model, the input transformation, and the VAE-decoder in the optimization pipeline. In practice, the time consumption is significantly lower than the upper bound due to the adaptive iteration mechanism, and the lower bound is characterized by:

Proposition B.1 (Lower Bound of the Diffusion Step). *For ResAdv-DDIM with the total sampling step T , the parameter of approximate iterations in g as K , the timestep of start adversarial optimization as t_s , the lower bound of the diffusion step is:*

$$\left(\frac{K \cdot t_s}{T} + 3 \right) \cdot \frac{t_s}{2} + T, \quad (19)$$

if K and t_s are set as $K|T$ and $\frac{T}{K}|t_s$.

Proof. From $K|T$ and $\frac{T}{K}|t_s$, let $T = K \cdot d$ where $d \in \mathbb{Z}^+$, $t_s = d \cdot m$ where $m \in \mathbb{Z}^+$. Under the optimal implementation, the forward process in the early-stop judgment and the optimization step are reused. Consider the case of always early-stopping, the total number of approximate iterations in g is $\sum_{t=1}^{t_s} \lceil \frac{t}{\lfloor T/K \rfloor} \rceil$, we have:

$$\sum_{t=1}^{t_s} \lceil \frac{t}{\lfloor T/K \rfloor} \rceil = \sum_{t=1}^{dm} \left\lceil \frac{t}{d} \right\rceil = \sum_{k=1}^m \sum_{t=1}^d \left\lceil \frac{kd + t - d}{d} \right\rceil = \sum_{k=1}^m (d \cdot k) = \frac{dm(m+1)}{2}. \quad (20)$$

By combining the total denoising step of x_{adv} and x_{exemplar} generation, the total step is:

$$\frac{dm(m+1)}{2} + t_s + T = \frac{t_s}{2} \cdot \left(\frac{K \cdot t_s}{T} + 1 \right) + t_s + T = \left(\frac{K \cdot t_s}{T} + 3 \right) \cdot \frac{t_s}{2} + T. \quad (21)$$

This completes the proof. \square

B.2 3D Object Generation

Base 3D Generation Method (Trellis) : Representing 3D data as matrices is inefficient and impractical due to computational problems. Our method is developed based on the *Trellis* model, which bridges the gap between 3D structure and the diffusion process with the structured latent (SLAT). The overall generation process could be formulated as:

$$\begin{aligned} z_0^{\text{slat}} &= \text{Diffusion Sampling}(\epsilon_{\text{slat}}, z_t^{\text{slat}}, \text{Text}), z_T^{\text{slat}} \sim \mathcal{N}(0, I) \\ \mathbf{pos} &= \text{Coords}(\mathcal{D}_{\text{slat}}(z_0^{\text{slat}})), \text{Coords} : \mathbb{R}^{b \times h \times w \times d} \rightarrow \mathbb{R}^{b \times n \times 2} \\ z_0 &= \text{Diffusion Sampling}(\epsilon, z_T, \mathbf{pos}, \text{Text}), z_T \sim \mathcal{N}(0, I) \\ \text{Model}_{\text{GS}} &= \mathcal{D}_{\text{GS}}(z_0, \mathbf{pos}), \mathcal{D}_{\text{GS}} : \{(z^i, pos^i)\}_{i=1}^L \rightarrow \{\{(x_i^k, c_i^k, s_i^k, \alpha_i^k, r_i^k)\}_{k=1}^K\}_{i=1}^L \\ x &= \text{Renderer}_{\text{GS}}(\text{Model}_{\text{GS}}, \text{Camera}). \end{aligned} \quad (22)$$

where z_0 and z_0^{slat} are latents sampled by the diffusion model, and are represented by sparse and dense tensors, respectively. $\mathcal{D}_{\text{slat}}$ is the coarse structure decoder, Coords transforms the voxel to point positions \mathbf{pos} , \mathcal{D}_{GS} is the refined structure decoder that decodes each vertex into multiple Gaussian points, and $\text{Renderer}_{\text{GS}}$ renders the Gaussian model to 2D images x with the camera parameter. Since the refined position is also encoded in z_0 , we implement the proposed ResAdv-DDIM with the noise estimation network ϵ that models the refined structure.

$\text{Renderer}_{\text{GS}}$ is the Gaussian renderer proposed in Gaussian splatting. Due to its advantage in optimization, we select this representation as the intermediate 3D data structure for gradient estimation. Specifically, the 3D Gaussian with center point p can be expressed as the Gaussian function:

$$G(p) = e^{-\frac{1}{2} p^T \Sigma^{-1} p} \quad (23)$$

To get a 2D projection image x from a 3D Gaussian point in a world coordinate with a viewing transformation W , the 2D covariance matrix Σ' as:

$$\Sigma' = JW\Sigma W' J', \quad (24)$$

where J is the Jacobi affine approximation matrix of the transformation matrix. To render the entire Gaussian model, the color of every pixel in the 2D image x is computed by blending N ordered points overlapping the pixel using the following equation:

$$C = \sum_{i \in N} c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (25)$$

where c_i is the color of each point evaluated by spherical harmonics (SH) color coefficients and (α_i) is determined by a 2D Gaussian with covariance Σ and optimizable per-point opacity.

Residual Approximation with EoT We adapt the expectation-over-transformation to bridge the 2D and 3D adversarial generation. By integrating the renderer, the residual approximation g and the adversarial optimization are represented as:

$$\begin{aligned} g_\theta^{(t, \text{Text})}(z_t, \mathbf{pos}, \text{Camera}) := \\ \text{Renderer}_{\text{GS}} \left(\mathcal{D}_{\text{GS}}(f_{\theta, \Delta T_1}^{(\Delta T_1, \text{Text})} \circ \dots \circ f_{\theta, \Delta T_k}^{(t, \text{Text})}(z_t, \mathbf{pos}), \mathbf{pos}), \text{Camera} \right) \\ z_{t-\Delta T} := f_{\theta, \Delta T}^{(t, \text{Text})} \left(\arg \max_{z_t} \mathbb{E}_{\text{Camera} \sim P_{\text{Cam}}} \left[\mathcal{L}_{\text{ATK}}(\mathcal{M}(g_\theta^{(t, \text{Text})}(z_t, \text{Camera}, \mathbf{pos}))) \right] \right) \end{aligned} \quad (26)$$

The camera settings are sampled based on the original configuration in *Trellis* framework that defines P_{Cam} , *i.e.*:

$$\begin{aligned} \Delta_{\text{yaw}} &\sim \mathcal{U} \left(-\frac{\pi}{4}, \frac{\pi}{4} \right), \\ \Delta_{\text{pitch}} &\sim \mathcal{U} \left(-\frac{\pi}{4}, \frac{\pi}{4} \right), \\ \mathbf{eye} &= 2 \cdot \begin{bmatrix} \sin(\theta_{\text{yaw}} + \Delta_{\text{yaw}}) \cos(\theta_{\text{pitch}} + \Delta_{\text{pitch}}) \\ \cos(\theta_{\text{yaw}} + \Delta_{\text{yaw}}) \cos(\theta_{\text{pitch}} + \Delta_{\text{pitch}}) \\ \sin(\theta_{\text{pitch}} + \Delta_{\text{pitch}}) \end{bmatrix}, \\ \mathbf{R} &= \text{LookAt}(\text{From}=\mathbf{eye}, \text{To}=(0, 0, 0), \text{UpAxis}=Z), \\ \text{Extrinsics} &= \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix}, \\ \text{Camera} &= [\text{Extrinsics}, \text{Intrinsics}_{fov=40^\circ}], \end{aligned} \quad (27)$$

where the eye position \mathbf{eye} is sampled on the sphere with $r = 2$, and the camera is toward the coordinate origin. Based on the practice of expectation-over-transformation(EoT) [26], the inner-loop optimization is performed by averaging the gradient from different sampled cameras:

$$\begin{aligned} \text{Cam}[1, 2, \dots, E] &\leftarrow \text{Sample}(P_{\text{Cam}}). \\ \text{grad} &= \frac{1}{E} \sum_{i=1}^E \nabla_{z_t} \mathcal{L}_{\text{ATK}}(\mathcal{M} \circ g_\theta^{(t, \text{Text})}(z_t, \mathbf{pos}, \text{Camera})) \\ v_z &\leftarrow \beta v_x + (1 - \beta) \cdot \text{grad} \\ z_t &\leftarrow z_t + s * v_z, \end{aligned} \quad (28)$$

where E is the step of EoT with the default value 1. For the optimization-related configurations, we maintained most of the parameters in 2D generation: $\epsilon = 10$, $K = 4$, $M = m = 30/E$, $\xi_1 = \xi_2 = 0.01$, $\beta = s = 0.5$ and $t_s = 0.75T$. The target label delay update mechanism is not applied. Diffusion-related configurations are coherent with the original pipeline, *i.e.*, $T = 50$.

Overall time Consumption We conducted a timing analysis on a single GPU using 100 labels from the ImageNet-Label dataset. The results are shown in Tabel 7. Gradient computation and residual approximation account for the majority of the computational cost. Meanwhile, *EoT* sampling is implemented via multi-view rendering using the Gaussian splatting renderer, while the rendering parameter computation does not occupy CUDA execution time. The variation in generation time arises from the adaptive number of optimization steps and differences in the size of sparse tensors in the Trellis framework. Specifically, in single iteration, fluctuation in computation speed is primarily caused by inconsistent data sizes in the Sparse Tensor. For the entire optimization, the fluctuation in

Table 7: GPU Time Consumption for Adversarial Optimization

Process	Single Iteration (ms)	Entire Optimization (ms)
Denoise Sampling	123.321 ± 36.878	2994.933 ± 269.841
Residual Approximation	147.256 ± 76.243	7040.238 ± 8915.112
Gaussian Decoding	21.547 ± 11.868	1029.394 ± 1185.636
Gaussian Rendering	28.556 ± 7.706	1364.221 ± 1732.996
Surrogate Model	6.777 ± 0.745	323.782 ± 423.503
Backward Process	208.471 ± 78.039	6261.577 ± 12119.871

entity, physical_entity, object, ungulate, whole, animal, organism, vertebrate, vascular_plant, instrumentality, mammal, placental, carnivore, vehicle, herb, self-propelled_vehicle, amphibian, canine, domestic_animal, electronic_equipment, device, container, covering, conveyance, commodity, monkey, abstraction, consumer_goods, structure, invertebrate, artifact, ruminant, invertebrate, matter, wheeled_vehicle, arthropod, causal_agent, reptile, equipment, implement, even-toed_ungulate, garment, game, diapsid, primate, protective_covering, relation, restraint, ape, natural_object, psychological_feature, geologicalFormation, attribute, starches, obstruction, aquatic_vertebrate, old_world_monkey, process, barrier, new_world_monkey, substance, communication, establishment, feline, tool, clothing, food, solid, piece_of_cloth, brass, screen, shelter, grouse, machine, vessel, craft, arachnid, fabric, durables, thing, place_of_business, reproductive_structure, plant, event, material, fastener, woody_plant, measure, home_appliance, mechanism, seafood, cognition, part, organ, group, game_equipment, shape, rodent, military_vehicle, area, mechanical_device, substance, nutrient, amphibian, salamander, support, produce, natural_elevation, mollusk, crustacean, aquatic_mammal, signal, indefinite_quantity, act, public_transport, hand_tool, medium, box, state, kitchen_appliance, edible_fruit, toiletry, shellfish, ware, utensil, fur, foodstuff, cloak, big_cat, footwear, ball, instrument, person, measuring_instrument, sports_equipment, stick, worker, insect, computer, lepidopterous_insect, vine

Figure 14: Overly polysemous tags filtered out

computation speed is also attributed to the adaptive optimization steps. Specifically, the average count of optimization steps for the inner optimization is: 47.810 (std=63.566). Attributed to the proposed early stopping mechanism, easier attack cases require fewer iterations and thus less time. The large variance in generation time reflects significant differences in attack difficulty across tasks, which highlights the effectiveness of our adaptive step design.

B.3 Construction and Analysis of Semantic-Abstracted Evaluation Task

Abstracted Label Set Construction As described in the main paper, we construct the coarse label set based on three steps: (1) hyponymic graph based on the hyponymic relation defined by *WordNet*, (2) select the proper abstraction level, (3) define the attack goal as the evasion attack on the abstracted label. Specifically, the hyponymic graph is defined as a directed graph $G = (V, E)$, $E \subset V \times V$, where each vertex V represents a word, each edge $e : v_1 \rightarrow v_2$ denotes that the word v_2 is a hypernym of v_1 . We denote the ImageNet label set as the vertex set $\mathbb{L} \subset V$, and construct the subgraph $G' = (V', E')$ as

$$V' = \text{TC}_G(L) := \left\{ v \in V \mid \exists u \in L, \exists k \in \mathbb{N}, (u \xrightarrow{k} v) \in E^+ \right\}$$

$$G' := (V', E') \quad \text{where} \quad E' = \{(u, v) \in E \mid u, v \in L'\} \quad (29)$$

Where TC denotes the transitive closure and G' is the induced subgraph. Next, we select the abstraction level. We first annotated and filtered out the following over-coarse labels from the vertex set V' , resulting in the label set L' . We also filter out polysemous tags listed in Figure 14.

Then, for each linear path, we select the node with the lowest height (or has more than one direct hyponym) as a candidate label, constraining the upper bound of the abstracting level, and eliminate

descendant labels of the candidate labels to constrain the lower bound of the abstracting level. These labels are represented as the AbstractedLabel. We construct the SemanticAE generation task by evading the AbstractedLabel, which is formulated as:

$$\begin{aligned}
& \text{find } x_{\text{adv}} \in \mathcal{S}(\text{Text}) \text{ s.t. } \mathcal{M}(x_{\text{adv}}) \in A_{\text{Text}}, \\
& \text{Text := "Realistic image of [AbstractedLabel], specifically, [label]",} \\
& A_{\text{Text}} := \{\text{label}_{\text{Adv}} \mid \text{AbstractedLabel} \notin \text{Ancestors}(\text{label}_{\text{Adv}})\}, \\
& \text{AbstractedLabel} \in \{c \in \mathbb{L}' \mid \exists l \in \mathbb{L} \text{ s.t. } c \in \text{Ancestors}(l) \wedge \text{Count}_{\text{Children}}(c) > 0\},
\end{aligned} \tag{30}$$

For simplicity, we use the term $c \in \text{Ancestors}(l)$ to denote there exists a path from l to c , and the term $\text{Count}_{\text{Children}}(c) > 0$ to denote the in-degree of $c > 0$. We select the abstracted label with more than 3 hyponym Image labels as the label set in the experiment. The abstracted labels and the corresponding hyponym imangenet labels are shown in Table 9.

Discussions on the ASR_{relative} metric. We design ASR_{relative} to match the goal of facilitating the red-teaming framework, as described in Section 3.1 of the main paper. The following proposition describes the relation between the ASR_{relative} evaluation metrics of SemanticAE generator and the application scenario of the multi-round reject-sampling-based data generation pipeline.

Proposition B.2 (ASR_{relative} characterize the upper-bound probability of the successful attack). *For any adversarial sampling algorithm K that generates the SemanticAE with $ASR_{\text{relative}} = p$ on the evaluated black-box model \mathcal{M}_T and the surrogate model \mathcal{M}_S , there exists an attack algorithm that achieves the successful attack with at least probability $p \cdot (1 - \epsilon)$ on \mathcal{M}_T , the average generation times less than $1/p_s$, and the maximal generation times $\left\lceil \frac{\log(\epsilon)}{\log(1-p_s)} \right\rceil$, for any $0 < \epsilon < 1$, if the following assumption holds true:*

1. **Non-Positive Attack:** For the sample generated from the instruction Text and the generation algorithm does not perform adversarial optimization towards misleading \mathcal{M} towards the labels corresponding to Text, the correct classification leads to semantic alignment, i.e., $\mathcal{M}(x) \in L_{\text{Text}} \rightarrow (x + \delta) \in \mathcal{S}(\text{Text})$, where L_{Text} is the correct label corresponding to Text and $L_{\text{Text}} = \overline{A_{\text{Text}}}$, and δ is a small perturbation with $\delta > \langle x_{\text{exemplar}}, x_{\text{adv}} \rangle$.
2. The sampling algorithm can generate a sample satisfying $\mathcal{M}(x) \in L_{\text{Text}}$ at least a probability of p_s by only accessing the instruction Text.
3. **Adjacency Assumption:** $P(\mathcal{M}_T(x_{\text{adv}}) \in L_{\text{Text}} \mid \mathcal{M}_T(x_{\text{exemplar}}) \in L_{\text{Text}}) > P(\mathcal{M}_T(x_{\text{adv}}) \in L_{\text{Text}} \mid \mathcal{M}_S(x_{\text{exemplar}}) \in L_{\text{Text}})$, where \mathcal{M}_S is the surrogate model, \mathcal{M}_T is the target model.
4. **Consistency Assumption:** $ASR_{\text{relative}} = P(\mathcal{M}_T(x) \in A_{\text{Text}} \mid \mathcal{M}_T(x) \in L_{\text{Text}})$ for the instruction Text given in the attack scenario.

Proof. We construct the Las Vegas-style sampling algorithm with the surrogate model \mathcal{M}_S . The sampling algorithm is re-run if $\mathcal{M}_S(x_{\text{exemplar}}) \notin L_{\text{Text}}$, until it reaches the upper-bound iteration $\left\lceil \frac{\log(\epsilon)}{\log(1-p_s)} \right\rceil$.

For the case of $\mathcal{M}_S(x_{\text{exemplar}}) \in L_{\text{Text}}$, based on assumption (1), $x_{\text{adv}} \in \mathcal{S}(\text{Text})$. Based on assumptions (3) and (4), we have

$$\begin{aligned}
P(\mathcal{M}_T(x_{\text{adv}}) \in A_{\text{Text}}) &= P(\mathcal{M}_T(x_{\text{adv}}) \in A_{\text{Text}} \mid \mathcal{M}_S(x_{\text{exemplar}}) \in L_{\text{Text}}) \\
&= 1 - P(\mathcal{M}_T(x_{\text{adv}}) \in L_{\text{Text}} \mid \mathcal{M}_S(x_{\text{exemplar}}) \in L_{\text{Text}}) \\
&> 1 - P(\mathcal{M}_T(x_{\text{adv}}) \in L_{\text{Text}} \mid \mathcal{M}_T(x_{\text{exemplar}}) \in T_{\text{Text}}) \quad (\text{by Assumption (3)}) \\
&= P(\mathcal{M}_T(x_{\text{adv}}) \in A_{\text{Text}} \mid \mathcal{M}_T(x_{\text{exemplar}}) \in T_{\text{Text}}) \\
&= ASR_{\text{relative}} = p \quad (\text{by Assumption (4)})
\end{aligned} \tag{31}$$

Therefore, with probability p , x_{adv} is a successful attack. Since $P(\mathcal{M}_S(x_{\text{exemplar}}) \in L_{\text{Text}}) > p_s$, the final success attack rate is :

$$P_{\text{final}} = p \cdot (1 - (1 - p_s)^{\left\lceil \frac{\log(\epsilon)}{\log(1-p_s)} \right\rceil}) > p \cdot (1 - (1 - p_s)^{\log_{1-p_s}\epsilon}) = p \cdot (1 - \epsilon) \tag{32}$$

The expected execution time of the sampling is:

$$\mathbb{E}[T] = \sum_{k=1}^{\lceil \frac{\log(\epsilon)}{\log(1-p_s)} \rceil} (1-p_s)^{k-1} p_s < \sum_{k=1}^{\infty} (1-p_s)^{k-1} p_s = p_s \cdot \frac{1}{(1-(1-p_s))^2} = \frac{1}{p_s}$$

This completes the proof. \square

We acknowledge that the evaluation might still not be adequate as it requires the assumption of *Non-Positive Attack*. However, the defect of the original non-reference evaluation is already shown in our experiments (detailed in Appendix D.4).

C Detailed Experiment Settings

C.1 2D Experiment Settings

Baseline Descriptions Our baseline method is constructed based on the categorization of the transfer attack method in Appendix A.3 and based on the discussions in Appendix A.2. Since the proposed module belongs to the intersection of the diffusion-based adversarial attack and the optimization methodology, we select the classical MI-FGSM [3], which is the base method of recent transfer attacks, and three diffusion-based adversarial optimization methods that are recently proposed and are suitable for SemanticAE , including AdvDiff [13], VENOM [14] and SD-NAE [12]. SD-NAE applies the gradient back-propagation optimization over the full diffusion steps, with the optimization formula as follows. It belongs to the first optimization paradigm discussed in Section 3.2

$$\max_{\text{TextEmbedding}} \mathcal{L}_{\text{ATK}}(\mathcal{M}(\underbrace{f_{\theta, \Delta T}^{\text{TextEmbedding}} \circ f_{\theta, \Delta t}^{\text{TextEmbedding}} \circ \dots \circ f_{\theta, \Delta t}^{\text{TextEmbedding}}(x_T)}_{T/\Delta T \text{ times}})), \quad (33)$$

AdvDiff and VENOM alter the latent embedding x_t during the diffusion denoising process, and belong to the second optimization paradigm discussed in Section 3.2. Compared to AdvDiff, VENOM applied the additional momentum mechanism and the early-stopping mechanism to the optimization process. Also, it tries to resample x_T based on the adversarial direction of x_0 when the optimization fails on the surrogate model. Our method, besides the technical improvement discussed in the main paper, adds a more fine-grained scheduling of optimization steps based on the new optimization paradigm. And we do not add the resample mechanism since it is not redundant if our method is used as a sampling module in an attack/data generation system that could perform reject sampling based on the application scenarios.

We integrate an input-transformation-based method (DeCoWA [41]) as the surrogate model to evaluate the collaboration capability of our method and other methods. We set the *num_warping*=2 for diffusion-based and our attacks, and *num_warping*=10 for MI-FGSM. The latter setting is consistent with the overall attack pipeline evaluated in the *DeCoWA* paper.

For the diffusion baselines, AdvDiff adapts the classifier guidance and takes the image class as the input of the diffusion model, while VENOM and SD-NAE adapt the classifier-free guidance. SD-NAE alters the selected text embedding by solving the maximization problem described in the main paper, and therefore $\max \mathcal{L}_{\text{ATK}}(\mathcal{M}(f_{\theta, \Delta T} \circ f_{\theta, \Delta t} \circ \dots \circ f_{\theta, \Delta t}(x_T)))$, Advdiff adapts the approximated optimization $'_{t-\Delta T} = f_{\theta, \Delta T}(\arg \max_{x'_t} \mathcal{L}_{\text{ATK}}(\mathcal{M}(x'_t)))$, and VENOM introduces conditional optimization and momentum mechanisms on it to stabilize the optimization. The diffusion model applied in VENOM and SD-NAE is *bguisard/stable-diffusion-nano-2-1*, and *latent-diffusion/cin256-v2* is applied for AdvDiff. We implemented our code based on the baselines VENOM and SD-NAE, and adapted the same diffusion model for consistent evaluation. In principle, our method can scale to larger models and is evaluated on the *Trellis* 3D generation model.

Surrogate and Target Models As described in the main paper, we adapt the target model set as $\mathcal{T} = \{\text{ResNet50} [45], \text{ViT-B/16} [46], \text{ConvNeXt-T} [47], \text{ResNet152}, \text{InceptionV3} [48], \text{Swin-Transformer-B} [49]\}$, and the surrogate model set as {ResNet50, ViT-B/16, ConvNext-T, ResNet50+DeCoWA}.

Loss Function We employ the same loss function as the original implementation for the baseline methods. For our model, we set the loss function in both 2D and 3D SemanticAE generation task as:

$$\mathcal{L}_{\text{Atk}}(\text{logits}, A_{\text{Text}}, L_{\text{Tar}}) = \text{LogSoftMax}(\text{logits})[L_{\text{Tar}}] - \frac{1}{|A_{\text{Text}}|} \sum_{i \in A_{\text{Text}}} \text{LogSoftMax}(\text{logits})[i], \quad (34)$$

Where L_{Tar} is the currently selected label (the label with the highest confidence in the set A_{Text} , and log softmax denotes $\log \frac{e^x}{\sum e^x}$.

Image Quality Assessment Metrics To evaluate semantic constraints, the pairwise semantic metric is proposed to measure the similarity between x_{exemplar} and x_{adv} , which defines as follows in the main body of our work:

$$\text{SemanticDiff}_{\mathcal{S}} = \langle x_{\text{exemplar}}, x_{\text{adv}} \rangle_{\mathcal{S}}, \quad (35)$$

\mathcal{S} is a visual similarity metric; we employed LPIPS and MS-SSIM for evaluation. Parameters of the evaluation metrics are adapted as common practice. For MS-SSIM, we adapt $kernelseize = 11$ and $\sigma = 1.5$. For LPIPS, we adapt *AlexNet* as the local feature extractor. For Clip_Q , we use the implementation of *piq* [77] and adapt *openai/clip-vit-base-patch16* as the image embedding extractor.

C.2 3D Experiment Settings

This section details the comprehensive framework established for the evaluation of 3D video generation models, encompassing dataset preparation, video synthesis methodologies for both benign and adversarial examples, and the metrics employed for performance assessment.

Dataset Preparation The ImageNet dataset, while extensive, contains labels with fine-grained semantic distinctions that can be challenging for text-to-3D video generation models to differentiate effectively. A coarse-graining procedure was applied to the Original Imagenet labels to address this.

Specifically, a predefined mapping, detailed in Table 9, was utilized to merge semantically similar labels. This process involved replacing the Original ImageNet labels with their corresponding Abstracted Labels. The resultant collection of these processed Abstracted Labels served as the prompt dataset for the subsequent video generation tasks. Let \mathbb{L} be the set of Original ImageNet labels and $\mathcal{T}_{\text{coarse}}$ be the set of Abstracted Labels. The mapping function $M : \mathbb{L} \rightarrow \mathcal{T}_{\text{coarse}}$ transforms each Original ImageNet label to its Abstracted Label. The set of prompts used for generation is $\mathcal{P} = \{t | t \in \mathcal{T}_{\text{coarse}}\}$.

3D Video Generation Two categories of video samples were generated: clean samples and adversarial examples. Clean video samples were synthesized using the TRELLIS[39] model. For each prompt $p \in \mathcal{P}$, a corresponding clean video V_{clean} was generated. Adversarial examples were generated based on the TRELLIS[39] model (version *TRELLIS-text-base*), employing a ResNet-50 model, pre-trained on ImageNet, as the surrogate model for guiding the adversarial attack. The generation of these adversarial examples was performed using our proposed methodology. During each generation instance, an abstracted text label $p \in \mathcal{P}$ was used as the input prompt. The target label for the attack was set to any Original ImageNet label $l_{\text{target}} \in \mathbb{L}$ such that $M(l_{\text{target}}) = p$. A constant perturbation strength, denoted as ϵ , was maintained at 10.0 across all adversarial generation processes. All videos are captured by the original rendering pipeline, *i.e.*, a camera surrounding the object as shown in the supplementary videos.

Evaluation Metrics The evaluation of the generated videos involved a frame-by-frame analysis using a pre-trained ResNet50 classifier.

For a given video V , consisting of N frames $\{f_1, f_2, \dots, f_N\}$, each frame f_i was individually classified by the ResNet-50 model. This yields a sequence of Original Imagenet labels, $c_i = \text{ResNet50}(f_i)$. Each c_i was then mapped to its Abstracted Label $t_i = M(c_i)$. The overall model prediction for the video V , denoted as P_V , was determined by the mode of these frame-level Abstracted Label predictions:

$$P_V = \text{mode}(\{t_1, t_2, \dots, t_N\})$$

where $\text{mode}(\cdot)$ returns the most frequently occurring element in the set.

A video V was deemed correctly classified if its model prediction P_V matched its ground truth Abstracted Label G_V . It was observed that, under certain parameter configurations (particularly for varying K), a minority of video generation attempts might fail. To ensure a rigorous and controlled comparison, a data curation step was implemented. The intersection of successfully generated videos across all conditions – clean samples (V_{clean}) and all sets of adversarial examples ($V_{adv}^{(0)}, V_{adv}^{(1)}, V_{adv}^{(2)}, V_{adv}^{(3)}, V_{adv}^{(4)}$) – was taken. Only videos present in this intersection were considered for the final evaluation. This ensures that performance metrics are calculated over an identical set of video instances, thus isolating the impact of the varied adversarial generation parameters.

Following the procedures outlined above, the Accuracy (ACC) and Attack Success Rate (ASR) were calculated. The specific mathematical formulations ASR are provided in the main body of this work.

D Detailed Results and Discussions

D.1 Transfer Attack Analysis

Experimental settings We select VENOM, MI-FGSM, SD-NAE, and AdvDiff as baseline methods, employing four surrogate models: ResNet50, DeCoWa, ConvNext-T, and ViT-B/16. Six target models are evaluated: ResNet50, ResNet152, ConvNext-T, ViT-B/16, Swin-B, and InceptionV3. For Abstracted Label tasks, our method adopts four different perturbations $\epsilon = \{2, 2.5, 3, 4\}$, while Original Imagenet label tasks use $\epsilon = \{1.5, 2, 2.5, 3\}$. For each surrogate-target model pair, adversarial examples are crafted using both our method and baselines. Evaluation metrics include Attack Success Rate (ASR), Accuracy (ACC) and LPIPS (lower values indicate better perceptual quality).

Data presentation Due to the inconsistency between different surrogate-target model pairs, we chose to present the experimental results using subplots. The x-axis represents the selected surrogate models, and the y-axis represents the attacked target models, with a total of 24 subplots. The experimental results are shown in Figure 20, 21, 22, 23. Figure 20, 22 display the relationship between LPIPS and ASR/ACC for examples generated by different methods in the abstracted label task. Figure 21, 23 show the relationship between LPIPS and ASR/ACC for examples generated by different methods in the original Imagenet label task. In the figures, data points of different colors represent different methods. The data points of our method under varying perturbation strengths are connected by lines, and its trend is fitted with a black dashed line. As a supplement, the numerical results with the standard deviation of the final metric over 6 random seeds are shown in Table 11 and Table 11, demonstrating the stability of the results.

Analysis Based on the observations from Figure 20, 21, the following conclusions can be derived:

- In our method, as the perturbation parameter ϵ gradually increases, ASR also rises, but the LPIPS increases as well. The relationship between ASR and the Natural logarithm of LPIPS essentially forms a straight line with a positive slope.
- The adversarial examples produced by SD-NAE (yellow markers) exhibit higher LPIPS in almost all cases, indicating that this method introduces more noticeable perturbations. However, compared to other baseline methods, SD-NAE does not always achieve a higher attack success rate.
- The adversarial examples produced by MI-FGSM (purple markers) have lower LPIPS than SD-NAE, but still significantly higher than our proposed attack method.
- The LPIPS of adversarial examples generated by VENOM (green markers) is comparable to our method. However, at similar LPIPS levels, our method achieves a higher ASR in most cases.
- AdvDiff (blue markers) produces adversarial examples with the lowest LPIPS, indicating better stealthiness. However, its attack success rate is significantly lower than all other methods.

Conclusion To summarize, in all 48 settings of this study, except for two cases (the first row, first column and first row, fourth column in Figure 20) where our method was slightly worse than the VENOM method, the remaining ASR-LPIPS curves of our method were all located above and to the

left of the comparison methods. In Figure 22, 23, all the methods, except for the VENOM method in the two subfigures of Figure 22(the first row and first column, and the first row and fourth column), are on the ACC-LPIPS curves of our method. This means our method achieved higher attack success rates at the same LPIPS levels. These results clearly demonstrate that our method outperforms the comparison methods in most cases (**46/48**). In addition, these exceptions both occurred in white-box attack scenarios, which were not our primary focus. VENOM achieved higher white-box accuracy by repeatedly resampling through rejection, while we treated this as a module separate from the Adversarial Sampling algorithm.

Therefore, our **InSUR** framework shows significant superiority in adversarial transferability.

D.2 Attack Robustness Analysis

Experimental settings : This experiment discusses the performance of our method compared to baseline methods when facing adversarial defenses. The defense methods used are JPEG and DiffPure. JPEG applies lossy compression to images (with a quality factor of 75 in this experiment), while DiffPure removes adversarial perturbations by feeding samples into a diffusion model for regeneration.

Data presentation : Figure 24, 25 show the results of our method and comparison methods on both defended and undefended models. Solid dots represent undefended results, while hollow dots represent defended results, with arrows indicating the impact of applying defenses.

Analysis :

- **MI-FGSM**: JPEG is a rule-based defense method, whereas DiffPure is a defense based on the in-manifold assumption. Due to the lack of in-manifold regularization, MI-FGSM-generated adversarial examples are less robust against DiffPure, showing a noticeable performance drop. This can be seen in Figure 24 vs. Figure 25.
- **AdvDiff**: As diffusion-based adversarial example generation method, it exhibits strong robustness against DiffPure (also diffusion-based). Figure 25 shows that their attack success rates (ASR) even increase after DiffPure defense, though they still underperform our method.
- **SD-NAE**: SD-NAE is also a diffusion-based adversarial example generation method, so it exhibits strong robustness against DiffPure. Besides, SD-NAE applies perturbations through text embeddings, which results in better performance in transfer attacks. However, this does not necessarily indicate an advantage, as the visualization results suggest that it may lead to uncontrollable semantic deviations. The high transfer attack success rate could be attributed to inherent changes in global semantics. Additionally, its optimization for white-box attacks is inadequate. For undefended white-box models, the success rate is lower than that of our method and VENOM.

Conclusion : Our method shows a decrease in ASR in most cases when facing JPEG and DiffPure defenses, but this effect is limited, and in the face of JPEG defense, our method is the only one able to keep the attack success rate ASR consistently above 80% when the target model is the same as the surrogate model (see Figure 24). When the DiffPure defense leads to a decrease in the attack success rate of our method and an increase in the attack success rate of the SD-NAE method, we are still able to ensure that at least one data point of our method outperforms the SD-NAE (e.g., the subplot in the first column of the fifth row of Figure 25).

In summary, our method enhances the optimization exploration capability for adversarial examples while maintaining In-Manifold Regularization(outperforms the comparison methods all cases (**48/48**)). This ensures that our approach remains effectively aggressive against defense methods.

D.3 On the Role of Residual-driven Attacking Direction Stabilization in 2D and 3D SemanticAE Applications

Boosting Multi-diffusion-step Regularized Adversarial Optimization Recall that, to solve the problem of **collaborating the adversarial optimization with the diffusion model for better transfer attacks and robust attacks**, we introduce the residual approximation in *ResAdv-DDIM*.

We further analyze it in the more refined evaluation tasks. The experiment is conducted with the surrogate model of **ViT-B/16** with 6 different target models on the abstracted label evasion task. The parameters of the maximal approximate iteration K and the $t_s/T = 0.25, 0.5, 0.75$ have been evaluated. The results are shown in Figure 15. Since altering the parameter changes the behavior of the semantic alignment, the semantic difference measurement is included, and the results are plotted in a figure. $K = 0$ indicates the setting without the residual approximation.

The figure shows that:

- When $t_s/T = 0.25$, both LPIPS and ASR are higher, and the introduction of the residual approximation lets the balance point between LPIPS and ASR offset, or more biased towards LPIPS optimization. This is due to (1) since the white box results, shown in the upper left figure, indicate that the successful ASR has already been achieved, and there is no need to improve the performance of adversarial optimization. Therefore, improvement is on LPIPS. (2) Since the residual approximation is not designed for regularizing transfer attacks, the transfer attack performance declines in ResNet and Inception models.
- When $t_s/T \in \{0.5, 0.75\}$, the result is different. Specifically, (1) the white-box results indicate that the adversarial optimization is non-optimal. (2) The residual approximation improves the performance of the white-box attacks and significantly improves the performance of transfer attacks. Improvements exist in both LPIPS and ASR.
- Simply increasing the step (the blue arrows) of undergoing diffusion steps of adversarial optimization may not improve adversarial transferability. However, it may decrease adversarial optimization performance under the same adaptive optimization mechanism, as the diffusion process may purify the adversarial optimization.
- Increasing t_s together with the residual approximation solves the collaboration problem between adversarial optimization and diffusion purification, leading to highly transferable SemanticAE.

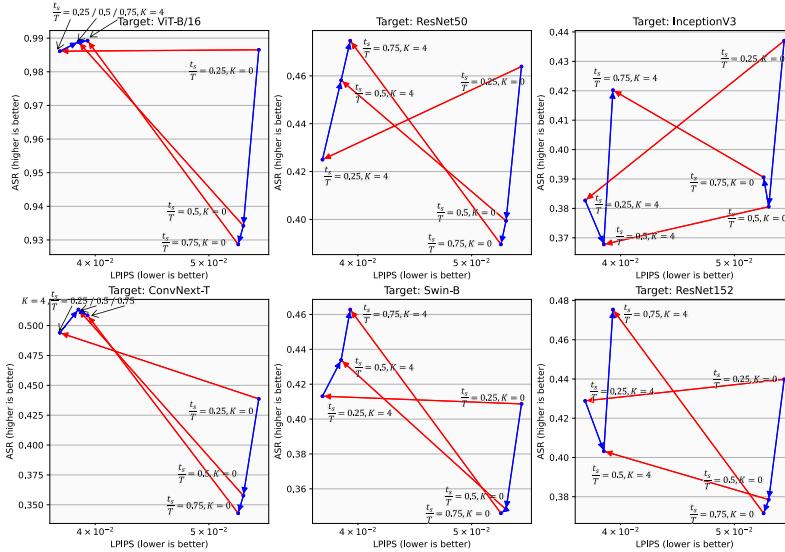


Figure 15: Parameter analysis with surrogate model ViT-B/16. Red arrows denote the performance before and after adding the residual approximation for $g(x)$. In each setting of $\frac{t_s}{T} \geq 0.5$, applying residual approximation achieves significant improvements (The upper left corner indicates a strictly superior direction). **This result is coherent with the design goal of ResAdv-DDIM.**

Collaboration with EoT Different from 2D optimization, the global gradient of 3D models cannot be obtained within a single iteration. Therefore, the EoT optimization is applied, accumulating the gradients across different perspectives. This further challenges the adversarial optimization capability. However, as shown in Table 8, with the residual approximation ($K > 0$), our method can collaborate well with *EoT* with different numbers of EoT steps, resulting in different gradient optimization step

sizes. The visualized comparison is shown in Figure 16 and the supplementary videos. The attack performance on different views is significantly higher than without residual approximation ($K = 0$), showing that the necessity of the residual approximation under diffusion+EoT generation pipeline.

Table 8: Comparison between residual approximation and expectation over transformation (EoT). The total iteration (EoT step * gradient descent step) is consistent. The gradient optimization stepsize is larger if the EoT step is higher.

EoT step	5	3	1	1	1	1	1
Residual approximation step (K)	4	4	4	3	2	1	0
ASR	0.922	0.913	0.922	0.912	0.912	0.902	0.451



Figure 16: 3D Visual Results Ablation.

D.4 On the Potential Adversarial Transferability to Semantic Evaluator

We analyze the character between and find the clue of the reference-free semantic evaluation metric being transferred in SemanticAE evaluation. Our experiment is based on hypothesis testing. Specifically, we evaluate the results from our method under the setting of Appendix D.1, which is 16 rounds of generation for each of the original ImageNet label evasion task and the abstracted label evasion task. Additionally, we evaluate the clip-score [78] metric with the settings:

$$Prompt = \begin{cases} \text{ImageNetLabel} & Task = \text{ImageNet} \\ \text{AbstractedLabel, ImageNetLabel} & Task = \text{Abstracted} \end{cases} \quad (36)$$

The backbone of the clip-score is *ViT-B/32*. Then we apply the linear regression on the clip-score, as the clip Semantic metric, and the LPIPS score on the factor of **whether the surrogate selects ViT-B/16**. The results is shown in Figure 17, Figure 18, and Figure 19. The following hypothesis is rejected with high confidence ($p < 0.02$):

Under the same generation settings, the CLIP semantic evaluation results are independent of the surrogate model selection.

Due to the high p-value associated with LPIPS, its correlation with surrogate model selection is relatively low. Although there are differences in model configurations and training tasks, both clip-score and the surrogate model adopt the **ViT** architecture. Therefore, we have reason to believe that the transfer attack has affected the evaluation of semantic similarity metrics.

We believe that the potential adversarial transferability to the deep-learning-based semantic evaluation model is difficult to tackle within the models, especially for the potential adversarial example that could perform weak-to-strong attacks. As discussed in Appendix B.3, our exemplar-based evaluation

task provides an alternative way to show the semantic alignment, avoids the requirement of non-referencing semantic evaluation, and directly shows the adversarial capability of the generated data.

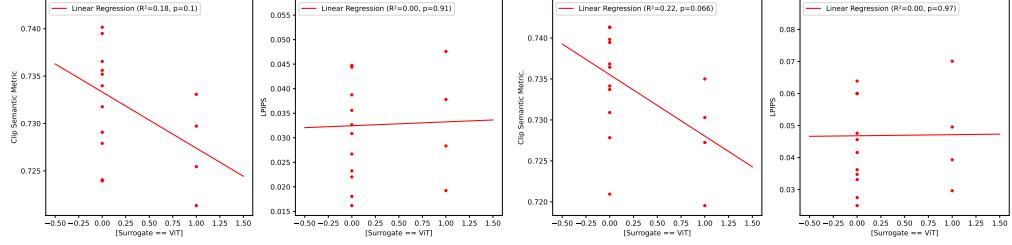


Figure 17: Results on the Original ImageNetFigure 18: Results on the Abstracted Label Evasion SemanticAEs

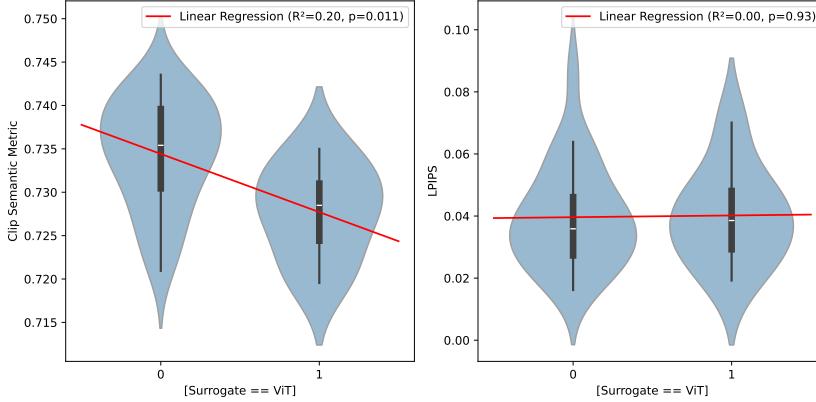


Figure 19: Combined distribution of the semantic metric on the ImageNet and abstracted label evasion SemanticAEs. With high confidence ($p < 0.02$), the CLIP semantic metric is affected by the attack transferability.

E Results Visualization

E.1 2D Visualized Results

To intuitively visualize the samples generated under the semantic constraints of original ImageNet label and abstracted label, we select SemanticAEs x_{adv} and corresponding nearby samples x_{exemplar} of seven original ImageNet labels (monarch, castle, goldfinch, aircraft carrier, vase, tiger, jellyfish) and six abstracted labels (aircraft, bag, beetle, bird, boat, fish) for visualization, as shown in Figure 26 and 27. The surrogate model is DeCoWa + ResNet. For our attack method. We set different parameters $\epsilon = \{1.5, 2, 2.5, 3\}$ and $\epsilon = \{2, 2.5, 3, 4\}$, respectively, for 2D ImageNet-label evasion attacks and 2D abstracted-label evasion attacks, which determines the strength of the semantic constraint. We also present the generation results of the baselines. For a single image I , we report confidence and mark the classification label y with different colors employing ResNet50 and ViT-B/16 as target models. In 2D ImageNet-label evasion attacks, Green is for the same classification label y and ImageNet label corresponding to the semantic constraint; otherwise, red. In 2D Abstracted-label evasion attacks, Green indicates that the classification label y belongs to the abstracted label corresponding to the semantic constraint. Therefore, a successful attack is defined as follows: for a pair of x_{adv} and x_{exemplar} , x_{exemplar} is classified correctly (green), while x_{adv} is classified incorrectly (red).

Analyzing 2D Visualized Results Based on the observations of Figure 26 and 27, the following conclusions can be derived:

- For our attack method, as the parameter ϵ gradually increases, signifying a progressive weakening of the semantic constraint strength, the number of successful attacks correspondingly increases. However, this also results in unnaturalness of SemanticAEs. For instance, the x_{adv} with the original ImageNet label "castle" has a castle with a blurred top when $\epsilon = 3$, and the x_{adv} with the abstracted label "ship" exhibits a blurry mast when $\epsilon = 4$.
- SemanticAEs generated by SD-NAE disturb more global semantics, which leads to semantic drift and dissimilarity between the x_{exemplar} and x_{adv} . It is consistent with the results shown in Appendix D.1: SD-NAE has a lower MS-SSIM score and a higher LPIPS score, with the same surrogate model, compared with other attack methods.
- As for MI-FGSM, noticeable noise can be observed in the background area of the SemanticAEs, which are not natural. In the main paper, MI-FGSM corresponds to a lower CLIP-QAI score related to noisiness.
- AdvDiff generates adversarial examples with artifacts at the edges in ImageNet-label evasion attacks, such as the x_{adv} of label "castle" and "aircraft carrier", which is a manifestation of low image quality. In abstracted-label evasion attacks, SemanticAEs hardly attack target models successfully.
- The adversarial examples of VENOM are natural. However, they struggle to effectively attack the ViT-B/16 model.

In conclusion, our method generates local in-manifold patterns to achieve strong attacks.

E.2 3D Visualized Results

Image Visualization To qualitatively substantiate the efficacy of our proposed methodology, we conducted a visual comparison of the generated video samples. Initially, ten unique labels were randomly selected from the complete set of Abstracted Labels. Subsequently, for each of these selected labels, the corresponding clean video sample and the adversarial video sample generated with $K = 4$ were chosen. From each of these videos, five frames were extracted at equidistant intervals. The visual results of these extracted frames are presented in Figure 28. A comparative analysis of each pair of clean and adversarial examples reveals that the adversarial counterparts maintain a high degree of visual similarity to the benign videos. Despite this perceptual resemblance, the adversarial examples demonstrate a high probability of inducing misclassification by the ResNet-50 model, thereby underscoring the effectiveness of our approach in generating robust yet inconspicuous adversarial attacks.

Furthermore, we performed a comparative analysis of adversarial examples generated under varying K parameter settings to demonstrate the rationale behind our parameter selection. Specifically, we selected the same video instance from the clean samples, the adversarial examples generated with $K = 0$, and those generated with $K = 4$. A single frame was extracted from each of these three video versions for visualization, as depicted in Figure 16. It is observable from the figure that when $K = 0$, the classification outcome for the adversarial example paradoxically improves compared to the clean sample, signifying a failure of the adversarial attack. Conversely, for $K = 4$, the adversarial example exhibits visual characteristics closely resembling those of the $K = 0$ sample. However, its efficacy in deceiving the classifier is significantly enhanced. This comparative visualization corroborates the superiority of our chosen parameter configuration ($K = 4$) in achieving a strong attack effect while preserving visual quality.

Video Visualization We performed a comparative analysis of adversarial examples generated under varying K parameter settings to demonstrate the rationale behind our parameter selection. Specifically, we selected three video instances, identified by their primary content as the *forklift* video, the *llama* video, and the *volcano* video. For each of these videos, we considered the clean sample, the adversarial example generated with $K = 0$ (without residual approximation), and that generated with $K = 4$.

It is observable from our quantitative analysis that the outcomes for $K = 0$ varied across the different videos. For the *forklift* video, the classification accuracy of the adversarial example generated with $K = 0$ paradoxically improved compared to its clean sample, signifying a failure of the adversarial attack under this specific setting. In contrast, for both the *llama* and *volcano* videos, the adversarial examples generated with $K = 0$ achieved lower classification accuracies than their respective clean samples, indicating some attack effect. However, their accuracies were still notably higher than those

of the adversarial examples generated with $K = 4$. This demonstrates that for the *llama* and *volcano* videos, while $K = 0$ initiated an attack, its deceiving capability was considerably weaker than that of $K = 4$.

Conversely, for $K = 4$, the adversarial examples for all three videos (*forklift*, *llama*, and *volcano*) exhibit visual characteristics closely resembling those of the $K = 0$ samples. However, their efficacy in deceiving the classifier is significantly enhanced across all instances. This comparative visualization and analysis corroborate the superiority of our chosen parameter configuration ($K = 4$) in achieving a strong attack effect while preserving visual quality.

F Boarder Impacts

While our goal is to catalyze the development of the red-teaming framework and trustworthy AI, we acknowledge that the proposed technology might be misused, including: (1) extending the proposed transfer attack improvement methods to jailbreak multi-modal LLMs. (2) extending the proposed 3D attack methods to generate physical adversarial examples to attack biometric authentication systems. However, our framework is not directly designed for these scenarios and requires further integration.

To protect from potential attacks in applications, we suggest developing the following closed-loop framework as a complement to traditional defense methods tested in Appendix D:

- Collect data generated by the proposed **InSUR** framework.
- Annotate the data with human feedback or rule-based models.
- Improve the alignment of the multi-modal models through fine-tuning on the dataset.

As a tool for data generation, we believe our framework is more beneficial for the model holder. For responsibility, we will release the code of **InSUR** framework *after* the paper is published for reference.

Table 9: Abstracted Label Mapping from ImageNet Numerical IDs

Abstracted Label	Original ImageNet Label IDs
dog	151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268
bird	7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146
musical_instrument	401, 402, 420, 432, 486, 494, 513, 541, 546, 558, 566, 577, 579, 593, 594, 641, 642, 683, 684, 687, 699, 776, 822, 875, 881, 889
furnishing	423, 431, 453, 493, 495, 516, 520, 526, 532, 548, 553, 559, 564, 648, 703, 736, 741, 765, 794, 831, 846, 854, 857, 861, 894
motor_vehicle	407, 408, 436, 468, 511, 555, 569, 573, 575, 609, 627, 656, 661, 665, 675, 717, 734, 751, 803, 817, 864, 867
snake	52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68
fish	0, 1, 2, 3, 4, 5, 6, 389, 390, 391, 392, 393, 394, 395, 396, 397
building	410, 425, 449, 497, 498, 580, 598, 624, 663, 668, 698, 727, 762, 832
saurian	38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48
ball	429, 430, 522, 574, 722, 747, 768, 805, 852, 890
headdress	433, 439, 452, 515, 518, 560, 667, 793, 808
beetle	300, 301, 302, 303, 304, 305, 306, 307
timepiece	409, 530, 531, 604, 704, 826, 835, 892
shop	415, 424, 454, 467, 509, 788, 860, 865
dish	925, 926, 933, 934, 962, 963, 964, 965
cat	281, 282, 283, 284, 285, 286, 287
weapon	413, 456, 471, 657, 744, 763, 764
bottle	440, 720, 737, 898, 899, 901, 907
fungus	991, 992, 993, 994, 995, 996, 997
spider	72, 73, 74, 75, 76, 77
ship	403, 510, 628, 724, 833, 913
boat	472, 554, 576, 625, 814, 914
turtle	33, 34, 35, 36, 37
bag	414, 636, 728, 748, 797
housing	500, 660, 663, 698, 915
crab	118, 119, 120, 121
wolf	269, 270, 271, 272
fox	277, 278, 279, 280
bear	294, 295, 296, 297
aircraft	404, 405, 417, 895
armor	465, 490, 524, 787
wheel	479, 694, 723, 739
fence	489, 716, 825, 912
personal_computer	527, 590, 620, 681
roof	538, 853, 858, 884
overgarment	568, 617, 735, 869
skirt	601, 655, 689, 775
bread	930, 931, 932, 962
squash	939, 940, 941, 942

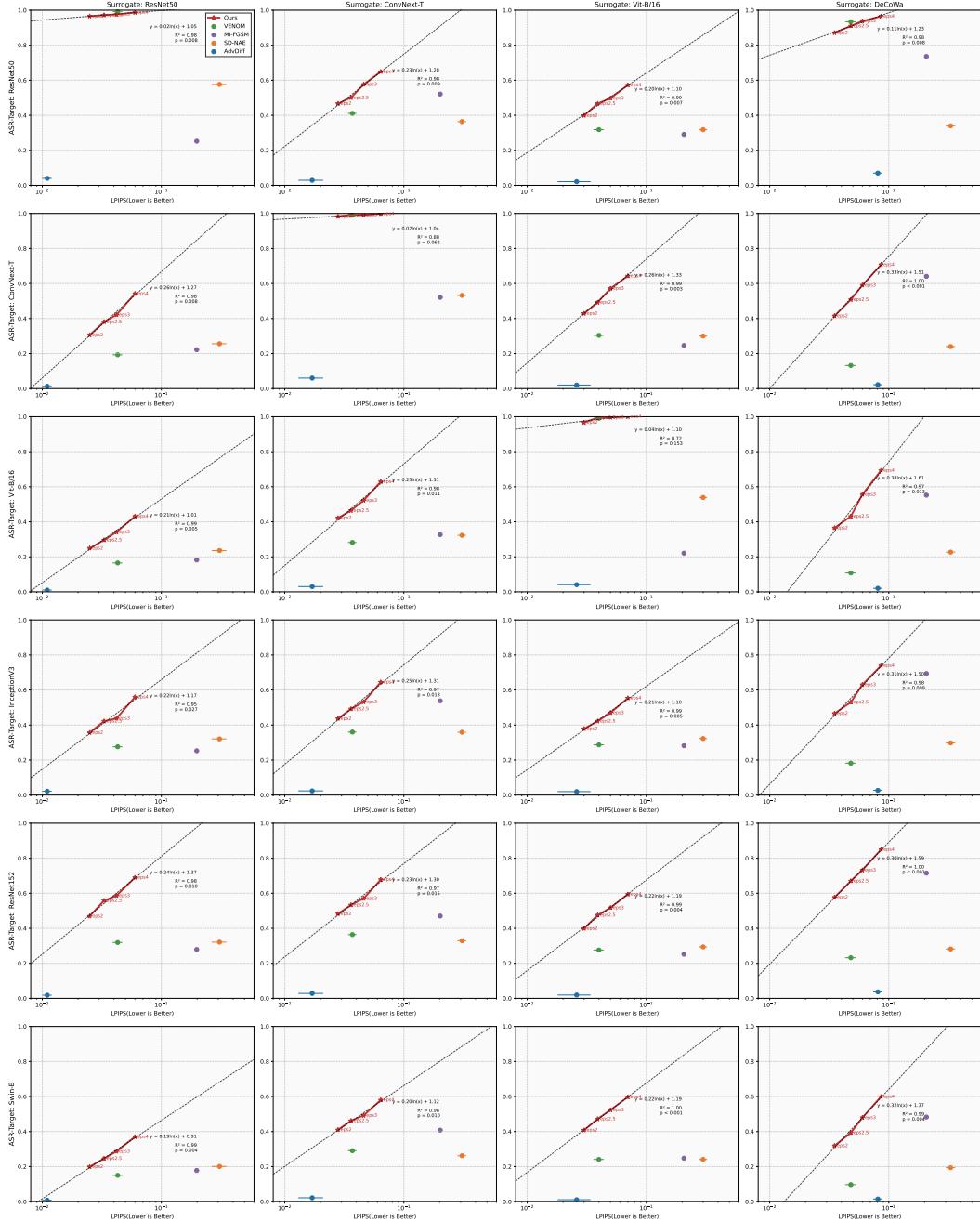


Figure 20: ASR of Abstracted Label

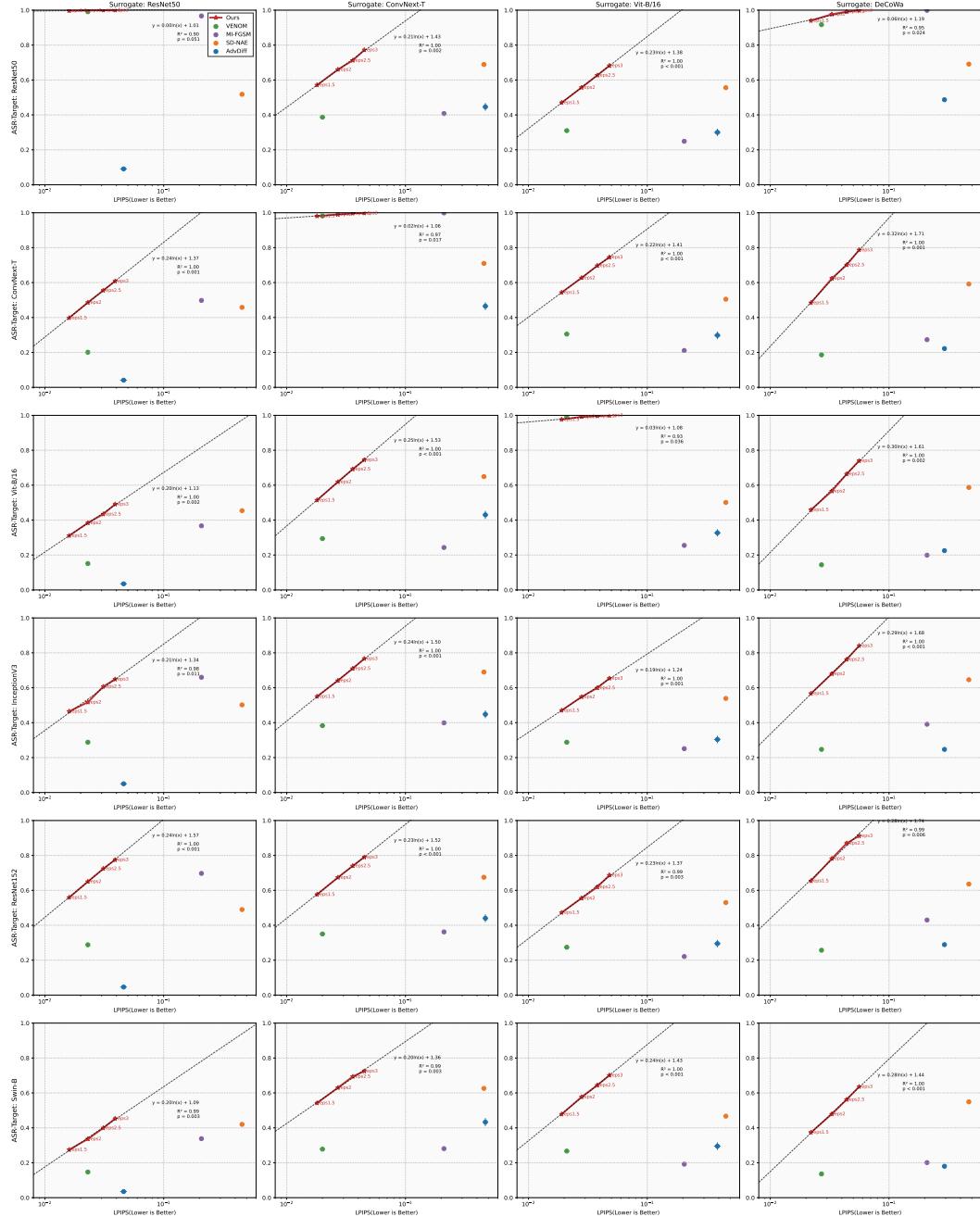


Figure 21: ASR of Original Imagenet label

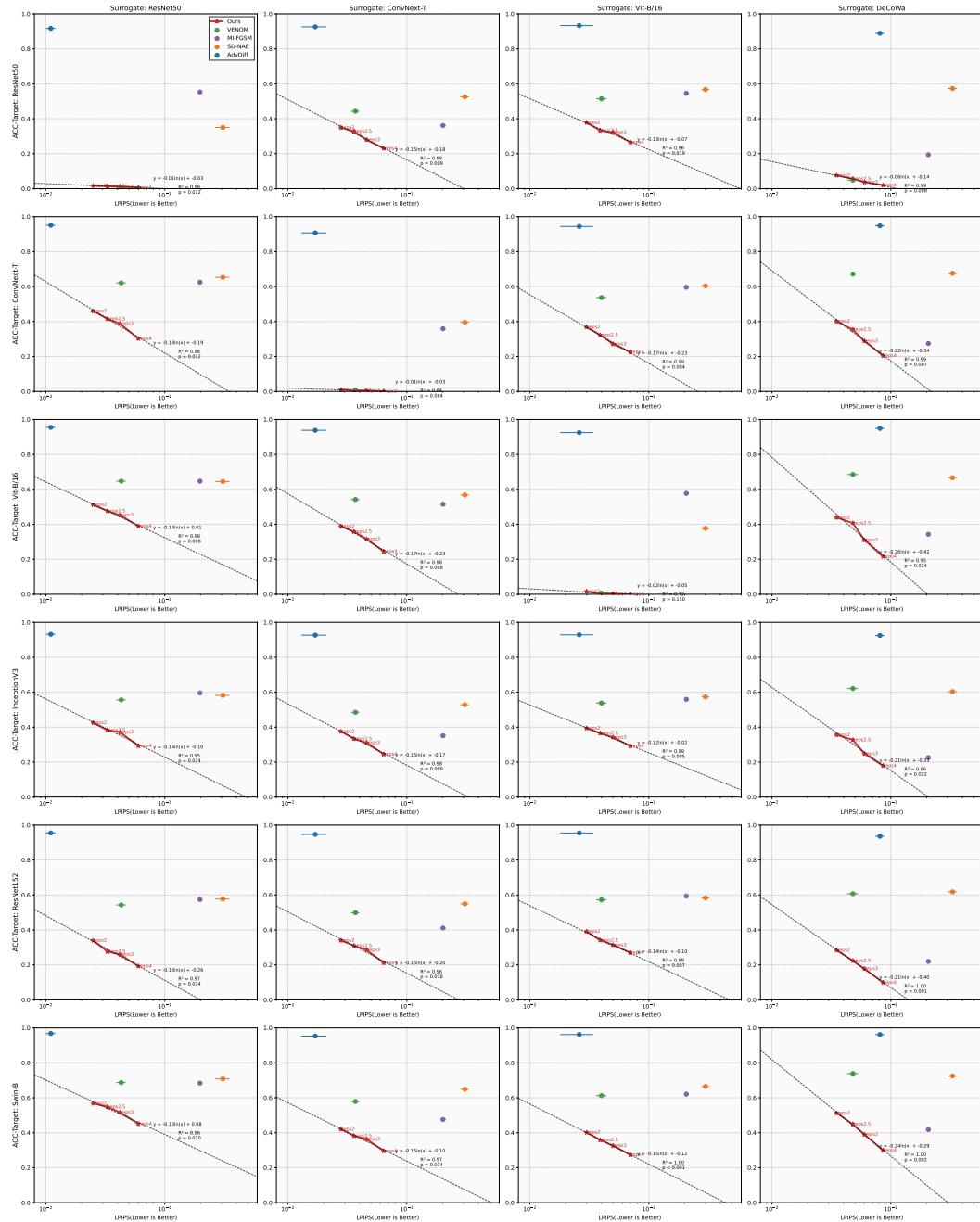


Figure 22: ACC of Abstracted Label

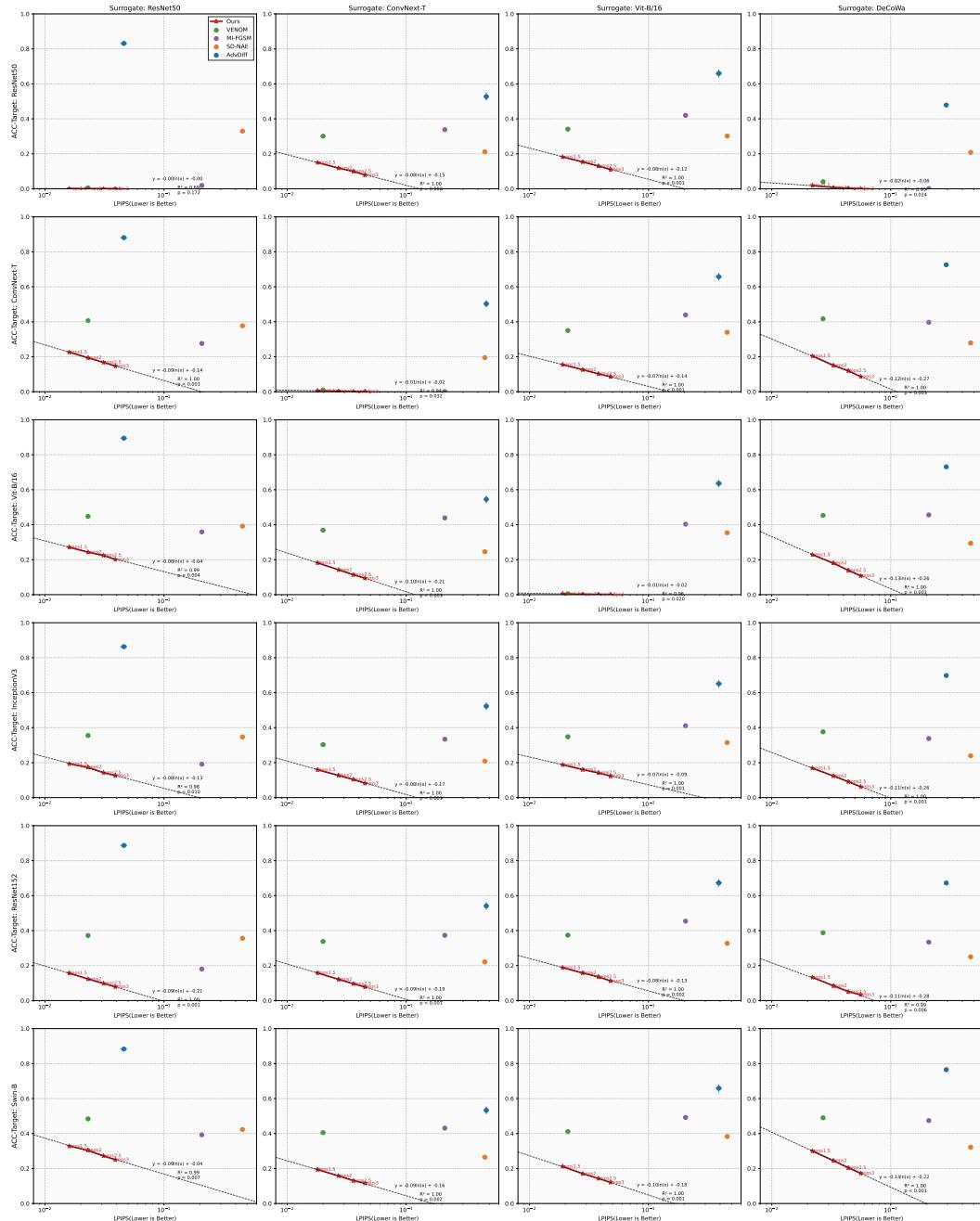


Figure 23: ACC of Original Imagenet label

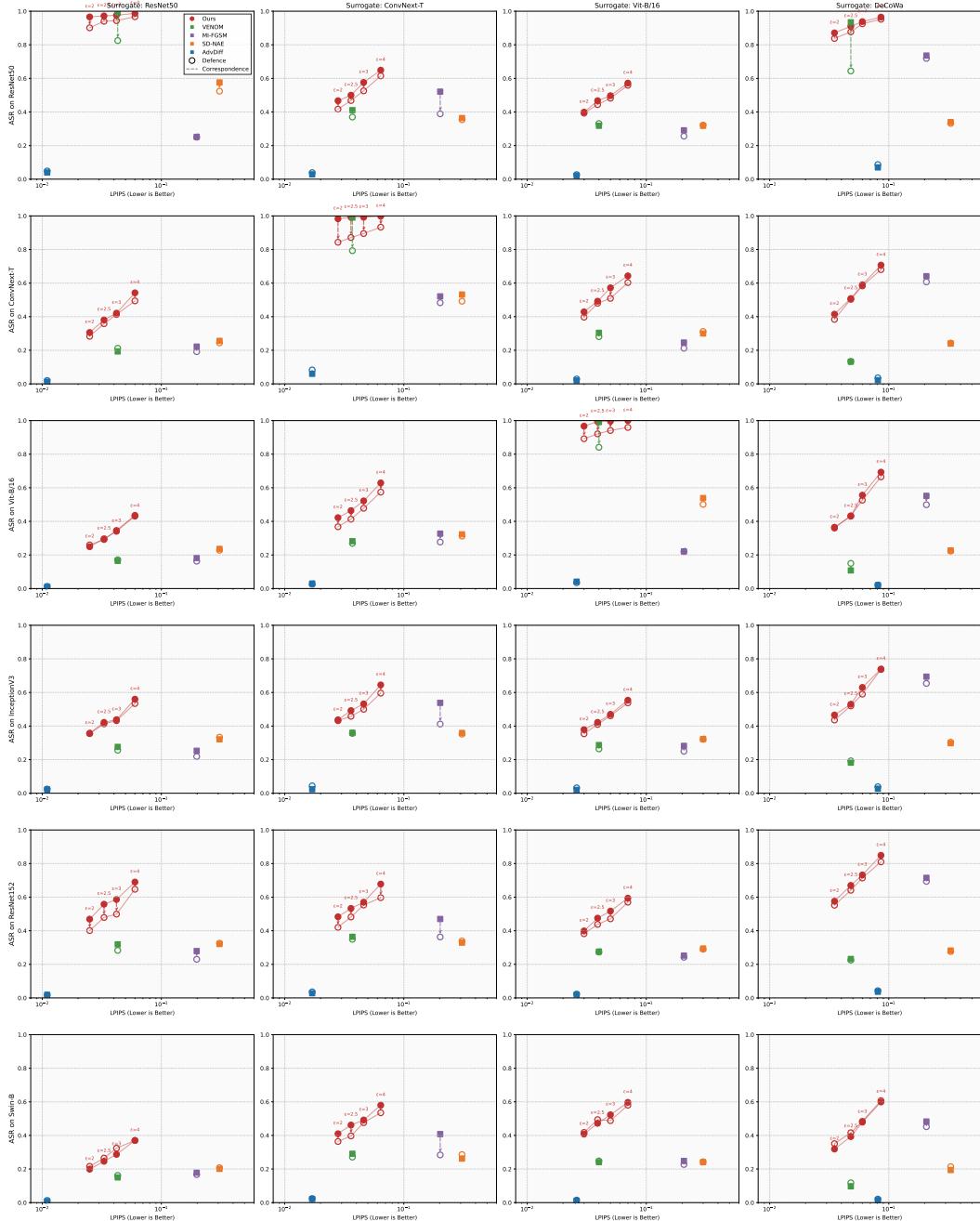


Figure 24: ASR on JPEG Defense

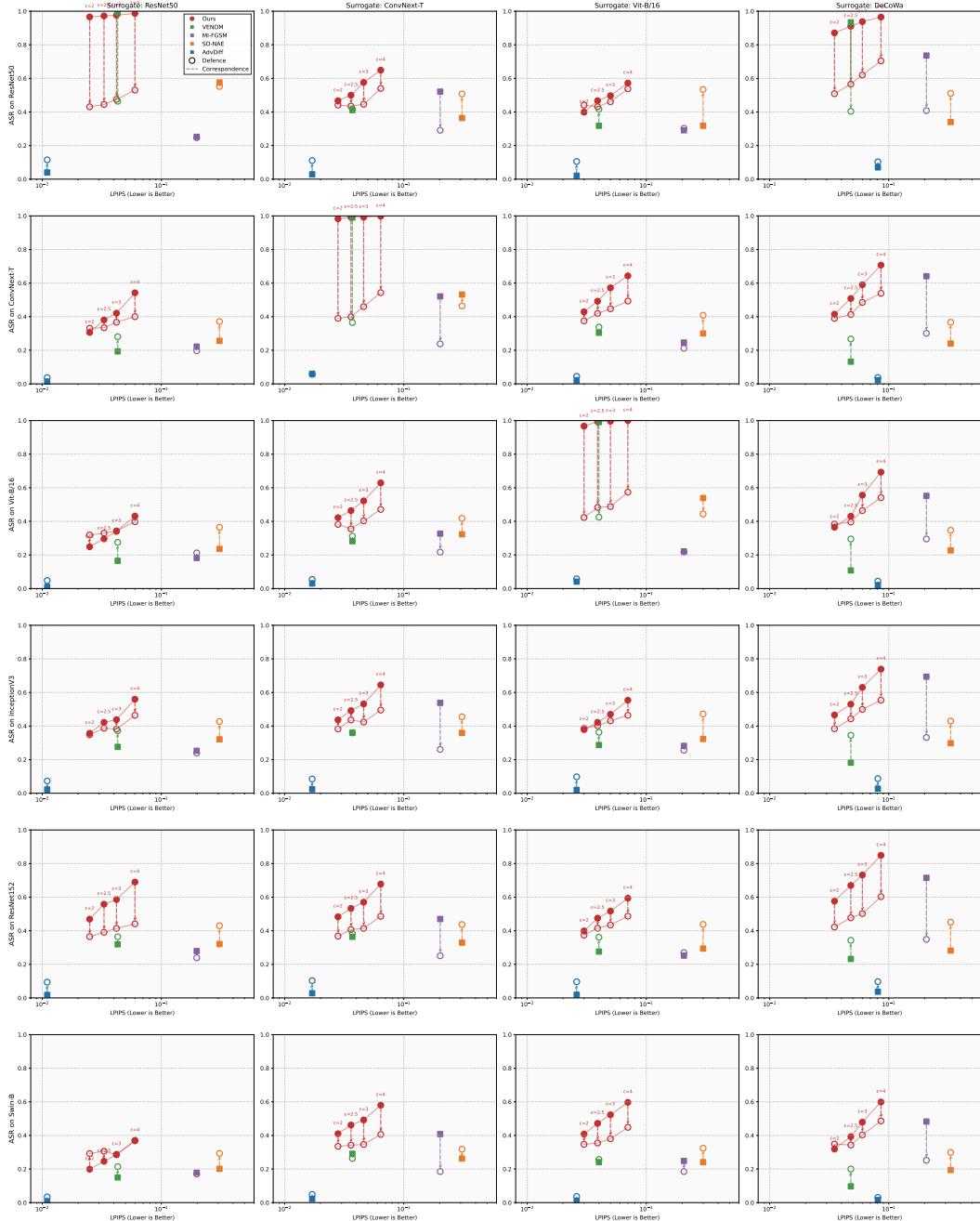


Figure 25: ASR on DiffPure Defense

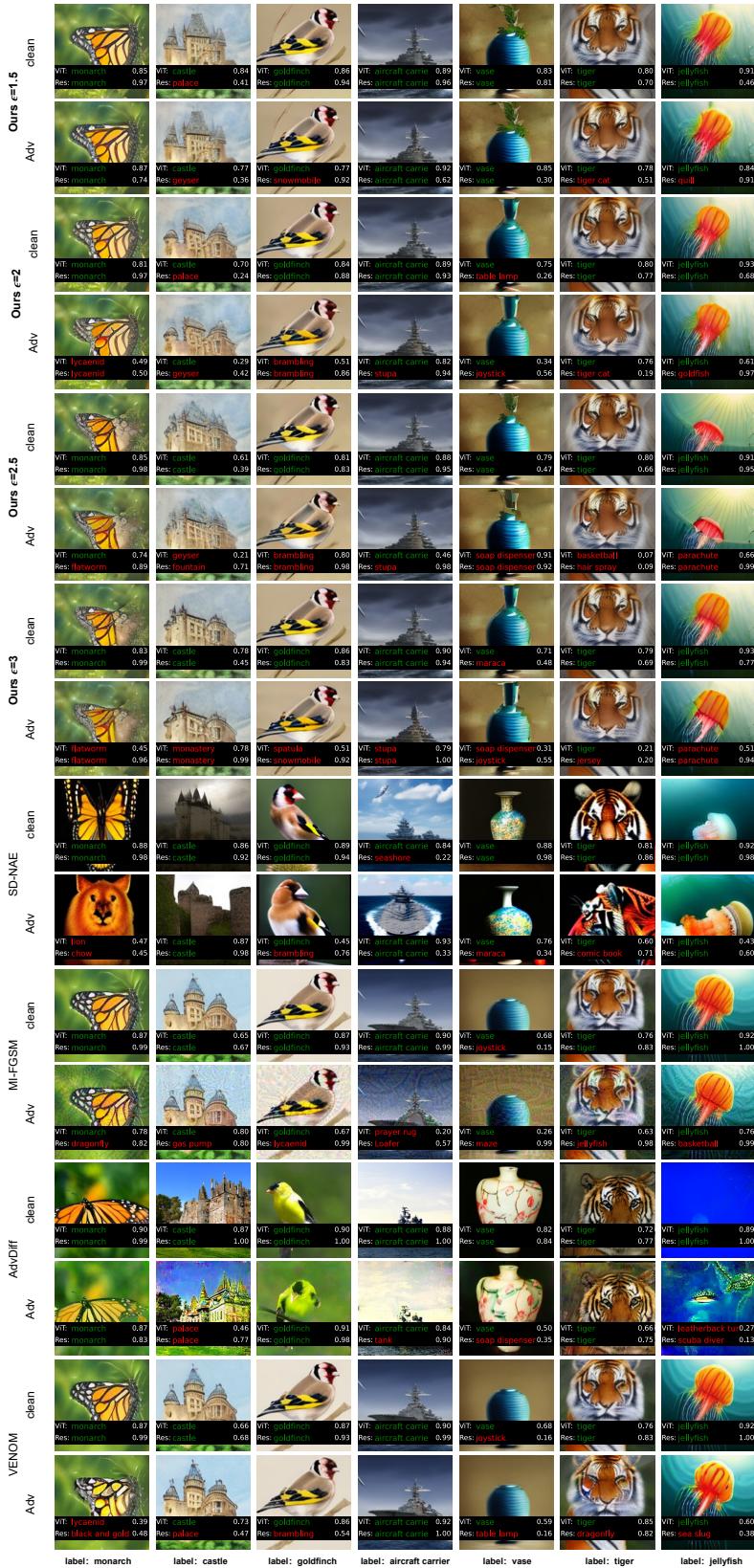


Figure 26: Visualization of 2D ImageNet-label Evasion Attacks. Surrogate model is ResNet50+DeCoWA.



Figure 27: Visualization of 2D Abstracted-label Evasion Attacks. Surrogate model is ResNet50+DeCoWA.

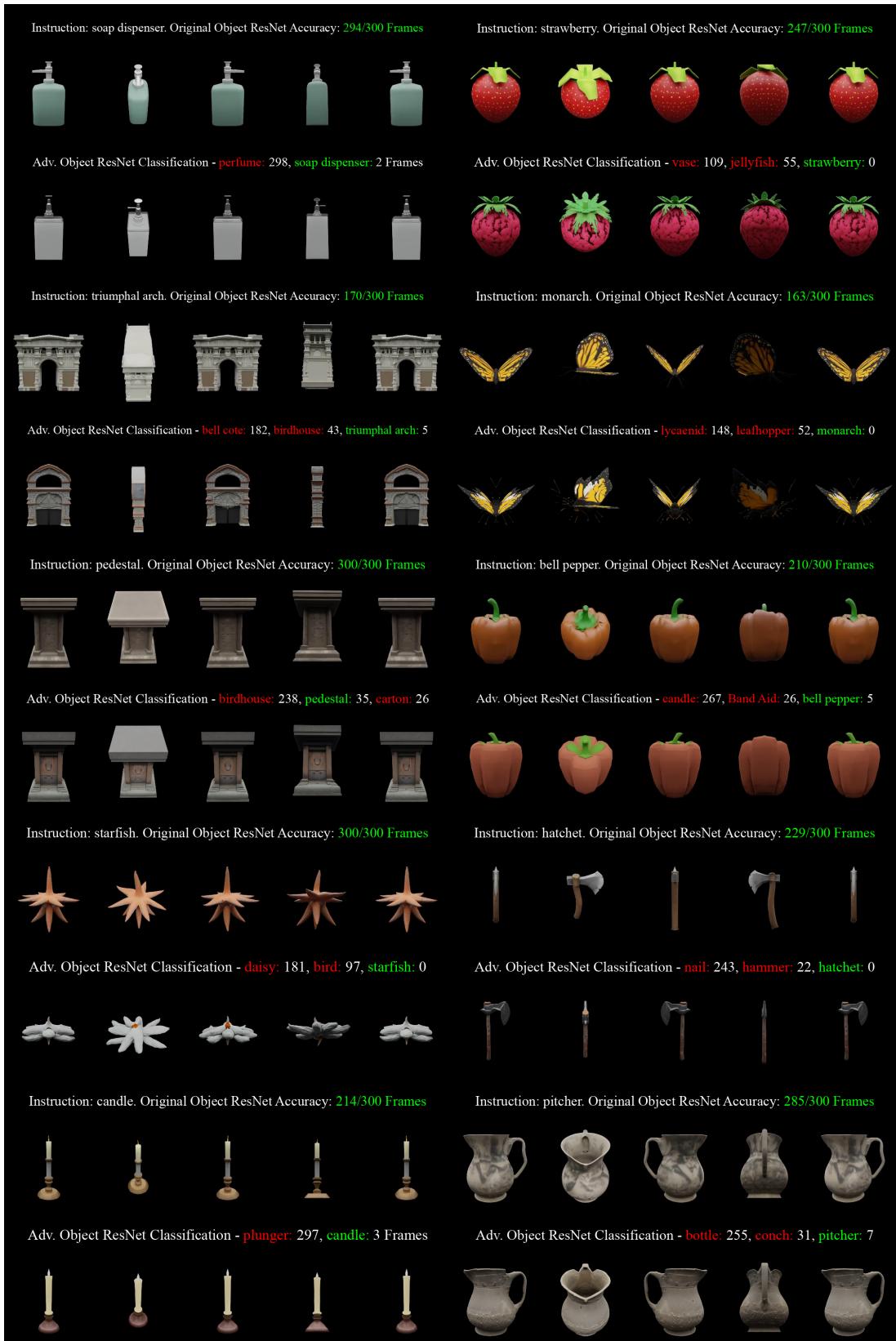


Figure 28: 3D Visual Results. Surrogate Model is ResNet50.

Table 10: Abstracted label transfer attack results

Table 11: Original ImageNet label transfer attack results.

method	Surrogate	Metrics						ViT-B/16					
		Clip/QA	LPIPS	MSSSIM	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR
ours _{c=1.5}	ConvNext-T	0.822 \pm 0.003	0.018 \pm 0.000	0.576 \pm 0.001	0.158 \pm 0.001	0.550 \pm 0.003	0.159 \pm 0.002	0.542 \pm 0.003	0.194 \pm 0.001	0.571 \pm 0.002	0.150 \pm 0.001	0.981 \pm 0.001	0.514 \pm 0.002
	ResNet-DecoWa	0.821 \pm 0.003	0.022 \pm 0.000	0.973 \pm 0.001	0.654 \pm 0.001	0.133 \pm 0.001	0.567 \pm 0.002	0.168 \pm 0.001	0.374 \pm 0.002	0.300 \pm 0.001	0.940 \pm 0.002	0.019 \pm 0.001	0.484 \pm 0.001
	ResNet50	0.821 \pm 0.004	0.016 \pm 0.000	0.979 \pm 0.000	0.559 \pm 0.002	0.157 \pm 0.001	0.465 \pm 0.002	0.193 \pm 0.001	0.275 \pm 0.000	0.328 \pm 0.001	0.97 \pm 0.000	0.001 \pm 0.000	0.228 \pm 0.003
	Vit-B/16	0.820 \pm 0.005	0.019 \pm 0.000	0.977 \pm 0.001	0.473 \pm 0.004	0.188 \pm 0.002	0.470 \pm 0.004	0.188 \pm 0.002	0.478 \pm 0.001	0.212 \pm 0.001	0.470 \pm 0.005	0.182 \pm 0.002	0.155 \pm 0.002
ours _{c=2}	ConvNext-T	0.819 \pm 0.002	0.027 \pm 0.000	0.968 \pm 0.000	0.674 \pm 0.001	0.121 \pm 0.001	0.641 \pm 0.003	0.127 \pm 0.002	0.630 \pm 0.002	0.158 \pm 0.001	0.660 \pm 0.001	0.118 \pm 0.001	0.970 \pm 0.001
	ResNet-DecoWa	0.815 \pm 0.004	0.033 \pm 0.001	0.960 \pm 0.001	0.782 \pm 0.001	0.084 \pm 0.000	0.680 \pm 0.003	0.124 \pm 0.001	0.479 \pm 0.002	0.245 \pm 0.001	0.977 \pm 0.002	0.007 \pm 0.001	0.619 \pm 0.002
	ResNet50	0.819 \pm 0.004	0.023 \pm 0.000	0.970 \pm 0.000	0.650 \pm 0.002	0.123 \pm 0.001	0.517 \pm 0.001	0.174 \pm 0.001	0.346 \pm 0.002	0.305 \pm 0.000	0.999 \pm 0.000	0.004 \pm 0.000	0.386 \pm 0.002
	Vit-B/16	0.817 \pm 0.004	0.028 \pm 0.001	0.967 \pm 0.000	0.555 \pm 0.004	0.158 \pm 0.001	0.548 \pm 0.003	0.160 \pm 0.001	0.577 \pm 0.001	0.170 \pm 0.001	0.557 \pm 0.002	0.027 \pm 0.003	0.126 \pm 0.000
ours _{c=2.5}	ConvNext-T	0.817 \pm 0.003	0.036 \pm 0.000	0.958 \pm 0.001	0.740 \pm 0.002	0.096 \pm 0.001	0.710 \pm 0.003	0.104 \pm 0.001	0.694 \pm 0.002	0.130 \pm 0.001	0.712 \pm 0.002	0.100 \pm 0.001	0.997 \pm 0.000
	ResNet-DecoWa	0.810 \pm 0.004	0.044 \pm 0.001	0.947 \pm 0.001	0.870 \pm 0.001	0.091 \pm 0.001	0.762 \pm 0.001	0.140 \pm 0.001	0.562 \pm 0.002	0.204 \pm 0.001	0.991 \pm 0.001	0.003 \pm 0.000	0.139 \pm 0.001
	ResNet50	0.815 \pm 0.006	0.031 \pm 0.000	0.961 \pm 0.000	0.724 \pm 0.001	0.098 \pm 0.001	0.607 \pm 0.003	0.142 \pm 0.001	0.399 \pm 0.002	0.272 \pm 0.001	1.00 \pm 0.000	0.000 \pm 0.000	0.664 \pm 0.002
	Vit-B/16	0.815 \pm 0.004	0.038 \pm 0.000	0.986 \pm 0.000	0.619 \pm 0.004	0.136 \pm 0.002	0.599 \pm 0.003	0.142 \pm 0.001	0.644 \pm 0.001	0.143 \pm 0.001	0.677 \pm 0.004	0.130 \pm 0.001	0.996 \pm 0.001
ours _{c=3}	ConvNext-T	0.814 \pm 0.003	0.045 \pm 0.000	0.947 \pm 0.001	0.788 \pm 0.001	0.076 \pm 0.001	0.767 \pm 0.002	0.082 \pm 0.001	0.725 \pm 0.001	0.167 \pm 0.001	0.772 \pm 0.001	0.001 \pm 0.000	0.094 \pm 0.001
	ResNet-DecoWa	0.807 \pm 0.004	0.056 \pm 0.000	0.934 \pm 0.001	0.912 \pm 0.001	0.034 \pm 0.001	0.841 \pm 0.001	0.062 \pm 0.001	0.636 \pm 0.002	0.173 \pm 0.001	0.999 \pm 0.000	0.000 \pm 0.000	0.745 \pm 0.002
	ResNet50	0.812 \pm 0.005	0.039 \pm 0.001	0.952 \pm 0.001	0.775 \pm 0.003	0.079 \pm 0.001	0.649 \pm 0.002	0.126 \pm 0.001	0.453 \pm 0.002	0.250 \pm 0.001	1.000 \pm 0.000	0.000 \pm 0.000	0.147 \pm 0.001
	Vit-B/16	0.812 \pm 0.003	0.048 \pm 0.001	0.947 \pm 0.001	0.667 \pm 0.003	0.112 \pm 0.001	0.654 \pm 0.002	0.123 \pm 0.001	0.701 \pm 0.002	0.121 \pm 0.001	0.682 \pm 0.001	0.110 \pm 0.001	0.45 \pm 0.002
MI-FGSM	ConvNext-T	0.543 \pm 0.004	0.211 \pm 0.001	0.877 \pm 0.000	0.373 \pm 0.001	0.334 \pm 0.001	0.334 \pm 0.002	0.334 \pm 0.002	0.281 \pm 0.002	0.431 \pm 0.001	0.474 \pm 0.001	0.998 \pm 0.001	0.338 \pm 0.001
	ResNet-DecoWa	0.548 \pm 0.003	0.209 \pm 0.002	0.880 \pm 0.001	0.430 \pm 0.003	0.340 \pm 0.001	0.660 \pm 0.001	0.191 \pm 0.001	0.358 \pm 0.002	0.201 \pm 0.002	0.497 \pm 0.001	0.001 \pm 0.000	0.273 \pm 0.002
	ResNet50	0.535 \pm 0.003	0.208 \pm 0.002	0.869 \pm 0.001	0.667 \pm 0.001	0.180 \pm 0.000	0.660 \pm 0.001	0.190 \pm 0.001	0.353 \pm 0.001	0.397 \pm 0.001	0.496 \pm 0.001	0.119 \pm 0.001	0.456 \pm 0.001
	Vit-B/16	0.539 \pm 0.003	0.205 \pm 0.001	0.836 \pm 0.001	0.221 \pm 0.002	0.454 \pm 0.001	0.251 \pm 0.002	0.411 \pm 0.001	0.192 \pm 0.002	0.492 \pm 0.001	0.420 \pm 0.001	0.111 \pm 0.002	0.359 \pm 0.002
VENOM	ConvNext-T	0.796 \pm 0.002	0.020 \pm 0.001	0.978 \pm 0.001	0.350 \pm 0.002	0.338 \pm 0.003	0.383 \pm 0.002	0.303 \pm 0.001	0.278 \pm 0.001	0.405 \pm 0.001	0.387 \pm 0.002	0.301 \pm 0.001	0.294 \pm 0.002
	ResNet-DecoWa	0.805 \pm 0.002	0.027 \pm 0.001	0.968 \pm 0.002	0.257 \pm 0.001	0.388 \pm 0.002	0.372 \pm 0.002	0.376 \pm 0.001	0.136 \pm 0.001	0.490 \pm 0.002	0.917 \pm 0.001	0.040 \pm 0.001	0.144 \pm 0.001
	ResNet50	0.795 \pm 0.003	0.023 \pm 0.000	0.972 \pm 0.001	0.288 \pm 0.001	0.355 \pm 0.001	0.372 \pm 0.002	0.247 \pm 0.001	0.147 \pm 0.001	0.484 \pm 0.001	0.991 \pm 0.001	0.005 \pm 0.000	0.151 \pm 0.001
	Vit-B/16	0.796 \pm 0.002	0.021 \pm 0.001	0.977 \pm 0.001	0.274 \pm 0.001	0.374 \pm 0.001	0.288 \pm 0.002	0.348 \pm 0.001	0.267 \pm 0.001	0.411 \pm 0.002	0.341 \pm 0.000	0.350 \pm 0.001	0.005 \pm 0.000
SD-NAE	ConvNext-T	0.848 \pm 0.009	0.048 \pm 0.002	0.942 \pm 0.018	0.655 \pm 0.001	0.221 \pm 0.001	0.690 \pm 0.002	0.209 \pm 0.001	0.626 \pm 0.002	0.255 \pm 0.001	0.710 \pm 0.001	0.195 \pm 0.001	0.649 \pm 0.002
	ResNet-DecoWa	0.845 \pm 0.013	0.470 \pm 0.018	0.421 \pm 0.018	0.250 \pm 0.001	0.646 \pm 0.001	0.259 \pm 0.001	0.549 \pm 0.001	0.423 \pm 0.002	0.691 \pm 0.002	0.208 \pm 0.002	0.587 \pm 0.001	0.294 \pm 0.001
	ResNet50	0.841 \pm 0.011	0.457 \pm 0.020	0.433 \pm 0.017	0.502 \pm 0.001	0.347 \pm 0.001	0.539 \pm 0.001	0.420 \pm 0.001	0.518 \pm 0.001	0.423 \pm 0.001	0.585 \pm 0.002	0.377 \pm 0.001	0.454 \pm 0.001
	Vit-B/16	0.844 \pm 0.009	0.459 \pm 0.016	0.441 \pm 0.016	0.530 \pm 0.002	0.327 \pm 0.001	0.467 \pm 0.001	0.383 \pm 0.001	0.505 \pm 0.002	0.302 \pm 0.001	0.501 \pm 0.001	0.340 \pm 0.001	0.354 \pm 0.001
AdvDiff	ConvNext-T	0.636 \pm 0.006	0.471 \pm 0.025	0.312 \pm 0.011	0.440 \pm 0.022	0.541 \pm 0.022	0.448 \pm 0.022	0.523 \pm 0.022	0.433 \pm 0.022	0.533 \pm 0.022	0.446 \pm 0.022	0.457 \pm 0.022	0.545 \pm 0.022
	ResNet-DecoWa	0.597 \pm 0.004	0.293 \pm 0.003	0.761 \pm 0.003	0.289 \pm 0.001	0.672 \pm 0.001	0.247 \pm 0.001	0.698 \pm 0.001	0.180 \pm 0.001	0.479 \pm 0.001	0.487 \pm 0.001	0.147 \pm 0.001	0.225 \pm 0.001
	ResNet50	0.634 \pm 0.004	0.046 \pm 0.003	0.929 \pm 0.001	0.046 \pm 0.001	0.857 \pm 0.001	0.050 \pm 0.001	0.863 \pm 0.001	0.035 \pm 0.002	0.91 \pm 0.001	0.931 \pm 0.001	0.041 \pm 0.002	0.881 \pm 0.001
	Vit-B/16	0.638 \pm 0.004	0.390 \pm 0.025	0.430 \pm 0.010	0.295 \pm 0.021	0.673 \pm 0.021	0.304 \pm 0.021	0.651 \pm 0.020	0.300 \pm 0.021	0.659 \pm 0.022	0.327 \pm 0.021	0.660 \pm 0.021	0.327 \pm 0.021