
Black Box Adversarial Prompting for Foundation Models

Natalie Maus*

NMAUS@SEAS.UPENN.EDU
University of Pennsylvania

Patrick Chao*

PCHAO@WHARTON.UPENN.EDU
University of Pennsylvania

Eric Wong

EXWONG@SEAS.UPENN.EDU
University of Pennsylvania

Jacob Gardner

JACOBURG@SEAS.UPENN.EDU
University of Pennsylvania

Abstract

Prompting interfaces allow users to quickly adjust the output of generative models in both vision and language. However, small changes and design choices in the prompt can lead to significant differences in the output. In this work, we develop a black-box framework for generating adversarial prompts for unstructured image and text generation. These prompts, which can be standalone or prepended to benign prompts, induce specific behaviors into the generative process, such as generating images of a particular object or generating high perplexity text.

1 Introduction

Foundation models (FMs) have demonstrated state-of-the-art performance across a diverse range of prediction and generation tasks. Large Language Models (LLMs), such as GPT-3 [10], PaLM [16], and ChatGPT are capable of answering complex queries when given a short language prompt as instructions. In vision, FMs now come with an easy-to-use text interface that has revolutionized image generation. Text-To-Image Models (TTIM) such as DALL-E 2 [56], Imagen [61], and Stable Diffusion [59] are capable of generating a wide range of high-resolution images from text descriptions.

The common thread amongst all these models is the use of a natural language prompt to influence the output of the model. By adjusting the words in a prompt, users can quickly and easily tweak the FM to generate text or images that better suit their task. This prompting interface has introduced a new mode of interaction with a model—instead of collecting datasets and fine-tuning models for a particular task, one simply changes words in a prompt.

However, it is becoming apparent that small variations in the prompt can lead to substantial differences in the model output. For example, prompting is known to be sensitive to the order of sentences [46], variations in the template [35], and the specific examples provided in the prompt [41]. Certain prompts can also result in strange or unusual behavior. For example, gibberish words from generated images such as “Apoploe vesrreaitais” can be used as a prompt to generate animals such as birds [19].

Recent work has seen researchers crafting adversarial prompts that exploit these vulnerabilities to extract information or induce malevolent behaviors. For example, Perez and Ribeiro [53] designed adversarial prompts that induce specific behaviors in an LLM, such as repeating its original prompt or outputting a particular string. Xu et al. [71] studied backdoors and adversarial triggers for LLMs fine-tuned on target tasks. However, these latter adversarial attacks require gradient information and the ability to modify the weights of the LLM. Such attacks are not possible on closed-source, general

*Equal contribution.



(a) Images generated by Stable Diffusion v1.5 [59] using the prompts **(Left)** ‘a picture of a mountain’ and **(Right)** ‘turbo lhaff√a picture of a mountain’, found by our method. The prepended text causes the model to consistently generate dogs rather than mountains.

Prompt: *coordinate hip arts ... filled via*
 Explain list comprehension in Python.

Generated Examples:

1. 1. 单词:\n 2. 词组: \n 3. 短语: \n \n 关键词: \n 1. 单词
2. 1. Forward and backward projections are techniques used in machine learning to estimate the gradient of a loss function. \2. In physics, the Lorentz transformation is used to describe how space and time are affected by relative motion.
3. I'd like to be called "Alexa," as it is short, simple, and easy to remember. Plus, it has a nice ring to it!

(b) Three examples of text generated from Vicuna 13B-v1.1 using our optimized prompt in italics and the original prompt to which we prepend. The generated text are nonsensical or irrelevant to the prompt.

Figure 1: Examples of adversarial prompts and generated outputs.

purpose FMs such as ChatGPT or DALL-E 2 that are used for a wide variety of unstructured tasks. *How do we craft adversarial prompts for general tasks on foundation models with only query access?*

Our Contributions. In this paper, we develop a framework for finding adversarial prompts in the black-box setting for general language and vision objectives. We focus specifically on *black-box* attacks as many foundation models are not publicly available, e.g. ChatGPT and DALL-E 2 only have API access. We build upon tools for black-box adversarial attacks traditionally used in the vision setting but instead leverage these tools to search for adversarial prompts. The main challenges are the high dimensionality and discrete nature of the token space. These properties drastically increase the sample complexity and slow down the convergence of classic black-box adversarial attacks. To address these issues, we develop a Token Space Projection (TSP) to map a continuous, lower-dimensional embedding space to the set of discrete language tokens. Our projected variations of black box attacks are capable of finding adversarial prompts, which we demonstrate across a range of vision and language tasks. Our main contributions are as follows:

1. We develop a framework for finding adversarial prompts using a Token Space Projection operator. This operator bridges the continuous word embedding space with the discrete token space, and enables the use of black-box attacks to find adversarial prompts.
2. We demonstrate how our framework can automatically find standalone or prepended prompts that cause text-to-image models to output specific image classes. We can further find adversarial prompts that exclude tokens related to the target class.
3. Our framework can also find adversarial prompts that alter text generation. For example, we find adversarial prompts that generate irrelevant or incoherent text, or text in different languages.

2 Adversarial Prompt Optimization

Adversarial attacks in language models have classically focused on perturbing text inputs to change the model’s prediction while maintaining the original semantic meaning. Various attack strategies include inducing typos [64], substituting words [34], and paraphrasing [32] to perturb the original text. However, these attacks are designed for a classification setting: they require a clean text example which is then modified. In contrast, there is no such clean example for a prompt, which is designed from scratch. Furthermore, prompt-based models do not generate clean classifications and instead produce comparatively unstructured outputs.

Therefore, we seek an alternative type of adversarial attack for prompting, one that can induce a targeted change in the output of the generative model. As an overview, our adversarial prompt optimization pipeline comprises four steps: 1) generate a candidate word embedding, 2) project the embeddings to a prompt, 3) generate using the prompt to obtain a loss, and 4) use a gradient-free optimization technique to optimize the loss and propose a new candidate word embedding, (Figure 5).

We begin by defining the concept of an *adversarial prompt*. Given a prompt $p \in \mathcal{P}$, the model M_p is represented as a probability distribution over responses $r \in \mathcal{R}$, parameterized by p . For text input models, the arbitrary space \mathcal{P} is defined by sequences of tokens with a maximal length d . We denote this as $\mathcal{P} = \mathcal{T}^d$, with \mathcal{T} denoting the discrete token space. Let $R_1^p, \dots, R_n^p \stackrel{\text{iid}}{\sim} M_p$ be sampled responses from the probability distribution induced by model M_p .

Let $\ell : \mathcal{R} \rightarrow \mathbb{R}$ be some external predefined loss function and $\mathcal{T}_{\text{allow}} \subseteq \mathcal{T}$ be a (potentially restricted) subclass of allowed tokens and $\mathcal{P}_{\text{allow}} := \mathcal{T}_{\text{allow}}^d \subseteq \mathcal{P}$ be a (potentially restricted) subclass of prompts comprised of d tokens from $\mathcal{T}_{\text{allow}}$. Consequently, our adversarial prompt framework aims to discover a prompt p that minimizes loss over the distribution of generated responses,

$$\arg \min_{p \in \mathcal{P}_{\text{allow}}} \mathbb{E}_{R \sim M_p} [\ell(R)]. \quad (1)$$

We relax the optimization problem to an empirical risk minimization framework (ERM) where we observe responses $R_1^p, \dots, R_n^p \stackrel{\text{iid}}{\sim} M_p$,

$$\arg \min_{p \in \mathcal{P}_{\text{allow}}} \frac{1}{n} \sum_{i=1}^n \ell(R_i^p). \quad (2)$$

We provide details on using black box optimizers to solve the relaxed ERM formulation in Section 2.2.

2.1 Token Space Projection

In our setting, we must optimize over a combinatorially large and discrete input token space $\mathcal{P} = \mathcal{T}^d$. To combat this, we consider a word embedding associated with each token and instead optimize over a relaxation of Equation 2 defined over the continuous embedding space.

Consider an embedding $g : \mathcal{T} \rightarrow \mathcal{E} \subset \mathbb{R}^m$, where for each token $t \in \mathcal{T}$, we have a unique associated embedding vector $g(t) \in \mathbb{R}^m$, where m is the embedding dimension. Any prompt $p = (t_1, \dots, t_d) \in \mathcal{P}$ comprised of d tokens can be represented as the embedding $g(p) := (g(t_1), \dots, g(t_d)) \in \mathcal{E}^d \subset \mathbb{R}^{m \times d}$. We note that the embedding g does not need to be the same embedding in the generative model, and we may utilize smaller embeddings for ease in optimization.

We could immediately optimize over the embedding space \mathcal{E}^d instead of \mathcal{P} and obtain adversarial examples, however, we assume only access to a black box generative model with text input functionality and thus are only interested in usable text prompts. To ensure we find text prompts, we project the candidate embeddings e to tokens, and evaluate the loss function with the projected tokens. The projection function $\text{proj}_{\mathcal{S}} : \mathcal{E} \rightarrow \mathcal{S}$ chooses the closest token in Euclidean distance in a set \mathcal{S} ,

$$\text{proj}_{\mathcal{S}}(e) := \arg \min_{t \in \mathcal{S}} \|g(t) - e\|_2. \quad (3)$$

Denote by $\mathbf{e} = (e_1, \dots, e_d)$ the embedding matrix corresponding to the d tokens. With slight abuse of notation, let $\text{proj}_{\mathcal{S}}(\mathbf{e}) = (\text{proj}_{\mathcal{S}}(e_1), \dots, \text{proj}_{\mathcal{S}}(e_d))$, mapping embedding matrices to prompts. In our setting, we consider $\text{proj}_{\mathcal{T}_{\text{allow}}}(e_i)$, projecting embeddings to the closest tokens in the restricted subset $\mathcal{T}_{\text{allow}}$ and, using the above abuse of notation, its prompt extension $\text{proj}_{\mathcal{P}_{\text{allow}}}(\mathbf{e})$. For notational convenience, let $q := \text{proj}_{\mathcal{P}_{\text{allow}}}$.

For $R_i^{q(\mathbf{e})} \stackrel{\text{iid}}{\sim} M_{q(\mathbf{e})}$, our final relaxation of our optimization problem is

$$\arg \min_{\mathbf{e} \in \mathcal{E}^d} \frac{1}{n} \sum_{i=1}^n \ell(R_i^{q(\mathbf{e})}). \quad (4)$$

Given the optimum \mathbf{e}^* , the final designed prompt is obtained by $\text{proj}_{\mathcal{P}_{\text{allow}}}(\mathbf{e}^*)$.

2.2 Optimization Approaches

We address the minimization problem in Equation 4 as a black-box optimization, naturally leading to the application of standard techniques to find optimal prompts. We utilize two black-box optimization techniques, the Square Attack algorithm [3], typically employed for generating black-box adversarial examples, and recent advancements in high-dimensional Bayesian optimization [21].

2.2.1 Square Attack

The first black box optimization algorithm we consider is the Square Attack algorithm [3]. Although originally created for attacking image classifiers, its simplicity and competitive query complexity make it a promising candidate for finding adversarial prompts.

The Square Attack algorithm is an iterative algorithm to optimize an arbitrary high dimensional black-box f with solely function evaluations. For simplicity, assume $f : \mathbb{R}^d \rightarrow \mathbb{R}$ and we would like to find an input X that achieves a low value for ℓ .

Let x_t denote the candidate vector at iteration t and x_0 be a chosen initialization. The iteration update for the Square Attack comprises three steps:

1. **Subset Selection:** Select a subset $S \subseteq [d]$ of the indices of x_t to update. Let $x_t^{(S)}$ be the corresponding subsetting vector of x_t .
2. **Sample Values:** Sample $v_1, \dots, v_k \in \mathbb{R}^{|S|}$ and create k new candidate vectors modifying only the entries in S to v_1, \dots, v_k : $x_{t,i}^{(S)} := v_i$ and $x_{t,i}^{(\bar{S})} := x_{t,i}^{(\bar{S})}$ for $i \in [d]$, where \bar{S} is the complement of S .
3. **Update:** Update x_{t+1} with the vector attaining the lowest value among $f(x_{t,1}), \dots, f(x_{t,k})$ and the original $f(x_t)$.

In our implementation, for the Subset Selection step, we choose a random subset S where $|S| \approx d/10$. For the Sample Values step, we choose a random sample using a normal distribution centered around the previous value $x_t^{(S)}$ and identity covariance matrix weighted by a constant term c and the standard deviation of the previous iteration’s function evaluations. The constant term c is a proxy for the ‘step size’, we select c by dividing the average distance between embedding vectors by 10 (we choose $c = 0.1$). We find that normalizing by the standard deviation improves performance, since we are less likely to take large steps close to convergence. More precisely,

$$v_i \stackrel{\text{iid}}{\sim} \frac{c}{\text{stdev}(f(x_{t,1}^{(S)}), \dots, f(x_{t,k}^{(S)}))} \cdot \mathcal{N}(x_t^{(S)}, \mathbf{I}).$$

We provide explicit algorithm details in Algorithm 1. In our setting, we choose ℓ for our optimization function f and X is our flattened word embedding in $\mathbb{R}^{m \cdot d}$.

2.2.2 Bayesian Optimization

Sample efficiency is critical when targeting large foundation models, as evaluating the loss in Equation 4, even for $n = 5$, can take several minutes. Bayesian Optimization (BO) is a general purpose method for solving noisy black-box optimization problems [63, 24], especially when sample efficiency is desired.

In Bayesian optimization, one is given a (possibly initially empty) set of function evaluations $\mathcal{D} : \{(x_1, y_1), \dots, (x_t, y_t)\}$, with y_i a noisy observation of $f(x_i)$. A probabilistic surrogate model—commonly a Gaussian process (GP) [58]—is trained on this dataset to obtain a predictive model of the objective function, $p(y^* | x^*, \mathcal{D})$.

An *acquisition function* is used which leverages this predictive posterior to find the most promising candidates in the search-space to evaluate next, efficiently trading off exploration and exploitation. When new data is acquired, the surrogate model is updated and becomes progressively more accurate. By sequentially selecting candidates to evaluate in this manner, BO can reduce the number of evaluations needed to optimize expensive black-box functions. For a more in-depth introduction to Bayesian optimization, see [26].

Trust Region Bayesian Optimization (TuRBO). Since we may have a token embedding dimension up to $m = 768$ and optimize over $d = 4$ tokens, we may search over a large $768 \times 4 = 3072$ -dimensional space for optimal prompts, a dimensionality well out of reach for traditional Bayesian optimization methods. Eriksson et al. [22] propose Trust Region Bayesian Optimization (TuRBO) which has enabled the use of BO on much higher dimensional functions without making additional assumptions about the search space.

TuRBO mitigates the curse of dimensionality which typically plagues BO algorithms in high-dimensional spaces by dynamically limiting the search space to be within a hyper-rectangular *trust region* so that BO avoids over-exploring the exponentially large search space \mathcal{X} .

The trust region is a hyper-rectangular subset of the input space \mathcal{X} centered at the best point found by the optimizer—the *incumbent*— \mathbf{x}_i^+ and has a side-length $\beta_i \in [\beta_{min}, \beta_{max}]$. If a local optimizer improves upon its own incumbent ρ_{succ} times in a row, β_i is increased to $\min(2\beta_i, \beta_{max})$. Similarly, when a local optimizer fails to make progress ρ_{fail} times in a row, the length β_i is reduced to $\beta/2$. If $\beta_i < \beta_{min}$, that local optimizer is restarted.

3 Losses and Threat Models

The provided formulation for adversarial prompts is flexible and can accommodate a variety of prompt classes and loss functions. The specific choice of prompt class \mathcal{P}_{allow} determines the so-called *threat model* from adversarial examples, or the space of tokens that the adversary is allowed to use to construct a prompt. On the other hand, the loss function ℓ specifies the targeted behavior that the adversary wishes to induce in the generated outputs. Specifying both the allowable prompt class and loss function fully defines both the threat model and the goal of the adversarial prompt.

3.1 Adversarial Targets

In contrast to the standard classification setting for adversarial examples, adversaries could have an arbitrarily complex goal when searching for prompts. To actualize this unrestricted space, we consider the following two candidate loss functions that provide concrete, measurable goals for the adversary.

Classifier Loss. One way to provide a target for adversarial prompts is to use a pretrained classifier. An adversary can then try to find a prompt that pushes the generated text or images toward a particular class. Specifically, the classifier loss is the negative log probability of a class CLS (e.g. cat, dog, car, etc.) from a classifier g :

$$\ell(R) = -\log g(R)_{CLS}. \quad (5)$$

With this loss, the adversary aims to force the model M_p to generate images of the target class CLS. For example, for text-to-image models, an adversary could use an ImageNet classifier to find adversarial prompts that generate images of dogs. For text-to-text models, an adversary could use a sentiment prediction model to prompt the model into generating text with positive or negative sentiment.

Feature Loss. We also consider a second loss function that captures more granular information using features in the generated outputs. For example, when R is generated text, $\ell(R)$ could be the *perplexity* of the text, computed using a separate language model. An adversary would then use this loss to find a prompt that pushes the model into generating text with high perplexity, perhaps nonsense or gibberish text.

3.2 Threat Models

Similar to adversarial examples for discriminative models, an adversarial prompt needs to have a well-defined threat model that cleanly states what an adversary is and is not allowed to do. However, existing threat models in language typically require a clean example that serves as a center of allowable perturbations for the adversary (i.e. by perturbing characters and words or paraphrasing). These threat models do not directly carry over to the prompting setting, as there is no clean example to center around.

Thus in this paper, we consider the following three types of threat models to restrict the adversary.

1. **Unrestricted Prompts:** The adversary is completely unrestricted ($\mathcal{P}_{\text{allow}} = \mathcal{P}$) and can use any sequence of tokens.
2. **Restricted Prompts:** The adversary is restricted to prompts that only use a subset of tokens $\mathcal{T}_{\text{allow}} \subset \mathcal{T}$. This restricted prompt space ($\mathcal{T}_{\text{allow}}^d = \mathcal{P}_{\text{allow}} \subset \mathcal{P}$) eliminates prompts that trivially achieve low loss, such as the prompt containing only the class name. Specifically in the Classifier Loss setting, $\mathcal{T}_{\text{allow}}$ excludes tokens that are similar to the desired class CLS. For example, if we would like to generate images of dogs, then $\mathcal{T}_{\text{allow}}$ would exclude tokens such as puppy, labrador, and emojis of dogs.
3. **Restricted Prepending Prompts:** The adversary is restricted to a subset of tokens and must prepend them to a predefined prompt p' . More specifically, $\mathcal{P}_{\text{pre}} = \{\text{concat}(p_1, p') : p_1 \in \mathcal{P}_{\text{allow}}\}$ where $\mathcal{P}_{\text{allow}}$ is the set of restricted prompts in Task 2. For example, the predefined prompt could be $p' = \text{a picture of a CLS}$. The adversary prepends tokens to cause the model to generate images of a different class CLS_2 (as specified in the loss function), while excluding tokens related to CLS_2 as in the previous task.

These three threat models increase in constraint and difficulty for the adversary. In particular, the third task has opposing text in the prompt by construction that an adversarial prompt must overcome while not using related words.

4 Experiments

Utilizing the Token Space Projection discussed in Section 2.1, we apply the two black-box optimization methods discussed in Section 2.2 to each task introduced in Section 3.2. For the text-to-image tasks, we use embeddings from CLIP ViT-L/14 tokenizer [55] where $m = 768$, and optimize over $d = 4$ tokens. For the text-to-text tasks, we use embeddings from the L-2 H-128 BERT model [68] where $m = 128$, and optimize over $d = 6$ tokens. Code to reproduce all results is available at https://github.com/DebugML/adversarial_prompting. Implementation details and hyperparameters are detailed in Appendix E.

4.1 Image Generation

We consider the optimization Tasks 1, 2, and 3 as described in Section 3 using the Classifier Loss. We use the Stable Diffusion v1.5 model [59] to generate 512x512 images, using a classifier free guidance scale of 7.5 [29] and 25 inference steps with the default scheduler. For our classifier loss function, we use the standard ResNet18 from TorchVision [48] and use the negative log class probability for the desired class.

For Tasks 1 and 2, we allow the black-box optimization methods to query a total of 5,000 prompts before terminating. If the optimizer fails to make progress for 1,000 prompt consecutive queries, we terminate early. Since Task 3 is more difficult, we allow 10,000 prompts to be queried and terminate early after failing to make progress for 3,000 consecutive queries.

4.1.1 Quantitative Evaluation Metrics

In addition to example images demonstrating success, we measure the *success* of optimized prompts quantitatively in two ways.

Most Probable Class Success (MPC Success). The optimizer achieves *MPC success* if it finds a prompt for which the target class is the largest probability class for the majority of 5 images generated according to a pretrained ResNet18 classification model.

Out-Performing Simple Baseline Prompts (OPB Success). Denote by CLS the target ImageNet class. We compute the loss using the following two baseline prompts: CLS and a picture of a CLS. The optimizer achieves *OPB success* if the optimizer finds a prompt that achieves lower loss than both of these baseline prompts.

Table 1: Success rates for TuRBO and Square Attack for all image generation tasks. We report the proportion of these classes for which each method achieved success according to our two measures of success discussed in Section 4.1.1. Both refers to MPC and OPB. Single refers to single/regular image-net classes. Higher refers to image-net classes further up in the image-net class hierarchy.

Task	Metric	TuRBO (Single)	Square (Single)	TuRBO (Higher)	Square (Higher)
Task 1	MPC	18.5%	1.5%	42.0%	12.7%
	OPB	36.8%	13.5%	61.3%	23.3%
	Both	10.5%	0.0%	36.0%	7.3%
Task 2	MPC	58.3%	8.3%	60.0%	0.0%
	OPB	83.3%	58.3%	75.0%	70.0%
	Both	50.0%	8.3%	55.0%	10.0%
Task 3	MPC	25.0%	0.0%	37.3%	0.0%

Goal Class: lizard

louisiana argonhilton deta

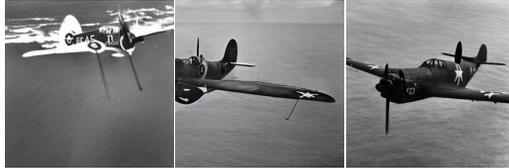


louisiana argon hilton deta



Goal Class: aircraft

pegasus yorkshire wii taken a picture of the ocean



pegasus yorkshirewii taken



(a) Task 2 generated images with the class lizard.

(b) Task 3 generated images with the class aircraft.

Figure 2: Sample images generated using adversarial prompts. In both figures, the first row contains images generated by the optimized prompt and are classified by the ResNet18 classifier as the target class. The second row contains images generated from each individual token in the optimized prompt, none of which are classified as the target class. See Appendix F for more examples.

4.1.2 Unrestricted Prompts

We consider a set of 400 target classes randomly selected from ImageNet classes that have a single token class name (e.g., excluding classes like brown bear). We report success rates for these 400 classes in Table 1. Both methods find prompts that are successful under both definitions from many target classes, with TuRBO consistently achieving a high rate of success.

Higher-Level Classes. In addition to the 400 regular ImageNet target classes, we randomly select a set of 150 higher-level classes from the ImageNet hierarchy. Higher-level ImageNet classes encompass some larger number of standard classes (i.e. wolf is a higher-level class encompassing the base-classes: timber wolf, white wolf, red wolf, coyote). For each higher-level class, the optimizer seeks a prompt that generates images that fall in any of the associated base-classes. The probability that an image falls in a higher-level class is the maximum over the probabilities of all associated base-classes.

Success rates for the 150 higher-level classes are reported in Table 1. Optimizing for higher-level classes is more straightforward since the optimizer does not need to differentiate between highly related classes such as timber wolf and white wolf. We observe that success rates are therefore uniformly higher in this setting.

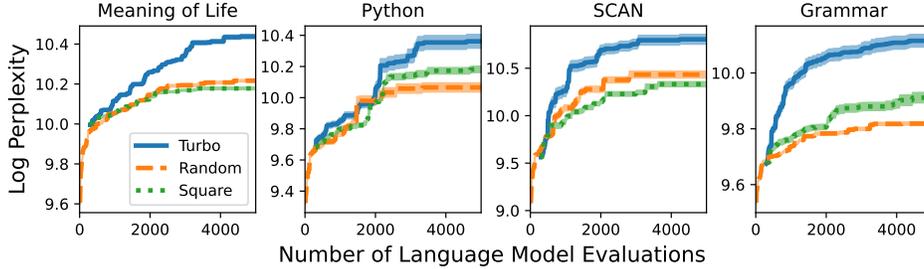


Figure 3: Plots of the maximum log perplexity with respect to the number of Vicuna 13B-v1.1 evaluations. The values are averaged over 20 random seeds, with the one standard error interval shaded. The full seed prompts are in Table 4.

4.1.3 Restricted Prompts

In Section 4.1.2, the optimization methods used an unconstrained vocabulary when designing prompts. Thus, when generating pictures of a wolf, the optimizers were allowed to use the token `wolf`, as well as other highly related words. In this section, we make this task significantly more challenging by forbidding the optimizers from using certain *high-similarity tokens* in the output prompt.

We randomly select 12 target regular ImageNet classes and another 20 target higher-level ImageNet classes. Since the restricted prompt optimization task is more difficult, we select these target classes from the subset of classes for which MPC success was achieved on the unrestricted task. The 12 selected regular ImageNet classes are listed in Table 2. The 20 selected higher-level ImageNet classes can be found in Table 3.

We define a high-similarity token (HST) as a token t where the prompt $p = t$ achieves low loss. We exclude all tokens which achieve a log probability of the target class higher than -3.0 . In order to determine a comprehensive list of HSTs for each target class, we compute the loss for each of the 50k tokens available to the model. Precomputing the loss for these 50k tokens is computationally prohibitive—thus we restrict our experiments to 32 target classes. We report results for optimization without access to HSTs for the 32 selected ImageNet classes in Table 1. Example optimal prompts with generated images can be found in Figure 2a, with additional results in Appendix F.

4.1.4 Restricted Prepending Prompts

We consider the prepending task introduced in Task 3 of Section 3.2 on a subset of the 32 classes used in Section 4.1.3. We prepend prompts onto the following three strings: a picture of a dog, a picture of a mountain, and a picture of the ocean. We provide quantitative results for all three variations of this task in Table 1. Because the prompts `CLS` and a picture of a `CLS` are no longer necessarily strong baselines when prepended to e.g. a picture of a dog, we no longer report OPB success for this task. Example adversarial prompts with generated images can be found in Figure 2b, with additional results in Appendix F.

4.2 Text Generation

With the rise in the use of generative text models, it’s imperative to embed safety measures, ensuring their outputs align with human intentions while avoiding potential risks. We consider the relatively innocuous task of optimizing prompts for large language models to increase the perplexity of the output. Perplexity quantifies the overall level of likelihood of a sequence, measured by the negative log-likelihood as measured by a language model – we use GPT-2 [54]. For a particular prompt, our objective is the average perplexity of the middle three of five generated outputs.

We use the Stable Vicuna 13B-v1.1 model [15]—a state of the art freely available LLM trained with human feedback—for text generation with a maximum generated length of 50 tokens. We optimize perplexity with a maximum budget of 5,000 prompts, terminating early if no progress is made for 2,000 prompts in a row.

We optimize over prompts to prepend to a seed prompt (which we denote as ‘Meaning of Life’, ‘Python’, ‘SCAN’, ‘Grammar’), and replicate these experiments 20 times per seed prompt. To choose

Prompt: *intent oder ranch delegate source every* Explain list comprehension in Python.

Generated Examples:

1. Willkommen auf meinem YouTube-Kanal! Hier ist mein erstes Video, in dem ich dir erkläre, was Intent oder Ranch ist und wie man es in Python verwenden kann.
2. Willkommen auf meinem GPT-3-Konto! Ich bin ein künstlicher Intelligenz-Sprachmodell und kann dir bei deinen Fragen helfen.
3. Willkommen auf meinem ZDF-Zukunftsforum-Kanal! Hier werden wir uns mit den wichtigsten Themen der Zukunft auseinandersetzen und diskutieren, welche Entwicklungen unser Leben

Figure 4: Three examples of text generated from Vicuna 13B-v1.1 using our optimized prompt in italics and the original prompt to which we prepend. The generated text are nonsensical or irrelevant to the prompt.

seed prompts, we randomly selected prompts from the GPT4ALL and PromptSource datasets [2, 4] – see Table 4 in Appendix G for a list of all prompts. We use the Square Attack and TuRBO optimizers from Section 2.2 and also include a random search baseline, where we simply randomly sample tokens. We provide the performance of these approaches on four seed prompts in Figure 3, and more results in Appendix G.

Figure 3 demonstrates our ability to increase the perplexity of the generated text, often by an order of magnitude. We also provide an example of the optimized adversarial prompt for the ‘Python’ seed prompt in Figure 4. This serves as a preliminary evidence that we may successfully optimize adversarial prompts to a desired criterion.

5 Background and Related Work

Adversarial Attacks. The field of adversarial examples has a long history of developing methodologies for finding input perturbations that can change the output of a classifier [7], starting from early text-based attacks to circumvent spam filters [18, 44, 45]. Gradient-ascent based approaches have long been used to generate small perturbations with bounded ℓ_p noise [6, 65, 27, 67, 47], as well as perceptually similar attacks [38]. In computer vision, visually realistic adversarial attacks have been produced on glasses [62], clothing [70], printable patches [23, 13, 43, 66], and more [20, 69]. In language, word embedding [49], word substitution [34], semantically similar [1], and BERT-based [40] attacks were developed to flip the output of a classifier. Most recently, the researchers have begun designing prompt attacks on LLMs such as prompt injection [53] and prompt backdoors [71] that work in specific scenarios. Concurrently, Yang et al. [72] developed a reinforcement learning-based framework for circumventing safety filters on text-to-image generative models.

Of particular relevance are black box adversarial attack frameworks, those which require only query-access to the model being attacked. Early attacks required training a substitute model to generate adversarial examples [52]. Other works have developed zeroth order optimization techniques to avoid training a substitute model [51, 12], with later improvements to reliability and query complexity [8, 30, 14, 3]. Other variants employ tools ranging from combinatorial optimization [50], bandits [31], and confidence scores [28]. Most of these frameworks have been demonstrated primarily in finding ℓ_p adversarial examples in the vision setting.

Prompting. Recent years have seen a prevalence of prompting [42] as a way to interact with models instead of fine-tuning with data. With a large enough of a language model, it became possible to adapt LLMs to new tasks with just a natural language instruction and a few examples as input [9]. Since then, a number of specialized LLMs have been developed for various tasks including code generation [11] and storytelling [74], Prompting has also seen recent use as an interface for guiding zero-shot image generation [57, 60] and image editing [36, 17]. By conditioning image generation on provided text, users can generate high quality images with variable content/style.

6 Conclusion

Broader Impacts. This paper presents a systematic and automated mechanism for manipulating generative models into producing unintended results. We emphasize that any such method, ours or similar, is inherently a double-edged sword. It can be used to bypass alignment safeguards, effectively ‘jailbreaking’ these models. Yet, these methods may facilitate a more systematic alignment process, an alternative to current hand-crafted efforts. Given the relative ease of manually crafting ‘jailbreak’ prompts, we posit that the advantages of an automated approach for investigating unintended behavior in generative models currently outweigh the potential drawbacks.

As closed-source foundation models grow in popularity, it is critical to understand how they react to unexpected prompts. In this paper, we demonstrate that it is possible to efficiently and consistently manipulate the generated output of these foundation models with adversarial prompts. The attack is effective even when using only a small, four token augmentation to otherwise benign prompts. In the future, these techniques may evolve to allow practitioners greater flexibility in prompting via backend prompt engineering. For example, it may be possible to encourage or discourage certain types of outputs by augmenting a user specified prompt with a small handful of tokens found via our adversarial prompting framework.

References

- [1] Moustafa Alzantot, Yash Sharma, Ahmed Elgohary, Bo-Jhang Ho, Mani Srivastava, and Kai-Wei Chang. “Generating natural language adversarial examples”. In: *arXiv preprint arXiv:1804.07998* (2018).
- [2] Yuvanesh Anand, Zach Nussbaum, Brandon Duderstadt, Benjamin Schmidt, and Andriy Mulyar. *GPT4All: Training an Assistant-style Chatbot with Large Scale Data Distillation from GPT-3.5-Turbo*. <https://github.com/nomic-ai/gpt4all>. 2023.
- [3] Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. “Square attack: a query-efficient black-box adversarial attack via random search”. In: *European Conference on Computer Vision*. Springer, 2020, pp. 484–501.
- [4] Stephen H. Bach et al. *PromptSource: An Integrated Development Environment and Repository for Natural Language Prompts*. 2022. arXiv: 2202.01279 [cs.LG].
- [5] Maximilian Balandat, Brian Karrer, Daniel R. Jiang, Samuel Daulton, Benjamin Letham, Andrew Gordon Wilson, and Eytan Bakshy. “BoTorch: A Framework for Efficient Monte-Carlo Bayesian Optimization”. In: *Advances in Neural Information Processing Systems 33*. 2020.
- [6] Battista Biggio, Iginio Corona, Davide Maiorca, Blaine Nelson, Nedim Šrndić, Pavel Laskov, Giorgio Giacinto, and Fabio Roli. “Evasion attacks against machine learning at test time”. In: *Joint European conference on machine learning and knowledge discovery in databases*. Springer, 2013, pp. 387–402.
- [7] Battista Biggio and Fabio Roli. “Wild patterns: Ten years after the rise of adversarial machine learning”. In: *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*. 2018, pp. 2154–2156.
- [8] Wieland Brendel, Jonas Rauber, and Matthias Bethge. “Decision-based adversarial attacks: Reliable attacks against black-box machine learning models”. In: *arXiv preprint arXiv:1712.04248* (2017).
- [9] Tom Brown et al. “Language models are few-shot learners”. In: *Advances in neural information processing systems 33* (2020), pp. 1877–1901.
- [10] Tom B. Brown et al. *Language Models are Few-Shot Learners*. 2020. DOI: 10.48550/ARXIV.2005.14165. URL: <https://arxiv.org/abs/2005.14165>.
- [11] Mark Chen et al. “Evaluating large language models trained on code”. In: *arXiv preprint arXiv:2107.03374* (2021).
- [12] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. “Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models”. In: *Proceedings of the 10th ACM workshop on artificial intelligence and security*. 2017, pp. 15–26.

- [13] Shang-Tse Chen, Cory Cornelius, Jason Martin, and Duen Horng Polo Chau. “Shapeshifter: Robust physical adversarial attack on faster r-cnn object detector”. In: *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer. 2018, pp. 52–68.
- [14] Minhao Cheng, Thong Le, Pin-Yu Chen, Jinfeng Yi, Huan Zhang, and Cho-Jui Hsieh. “Query-efficient hard-label black-box attack: An optimization-based approach”. In: *arXiv preprint arXiv:1807.04457* (2018).
- [15] Wei-Lin Chiang et al. *Vicuna: An Open-Source Chatbot Impressing GPT-4 with 90%* ChatGPT Quality*. Mar. 2023. URL: <https://vicuna.lmsys.org>.
- [16] Aakanksha Chowdhery et al. *PaLM: Scaling Language Modeling with Pathways*. 2022. arXiv: 2204.02311 [cs.CL].
- [17] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. “Diffedit: Diffusion-based semantic image editing with mask guidance”. In: *arXiv preprint arXiv:2210.11427* (2022).
- [18] Nilesh Dalvi, Pedro Domingos, Mausam, Sumit Sanghai, and Deepak Verma. “Adversarial classification”. In: *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*. 2004, pp. 99–108.
- [19] Giannis Daras and Alexandros G Dimakis. “Discovering the Hidden Vocabulary of DALLE-2”. In: *arXiv preprint arXiv:2206.00169* (2022).
- [20] Logan Engstrom, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry. “A rotation and a translation suffice: Fooling cnns with simple transformations”. In: *arXiv preprint arXiv:1712.02779* 1.2 (2017), p. 3.
- [21] David Eriksson and Martin Jankowiak. “High-dimensional Bayesian optimization with sparse axis-aligned subspaces”. In: *Uncertainty in Artificial Intelligence*. PMLR. 2021, pp. 493–503.
- [22] David Eriksson, Michael Pearce, Jacob Gardner, Ryan D Turner, and Matthias Poloczek. “Scalable Global Optimization via Local Bayesian Optimization”. In: *Advances in Neural Information Processing Systems*. Ed. by H. Wallach, H. Larochelle, A. Beygelzimer, F. d Alché-Buc, E. Fox, and R. Garnett. Vol. 32. 2019. URL: <https://proceedings.neurips.cc/paper/2019/file/6c990b7aca7bc7058f5e98ea909e924b-Paper.pdf>.
- [23] Kevin Eykholt, Ivan Evtimov, Earlene Fernandes, Bo Li, Amir Rahmati, Florian Tramer, Atul Prakash, Tadayoshi Kohno, and Dawn Song. “Physical Adversarial Examples for Object Detectors”. In: *CoRR* (2018).
- [24] Peter I Frazier. “A tutorial on Bayesian optimization”. In: *arXiv preprint arXiv:1807.02811* (2018).
- [25] Jacob R Gardner, Geoff Pleiss, David Bindel, Kilian Q Weinberger, and Andrew Gordon Wilson. “Gpytorch: Blackbox matrix-matrix Gaussian process inference with gpu acceleration”. In: *arXiv preprint arXiv:1809.11165* (2018).
- [26] Roman Garnett. *Bayesian Optimization*. to appear. Cambridge University Press, 2023.
- [27] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. “Explaining and harnessing adversarial examples”. In: *arXiv preprint arXiv:1412.6572* (2014).
- [28] Chuan Guo, Jacob Gardner, Yurong You, Andrew Gordon Wilson, and Kilian Weinberger. “Simple black-box adversarial attacks”. In: *International Conference on Machine Learning*. PMLR. 2019, pp. 2484–2493.
- [29] Jonathan Ho and Tim Salimans. “Classifier-Free Diffusion Guidance”. In: *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*. 2021. URL: <https://openreview.net/forum?id=qw8AKxfYbI>.
- [30] Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. “Black-box adversarial attacks with limited queries and information”. In: *International Conference on Machine Learning*. PMLR. 2018, pp. 2137–2146.
- [31] Andrew Ilyas, Logan Engstrom, and Aleksander Madry. “Prior convictions: Black-box adversarial attacks with bandits and priors”. In: *arXiv preprint arXiv:1807.07978* (2018).

- [32] Mohit Iyyer, John Wieting, Kevin Gimpel, and Luke Zettlemoyer. “Adversarial example generation with syntactically controlled paraphrase networks”. In: *arXiv preprint arXiv:1804.06059* (2018).
- [33] Martin Jankowiak, Geoff Pleiss, and Jacob R. Gardner. “Parametric Gaussian Process Regressors”. In: *Proceedings of the 37th International Conference on Machine Learning*. ICML’20. JMLR.org, 2020.
- [34] Robin Jia, Aditi Raghunathan, Kerem Göksel, and Percy Liang. “Certified robustness to adversarial word substitutions”. In: *arXiv preprint arXiv:1909.00986* (2019).
- [35] Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. “How can we know what language models know?” In: *Transactions of the Association for Computational Linguistics* 8 (2020), pp. 423–438.
- [36] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. “Imagic: Text-based real image editing with diffusion models”. In: *arXiv preprint arXiv:2210.09276* (2022).
- [37] Andreas Köpf et al. *OpenAssistant Conversations – Democratizing Large Language Model Alignment*. 2023. arXiv: 2304.07327 [cs.CL].
- [38] Cassidy Laidlaw, Sahil Singla, and Soheil Feizi. “Perceptual adversarial robustness: Defense against unseen threat models”. In: *arXiv preprint arXiv:2006.12655* (2020).
- [39] Brenden M. Lake and Marco Baroni. “Generalization without Systematicity: On the Compositional Skills of Sequence-to-Sequence Recurrent Networks”. In: *ICML*. 2018. URL: <https://arxiv.org/pdf/1711.00350.pdf>.
- [40] Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. “Bert-attack: Adversarial attack against bert using bert”. In: *arXiv preprint arXiv:2004.09984* (2020).
- [41] Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. “What Makes Good In-Context Examples for GPT-3?” In: *arXiv preprint arXiv:2101.06804* (2021).
- [42] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. “Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing”. In: *ACM Computing Surveys* 55.9 (2023), pp. 1–35.
- [43] Xin Liu, Huanrui Yang, Ziwei Liu, Linghao Song, Hai Li, and Yiran Chen. “Dpatch: An adversarial patch attack on object detectors”. In: *arXiv preprint arXiv:1806.02299* (2018).
- [44] Daniel Lowd and Christopher Meek. “Adversarial learning”. In: *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*. 2005, pp. 641–647.
- [45] Daniel Lowd and Christopher Meek. “Good Word Attacks on Statistical Spam Filters.” In: *CEAS*. Vol. 2005. 2005.
- [46] Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. “Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity”. In: *arXiv preprint arXiv:2104.08786* (2021).
- [47] Pratyush Maini, Eric Wong, and Zico Kolter. “Adversarial robustness against the union of multiple perturbation models”. In: *International Conference on Machine Learning*. PMLR. 2020, pp. 6640–6650.
- [48] TorchVision maintainers and contributors. *TorchVision: PyTorch’s Computer Vision library*. <https://github.com/pytorch/vision>. 2016.
- [49] Takeru Miyato, Andrew M Dai, and Ian Goodfellow. “Adversarial training methods for semi-supervised text classification”. In: *arXiv preprint arXiv:1605.07725* (2016).
- [50] Seungyong Moon, Gaon An, and Hyun Oh Song. “Parsimonious black-box adversarial attacks via efficient combinatorial optimization”. In: *International Conference on Machine Learning*. PMLR. 2019, pp. 4636–4645.
- [51] Nina Narodytska and Shiva Prasad Kasiviswanathan. “Simple black-box adversarial perturbations for deep networks”. In: *arXiv preprint arXiv:1612.06299* (2016).

- [52] Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram Swami. “Practical black-box attacks against machine learning”. In: *Proceedings of the 2017 ACM on Asia conference on computer and communications security*. 2017, pp. 506–519.
- [53] Fábio Perez and Ian Ribeiro. “Ignore Previous Prompt: Attack Techniques For Language Models”. In: *arXiv preprint arXiv:2211.09527* (2022).
- [54] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. “Language Models are Unsupervised Multitask Learners”. In: 2019.
- [55] Alec Radford et al. “Learning Transferable Visual Models From Natural Language Supervision”. In: *International Conference on Machine Learning*. 2021.
- [56] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. “Hierarchical Text-Conditional Image Generation with CLIP Latents”. In: (2022). URL: <https://arxiv.org/abs/2204.06125>.
- [57] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. “Zero-shot text-to-image generation”. In: *International Conference on Machine Learning*. PMLR. 2021, pp. 8821–8831.
- [58] Carl Edward Rasmussen. “Gaussian processes in machine learning”. In: *Summer School on Machine Learning*. Springer. 2003, pp. 63–71.
- [59] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. “High-Resolution Image Synthesis With Latent Diffusion Models”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2022, pp. 10684–10695.
- [60] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. “High-resolution image synthesis with latent diffusion models”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 10684–10695.
- [61] Chitwan Saharia et al. “Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding”. In: (2022). URL: <https://arxiv.org/abs/2205.11487>.
- [62] Mahmood Sharif, Sruti Bhagavatula, Lujio Bauer, and Michael K. Reiter. “Accessorize to a Crime: Real and Stealthy Attacks on State-of-the-Art Face Recognition”. In: *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, Vienna, Austria, October 24-28, 2016*. 2016, pp. 1528–1540.
- [63] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. “Practical Bayesian optimization of machine learning algorithms”. In: *Advances in neural information processing systems 25* (2012).
- [64] Lichao Sun, Kazuma Hashimoto, Wenpeng Yin, Akari Asai, Jia Li, Philip Yu, and Caiming Xiong. “Adv-bert: Bert is not robust on misspellings! generating nature adversarial samples on bert”. In: *arXiv preprint arXiv:2003.04985* (2020).
- [65] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. “Intriguing properties of neural networks”. In: *arXiv preprint arXiv:1312.6199* (2013).
- [66] Simen Thys, Wiebe Van Ranst, and Toon Goedemé. “Fooling automated surveillance cameras: adversarial patches to attack person detection”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. 2019.
- [67] Florian Tramèr and Dan Boneh. “Adversarial training and robustness for multiple perturbations”. In: *Advances in Neural Information Processing Systems*. 2019, pp. 5858–5868.
- [68] Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. “Well-Read Students Learn Better: On the Importance of Pre-training Compact Models”. In: *arXiv: Computation and Language* (2019).
- [69] Eric Wong, Frank R Schmidt, and J Zico Kolter. “Wasserstein adversarial examples via projected sinkhorn iterations”. In: *arXiv preprint arXiv:1902.07906* (2019).
- [70] Zuxuan Wu, Ser-Nam Lim, Larry S Davis, and Tom Goldstein. “Making an invisibility cloak: Real world adversarial attacks on object detectors”. In: *European Conference on Computer Vision*. Springer. 2020, pp. 1–17.

- [71] Lei Xu, Yangyi Chen, Ganqu Cui, Hongcheng Gao, and Zhiyuan Liu. “Exploring the Universal Vulnerability of Prompt-based Learning Paradigm”. In: *arXiv preprint arXiv:2204.05239* (2022).
- [72] Yuchen Yang, Bo Hui, Haolin Yuan, Neil Gong, and Yinzhi Cao. “SneakyPrompt: Evaluating Robustness of Text-to-image Generative Models’ Safety Filters”. In: *arXiv preprint arXiv:2305.12082* (2023).
- [73] Ryoma Yoshimura, Masahiro Kaneko, Tomoyuki Kajiwara, and Mamoru Komachi. “SOME: Reference-less Sub-Metrics Optimized for Manual Evaluations of Grammatical Error Correction”. In: *Proceedings of the 28th International Conference on Computational Linguistics*. Barcelona, Spain (Online): International Committee on Computational Linguistics, Dec. 2020, pp. 6516–6522. DOI: 10.18653/v1/2020.coling-main.573. URL: <https://aclanthology.org/2020.coling-main.573>.
- [74] Ann Yuan, Andy Coenen, Emily Reif, and Daphne Ippolito. “Wordcraft: story writing with large language models”. In: *27th International Conference on Intelligent User Interfaces*. 2022, pp. 841–852.
- [75] Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. “SWAG: A Large-Scale Adversarial Dataset for Grounded Commonsense Inference”. In: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2018.

Appendix

A Limitations

The main limitations on our work are the evaluation of the open source models and necessity of real valued loss functions. Our approaches apply to arbitrary vision and language models, and we formulate our attacks as black box optimizations. However, due to the cost of systematic evaluation on black-box models, we perform our experiments on open source models, as systematically evaluating our approach on DALLE-2 or ChatGPT may become prohibitively expensive. Furthermore, our optimization requires some form of real-valued signal to ensure improvement. For possible tasks of interest, there may be a binary loss, which would possibly provide too sparse of a signal to optimize.

B Compute

These experiments were evaluated on 4 Nvidia Tesla A100s. Each individual run takes about 30-60 minutes on a single GPU, for a total of about 1000 GPU hours for all experiments.

C Adversarial Prompt Optimization Pipeline

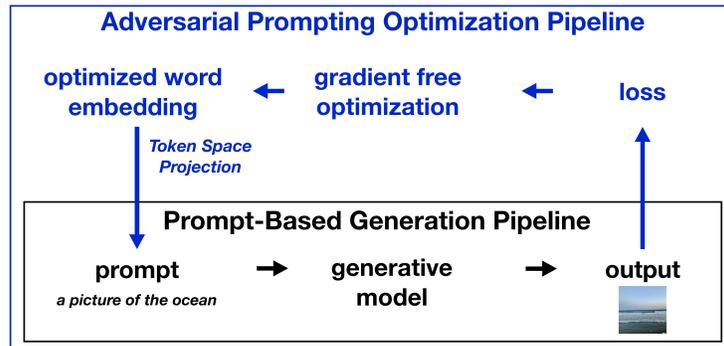


Figure 5: Adversarial prompting pipeline for a prompt-based generative model.

D Square Attack Algorithm

We include an algorithmic description of the Square Attack described in Section 2.

D.1 Selected ImageNet Classes for Tasks 2 and 3

Table 2 lists the 12 ImageNet selected for optimizing without access to HSTs and for the prepending task. Table 3 lists the 20 classes selected that were higher level in the ImageNet hierarchy.

Algorithm 1 Square Attack Algorithm

Input: loss function $\ell : \mathbb{R}^d \rightarrow \mathbb{R}$

```
1: Initialize  $x_0$ 
2:  $\sigma \leftarrow 1$ 
3: for  $t = 0, \dots, T - 1$  do
4:   Select random subset  $S \subseteq [d]$  with  $|S| = d/10$ 
5:    $x_{t,1}, \dots, x_{t,k} \leftarrow x_t$ 
6:   Sample  $v_1, \dots, v_k \stackrel{\text{iid}}{\sim} \frac{1}{10\sigma} \mathcal{N}(x_0^{(S)}, I)$ 
7:    $x_{t,i}^{(S)} \leftarrow v_i$ 
8:   Compute and store  $q_t \leftarrow (f(x_{t,1}), \dots, f(x_{t,k}))$ 
9:    $\sigma \leftarrow \text{stdev}(q_t)$ 
10:  if  $f(x_t) \geq \min(q_t)$  then
11:     $x_{t+1} \leftarrow \arg \min q_t$ 
12:  else
13:     $x_{t+1} \leftarrow x_t$ 
14:  end if
15: end for
16: Return  $x_T$ 
```

Table 2: The 12 regular ImageNet classes randomly selected for more difficult optimization tasks (optimizing without access to high-similarity tokens and pre-pending tasks).

CLASSES

TABBY, SPORTSCAR, VIOLIN, BALLPLAYER,
LIBRARY, LIPSTICK, PINWHEEL, MASK,
CHURCH, GOBLET, GOLDFINCH, PICKUP

Table 3: The 20 higher-level ImageNet classes selected for more difficult optimization tasks (optimizing without access to high-similarity tokens and prepending tasks).

CLASSES

AIRCRAFT, BIG-CAT, BIRD, DOG, PERSON,
TRUCK, STORE, LIZARD, HEADDRESS, TOI-
LETRY, SEAT, MUSICAL-INSTRUMENT, PIANO,
GYMNASTIC-APPARATUS, BALL, EQUINE,
SWIMSUIT, FRUIT, DOMESTIC-CAT, BUS

E Implementation Details and Hyperparameters

We use the Hugging Face API to load in and apply our generative models, the details of which we describe in the relevant sections below. We use the open source BoTorch [5, 25] implementation of TuRBO. Code to reproduce all results is available at https://github.com/DebugML/adversarial_prompting. For the TuRBO optimization method, all trust region hyper-parameters are set to the TuRBO defaults as used in [22].

Since we consider large numbers of function evaluations for several tasks, we use an approximate Gaussian process (GP) surrogate model. In particular, we use a Parametric Gaussian Process Regressor (PPGPR) [33] with a deep kernel (a GP with several linear layers between the search space and the GP kernel). We use a deep kernel with 3 linear layers with 256, 128, and 64 nodes. We use 100 random points to initialize optimization across all methods. To account for the variability in the generative model outputs, we generate 5 unique outputs for each prompt queried and compute the average loss.

F Additional Examples

We provide additional examples of optimal prompts found by our methods for the image generation tasks and the resultant images generated by those prompts. See additional examples for Task 2 in Figures 6 to 9. See additional examples for Task 3 in Figures 10 to 13.

Goal Class: lipstick

eyebrows octavia assumptions 🙌



eyebrows

octavia

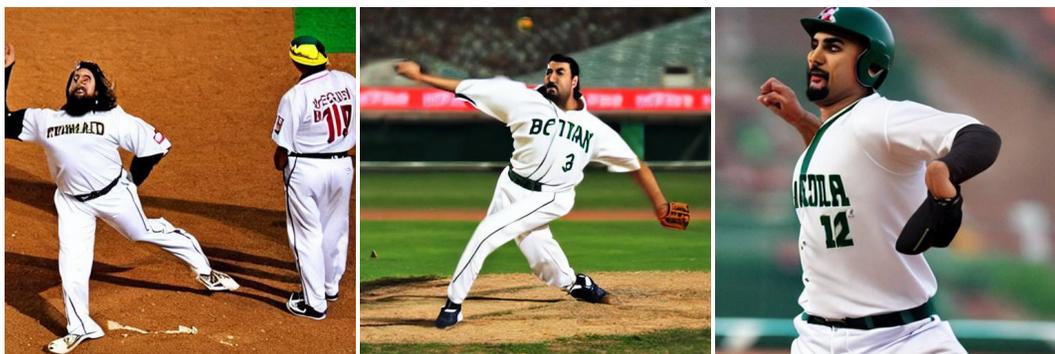
assumptions 🙌



Figure 6: Visualizing images generated by an optimal prompt found by TuRBO for the target higher-level ImageNet class: lipstick. Top row shows three images generated by the optimal prompt. All three images are classified by the ResNet as lipstick, thus indicating this was a successful prompt. The bottom row shows images generated from each individual token in the optimal prompt. None of the individual tokens are able to generate the target class.

Goal Class: ballplayer

mohammed 🍌 sal threw



mohammed



sal

threw



Figure 7: Visualizing images generated by an optimal prompt found by TuRBO for the target regular ImageNet class ballplayer. Top row shows three images generated by the optimal prompt. All three images are classified by the ResNet as ballplayers, thus indicating this was a successful prompt. The bottom row shows images generated from each individual token in the optimal prompt. None of the individual tokens are able to generate the target class.

Goal Class: piano

wider organist cart



wider

organi

st

cart



Figure 8: Visualizing images generated by an optimal prompt found by TuRBO for the target higher-level ImageNet class: piano. Top row shows three images generated by the optimal prompt. All three images are classified by the ResNet as pianos, thus indicating this was a successful prompt. The bottom row shows images generated from each individual token in the optimal prompt. None of the individual tokens are able to generate the target class.

Goal Class: toiletry

artist led whitening oil



artist

led

whitening

oil



Figure 9: Visualizing images generated by an optimal prompt found by TuRBO for the target higher-level ImageNet class: toiletry. Top row shows three images generated by the optimal prompt. All three images are classified by the ResNet as toiletries, thus indicating this was a successful prompt. The bottom row shows images generated from each individual token in the optimal prompt. None of the individual tokens are able to generate the target class.

Goal Class: dog

turbo lhaff✓a picture of a mountain

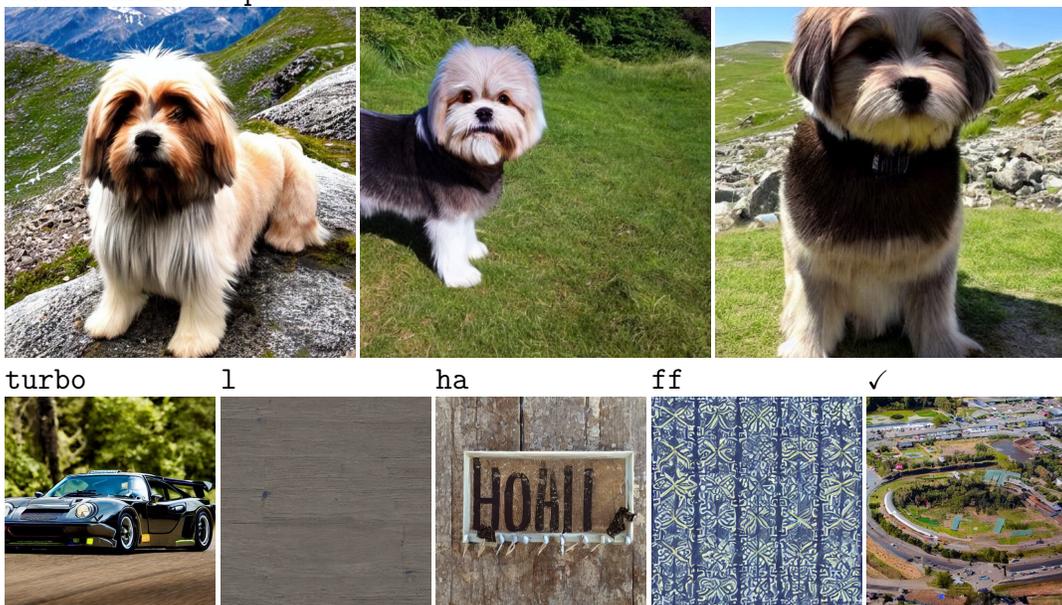


Figure 10: Visualizing images generated by an optimal prompt found by our method for the higher-level ImageNet class: dog. Top row shows three images generated by the optimal prompt. All three images are classified by the Resnet18 classifier as dogs (rather than mountains). The bottom row shows images generated from each individual token in the optimal prompt. None of the individual tokens are able to generate the target class.

Goal Class: person

cllr mods a picture of a dog



cllr

mods



Figure 11: Visualizing images generated by an optimal prompt found by our method for the higher-level ImageNet class: person. Top row shows three images generated by the optimal prompt. All three images are classified by the ResNet18 classifier as people (rather than dogs). The bottom row shows images generated from each individual token in the optimal prompt. None of the individual tokens are able to generate the target class.

Goal Class: sportscar

jaguar fp euphoria idan a picture of the ocean

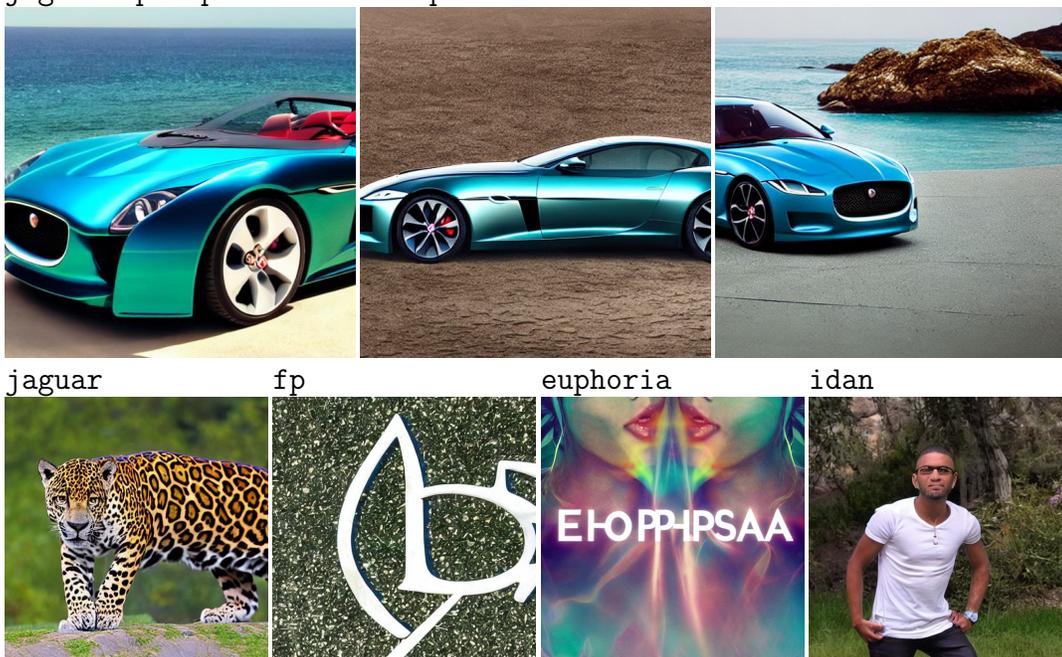


Figure 12: Visualizing images generated by an optimal prompt found by our method for the single ImageNet class: sportscar. Top row shows three images generated by the optimal prompt. All three images are classified by the ResNet18 classifier as sportscars (rather than oceans). The bottom row shows images generated from each individual token in the optimal prompt. None of the individual tokens are able to generate the target class.

Goal Class: ballplayer

fiji players hormoncaine a picture of a mountain



fiji



players



hormon



caine



Figure 13: Visualizing images generated by an optimal prompt found by TuRBO for the target higher-level ImageNet class: ballplayer. Top row shows three images generated by the optimal prompt. All three images are classified by the ResNet as ballplayers, thus indicating this was a successful prompt. The bottom row shows images generated from each individual token in the optimal prompt. None of the individual tokens are able to generate the target class.

G Text Generation

G.1 Seed Prompt Details

Origin	Reference Name	Full Prompt
SCAN, template_jump_around_right, test, affirmative_bottom, 566 [39]	SCAN	‘run opposite left after jump around right thrice\n\n Given the commands above, produce the corresponding correct sequence of actions. The actions should be comma-separated.’
2023 AMC8 Problem 21	Math	‘Bob writes the numbers 1, 2, ..., 9 on separate pieces of paper, one number per paper. He wishes to divide the papers into 3 piles of three papers so that the sum of the numbers in each pile will be the same. In how many ways can this be done?’
tmu_gfm_dataset, train, correct-sentence, 329 [73]	Grammar	‘Grammatically improve the below text. Note that the original meaning has to be preserved and also it should sound natural.\n\n Text: People needs a safe environment to live in, and also needs a private environment to stay independent.’
GPT4All [2]	Python	‘Explain list comprehension in Python.’
SWAG, regular, test, how_ends, 8686 [75]		‘From behind, someone swings the poker at his head. He...\n How does the description likely end?\n (a): falls face down on the carpet.\n (b): glances at all of his heads.\n (c): looks grandma in the eye.\n (d): stares at the phone.\n’
GPT4All [2]	Poem	‘Compose a 10 line rhyming poem about cats.’
OASST1 Validation [37]	Meaning of Life	‘What is the meaning of life, written in the style of Dr. Seuss poem?’
GPT4ALL [2]	Slogan	‘Design a slogan for a restaurant that serves traditional Taiwanese dishes.’

Table 4: Table of the seed prompts considered, where the prompt was obtained from, and our reference name. For prompts from an explicit dataset, we provide specific identifying details for the sample, such as the exact sample number. For some of the prompts, we modify the wording slightly to avoid using prompts the model has been trained on.

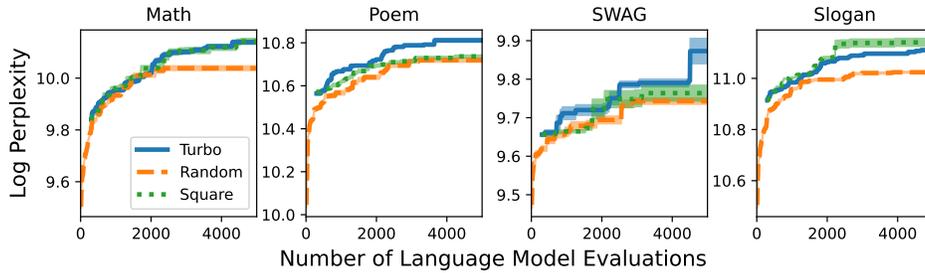


Figure 14: Plots of the maximum log perplexity with respect to the number of language model evaluations. The values are averaged over 20 random seeds, with the standard error shaded. The full seed prompts are in Table 4.

G.2 More Seed Prompt Generations

We provide more examples of generated text outputs in Figures 15 to 18 and the performance on other seed prompts in Figure 14. We see that we are consistently able to optimize desired loss functions and obtain nonsensical outputs.

Prompt: *oral geschichte ' efficient ding once* Grammatically improve the below text. Note that the original meaning has to be preserved and also it should sound natural.
 \n\nText: People needs a safe environment to live in, and also needs a private environment to stay independent.

Generated Examples:

1. Menschen brauchen eine sichere Umgebung, in der sie leben können, und auch eine private Umgebung, in der sie unabhängig bleiben können.
2. Menschen brauchen eine sichere Umgebung, um in ihr zu leben, und auch eine eigene Umgebung, um unabhängig zu bleiben.
3. Menschen brauchen eine sichere Umgebung, in der sie leben können, und auch eine private Umgebung, um unabhängig zu bleiben.

Figure 15: Grammar Prompt: Three examples of text generated from Vicuna 13B-v1.1 using our optimized prompt in italics and the original prompt to which we prepend. The generated text are nonsensical or irrelevant to the prompt. Note that while 'geschichte' is 'story' in German, without the apostrophe, the model will not consistently generate German text.

Prompt: *cod irs theorem fit lay pg* Bob writes the numbers 1, 2, ..., 9 on separate pieces of paper, one number per paper. He wishes to divide the papers into 3 piles of three papers so that the sum of the numbers in each pile will be the same. In how many ways can this be done?

Generated Examples:

1. This problem is known as the “Coding Theorem” or “Codirras-Wagner Theorem”. It can be solved using combinatorics, specifically permutations and combinations. Let’s denote the three piles of papers as
2. This problem is an example of the “Coffee and Cake” problem, also known as the “Irregular Rearrangement” problem. The solution to this problem is to use the “Codirisu Theorem”, which states
3. This problem is a variation of the so-called “Coffee Shop” problem, which is a well-known problem in combinatorics. The solution to this problem is to use the “Coffee Shop” solution, which

Figure 16: Math Prompt: Three examples of text generated from Vicuna 13B-v1.1 using our optimized prompt in italics and the original prompt to which we prepend. The problems and theorems, e.g. “Coffee and Cake problem” and “Codirras-Wagner Theorem”, are not real problems or theorems.

Prompt: *ben post ole fin quick marine* run opposite left after jump around right thrice\n\n Given the commands above, produce the corresponding correct sequence of actions. The actions should be comma-separated.

Generated Examples:

1. Human, bench press 50kg, 3 sets of 8 reps, bodyweight squats 20 reps, rowing machine 500m, push-ups 10 reps, 1-
2. Human, bench press 100 pounds for 3 sets of 8 reps, bicep curls 20 pounds for 3 sets of 10 reps, tricep dips 40 pounds for
3. Human, bench press 50 kg for 3 sets of 10 reps, bicep curl 20 kg for 3 sets of 10 reps, tricep dips for 3 sets of

Figure 17: SCAN Prompt: Three examples of text generated from Vicuna 13B-v1.1 using our optimized prompt in italics and the original prompt to which we prepend. The generated text comprises exercise routines rather than relevant actions from the prompt.

Prompt: *default since levant rico ny levant* run opposite left after jump around right thrice\n\n Given the commands above, produce the corresponding correct sequence of actions. The actions should be comma-separated.

Generated Examples:

1. Human, default, default,
2. Human, default, default,
3. Human, default, default,

Figure 18: SCAN Prompt: Three examples of text generated from Vicuna 13B-v1.1 using our optimized prompt in italics and the original prompt to which we prepend. The generated text is essentially only the word ‘default’.