# CLIPDraw: Exploring Text-to-Drawing Synthesis through Language-Image Encoders

**Kevin Frans**[1,2], **L. B. Soros**[1] **and Olaf Witkowski**[1,3,4]

[1]Cross Labs, Cross Compass Ltd., Tokyo, Japan
[2]Massachusetts Institute of Technology, Cambridge, MA, USA
[3]Earth-Life Science Institute, Tokyo Institute of Technology, Japan
[4]College of Arts and Sciences, University of Tokyo, Japan
kvfrans@csail.mit.edu

"A drawing of a cat".　　"Horse eating a cupcake".　　"A 3D rendering of a temple".　　"Family vacation to Walt Disney World".　　"Self".

**Various drawings synthesized by CLIPDraw**, along with the corresponding description prompts used. CLIPDraw synthesizes images from text by performing gradient descent over a set of RGBA Bézier curves, with the goal of minimizing cosine distance between the CLIP encodings of generated images and description prompts. CLIPDraw does not require learning a new model, and can generally synthesize images within a minute on a typical GPU.

## Abstract

This work presents CLIPDraw, an algorithm that synthesizes novel drawings based on natural language input. CLIPDraw does not require any training; rather a pre-trained CLIP language-image encoder is used as a metric for maximizing similarity between the given description and a generated drawing. Crucially, CLIPDraw operates over vector strokes rather than pixel images, a constraint that biases drawings towards simpler human-recognizable shapes. Results compare between CLIPDraw and other synthesis-through-optimization methods, as well as highlight various interesting behaviors of CLIPDraw, such as satisfying ambiguous text in multiple ways, reliably producing drawings in diverse artistic styles, and scaling from simple to complex visual representations as stroke count is increased. Code for experimenting with the method is available at: **https://colab.research.google.com/github/kvfrans/clipdraw/blob/main/clipdraw.ipynb**

## 1. Introduction

As humans, when we hear a description of a scene, it'e easy to picture what it may look like in our heads. Conversely, when we construct a mental image, it's easy to describe that scene in words. At some level, humans have a deeply coupled representation for both textual and visual structures that is key to understanding our everyday world.

The recent introduction of CLIP (Radford et al., 2021), a dual language-image encoder, is a large step towards unifying textual and visual information. In a CLIP model, both text and images are mapped onto the same representational space, thus enabling the similarity between images and textual descriptions to be measured. When trained on large amounts of data, CLIP representations have been shown to solve a robust range of image-based recognition tasks.

This work presents *CLIPDraw*, an algorithm that synthesizes novel drawings based on natural language input. CLIPDraw does not require any training; rather a pre-trained CLIP model is used as a metric for maximizing similarity between the given description and a generated drawing. Rather than photorealistic images, CLIPDraw aims to syn-

thesize simple drawings that nevertheless match the prompt. Thus, CLIPDraw optimizes a set of vector strokes rather than pixel images, a constraint that biases drawings towards simple human-recognizable shapes.

The aim of this work is to present CLIPDraw a testbed for exploring language-image relationships and synthesizing AI-assisted artwork, as well as to showcase various nuances of the method. Results compare between CLIPDraw and other optimization-based text-to-image methods, along with highlighting several interesting behaviors:

- By adjusting descriptive adjectives, such as "watercolor" or "3D rendering", CLIPDraw produces drawings of vastly different styles.

- CLIPDraw often matches the description prompt in creative ways, such as writing words from the prompt inside the image itself, or interpreting ambiguous nouns in multiple ways.

- Running CLIPDraw with a low stroke count results in cartoonish drawings, while high stroke counts tend to result in realistic renderings.

- By giving CLIPDraw abstract prompts, such as "happiness" or "self", we can examine what visual concepts the CLIP model associates with them.

- CLIPDraw behavior can be further controlled through the use of negative prompts, such as "a messy drawing", to encourage the opposite behavior.

## 2. Related Work

**Text-to-Image Synthesis.** This work greatly draws from the field of text-to-image synthesis, whose primary aim is to generate images that correctly match a given textual description. In recent years, focus has been on methods that aim to learn a direct text-to-image mapping function, often in the form of a conditional GAN (Goodfellow et al., 2014a; Mirza & Osindero, 2014; Reed et al., 2016; Frolov et al., 2021). Commonly used datasets include Oxford-120 Flowers (Nilsback & Zisserman, 2008), CUB-200 Birds (Wah et al., 2011), and COCO (Lin et al., 2014), all of which are comprised of natural images along with a set of captions describing them. While GAN-based methods have enabled considerable progress towards photorealistic image synthesis, strong autoregressive models have achieved similar quality results (Oord et al., 2017; Chen et al., 2020), with the recent DALL-E model (Ramesh et al., 2021) showcasing the benefit of scaling text-to-image synthesis networks to a large capacity. In comparison to text-to-image generative models, which require large amounts of training, this work follows the framework of *synthesis through optimization*, in which images are generated through evaluation-time optimization against a given metric.

**Synthesis Through Optimization.** Instead of directly learning an image generation network, an alternative method of image synthesis is to optimize towards a matching image during evaluation time. This framework is often referred to as *activation maximization* (Erhan et al., 2009; Nguyen et al., 2016; Mordvintsev et al., 2015), where a random image is optimized through backpropogation to increase certain neuron activations of a pretrained network. Activation-maximization methods have produced highly realistic images, however it is a challenge to understand the meaning of a neuron activation. CLIPDraw builds off a set of methods where rather than maximizing an activation, the objective is to minimize the distance between the produced image and a given description phrase, as defined by a powerful CLIP language-image encoder (Fernando et al., 2021; Murdock, 2021; Galatolo et al., 2021). A key issue in synthesis through optimization is that the produced images often leave the space of natural images (Nguyen et al., 2015), or fool the system through adversarial means (Goodfellow et al., 2014b), thus a body of work aims to discover 'natural image priors' to constrain which images may be produced (Nguyen et al., 2016; 2017). While a typical solution is to constrain optimization to the generative space of a GAN, this setup can be expensive to evaluate, and only allows synthesis of images producible by the GAN generator. Because CLIPDraw focuses on synthesizing drawings rather than realistic pictures, CLIPDraw instead limits optimization to a set of vector curves. This constraint results in stroke-based images, which must capture larger features such as shapes and outlines, rather than fine-grained textures.

**Vector Graphics.** This work builds largely from work by Li et al. (2020), which introduces a differentiable renderer for vector graphics. Image generation methods that operate over vector images have traditionally required a vector-based dataset, however recent work has shown how differentiable renderers can be used to bypass this limitation (Reddy et al., 2021; Shen & Chen, 2021). CLIPDraw uses a differentiable renderer as a representation for generating drawings; namely a set of RGBA Bézier curves are optimized rather than a matrix of pixels.

## 3. Method

The objective of CLIPDraw is to synthesize a drawing that matches a given description prompt (see Front Figure). Specifically, a pre-trained CLIP model is used as a judge. A CLIP model contains two networks – an image encoder and a textual encoder – which both map their respective inputs into a shared encoding space of a 512-length vector. Similarity is measured via the cosine distance between two encodings. Thus, the goal of CLIPDraw is to produce an image which, when encoded via CLIP, matches the CLIP encoding of the given description prompt.
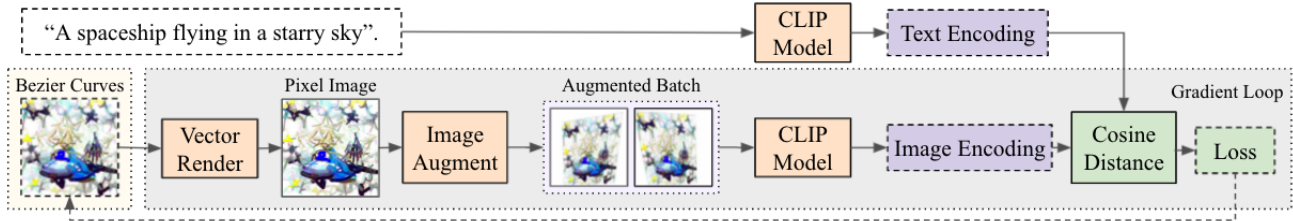
*Figure 1.* **CLIPDraw iteratively synthesizes images through evaluation-time gradient descent.** Starting from a random set of Bézier curves, the position and colors of the curves are optimized so that the generated drawings best match the given description prompt. Before being passed into the CLIP encoder, drawings are augmented into multiple perspective-shifted copies.
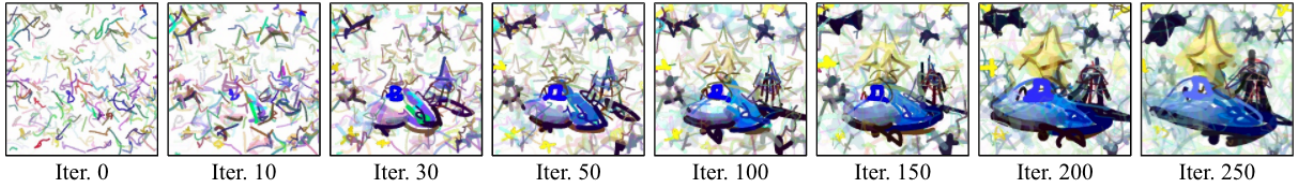


| Iter. 0 | Iter. 10 | Iter. 30 | Iter. 50 | Iter. 100 | Iter. 150 | Iter. 200 | Iter. 250 |

*Figure 2.* **A typical CLIPDraw run gradually forms messy curves into concrete shapes.** In this example, the drawing first develops a background of star-shaped structures, which eventually develop into a large spaceship. Near the later iterations, more pronounced stars appear, in addition to a Darth Vader-like figure riding the spaceship.

Drawings in CLIPDraw are represented by a set of differentiable RGBA Bézier curves, following the method by Li et al. (2020). Each curve is parametrized by between 3 to 5 control points, along with thickness and an RGBA color vector. Drawings initially begin with curves randomly distributed throughout the image, with a white background as the default color. During optimization, the number of curves and control points is fixed, however the positions of the points along with the thickness and color vectors can be optimized via gradient descent.

The CLIPDraw algorithm (Algorithm 1) works by running evaluation-time gradient descent, as shown in Figure 1. First, the description phrase is encoded via the CLIP model, and a random set of $N$ Bézier curves are initialized. During each iteration, the curves are rendered to a pixel image via the differentiable renderer, and the resulting image is then duplicated $D$ times and augmented by a random perspective shift and random crop-and-resize. The resulting batch of augmented images is passed into the CLIP image encoder, and the cosine distances to the description phrase are summed to form the loss value. Because all operations are differentiable, gradient descent can be run through the entire loop, optimizing the parameters of the curves to decrease loss. This procedure is repeated $I$ times, until convergence.

The goal of the image augmentation is to force drawings to remain recognizable when viewed through various distortions. Without image augmentation, synthesis-through-optimization methods often result in adversarial

---

**Algorithm 1** CLIPDraw

> **Input:** Description Phrase $desc$; Iteration Count $I$; Curve Count $N$; Augment Size $D$; Pre-trained CLIP model.
> **Begin:**
> Encode Description Phrase. *EncPhr = CLIP(desc)*
> Initialize Curves. *Curves$_{0..N}$ = RandomCurve()*
> **for** $i = 0$ **to** $I$ **do**
>     Render Curves to Pixels. *Pixels = DiffRender(Curves)*
>     Augment the Image. *AugBatch$_{0..D}$ = Augment(Pixels)*
>     Encode Image. *EncImg = CLIP(AugBatch)*
>     Compute Loss. *Loss = −CosineSim(EncPhr, EncImg)*
>     Backprop. *Curves ← Minimize(Loss)*
> **end for**

---

images that fulfill the numerical objective but are unrecognizable to humans. This work specifically uses the `torch.transforms.RandomPerspective` and `torch.transforms.RandomResizedCrop` functions in sequence. Note that the specific details of the augmentation were not the focus of this work, thus a more robust augmentation function may exist and is left to future research.

Figure 2 showcases the gradual synthesis of a typical CLIPDraw drawing. Note that while the optimization process is largely deterministic, there is randomness in the initial curves and image augmentations, thus multiple runs of CLIPDraw can result in different drawings.
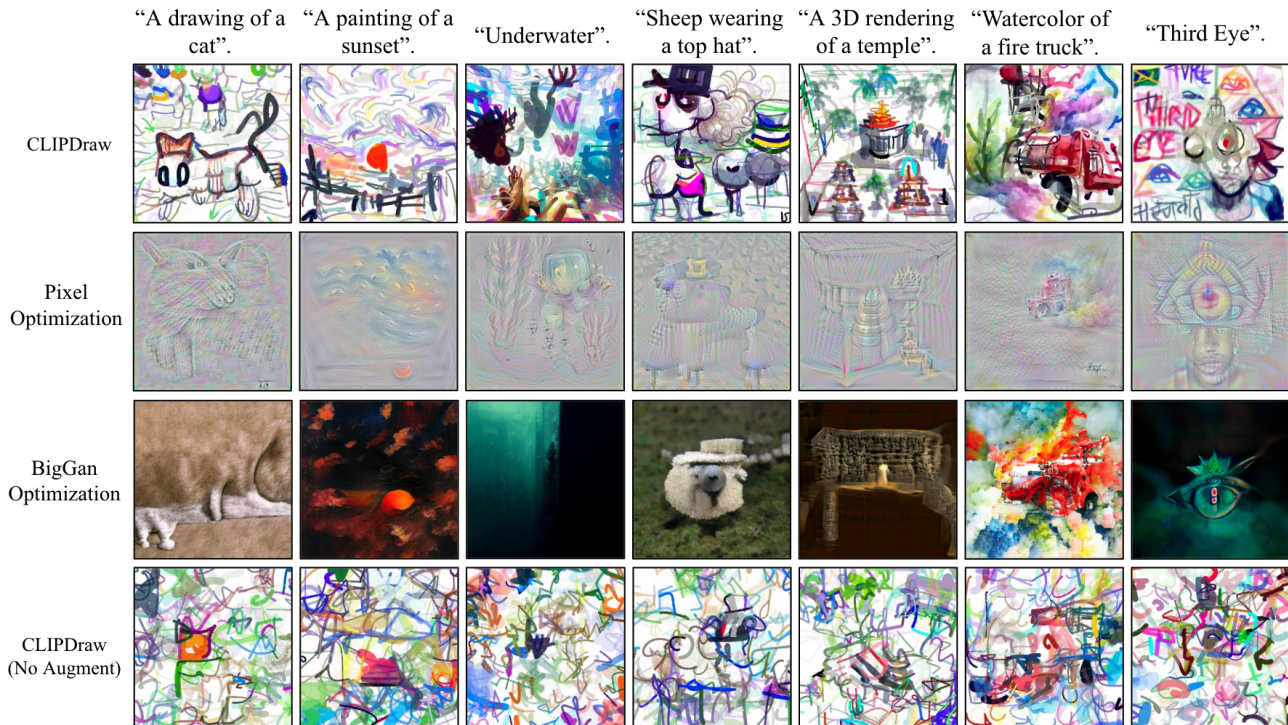
*Figure 3.* **Images synthesized via various synthesis-through-optimization methods**, with the goal of matching a given CLIP-encoded description phrase. CLIPDraw can produce a diverse set of human-recognizable drawings based on simple strokes and shapes.

## 4. Results

In the following sections, various interesting behaviors of CLIPDraw are highlighted through a variety of examples. With the exception of Section 4.1, example images are picked to best convey the behavior in consideration. Focus is placed on qualitative observations, unusual behavior, or recurring trends in CLIPDraw image synthesis.

### 4.1. How does CLIPDraw compare to other synthesis-through-optimization methods?

Compared to methods that learn a direct generative model, optimization-based synthesis methods do not require prior training. Instead, images are generated through an evaluation-time optimization loop, aiming to maximize a given objective. This work specifically focuses on synthesizing images that match the CLIP encoding of a description prompt. The following methods are compared:

- **CLIPDraw**, in which drawings are produced by a set of RGBA Bézier curves. The control points, thickness, and colors of the curves can all be adjusted.

- **Pixel Optimization**, which instead optimizes a 224x224x3 matrix of RGB pixel colors. Otherwise,

all algorithmic aspects are the same as CLIPDraw, including image augmentation.

- **BigGan Optimization**, in which images are produced using a pre-trained BigGAN generator. The weights of the generator are frozen; only the latent $Z$ vectors are optimized. Samples are generated using the method by Murdock (2021).

- **CLIPDraw (No Augment)**, which is identical to CLIPDraw, except no image augmentation is applied to the synthesized drawings.

In Figure 3, various methods are run on the same CLIP matching objective. Each method is run for 250 steps of gradient descent. In the CLIPDraw results, stroke count is 256, and 8 duplicates are used during image augmentation. CLIPDraw tends to result in a diverse set of human-recognizable drawings based on simple strokes and shapes. On the other hand, Pixel Optimization creates interesting textures but fails to compose colors and shapes. BigGan Optimization can synthesize high-resolution images, but is constrained to the set of images its generator can produce, thus it fails in out-of-distribution prompts. CLIPDraw without image augmentation produces images that score high numerically, but are nonsense when viewed by humans.
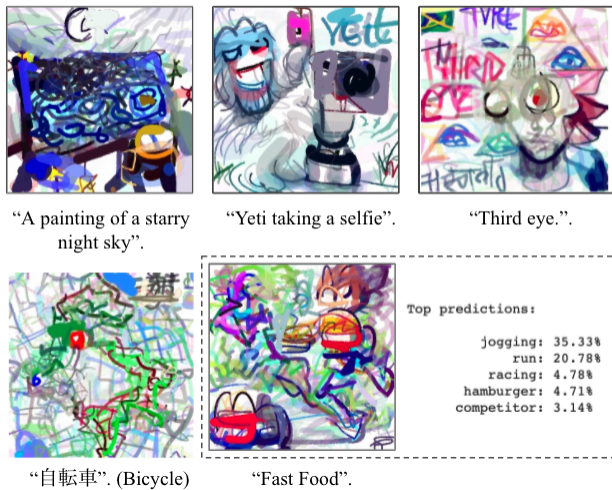
"A painting of a starry night sky".    "Yeti taking a selfie".    "Third eye.".

"自転車". (Bicycle)    "Fast Food".

Top predictions:

jogging: 35.33%
run: 20.78%
racing: 4.78%
hamburger: 4.71%
competitor: 3.14%

*Figure 4.* **CLIPDraw often matches the description prompt through a variety of creative techniques**, such as forming letters inside the images, or interpreting ambiguous words in multiple ways.

## 4.2. What kinds of visual techniques does CLIPDraw use to satisfy the textual description?

CLIPDraw often results in drawings that match their description prompts in multiple, unexpected ways, as shown in Figure 4. A prime example is the prompt for "a painting of a starry night sky". In this drawing, the main background features a sky with a prominent moon and a few scattered stars. The drawing itself is rendered in a painterly style, however the drawing also features an actual painting canvas and painter. Inside the canvas, black and blue swirls resemble Van Gogh's 1889 "The Starry Night".

Another interesting behavior of CLIPDraw is its tendency to write words in the drawing itself. In "Yeti taking a selfie", letters resembling "Yeti" can be seen in the top-right corner. In "Third Eye", again words resembling "third" and "eye" are scattered throughout the image.

At times, the drawings contain symbols that do not literally contain the description, but are tangentially associated, such as the prompt "自転車" (*bicycle* in Japanese) resembling a Google Maps screenshot with a Japanese-like character in the corner.

The ambiguity of prompts also presents intriguing results. In the prompt "Fast Food", a McDonald's logo along with a set of hamburgers is shown. However, also present are two joggers in a footrace, providing another interpretation of the phrase "fast". Included in Figure 4 are the top words predicted by CLIP as being closest to the image, showing that CLIP recognizes both "jogging" and "hamburger" as terms related to the synthesized drawing of "fast food".



"A drawing of a cat".    "A realistic photograph of a cat".    "A cat as 3D rendered in Unreal Engine".

"A Japanese woodblock print of one cat".    "A watercolor painting of a cat".    "A 3D wireframe model of a cat".
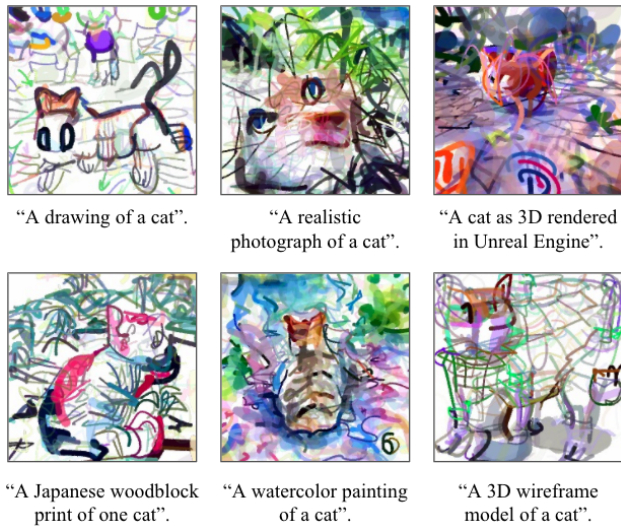
*Figure 5.* **By adjusting descriptive adjectives, CLIPDraw can produce drawings of vastly different styles.** Styles vary not only in the texture of the images, but showcase different representations of the underlying content, such as a cartoonish cat when prompted for a "drawing", versus a cat in perspective when prompted for a "3D wireframe model".

## 4.3. Can CLIPDraw reliably produce drawings in different styles?

A useful feature of CLIPDraw is its ability to adjust not just the content of its drawings, but also the styles, based on the description prompts given. Part of this flexibility is due to the robustness of curve-based images: in comparison to methods that use a pre-trained GAN generator, CLIPDraw drawings are not limited to the space of natural images. Thus, a variety of styles can be produced, and these styles are easily explorable through text.

As shown in Figure 5, a synthesized image of a cat can look vastly different depending on the descriptor words included. When asked for a "drawing of a cat", CLIPDraw synthesized a cartoonish depiction of a cat, comprised mostly of an outline and simple face. A "realistic photograph" features more detailed shading, while a "cat as 3D rendered in Unreal Engine" showcases complex lighting along with a depth-based blurring. Further styles feature a bias towards certain colors, such as the reds and greens of Japanese woodblock prints, or the multi-color blends of watercolors.

An interesting result is that adjusting descriptive adjectives not only changes the textures of the drawings, akin to Style Transfer methods (Gatys et al., 2015), but also changes its structural representation of the underlying content. For example, prompting for "a drawing" produces a flat cartoonish cat, while prompts like "a 3D wireframe" produce a cat in perspective, with depth and shadows.
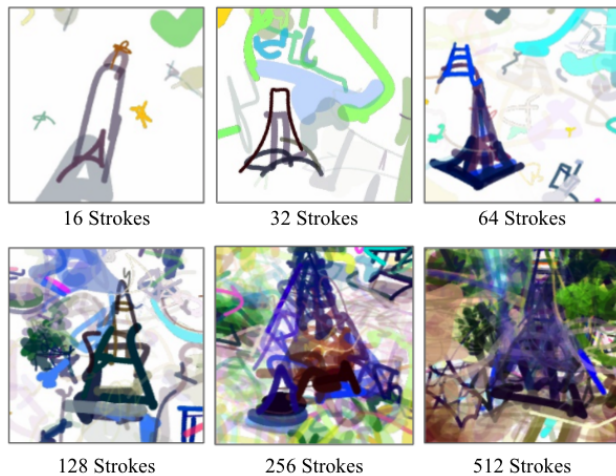
16 Strokes    32 Strokes    64 Strokes

128 Strokes    256 Strokes    512 Strokes



"Self".    "Happiness".    "Translation".

"The space between infinity".    "What do you look like, CLIPDraw?".    "Enlightenment".

*Figure 6.* **When stroke count is increased, CLIPDraw produces drawings of increasing realism.** Shown are multiple runs of the prompt "The Eiffel Tower" with a range of stroke counts. Low-stroke drawings opt for an abstract cartoonish representation, while high-stroke drawings capture 3D depth, background content, and complex shading.

*Figure 7.* **Abstract prompts grant insight into how CLIP relates visual concepts.** Synthesized images often contain symbols that indirectly relate to the description prompt through a cultural connection.

## 4.4. How does the stroke count affect what drawings CLIPDraw produces?

A crucial parameter of CLIPDraw is the number of strokes in each drawing. When the stroke count is low, CLIPDraw tends to produce cartoonish or abstract representations of the given description prompt. As stroke counts are increased, drawings become more detailed and incorporate additional features. Figure 6 showcases results for "The Eiffel Tower" on various stroke counts. In the 16-stroke example, the tower is drawn as only a few straight lines. The 32 and 64-stroke examples begin to display signs of 3D structure, such as a square base and a triangle-like scaffolding. Higher stroke count images begin more details on the Eiffel Tower itself, along with additional features such as background colors and complex lighting.

A common thread in synthesis-through-optimization methods is that pure optimization leads to undesirable results; it is also necessary to constrain optimization to a suitable space of images, such as the natural images generated by a GAN, or any image made of strokes in the case of CLIP-Draw. Limiting stroke count furthers this constraint. When optimizing within the space of 16-stroke images, it is hard to achieve details or textures, thus synthesized drawings will reveal the most basic forms that make up a visual concept. As a tool for AI-assisted art, the stroke-count parameter presents an easy way of adjusting between "simple" and "complex". This behavior may be useful in applications such as UI or icon design, where simplicity is important.

## 4.5. What happens if abstract words are given as a description prompt?

When given an abstract prompt without a literal interpretation, CLIPDraw must utilize cultural connections to come up with visual concepts that relate to the description. Often, this results in drawings that contain symbols relating to the given prompt, such as in Figure 7 with "Happiness" containing smiling faces and fireworks, "Translation" showcasing English and Japanese-like characters, and "Enlightenment" featuring a prominent monk-like figure.

At times, synthesized drawings demonstrate concepts through more complex relationships. In the prompt "Self", the resulting drawing features a body with multiple heads, evoking e.g. the idea that a person's self may contain multiple outward personalities. When asked "What do you look like, CLIPDraw?", the synthesized drawing contains a smiling face followed by text resembling "CLIPDRAW". Finally, "The space between infinity" presents a dream-like landscape with an infinity symbol under a galaxy-filled sky.

The ability for CLIPDraw to answer abstract prompts through related concepts presents a potential tool for exploration of human culture. As CLIP was trained on vast amounts of human data, cultural connections featured in synthesized drawings can grant insight into the typical connections humans may make. As a tool for AI-assisted art, this ability is also useful, as artists who want to evoke certain emotions will employ concepts that are culturally related to that emotion, akin to the smiles and fireworks in the drawing of "Happiness" above.

+1.0 "Hashtag".

+1.0 "A torii gate."

+1.0 "A realistic painting of a sailboat".

+1.0 "Hashtag".
-0.3 "Words and text".
-0.3 "Many ugly, messy, drawings".

+1.0 "A torii gate."
-0.3 "Purple".
-0.3 "Red".

+1.0 "A realistic painting of a sailboat".
-0.3 "A badly drawn sketch of many sailboats".
-0.3 "Many ugly, messy sailboats".

*Figure 8.* **CLIPDraw behavior can be adjusted through negative description prompts.** Negative prompts discourage synthesized drawings from matching with them, presenting a tool for fine-tuning results.

### 4.6. Can synthesized drawings be fine-tuned via additional negative prompts?

A common pain point in AI-assisted image synthesis is that is hard to control what the AI will produce. One potential solution in CLIP-based methods is to introduce negative prompts. In this setup, the optimization objective is to minimize cosine distance between the CLIP-encoded drawing and the description prompt, while maximizing distance between the drawing and a set of negative prompts.

Figure 8 showcases how negative prompts can fine-tune the behavior of CLIPDraw. Presented are pairs of drawings synthesized from the same random initialization, with the bottom row utilizing additional negative prompts, weighted on a 0.3:1 scale. In the "Hashtag" example, the original drawing contains many instances of the word "hashtag" written out. By penalizing "Words and text", the bottom example contains fewer words, instead featuring a set of faces typically seen in social media. A drawing of "a torii gate" originally results in a purple and red structure, however by penalizing "Purple" and "Red" the main color of the drawing switches to green. Lastly, the original drawing for "a realistic painting of a sailboat" features many sailboats on an ocean, and penalizing the phrase "many sailboats" results in a drawing featuring only a singular sailboat in the center.

In general, while negative prompts present a richly semantic way to fine-tune image synthesis in CLIP-based methods, it remains tricky to locate prompts that consistently encour-

age the intended behavior. Many times, negative prompts show negligible effect on the resulting drawing. During experiments, a cure-all negative prompt such as "a low-quality drawing", with the goal of consistently improving drawing quality, was unable to be found. Further work remains on how to best influence CLIP-based synthesis-through-optimization methods through additional objectives, whether negative or positive.

## 5. Discussion

This work presents CLIPDraw, a text-to-drawing synthesis method based on the CLIP language-image encoder. CLIP-Draw does not require any model training; rather, drawings are synthesized through iterative optimization during evaluation time. CLIPDraw is not the first method to utilize evaluation-time optimization for image synthesis; in fact, many recent works have used CLIP as an objective as well. However, by constraining image synthesis to images made up of RGBA Bézier curves, CLIPDraw biases towards simple drawings of human-recognizable concepts. The focus of this paper is to examine the nuances of CLIPDraw behavior, and experiments focus on specific questions and observations about synthesized drawings.

### 5.1. Limitations

The CLIPDraw method presented comes with various limitations. First, CLIPDraw inherently biases towards drawings rather than photorealistic images. Thus, synthesizing high-resolution images is a challenge, and CLIPDraw will often fall short of methods that incorporate a high-functioning generative model. This problem is related to a classic pitfall in synthesis-through-optimization methods, which is that the numerical objective used is not necessarily what is desired, e.g. an image may very closely match the CLIP objective, while looking messy and ugly to a human. Thus it is important to introduce auxiliary objectives or constraints. In the case of CLIPDraw, drawings are constrained by the Bézier curve representation, however stricter constraints such as fooling a GAN discriminator may improve synthesis quality.

A second limitation has to do with using CLIP encodings as a synthesis objective. While CLIP provides a rich textual representation for describing an image, in comparison to coarser neuron-activation objectives, it still remains a challenge to specify details. For example, it is hard to tell CLIPDraw to move a sailboat to the other side of the image. This work explored negative prompts as a possible direction towards more fine-grained adjustments, however a consistently satisfying method was hard to locate. A promising path in future research can lie in how to correctly adjust synthesized images, or introduce finer detail into description prompts via additional objectives.

## 5.2. Ethics and Social Biases

An important concept to keep in mind when dealing with human data is the existence of inherent social biases contained within. The pre-trained CLIP model is trained on a large corpus of online data, so its representations may include connections or biases that are undesirable. As CLIPDraw does not learn a new model, but instead optimizes based on CLIP itself, the bias studies presented in the CLIP paper (Radford et al., 2021) are highly relevant for CLIPDraw as well. This work does not specifically present additional social bias studies, however it is important to keep in mind that these biases exist when applying CLIPDraw to real-world use cases. As an example, in Section 4.5, a use case for CLIPDraw is mentioned as a tool for exploring visual connections in human culture. It is crucial to recognize that symbols or connections that are formed by CLIPDraw are not necessarily reflective of human culture, but rather are artifacts of the data used to train the original CLIP model. Thus, while CLIPDraw can be used to synthesize drawings that utilize cultural connections to evoke emotions or abstract concepts, it remains the duty of the user to ensure that the final product is up to desired standards. This caveat is especially key in automated setups, where AI-assisted art may produce unwanted results when run without a human in the loop.

## 5.3. Future Work

Overall, the aim of this work is to introduce CLIP-Draw as an easily accessible starting point to experiment with natural language image synthesis. Due to its focus on drawings rather than photorealistic rendering, CLIPDraw presents a straightforward method to examine language-image relationships without the overhead of realism. The presented CLIPDraw implementation can generally synthesize images within a minute on a typical GPU. Thus, interactable source code for experimenting the methods is made available at this Colab notebook: **https://colab.research.google.com/github/ kvfrans/clipdraw/blob/main/clipdraw.ipynb**

The results presented in this paper aim to describe various interesting behaviors of CLIPDraw, however they are by no means exhaustive. CLIP-based text-to-image synthesis remains a field with many promising directions, and we hope others will use this work as a backboard for additional research into the nuances of synthesis-through-optimization, or as a practical tool for AI-assisted art and other interactive visual applications.

## References

Chen, M., Radford, A., Child, R., Wu, J., Jun, H., Luan, D., and Sutskever, I. Generative pretraining from pixels. In *International Conference on Machine Learning*, pp. 1691–1703. PMLR, 2020.

Erhan, D., Bengio, Y., Courville, A., and Vincent, P. Visualizing higher-layer features of a deep network. *University of Montreal*, 1341(3):1, 2009.

Fernando, C., Eslami, S., Alayrac, J.-B., Mirowski, P., Banarse, D., and Osindero, S. Generative art using neural visual grammars and dual encoders. *arXiv preprint arXiv:2105.00162*, 2021.

Frolov, S., Hinz, T., Raue, F., Hees, J., and Dengel, A. Adversarial text-to-image synthesis: A review. *arXiv preprint arXiv:2101.09983*, 2021.

Galatolo, F. A., Cimino, M. G., and Vaglini, G. Generating images from caption and vice versa via clip-guided generative latent space search. *arXiv preprint arXiv:2102.01645*, 2021.

Gatys, L. A., Ecker, A. S., and Bethge, M. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576*, 2015.

Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial networks. *arXiv preprint arXiv:1406.2661*, 2014a.

Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014b.

Li, T.-M., Lukáč, M., Gharbi, M., and Ragan-Kelley, J. Differentiable vector graphics rasterization for editing and learning. *ACM Transactions on Graphics (TOG)*, 39 (6):1–15, 2020.

Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. Microsoft coco: Common objects in context. In *European conference on computer vision*, pp. 740–755. Springer, 2014.

Mirza, M. and Osindero, S. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.

Mordvintsev, A., Olah, C., and Tyka, M. Inceptionism: Going deeper into neural networks, 2015. URL https://research.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html.

Murdock, R. The big sleep: Bigganxclip, 2021. URL https://colab.research.google.com/drive/1NCceX2mbiKOSlAd_o7IU7nA9UskKN5WRl.

Nguyen, A., Yosinski, J., and Clune, J. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 427–436, 2015.

Nguyen, A., Dosovitskiy, A., Yosinski, J., Brox, T., and Clune, J. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. *arXiv preprint arXiv:1605.09304*, 2016.

Nguyen, A., Clune, J., Bengio, Y., Dosovitskiy, A., and Yosinski, J. Plug & play generative networks: Conditional iterative generation of images in latent space. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4467–4477, 2017.

Nilsback, M.-E. and Zisserman, A. Automated flower classification over a large number of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, pp. 722–729. IEEE, 2008.

Oord, A. v. d., Vinyals, O., and Kavukcuoglu, K. Neural discrete representation learning. *arXiv preprint arXiv:1711.00937*, 2017.

Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021.

Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., Chen, M., and Sutskever, I. Zero-shot text-to-image generation. *arXiv preprint arXiv:2102.12092*, 2021.

Reddy, P., Gharbi, M., Lukac, M., and Mitra, N. J. Im2vec: Synthesizing vector graphics without vector supervision. *arXiv preprint arXiv:2102.02798*, 2021.

Reed, S., Akata, Z., Yan, X., Logeswaran, L., Schiele, B., and Lee, H. Generative adversarial text to image synthesis. In *International Conference on Machine Learning*, pp. 1060–1069. PMLR, 2016.

Shen, I.-C. and Chen, B.-Y. Clipgen: A deep generative model for clipart vectorization and synthesis. *IEEE Transactions on Visualization and Computer Graphics*, pp. 1–1, 2021. doi: 10.1109/TVCG.2021.3084944.

Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. The caltech-ucsd birds-200-2011 dataset. 2011.