

# A Taxonomy of Prompt Modifiers for Text-To-Image Generation

JONAS OPPENLAENDER, University of Jyväskylä, Finland

Text-guided synthesis of images has become enormously popular and online communities dedicated to text-to-image generation and art generated with Artificial Intelligence (AI) have emerged. While deep generative models can synthesize high-quality images and artworks from simple descriptive text prompts, practitioners of text-to-image generation typically seek to control the generative model's output by adding short key phrases ("modifiers") to the prompt. This paper identifies six types of prompt modifiers used by practitioners in the online text-to-image community based on a 3-month ethnographic study. The novel taxonomy of prompt modifiers provides researchers a conceptual starting point for investigating the practice of text-to-image generation, but may also help practitioners of AI generated art improve their images. We further outline how prompt modifiers are applied in the practice of "prompt engineering." We discuss research opportunities of this novel creative practice in the field of Human-Computer Interaction (HCI). The paper concludes with a discussion of broader implications of prompt engineering from the perspective of Human-AI Interaction (HAI) in future applications beyond the use case of text-to-image generation and AI generated art.

Additional Key Words and Phrases: prompt engineering, text-to-image generation, human-AI interaction, AI generated art

## 1 INTRODUCTION

Text-to-image generation has become widely popular both in academia and as a new creative practice among practitioners of "AI art." Based on deep learning, text-to-image generation systems can generate digital images from short descriptive texts (called *prompts*, such as "*an oil painting of a beautiful landscape at dawn*"). To be effective, the textual input prompts need to be given in a certain format in order to, for instance, generate images with a certain style. This is commonly achieved by adding keywords and key phrases to the prompt (so-called "prompt modifiers"). Examples of images synthesized from textual prompts are depicted in Figure 1. Given the quality of these images, it is not surprising that an enthusiastic online community around this novel text-based way of creating images and art has developed. Within this community, the practice and skill of writing prompts is known by the term "prompt engineering" due to its iterative and experimental nature [28]. Prompt engineering is an emerging research area in the field of Human-Computer Interaction (HCI) concerned with how to phrase input prompts for deep generative models and – from a broader perspective – how humans can effectively interact with artificial intelligence.

The learning curve of prompt engineering can still be steep. Some prompt modifiers used within the community of practitioners are not intuitive and from looking at an image, it is impossible to tell the input prompt used to synthesize the image. On social media, many artists do not share their complete prompts for their artworks and it is often not clear how these artworks were created. Therefore, prompt engineering is a non-intuitive skill that is learned from extensive experimentation and trial and error [28]. A growing number of resources in the gray and scholarly literature present systematic experimentation on the effect of different prompt modifiers [12, 15, 28, 38, 54]. Online databases have been created in which users can explore artworks, prompts, and prompt modifiers [e.g. 1, 2, 27, 36, 57]. These resources and guides are part of a growing online ecosystem around text-to-image generation [37].

While guides, resources, and datasets about prompting are available, there is still a gap in our understanding of prompt modifiers. No previous study has investigated different types of prompt modifiers. With a specific focus on digital art generated with text-to-image systems, this paper contributes a taxonomy of prompt modifiers used by practitioners in the text-to-image community, based on an ethnographic study of the community's prompt engineering practices. The work is

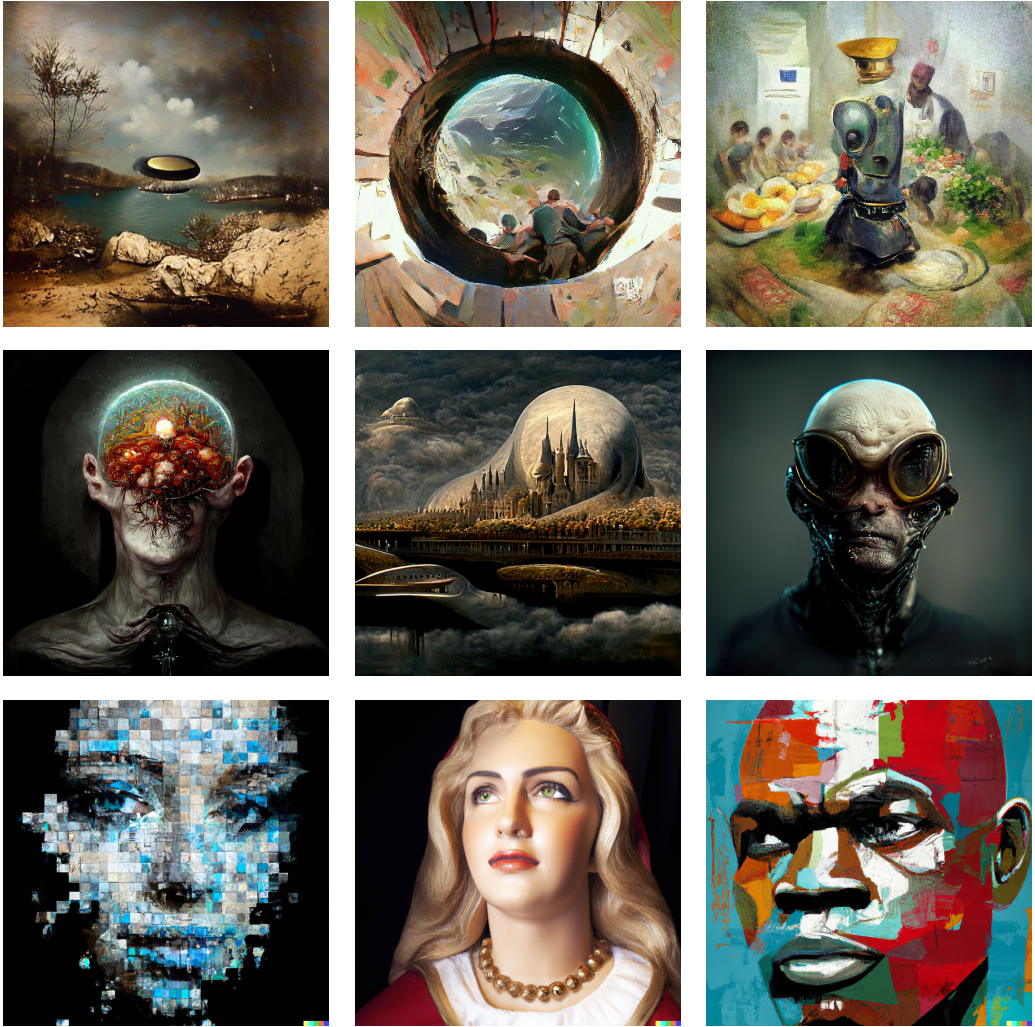


Fig. 1. Selected images generated with text-to-image generation using VQGAN-CLIP (top), Midjourney.com (middle), and DALL-E 2 (bottom).

based on a three-month online ethnography which analyzed how prompt modifiers are being applied in prompt writing. This paper contributes toward a better understanding of prompt engineering as a practice within HCI in order to inform the HCI research community on the emerging practice of prompt engineering within the broader context of human interactions with artificial intelligence. This paper aims to enhance the theoretical understanding of how people write prompts and use prompts modifiers. Through understanding prompt writing, we can pave the way towards a broader and unified theory of prompt engineering which the HCI literature is currently missing. The paper also touches on how the technology behind text-to-image systems and the practice of prompt engineering has broader implications in research on HCI and Human-centered AI (HCAI).

The paper is structured as follows. We first provide a brief introduction into text-to-image synthesis and prompt engineering in Section 2. After describing the methodological approach in

Section 3, a taxonomy of six different types of prompt modifiers is presented in Section 4. It is demonstrated how these prompt modifiers are applied in the context of prompt engineering in Section 5. The paper concludes with a discussion of opportunities for future research on text-to-image generation and the broader implications beyond AI generated art (sections 6 and 7).

## 2 BACKGROUND

This section discusses the evolution of text-to-image generation, particularly highlighting the role of OpenAI’s CLIP model. It then delves into the concepts of “prompt engineering,” a creative practice for controlling image generation, and “prompt modifiers,” keywords used to refine the image output. The section also underscores the contribution of the online community in advancing these creative practices.

### 2.1 Text-to-Image Generation

The field of image synthesis using deep learning has seen an unprecedented growth with the breakthrough development of multimodal models trained on large amounts pairs of images and text scraped from the World Wide Web. The development was initially spurred by OpenAI’s multimodal model CLIP [44]. CLIP is a contrastive language-vision model trained in an unsupervised way to perform zero-shot classification of images. CLIP provides a convenient way to transform both text and images into a common vector-based representation. When used as a discriminator component in text-conditioned generative systems, CLIP can “guide” the image generation process. CLIP was originally a part of OpenAI’s DALL-E architecture [46], a text-to-image system that was never released in its entirety. However, OpenAI did release the weights of the CLIP model. This resulted in a vast number of open source implementations of text-to-image systems, first as CLIP-guided generative adversarial networks (e.g., VQGAN-CLIP by Crowson et al. [9]) and later as diffusion based image generation systems, such as CLIP Guided Diffusion [8] and Latent Diffusion [49].

This paper investigates text-to-image generation from the lens of Human-Computer Interaction (HCI). In order to generate images from text, one not only has to choose the right words to make the text-to-image system generate the desired images, one also has to add different keywords and key phrases to control the style and quality of the image generation. This creative practice of writing effective prompts is sometimes referred to as “prompt engineering.” This paper investigates what (and how) different types of prompt modifiers are being applied in prompt engineering.

### 2.2 Prompt Engineering and Prompt Modifiers

Prompt engineering [28] – also referred to as “prompt design” [35], “prompt programming” [47], and “prompting” [2] for short – is the practice of writing textual inputs for generative systems. In the context of text-to-image generation, “carefully selected and composed sentences are used to achieve a certain visual style in the synthesized image” [50]. The practice has seen an ideal application ground in AI generated art, but it is not limited to text-to-image generation. The term prompt engineering was originally coined to denote the practice of writing textual inputs for the language model GPT-3 [28]. This autoregressive language model requires context to produce relevant text as output. Templates have been developed to optimally provide textual inputs to GPT-3. OpenAI’s documentation, for instance, lists 49 “recipes” on how to phrase input prompts for their language model.<sup>1</sup> Templating languages and interfaces have been developed to advance the field of prompting [2]. Using such recipes and tools, the output of the language model can be adapted to different down-stream tasks, such as correcting grammar, summarizing text, answering questions, generating product names, or acting as a chat bot.

<sup>1</sup>See <https://beta.openai.com/examples>.

Templates have also emerged for writing input prompts for text-to-image systems, particularly in the online community around AI generated art. For instance, the “Traveler’s Guide to the Latent Space” recommends the following prompt template [54]:

*[Medium][Subject][Artist(s)][Details][Image repository support]*

Similar templates are being followed in many resources originating from within the online community, such as the DALL-E Prompt Book [38]. Figure 2 provides an example of a typical textual input prompt and the resulting AI generated image.



Fig. 2. Digital artwork generated with DISCO Diffusion from the input prompt “A beautiful painting of a singular lighthouse, shining its light across a tumultuous sea of blood by greg rutkowski and thomas kinkade, Trending on artstation.” This prompt is part of the default configuration settings in the DISCO Diffusion notebook.<sup>2</sup>

Prompt engineering is not a hard science as found in the fields of science, technology, engineering, and mathematics (STEM). Rather, it is a term that originates from within the online community of practitioners of text-to-image generation. The term reflects the community’s self-understanding, similar to the terms “AI art” and “AI artist” which also originate from within the community. Due to the rise in popularity of text-to-image systems, practitioners of AI art include not only technology-savvy developers and early-adopting hobbyists, but also artists, professionals, semi-professionals, and “Pro-Ams” [23] with or without commercial interests. In the remainder of this paper, we will refer to the members of the online text-to-image community as *practitioners*.

Prompt engineering resembles a conversation with the text-to-image system. A practitioner typically will run a prompt, observe the outcome, and adapt the prompt to improve the outcome. Prompt engineering, thus, is iterative and practitioners formulate prompts as probes into the generative models’ latent space. The online community quickly found that the aesthetic qualities and subjective attractiveness of images can be improved by adding certain keywords and key phrases to the textual input prompts. The terms may be referred to by a number of different names, such as “style phrases,” “clarifying keywords” [39], or “vitamin phrases” [42]. In this paper, we refer

<sup>2</sup>See <https://github.com/alembics/disco-diffusion>.



to them as *prompt modifiers*. By adding a prompt modifier to a textual input, one seeks to direct the text-to-image system in certain directions, hence “modifying” the resulting image.

In practice, prompt modifiers are applied through experimentation or based on best practices learned from experience or online resources. An example of an iterative application of prompt modifiers can be seen in Figure 3. Knowing what prompt modifiers work best for a given subject term is often the result of the practitioner’s iterative experimentation, research in online communities, and the use of online tools and resources created for supporting the practice of prompt engineering [37].



Text prompts:

- a) “*ufo landing*”
- b) “*ufo landing, daguerreotype*”
- c) “*ufo landing, daguerreotype, trending on /r/art*”
- d) “*ufo landing, daguerreotype, by greg rutkowski, trending on /r/art*”

Fig. 3. Example of iterative prompt engineering for generating an image. Images generated with VQGAN-CLIP by Crowson et al. [9] with 175 iterations, CLIP model ViT-B/32, VQGAN model wikiart\_16384, and seed 6087304447281500163.

### 3 METHOD

In this research, a dual-methodological approach was adopted, leveraging both autoethnographic (Section 3.1) and online ethnographic studies (Section 3.2), to delve into the nuanced aspects of prompt engineering and text-to-image art generation. Understanding the intricacies of prompt engineering, an acquired skill cultivated through iterative experimentation, necessitates a hands-on, experiential approach. Hence, autoethnography provided a fitting method to gain an intimate, practitioner’s perspective. By conducting an autoethnographic study, the author was able to engage with the process of text-to-image synthesis personally, thereby capturing its nuances from a first-hand perspective. However, the complex and communal nature of this emerging field necessitated a broader perspective. To capture the collective wisdom and shared practices within the field, the author further complemented the autoethnographic approach with an online ethnography of the text-to-image art community. This approach allowed the author to glean insights from the shared experiences and resources of the broader community of practitioners active in online spaces, primarily on Twitter. The synthesis of these two complementary approaches aimed to provide a comprehensive understanding of prompt engineering, bridging the gap between individual experience and collective knowledge.

#### 3.1 Autoethnographic Research on Prompt Engineering

Prompt engineering is learned through iterative experimentation akin to “brute-force trial and error” [28]. Therefore, prompt engineering is an acquired skill that is associated with a learning

curve. The skill can be learned from community-provided resources, such as written guides and reports of systematic experimentation, or from prompts shared on social media, such as online communities dedicated to text-to-image art [37]. However, to appreciate and understand the craft of prompt engineering and text-to-image generation, one has to apply the knowledge and experiment with different input prompts. Autoethnography research is, therefore, an appropriate method to learn about prompt engineering.

The author conducted a 3-month autoethnographic study [10, 11, 13] between October 2021 and December 2021. This personal ethnography [7] allowed the author to get a “practitioner’s perspective” [11] of text-to-image generation by “learning from self-use” [33]. The author experimented with text-to-image synthesis and created digital images with a text-to-image system using notebooks hosted on Google’s Colaboratory (Colab).<sup>3</sup> The author started on average at least one Colab session every work day between October 4 and December 31, 2021. The free tier of Google Colab was used in all sessions. This limited the overall working time to about 2 hours per day, depending on the computational power of the assigned resources and whether penalties were incurred the previous day. VQGAN–CLIP by Crowson et al. [9] was chosen as text-to-image system using a notebook titled “VQGAN and CLIP (z + quantize method with augmentations)”.<sup>4</sup> This VQGAN–CLIP notebook was originally created by Katherine Crowson, with “modifications by Eleiber # 8347” and a “friendly interface” by “Abulafia # 3734” and further modifications by Justin John. VQGAN–CLIP was selected for several reasons. First, VQGAN–CLIP was one of the first text-to-image systems that experienced widespread popularity in the emerging text-to-image art community in 2021. This made VQGAN–CLIP instrumental to the growth of the community [9]. Second, the system can be executed on Google’s Colaboratory (Colab) free of charge. The system requires less memory than later systems, and it is therefore less likely that image generation will fail due to insufficient memory. Third, the VQGAN–CLIP notebook on Colab is very accessible and straight-forward to use, with only a small number of configuration parameters (cf. Figure 3). Last, the system is deterministic. Consecutive runs with the same configuration parameters will produce exactly the same images which makes the images reproducible. This is not the case with some of the later systems which make use of non-deterministic algorithms. The author generated 885 images in the course of this study.

The autoethnographic research was not conducted from scratch. Rather, it was informed by learning from the community on social media. To this end, the autoethnographic research was complemented with an online ethnography of the text-to-image art community on Twitter and a study of online community resources, described in the following section.

### 3.2 Ethnographic Study of the Text-to-Image Art Community

An ethnographic study of prompt engineering was conducted on Twitter (see Section 3.2.1). The aim of this social media ethnography [40, 41] was to learn more about the textual prompts used in the community of practitioners of text-to-image art. Insights derived from the study of this community were used in the autoethnographic experimentation with the text-to-image system. The research was complemented with a review of the literature (Section 3.2.2).

*3.2.1 Twitter community.* A dedicated online community around text-to-image generation with specific focus on AI generated art – which practitioners in the online community sometimes refer to as “AI art” [30] – has emerged. Social media services, such as Twitter, are a well-suited outlet for practitioners in this community to post and share images and experiences.

<sup>3</sup><https://colab.research.google.com>

<sup>4</sup>[https://colab.research.google.com/github/justinjohn0306/VQGAN-CLIP/blob/main/VQGAN%2BCLIP\(Updated\).ipynb](https://colab.research.google.com/github/justinjohn0306/VQGAN-CLIP/blob/main/VQGAN%2BCLIP(Updated).ipynb)

During the 3-month period of research, the author took the role of “participant-as-observer” [20] by engaging with the text-to-image art community on Twitter, participating in discussions, and posting images created with the text-to-image system. The author followed posts on Twitter to learn about different prompts used in the text-to-image art community. To this end, the author followed trending hashtags, such as #vqganclip, #VQGAN, #clipguideddiffusion, #digitalart, #AIArt, and #generativeart. Not every practitioner of text-to-image art shares their prompts on Twitter. Especially if commercial interests are involved – e.g., selling the art as non-fungible tokens (NFTs) – practitioners may keep their prompts a secret. The research material, therefore, was sparse at the time of conducting the study. However, some practitioners are more liberal in sharing their prompts. It is the posts from this group of Twitter users that informed this research (e.g., posts by Katherine Crowson (@RiversHaveWings), Hannah Johnston (@hannahjdotca), @nshepperd1, and John David Pressman (@jd\_pressman), to name but a few).

*3.2.2 Review of community resources.* In parallel to the research on the online community, a review of the literature was conducted, with specific focus on text-to-image generation and the practice of prompt engineering for digital art. With the exception of Liu and Chilton’s design guidelines for prompt engineering [28, 43], there still is little scholarly literature on the practice of text-to-image generation for AI generated art in the field of HCI. Therefore, the literature review primarily focused on sources in the gray literature, such as community-provided resources, documents, guides, experiment reports, blog posts, articles on the Web.

### **3.3 Inductive Development of the Taxonomy**

The taxonomy was developed inductively from pieces of information found during the research. Due to the relative scarcity of this material at the time of writing, the development of the taxonomy was conducted iteratively, as follows. A list of potential candidates for prompt modifiers was inductively compiled and grouped. This list was subject to continual reinterpretation when novel instances of prompts were encountered. Whenever a candidate for a novel type of prompt modifier was found in a post on Twitter or the literature, the author revisited the list of prompt modifiers. Therefore, the resulting taxonomy was iteratively and inductively revised and expanded when new types of prompt modifiers were encountered. After some weeks of collecting data this way, the list of prompt modifiers and taxonomy did no longer grow, even if instances of novel and atypical prompts were encountered. This indicates the completeness of the developed taxonomy.

The findings were documented in a PowerPoint presentation with text and images to produce an evocative and aesthetic description of the ethnographic research. This iteration also served as verification of the correctness of the taxonomy. The author’s creation of and engagement with the presentation acted as a daily conversation with the research material. This allowed the author to concurrently and iteratively develop and articulate an understanding of the subject matter both visually and textually. At the end of the research period, the author engaged in a summative analysis [11] of the research material to review the completeness and consistency of the taxonomy.

### **3.4 Self-Disclosure**

While the author has experimented with text-to-image systems and produced digital artworks with these systems, the author is not an artist. The author’s background is in Computer Science with focus on Human-Computer Interaction (HCI) and Social Computing. The research was conducted not from a technical lens, but a human-centered lens [21]. The author’s specific interest in prompt engineering is the text-based interactions of users with text-to-image systems and the novel creative practices that arise from these systems.

#### 4 TAXONOMY OF PROMPT MODIFIERS

This research points towards there being six different types of prompt modifiers (subject terms, image prompts, style modifiers, quality boosters, repeating terms, and magic terms) used by practitioners in the text-to-image art community (see summarized in Table 1). This taxonomy reflects the practitioner’s comprehension of prompt modifiers, a knowledge that was instrumental in classifying these modifiers into six distinct categories.

Table 1. Taxonomy of prompt modifiers.

Modifier	Description
Subject term	Denotes the subject
Style modifier	Indicates an artistic style
Image prompt	Indicates a style or subject via an image
Quality booster	A term intended to improve the quality of the image
Repeating term	Repetition of subject terms or style terms with the intention of strengthening this subject or style
Magic term	A term that is semantically different from the rest of the prompt with the intention to produce surprising results

**Subject terms** indicate the desired subject to the text-to-image system (e.g., “*a landscape*” or “*an old car in a meadow*”). While it is possible to generate images without subject terms, the subject is essential for controlling the image generation process. On the other hand, since text-to-image systems were trained on images in context of their descriptive text, subject terms can, in some cases, have less control over the outcome. One such case is the artist Zdzisław Beksiński who developed a unique and recognizable style but never provided titles for his artworks. For this reason, early text-to-image systems, such as VQGAN–CLIP, struggled to reliably reproduce specific subjects in images generated to resemble Beksiński’s artworks.

**Style modifiers** can be added to a prompt to produce images in a certain style. Style modifiers will consistently reproduce a characteristic style (e.g., “oil painting”) or artistic medium (e.g., “mixed media”). For instance, the modifier “*by Francisco Goya*” will generate digital images in the recognizable style of the late Spanish painter. Other examples of this type of modifier include, but are not limited to, “*oil on canvas*,” “*#pixelart*,” “*hyperrealistic*,” “*abstract painting*,” “*surreal*,” “*Cubism*” or “*cubist*,” “*cabinet card*,” “*in the style of a cartoon*,” “*by Claude Lorrain*,” and “*in the style of Hudson River School*,” to name but a few. As can be seen from the above list, style modifiers can include information about art periods, schools, and styles, but also art materials and media, techniques, and artists. When it comes to the latter, modifiers such as “*by Greg Rutkowski*” and “*by James Gurney*” have become popular in the community of text-to-image art as a means to produce images in a certain style and quality.

**Image prompts** act similar to subject terms and style modifiers in that they provide the text-to-image system a (visual) target for the synthesis of the image (both in terms of style and subject). Image prompts are typically specified as one or several urls that are added to the textual input prompt or provided in a separate array. Image prompts are different from “initial images” which were investigated by Qiao et al. [43]. Whereas an image prompt can consist of multiple images, there can only be one initial image. This initial image can be specified as a starting point for the image generation, for instance, for the purpose of enhancing or distorting the initial image. This is made possible because of the iterative nature of the image generation process which typically starts with an image filled with random noise (such as Perlin noise).



**Quality boosters** can be added to a prompt to increase aesthetic qualities and the level of detail in images. Examples of this type of modifier are the terms “*trending on artstation*,” “*award-winning*,” “*masterpiece*,” “*highly detailed*,” “*awesome*,” “*#wow*,” “*epic*,” and “*rendered in Unreal Engine*.” This type of modifier can also be applied in the form of “extra fluff” added to the prompt. Verbosity in the prompt may boost the amount of details and overall quality of the generated image, at the expense of the subject becoming less controllable. For instance, the prompt “*painting of an exploding heart*” could potentially be improved by appending the modifiers “*highly detailed, eclectic, fiery, vfx, rendered in octane, postprocessing, 8k*.”

**Repeating terms** can strengthen the associations formed by the generative system. For instance, the prompt “*space whale. a whale in space*”<sup>5</sup> by @nshepperd1 will likely produce subjectively better results than either of the two subject terms alone. The use of different phrasing and synonyms will cause the text-to-image system to more reliably activate regions in the neural network’s latent space that are associated with the subject terms. This is not only an imagined effect. The prompt “*a very very very very very beautiful landscape*” will, for instance, likely produce a better image than a prompt without repeating terms. Technically, this is due to likelihood-maximizing language models becoming stuck in positive feedback loops from repeated phrases [24].

**Magic terms** introduce randomness to the image that can lead to surprising results. For instance, Twitter user @jd\_pressman added the magic term “*control the soul*” to the prompt “*orchestra conductor leading a chorus of sound wave audio waveforms swirling around him on the orchestral stage*”.<sup>6</sup> The term was added to – in Pressman’s words – produce “more magic, more wizard-ish imagery”.<sup>7</sup> Magic terms, thus, introduce an element of unpredictability and surprise to the resulting images, often with the intention of increasing the variation in the output. Magic terms can refer to terms that are semantically distant to the main subject of the prompt, or they can refer to non-visual qualities, such as the sense of touch (somatosensory), sense of hearing (auditory), sense of smell (olfactory), and sense of taste (gustatory) (e.g., “*feed the soul*” and “*feel the sound*”).

In summary, prompt modifiers come in a variety of types and can take different forms. They can, for instance, be added as hash tags (e.g., “*#wow*”), attribution phrases (e.g., “*by [artist]*”), or more complex composite statements (e.g., “*in the style of [artist]*”). Further, not every part of a prompt has the same importance and there are specific affordances of text-to-image systems that are being used in the practice of prompt engineering, as described in the following section.

## 5 PROMPT ENGINEERING IN PRACTICE

This section provides an overview of how the different types of prompt modifiers are being applied in the practice of prompt engineering with specific focus on the generation of static images from either textual or visual input prompts. We specifically focus on demonstrating and explaining the iterative process of text-image generation with its iterative different steps (as described in Table 1).

The first step in iterative prompt design is to denote the **subject** with one or several terms. While images can be generated from random text or even single characters and emojis [37], the subject term is fundamental to the controlled generation of digital images. Consequently, a prompt typically contains at least one subject term. Any other parts of the prompt are optional. It is, for instance, possible to generate artworks with the prompt “*car*.” In practice, however, practitioners use modifiers to improve the resulting images and to exercise more control over the image creation process.

<sup>5</sup><https://twitter.com/nshepperd1/status/1456584388037148678>

<sup>6</sup>[https://twitter.com/jd\\_pressman/status/1457171648293924867](https://twitter.com/jd_pressman/status/1457171648293924867)

<sup>7</sup>[https://twitter.com/jd\\_pressman/status/1457445367125921793](https://twitter.com/jd_pressman/status/1457445367125921793)

**Modifiers** are typically added with the intention to either modify the style of the generated image or boost its quality. As mentioned in Section 4, style modifiers and quality boosters do not form a disparate set. Rather, the two types of modifiers can have overlapping effects and the difference between the two types of prompt modifiers is sometimes not fully apparent. For instance, the modifier “*by Greg Rutkowski*” exhibits this property. Greg Rutkowski<sup>8</sup> is a contemporary illustrator and concept artist who has been embraced by the text-to-image art community in their practice of prompt engineering. Images generated with the modifiers “*by greg rutkowski*” or “*in the style of greg rutkowski*” are of high quality, texture-rich, and contain a high amount of details. As such, this modifier is often not used as a style modifier – as one would expect –, but as a quality booster in the community, even though a trained eye may tell by the style of the image that the prompt modifier was being used.

Once a style modifier has been added, the style can be reinforced and “solidified” without losing expressivity. **Solidifiers** (in the form of repeating terms) can be applied to any of the other types of modifiers (subjects, style modifiers, and quality boosters), although they are most commonly applied to subject terms. Image prompts are a special case in that they can carry both information about the subject and style because of their visual nature. If the textual prompt is aligned with the image prompt, the image prompt can also act as a solidifier. On the other hand, if several images that are different from each other are added to the prompt, the image prompts will contribute to variation in the output. Last, **magic terms** may be optionally added to increase the chance of surprising results. The use of magic terms will result in more variation in the output, while maintaining the overall style.

Each of the six types of prompt modifiers can be assigned **weights**. Weighted terms can be negative to exclude subjects and styles from being generated. For instance, VQGAN–CLIP tends to generate heart-shaped objects when the prompt contains the word “*love*.” By adding a negative weight to the prompt (e.g., “*heart:-1*”), the system can be instructed not to activate the corresponding latents in its neural network. The resulting image is thus free from heart-shaped objects. Weighted terms can also be used to seamlessly mix styles. For instance, Twitter user @c0y0te6 mixed the styles of two artists in the prompt “*a painting of a high prestess [sic] summoning a demon by Ralph McQuarrie:75 | by Zdzislaw Beksiński:25*”.<sup>9</sup> The style of Ralph McQuarrie is, in this case, given precedence over the style of Zdzisław Beksiński (with a ratio of 3:1).

Table 2 summarizes the iterative nature of prompt writing (c.f. Figure 3). Subject terms are most important for the controlled generation of images and usually written as first step. Modifiers and solidifiers are then added to the prompt, either iteratively (image after image) or from learned experience. Last, weights can be applied to exclude or mix subjects and styles.

Table 2. The iterative practice of prompt writing.

Step	Purpose	Prompt modifier	Importance
1	Define	subject term, initial images, image prompt	required
2	Modify	style modifier, quality booster, initial images, image prompt	optional
3	Solidify	repeating terms, initial images	optional
4	Vary	magic terms, initial images	optional
5	Mix/Exclude	mixing and exclusion	optional

<sup>8</sup><https://www.artstation.com/rutkowski>

<sup>9</sup><https://twitter.com/c0y0te6/status/1481780797858275329>

## 6 DISCUSSION

The availability and accessibility of text-to-image generation as a new creative practice and artistic medium [29], paired with a specific bundle of technologies and resources that support the ecosystem of this “emerging art scene” [37, 55], have resulted in an explosion of AI generated artworks being shared online. The application of prompt modifiers is key to the emerging creative practice called prompt engineering. The taxonomy of six different types of prompt modifiers represents an initial work to bringing structure to the creation process and research in text-to-image systems. The taxonomy of prompt modifiers is reified for the sparse HCI literature around prompt engineering as a logical building block in this emerging field of research.

Midjourney has over 15 million members at the time of writing,<sup>10</sup> and open source systems, such as Stable Diffusion, are available for execution on cloud or local hardware. Today, everyone is able to synthesize digital images and artworks from natural language using free or relatively inexpensive means, with implications for productivity and creativity [37]. Gartner estimated in 2021 that by the year 2024, 80% of technology products and services will be built by people who are not technology professionals [18]. Increasingly, deep generative models will be used by laypeople without technical expertise and skills. Interaction with opaque deep learning models will increasingly become more common in future use cases and applications of artificial intelligence. Therefore, prompt engineering is an emerging and important research area in the field of Human-Computer Interaction (HCI). However, with the exception of the design guidelines by Liu and Chilton [28] and Qiao et al. [43], the scholarly literature in the field of HCI on prompt engineering still resembles a cottage industry, with concepts and structures yet to emerge. Meanwhile, many resources started to emerge from within the online community, such as Smith’s “Traveler’s Guide to the Latent Space” [54] and Parsons’s “DALL-E Prompt Book” [38]. Drawing on gray literature, such as the above, and extensive auto-ethnographic research, this work provides a taxonomy of prompt modifiers as a starting point for systematizing the practice of prompt engineering for text-to-image generation. The subsequent discussion will examine the broader implications of prompt engineering for human-AI interaction.

### 6.1 Broader Implications for Human-AI Interaction

Research on prompt engineering has broader implications and is not only limited to the field of text-to-image synthesis and AI generated art, but also relevant to the interaction of humans with deep learning models and artificial intelligence in general.

*6.1.1 AI and the future of creative work.* There is much potential for deep learning to disrupt and transform entire sectors of the creative economy. Recently, there has been an interest into developing generative systems that are able to synthesize more complex outcomes. For instance, systems for text-to-video generation have been presented by Hong et al. [25], Ho et al. [22], Singer et al. [53], and Villegas et al. [56]. Low-code and no-code tools for creating online products and experiences will become increasingly common in the future. Declarative machine learning systems may – as a next wave of machine learning – bring machine learning to non-coders [32]. This technology will extend the currently rather narrow focus of prompt engineering on language models and text-to-image synthesis to more broader application domains. In the future, we may see deep generative models with generative capabilities that transcend what we can imagine today. Deep generative models could, for instance, create entire interactive story-driven worlds and games from short text prompts.

Such powerful AI-based systems will have implications for the future of creative work. Artificial intelligence will not only transform the way we interact with computers and perform work online,

---

<sup>10</sup>See <https://discord.com/servers>.

but also the content of our work and the human agency in the work. An example of an application that has such transformative potential is OpenAI’s Codex [5, 59]. Codex is a large language model that interprets commands in natural language and generates programming code. In the future, instead of typing code, we will be able to describe a software and its expected outputs in natural language. Pre-trained generative models, such as Codex, BLOOM [3], or other “foundation models” [4], will then generate executable software code based on the human’s spoken or written input prompts. This technology has already found application in GitHub’s CoPilot<sup>11</sup>, an “AI pair programmer” assisting its users in auto-completing programming code. In academia, researchers increasingly rely on language models as creativity support tools for writing academic papers [26]. The change in the agency of humans and computers brought by generative models will be transformative to creative work, such as software development and research.

*6.1.2 Beyond text-to-image generation.* The use case of art generated with text-to-image systems discussed in this paper is but one of many application areas of prompt engineering, with implications for the future of creative work and Human-AI Interaction (HAI) in general. The latter can be viewed from many different perspectives, such as human-centered AI [52], human-AI partnerships [45], and human-AI cooperation [6]. Irrespective of the term used to describe our relationship with AI, we will increasingly interact with opaque models through prompts in natural language.

Research on how to design prompts is therefore timely and important. Increasingly, we see use-facing applications being powered by foundation-scale models. The emergent properties of these models make it possible to use them for a vast number of different use cases and applications. Internally, such applications are often enabled by prompt engineering. For instance, tool-augmented language models [31, 48] internally use prompts to enable the language model to use external tools. Research on prompt engineering, thus, will advance our understanding of how people can effectively interact with and employ machine learning models for solving complex tasks.

With these considerations in mind, we turn our attention to specific opportunities and challenges of prompt engineering within the field of Human-Computer Interaction (HCI).

## 6.2 Opportunities for Research on Prompt Engineering in HCI

This section discusses opportunities for future research on prompt engineering in the field of HCI. Specifically, we touch on the community dynamics surrounding text-to-image art creation, novel workflows and techniques employed by practitioners, and the embedded biases in AI-driven systems. Additionally, we explore the relevance of prompt engineering for research on computational aesthetics and human-AI alignment.

*6.2.1 Social aspects of prompt engineering.* There are social components to the use of text-to-image generation systems. Prompt engineers face an interesting challenge: Because text-to-image systems were trained on images and text scraped from the Web, users of text-to-image systems need to imagine and predict how other people described and reacted to images posted on the Web. Describing an image in detail is often not enough to achieve optimal results – one has to imagine the image as if it already existed on the Web.

Another social aspect in prompt engineering are the dedicated communities that came into existence only recently. Practitioners of text-to-image art are producing artworks in shared Discord-based chat rooms, such as on Midjourney.<sup>12</sup> These dedicated communities offer a rich set of social features worth investigating more closely in HCI research. For instance, members on Midjourney have their own profiles that bundle the members’ successful creations together with the prompts

<sup>11</sup>[copilot.github.com](https://copilot.github.com)

<sup>12</sup><https://www.midjourney.com>

used to create the images. Midjourney introduced a 2D map in which members can explore other members based on the similarity of their prompts. Midjourney also has dedicated “group jam” sections in which members can iterate on and further develop other members’ works and there is a “theme of the day” section. Long running threads are quite common in this community. Community-learning is an interesting area of research in this regard. How do members receive and seek inspiration in the community? How do novices learn the craft of prompt engineering and is there learning taking place in the community as a whole?

Future work could explore and ethnographically investigate the online community around text-to-image art and its prompt engineering practices in more detail, using the taxonomy presented in this paper as a conceptual starting point or framework.

**6.2.2 Human-AI co-creation.** While the heart piece of prompt engineering is prompt writing, prompt engineering is only a starting point in some practitioners’ creative work flows. Novel creative practices are emerging. For instance, practitioners may develop complex work flows for creating their artworks (e.g., generating initial images with one text-to-image system as a source for inspiration, then continuing on another text-to-image system before finalizing the images in a photo editor). The different affordances of text-to-image systems still need to be reified and systematized in the HCI community. For instance, some text-to-image systems enable the creation of zooming animations, others can complete parts of images which is called image inpainting [60] and outpainting.<sup>13</sup> These novel creative practices offer a level of interactivity beyond mere generation of static images from textual input prompts. Further, practitioners may make certain idiosyncratic choices when they create text-based generative art (e.g., selecting certain numerical values as seed for the model or adapting the canvas size to certain subject terms). Some of these choices may fall into the realm of folk theories [14, 19] – that is, causal attributions that may or may not be true –, while other choices may be based on the practitioner’s experimentation and experience with prompt engineering. Future work could investigate these creative practices, work flows, strategies, and beliefs adopted by practitioners in the text-to-image art community. The emerging research field also offers an opportunity for HCI researchers to make technical contributions [58] in the form of creativity support tools, user interfaces, and interactive experiences to support text-to-image generation, to teach novices the practice of prompt engineering, and to advance the emerging AI generated art ecosystem. Research in this space could make a timely contribution to a novel computational medium and an emerging digital art form.

**6.2.3 Bias in image generation systems.** Another interesting area for future work is bias encoded in text-to-image generation systems. It has been shown, for instance, that the CLIP model contains bias<sup>14</sup> and some text-to-image systems prompted with “*princess*” will produce images of women with light skin color, reflecting the bias in the training data toward Western, educated, industrialized, rich and democratic (WEIRD) subjects.<sup>15</sup> OpenAI recently announced that bias was reduced in their DALL-E 2 model [34], but at the cost of potentially reducing signal-to-noise of the generated images.<sup>16</sup>

Responsible deployment of large models and the potential risks are two concerns often listed for not fully releasing a model. While organizations such as OpenAI and Google can be commended for trying to be responsible with their powerful systems, these organizations act paternalistic and impose their value and belief system onto their users which is another source of bias. DALL-E 2, in particular, can be a source of frustration for its users who are often faced with content policy

<sup>13</sup>See, for instance, <https://twitter.com/adampickard/status/1551584412659335168>.

<sup>14</sup>See <https://twitter.com/RiversHaveWings/status/1432100170645180416>.

<sup>15</sup>See <https://twitter.com/EMostaque/status/1495323912951021568>.

<sup>16</sup>See <https://twitter.com/minimaxir/status/1549070583035416576>.

notices for terms relating to war or sexual content (with a threat of account closure if the warning is incurred too often). Pressman et al. recently raised an important point: Humans are sexual beings and the androgynous values imposed on text-to-image systems with the intent of making them “safe-for-work” deprives users of “a key component of human aesthetic values and experience” [42].

*6.2.4 Computational aesthetics and Human-AI alignment.* The goal of making computers evaluate and understand aesthetics is much older than text-to-image generation [17]. Recently, there is renewed research on neural image assessment and computational aesthetics. State-of-the-art text-to-image systems increasingly consider human aesthetics in an attempt to produce better images [51]. Prompts are a vast resource for research on computational aesthetics, as they encapsulate a person’s stated intent. This intent, however, is likely only partially explicit. Research on prompt engineering, therefore, also relates to research on human-AI alignment [16]. This research area is concerned with teaching artificial intelligence to understand human values. Prompts for text-to-image generation systems could form an interesting study resource for this kind of research.

## 7 CONCLUSION

This research contributes to the academic understanding of text-to-image generation by proposing a novel taxonomy of six types of prompt modifiers: subject terms, image prompts, style modifiers, quality boosters, repeating terms, and magic terms. The taxonomy of prompt modifiers lays the foundation for future structured investigations into prompt engineering for text-to-image generation and AI generated art. Moreover, the taxonomy highlights the unique affordances of text-to-image systems, providing a clearer understanding of the design (“engineering”) of prompts for image generation.

In the practice of prompt engineering for generating static images from textual or visual inputs, subject terms are fundamental to the controlled creation of images. Practitioners often use prompt modifiers to improve image quality and exercise greater control over the creation process. Modifiers either modify the image style or enhance its quality, and these two types can overlap in their effects. For example, the modifier “by Greg Rutkowski” is typically used by practitioners as a quality booster rather than a style modifier, despite the artist’s distinct style. Solidifiers can also be used to reinforce a chosen style or subject without loss of expressivity. Image prompts, due to their visual nature, can carry information about both subject and style. The use of magic terms can increase output variation while maintaining style. Additionally, prompt modifiers can be assigned weights to control image generation further. Negative weights can exclude certain subjects or styles, while positive weights can be used to mix styles. The process of prompt writing is iterative, starting with subject terms, followed by the addition of modifiers and solidifiers, and finally applying weights for precise control.

This work has illuminated the burgeoning field of prompt engineering, which is central to the emerging practice of text-to-image synthesis and AI-generated art. The development of a taxonomy of six types of prompt modifiers is a stepping stone to bring structure to this area of study. Future research will be critical in addressing several key areas, including the ethical and societal implications of AI-generated creative work, the social aspects of prompt engineering, the co-creation process between humans and AI, potential bias in image generation systems, and the alignment of AI with human values. As non-technical users increasingly interact with complex AI models, the need for HCI research in prompt engineering will only continue to grow. The exploration of these topics will not only advance our understanding of how people can effectively interact with machine learning models, but also inform the design of future AI-driven systems and contribute to the development of a novel digital art form.



## REFERENCES

- [1] ArtHub. 2022. arthub.ai. <https://arthub.ai/> [Accessed Nov. 9, 2022].
- [2] Stephen Bach, Victor Sanh, Zheng Xin Yong, Albert Webson, Colin Raffel, Nihal V. Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Fevry, Zaid Alyafeai, Manan Dey, Andrea Santilli, Zhiqing Sun, Srulik Ben-david, Canwen Xu, Gunjan Chhablani, Han Wang, Jason Fries, Maged Al-shaibani, Shanya Sharma, Urmish Thakker, Khalid Almubarak, Xiangru Tang, Dragomir Radev, Mike Tian-jian Jiang, and Alexander Rush. 2022. PromptSource: An Integrated Development Environment and Repository for Natural Language Prompts. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*. Association for Computational Linguistics, Dublin, Ireland, 93–104. <https://doi.org/10.18653/v1/2022.acl-demo.9>
- [3] BigScience Initiative. 2022. Introducing The World’s Largest Open Multilingual Language Model: BLOOM. (2022). <https://bigscience.huggingface.co/blog/bloom> [Accessed Nov. 9, 2022].
- [4] Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel, Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Munyikwa, Suraj Nair, Avani Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Rob Reich, Hongyu Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher Ré, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. 2021. *On the Opportunities and Risks of Foundation Models*. Technical Report. Stanford University. <https://crfm.stanford.edu/assets/report.pdf>
- [5] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating Large Language Models Trained on Code. (2021). arXiv:2107.03374 [cs.LG] [Preprint]. Available at: <https://arxiv.org/abs/2107.03374> [Accessed Nov. 9, 2022].
- [6] Jacob W. Crandall, Mayada Oudah, Tennom, Fatimah Ishowo-Oloko, Sherief Abdallah, Jean-François Bonnefon, Manuel Cebrian, Azim Shariff, Michael A. Goodrich, and Iyad Rahwan. 2018. Cooperating with machines. *Nature Communications* 9, 1 (2018), 12 pages. <https://doi.org/10.1038/s41467-017-02597-8>
- [7] Lyall Crawford. 1996. Personal ethnography. *Communication Monographs* 63, 2 (1996), 158–170. <https://doi.org/10.1080/03637759609376384>
- [8] Katherine Crowson. 2021. CLIP Guided Diffusion HQ 256x256. (2021). [https://colab.research.google.com/drive/12a\\_Wrfi2\\_gwwAuN3VvMTwVMz9TfqctNj](https://colab.research.google.com/drive/12a_Wrfi2_gwwAuN3VvMTwVMz9TfqctNj) [Accessed Nov. 9, 2022].
- [9] Katherine Crowson, Stella Biderman, Daniel Kornis, Dashiell Stander, Eric Hallahan, Louis Castricato, and Edward Raff. 2022. VQGAN-CLIP: Open Domain Image Generation and Editing with Natural Language Guidance. In *Computer Vision – ECCV 2022*, Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner (Eds.). Springer Nature, Cham, Switzerland, 88–105.
- [10] Norman K. Denzin and Yvonna S. Lincoln. 2017. *The SAGE Handbook of Qualitative Research* (5th ed.). SAGE, Thousand Oaks, CA.
- [11] Margot Duncan. 2004. Autoethnography: Critical Appreciation of an Emerging Art. *International Journal of Qualitative Methods* 3, 4 (2004), 28–39. <https://doi.org/10.1177/160940690400300403>
- [12] Remi Durant. 2021. Artist Studies by @remi\_durant. (2021). <https://remidurant.com/artists/> [Accessed Nov. 9, 2022].
- [13] Carolyn Ellis, Tony E. Adams, and Arthur P. Bochner. 2011. Autoethnography: An Overview. *Historical Social Research / Historische Sozialforschung* 36, 4 (138) (2011), 273–290. <http://www.jstor.org/stable/23032294>

- [14] Motahhare Eslami, Karrie Karahalios, Christian Sandvig, Kristen Vaccaro, Aimee Rickman, Kevin Hamilton, and Alex Kirlik. 2016. First I "like" It, Then I Hide It: Folk Theories of Social Feeds. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems (CHI '16)*. Association for Computing Machinery, New York, NY, 2371–2382. <https://doi.org/10.1145/2858036.2858494>
- [15] Harmeeet Gabha. 2022. Disco Diffusion 70+ Artist Studies. (2022). <https://weirdwonderful.ai/resources/disco-diffusion-70-plus-artist-studies/> [Accessed Nov. 9, 2022].
- [16] Iason Gabriel. 2020. Artificial Intelligence, Values, and Alignment. *Minds and Machines* 30, 3 (2020), 411–437. <https://doi.org/10.1007/s11023-020-09539-2>
- [17] Philip Galanter. 2012. *Computational Aesthetic Evaluation: Past and Future*. Springer Berlin Heidelberg, Berlin, Heidelberg, 255–293. [https://doi.org/10.1007/978-3-642-31727-9\\_10](https://doi.org/10.1007/978-3-642-31727-9_10)
- [18] Gartner. 2021. Gartner Says the Majority of Technology Products and Services Will Be Built by Professionals Outside of IT by 2024. Press release. (14 June 2021). <https://www.gartner.com/en/newsroom/press-releases/2021-06-10-gartner-says-the-majority-of-technology-products-and-services-will-be-built-by-professionals-outside-of-it-by-2024> [Accessed Nov. 9, 2022].
- [19] Susan A. Gelman and Cristine H. Legare. 2011. Concepts and folk theories. *Annual Review of Anthropology* 40 (2011), 379–398. <https://doi.org/10.1146/annurev-anthro-081309-145822>
- [20] Raymond L. Gold. 1958. Roles in Sociological Field Observations. *Social Forces* 36, 3 (1958), 217–223. <http://www.jstor.org/stable/2573808>
- [21] Mark Guzdial. 2013. Human-Centered Computing: A New Degree for Licklider’s World. *Commun. ACM* 56, 5 (may 2013), 32–34. <https://doi.org/10.1145/2447976.2447987>
- [22] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. 2022. Imagen Video: High Definition Video Generation with Diffusion Models. (2022). [Preprint]. Available at: <https://arxiv.org/abs/2210.02303> [Accessed Nov. 14, 2022].
- [23] Michaela Hoare, Steve Benford, Rachel Jones, and Natasa Milic-Frayling. 2014. Coming in from the Margins: Amateur Musicians in the Online Age. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '14)*. Association for Computing Machinery, New York, NY, 1295–1304. <https://doi.org/10.1145/2556288.2557298>
- [24] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In *Proceedings of the International Conference on Learning Representations (ICLR '20)*. 16 pages.
- [25] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. 2022. CogVideo: Large-scale Pretraining for Text-to-Video Generation via Transformers. (2022). <https://doi.org/10.48550/ARXIV.2205.15868> [Preprint]. Available at: <https://arxiv.org/pdf/2205.15868v1.pdf> [Accessed Nov. 9, 2022].
- [26] Matthew Hutson. 2022. Could AI help you to write your next paper? *Nature* 611 (2022), 192–193.
- [27] Lexica.art. 2022. Lexica.art. <https://lexica.art/> [Accessed Nov. 9, 2022].
- [28] Vivian Liu and Lydia B Chilton. 2022. Design Guidelines for Prompt Engineering Text-to-Image Generative Models. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems (CHI '22)*. Association for Computing Machinery, New York, NY, Article 384, 23 pages. <https://doi.org/10.1145/3491102.3501825>
- [29] Jon McCormack, Camilo Cruz Gambardella, Nina Rajcic, Stephen James Krol, Maria Teresa Llano, and Meng Yang. 2023. Is Writing Prompts Really Making Art? <https://doi.org/10.48550/ARXIV.2301.13049>
- [30] Jon McCormack, Toby Gifford, and Patrick Hutchings. 2019. Autonomy, Authenticity, Authorship and Intention in Computer Generated Art. In *Computational Intelligence in Music, Sound, Art and Design*, Anikó Ekárt, Antonios Liapis, and María Luz Castro Pena (Eds.). Springer International Publishing, Cham, 35–50.
- [31] Grégoire Mialon, Roberto Dessi, Maria Lomeli, Christoforos Nalmpantis, Ram Pasunuru, Roberta Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, Edouard Grave, Yann LeCun, and Thomas Scialom. 2023. Augmented Language Models: a Survey. <https://doi.org/10.48550/arXiv.2302.07842> arXiv:2302.07842 [cs.CL]
- [32] Piero Molino and Christopher Ré. 2021. Declarative Machine Learning Systems. *Commun. ACM* 65, 1 (dec 2021), 42–49. <https://doi.org/10.1145/3475167>
- [33] Carman Neustaedter and Phoebe Sengers. 2012. Autobiographical Design in HCI Research: Designing and Learning through Use-It-Yourself. In *Proceedings of the Designing Interactive Systems Conference (DIS '12)*. Association for Computing Machinery, New York, NY, 514–523. <https://doi.org/10.1145/2317956.2318034>
- [34] OpenAI. 2022. Reducing Bias and Improving Safety in DALL·E 2. (18 July 2022). <https://openai.com/blog/reducing-bias-and-improving-safety-in-dall-e-2/> [Accessed Nov. 9, 2022].
- [35] OpenAI. nd.. Completion – OpenAI API. (nd.). <https://beta.openai.com/docs/guides/completion> [Accessed Nov. 9, 2022].
- [36] OpenArt.ai. 2022. OpenArt.ai. <https://openart.ai/> [Accessed Nov. 9, 2022].
- [37] Jonas Oppenlaender. 2022. The Creativity of Text-to-Image Generation. In *Proceedings of the 25th International Academic Mindtrek conference (Academic Mindtrek '22)*. ACM, 11 pages pages. <https://doi.org/10.1145/3569219.3569352>

- [38] Guy Parsons. 2022. *The DALL-E 2 Prompt Book*. <https://dallery.gallery/wp-content/uploads/2022/07/The-DALL-E-2-prompt-book-v1.01.pdf> [Accessed Nov. 9, 2022].
- [39] Nikita Pavlichenko and Dmitry Ustalov. 2022. Best Prompts for Text-to-Image Models and How to Find Them. (2022). <https://doi.org/10.48550/ARXIV.2209.11711> [Preprint]. Available at: <https://arxiv.org/abs/2209.11711> [Accessed Nov. 9, 2022].
- [40] Sarah Pink, Heather Horst, John Postill, Larissa Hjorth, Tania Lewis, and Jo Tacchi. 2016. *Digital Ethnography: Principles and Practice*. SAGE, London, UK.
- [41] John Postill and Sarah Pink. 2012. Social Media Ethnography: The Digital Researcher in a Messy Web. *Media International Australia* 145, 1 (2012), 123–134. <https://doi.org/10.1177/1329878X1214500114>
- [42] John David Pressman, Katherine Crowson, and Simulacra Captions Contributors. 2022. *Simulacra Aesthetic Captions*. Technical Report Version 1.0. Stability AI. url <https://github.com/JD-P/simulacra-aesthetic-captions>.
- [43] Han Qiao, Vivian Liu, and Lydia Chilton. 2022. Initial Images: Using Image Prompts to Improve Subject Representation in Multimodal AI Generated Art. In *Creativity and Cognition (C&C '22)*. Association for Computing Machinery, New York, NY, 15–28. <https://doi.org/10.1145/3527927.3532792>
- [44] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *Proceedings of the 38th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 8748–8763. <https://proceedings.mlr.press/v139/radford21a.html>
- [45] Sarvapali D. Ramchurn, Sebastian Stein, and Nicholas R. Jennings. 2021. Trustworthy human-AI partnerships. *iScience* 24, 8 (2021), 13 pages. <https://doi.org/10.1016/j.isci.2021.102891>
- [46] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-Shot Text-to-Image Generation. In *Proceedings of the 38th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 8821–8831.
- [47] Laria Reynolds and Kyle McDonell. 2021. Prompt Programming for Large Language Models: Beyond the Few-Shot Paradigm. In *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems (CHI EA '21)*. Association for Computing Machinery, New York, NY, Article 314, 7 pages. <https://doi.org/10.1145/3411763.3451760>
- [48] Toran Bruce Richards. 2023. Significant-Gravitas/Auto-GPT GitHub repository. <https://github.com/Significant-Gravitas/Auto-GPT>.
- [49] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2021. High-Resolution Image Synthesis with Latent Diffusion Models. (2021). arXiv:2112.10752 [cs.CV] [Preprint]. Available at: <https://arxiv.org/abs/2112.10752> [Accessed Nov. 9, 2022].
- [50] Robin Rombach, Andreas Blattmann, and Björn Ommer. 2022. Text-Guided Synthesis of Artistic Images with Retrieval-Augmented Diffusion Models. (2022). [Preprint]. Available at: <https://arxiv.org/abs/2207.13038> [Accessed Nov. 9, 2022].
- [51] Christoph Schuhmann. 2022. LAION-Aesthetics. <https://laion.ai/blog/laion-aesthetics/> <https://laion.ai/blog/laion-aesthetics/> [Accessed Nov. 11, 2022].
- [52] Ben Shneiderman. 2020. Human-Centered Artificial Intelligence: Reliable, Safe & Trustworthy. *International Journal of Human-Computer Interaction* 36, 6 (2020), 495–504. <https://doi.org/10.1080/10447318.2020.1741118>
- [53] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. 2022. Make-A-Video: Text-to-Video Generation without Text-Video Data. (2022). <https://doi.org/10.48550/ARXIV.2209.14792> [Preprint]. Available at: <https://arxiv.org/abs/2209.14792> [Accessed Nov. 14, 2022].
- [54] Ethan Smith. 2022. A Traveler’s Guide to the Latent Space. (2022). <https://sweet-hall-e72.notion.site/A-Traveler-s-Guide-to-the-Latent-Space-85efba7e5e6a40e5bd3cae980f30235f> [Accessed Nov. 9, 2022].
- [55] Charlie Snell. 2021. Alien Dreams: An Emerging Art Scene. (2021). <https://ml.berkeley.edu/blog/posts/clip-art/> [Accessed Nov. 9, 2022].
- [56] Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang, Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. 2022. Phenaki: Variable Length Video Generation from Open Domain Textual Descriptions. (2022). <https://openreview.net/forum?id=vOEXS39nOF> [Accessed Nov. 14, 2022].
- [57] Zijie J. Wang, Evan Montoya, David Munechika, Haoyang Yang, Benjamin Hoover, and Duen Horng Chau. 2022. DiffusionDB: A Large-scale Prompt Gallery Dataset for Text-to-Image Generative Models. (2022). <https://doi.org/10.48550/ARXIV.2210.14896> [Preprint]. Available at: <https://arxiv.org/abs/2210.14896> [Accessed Nov. 9, 2022].
- [58] Jacob O. Wobbrock and Julie A. Kientz. 2016. Research Contributions in Human-Computer Interaction. *Interactions* 23, 3 (2016), 38–44. <https://doi.org/10.1145/2907069>
- [59] Wojciech Zaremba and Greg Brockman. 2021. OpenAI Codex. (2021). <https://openai.com/blog/openai-codex> [Accessed Nov. 9, 2022].

- [60] Lisai Zhang, Qingcai Chen, Baotian Hu, and Shuoran Jiang. 2020. *Text-Guided Neural Image Inpainting*. Association for Computing Machinery, New York, NY, 1302–1310. <https://doi.org/10.1145/3394171.3414017>