
Discovering the Hidden Vocabulary of DALLE-2

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We discover that DALLE-2 seems to have a hidden vocabulary that can be
2 used to generate images with absurd prompts. For example, it seems that
3 Apoploe vesrreaitais means birds and Contarra ccetnxiams luryca
4 tanniounons (sometimes) means bugs or pests. We find that these prompts
5 are often consistent in isolation but also sometimes in combinations. We present
6 our black-box method to discover words that seem random but have some corre-
7 spondence to visual concepts. This creates important security and interpretability
8 challenges.



Figure 1: Images generated with the prompt: “Apoploe vesrreaitais eating Contarra ccetnxiams luryca tanniounons”. We discover that DALLE-2 has its own vocabulary where Apoploe vesrreaitais means birds and Contarra ccetnxiams luryca tanniounons (sometimes) means bugs. Hence, this prompt means “Birds eating bugs”.

9 1 Introduction

10 DALLE [1] and DALLE-2 [2] are deep generative models that take as input a text caption and
11 generate images of stunning quality that match the given text. DALLE-2 uses Classifier-Free
12 Diffusion Guidance [3] to generate high quality images. The conditioning is the CLIP [4] text
13 embeddings for the input text.

14 A known limitation of DALLE-2 is that it struggles with text. For example, text prompts such as:
15 “An image of the word airplane” often lead to generated images that depict gibberish text.
16 We discover that this produced text is not random, but rather reveals a hidden vocabulary that the

17 model seems to have developed internally. For example, when fed with this gibberish text, the model
18 frequently produces airplanes.

19 Some words from this hidden vocabulary can be learned and used to create absurd prompts that
20 generate natural images. For example, it seems that `Apoploe vesrreaitais` means birds and
21 `Contarra ccetnxniam luryca tanniounons` (sometimes) means bugs or pests. We found that
22 we can generate images of cartoon birds with prompts like `An image of a cartoon apoploe`
23 `vesrreaitais` or even compose these terms to create birds eating bugs as shown in Figure 1.

24 **2 Discovering the DALLE-2 Vocabulary**

25 We provide a simple method to discover words of the DALLE-2 vocabulary. We use (in fact, we only
26 have) query access to the model, through the API. We describe the method with an example. Assume
27 that we want to find the meaning of the word: `vegetables`. Then, we can prompt DALLE-2 with
28 one of the following sentences (or a variation of those):

- 29 • `“A book that has the word vegetables written on it.”`
- 30 • `“Two people talking about vegetables, with subtitles.”`
- 31 • `“The word vegetables written in 10 languages.”`

32 For each of the above prompts, DALLE-2 usually creates images that have some text written text
33 on it. The written text often seems gibberish to humans, as mentioned in the original DALLE-2
34 paper [2] and also in the preliminary evaluation of the system by [5]. However, we make the
35 surprising observation that this text is not as random as it initially appears. In many cases, it is
36 strongly correlated to the word we are looking to translate. For example, if we prompt DALLE-2 with
37 the text: `“Two farmers talking about vegetables, with subtitles.”`, we get the image
38 shown in Figure 2(a). We parse the text that appears in the images and we prompt the model with it as
39 shown in Figure 2(b), (c). It seems that `Vicootes` means vegetables and `Apoploe vesrreaitais`
40 means birds. It appears that the farmers are talking about birds that interfere with their vegetables.

41 We note that this simple method doesn’t always work. Sometimes, the generated text gives random
42 images when prompted back to the model. However, we found that with some experimentation
43 (selecting some words, running different produced texts, etc.) we can usually find words that appear
44 random and are correlated with some visual concept (at least under some contexts). We encourage
45 the interested readers to refer to the Limitations Section for more information.

46 **3 A Preliminary Study of the Discovered Vocabulary**

47 **Compositionality.** From the previous example, we learned that `Apoploe vesrreaitais` seems to
48 mean birds. By repeating the experiment with the prompt about farmers, we also learn that: `Contarra`
49 `ccetnxniam luryca tanniounons` may mean pests or bugs. An interesting question is whether
50 we can compose these two concepts in a sentence, as we could do in an ordinary language. In Figure
51 1, we illustrate that this is possible, at least sometimes. The sentence: `“Apoploe vesrreaitais`
52 `eating Contarra ccetnxniam luryca tanniounons”` gives images in which birds are eating
53 bugs. We found that this happens for some, but not all of the generated images.

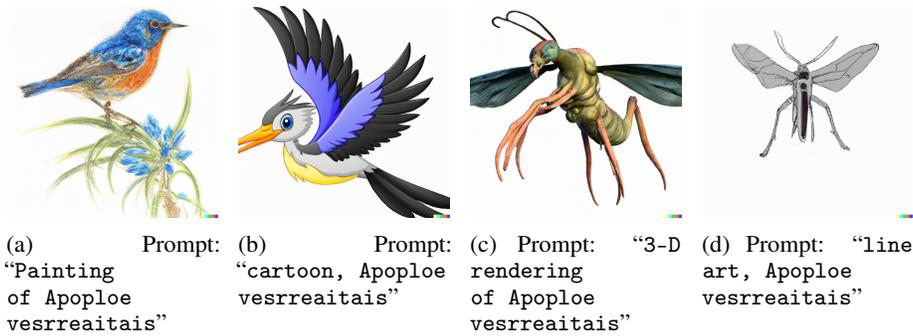
54 **Style Transfer.** DALLE-2 is capable of generating images of some concept under different styles
55 that can be specified in the prompt [2]. For example, one might ask for a photorealistic image
56 of an apple or a line-art showing an apple. We test whether the discovered words, (e.g. `Apoploe`
57 `vesrreaitais`) correspond to visual concepts that can be transformed into different styles, depending
58 on the context of the prompt. The results of this experiment are shown in Figure 3. It seems that the
59 prompt sometimes leads to flying insects as opposed to birds.

60 **Text’s consistency with the caption and the generated image.** Recall the example with the
61 farmers. The prompt was: `“Two farmers talking about vegetables, with subtitles.”`.
62 From this example, we discovered the word `vegetables`, but also the word `birds`. It is very plausible
63 that two farmers would be talking about birds and hence this opens the very interesting question of
64 whether the text outputs of DALLE-2 are consistent with the text conditioning and the generated



(a) Image generated with the prompt: “Two farmers talking about vegetables, with subtitles.” (b) Image generated with the prompt: “Vicootes.” (c) Image generated with the prompt: “Apoploe vesrreaitais.”

Figure 2: Illustration of our method for discovering words that seem random but can be understood by DALLE-2. We first query the model with the prompt: “Two farmers talking about vegetables, with subtitles.”. The model generates an image with some gibberish text on it. We then prompt the model with words from this generated image, as shown in (b), (c). It seems that Vicootes means vegetables and Apoploe vesrreaitais means birds. Possibly farmers are talking about birds that interfere with their vegetables.



(a) Prompt: “Painting of Apoploe vesrreaitais” (b) Prompt: “cartoon, Apoploe vesrreaitais” (c) Prompt: “3-D rendering of Apoploe vesrreaitais” (d) Prompt: “line art, Apoploe vesrreaitais”

Figure 3: Illustration of DALLE-2 generations for Apoploe vesrreaitais under different styles. The visual concept of “something that flies” is maintained across the different styles.

65 image. Our initial experiments show that sometimes we get gibberish text that translates to visual
 66 concepts that match the caption that created the gibberish text in the first place. For example, the
 67 prompt: “Two whales talking about food, with subtitles.” generates an image with the
 68 text “Wa ch zod ahaakes rea.” (or at least something close to that). We feed this text as prompt
 69 to the model and in the generated images we see seafood. This is shown in Figure 4. It seems that
 70 the gibberish text indeed has a meaning that is sometimes aligned with the text-conditioning that
 71 produced it.

72 4 Security and Interpretability Challenges

73 There are many interesting directions for future research. It was not clear to us if some of the
 74 gibberish words are misspellings of normal words in different languages, but we could not find any
 75 such examples in our search. For many of the prompts, the origins of these words remains confusing
 76 and some of the words were not as consistent as others in our preliminary experiments. Another
 77 interesting question is if Imagen [6] has a similar hidden vocabulary given that it was trained with a
 78 language model as opposed to CLIP. We conjecture that our prompts are adversarial examples for
 79 CLIP’s [4] text encoder, i.e. the vector representation of Apoploe vesrreaitais is close to the
 80 representation of bird. It is natural to use other methods (e.g. white box) of adversarial attacks on
 81 CLIP to generate absurd text prompts that produce target images in DALLE2.



Figure 4: Left: Image generated with the prompt: “Two whales talking about food, with subtitles.”. Right: Images generated with the prompt: “Wa ch zod ahaakes rea.”. The gibberish text, “Wa ch zod ahaakes rea.”, produces images that are related to the caption and the visual output of the first image.

82 **Robustness and Limitations.** One of the central questions is how consistent this method is.
 83 For example, our preliminary study shows that prompts like *Contarra cctetxniam luryca*
 84 *tanniounons* sometimes produces bugs and pests (about half of the generated images) and sometimes
 85 different images, mostly animals. We found that *Apoploe vesrreaitais* is much more robust
 86 and can be combined in various ways as we show. We also want to emphasize that finding other
 87 robust prompts is challenging and requires a lot of experimentation. In our experiments we tried
 88 various ways of making DALLE generate images, selected parts of the generated text and tested its
 89 consistency. However, even if this method works for a few gibberish prompts (that are hard to find),
 90 this is still a big interpretability and security problem. If a system behaves in wildly unpredictable
 91 ways, even if this happens rarely and under unexpected conditions like gibberish prompts, this is still
 92 a significant concern, especially for some applications.

93 The first security issue relates to using these gibberish prompts as backdoor adversarial attacks or
 94 ways to circumvent filters. Currently, Natural Language Processing systems filter text prompts that
 95 violate the policy rules and gibberish prompts may be used to bypass these filters. More importantly,
 96 absurd prompts that consistently generate images challenge our confidence in these big generative
 97 models. Clearly more foundational research is needed in understanding these phenomena and creating
 98 robust language and image generation models *that behave as humans would expect*.

99 5 Conclusions and Future Work

100 In this work, we showed that the state-of-the-art text-conditional generative model DALLE-2 has a
 101 hidden vocabulary that be used to generate images with prompts that cannot be parsed by humans. We
 102 developed a suprisingly simple method that, given only query access to the model, sometimes help us
 103 extract gibberish words that correspond to consistent visual concepts. Recently, powerful open-source
 104 text-to-image models have been released for everyone to use [7]. In the future, we plan to explore
 105 more powerful methods (that use access to the weights) to discover gibberish text that corresponds to
 106 concepts of interest. We are also interested in understanding how this hidden vocabulary is shaped.

107 **References**

- 108 [1] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark
109 Chen, and Ilya Sutskever. Zero-shot text-to-image generation, 2021.
- 110 [2] Aditya Ramesh, Prfulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical
111 text-conditional image generation with clip latents, 2022.
- 112 [3] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop
113 on Deep Generative Models and Downstream Applications*, 2021.
- 114 [4] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
115 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.
116 Learning transferable visual models from natural language supervision, 2021.
- 117 [5] Gary Marcus, Ernest Davis, and Scott Aaronson. A very preliminary analysis of dall-e 2, 2022.
- 118 [6] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed
119 Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim
120 Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image
121 diffusion models with deep language understanding, 2022.
- 122 [7] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
123 resolution image synthesis with latent diffusion models, 2021.