
Prompt Recovery for Image Generation Models: A Comparative Study of Discrete Optimizers

Anonymous authors

Paper under double-blind review

ABSTRACT

Recovering natural language prompts for image generation models, solely based on the generated images is a difficult discrete optimization problem. In this work, we present the first head-to-head comparison of recent discrete optimization techniques for the problem of *prompt inversion*. Following prior work on prompt inversion, we use CLIP’s Radford et al. (2021) text-image alignment as an inexpensive proxy for the distribution of prompt-image pairs, and compare several discrete optimizers against BLIP2’s image captioner (Li et al., 2023) and PRISM (He et al., 2024) in order to evaluate the quality of discretely optimized prompts across various metrics related to the quality of inverted prompts and the images that they generate. We find that while the discrete optimizers effectively minimize their objectives, CLIP similarity between the inverted prompts and the ground truth image acts as a poor proxy for the distribution of prompt-image pairs – responses from well-trained captioners often lead to generated images that more closely resemble those produced by the original prompts. This finding highlights the need for further investigation into inexpensive methods of modeling the relationship between the prompts for generative models and their output space.

1 INTRODUCTION

Images generated by AI models are flooding the internet and the models and prompts used for generation often feel like components of alchemy. Naturally, the academic community aims to better understand the mechanisms at play in these systems, with one crucial focus on inverting the generative process by recovering prompts from images (Fan et al., 2024).

As the name suggests, given only the output image from a generative model, prompt inversion methods seek to find the prompt that generated the image. While the goal is clear, this problem leads to a wide variety of approaches. For example, given an initial prompt estimate, by strategically refining the text input Sohn et al. (2023) have found that the image generation process can be effectively controlled. Yet, this refining process can be difficult, leading to Wen et al. (2023) leveraging CLIP’s (Radford et al., 2021) embedding space to directly optimize natural language inputs to be close to target images. Similarly, Mahajan et al. (2024) have proposed an inversion technique that backpropagates through intermediate steps of the diffusion process instead of relying on CLIP embeddings. In contrast, it has been found that training a captioner on a dataset of prompt-image pairs effectively learns the prompt distribution well enough to act as an inverter.

While prompt inversion is an interesting task in its own right, there are two practical motivations to develop strong methods. First, those interested in better controlling the output of image generation models may want to find a prompt from an image as a starting point for their own prompt engineering. Second, by extracting prompts from images, one might better understand the various associations these text-to-image models have and debug them to avoid explicit content. Yet, to date, there is no standardized comparison of these methods for prompt inversion from image generation models.

054 In this work, we benchmark several approaches to prompt inversion and attempt to shed light
055 on three primary questions: 1) How closely do the distribution of images generated by each
056 approach align with the distribution generated by the original prompt? 2) Do text-image
057 alignment models like CLIP successfully act as a proxy for the joint distribution of prompt-
058 image pairs? 3) Does the ability for discrete optimization methods to search the input space
059 allow them to outperform the ingrained knowledge of learned models? We investigate these
060 questions by comparing four generic methods of discrete optimization within the prompt
061 space (Zou et al., 2023; Wen et al., 2023; Zhu et al., 2023; Andriushchenko, 2023) against two
062 methods that directly sample from the space (Li et al., 2023; He et al., 2024) and evaluate
063 each across several metrics. We follow our analysis with a discussion on the efficacy of
064 different approaches to approximating and searching within prompt spaces.

066 2 RELATED WORK

069 To best situate this paper among prior work, we discuss several motivations for executing
070 prompt inversion, other domains where discrete optimization is relevant, and the overall
071 goal of our work in contrast to related papers.

072 Here, we focus on discrete optimization methods for recovering image prompts, but discrete
073 optimization over natural language has several other applications including jailbreaking
074 LLMs (Andriushchenko, 2023; Zou et al., 2023; Zhu et al., 2023) and measuring memoriza-
075 tion (Schwarzschild et al., 2024; Kassem et al., 2024). In particular, Zou et al. (2023) propose
076 a method to find adversarial prompts for LLMs that break their safety alignment. We ex-
077 periment with this optimizer for image generation inversion as Zou et al. (2023) compare
078 their optimizer to PEZ in their work, but only with the goal of jailbreaking LLMs.

079 Whereas prompt optimization strategies in the text generation space have specific goals,
080 such as generating targeted strings, the image generation space has struggled with tractable
081 options for aligning prompts and generated images. While we follow prior work in using
082 CLIP as a proxy model, this choice is primarily driven by the practical challenges of directly
083 optimizing prompts through backpropagation in the diffusion process. The computational
084 requirements for coarse-grained exploration and fine-grained search that discrete optimiza-
085 tion often calls for (Parker and Rardin, 2014) would entail generating multiple full images
086 for every candidate prompt at every step. Following the search parameters recommended
087 by (Zou et al., 2023), optimizing a single prompt would require generating a minimum of
088 512 images per step, which quickly becomes prohibitive over numerous iterations.

089 Mahajan et al. (2024) have attempted to address this burden by focusing on the similarity
090 between predicted noise residuals at specific diffusion timesteps, rather than generating full
091 images. However, in alignment with prior work on noise inversion (Song et al., 2020; Mokady
092 et al., 2023) the authors find that prompts only have strong influence on the generated
093 image during a narrow range of timesteps. At early timesteps, the image becomes largely
094 “locked in,” so even substantial changes to the prompt have little effect. In contrast, at
095 later timesteps, the stochasticity of the diffusion process leads to large variations in the final
096 image, even when the correct prompt is used. This unpredictability makes it difficult to
097 rely on noise residual comparisons for consistent prompt inversion. Thus, prompt inversion
098 methods (Wen et al., 2023; Williams and Kolter, 2024) often rely on deterministic proxy
099 models like CLIP, which offer a more stable and efficient alternative. Through CLIP’s text-
100 image alignment, we can more reliably approximate the prompt-image distribution without
having to address the risk of stochasticity producing divergent results.

101 Importantly, direct discrete optimization is not the only method for finding viable prompts.
102 Several approaches focus on using black box models to sample prompts. Both Zhang et al.
103 (2024) and He et al. (2024) use pretrained language models to extract prompts for given
104 a output across text generation and image generation tasks respectively. Moreover, as we
105 show in this work, even a simple captioner that has not been finetuned for prompt generation
106 often outperforms discrete optimization methods. In fact, Reade et al. (2023) have found
107 that a captioner fine-tuned on pairs of prompts and the images that they generate can
effectively sample prompts that are exceptionally similar to the ground truth.

108 Despite this performance, we focus on the discrete case as direct discrete optimization can be
 109 beneficial for better understanding the behavior of the image generation models. Similarly to
 110 prior work on counterfactual explanations (Verma et al., 2020), directly optimizing inputs for
 111 a desired outputs helps to better understand the decision boundaries of classifiers (Ribeiro
 112 et al., 2016). While generative models are not classifiers with explicit decisions and accuracy
 113 metrics, they are constantly making decisions on their representations based on the prompts.
 114 From background color to subject ethnicity, discrete optimization methods may provide a
 115 useful understanding of the relationship between prompts and images (Williams et al., 2024).

116 We emphasize solidifying ways of comparing discrete optimizers for image generation tasks.
 117 Even with the rise of novel discrete optimization methods, standard prompt recovery com-
 118 parisons over images are missing. We focus on a holistic benchmark on not only the similarity
 119 between prompt and image, but also the similarity among images generated by the inverted
 120 prompts which to the best of our knowledge has not been standardized in this setting.

122 3 SELECTED ALGORITHMS

124 To best introduce the optimizers we study, it is critical to pose the prompt inversion problem
 125 formally. Consider a tokenizer that maps from natural language to sequences of integer
 126 valued tokens corresponding to a list of indices in a vocabulary of tokens \mathbb{T} . Let $x \in \mathbb{T}^s$ be a
 127 length s sequence of tokens. Next, let $\mathbf{E} \in \mathbb{R}^{|\mathbb{T}| \times d}$ be a matrix whose rows are d -dimensional
 128 embedding vectors, one for each token in the vocabulary. To embed a sequence x , we can
 129 define $\mathbf{X} \in \{0, 1\}^{s \times |\mathbb{T}|}$ s.t. $\sum_{i=1}^T \mathbf{X}_{j,i} = 1 \forall j \in \{1, \dots, M\}$ to be a matrix of one-hot encoded
 130 rows for the integers in the sequence x . The product $\mathbf{X}\mathbf{E}$ defines an $s \times d$ embedding of x .

131 Prompt inversion techniques seek to find the sequence of tokens x , or equivalently their
 132 corresponding one-hot encodings \mathbf{X} , that solve $\mathcal{M}^{-1}(Y)$, where \mathcal{M} is a stochastic generative
 133 model that maps a sequence of tokens x to an image Y . Typically we express the solution as
 134 the minimizer of some loss function \mathcal{L} , or the solution to the following optimization problem.

$$135 \quad \underset{\mathbf{X} \in \{0,1\}^{s \times |\mathbb{T}|}}{\operatorname{argmin}} \quad \mathcal{L}(\mathcal{M}(\mathbf{X}\mathbf{E}), Y) \quad \text{s.t.} \quad \sum_{i=1}^{|\mathbb{T}|} X_{j,i} = \mathbf{1} \quad \forall j \in \{1, \dots, s\} \quad (1)$$

138 Khashabi et al. (2021) show that embeddings in \mathbb{R}^d outside of the discrete set of the rows of
 139 \mathbf{E} have little meaning to the generative model \mathcal{M} . As a consequence, most prompt inversion
 140 methods focus on strategies for discrete optimization within the embedding table \mathbf{E} ; we call
 141 this ‘hard prompting’ in a discrete space rather than ‘soft prompting’ in a continuous space.
 142 While the gradient exists with respect to the entries of the input $\mathbf{X}\mathbf{E} \in \mathbb{R}^{s \times d}$, continuous
 143 descent-based methods risk finding minima outside of \mathbb{T}^s , leaving us without hard tokens.

144 Moreover, computing the gradient through the full generation model \mathcal{M} may be too expen-
 145 sive (for example when \mathcal{M} is a diffusion model forward passes may take multiple seconds),
 146 but as emphasized above, prior work often uses CLIP (Radford et al., 2021) to encode im-
 147 ages and text in a shared latent space. Some of the methods we examine operate wholly
 148 within CLIP’s latent space to compute the loss between the prompt and the target image.
 149 These methods approximate Equation (1) by solving the following problem where $\mathcal{L}_{\text{CLIP}}$ is
 150 a similarity loss defined over CLIP embeddings.

$$151 \quad \underset{\mathbf{X} \in \{0,1\}^{s \times |\mathbb{T}|}}{\operatorname{argmin}} \quad \mathcal{L}_{\text{CLIP}}(\mathbf{X}\mathbf{E}, Y) \quad \text{s.t.} \quad \sum_{i=1}^{|\mathbb{T}|} X_{j,i} = \mathbf{1} \quad \forall j \in \{1, \dots, s\} \quad (2)$$

155 3.1 PEZ

156 The first approach we consider is PEZ (Wen et al., 2023), a version of projected gradient
 157 descent where descent steps are made in the continuous embedding space. The gradients of
 158 the objective in Equation (2) are evaluated at points in embedding space corresponding to
 159 real tokens, but the trajectory of the iterates may deviate from the discrete token set.

160 More formally, let $\text{Proj}_{\mathbf{E}}(\cdot)$ be an operator that projects vectors (or matrices row-wise) from
 161 \mathbb{R}^d to their nearest row-vector of \mathbf{E} , and let $\xi_i \in \mathbb{R}^{s \times d}$ be a soft prompt. As an iterative

162 gradient-based optimizer, PEZ produces a sequence of iterates $[\xi_0, \xi_1, \dots, \xi_n]$ as it solves
 163 the minimization problem in Equation (2). To update from ξ_i to ξ_{i+1} , PEZ computes the
 164 gradient of the loss at the hard prompt $\text{Proj}_{\mathbf{E}}(\xi_i)$ and takes a step in the direction of this
 165 gradient from the soft prompt ξ_i and then calls $\text{Proj}_{\mathbf{E}}(\xi_i)$ to project back to the space of
 166 hard prompts. Thus, PEZ gives a fast, lightweight method of discrete optimization while
 167 still using gradient-based descent to approximately solve the problem in Equation (1). For
 168 more information, see Algorithm 1 as described by Wen et al. (2023). For a single image,
 169 we run PEZ over the CLIP loss for 3000 steps and return the prompt that maximizes the
 170 CLIP similarity between the image embedding and the text embedding of the prompt.

171 3.2 GREEDY COORDINATE GRADIENTS

172 Greedy Coordinate Gradients (GCG) (Zou et al., 2023) is an alternative method for opti-
 173 mizing over the discrete vocabulary using the gradients of the objective with respect to the
 174 matrix \mathbf{X} in Equation (2). In particular, we compute the gradient of the loss with respect
 175 to \mathbf{X} , which is a matrix of the same shape that approximately ranks token swaps. As each
 176 entry in a given row of \mathbf{X} corresponds to a token in the vocabulary, each row i in its gradi-
 177 ent relays to us how influential changing the token x_i to each other token in the vocabulary
 178 might be in lowering the loss. More formally, we compute $\nabla_{\mathbf{X}} \mathcal{L}_{\text{CLIP}}(\mathcal{M}(\mathbf{X}\mathbf{E}), Y)$, then, just
 179 as gradient descent methods takes steps in the opposite direction of the gradient, we select
 180 a random batch of candidate swaps from the top k largest entries of the *negative* gradient.
 181 A given swap corresponds to a single token change in x and we directly compute the loss
 182 for each of these candidates and greedily accept the best one as our new iterate. As done
 183 with PEZ, we run GCG over the CLIP loss for 3000 steps, returning the best prompt as
 184 determined by CLIP similarity between the image embedding and the prompt’s embedding.

185 3.3 AUTODAN

186 AutoDAN (Zhu et al., 2023) was proposed as a method of finding human-readable adversarial
 187 attacks on aligned language models. The optimizer solves Eq. (1) by iteratively optimizing
 188 a single token appended to the current prompt. Given an initial prefix, e.g., “Image of a”,
 189 the algorithm searches for the token that follows ‘a’ that minimizes the objective function.
 190 The optimizer incorporates a ‘readability’ objective based on the log probability of the next
 191 token given an underlying language model. Similarly to GCG, AutoDAN employs a coarse-
 192 to-fine search strategy by appending an initial token, \hat{x} to the current iterate x , and scores
 193 each token in the vocabulary according to the following scoring function:

$$194 \text{score}(x_i) = -(\nabla_{\hat{x}} \mathcal{L}([x, \hat{x}]E)) + \log(p(x_i|x)) \quad (3)$$

195 The algorithm selects the top k scoring tokens and performs a fine-grained search by com-
 196 puting the exact loss over each, taking the token that minimizes the loss, \mathcal{L} . This minimizing
 197 token is then appended to x , giving $x_{t+1} = [x_t \ x_i^*]$.

198 AutoDAN was originally designed for text-to-text language models, where the log proba-
 199 bility, $\log(p(x_i|x))$ was directly available. However, in this review, we use CLIP to deter-
 200 mine the quality of the prompt, which does not inherently compute the log probability.
 201 We thus use FUSE (Williams and Kolter, 2024), a recently proposed approach for solving
 202 multi-objective problems across models and tokenizers. FUSE approximates the jacobian
 203 of a mapping between the two models and uses the embeddings of a text-to-text language
 204 model, such as GPT2 to compute both the log probability, $\log(p(x_i|x))$, and the gradient,
 205 $\nabla_{x_{\text{GPT}}} \mathcal{L}_{\text{CLIP}}(f(x_{\text{GPT}}))$, where f maps from GPT’s embeddings to CLIP’s embeddings.
 206 This allows us to apply a language prior when optimizing a prompt with CLIP. We addi-
 207 tionally explore the scenario in which we do not use a language prior, by reverting to the
 208 standard case in which we fix $p(x_i|x) = \frac{1}{|\mathbb{V}|}$. In our experiments, we run AutoDAN for 16
 209 steps, which enforces a a maximum token length of 16 due to one-by-one generation of new
 210 tokens. We also utilize a beam search with a beam width of 5.

211 3.4 RANDOM SEARCH

212 Andriushchenko (2023) suggests that such sophisticated strategies may not be critical for
 213 prompt optimization—given enough time, random searches can perform adequately in a

216 variety of settings. Thus, we explore a variant of random search (Rastrigin, 1963). While
217 random search traditionally selects random candidates from within a ball around the
218 current iterate, this approach does not directly map to hard prompting. Due to the curse
219 of dimensionality, true random samples around these high-dimensional embedding spaces
220 are sampled from a ball of with negligible volume around the initial embedding; a nearest
221 neighbor projection would often fail to return a new candidate.

222 In order to address this limitation, we randomly sample from new tokens from the l_0 ball
223 around each element in the sequence \mathbf{XE} . At every iteration, we select a batch of candi-
224 dates and greedily accept the best single-token replacement as the next iterate. We compare
225 the prompt found by Random Searching over the same number of steps as done for PEZ
226 and GCG, determining the best prompt by CLIP similarity in the same way.

228 3.5 PRISM

229 PRISM, proposed by He et al. (2024), highlights that text-to-image generation is not a
230 one-to-one mapping – multiple prompts can describe the same image, and many images can
231 correspond to the same prompt. Rather than relying on discrete token space optimization,
232 PRISM optimizes over a distribution of prompts. Inspired by LLM jailbreaking methods (eg.
233 Chao et al., 2023), PRISM leverages in-context learning in vision-language models (VLMs)
234 to iteratively refine the prompt distribution. This process incorporates the history of refer-
235 ence images, generated prompts, output images from an anchor text-to-image model, and
236 evaluations from a VLM judge, using techniques similar to chain-of-thought (Wei et al.,
237 2022) and textual gradient (Pryzant et al., 2023). After K iterations across N parallel
238 streams, the best-performing prompt is selected using the same VLM judge. In our ex-
239 periments, we use GPT-4-o-mini as the VLM and Stable Diffusion XL-Turbo (Sauer et al.,
240 2023) as the anchor text-to-image model, following He et al. (2024)’s setup with $N = 6$ and
241 $K = 5$. To ensure fair comparisons, we limit the generated prompts to 20 tokens.

243 3.6 CAPTIONING

244 Lastly, we use automated image captions as a proxy for the inverted prompts. Given that a
245 prompt for an image generation model likely encodes information about the setting of the
246 desired image, its subject, its quality, and other properties, we assume that captioning an
247 image provides a human-readable token sequence with some or all of these same properties
248 necessary to generate the image. Moreover, as captioners are typically autoregressive, they
249 have the potential to return an approximate inversion much faster than other methods.

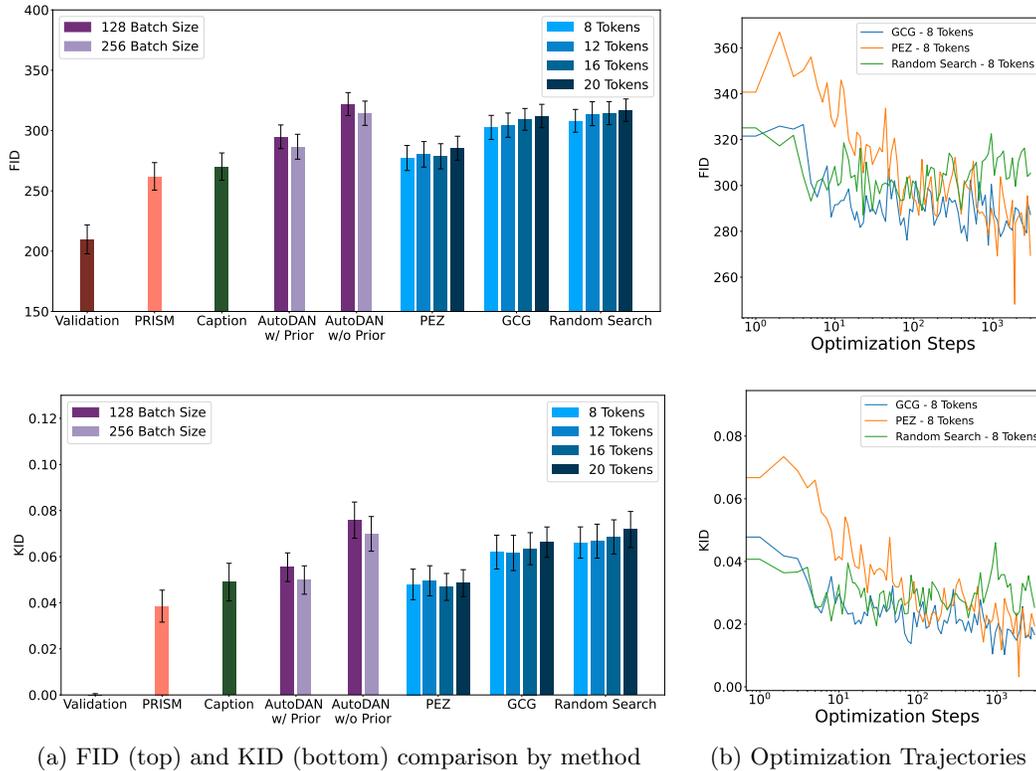
250 Here, we focus on a single model, BLIP-2 (Li et al., 2023). This model is a generic and
251 compute-efficient vision-language pre-training (VLP) method. VLP techniques aim to learn
252 multimodal foundation models on a variety of vision-language tasks. BLIP-2 leverages a
253 trainable module, the Q-former, in order to bridge the gap between a frozen image encoder
254 and a frozen LLM, facilitating image-text matching tasks, image-grounded text generation
255 tasks, and image-text contrastive learning tasks. We prioritize BLIP-2’s image-grounded
256 text generation as the frozen CLIP-style encoder aligns well with the above prompt inversion
257 methods, all of which use frozen CLIP encoders.

259 4 EVALUATION

260 For each optimizer detailed above, we assess their performance across several criteria. Con-
261 sidering the stochastic nature of image generation, we measure the effectiveness of an in-
262 verted prompt by asking the following questions.

- 263 1. How similar (FID (Heusel et al., 2017), KID (Bińkowski et al., 2018)) are images generated
264 with the inverted prompt to images generated by the original prompt?
- 265 2. How well (CLIP (Hessel et al., 2021)) do the inverted prompt and original image align?
- 266 3. How well (Text Embedding Similarity (Reade et al., 2023)) does the semantic content of
267 the inverted prompt align with the semantic content of the original prompt?
268
269

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294



295 Figure 1: Comparison between images generated by inverted prompts and images generated
296 by the original prompts.
297

298
299 We address the stochasticity inherent to the image generation process by averaging the
300 performance of each method across several images generated by the original prompt and
301 the inverted prompts. First, we randomly sample 100 prompts from an existing dataset of
302 prompts¹ used by Stable Diffusion (Rombach et al., 2021).² Given each of these prompts,
303 we generate 10 baseline images for each baseline prompt, and invert each according for all
304 of the 7 methods considered here. Once we have found an inverted prompt for each baseline
305 image, we generate 2 images for each inverted prompt, and finally compute our metrics
306 across the 10 baseline prompts and images and the 20 images based on the 10 inverted
307 prompts. In addition, we choose 75 log scaled time points within the 3000 optimization
308 steps used for PEZ, GCG, and Random Search and repeat our full analysis on a subset of
309 DiffusionDB prompts in order to better understand the convergence of each method.

310 5 EMPIRICAL RESULTS

311
312
313 In this section we present quantitative and qualitative results comparing each method.
314 Across several metrics, we see the quantitative rankings are consistent, but we find upon
315 qualitative examination that these numeric rankings show only a partial picture. Examining
316 the images and the recovered prompts themselves we see trade offs across methods.

317 5.1 QUANTITATIVELY RANKING METHODS

318
319 **Image to Image Comparisons** For image to image comparisons (Figure 1), we analyze
320 images generated by the best early-stopped prompt for each method and the convergence
321

322 ¹We use samples from the Poloclub DiffusionDB dataset of prompt-image pairs (Wang et al.,
323 2022) to find our evaluation prompts.

²All images are generated with StableDiffusion 2-1: stabilityai/stable-diffusion.

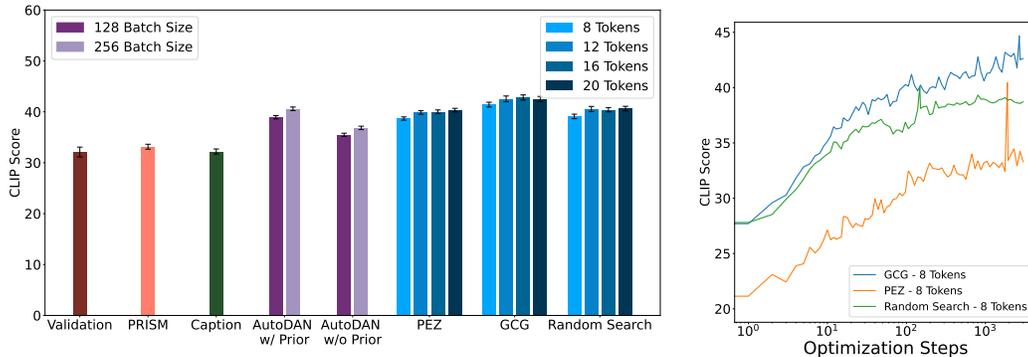


Figure 2: CLIP Similarity between the inverted prompt and images generated by the original prompt. This CLIP Similarity is the objective that each optimizer is maximizing.

rates across our considered image similarity metrics for each algorithm. Our validation set, which consists of the ground truth prompts has an FID of 209.78 and a KID of -0.001 . The KID score in particular tells us that the closer any algorithm gets to a KID of 0, the more similar that prompt will be to the ground truth, whereas, while the ordering may be consistent with FID scores, it is possible that using FID rather than KID may incorrectly show that a method improves over the validation set.

We find that generating images from PRISM prompts provide the most similar images to those generated by the original prompt, with those images generated by BLIP-2 and PEZ as close seconds; PRISM gives average FID and KID values of 262.015 and 0.0385 respectively, while the captioner generates images with average FID and KID values of 270.085 and 0.0489. PEZ follows these with average FID and KID values of 280.392 and 0.0482. In addition, we see a significant gap in performance between AutoDAN with a prior and AutoDAN without a prior, where the former performs much more similarly to the captioners and the latter performing in line with GCG and a Random Search.

Analyzing the objective trajectory over the course of optimization reveals interesting trade-offs. We used a small validation set to determine the number of steps for all algorithms to converge for the given prompts and images used in this study. We determined that all optimizers stop receiving meaningful improvements after 3000 steps. We observe that GCG and a Random Search find a prompt comparable to their best early-stopped prompt within the first 25 steps and then struggle to descend further, analogous to applying too high of a learning rate to optimization problems. On the other hand, PEZ has a slower convergence, but it descends consistently across all steps until it finds prompts that improve over both the GCG and Random Search prompts. Moreover, as PEZ uses a single forward and backward pass, it requires much less time to run than the comparison methods. In other words, PEZ finds prompts that generate images more similar to the ground truth in much less time than all other optimizers considered here, except for the BLIP-2 Captioner.

Text to Image Comparisons When we focus on the alignment between the text and images we see an interesting trend emerge. We first compare the CLIP similarity between the inverted prompts and the original image (in the top of Figure 2). Note that this is the optimization objective used across all optimizers.

We find that all optimizers do a good job maximizing their objective. While AutoDAN without a language prior performs the worst over all optimizers, it still does a better job of maximizing the CLIP similarity over the validation set, PRISM, and the BLIP2 Captioner. Optimizing the objective with GCG and AutoDAN with a language prior performs the best over the discrete optimizers, with PEZ coming a close third. The contrast between the performance of each optimizer on their objective and their relative lack of performance across the image-to-image and text-to-text metrics suggests that the CLIP objective is acting as a poor proxy for finding prompts for generative image models. While there may be room for improvement over the CLIP objective for this task, this comparison allows us to take

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

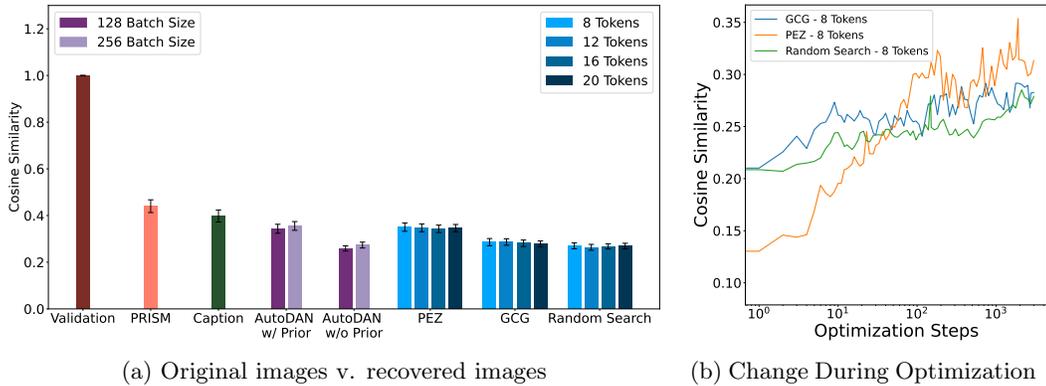


Figure 3: Cosine Similarity between text embeddings for the original and inverted prompts. Based on the metric used by Reade et al. (2023)

a better look at the convergence rates of all optimizers on their objective. Just as in the image-to-image comparison, GCG and Random Search quickly find a good prompt (within 20 steps) and then very slowly improve from there.

Yet PEZ follows a much more gradual curve, with sharp peaks when new optima are found. As these are log scaled in their x-axes, we do not see all peaks except for the early stopped result. The average prompt found with PEZ is much lower than the comparison methods, but the peaks are in line with the other methods. Additionally, GCG and Random Search again very quickly within the first 20 steps and then very slowly update from there. This overreliance on early-stopping may be a weakness for PEZ. Rather than oscillating tightly around the optima, PEZ oscillates wildly around high quality prompts. In essence, PEZ may better explore the prompt space, while methods incorporating fine-grained search (such as GCG) are more adept at exploiting it.

Text to Text Comparisons Lastly, we compare the similarity in the text of the found prompts to the ground truth prompts. Figure 3 shows the cosine similarity between the text embeddings³ of the found prompts and the ground truth prompts.

Just as in the image-to-image case, we find that using responses from PRISM as the inverted prompt outperforms all of our comparisons, with a cosine similarity of 0.440 to the original prompt; the BLIP-2 captioner comes in second with a cosine similarity of 0.397 to the original prompt. AutoDAN with a language prior and PEZ follow behind with respective similarities of 0.355 and 0.346 averaged across all lengths. GCG, Random Search, and AutoDAN without a language prior remain clustered together in terms of their performance. Moreover, when looking at their convergence rate, we see the same story as above. GCG and Random Search very quickly ascend, and while PEZ ascends more slowly it eventually exceeds GCG and Random Search in their performance within the first 100 optimizations steps of its allowed 3000 steps.

Hyperparameter Choice Each optimizer has various hyperparameters that can be adjusted. Intuitively, it may seem that the number of free tokens (tokens that can be optimized) is a particularly relevant parameter. However, this is not always the case. In line with the findings by Wen et al. (2023), PEZ’s performance remains fairly independent of the number of free tokens up to a certain point. Given that we allowed no fewer than 8 free tokens, there may be a performance drop-off if we further decrease the number of tokens. However, PEZ’s performance across all metrics does not appear to be significantly dependent on its free tokens. For context, the BLIP-2 captioner, which does not have a fixed length, can serve as a benchmark for reasonable prompt length. Its captions average 10.6 tokens using CLIP’s tokenizer, and 5.8 tokens after removing stop-words. Similarly, we see that neither

³Embeddings were computed using `sentence-transformers/all-MiniLM-L6-v2` to be in line with (Reade et al., 2023)

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

	Original Prompt	a friendly goblin with a big ((human)) nose and wide eyes, dark hairs, big ears, covered in branches and moss, portrait by daniel docciu and dave dorman and jeff easley
	PRISM	Friendly green goblin face, smiling, tangled branches, soft forest background
	BLIP-2	green troll in a tree with leaves and branches on it's head and a smiley face on it's face
	AutoDAN w/ Prior	character from Magic Treefolk depicting Green Elf head with smile during 2015 promotional image walkthrough art group image
	AutoDAN w/o Prior	animated jester troll grass wordpress goblin frightening branch artwork today)) oman ly 6 jester head. _ reid /
	PEZ	newmlb mtgnflrevealed loki revealed reveal).. goofy smiling scary creepy ytree ! orc =)) arbormates
	GCG	typically greener donny recent reported atrist..... frightening cohen substantially ∴ . kal ears googmirrowickedcriticalrole trees believes
	Random Search	cantenzenegger oaks]ñodo grin grassy goblin ytless...(! ı.; sends iconforeveryimp huskdns

Table 1: Example images and corresponding 20-token prompts. Each image is generated by the *original prompt* and we show examples of the inversion result from each method. Other than AutoDAN with the language prior applied, no discrete optimizer produces more human-readable prompts than another despite the quantitative differences in their performance.

GCG nor Random Search dependent on the number of free tokens. With no metric showing a statistically significant improvement for when adding more or fewer optimizable tokens.

AutoDAN, like the BLIP-2 captioner, can return a variable number of tokens due to early stopping. However, we use AutoDAN with a beam search, where each individual step uses a much smaller batch size for its fine-grained search, unlike GCG and Random Search. The latter have a 512 batch size, while AutoDAN with 4 beams and 128 tokens evaluates the same number of tokens for the fine-grained search. Increasing AutoDAN to 256 tokens, effectively doubling the number of tokens it searches over compared to GCG and Random Search, results in a small improvement across the board. Based on the optimal batch sizes described by Zou et al. (2023), further improvements in AutoDAN might be achieved by allowing a 512 batch size. However, there are likely diminishing returns, especially as computation time increases with larger batch sizes.

5.2 QUALITATIVELY ASSESSING INVERTED PROMPTS

In the quantitative evaluation above, we show that PRISM and the captioner return prompts that may be better across several metrics compared to searching for a prompt via discrete optimization. Here, we show a qualitative example (Table 1) of an image generated by one of the ground truth prompts and the different results that each method find. Other than AutoDAN with the language prior applied, no discrete optimizer produces human-readable prompts despite the quantitative similarities in their performance. We thus separate this subsection into natural language and keyword-based prompts that without a language prior.

PRISM provides prompts that are exceptionally more detailed than the comparison methods; opting for short descriptive clauses rather than the full sentences that BLIP-2 uses. As described above, when the length of a prompt is limited, the additional stop words required by full-sentence prompts reduce the number of concepts that can be included in a prompt,

486 significantly affecting the final image. On the other hand, AutoDAN’s language prior seems
487 to find natural language prompts that evoke the imagery described by the image, “char-
488 acter from Magic Treefolk...” does not describe anything from the image (to the best of
489 our knowledge there is not media called Magic Treefolk), but if such media existed then we
490 would not be surprised to find that something called Magic Treefolk included depictions of
491 goblins or other forest critters.

492 When comparing each recovered prompt to the original prompt, there is often a significant
493 amount of information lost during the generation process that is unrecoverable. Both Ran-
494 dom Search and PEZ capture basic information such as “trees” or “green”. These methods
495 try to include the single tokens that encode as much information as possible. Similarly
496 to the “Magic Treefolk” example above, GCG uses the token “criticalrole” for a similar
497 purpose. Critical Role [Mercer \(2015–present\)](#), a ‘Dungeons & Dragons’-based web series,
498 embeds a relationship between the prompt and creatures found in Dungeons & Dragons
499 through a single token. Moreover, without the need for a language prior, it does not need to
500 waste tokens fitting ‘criticalrole’ into a coherent sentence. Yet, it may cause an overreliance
501 on these ‘keyword’ tokens and allow unrelated tokens such as ‘goog’ to be included in a
502 prospective prompt. This comparison may shed light on why PEZ outperforms GCG and
503 random search, as PEZ appears to stay more on topic. PEZ includes the tokens “loki”,
504 “tree”, “arbor”, “scary”, “goofy”, “orc” and “smiling”, while GCG and Random Search
505 do not provide significantly more specificity than “green”, “trees”, and “criticalrole”; and
506 “oaks”, “goblin” and “grassy” respectively. At its core, PEZ is a projected gradient de-
507 scent method, using common optimizers, such as SGD or Adam with a weight decay. This
508 approach encourages some form of regularization in its optimization, that discourages the
509 one-and-done approach that GCG and Random Search seem to use, where they discourage
510 repeating the same general concepts or tokens in a prompt.

511 6 DISCUSSION

512 Our results prompt discussion on the practical implications, the limitations, and the future
513 directions related to prompt inversion. To begin, someone interested in finding good prompts
514 from images can conclude from our work that image captioning tools are a good approach.
515 They are fast, as the heavy lifting is done ahead of time in training these models rather
516 than optimizing anything per image in deployment. They also best capture natural sounding
517 language, a goal that discrete optimizers might better incorporate as these tools mature.

518 The limitations of our work center mostly on the fact that the diffusion and image-text
519 embedding space is so heavily driven by only a few models. As the set of state-of-the-art
520 large text-prompted image generations models grows, the trends we report may no longer
521 hold. In the same vein, small variations in the optimization strategies could have large
522 impacts on these results. In short, like any empirical benchmark results, our findings are
523 subject to change as the field progresses.

524 Here, we enumerate several questions and quirks arising from our work that warrant further
525 investigation. First, [Zou et al. \(2023\)](#) report that GCG is effective at jailbreaking LLMs
526 and PEZ is not. This stands in stark contrast to these two methods relative performance
527 at prompt inversion. Why might optimizing over natural language be so different in these
528 settings? This could be a difference in the particular models or in the loss landscapes. Sec-
529 ond, GCG and random search perform so similarly begging the question why does gradient
530 information make so little difference? The intuition that the gradient signal is informative
531 comes from observing the success of PEZ, so why is the combination of search and gradient-
532 based optimization in GCG leave it so similar to random search alone? Finally, we posit
533 that there is a lot of room for improvement. In other words, prompt inversion is far from
534 solved, and it makes for a great test bed for new discrete optimization approaches.

535 REFERENCES

536 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini
537 Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning

540 transferable visual models from natural language supervision. In *International conference*
541 *on machine learning*, pages 8748–8763. PMLR, 2021.

542

543 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-
544 image pre-training with frozen image encoders and large language models. *arXiv preprint*
545 *arXiv:2301.12597*, 2023.

546 Yutong He, Alexander Robey, Naoki Murata, Yiding Jiang, Joshua Williams, George J
547 Pappas, Hamed Hassani, Yuki Mitsufuji, Ruslan Salakhutdinov, and J Zico Kolter. Au-
548 tomated black-box prompt engineering for personalized text-to-image generation. *arXiv*
549 *preprint arXiv:2403.19103*, 2024.

550

551 Zezhong Fan, Xiaohan Li, Chenhao Fang, Topojoy Biswas, Kaushiki Nag, Jianpeng Xu, and
552 Kannan Achan. Prompt optimizer of text-to-image diffusion models for abstract concept
553 understanding. *arXiv preprint arXiv:2404.11589*, 2024.

554 Kihyuk Sohn, Nataniel Ruiz, Kimin Lee, Daniel Castro Chin, Irina Blok, Huiwen Chang,
555 Jarred Barber, Lu Jiang, Glenn Entis, Yuanzhen Li, et al. Styledrop: Text-to-image
556 generation in any style. *arXiv preprint arXiv:2306.00983*, 2023.

557

558 Yuxin Wen, Neel Jain, John Kirchenbauer, Micah Goldblum, Jonas Geiping, and Tom
559 Goldstein. Hard prompts made easy: Gradient-based discrete optimization for prompt
560 tuning and discovery. *arXiv preprint arXiv:2302.03668*, 2023.

561 Shweta Mahajan, Tanzila Rahman, Kwang Moo Yi, and Leonid Sigal. Prompting hard or
562 hardly prompting: Prompt inversion for text-to-image diffusion models. In *Proceedings*
563 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6808–
564 6817, 2024.

565

566 Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. Universal and transferable
567 adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023.

568 Sicheng Zhu, Ruiyi Zhang, Bang An, Gang Wu, Joe Barrow, Zichao Wang, Furong Huang,
569 Ani Nenkova, and Tong Sun. Autodan: Automatic and interpretable adversarial attacks
570 on large language models. *arXiv preprint arXiv:2310.15140*, 2023.

571 Maksym Andriushchenko. Adversarial attacks on gpt-4 via simple random search. 2023.

572

573 Avi Schwarzschild, Zhili Feng, Pratyush Maini, Zachary C Lipton, and J Zico Kolter. Re-
574 thinking llm memorization through the lens of adversarial compression. *arXiv preprint*
575 *arXiv:2404.15146*, 2024.

576

577 Aly M Kassem, Omar Mahmoud, Niloofar Mireshghallah, Hyunwoo Kim, Yulia Tsvetkov,
578 Yejin Choi, Sherif Saad, and Santu Rana. Alpaca against vicuna: Using llms to uncover
579 memorization of llms. *arXiv preprint arXiv:2403.04801*, 2024.

580

581 R Gary Parker and Ronald L Rardin. *Discrete optimization*. Elsevier, 2014.

582

583 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models.
584 *arXiv preprint arXiv:2010.02502*, 2020.

585

586 Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text
587 inversion for editing real images using guided diffusion models. In *Proceedings of the*
588 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6038–6047,
589 2023.

589

590 Joshua Nathaniel Williams and J. Zico Kolter. Fuse-ing language models: Zero-shot adapter
591 discovery for prompt optimization across tokenizers, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2408.04816)
592 [abs/2408.04816](https://arxiv.org/abs/2408.04816).

593

593 Collin Zhang, John X Morris, and Vitaly Shmatikov. Extracting prompts by inverting llm
outputs. *arXiv preprint arXiv:2405.15012*, 2024.

594 Walter Reade, Will Cukierski, and Ashley Chow. Stable diffusion - image to prompts, 2023.
595 URL <https://kaggle.com/competitions/stable-diffusion-image-to-prompts>.
596

597 Sahil Verma, Varich Boonsanong, Minh Hoang, Keegan Hines, John Dickerson, and Chirag
598 Shah. Counterfactual explanations and algorithmic recourses for machine learning: A
599 review. *ACM Computing Surveys*, 2020.

600 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Model-agnostic interpretability
601 of machine learning. *arXiv preprint arXiv:1606.05386*, 2016.
602

603 Joshua N Williams, Molly FitzMorris, Osman Aka, and Sarah Laszlo. Drawl: Understanding
604 the effects of non-mainstream dialects in prompted image generation. *arXiv preprint*
605 *arXiv:2405.05382*, 2024.

606 Daniel Khashabi, Shane Lyu, Sewon Min, Lianhui Qin, Kyle Richardson, Sean Welleck,
607 Hannaneh Hajishirzi, Tushar Khot, Ashish Sabharwal, Sameer Singh, et al. Prompt
608 waywardness: The curious case of discretized interpretation of continuous prompts. *arXiv*
609 *preprint arXiv:2112.08348*, 2021.

610 LA Rastrigin. The convergence of the random search method in the extremal control of a
611 many parameter system. *Automaton & Remote Control*, 24:1337–1342, 1963.
612

613 Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J Pappas, and
614 Eric Wong. Jailbreaking black box large language models in twenty queries. *arXiv preprint*
615 *arXiv:2310.08419*, 2023.

616 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V
617 Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language
618 models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
619

620 Reid Pryzant, Dan Iter, Jerry Li, Yin Tat Lee, Chenguang Zhu, and Michael Zeng. Au-
621 tomatic prompt optimization with” gradient descent” and beam search. *arXiv preprint*
622 *arXiv:2305.03495*, 2023.

623 Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffu-
624 sion distillation. *arXiv preprint arXiv:2311.17042*, 2023.
625

626 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp
627 Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equi-
628 librium. *Advances in neural information processing systems*, 30, 2017.

629 Mikołaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying
630 mmd gans. *arXiv preprint arXiv:1801.01401*, 2018.
631

632 Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A
633 reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*,
634 2021.

635 Zijie J. Wang, Evan Montoya, David Munechika, Haoyang Yang, Benjamin Hoover, and
636 Duen Horng Chau. Large-scale prompt gallery dataset for text-to-image generative mod-
637 els. *arXiv:2210.14896 [cs]*, 2022. URL <https://arxiv.org/abs/2210.14896>.
638

639 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer.
640 High-resolution image synthesis with latent diffusion models, 2021.

641 Matthew Mercer. Critical role, 2015–present. URL <https://www.criticalrole.com>. Geek
642 & Sundry / Critical Role Productions.
643
644
645
646
647