

---

# Exploring Semantic-constrained Adversarial Example with Instruction Uncertainty Reduction

---

Jin Hu<sup>1,2</sup> Jiakai Wang<sup>2\*</sup> Linna Jing<sup>1</sup> Haolin Li<sup>3</sup> Haodong Liu<sup>3</sup>  
Haotong Qin<sup>4</sup> Aishan Liu<sup>1</sup> Ke Xu<sup>1,2</sup> Xianglong Liu<sup>1,2</sup>

<sup>1</sup>State Key Laboratory of Complex & Critical Software Environment, Beihang University

<sup>2</sup>Zhongguancun Laboratory <sup>3</sup>School of Computer Science and Engineering, Beihang University

<sup>4</sup>Dept. of Information Technology and Electrical Engineering, ETH Zurich

{hujin,linnajing,lh142195,21373450,liuaishan,kexu,xlliu}@buaa.edu.cn  
wangjk@mail.zgclab.edu.cn haotong.qin@pbl.ee.ethz.ch

## 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, 3], transfer AEs [4, 5, 6], and naturalistic AEs [7, 8, 9], as well as the applications in evaluating autonomous driving [10, 11, 12] or LLM systems [13, 14].

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 [15, 16, 17], , but the de facto potential of SemanticAE has still not been fully released in performing transferable, adaptive, and effective attacks. In light of the recent advancements in language-driven multimodal intelligence and the increasing demand for alignment [18, 19, 20], we believe that 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* [21], *PGD* [22], *AutoAttack* [23], iterative optimization in the latent space of generative models, e.g., *DiffPGD* [24], *NAP* [25], *AC-GAN* [26], and training a neural network for adversarial attack generation, e.g., *AdvGAN* [27]. Adversarial examples may not be robust in the physical world. For physical attacks, expectation-over-transformation (EoT) [28] and 3D simulation [29] are proposed to bridge the digital-physical gap.

**Semantic-constrained Adversarial Example** [30] 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 [31, 32, 33], 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 diverse-distributed adversarial examples without referencing a static image. We formulate the

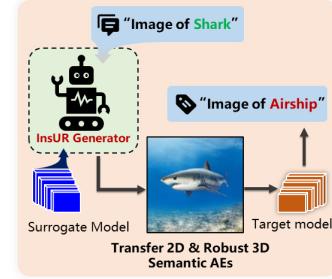


Figure 1: SemanticAEs are generated directly by instructions.

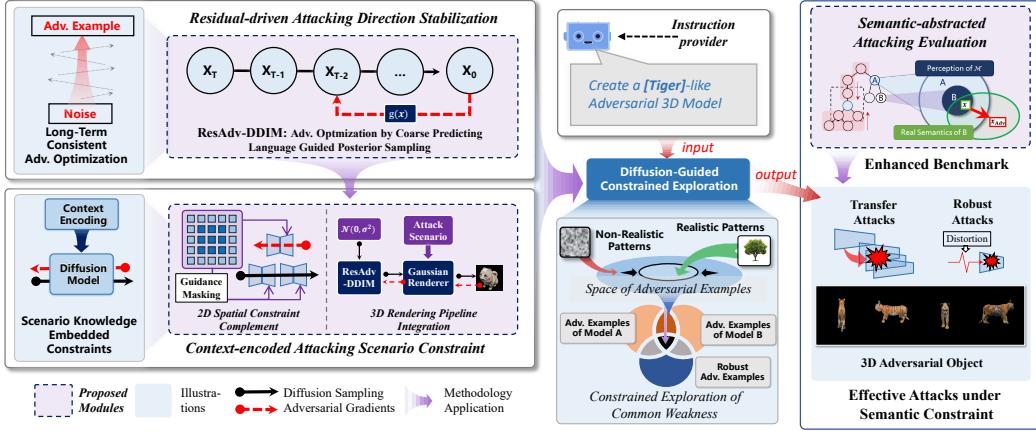


Figure 2: Overview of multi-dimensional instruction uncertainty reduction (**InSUR**) framework.

adversarial example generating task constrained by natural language’s semantics as the *semantically-constrained adversarial example generation* problem. Recent works([24, 33, 34, 16, 17]) 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. An extended technical background and related works are provided in Appendix A.

### 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  $x_{\text{adv}}$  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,  $A_{\text{Text}}$  defines the types of target model’s output that are conceptually antonyms of Text. 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  $x_{\text{adv}}$  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 [35]. 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. 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-abstacted attacking evaluation enhancement is proposed to facilitate further developments of SemanticAE generation.

### 3.2 Residual-driven Attacking Direction Stabilization with ResAdv-DDIM

**Semantically-Constrained Optimization Problem** Referring to the generative-model-based adversarial examples, we solve SemanticAE generation by maximizing the loss  $\mathcal{L}_{\text{ATK}}$  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 [36], as a constraint defined by Text [24, 15]. 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 [24] or directly altering the sampling process [16], 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, 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, limiting the advantage of multi-step diffusion models in transfer attacks.

Overall, such technical challenge originates from the conflict between (1) the accurate estimation of  $x_t$ , causing the robust problems of the language guidance defined by  $f$  and the computational problems, and (2) the approximated estimation of  $x_t$ , causing the non-optimality of attack optimization. We solve this problem by improve the approximation with the novel *ResAdv-DDIM* posterior sampler.

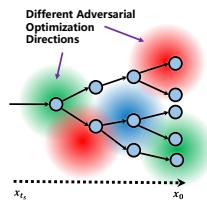


Figure 3:  
Inconsistent adv.  
direction problem.

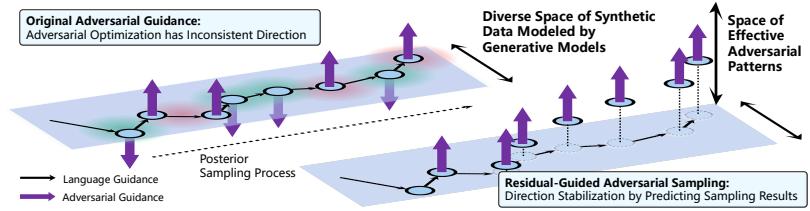


Figure 4: Residual-driven attacking direction stabilization. *ResAdv-DDIM* is designed for the efficient and thorough exploration of new adversarial patterns constrained by multi-step sampling processes.

**Residual-Guided Adversarial DDIM Sampler** Inspired by *Learning to Optimize* [37], we handle the challenge by **predicting a coarse sketch of the future-step denoising result**  $x_0$  for estimating the attack optimization direction with  $\mathcal{L}_{\text{ATK}}$ . 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.

Drawing inspiration from the DDIM sampler’s architectural design, we leverage DDIM’s multi-step posterior sampling capabilities to achieve a coarse prediction of  $x_0$ . Specifically, we define the generation process as:

$$g_\theta(x_t) = \underbrace{f_{\theta, \Delta T_1} \circ f_{\theta, \Delta T_2} \circ \dots \circ f_{\theta, \Delta T_k}}_{k \text{ times}, k \ll T/\Delta T}(x_t), \quad \text{where } \sum_{i=1}^k \Delta T_i = t \quad (4)$$

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

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 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  $x_{t_s}$ :

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

**(2) Stabilizing the temporary label selection in  $\mathcal{L}_{\text{ATK}}$ .** For evasion attacks, a typical construction of attack loss is maximizing the confidence of the highest confident label  $l_{\text{tar}}$  in the incorrect label set  $A_{\text{Text}}$ . To eliminate the guidance fluctuation further, we adapt the fixed  $l_{\text{tar}}$  after its initialization in the initial steps of attack optimization. **(3) Adaptive attack optimization iteration.** For the maximization problem defined in Eq. 4, we employ the early-stop mechanism with the following signal:

$$(i \leq n \wedge \arg \max_{l \notin A_{\text{Text}}} P(\mathcal{M}(x_t) = l) > \xi_1) \vee (i = 1 \wedge \arg \max_{l \notin A_{\text{Text}}} P(\mathcal{M}(x_t) = l) > \xi_2), \quad (6)$$

Where  $i$  is the current number of iterations,  $n$  is the upper-bound iteration number, and the probability is estimated by the confidence of  $\mathcal{M}$ 's prediction. The optimization is performed if and only if the probability of the unsuccessful attack is lower than the threshold defined by  $\xi_1$  and  $\xi_2$ . We set  $\xi_1 = 0.1$  and  $\xi_2 = 0.01$ ,  $n = 10$  for the first and last steps and  $n = 3$  for other steps. In addition, momentum optimization, introduced in [17, 38], is applied.

### 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. Detailed implementation, together with the ResAdv-DDIM sampler, is shown in Appendix B.

**Spatial Constraint Complement for 2D SemanticAE Generation** Effective SemanticAE generator shall leverage the optimization space of the image backgrounds and generate patterns that amplify the attack's effectiveness. However, the background generation automatically learned by the diffusion model is overly uniform, since the attack functionalities are not considered in the original training. To generate attack-related adversarial backgrounds, we introduce a guidance-masking-based method into the key guidance function  $f_\theta$  applied in both posterior sampling and adversarial optimization process in ResAdv-DDIM. Together with the deterministic DDIM,  $f_\theta$  is formulated as:

$$f_{\theta, \Delta T}(x_t) = \sqrt{\bar{\alpha}_{t-\Delta T}/\bar{\alpha}_t} (x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(x_t, t)) + \sqrt{1 - \bar{\alpha}_{t-\Delta T}} \cdot \epsilon_\theta(x_t, t), \quad (7)$$

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

where  $\alpha$  defines the noise ratio in the diffusion model,  $\epsilon_\theta$  is the noise estimating network, and  $M$  is the guidance masking that regularizes the spatial distribution of semantic guidance Text. Since adversarial optimization in ResAdv-DDIM incorporates the results of the coarse denoising steps, the optimization can automatically adapt the corrected guidance and generate attack-related backgrounds.

**Differentiable Rendering Pipeline Integration for 3D SemanticAE Generation** 3D Data is valuable for world modeling [39, 40]. 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 [40], we fill the gap between the 3D generation and 2D target models. As shown in Figure 5, we optimize latents with the 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* [41] framework, which is a recently published diffusion-based 3D access generation framework. For simplicity, we reformulate the *Trellis* sampling and rendering process as:

$$\begin{aligned} \mathbf{pos} &= \text{Coords}(\mathcal{D}_{\text{slat}}(\mathbf{z}_0^{\text{slat}})), \quad \text{Model}_{\text{GS}} = \mathcal{D}_{\text{GS}}(\mathbf{z}_0, \mathbf{pos}), \quad x = \text{Renderer}_{\text{GS}}(\text{Model}_{\text{GS}}, \text{Camera}), \\ \mathcal{D}_{\text{GS}} : \{(z^i, pos^i)\}_{i=1}^L &\rightarrow \{\{(x_i^k, c_i^k, s_i^k, \alpha_i^k, r_i^k)\}_{k=1}^K\}_{i=1}^L, \end{aligned} \quad (8)$$

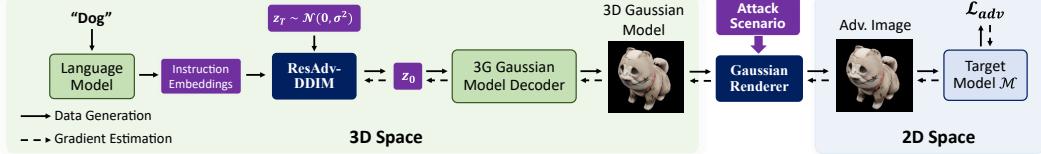


Figure 5: 3D optimization pipeline. Scenario knowledge is encoded with the Gaussian renderer.

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,  $\text{Coords}$  transforms the voxel to point positions  $\text{pos}$ .  $\mathcal{D}_{GS}$  is the refined structure decoder that decodes each vertex into multiple Gaussian points and  $\text{Renderergs}$  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  $g_\theta$ :

$$g_\theta(z_t, \text{pos}, \text{Camera}) := \text{Renderergs}(\mathcal{D}_{GS}(f_{\theta, \Delta T_1} \circ \dots \circ f_{\theta, \Delta T_k}(z_t, \text{pos}), \text{pos}), \text{Camera}) \\ z_{t-\Delta T} := f_{\theta, \Delta T}(\arg \max_{z_t} \mathbb{E}_{\text{Camera} \sim P_{\text{Cam}}} [\mathcal{L}_{\text{ATK}}(\mathcal{M}(g_\theta(z_t, \text{Camera}, \text{pos})))]) \quad (9)$$

We use the EoT method with gradient accumulation to optimize  $z_t$  for unknown camera positioning, i.e., samples Camera from  $P_{\text{Cam}}$  in each iteration. The scenario knowledge is encoded in  $P_{\text{Cam}}$  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.

### 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  $A_{\text{Text}}$ , 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.

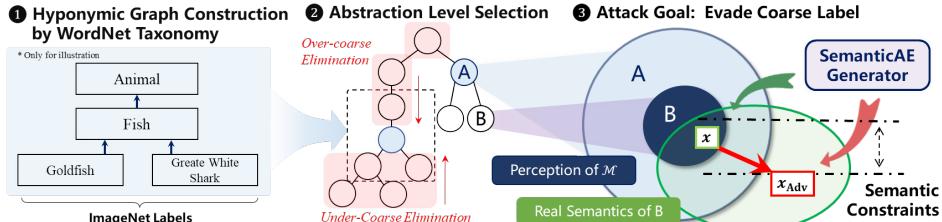


Figure 6: Construction of abstract label evasion evaluation task.

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  $A_{\text{Text}}$  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* [42] 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:

$$\text{Text} := \text{"Realistic image of [AbstractedLabel], specifically, [label]"}, \\ A_{\text{Text}} := \{\text{label}_{\text{Adv}} \mid \text{AbstractedLabel} \notin \text{Ancestors}(\text{label}_{\text{Adv}})\}, \quad (10)$$

$$\text{AbstractedLabel} \in \{c \in \mathbb{L}' \mid \exists l \in \mathbb{L} \text{ s.t. } c \in \text{Ancestors}(l) \wedge \text{Count}_{\text{Children}}(c) > 0\},$$

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 [17], 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  $x_{\text{adv}}$  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:

$$ASR_{\text{Relative}} = \frac{\sum_{i=1}^K \text{Attack Success}(x_{\text{adv}}^{(i)}) \wedge \text{Classification Correct}(x_{\text{exemplar}}^{(i)})}{K \cdot \text{Accuracy}(\mathcal{X}_{\text{exemplar}})} \in [0, 1], \quad (11)$$

$$\text{SemanticDiff}_S = \langle x_{\text{exemplar}}, x_{\text{adv}} \rangle_S,$$

where  $K$  is the amount of samples,  $S$  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,  $ASR_{\text{Relative}}$  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*( $x_{\text{exemplar}}$ ).

## 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 [38], DeCoWA [43]), and ③ Diffusion-based naturalistic AE generation methods (AdvDiff [16], SD-NAE [15], VENOM [17]). 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* [44] and *MSSSIM* [45], are been employed to measure the similarity between  $x_{\text{exemplar}}$  and  $x_{\text{adv}}$  under the proposed non-adversarial exemplar evaluation. Also, we supplement the non-reference image quality assessment *CLIP-IQA* [46] 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  $ASR_{\text{Relative}}$ , defined in Eq. 11, across diverse targets set  $\mathcal{T} = \{\text{ResNet50} [47], \text{ViT-B/16} [48], \text{ConvNeXt-T} [49], \text{ResNet152}, \text{InceptionV3} [50], \text{Swin-Transformer-B} [51]\}$ . The first three models are used individually as surrogates in experiments. The main paper presents the average  $ASR_{\text{relative}}$  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.

### 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  $\epsilon$  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. Tabel 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  $1.19 \times$**  average ASR and  **$1.08 \times$**  minimal ASR across **all**

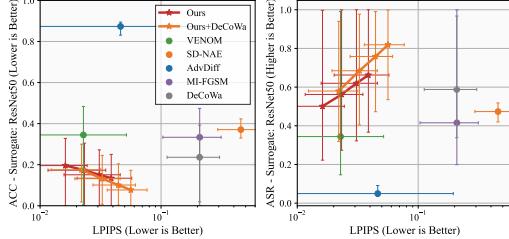


Figure 7: ImageNet label results.

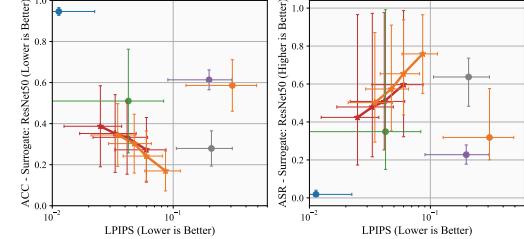


Figure 8: Abstracted label results.

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

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

**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  $\epsilon$ -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 ant targets model is ResNet50. The results are in Tabel 2. There is no 3D SemanticAE previously available. It shows that our results show the 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 ( $Iter_{max}$ ) in the adversarial guidance  $g$  in Eq 4 in the 2D abstracted label task and the 3D task. The performance improvement is significant and consistent compared to the result without future-step sampling prediction ( $Iter_{max}=0$ ). 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 diffusion steps provide better on-manifold regularization and better adversarial transfer-

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	<b>2.8%</b>	<b>92.2%</b>	<b>0.665</b>	<b>0.258</b>

Table 3: Ablation of residual approximation

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

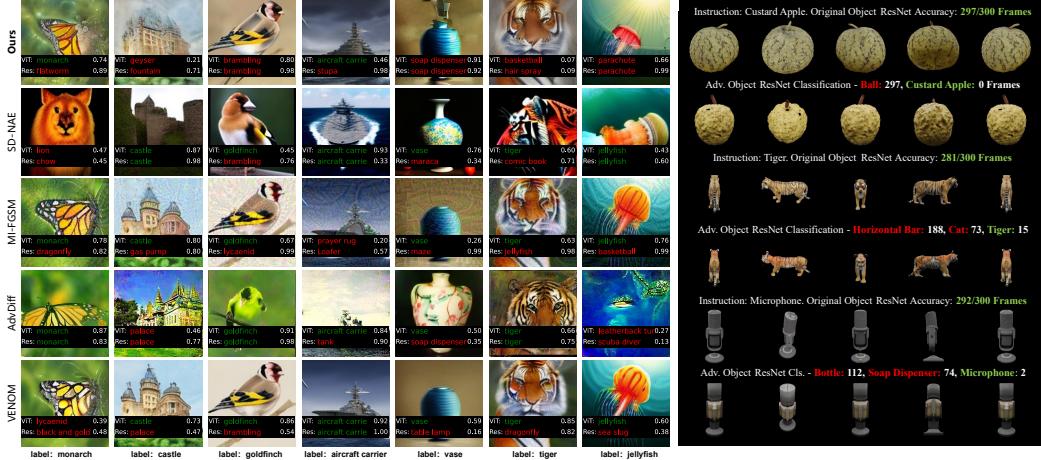


Figure 10: Comparison of 2D SemanticAEs.



Figure 11: 3D SemanticAEs

ability. Detailed results on more settings and further analysis are shown in the Appendices. Moreover, attributed to the adaptive iteration mechanism, more accurate estimation leads to lower optimization steps, and therefore, the increase in time consumption is sub-linear.  $\text{Clip}_Q$ , MSSSIM, and LPIPS are also slightly improved, which may be attributed to the enriched guidance diversity and the better collaboration between adversarial and language guidance.

Table 4: Ablation of guidance masking.

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

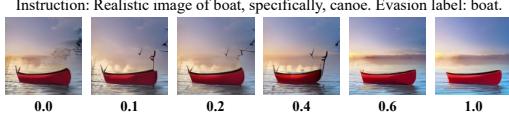


Figure 9: Vis. of different guidance masking. The value under images denotes  $M_{\text{edge}}/M_{\text{mid}}$ .

**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. 7, and setting to 1.0 represent removing the masking. The results are shown in Tabel 4 and Figure 9, indicating that (1) decreasing  $M_{\text{edge}}$  leads to the increase of unconditional guidance, which enriches the background diversity as shown in the figure, and may leads to the improvements of  $\text{Clip}_Q$ . (2) When the budget  $\epsilon$  is small ( $\epsilon = 2$ ) and the optimization space is relatively narrow, diversifying the background is beneficial. Note that  $\epsilon$  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.

#### 4.4 Visualization of Generated SemanticAEs

We selected the 2D / 3D samples with  $x_{\text{exemplar}}$  correctly classified and visualized in Fig. 10 and Fig. 11. 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 MSSSIM 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.

### 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.

## Acknowledgments and Disclosure of Funding

This work is supported by Zhongguancun Laboratory and State Key Laboratory of Complex & Critical Software Environment, Beihang University.

## References

- [1] Huali Ren, Teng Huang, and Hongyang Yan. Adversarial examples: attacks and defenses in the physical world. *International Journal of Machine Learning and Cybernetics*, 12(11):3325–3336, 2021.
- [2] Yanjie Li, Yiqian Li, Xuelong Dai, Songtao Guo, and Bin Xiao. Physical-world optical adversarial attacks on 3d face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 24699–24708, 2023.
- [3] Chengyin Hu, Weiwen Shi, Wen Yao, Tingsong Jiang, Ling Tian, Xiaoqian Chen, and Wen Li. Adversarial infrared curves: An attack on infrared pedestrian detectors in the physical world. *Neural networks*, 178:106459, 2024.
- [4] Wuping Ke, Desheng Zheng, Xiaoyu Li, Yuanhang He, Tianyu Li, and Fan Min. Improving the transferability of adversarial examples through neighborhood attribution. *Knowledge-Based Systems*, 296:111909, 2024.
- [5] Yaguan Qian, Kecheng Chen, Bin Wang, Zhaoquan Gu, Shouling Ji, Wei Wang, and Yanchun Zhang. Enhancing transferability of adversarial examples through mixed-frequency inputs. *IEEE Transactions on Information Forensics and Security*, 2024.
- [6] Rongyi Zhu, Zeliang Zhang, Susan Liang, Zhuo Liu, and Chenliang Xu. Learning to transform dynamically for better adversarial transferability. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24273–24283, 2024.
- [7] 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.
- [8] Bao Gia Doan, Minhui Xue, Shiqing Ma, Ehsan Abbasnejad, and Damith C Ranasinghe. Tnt attacks! universal naturalistic adversarial patches against deep neural network systems. *IEEE Transactions on Information Forensics and Security*, 17:3816–3830, 2022.
- [9] Zimu Wang, Wei Wang, Qi Chen, Qiufeng Wang, and Anh Nguyen. Generating valid and natural adversarial examples with large language models. In *2024 27th International Conference on Computer Supported Cooperative Work in Design (CSCWD)*, pages 1716–1721. IEEE, 2024.
- [10] 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.
- [11] 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.
- [12] 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.
- [13] Xilie Xu, Keyi Kong, Ning Liu, Lizhen Cui, Di Wang, Jingfeng Zhang, and Mohan Kankanhalli. An llm can fool itself: A prompt-based adversarial attack. *arXiv preprint arXiv:2310.13345*, 2023.
- [14] 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.

- [15] Yueqian Lin, Jingyang Zhang, Yiran Chen, and Hai Li. Sd-nae: Generating natural adversarial examples with stable diffusion. *arXiv preprint arXiv:2311.12981*, 2023.
- [16] Xuelong Dai, Kaisheng Liang, and Bin Xiao. Advdif: Generating unrestricted adversarial examples using diffusion models. In *European Conference on Computer Vision*, pages 93–109. Springer, 2024.
- [17] Hui Kuurila-Zhang, Haoyu Chen, and Guoying Zhao. Venom: Text-driven unrestricted adversarial example generation with diffusion models. *arXiv preprint arXiv:2501.07922*, 2025.
- [18] Taecheon Kim, Sangyun Chung, Damin Yeom, Youngjoon Yu, Hak Gu Kim, and Yong Man Ro. Mscotdet: Language-driven multi-modal fusion for improved multispectral pedestrian detection. *IEEE Transactions on Circuits and Systems for Video Technology*, 2024.
- [19] Jianzong Wu, Xiangtai Li, Chenyang Si, Shangchen Zhou, Jingkang Yang, Jiangning Zhang, Yining Li, Kai Chen, Yunhai Tong, Ziwei Liu, et al. Towards language-driven video inpainting via multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12501–12511, 2024.
- [20] Fukun Yin, Xin Chen, Chi Zhang, Biao Jiang, Zibo Zhao, Wen Liu, Gang Yu, and Tao Chen. Shapegpt: 3d shape generation with a unified multi-modal language model. *IEEE Transactions on Multimedia*, 2025.
- [21] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [22] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- [23] 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.
- [24] Haotian Xue, Alexandre Araujo, Bin Hu, and Yongxin Chen. Diffusion-based adversarial sample generation for improved stealthiness and controllability. *Advances in Neural Information Processing Systems*, 36:2894–2921, 2023.
- [25] Jiahui Lei, Congyue Deng, Bokui Shen, Leonidas Guibas, and Kostas Daniilidis. Nap: Neural 3d articulation prior. *arXiv preprint arXiv:2305.16315*, 2023.
- [26] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier gans. In *International conference on machine learning*, pages 2642–2651. PMLR, 2017.
- [27] Chaowei Xiao, Bo Li, Jun-Yan Zhu, Warren He, Mingyan Liu, and Dawn Song. Generating adversarial examples with adversarial networks. *arXiv preprint arXiv:1801.02610*, 2018.
- [28] Anish Athalye, Logan Engstrom, Andrew Ilyas, and Kevin Kwok. Synthesizing robust adversarial examples. In *International conference on machine learning*, pages 284–293. PMLR, 2018.
- [29] 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.
- [30] Yang Song, Rui Shu, Nate Kushman, and Stefano Ermon. Constructing unrestricted adversarial examples with generative models. *Advances in neural information processing systems*, 31, 2018.
- [31] 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.

- [32] 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 (ICCV)*, pages 7848–7857, October 2021.
- [33] Xinquan Chen, Xitong Gao, Juanjuan Zhao, Kejiang Ye, and Cheng-Zhong Xu. Advdiffuser: Natural adversarial example synthesis with diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4562–4572, 2023.
- [34] 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.
- [35] Yanbo Chen and Weiwei Liu. A theory of transfer-based black-box attacks: Explanation and implications. *Advances in Neural Information Processing Systems*, 36:13887–13907, 2023.
- [36] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [37] 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.
- [38] 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.
- [39] 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.
- [40] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023.
- [41] 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. *arXiv preprint arXiv:2412.01506*, 2024.
- [42] George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995.
- [43] 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.
- [44] 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.
- [45] Zhou Wang, Eero P Simoncelli, and Alan C Bovik. Multiscale structural similarity for image quality assessment. In *The Thirty-Seventh Asilomar Conference on Signals, Systems & Computers*, 2003, volume 2, pages 1398–1402. Ieee, 2003.
- [46] 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.
- [47] 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.

- [48] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [49] 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.
- [50] 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.
- [51] 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.
- [52] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- [53] 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.
- [54] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. In *International Conference on Machine Learning*, pages 32211–32252. PMLR, 2023.
- [55] Zhengyang Geng, Mingyang Deng, Xingjian Bai, J Zico Kolter, and Kaiming He. Mean flows for one-step generative modeling. *arXiv preprint arXiv:2505.13447*, 2025.
- [56] Viraj Prabhu, Sriram Yenamandra, Prithvijit Chattopadhyay, and Judy Hoffman. Lance: Stress-testing visual models by generating language-guided counterfactual images. *Advances in Neural Information Processing Systems*, 36:25165–25184, 2023.
- [57] 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.
- [58] 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 Srikumar, 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.
- [59] 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.
- [60] 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.
- [61] Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples and black-box attacks. *arXiv preprint arXiv:1611.02770*, 2016.
- [62] 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.
- [63] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. *arXiv preprint arXiv:2010.01412*, 2020.

- [64] 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.
- [65] 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.
- [66] 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.
- [67] 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 (ECCV)*, pages 52–67, 2018.
- [68] 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.
- [69] Leheng Li, LIAN Qing, and Ying-Cong Chen. Adv3d: Generating 3d adversarial examples in driving scenarios with nerf. 2023.
- [70] 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.
- [71] Sergey Kastrulyin, Jamil Zakirov, Denis Prokopenko, and Dmitry V Dylov. Pytorch image quality: Metrics for image quality assessment. *arXiv preprint arXiv:2208.14818*, 2022.
- [72] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Backgrounds</b>	<b>2</b>
<b>3</b>	<b>Methodology</b>	<b>3</b>
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
<b>4</b>	<b>Experiments</b>	<b>7</b>
4.1	Experiment Settings . . . . .	7
4.2	Overall Performance Evaluation . . . . .	7
4.3	Key Ablation Studies . . . . .	8
4.4	Visualization of Generated SemanticAEs . . . . .	9
<b>5</b>	<b>Conclusion</b>	<b>9</b>
<b>A</b>	<b>Technical Background and Related Works</b>	<b>17</b>
A.1	Diffusion models . . . . .	17
A.2	Generating Hard Samples from Language . . . . .	17
A.3	Transfer Attack Methodology . . . . .	18
A.4	3D Generation and Gaussian Splatting . . . . .	18
<b>B</b>	<b>Detailed Implementation of InSUR Framework</b>	<b>19</b>
B.1	2D Image Generation . . . . .	19
B.2	3D Object Generation . . . . .	21
B.3	Construction and Analysis of Semantic-Abstracted Evaluation Task . . . . .	22
<b>C</b>	<b>Detailed Experiment Settings</b>	<b>25</b>
C.1	2D Experiment Settings . . . . .	25
C.2	3D Experiment Settings . . . . .	25
<b>D</b>	<b>Detailed Results and Discussions</b>	<b>27</b>
D.1	Transfer Attack Analysis . . . . .	27
D.2	Attack Robustness Analysis . . . . .	28
D.3	On the Role of Residual-driven Attacking Direction Stabilization in 2D and 3D SemanticAE Applications . . . . .	28
D.4	On the Potential Adversarial Transferability to Semantic Evaluator . . . . .	29
<b>E</b>	<b>Results Visualization</b>	<b>32</b>
E.1	2D Visualized Results . . . . .	32
E.2	3D Visualized Results . . . . .	32
<b>F</b>	<b>Social Impacts</b>	<b>34</b>

## 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  $\beta$ . 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  $\epsilon_\theta$ . 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 [36] sampling procedure could be formulated as (deterministic version):

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

where  $\bar{\alpha} = \prod_{s=1}^t (1 - \beta_s)$ , and  $\Delta T$  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* [52].

**Discussions** In the application in content edit [53], 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 [54, 55]. 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 [56, 57, 58]. 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, [59] 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** [60], *i.e.*, creating efficient methodologies that utilize small models with limited capability to build the censorship framework for larger models. Also, we believe that the two methods could be implemented into an integrated red-team system. 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 [61]. Additionally, transferability can also be improved by ensembling multiple checkpoints along the training trajectory of a single model [62]. The Sharpness-Aware Minimization (SAM) method further boosts transferability by suppressing sharpness in the loss landscape [63]. 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 [64, 43]. This could be regarded as improving the diversity of the surrogate model. From the theoretical perspective, decreasing the cooperation within the adversarial pattern [65] and in-manifold constraint [35] is beneficial for transferability. In the adversarial example generation pipeline, the generative-model-based method, which this work focus on, could be regard as a replacement of gradient decent with the 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 [66] and meshes [67], ② implicit representations, such as Neural Radiance Fields (NeRFs) [39], ③ hybrid representations, such as 3D Gaussian [40]. 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  $\alpha$ , 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  $\alpha$ .

With the development of 3D representations, 3D adversarial examples for different 3D structures have also been advanced [68, 69, 70]. 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:

$$x_{t-\Delta T} = \sqrt{\bar{\alpha}_{t-\Delta T}/\bar{\alpha}_t} (x_t - \sqrt{1-\bar{\alpha}_t} \epsilon_\theta(x_t, t)) + \sqrt{1-\bar{\alpha}_{t-\Delta T}} \cdot \epsilon_\theta(x_t, t), \quad (13)$$

where  $\bar{\alpha} = \prod_{s=1}^t (1 - \beta_s)$  is the sampling parameter, and  $\Delta T$  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:

$$\begin{aligned} f_{\theta, \Delta T}^{(t, \text{Text})}(x_t) &= \sqrt{\frac{\bar{\alpha}_{t-\Delta T}}{\bar{\alpha}_t}} (x_t - \sqrt{1-\bar{\alpha}_t} \epsilon_\theta(x_t, t, \text{Text})) + \sqrt{1-\bar{\alpha}_{t-\Delta T}} \cdot \epsilon_\theta(x_t, t, \text{Text}), \\ \epsilon_\theta(x_t, t, \text{Text}) &:= (1-M) \cdot \epsilon_{\theta, \text{Unconditional}}(x_t, t) + M \cdot \epsilon_{\theta, \text{Conditional}}(x_t, t, \text{Text}). \end{aligned} \quad (14)$$

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

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

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:

$$\begin{aligned} g_\theta^{(t, \text{Text})}(x_t) &:= \underbrace{f_{\theta, \Delta T_1}^{(\Delta T_1, \text{Text})} \circ f_{\theta, \Delta T_2}^{(\Delta T_2 + \Delta T_1, \text{Text})} \circ \dots \circ f_{\theta, \Delta T_k}^{(t, \text{Text})}(x_t)}_{k \text{ times}}, \text{ where } \sum_{i=1}^k \Delta T_i = t, \\ k &:= \lceil \frac{t}{[T/K]} \rceil, \quad \Delta T_i(t) = \begin{cases} \lfloor T/K \rfloor, & 1 < i \leq k \\ t \mod [T/K], & i = 1 \end{cases}, \end{aligned} \quad (16)$$

where  $K$  is the maximal iteration number, corresponding to  $\text{Iter}_{\max}$  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  $\Delta T$ . 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:

$$x_{t-\Delta T} = \begin{cases} f_{\theta, \Delta T}^{(t, \text{Text})} \left( \arg \max_{x_t} \mathcal{L}_{\text{ATK}}(\mathcal{M} \circ g_\theta^{(t, \text{Text})}(x_t)) \right), & t \leq t_s, \\ f_{\theta, \Delta T}^{(t, \text{Text})}(x_t), & t > t_s. \end{cases} \quad (17)$$

**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 17, which is formulated as:

$$n := \begin{cases} M, & t = t_s \vee t < 4, \\ m, & \text{otherwise.} \end{cases}$$

$$\text{PerformOptimize} := \underbrace{(i \leq n \wedge \arg \max_{l \notin A_{\text{Text}}} P(\mathcal{M}(x_t) = l) > \xi_1)}_{\substack{\text{If confidence} \leq \xi_1, \text{ early stop. Maximal optimize iterations is } n.}} \vee$$

$$\underbrace{(i = 1 \wedge \arg \max_{l \notin A_{\text{Text}}} P(\mathcal{M}(x_t) = l) > \xi_2)}_{\substack{\text{If confidence} \leq \xi_2, \text{ do not optimize at the first step.}}}, \quad (18)$$

where  $n$  stands for the maximal adversarial optimization iteration in the single diffusion step. We set  $\xi_1 = 0.1$ ,  $\xi_2 = 0.01$ ,  $m = 3$  and  $M = 10$  for diversified strategy in different sampling steps, *i.e.*, the initial and final steps are set with a higher maximal iteration number ( $n = M$ ).

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,  $\mathcal{M}$ ,  $A_{\text{Text}}$ ,  $f$ ,  $\epsilon$ , optimization parameters  $T, K, t_s, t_k, \beta, s$   
Init  $x_T \sim \mathcal{N}(0, I)$ ,  $v_x \leftarrow 0$ ,  $x'_T \leftarrow x_T$ .  
**for**  $t = T, \dots, 1$  **do**

1	<b>if</b> $t \leq t_s$ <b>then</b> <b>for</b> $i = 1, 2, \dots, n$ <b>do</b> <b>if</b> $\neg \text{PerformOptimize (Eq. 18)}$ <b>then</b> <b>break</b> <span style="float: right;">/* Adaptive Iteration */</span> <b>if</b> $t = t_s \vee t < t_k$ <b>then</b> $L_{\text{Tar}} = \arg \max_{L \in A_{\text{Text}}} \mathcal{M}(g_{\theta}^{(t, \text{Text})}(x_t)) \logits[L]$ <span style="float: right;">/* Target Update */</span> $v_x \leftarrow \beta v_x + (1 - \beta) \nabla_{x_t} \mathcal{L}_{\text{ATK}}^{A_{\text{Text}}} (\mathcal{M} \circ g_{\theta}^{(t, \text{Text})}(x_t), L_{\text{Tar}})$ . <span style="float: right;">/* Momentum */</span> $x_t \leftarrow x_t + s * v_x$ . <span style="float: right;">/* Adversarial Optimization */</span>  <b>if</b> $t = t_s$ <b>then</b> $x'_{t-1} \leftarrow x_{t-1}$ <span style="float: right;">/* Determine Exemplar on Step <math>t_s</math> */</span> $x_{t-1} \leftarrow f_{\theta, 1}^{(t, \text{Text})}(x_t))(x_t)$ <span style="float: right;">/* DDIM Sampling Step for <math>x_{\text{adv}}</math> */</span> 2 $x'_{t-1} \leftarrow f_{\theta, 1}^{(t, \text{Text})}(x_t))(x'_t)$ . <span style="float: right;">/* DDIM Sampling Step for <math>x_{\text{exemplar}}</math> */</span> 3 $x_{t-1} \leftarrow x_{t-1} + \min\{\epsilon, \ x'_{t-1} - x_{t-1}\ _2\} \cdot \frac{x'_{t-1} - x_{t-1}}{\ x'_{t-1} - x_{t-1}\ _2}$ . <span style="float: right;">/* Semantic Constraint */</span> <b>return</b> $x_0, x'_0$ .
---	---

---

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  $\lceil \frac{t_s}{[T/K]} \rceil \times$  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  $K \mid T$  and  $\frac{T}{K} \mid t_s$ , 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  $x_{\text{adv}}$  and  $x_{\text{exemplar}}$  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  $z_0^{\text{slat}}$  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  $\mathbf{pos}$ ,  $\mathcal{D}_{\text{GS}}$  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  $\epsilon$  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  $\Sigma'$  as:

$$\Sigma' = JW\Sigma W^T J^T, \quad (24)$$

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  $(\alpha_i)$  is determined by a 2D Gaussian with covariance  $\Sigma$  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,  $g$  and the adversarial optimization is

represented as:

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

The camera settings are sampled based on the original configuration in *Trellis* framework that defines  $P_{\text{Cam}}$ , *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) [28], the inner-loop optimization is performed by averaging the gradient from different sampled cameras:

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

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.*,  $T = 50$ .

### 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

$$\begin{aligned} V' &= \text{TC}_G(L) := \left\{ v \in V \mid \exists u \in L, \exists k \in \mathbb{N}, (u \xrightarrow{k} v) \in E^+ \right\} \\ G' &:= (V', E') \quad \text{where } E' = \{(u, v) \in E \mid u, v \in L'\} \end{aligned} \quad (29)$$

Where  $\text{TC}$  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.

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

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 imagenet labels are shown in Table 6.

**Discussions on the ASR<sub>relative</sub> metric.** We design ASR<sub>relative</sub> 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<sub>relative</sub> evaluation metrics of SemanticAE generator and the application scenario of the multi-round reject-sampling based data generation pipeline.

**Proposition B.2** (*ASR<sub>relative</sub> characterize the upper-bound probability of the successful attack*). *For any adversarial sampling algorithm  $K$  that generates the SemanticAE with  $\text{ASR}_{\text{relative}} = p$  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  $L_{\text{Text}}$  is the correct label corresponding to Text and  $L_{\text{Text}} = \overline{A}_{\text{Text}}$ , and  $\delta$  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  $p$ ,  $x_{\text{adv}}$  is a successful attack. Since  $P(\mathcal{M}_S(x_{\text{exemplar}}) \in L_{\text{Text}}) > p_s$ , the final success attack rate is :

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

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 Since the proposed module belongs to the intersection of the diffusion-based adversarial attack and the optimization methodology, we select the classical MI-FGSM [38], 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 [16], VENOM [17] and SD-NAE [15]. We integrate an input-transformation-based method (DeCoWA [43]) 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  $t'_{t-\Delta T} = f_{\theta, \Delta T}(\arg \max_{x'_t} \mathcal{L}_{\text{ATK}}(\mathcal{M}(x'_t)))$ , 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  $\mathcal{T} = \{\text{ResNet50} [47], \text{ViT-B/16} [48], \text{ConvNeXt-T} [49], \text{ResNet152}, \text{InceptionV3} [50], \text{Swin-Transformer-B} [51]\}$ , and the surrogate model set as  $\{\text{ResNet50}, \text{ViT-B/16}, \text{ConvNext-T}, \text{ResNet50+DeCoWA}\}$ .

**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:

$$\mathcal{L}_{\text{Atk}}(\text{logits}, A_{\text{Text}}, L_{\text{Tar}}) = \text{LogSoftMax}(\text{logits})[L_{\text{Tar}}] - \frac{1}{|A_{\text{Text}}|} \sum_{i \in A_{\text{Text}}} \text{LogSoftMax}(\text{logits})[i], \quad (33)$$

Where  $L_{\text{Tar}}$  is the currently selected label (the label with the highest confidence in the set  $A_{\text{Text}}$ , 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  $x_{\text{exemplar}}$  and  $x_{\text{adv}}$ , which defines as follows in the main body of our work:

$$\text{SemanticDiff}_{\mathcal{S}} = \langle x_{\text{exemplar}}, x_{\text{adv}} \rangle_{\mathcal{S}}, \quad (34)$$

$\mathcal{S}$  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  $\text{Clip}_Q$ , we use the implementation of *piq* [71] 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 6, 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  $\mathcal{T}_{coarse}$  be the set of Abstracted Labels. The mapping function  $M : \mathbb{L} \rightarrow \mathcal{T}_{coarse}$  transforms each Original ImageNet label to its Abstracted Label. The set of prompts used for generation is  $\mathcal{P} = \{t | t \in \mathcal{T}_{coarse}\}$ .

**3D Video Generation** Two categories of video samples were generated: clean samples and adversarial examples. Clean video samples were synthesized using the TRELLIS[41] model. For each prompt  $p \in \mathcal{P}$ , a corresponding clean video  $V_{clean}$  was generated. Adversarial examples were generated based on the TRELLIS[41] 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  $l_{target} \in \mathbb{L}$  such that  $M(l_{target}) = p$ . A constant perturbation strength, denoted as  $\epsilon$ , 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 17, 18, 19, 20. Figure 17, 19 display the relationship between LPIPS and ASR/ACC for examples generated by different methods in the abstracted label task. Figure 18, 20 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 8 and Table 8, demonstrating the stability of the results.

**Analysis** :Based on the observations from Figure 17, 18, the following conclusions can be derived:

- In our method, as the perturbation parameter  $\epsilon$  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 17) 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 19, 20, all the methods, except for the VENOM method in the two subfigures of Figure 19(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 21, 22 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 21 vs. Figure 22.
- **AdvDiff**: As diffusion-based adversarial example generation method, it exhibits strong robustness against DiffPure (also diffusion-based). Figure 22 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 21). 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 22).

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 12. 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.

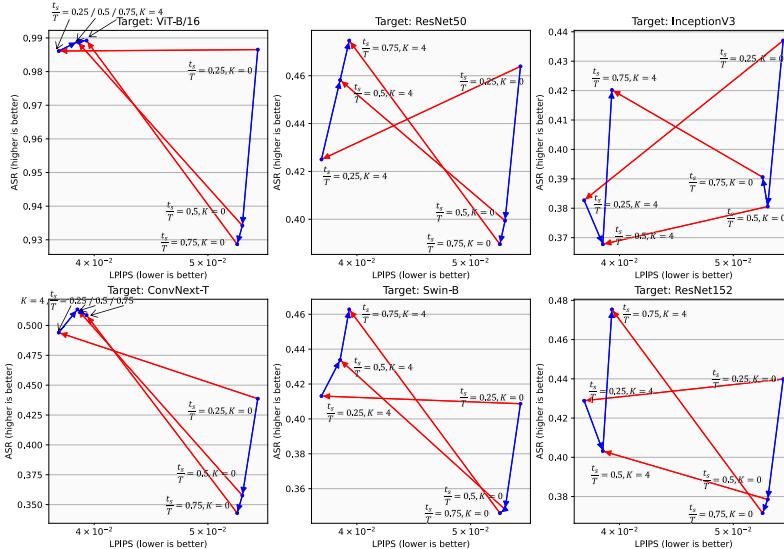


Figure 12: Parameter analysis with surrogate model ViT-B/16. Red arrows denote the performance before and after adding the residual approximation for  $g(x)$ . In each setting of  $\frac{t_s}{T} \geq 0.5$ , applying residual approximation achieves significant improvements (The upper left corner indicates 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 5, 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 13 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.

#### 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

Table 5: 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	<b>1</b>	1	1	1	1
Residual approximation step (K)	4	4	<b>4</b>	3	2	1	0
ASR	<b>0.922</b>	0.913	<b>0.922</b>	0.912	0.912	0.902	0.451



Figure 13: 3D Visual Results Ablation.

task. Additionally, we evaluate the clip-score [72] metric with the settings:

$$Prompt = \begin{cases} \text{ImageNetLabel} & Task = \text{ImageNet} \\ \text{AbstractedLabel}, \text{ImageNetLabel} & Task = \text{Abstracted} \end{cases} \quad (35)$$

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 14, Figure 15, and Figure 16. 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.

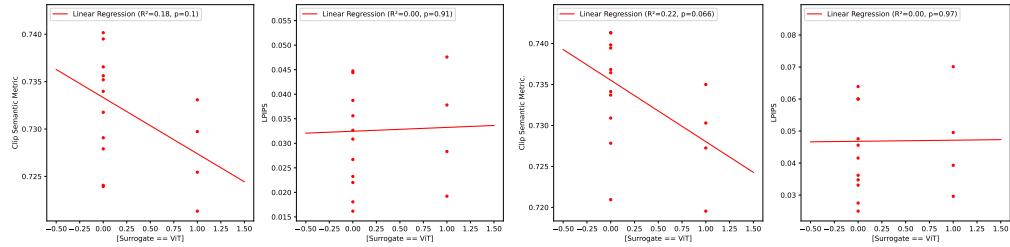


Figure 14: Results on the Original ImageNet Label Evasion SemanticAEs

Figure 15: Results on the Abstracted Label Evasion SemanticAEs

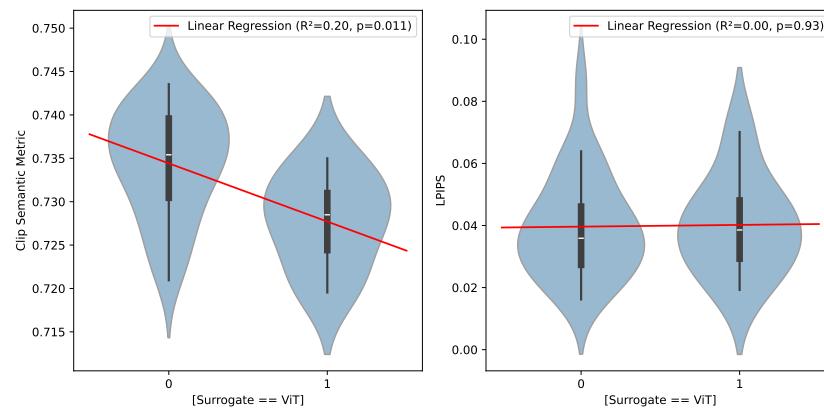


Figure 16: Combined distribution of the semantic metric on the ImageNet and abstracted label evasion SemanticAEs. With high confidence ( $p < 0.02$ ), the CLIP semantic metric is affected by the attack transferability.

## 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  $x_{\text{adv}}$  and corresponding nearby samples  $x_{\text{exemplar}}$  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 23 and 24. 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  $x_{\text{adv}}$  and  $x_{\text{exemplar}}$ ,  $x_{\text{exemplar}}$  is classified correctly (green), while  $x_{\text{adv}}$  is classified incorrectly (red).

**Analyzing 2D Visualized Results** Based on the observations of Figure 23 and 24, the following conclusions can be derived:

- For our attack method, as the parameter  $\epsilon$  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  $x_{\text{adv}}$  with the original ImageNet label "castle" has a castle with a blurred top when  $\epsilon = 3$ , and the  $x_{\text{adv}}$  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  $x_{\text{exemplar}}$  and  $x_{\text{adv}}$ . 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  $x_{\text{adv}}$  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 25. 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 13. 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 Social 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 6: 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

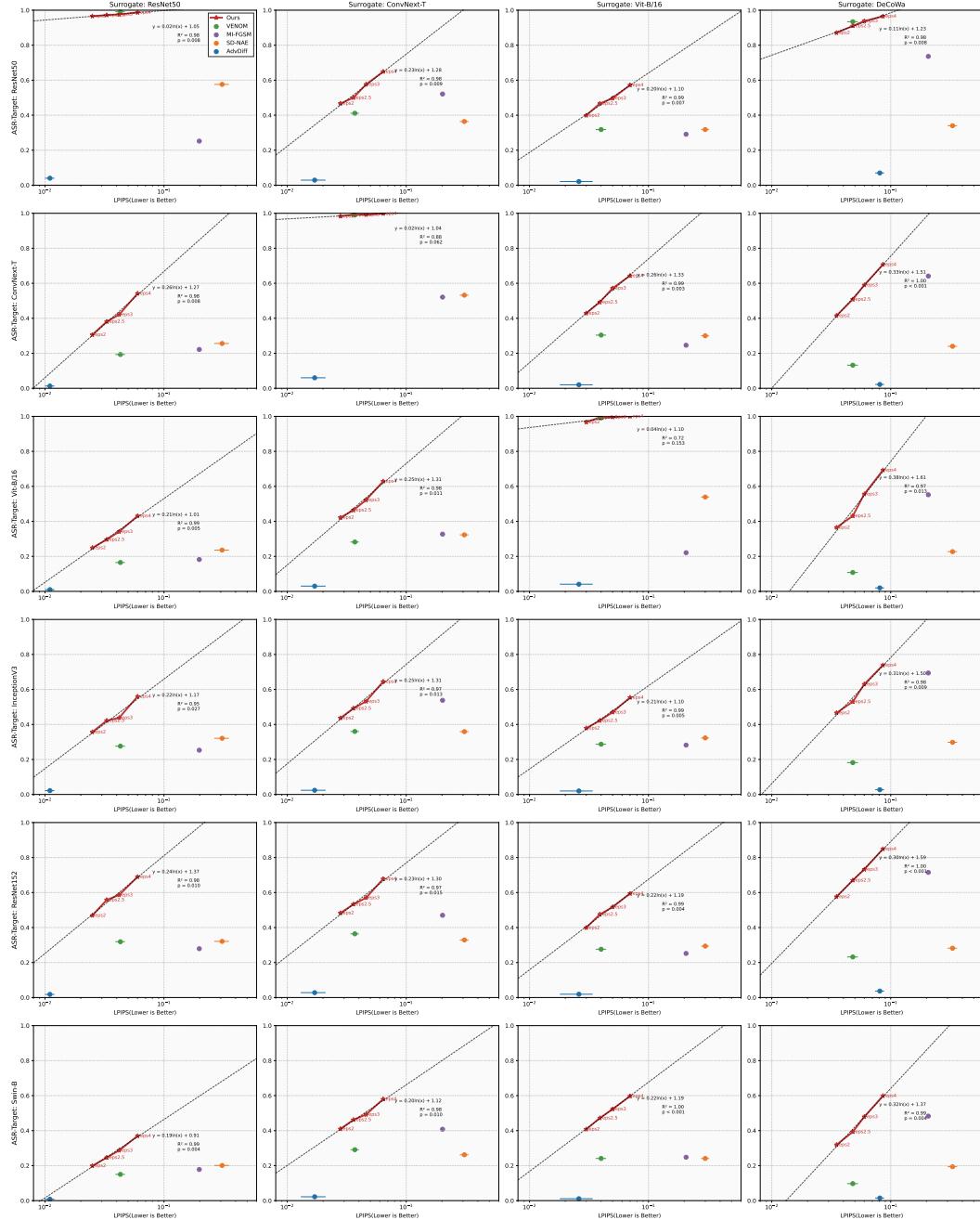


Figure 17: ASR of Abstracted Label

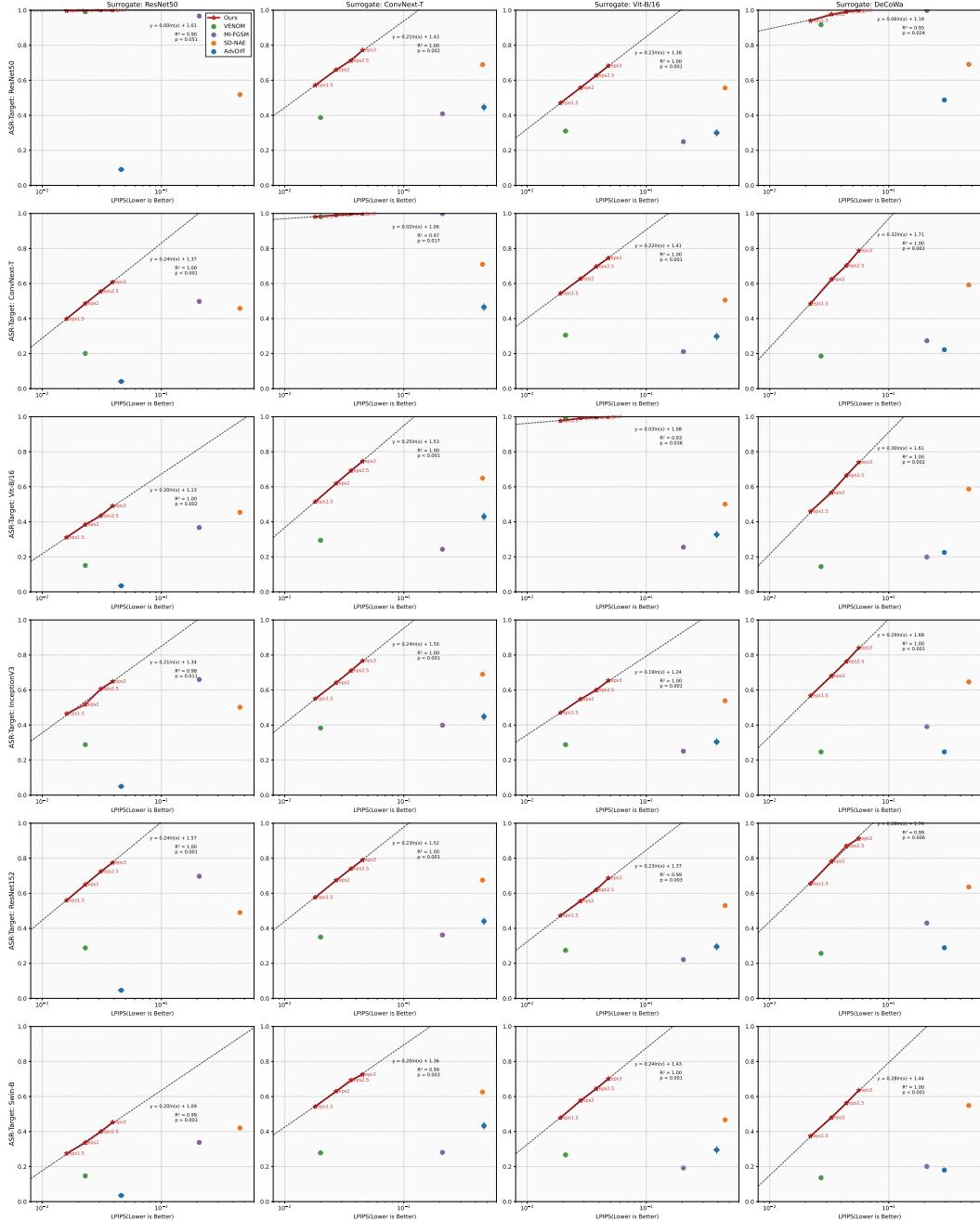


Figure 18: ASR of Original Imagenet label

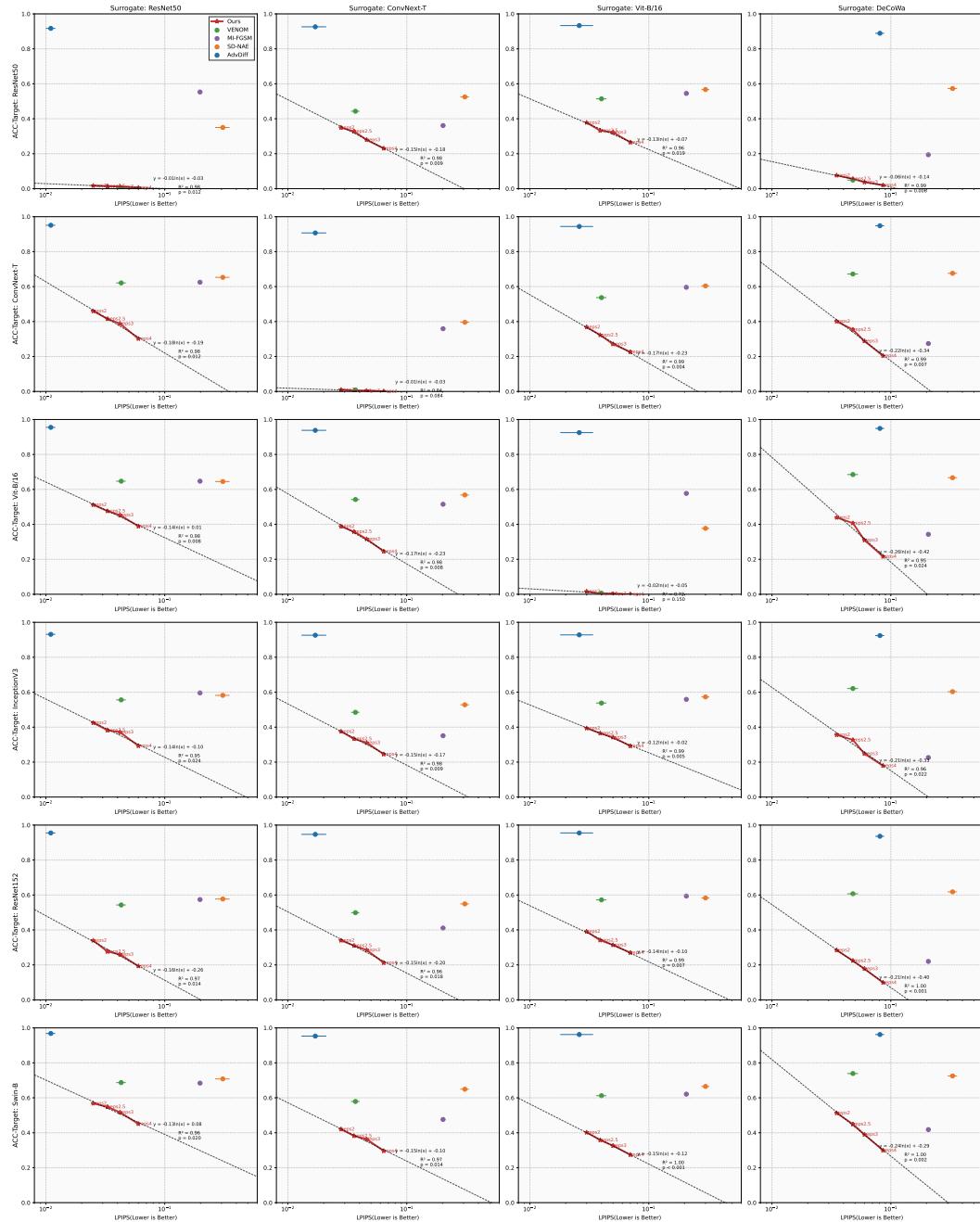


Figure 19: ACC of Abstracted Label

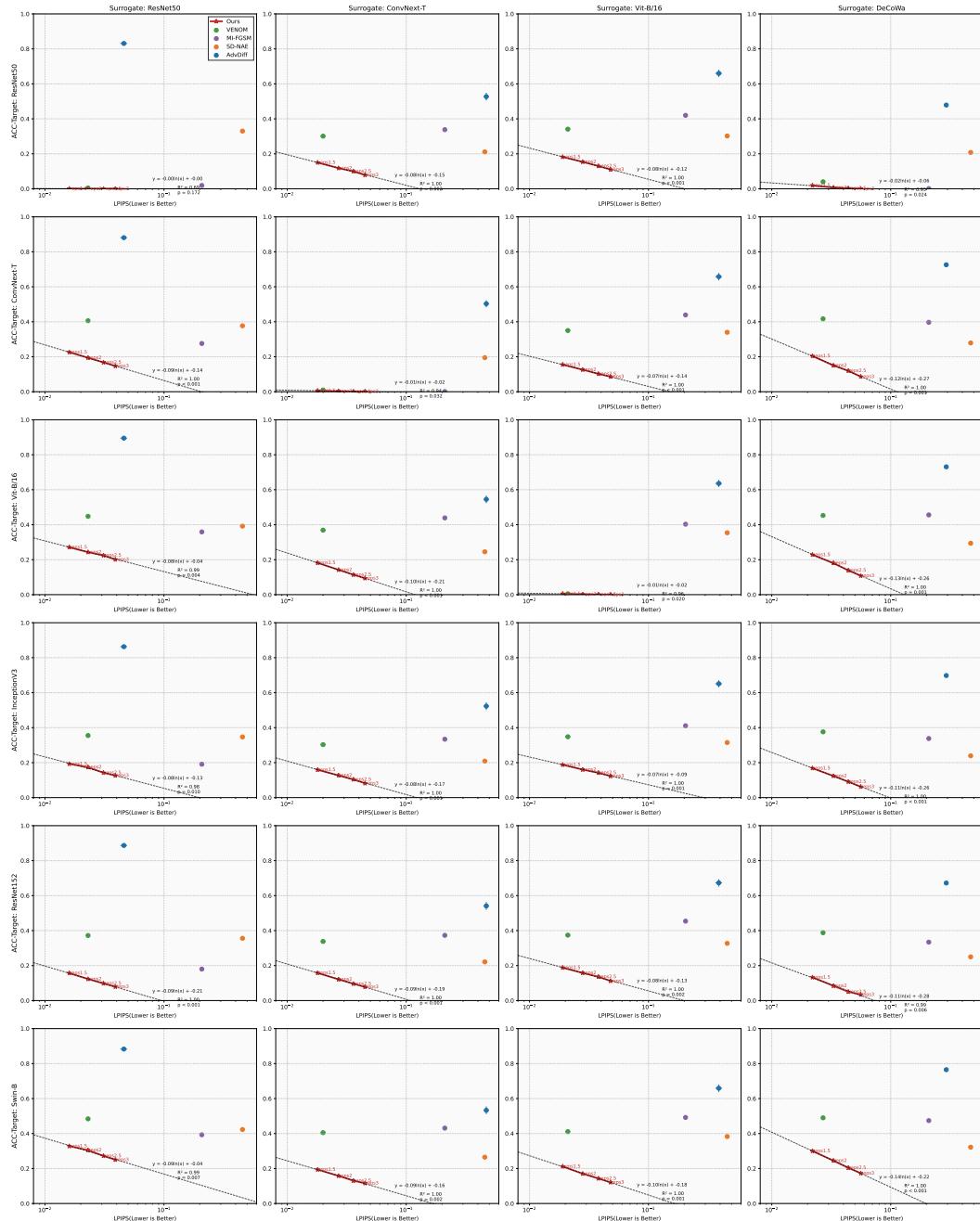


Figure 20: ACC of Original Imagenet label

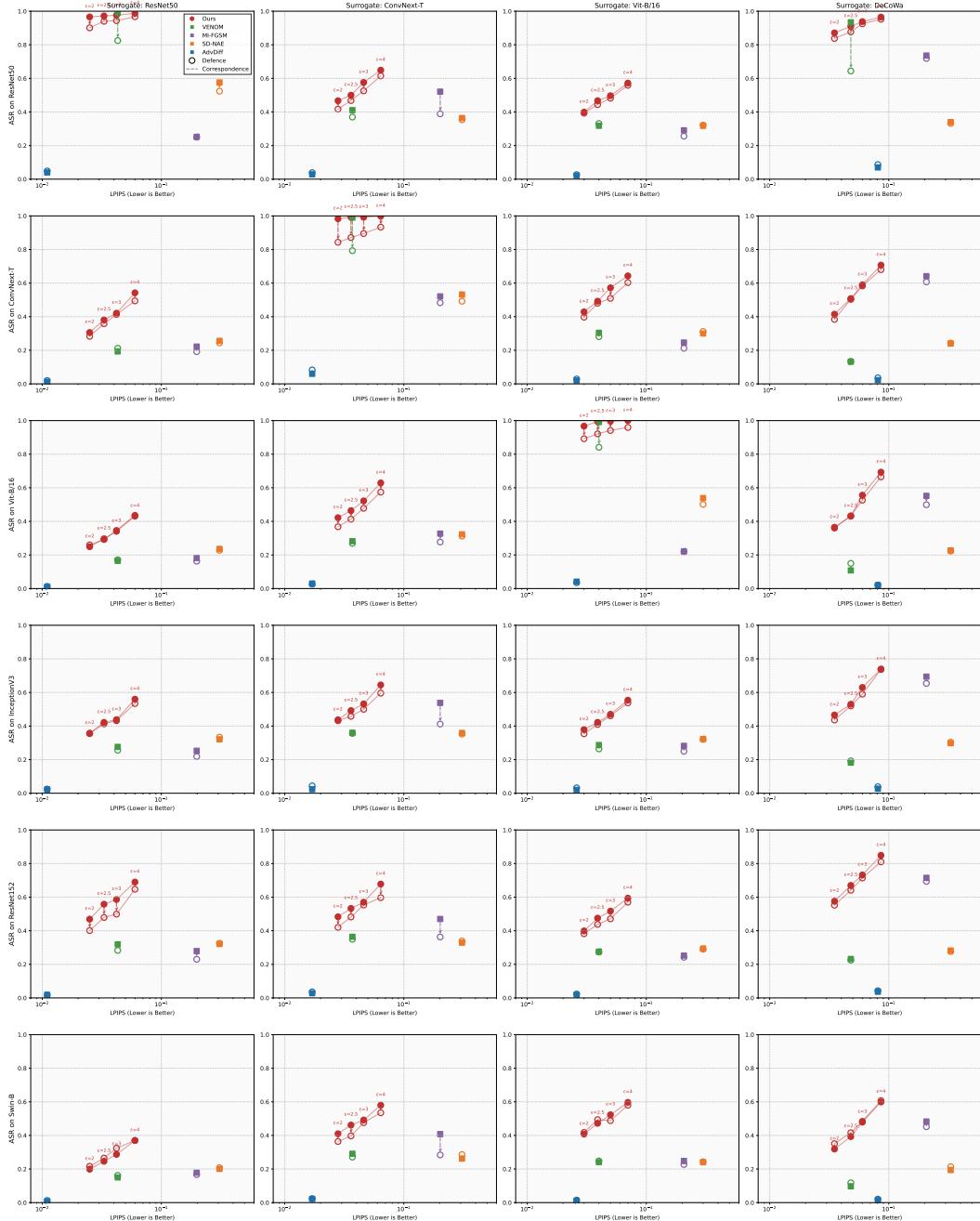


Figure 21: ASR on JPEG Defense

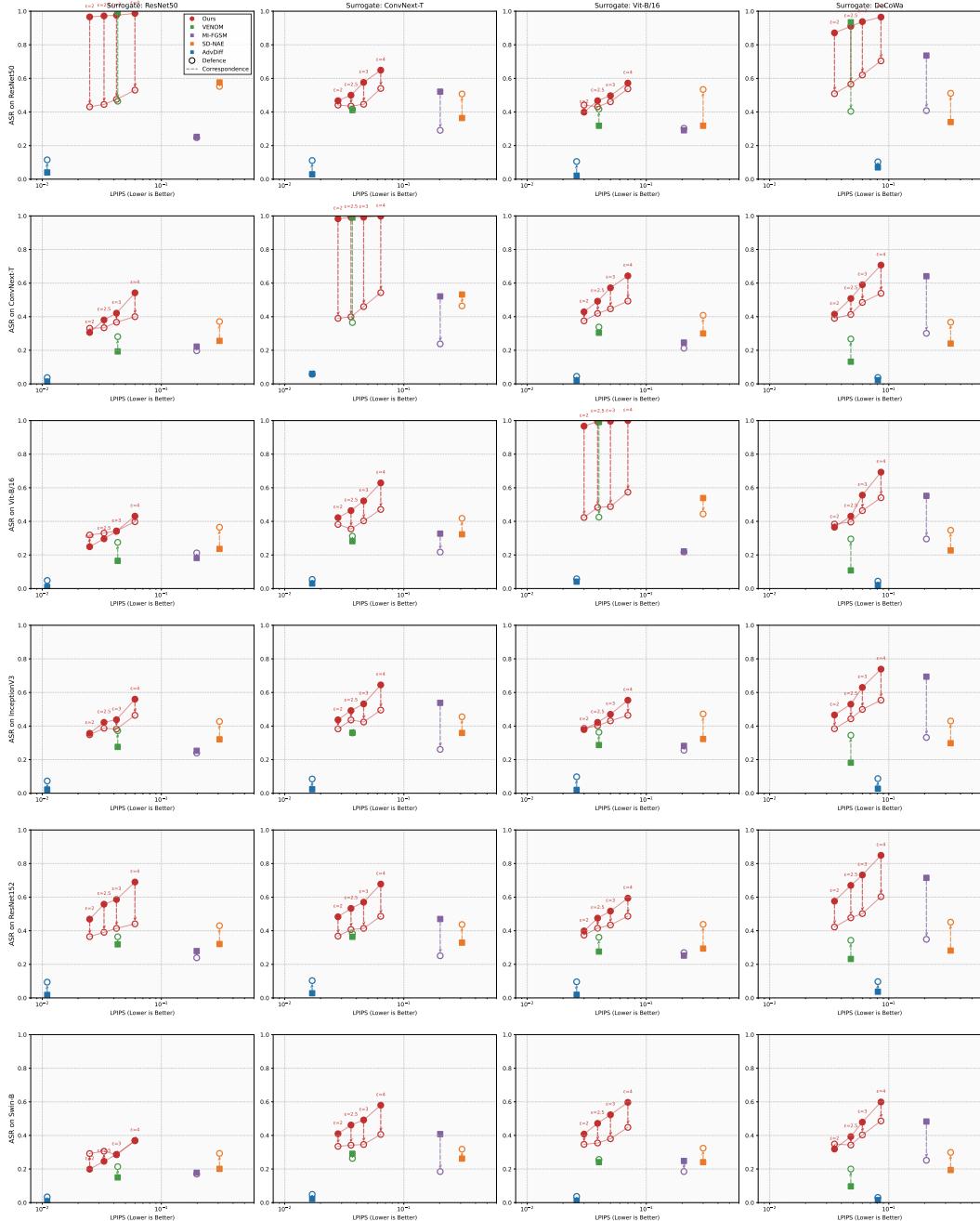


Figure 22: ASR on DiffPure Defense

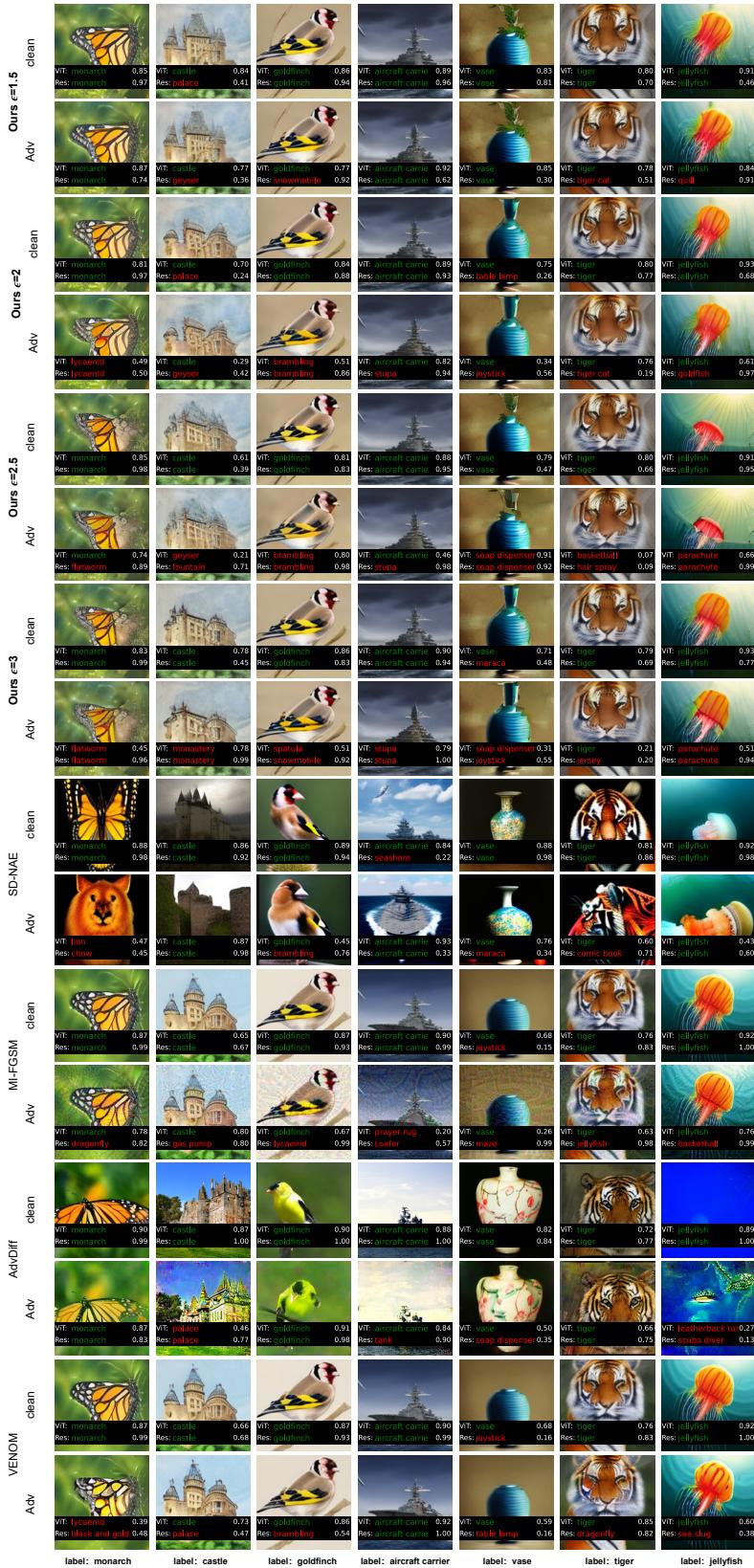


Figure 23: Visualization of 2D ImageNet-label Evasion Attacks. Surrogate model is ResNet50+DeCoWA.

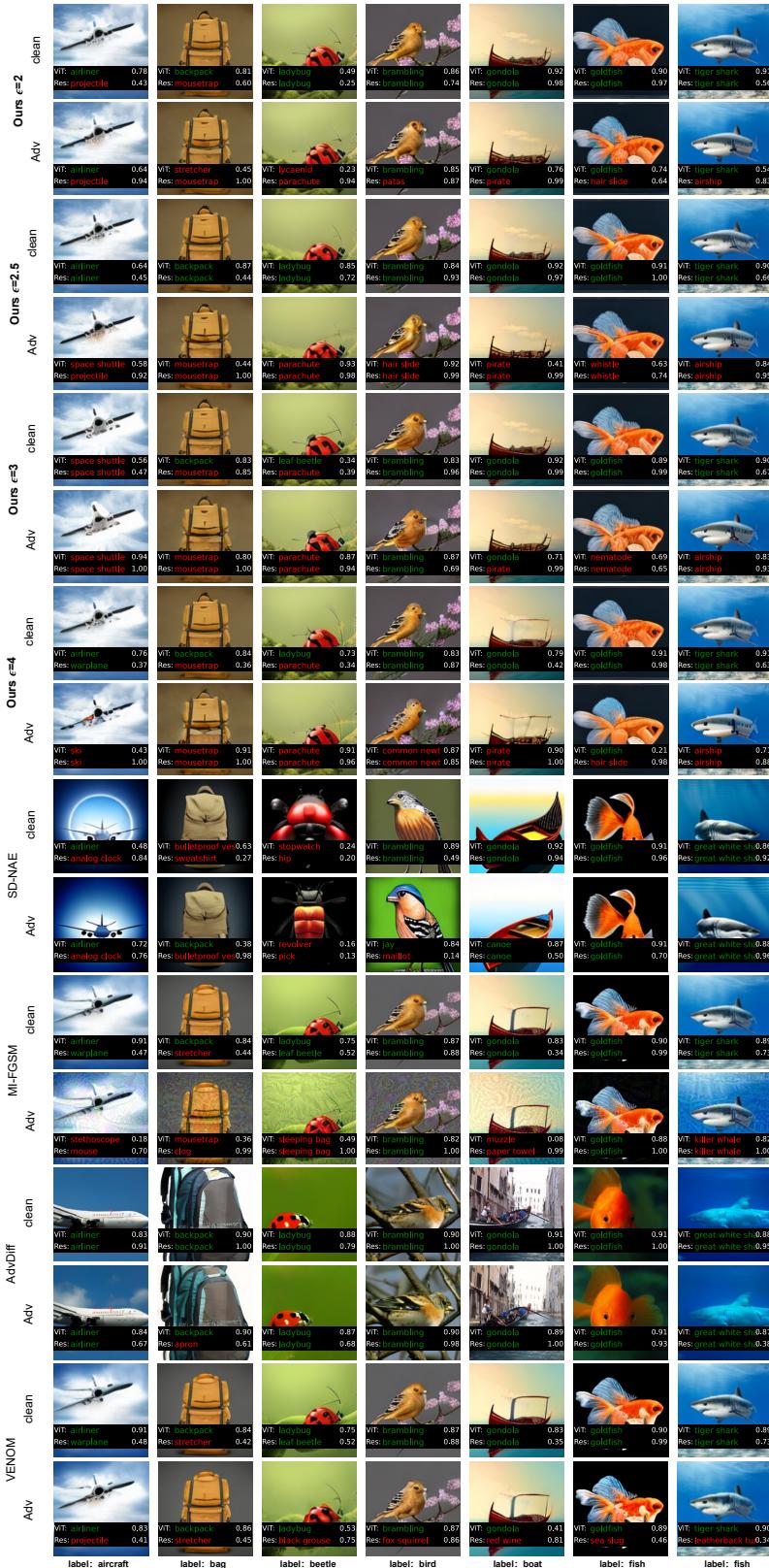


Figure 24: Visualization of 2D Abstracted-label Evasion Attacks. Surrogate model is ResNet50+DeCoWA.

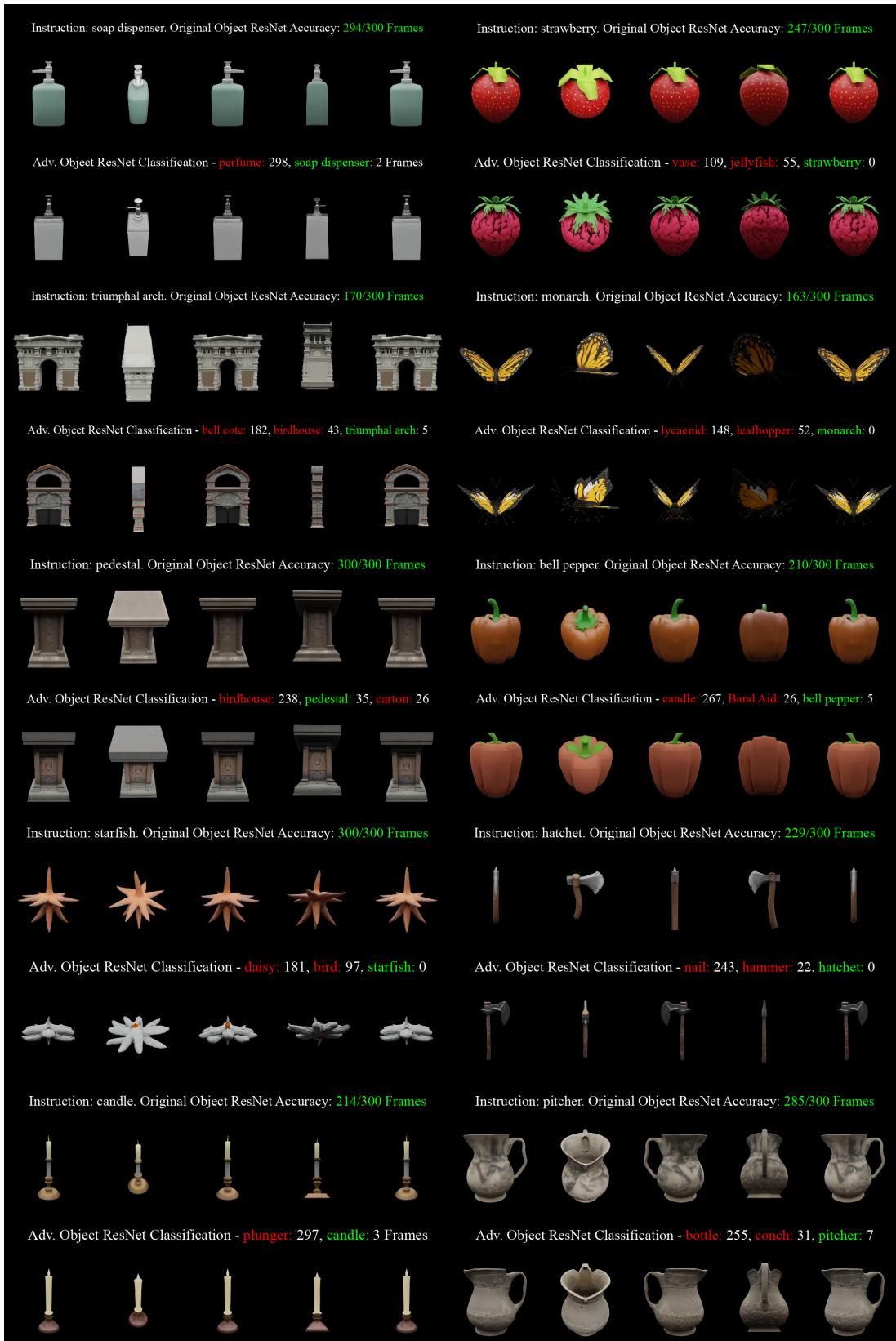


Figure 25: 3D Visual Results. Surrogate Model is ResNet50.

Table 7: Abstracted label transfer attack results.

method	Surrogate	Clip <sub>QA</sub>	Metrics			ResNet152			InceptionV3			ResNet50			ViT-B/16			ConvNext-T			Swin-B		
			LPIPS	MSSSIM	ASR	ACC																	
ours <sub>c=2</sub>	ResNet+DeCoVA	0.809 <sub>±0.009</sub>	0.035 <sub>±0.001</sub>	0.576 <sub>±0.001</sub>	0.284 <sub>±0.003</sub>	0.466 <sub>±0.002</sub>	0.356 <sub>±0.002</sub>	0.871 <sub>±0.003</sub>	0.076 <sub>±0.002</sub>	0.966 <sub>±0.001</sub>	0.017 <sub>±0.001</sub>	0.249 <sub>±0.005</sub>	0.365 <sub>±0.006</sub>	0.439 <sub>±0.005</sub>	0.401 <sub>±0.002</sub>	0.319 <sub>±0.005</sub>	0.513 <sub>±0.004</sub>	0.401 <sub>±0.002</sub>	0.319 <sub>±0.005</sub>	0.513 <sub>±0.004</sub>			
	ResNet50	0.813 <sub>±0.008</sub>	0.025 <sub>±0.001</sub>	0.968 <sub>±0.002</sub>	0.469 <sub>±0.003</sub>	0.340 <sub>±0.003</sub>	0.357 <sub>±0.006</sub>	0.425 <sub>±0.006</sub>	0.576 <sub>±0.001</sub>	0.016 <sub>±0.001</sub>	0.249 <sub>±0.005</sub>	0.512 <sub>±0.003</sub>	0.364 <sub>±0.005</sub>	0.406 <sub>±0.004</sub>	0.460 <sub>±0.003</sub>	0.199 <sub>±0.002</sub>	0.568 <sub>±0.002</sub>	0.201 <sub>±0.002</sub>	0.708 <sub>±0.002</sub>	0.653 <sub>±0.002</sub>			
	ViT-B/16	0.815 <sub>±0.012</sub>	0.030 <sub>±0.001</sub>	0.965 <sub>±0.001</sub>	0.399 <sub>±0.002</sub>	0.340 <sub>±0.003</sub>	0.379 <sub>±0.003</sub>	0.594 <sub>±0.003</sub>	0.594 <sub>±0.004</sub>	0.378 <sub>±0.004</sub>	0.967 <sub>±0.001</sub>	0.229 <sub>±0.004</sub>	0.368 <sub>±0.003</sub>	0.408 <sub>±0.005</sub>	0.422 <sub>±0.004</sub>	0.420 <sub>±0.004</sub>	0.410 <sub>±0.002</sub>	0.420 <sub>±0.002</sub>	0.420 <sub>±0.002</sub>	0.420 <sub>±0.002</sub>			
	ConnNext-T	0.814 <sub>±0.012</sub>	0.028 <sub>±0.001</sub>	0.967 <sub>±0.001</sub>	0.488 <sub>±0.006</sub>	0.340 <sub>±0.004</sub>	0.437 <sub>±0.005</sub>	0.488 <sub>±0.005</sub>	0.575 <sub>±0.004</sub>	0.349 <sub>±0.003</sub>	0.377 <sub>±0.004</sub>	0.530 <sub>±0.005</sub>	0.327 <sub>±0.003</sub>	0.422 <sub>±0.004</sub>	0.422 <sub>±0.004</sub>	0.420 <sub>±0.004</sub>	0.420 <sub>±0.004</sub>	0.420 <sub>±0.004</sub>	0.420 <sub>±0.004</sub>				
	ResNet+DeCoVA	0.808 <sub>±0.009</sub>	0.048 <sub>±0.001</sub>	0.943 <sub>±0.001</sub>	0.570 <sub>±0.003</sub>	0.224 <sub>±0.002</sub>	0.530 <sub>±0.005</sub>	0.327 <sub>±0.005</sub>	0.910 <sub>±0.002</sub>	0.056 <sub>±0.002</sub>	0.431 <sub>±0.001</sub>	0.206 <sub>±0.005</sub>	0.407 <sub>±0.001</sub>	0.058 <sub>±0.002</sub>	0.388 <sub>±0.003</sub>	0.449 <sub>±0.004</sub>	0.393 <sub>±0.004</sub>	0.449 <sub>±0.004</sub>	0.393 <sub>±0.004</sub>				
	ResNet50	0.808 <sub>±0.013</sub>	0.033 <sub>±0.001</sub>	0.958 <sub>±0.002</sub>	0.558 <sub>±0.002</sub>	0.276 <sub>±0.001</sub>	0.422 <sub>±0.007</sub>	0.382 <sub>±0.004</sub>	0.972 <sub>±0.001</sub>	0.014 <sub>±0.001</sub>	0.206 <sub>±0.005</sub>	0.476 <sub>±0.004</sub>	0.181 <sub>±0.004</sub>	0.415 <sub>±0.003</sub>	0.246 <sub>±0.002</sub>	0.548 <sub>±0.003</sub>	0.322 <sub>±0.003</sub>	0.426 <sub>±0.002</sub>	0.322 <sub>±0.003</sub>				
ours <sub>c=2.5</sub>	ViT-B/16	0.814 <sub>±0.012</sub>	0.036 <sub>±0.001</sub>	0.957 <sub>±0.001</sub>	0.533 <sub>±0.005</sub>	0.310 <sub>±0.002</sub>	0.492 <sub>±0.008</sub>	0.365 <sub>±0.002</sub>	0.500 <sub>±0.006</sub>	0.333 <sub>±0.006</sub>	0.500 <sub>±0.005</sub>	0.352 <sub>±0.006</sub>	0.464 <sub>±0.004</sub>	0.358 <sub>±0.003</sub>	0.994 <sub>±0.001</sub>	0.042 <sub>±0.001</sub>	0.421 <sub>±0.002</sub>	0.472 <sub>±0.002</sub>	0.382 <sub>±0.002</sub>				
	ConnNext-T	0.812 <sub>±0.008</sub>	0.060 <sub>±0.001</sub>	0.929 <sub>±0.002</sub>	0.732 <sub>±0.002</sub>	0.178 <sub>±0.001</sub>	0.630 <sub>±0.003</sub>	0.248 <sub>±0.003</sub>	0.938 <sub>±0.001</sub>	0.037 <sub>±0.001</sub>	0.576 <sub>±0.002</sub>	0.310 <sub>±0.002</sub>	0.590 <sub>±0.005</sub>	0.288 <sub>±0.003</sub>	0.479 <sub>±0.003</sub>	0.390 <sub>±0.003</sub>	0.517 <sub>±0.003</sub>	0.390 <sub>±0.003</sub>					
	ResNet+DeCoVA	0.808 <sub>±0.009</sub>	0.026 <sub>±0.001</sub>	0.949 <sub>±0.003</sub>	0.586 <sub>±0.002</sub>	0.314 <sub>±0.006</sub>	0.470 <sub>±0.004</sub>	0.342 <sub>±0.004</sub>	0.976 <sub>±0.001</sub>	0.012 <sub>±0.001</sub>	0.497 <sub>±0.005</sub>	0.332 <sub>±0.005</sub>	0.996 <sub>±0.001</sub>	0.002 <sub>±0.001</sub>	0.521 <sub>±0.004</sub>	0.523 <sub>±0.004</sub>	0.523 <sub>±0.004</sub>	0.523 <sub>±0.004</sub>					
	ResNet50	0.806 <sub>±0.011</sub>	0.044 <sub>±0.002</sub>	0.944 <sub>±0.001</sub>	0.517 <sub>±0.006</sub>	0.314 <sub>±0.004</sub>	0.438 <sub>±0.007</sub>	0.371 <sub>±0.003</sub>	0.971 <sub>±0.001</sub>	0.012 <sub>±0.001</sub>	0.497 <sub>±0.005</sub>	0.341 <sub>±0.005</sub>	0.996 <sub>±0.001</sub>	0.002 <sub>±0.001</sub>	0.521 <sub>±0.004</sub>	0.521 <sub>±0.004</sub>	0.521 <sub>±0.004</sub>	0.521 <sub>±0.004</sub>					
	ViT-B/16	0.811 <sub>±0.001</sub>	0.050 <sub>±0.002</sub>	0.950 <sub>±0.002</sub>	0.570 <sub>±0.007</sub>	0.285 <sub>±0.006</sub>	0.532 <sub>±0.006</sub>	0.311 <sub>±0.003</sub>	0.576 <sub>±0.008</sub>	0.279 <sub>±0.005</sub>	0.522 <sub>±0.004</sub>	0.316 <sub>±0.003</sub>	0.992 <sub>±0.001</sub>	0.004 <sub>±0.001</sub>	0.492 <sub>±0.006</sub>	0.363 <sub>±0.006</sub>	0.363 <sub>±0.006</sub>						
	ConnNext-T	0.808 <sub>±0.011</sub>	0.046 <sub>±0.002</sub>	0.946 <sub>±0.002</sub>	0.570 <sub>±0.007</sub>	0.270 <sub>±0.007</sub>	0.539 <sub>±0.003</sub>	0.317 <sub>±0.003</sub>	0.599 <sub>±0.002</sub>	0.270 <sub>±0.001</sub>	0.520 <sub>±0.002</sub>	0.317 <sub>±0.003</sub>	0.993 <sub>±0.001</sub>	0.004 <sub>±0.001</sub>	0.492 <sub>±0.006</sub>	0.363 <sub>±0.006</sub>	0.363 <sub>±0.006</sub>						
ours <sub>c=3</sub>	ResNet+DeCoVA	0.807 <sub>±0.012</sub>	0.086 <sub>±0.002</sub>	0.920 <sub>±0.002</sub>	0.639 <sub>±0.002</sub>	0.099 <sub>±0.001</sub>	0.344 <sub>±0.002</sub>	0.179 <sub>±0.003</sub>	0.965 <sub>±0.001</sub>	0.020 <sub>±0.001</sub>	0.431 <sub>±0.001</sub>	0.006 <sub>±0.001</sub>	0.431 <sub>±0.001</sub>	0.302 <sub>±0.003</sub>	0.570 <sub>±0.005</sub>	0.451 <sub>±0.005</sub>	0.451 <sub>±0.005</sub>						
	ResNet50	0.799 <sub>±0.010</sub>	0.060 <sub>±0.002</sub>	0.929 <sub>±0.003</sub>	0.690 <sub>±0.001</sub>	0.193 <sub>±0.001</sub>	0.560 <sub>±0.010</sub>	0.293 <sub>±0.005</sub>	0.987 <sub>±0.001</sub>	0.006 <sub>±0.001</sub>	0.576 <sub>±0.001</sub>	0.265 <sub>±0.001</sub>	1.000 <sub>±0.000</sub>	0.000 <sub>±0.000</sub>	0.643 <sub>±0.004</sub>	0.597 <sub>±0.005</sub>	0.274 <sub>±0.004</sub>						
	ViT-B/16	0.805 <sub>±0.012</sub>	0.070 <sub>±0.001</sub>	0.922 <sub>±0.002</sub>	0.594 <sub>±0.004</sub>	0.270 <sub>±0.003</sub>	0.554 <sub>±0.006</sub>	0.292 <sub>±0.003</sub>	0.572 <sub>±0.003</sub>	0.245 <sub>±0.005</sub>	0.649 <sub>±0.004</sub>	0.230 <sub>±0.003</sub>	0.629 <sub>±0.006</sub>	0.245 <sub>±0.006</sub>	0.580 <sub>±0.006</sub>	0.296 <sub>±0.006</sub>							
	ConnNext-T	0.808 <sub>±0.008</sub>	0.064 <sub>±0.002</sub>	0.926 <sub>±0.002</sub>	0.671 <sub>±0.006</sub>	0.220 <sub>±0.002</sub>	0.694 <sub>±0.012</sub>	0.226 <sub>±0.002</sub>	0.736 <sub>±0.002</sub>	0.194 <sub>±0.002</sub>	0.576 <sub>±0.003</sub>	0.221 <sub>±0.002</sub>	0.552 <sub>±0.003</sub>	0.343 <sub>±0.003</sub>	0.641 <sub>±0.003</sub>	0.483 <sub>±0.004</sub>	0.418 <sub>±0.003</sub>						
	ResNet+DeCoVA	0.874 <sub>±0.013</sub>	0.207 <sub>±0.008</sub>	0.870 <sub>±0.006</sub>	0.571 <sub>±0.002</sub>	0.220 <sub>±0.002</sub>	0.671 <sub>±0.012</sub>	0.220 <sub>±0.002</sub>	0.735 <sub>±0.002</sub>	0.194 <sub>±0.002</sub>	0.576 <sub>±0.003</sub>	0.221 <sub>±0.002</sub>	0.647 <sub>±0.002</sub>	0.222 <sub>±0.003</sub>	0.625 <sub>±0.004</sub>	0.178 <sub>±0.004</sub>	0.684 <sub>±0.004</sub>						
	ResNet50	0.551 <sub>±0.014</sub>	0.198 <sub>±0.009</sub>	0.885 <sub>±0.006</sub>	0.574 <sub>±0.003</sub>	0.279 <sub>±0.003</sub>	0.594 <sub>±0.003</sub>	0.253 <sub>±0.002</sub>	0.596 <sub>±0.002</sub>	0.252 <sub>±0.001</sub>	0.553 <sub>±0.003</sub>	0.182 <sub>±0.001</sub>	0.577 <sub>±0.001</sub>	0.221 <sub>±0.002</sub>	0.577 <sub>±0.002</sub>	0.248 <sub>±0.002</sub>							
mifgsm	ViT-B/16	0.521 <sub>±0.013</sub>	0.207 <sub>±0.008</sub>	0.860 <sub>±0.006</sub>	0.252 <sub>±0.004</sub>	0.594 <sub>±0.003</sub>	0.282 <sub>±0.004</sub>	0.559 <sub>±0.004</sub>	0.291 <sub>±0.002</sub>	0.545 <sub>±0.001</sub>	0.291 <sub>±0.002</sub>	0.545 <sub>±0.001</sub>	0.221 <sub>±0.002</sub>	0.546 <sub>±0.002</sub>	0.246 <sub>±0.002</sub>	0.621 <sub>±0.002</sub>							
	ConnNext-T	0.532 <sub>±0.012</sub>	0.202 <sub>±0.009</sub>	0.858 <sub>±0.005</sub>	0.470 <sub>±0.003</sub>	0.411 <sub>±0.003</sub>	0.538 <sub>±0.003</sub>	0.351 <sub>±0.002</sub>	0.521 <sub>±0.003</sub>	0.361 <sub>±0.002</sub>	0.521 <sub>±0.002</sub>	0.327 <sub>±0.002</sub>	0.515 <sub>±0.002</sub>	0.359 <sub>±0.002</sub>	0.408 <sub>±0.002</sub>	0.476 <sub>±0.001</sub>							
	ResNet+DeCoVA	0.799 <sub>±0.016</sub>	0.048 <sub>±0.005</sub>	0.944 <sub>±0.005</sub>	0.232 <sub>±0.004</sub>	0.607 <sub>±0.003</sub>	0.182 <sub>±0.002</sub>	0.621 <sub>±0.001</sub>	0.934 <sub>±0.002</sub>	0.049 <sub>±0.002</sub>	0.685 <sub>±0.004</sub>	0.132 <sub>±0.002</sub>	0.672 <sub>±0.002</sub>	0.097 <sub>±0.001</sub>	0.739 <sub>±0.002</sub>								
	ResNet50	0.799 <sub>±0.012</sub>	0.043 <sub>±0.004</sub>	0.951 <sub>±0.004</sub>	0.543 <sub>±0.004</sub>	0.572 <sub>±0.003</sub>	0.276 <sub>±0.004</sub>	0.556 <sub>±0.004</sub>	0.931 <sub>±0.001</sub>	0.006 <sub>±0.001</sub>	0.615 <sub>±0.003</sub>	0.647 <sub>±0.004</sub>	0.193 <sub>±0.003</sub>	0.621 <sub>±0.002</sub>	0.150 <sub>±0.001</sub>	0.687 <sub>±0.003</sub>							
	ViT-B/16	0.780 <sub>±0.016</sub>	0.040 <sub>±0.004</sub>	0.958 <sub>±0.004</sub>	0.276 <sub>±0.003</sub>	0.572 <sub>±0.002</sub>	0.287 <sub>±0.004</sub>	0.538 <sub>±0.004</sub>	0.918 <sub>±0.001</sub>	0.014 <sub>±0.001</sub>	0.514 <sub>±0.003</sub>	0.991 <sub>±0.001</sub>	0.006 <sub>±0.001</sub>	0.304 <sub>±0.002</sub>	0.537 <sub>±0.001</sub>	0.612 <sub>±0.002</sub>							
	ConnNext-T	0.785 <sub>±0.010</sub>	0.037 <sub>±0.003</sub>	0.961 <sub>±0.002</sub>	0.364 <sub>±0.004</sub>	0.499 <sub>±0.004</sub>	0.360 <sub>±0.003</sub>	0.485 <sub>±0.002</sub>	0.411 <sub>±0.003</sub>	0.443 <sub>±0.002</sub>	0.411 <sub>±0.002</sub>	0.282 <sub>±0.003</sub>	0.542 <sub>±0.002</sub>	0.090 <sub>±0.001</sub>	0.291 <sub>±0.003</sub>								
SD-NAE	ResNet+DeCoVA	0.782 <sub>±0.010</sub>	0.331 <sub>±0.030</sub>	0.568 <sub>±0.034</sub>	0.518 <sub>±0.002</sub>	0.521 <sub>±0.002</sub>	0.321 <sub>±0.002</sub>	0.532 <sub>±0.003</sub>	0.582 <sub>±0.002</sub>	0.503 <sub>±0.002</sub>	0.573 <sub>±0.003</sub>	0.526 <sub>±0.003</sub>	0.567 <sub>±0.004</sub>	0.645 <sub>±0.001</sub>	0.526 <sub>±0.004</sub>	0.604 <sub>±0.002</sub>	0.604 <sub>±0.002</sub>						
	ViT-B/16	0.787 <sub>±0.010</sub>	0.300 <sub>±0.023</sub>	0.609 <sub>±0.030</sub>	0.294 <sub>±0.002</sub>	0.583 <sub>±0.004</sub>	0.323 <sub>±0.004</sub>	0.578 <sub>±0.003</sub>	0.582 <sub>±0.004</sub>	0.564 <sub>±0.003</sub>	0.584 <sub>±0.003</sub>	0.564 <sub>±0.003</sub>	0.586 <sub>±0.003</sub>	0.377 <sub>±0.003</sub>	0.665 <sub>±0.003</sub>								
	ConnNext-T	0.782 <sub>±0.012</sub>	0.308 <sub>±0.026</sub>	0.603 <sub>±0.033</sub>	0.329 <sub>±0.002</sub>	0.549 <sub>±0.002</sub>	0.359 <sub>±0.003</sub>	0.582 <sub>±0.002</sub>	0.585 <sub>±0.003</sub>	0.564 <sub>±0.002</sub>	0.585 <sub>±0.002</sub>	0.565 <sub>±0.002</sub>	0.586 <sub>±0.002</sub>	0.396 <sub>±0.003</sub>	0.626 <sub>±0.002</sub>								
	ResNet+DeCoVA	0.629 <sub>±0.007</sub>	0.088 <sub>±0.007</sub>	0.934 <sub>±0.007</sub>	0.037 <sub>±0.001</sub>	0.936 <sub>±0.002</sub>	0.027 <sub>±0.001</sub>	0.924 <sub>±0.001</sub>	0.070 <sub>±0.002</sub>	0.889 <sub>±0.003</sub>	0.020 <sub>±0.001</sub>	0.949 <sub>±0.002</sub>	0.022 <sub>±0.001</sub>	0.948 <sub>±0.002</sub>	0.015 <sub>±0.001</sub>	0.962 <sub>±0.001</sub>							
	ResNet50	0.621 <sub>±0.007</sub>	0.062 <sub>±0.007</sub>	0.911 <sub>±0.001</sub>	0.992 <sub>±0.001</sub>	0.018 <sub>±0.001</sub>	0.955 <sub>±0.003</sub>	0.022 <sub>±0.001</sub>	0.931 <sub>±0.001</sub>	0.040 <sub>±0.002</sub>	0.917 <sub>±0.003</sub>	0.0											

Table 8: 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 <sub>e=1.5</sub>	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	0.398 $\pm$ 0.002
	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 $\pm$ 0.004	0.478 $\pm$ 0.001	0.212 $\pm$ 0.001	0.470 $\pm$ 0.005	0.182 $\pm$ 0.005	0.153 $\pm$ 0.002
ours <sub>e=2</sub>	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	0.348 $\pm$ 0.003	0.160 $\pm$ 0.001	0.577 $\pm$ 0.001	0.154 $\pm$ 0.004	0.486 $\pm$ 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	0.104 $\pm$ 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
ours <sub>e=2.5</sub>	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	0.120 $\pm$ 0.001	0.997 $\pm$ 0.000	0.001 $\pm$ 0.000	0.692 $\pm$ 0.002
	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	0.109 $\pm$ 0.001	0.562 $\pm$ 0.002	0.204 $\pm$ 0.001	0.991 $\pm$ 0.001	0.003 $\pm$ 0.000	0.664 $\pm$ 0.002
	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	0.555 $\pm$ 0.001
	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 $\pm$ 0.002	0.143 $\pm$ 0.001	0.677 $\pm$ 0.004	0.130 $\pm$ 0.001	0.697 $\pm$ 0.001
ours <sub>e=3</sub>	ConvNext-T	0.814 $\pm$ 0.003	0.045 $\pm$ 0.000	0.947 $\pm$ 0.001	0.793 $\pm$ 0.001	0.078 $\pm$ 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	0.745 $\pm$ 0.002
	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	0.740 $\pm$ 0.003
	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	0.491 $\pm$ 0.002
	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	0.455 $\pm$ 0.002
MI-FGSM	ConvNext-T	0.543 $\pm$ 0.004	0.211 $\pm$ 0.001	0.877 $\pm$ 0.000	0.334 $\pm$ 0.001	0.373 $\pm$ 0.002	0.334 $\pm$ 0.002	0.334 $\pm$ 0.002	0.281 $\pm$ 0.002	0.431 $\pm$ 0.004	0.438 $\pm$ 0.001	0.999 $\pm$ 0.001	0.439 $\pm$ 0.001
	ResNet-DecoWa	0.548 $\pm$ 0.003	0.209 $\pm$ 0.002	0.880 $\pm$ 0.001	0.430 $\pm$ 0.003	0.226 $\pm$ 0.002	0.372 $\pm$ 0.001	0.221 $\pm$ 0.001	0.201 $\pm$ 0.002	0.474 $\pm$ 0.001	0.999 $\pm$ 0.000	0.001 $\pm$ 0.000	0.456 $\pm$ 0.001
	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	0.191 $\pm$ 0.001	0.358 $\pm$ 0.002	0.393 $\pm$ 0.001	0.967 $\pm$ 0.001	0.019 $\pm$ 0.000	0.359 $\pm$ 0.002
	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 $\pm$ 0.002	0.240 $\pm$ 0.001	0.111 $\pm$ 0.001	0.439 $\pm$ 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 $\pm$ 0.001	0.383 $\pm$ 0.002	0.347 $\pm$ 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	0.372 $\pm$ 0.002	0.376 $\pm$ 0.001	0.490 $\pm$ 0.002	0.917 $\pm$ 0.001	0.040 $\pm$ 0.001	0.369 $\pm$ 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	0.355 $\pm$ 0.002	0.355 $\pm$ 0.001	0.147 $\pm$ 0.001	0.484 $\pm$ 0.001	0.991 $\pm$ 0.001	0.005 $\pm$ 0.000	0.448 $\pm$ 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	0.348 $\pm$ 0.002	0.348 $\pm$ 0.002	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	0.265 $\pm$ 0.001	0.711 $\pm$ 0.001	0.195 $\pm$ 0.001	0.649 $\pm$ 0.002
	ResNet-DecoWa	0.845 $\pm$ 0.013	0.047 $\pm$ 0.002	0.433 $\pm$ 0.017	0.490 $\pm$ 0.001	0.250 $\pm$ 0.001	0.502 $\pm$ 0.001	0.347 $\pm$ 0.001	0.420 $\pm$ 0.001	0.423 $\pm$ 0.001	0.518 $\pm$ 0.001	0.048 $\pm$ 0.001	0.294 $\pm$ 0.001
	ResNet50	0.844 $\pm$ 0.011	0.049 $\pm$ 0.009	0.441 $\pm$ 0.016	0.530 $\pm$ 0.002	0.327 $\pm$ 0.001	0.539 $\pm$ 0.001	0.467 $\pm$ 0.001	0.383 $\pm$ 0.001	0.536 $\pm$ 0.002	0.302 $\pm$ 0.001	0.505 $\pm$ 0.001	0.392 $\pm$ 0.001
	Vit-B/16	0.844 $\pm$ 0.016	0.047 $\pm$ 0.009	0.441 $\pm$ 0.014	0.530 $\pm$ 0.002	0.327 $\pm$ 0.001	0.539 $\pm$ 0.001	0.467 $\pm$ 0.001	0.383 $\pm$ 0.001	0.536 $\pm$ 0.002	0.302 $\pm$ 0.001	0.501 $\pm$ 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	0.466 $\pm$ 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	0.887 $\pm$ 0.001	0.050 $\pm$ 0.001	0.863 $\pm$ 0.001	0.035 $\pm$ 0.001	0.883 $\pm$ 0.002	0.091 $\pm$ 0.001	0.941 $\pm$ 0.001	0.081 $\pm$ 0.001	0.353 $\pm$ 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