# *Promptly Yours?* A Human Subject Study on Prompt Inference in AI-Generated Art

Khoi Trinh[1], Joseph Spracklen[2], Raveen Wijewickrama[2],
Bimal Viswanath[3], Murtuza Jadliwala[2], Anindya Maiti[1]

*khoitrinh@ou.edu, joseph.spracklen@my.utsa.edu, raveen.wijewickrama@utsa.edu,*
*vbimal@vt.edu, murtuza.jadliwala@utsa.edu, am@ou.edu*

[1] University of Oklahoma
[2] University of Texas at San Antonio
[3] Virginia Tech

## Abstract

The emerging field of AI-generated art has witnessed the rise of prompt marketplaces, where creators can purchase, sell, or share prompts to generate unique artworks. These marketplaces often assert ownership over prompts, claiming them as intellectual property. This paper investigates whether concealed prompts sold on prompt marketplaces can be considered as secure intellectual property, given that humans and AI tools may be able to approximately infer the prompts based on publicly advertised sample images accompanying each prompt on sale. Specifically, our survey aims to assess (i) how accurately can humans infer the original prompt solely by examining an AI-generated image, with the goal of generating images similar to the original image, and (ii) the possibility of improving upon individual human and AI prompt inferences by crafting human-AI combined prompts with the help of a large language model. Although previous research has explored the use of AI and machine learning for prompt inference (and also to protect against it), we are the first to include humans in the loop. Our findings indicate that while humans and human-AI collaborations can infer prompts and generate similar images with high accuracy, they are not as successful as using the original prompt.

## 1 Introduction

Artificial Intelligence (AI) has made remarkable strides in the domain of creative and artistic expression, enabling an easy and automated process for everyone to generate visually captivating and conceptually intriguing art work. Central to this creation process of AI-generated art and images are deep learning based text-to-image (txt2img) models for image generation that utilize text prompts as input (instructions) from users to generate unique and diverse image/art outputs. Some of the most popular open-source or commercially-available examples of such txt2img models include Midjourney (MidJourney (web)), DALL-E 2 (Ramesh et al. 2021), Stable Diffusion (Rombach et al. 2022) and GLIDE (Nichol et al. 2022).

At a high level, these models have two main components – a Language Model (e.g., CLIP[1]) and a Generative image model (e.g., Stable Diffusion (Rombach et al. 2022)).

The language model converts a given text prompt to a latent representation, which is then used to condition the generative image model to produce an image that captures the prompt description. Furthermore, these models are trained on vast datasets of text-image pairs, allowing them to understand and render complex visual concepts from textual descriptions with remarkable accuracy and creativity. For example, Stable Diffusion was trained on the publicly available LAION-5B dataset, containing 5 billion image-caption pairs, derived from data scraped from the Internet.

Text prompts serve as critical input instructions to txt2img models for generating high-quality text-conditioned images and it is non-trivial to deduce an appropriate prompt for the desired image, often requiring creativity and trial-and-error (Wang et al. 2023). This has resulted in the emergence of new prompt engineering jobs and prompt marketplaces for AI-generated art, where prompt engineers, artists, and enthusiasts can exchange and sell prompts that can generate custom high-quality art. Given the importance of selecting the right prompt for generating a desired image and the non-triviality in determining one, these text prompts are often treated as protected information by their creators. Prompt marketplaces often claim intellectual property rights over the prompts, asserting that they are valuable and original creations worthy of legal protection (PromptBase; Promptrr.io). Given the protected status of (input) text prompts, two research questions arise that are largely unexplored thus far: (i) how accurately can humans infer the input text prompt by just viewing the (AI-generated) image generated from that prompt? and (ii) can AI tools assist humans in more accurately inferring text prompts of a target (AI-generated) image?

To address these research questions, our study employs a human subject survey to assess how accurately users (participants) can infer prompts by just visually examining the AI-generated images. Our survey results provides valuable insights on how users' prediction (of prompts) performance, measured using well-defined metrics, varies with different prompt-related attributes and varying user demography and backgrounds. After combining responses from the survey (taken by human subject participants) with responses from an AI-based prompt inference model, we further re-evaluate and compare the overall inference accuracy in this human-

---

[1]https://openai.com/research/clip

AI collaborative setting. Our results show that although both human and combined human-AI efforts can accurately infer prompts and recreate images to a great extent, they fall short of the effectiveness achieved with the original prompts. Consequently, marketplaces for selling prompts and creators who offer prompts for AI-generated art can continue to maintain a viable business model.

## 2 Related Work and Research Goals

### 2.1 Prompt Inference in AI Art

Previous research in the literature has primarily explored AI/ML-based inference techniques to deduce (infer) prompts and has also proposed AI/ML-based methods to safeguard against such kind of prompt inference threats. Several prior works have employed machine learning algorithms to reverse-engineer the prompts used in the creation of artworks (CLIP Interrogator; Shen et al. 2023; Wu et al. 2022; Li et al. 2022), highlighting the potential vulnerability of prompt concealment. On the other hand, efforts have also been made to develop protective measures to secure the prompts from unauthorized access or replication (Shen et al. 2023; Struppek, Hintersdorf, and Kersting 2022; Zhai et al. 2023). These efforts involve strategies to thwart prompt inference, backdoor injections, and data poisoning, implementing rigorous dataset inspections, and employing anomaly detection. However, if humans are able to infer prompts with a high degree of accuracy, the effectiveness of these protective measures against prompt inference may be called into question. Furthermore, it is possible that AI-assisted prompt inference and prompt inference by humans can be effectively combined to significantly improve the overall inference accuracy, and needs to be further studied. While there are recent studies investigating the effectiveness of human-AI collaboration leading to more creative artworks (Lyu et al. 2022; Brade et al. 2023), currently we have little understanding of how a similar human-AI collaboration may work towards prompt inference of AI-generated art.

### 2.2 Challenges in Prompt Inference

Inferring prompts of AI-generated images is a complex task that both humans and AI models (CLIP Interrogator; MidJourney (web; Shen et al. 2023) can find challenging due to the complexities in visual content and the subtleties of language. Figure 1 exemplifies this by illustrating how varying the modifiers (in this case, "pixel art" and "dark colors") for a consistent subject (a cat) can lead to vastly different visual outcomes when combined with different modifiers. Even the addition of a single modifier can dramatically alter the resulting image, demonstrating how each element in a prompt contributes to the generated image. Conversely, omitting a crucial modifier could lead to a visual representation that misses the depth or context intended. When humans try to deduce the prompts for AI-generated images, they often rely on their subjective interpretation and understanding of the visual content, which can lead to varied conclusions on the subjects and modifiers used in the image generation. The understanding of how humans may infer prompts, specifically



| cat | cat, dark colors | cat, pixel art | cat, pixel art, dark colors |

Figure 1: Image generations using SDXL with prompts containing the same subject (cat) and different combinations of two modifiers (pixel art and dark colors).

the subjects and modifiers used in AI-generated art, with or without the aid of AI-based prompt inference tools, remains an unexplored research area and is the focus of this work.

### 2.3 Research Goals

Given the above intricacies in human and AI prompt inferences, it is not trivial to assess whether concealed prompts sold on prompt marketplaces can be considered as secure intellectual property, or are they vulnerable to prompt inference and replication. Moreover, there is a lack of a clearly defined measurement based on which a prompt inference can be deemed successful, especially for human prompt inference. To address this gap, our research goals are as follows:

- Determining the accuracy with which individuals can infer the original prompts from images created by a txt2img model, aiming to reproduce similar images. This is accomplished by designing and conducting a comprehensive human participant survey as outlined in Section 3 and Section 4. Additionally, exploring whether a participant with an art background (e.g. an art major in college, a graphic designer) would have an advantage over one without such a background. Evaluation is to be based on similarity between the original image and the images reproduced using the inferred prompt.

- Determining any improvement in prompt inference accuracy when human inferred prompts are combined with AI inferred ones, aiming to further improve the similarity of the reproduced images to the original images. This is achieved by integrating the survey responses of the human participants and AI inferred prompts as discussed in Section 4.3. Evaluation is to be based on similarity between the original image and the images reproduced using the combined prompt.

- Establishing robust thresholds of metrics for measuring the success of both standalone and combined human-AI inferences, as detailed in Section 4.2, ensuring a clear framework for assessment.

## 3 Survey Design and Participants

### 3.1 Prompt and Image Datasets

The images presented to the participants were generated from two distinct types of prompts: *Controlled* and *Uncontrolled*. Controlled prompts serve as a baseline to gauge

participants' ability to recognize common subjects and modifiers, allowing us to assess other variables associated with prompt inference, such as any differences between `txt2img` models and demographic-based differences. Conversely, uncontrolled prompts feature a much more diverse mix of subjects and modifiers and is used to construct more comprehensive and challenging inference tasks for participants.

**Controlled Dataset.** The controlled prompts set was constructed through an analysis of common subjects discovered on the dynamically changing homepage of Lexica[2], a popular prompt and AI art sharing platform. We selected Lexica for this data gathering task as they have web crawling-friendly terms of service. Our *Selenium* script identified 100579 different images and accompanying prompts by their specific HTML elements. Subsequent processing with *spaCy*, a Natural Language Processing (NLP) framework, enabled us to distill the subjects from these prompts. The five most frequently occurring subjects identified were *man*, *woman*, *astronaut*, *cat*, and *robot*. In parallel with subject selection, we curated a set of modifiers by referencing a Midjourney styles and keywords repository (Wulfken (web). The modifiers spanned various categories such as themes/genres, lighting, drawing and art medium, perspective, emotion/mood, colors and palettes, geography and culture, rendering/shading style, culminating in a list of 121 modifiers (details in Appendix E.5). These subjects and modifiers were then randomly sampled and combined, ranging from one to five modifiers per subject, to formulate 100 controlled prompts, 20 per subject (Appendix E.6). These 100 prompts were then used to generated images for Parts I and III of our survey (3), 25 for each of the four selected `txt2img` models.

**Uncontrolled Dataset.** The uncontrolled prompts set was constructed through a random sampling of 100 whole/complete prompts discovered on PromptHero (a platform for discovering and sharing AI-generated art and text prompts). There were no restrictions on what subject and modifiers appeared in these prompts, as long as they were non-explicit in nature. These 100 prompts were then used to generate images for Part II of our survey, 25 for each of the four selected `txt2img` models.

## 3.2 Analyzed `txt2img` Models

Our study focuses on four popular `txt2img` models, including two base models and two fine-tuned models:

- **MidJourney v5.0** (MidJourney (web).
- **Stable Diffusion XL** (SDXL) (Podell et al. 2023; Rombach et al. 2022).
- **DreamShaper XL** (Podell et al. 2023).
- **Realistic Vision v5** (Hugging Face b).

At the time we began this study in late 2023, these models had the highest number of prompts being shared or sold on platforms such as PromptHero, Promptrr.io, Prompti AI,

---

[2]https://lexica.art

PromptBase, and CivitAI. More details on these four models and how they were employed in our study can be found in Appendix B.

Besides the above, there are other popular `txt2img` models which we also attempted to include in our study, but could not for various reasons. For instance, OpenAI's popular DALL-E model utilizes a modified version of the GPT model by adapting its transformer-based architecture for image generation (Ramesh et al. 2021). However, DALL-E was omitted from our study because we were unable to lock in a specific version of the model for image creation. We were concerned that DALL-E's updates during the time-frame of our study could affect the consistency of the images produced initially and those generated later for comparison.

For our user study to assess the ability of humans to infer prompts from AI-generated images, we design and implement a custom survey tool. This tool dynamically loads and displays images (and its metadata) together with related question(s), and can capture participant responses to these questions (see Figure 22 and Appendices E.1 to E.3). We collected responses from 230 participants using this custom tool, of which 59 were recruited from two public university campuses (in the US), while the remaining 171 participants were recruited via the Amazon Mechanical Turk (MTurk) platform. This study has been approved by the Institutional Review Boards (IRBs) of the institutions involved. Before completing the main survey task, participants completed a preliminary questionnaire (Figure 16 in Appendix E) that requested demographic information and self-reported familiarity with generative AI tools. Below, we first summarize findings from the preliminary survey, followed by a detailed description of the main survey tasks completed by the participants.

## 3.3 Pre-Survey: Demographics

30.9% of the participants in our study are in the age group of 35-44 years, followed by 26.5% in 25-34 years and 20.4% in 18-24 years. A small number of remaining participants are 45 years or older. Figure 17a and Figure 17b (in Appendix C) illustrate the age and gender distribution of our participants, respectively. Most of our participants are male (59.1%), with the remaining females (38.7%) and those who identify as other gender categories (2.2%). 15.65% participants reported having an arts background or education, indicating a certain level of expertise or familiarity with creative fields, while the remaining 84.35% did not have such a background, representing a broader range of occupations and experiences.

## 3.4 Pre-Survey: Familiarity with Generative AI

We separated the level of familiarity of our participants with various generative AI tools into four distinct categories, namely, familiarity with the text, image, audio, and video generative AI tools (Figure 2). The predominant level of familiarity for text and audio tools is "*Slightly Familiar*," (41.7% for text and 32.6% for audio), suggesting a moderate acquaintance with these technologies. For image generation tools, the results indicate an evenly spread level of familiarity, with both "*Slightly Familiar*" and "*Somewhat Famil-*
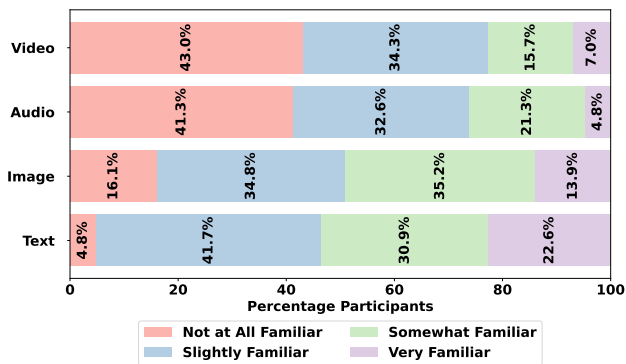
Figure 2: Participants' familiarity levels with different generative AI tools.

*iar*" categories capturing the largest proportions at roughly 35.0% each. For video generation tools, both "*Not At All Familiar*" and "*Slightly Familiar*" categories capture the largest proportions at 43.0% and 34.3%, respectively. "*Very Familiar*" holds smaller shares across video and audio tools (4.8% for audio and 7.0% for video), while for image and text generation tools "*Not at All Familiar*" remains the least represented (4.8% for text and 16.1% for image). This distribution demonstrates that our participants are moderately familiar with a broad range of generative AI tools.

## 3.5 Parts I-III: Prompt Inference

Our pool of participants, both recruited on campus and through MechTurk, participated in the main survey, comprised of three sequentially ordered parts (or phases), where in each part participants were tasked with either selecting or typing-in the most appropriate prompt (comprising of a subject and modifiers) for a series of AI-generated images. In Parts I and III, a selection of 100 images using a controlled set of subjects and modifiers, each methodically produced using the four image generation models outlined in Section 3.2, were displayed to the participants. Each of these models contributed an equal share of images to the controlled dataset. Consistent exposure to all subjects and modifiers in the controlled set was maintained by appropriately varying the subset of images each of the participants were exposed to.

In contrast, Part II expanded the scope of inference by using 100 images from an uncontrolled dataset. These images were generated from popular prompts found on online prompt sharing websites, offering a wide array of subjects and modifiers. This uncontrolled set also used the same four image generation models. Detailed information on the subjects and modifiers selected for the controlled dataset and the composition of the uncontrolled dataset can be found in Section 3.1.

**Part I: Subject and Modifier Selection in Controlled Dataset.** In this part, participants were presented with a sequence of five AI-generated images with controlled subjects and modifiers (i.e, from the controlled dataset). For each displayed image, they were required to first select the correct

subject from the options provided. Next, they had to choose the correct modifier(s) from the given options (between 1 to 5 checkbox options), with the possibility of selecting multiple modifiers if preferred. For each question, there will be one correct subject, with the rest of the options filled in randomly; at least one and at most five correct modifiers. For questions where there are less than five correct modifiers, the remaining modifiers are filled in randomly. Figure 19 shows an example of a question for this part. The goal of this task was to measure the participants' ability to discern and match the given images to their corresponding prompts accurately.

**Part II: Prompt Creation in Uncontrolled Dataset.** This second part challenged participants with another five AI-generated images, but this time the prompts were uncontrolled, i.e., no pre-selected options were provided for both subjects and modifiers. In this task, for each image participants had to type what they believed to be the most fitting subject and modifier(s) as a complete sentence. Figure 20 (in (in Appendix E.2) shows an example of a question for this part. This open-ended task assessed their creative inference abilities and how they translated their interpretation of the images into descriptive prompts.

**Part III: Prompt Creation in Controlled Dataset.** In this final task, participants were shown yet another set of five AI-generated images. Similar to Part II, they were instructed to generate sentence-like prompts for these images. However, in Part III the images were again based on controlled subjects and modifiers (i.e., from the controlled dataset). Figure 21 (in Appendix E.3) shows an example of a question for this part. The responses of the participants were then compared with the actual prompts and images displayed to them, providing a measure of how well the participants could infer and replicate the structured prompts that were initially used to generate the images displayed.

## 3.6 Part IV: Rating Similarity

In addition to parts I through III, a subset (all of our 171 MTurk participants) were asked to take an additional part. In this last part (Part IV) of the study, participants rated the similarity between pairs of images, *one shown to prior participants and one generated using their (prior participants') prompts*. Specifically, participants were presented with pairs of images, similar to the ones depicted in Figure 22 (in Appendix E.4). Each pair consisted of an original image from Part III and a newly generated image based on a prompt submitted by one of the first 25 (on campus) participants. The participants' task for each image pair was to assess and rate the similarity between the two images. They were provided with a Likert scale ranging from "Not at All Similar" to "Very Similar". This allowed participants to express their perceived degree of similarity, taking into account factors such as color scheme, composition, and subject representation. After making their selection, they would proceed to the next comparison until they had evaluated all five pairs assigned to them.

## 3.7 Survey Responses

In total, 230 on campus and MTurk participants took our survey between December 2023 and March 2024. The responses collected consists of subject and modifier selections from a fixed set of options in Part I, and whole prompts in Parts II and III of the survey. After completion of the survey and elimination of invalid responses, we recorded 1078 responses for Part I, 1145 responses for Part II, and 1141 responses for Part III. In addition, the participants who were selected for Part IV of our study gave us 855 responses. In Parts II and III, responses that were incoherent or irrelevant, such as the failure to specify at least one subject or modifier, were also excluded.

# 4 Experimental Setup

## 4.1 Metrics

Next, we define the metrics that we employ to assess the prompt inference accuracy of the survey participants in our experiments. By means of these metrics, we aim to compare and compute the discrepancy between the original prompts and participants' responses (inferred prompts), and between the original image (displayed to the participants) and images generated from their (inferred) prompts. A detailed description of these metrics can be found in Appendix D. These metrics provide a comprehensive assessment of both objective (e.g., image hash, perceptual similarity) and subjective (e.g., surveyed similarity ratings) aspects, ensuring a well-rounded evaluation of the participants' performance in the study.

- **MSQ Scores (Survey Part I).** Score between 0 (no correct selection) to 2 (all correct selections) for each question to evaluate multiple selection questions (MSQs).

- **Image Hash**[3]**.** Hamming distance between the hashes of two images. A lower image hash score (difference) implies that the two images are similar, and vice versa.

- **Perceptual Similarity (Zhang et al. 2018).** A lower perceptual similarity score implies that the two images are similar, and vice versa.

- **Image Embedding Similarity (CLIP Score) (Wang, Chan, and Loy 2023).** CLIP score between an image and a prompt (text) where a higher score indicates greater relevance or similarity. We employ two state-of-the-art CLIP models in our evaluation, OpenAI's ViT-L/14 Transformer[4] (L14) and ViT-B/32 Transformer[5] (B32).

- **Text Embedding Similarity (Semantic Similarity) (Reimers and Gurevych 2019).** A higher semantic similarity implies the two prompts are similar, and vice versa.

- **Surveyed Similarity Rating (from Survey Part IV).** Participants in Part IV of the study rated pairs of images on a Likert scale consisting of "Not at All Similar", "Slightly Similar", "Somewhat Similar", and "Very Similar."

---

[3]https://pypi.org/project/ImageHash/

[4]https://huggingface.co/openai/clip-vit-large-patch14

[5]https://huggingface.co/openai/clip-vit-base-patch32

## 4.2 Quantifying Successful Inferences

Given our study's focus on assessing the accuracy of human prompt inference is ultimately in understanding their ability to generate images resembling those presented to them, we now establish a quantifiable measure of *success threshold* based on the above metrics, specifically image hash, perceptual similarity, and CLIP scores. Considering the inherent variability in txt2img generations for identical prompts, achieving perfect scores (1 for CLIP score, 0 for image hash, and perceptual similarity) is implausible even with completely accurate prompt inference, as illustrated by Figures 18a and 18b scores.

Therefore, a more appropriate threshold for gauging successful inference involves analyzing the range of scores for images generated from the same prompt across different instances. Accordingly, to determine an effective *success threshold*, we assessed the scores generated from identical prompts across multiple image creations, using our 100 controlled dataset prompts (Appendix E.6). For each of these prompts, we generated two different images using SDXL, and calculated the image hash, perceptual similarity, and B32 and L14 CLIP scores. After calculating the averages of these scores, we determined the success thresholds, denoted by $\theta$, as follows: $\theta_{\text{hash}} = 26.83$ for image hash, $\theta_{\text{ps}} = 0.575$ for perceptual similarity, and $\theta_{\text{B32}} = 0.875$ and $\theta_{\text{L14}} = 0.851$ for B32 and L14 CLIP scores, respectively. Participants whose inferred prompts are able to generate images close to or better than these thresholds can be deemed successful in their inference.

## 4.3 Human-AI Combined Inference

We now detail our experimentation on the effectiveness of combining human inferred prompts with those produced by AI models towards accurately recreating AI-generated art. Specifically, we utilize prompt responses collected in Part III of our survey, which involved controlled dataset prompts, and AI-generated prompts obtained through the CLIP Interrogator, a model that analyzes an image and generates descriptive text prompts by leveraging OpenAI's CLIP model to match images and text representations (CLIP Interrogator). These human and AI prompt pairs are then consolidated into one combined prompt for comparative evaluation against only human inferred prompt, CLIP Interrogator prompt, and the success threshold, in Section 5.3.

To construct a combined prompt from each pair of human prompt and corresponding CLIP Interrogator prompt, we employ a large language model such as GPT-4 (OpenAI 2023), and instruct it to merge the two prompts into a succinct new prompt of up to 25 words without adding any extraneous information. An example of this process is illustrated below, where the first prompt was from a participant and the second prompt was generated by CLIP Interrogator:

**Instruction**: Combine these two prompts into a new prompt of 25 words, without any extra information added:
1. man dressed in steampunk in a steampunk factory
2. a man in a steampunk suit and top hat standing in front of a giant clock with gears, Bastien L. Deharme, steampunk, a character portrait, fantasy art
**GPT-4 Response**: a man in a steampunk suit and top hat stands in a factory, surrounded by giant clocks and gears, embodying a fantasy art portrait.

The combined prompts were capped at 25 words to prevent GPT-4 from inadvertently introducing additional keywords when given unrestricted length. Employing a large language model was considered more suitable than merely concatenating the two prompts to avoid repeating subjects and modifiers. Such repetitions could skew the `txt2img` generations by inadvertently emphasizing repeated subjects and modifiers over others.

## 5 Results and Analysis

### 5.1 Subject and Modifier Selection Evaluation

In Figure 3, the $y$-axis illustrates the Part I MSQ score distributions over the specified categories (models, subjects, and demographic), reflecting participants' ability to match AI-generated images with their corresponding subjects and modifiers where a higher score indicate a more accurate prompt matching. As seen in Figure 3a, Midjourney-generated images showed a relatively high score concentration towards the upper end, with a median score of 1.5 and an interquartile range (IQR) ranging from 1.33 to 1.9. Realistic Vision 5's generations demonstrated a median of 1.38 and an IQR from 1.2 to 1.6. DreamShaper XL's images displayed a median of 1.49 and an IQR from 1.28 to 1.75, suggesting slightly more consistency than Realistic Vision 5. Stable Diffusion XL matched Midjourney in median score but had a narrower IQR from 1.3 to 1.67.

Analyzing scores by subject in Figure 3b, cat-themed images outperformed others with a median of 1.5 and an IQR from 1.33 to 1.84, while astronaut images lagged behind with a median of 1.33 and an IQR from 1.23 to 1.58. Next, dissecting the impact of an art background, Figure 3c shows that participants with art-related education or employment achieved slightly higher median scores at 1.5 and an IQR of 1.32 to 1.61 compared to those without, who had a median of 1.48 and a similar IQR but with a broader range and more outliers. Participant demographics split by recruitment source, shown in Figure 3d, highlight differences in scoring patterns. MTurk workers registered a median score of 1.48 and an IQR of 1.31 to 1.6, while university-recruited participants posted a slightly higher median of 1.54 with an IQR from 1.42 to 1.67. Notably, the MTurk group exhibited a broader score range and more outliers, indicating variances in scoring behavior between the two groups. These results broadly imply that when given with options, participants accurately identified subjects and modifiers in AI-generated images, achieving high scores across different models.
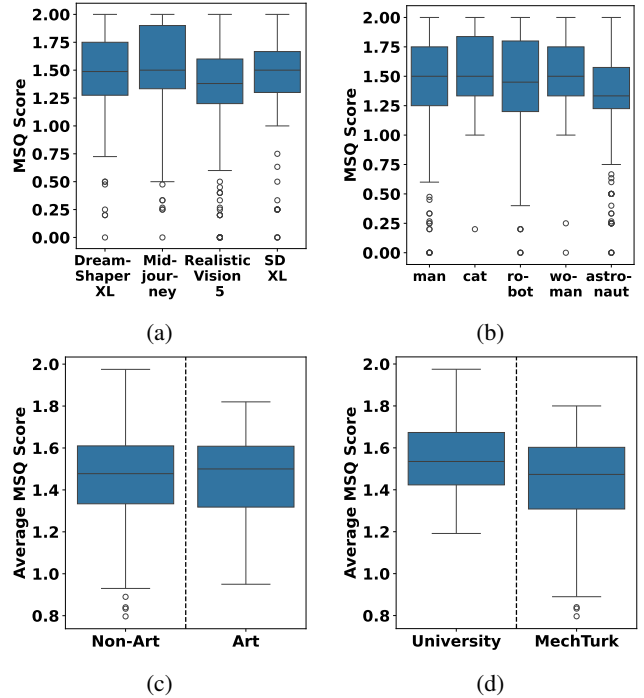


Figure 3: Analysis of multiple answer question (MSQ) responses from Part I, aimed at identifying disparities in (a) four distinct `txt2img` models, (b) various subjects depicted in the images, (c) the impact of participants' arts background, and (d) variations attributable to recruitment sources.

### 5.2 Human Inference Evaluation

We next utilize metrics such as the image hash, perceptual similarity, CLIP score, semantic similarity, and survey similarity ratings (outlined in Section 4.1) as applicable, to evaluate various factors affecting human inference accuracy. These metrics have ranges of 0 to 1 for CLIP score, semantic similarity, and perceptual similarity; and 0 to 64 for image hash score.

**Models.** Both the uncontrolled and controlled datasets (from Parts II and III, respectively) serve as our primary data sources for this analysis where we compare the images generated from prompts submitted by the participants using our four selected `txt2img` models. These images are compared to the original image from the corresponding correct prompt using the the similarity metrics.

Figure 4 shows the distribution of semantic similarity and CLIP scores obtained for the participant responses for survey Parts II and III. For the controlled dataset (Figure 4a), semantic similarity (a higher value represents a higher similarity between the original and the participant inferred prompt) reveals that the L14 model falls short of the B32 model's performance, with a notable difference in median scores (25.687% lower on average for L14). This pattern persists in the CLIP scores (where a higher value represents a higher similarity between the original and the inferred prompt gen-
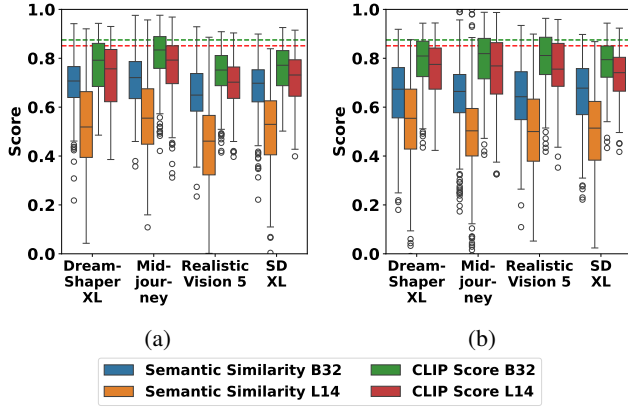
Figure 4: Semantic similarity and CLIP score for different `txt2img` models, in (a) controlled dataset, and (b) uncontrolled dataset. Success thresholds are depicted as dashed lines, green dashed line for the L14 model and red dashed line for the B32 model.
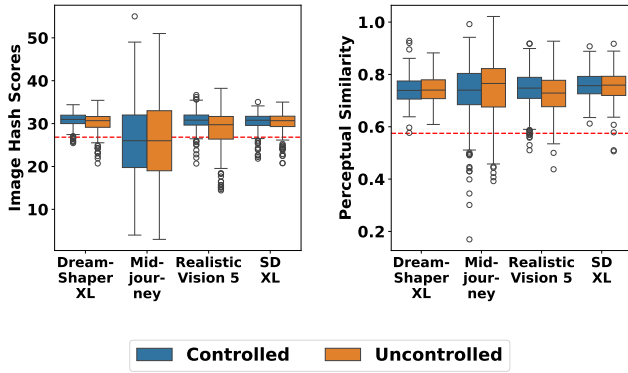


Figure 5: Perceptual similarity and image hash scores for images generated using the four different `txt2img` models.

erated image), where B32 outperforms L14, emphasizing the latter's lower median by 5.294% on average. Among the `txt2img` models, Midjourney stands out for its superior semantic similarity and CLIP scores, indicating a closer match to the inferred prompts, whereas Realistic Vision 5 consistently records the lowest score ranges across these metrics.

In the uncontrolled dataset (Figure 4b), we observe a similar trend, with the L14 model lagging behind the B32 model in both semantic similarity (the median is 22.06% lower on average) and CLIP scores (the median 5.99% lower on average) (Figure 4b). The majority of the CLIP scores observed fell below the success thresholds (dashed lines), $\theta_{L14} = 0.851$ and $\theta_{B32} = 0.875$, as denoted by the red and green dashed lines respectively in Figure 4. However, it is noteworthy that 162 images from the Midjourney model out of 2,286 total images were able to surpass these thresholds under B32 across both controlled and uncontrolled datasets.

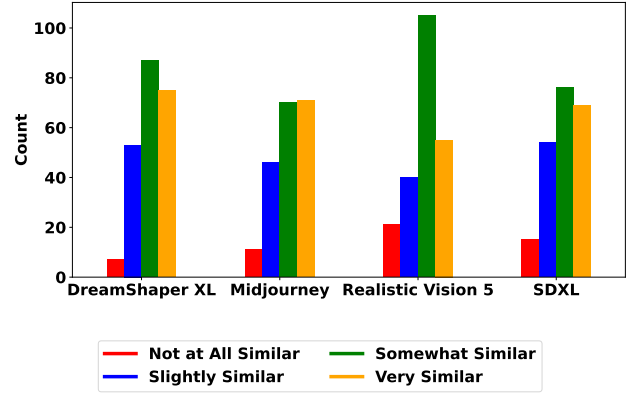The perceptual similarity and image hash scores, as



Figure 6: Part IV similarity rating distributions for images generated using the four different `txt2img` models.

shown in Figure 5, reveal the influence of dataset conditions (controlled vs. uncontrolled) on image generation in a somewhat counter-intuitive manner. Under controlled conditions, the hash scores are higher, with the median hash score across all models being on average 1.185% higher than that for the uncontrolled dataset (note that a lower image hash scores indicate greater similarity). Similarly, the median perceptual similarity is higher, although only marginally, at 0.329% on average compared to the uncontrolled dataset (note that a lower perceptual similarity score is more desirable, i.e., greater similarity). These findings indicate that a more structured approach to prompt inference results in the generation of images that considerably deviate from the original images compared to uncontrolled prompt inference.

In terms of success thresholds, a pattern similar to CLIP scores was observed in the image hash and perceptual similarity evaluations. The success thresholds are depicted in Section 5.2 and Section 5.2 by red dashed lines. Surveyed similarity ratings (Figure 6) provide direct insights into the perceived accuracy of images generated from human-inferred prompts. DreamShaper XL notably receives a majority of "Somewhat Similar" ratings (87), along with the highest count of "Very Similar" ratings at 75, suggesting that images generated with this model are likely to be more similar to the corresponding original image.

These analyses reveals how different `txt2img` models perform while interpreting prompts and corresponding image generations, emphasizing the challenges in reproducing AI art through prompt inference that aligns with human interpretations. Despite certain minor favorable observations, the overall performance across all three types of metrics did not reach the success thresholds, indicating the challenging nature in human prompt inferences.

**Subjects.** We next investigate the inference accuracy for different subjects in prompts using the controlled dataset. Figure 7 illustrates the semantic similarity and CLIP scores (where higher scores indicate a greater match), for five selected subjects. It is observed that the L14 model typically scores lower compared to the B32 model in both metrics. Subjects such as *robots* and *cats* achieve high CLIP scores,
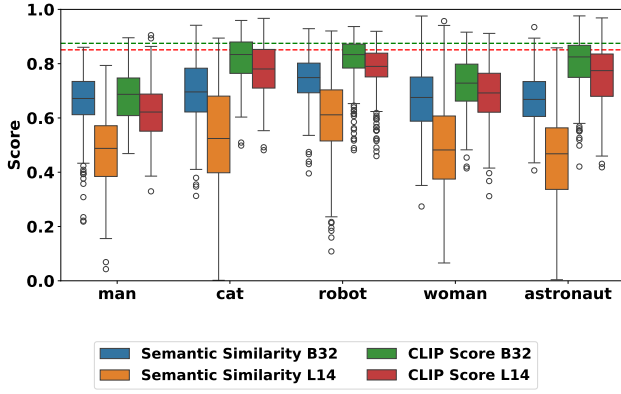
Figure 7: Semantic similarity and CLIP score for images containing different subjects, in controlled dataset.
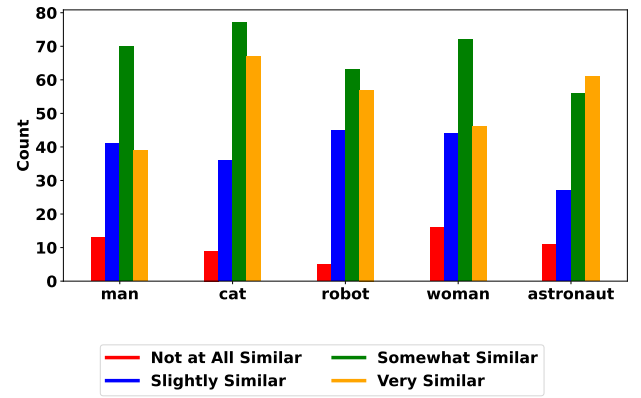


Figure 8: Perceptual similarity and image hash scores for images containing different subjects, in controlled dataset.



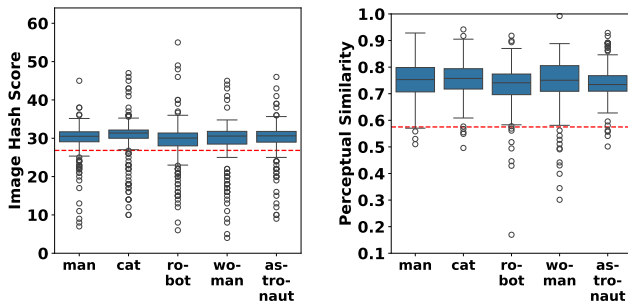Figure 9: Part IV similarity rating distributions for images containing the five different subjects.



Figure 10: Top 10 and bottom 10 modifiers by weighted average ratings of the overall image similarity, by participants in Part IV.

suggesting a strong alignment with their respective prompts. This observation prompts an inquiry into the nature of the subjects' representations: are they inherently distinct, consistently depicted by the models, or is it that humans are inherently better at inferring prompts related to these subjects? Inherent distinctness would imply that *robots* and *cats* possess unique, identifiable features that are readily recognized by the B32 and L14 models. Consistent depiction, on the other hand, would result in uniformity in representing these subjects across various image generations. Meanwhile, if humans are inherently better at inferring prompts for these particular subjects, this could also contribute to the observed accuracy. While *cats* may not always meet the desired benchmarks, in the top 25% of cases, their inference performance reaches the success threshold, with 66 inference instances above $\theta_{B32}$ and 61 instances above $\theta_{L14}$. Except for *man* and *woman*, all subjects exhibit at least 30 instances surpassing the thresholds.

Perceptual similarity and image hash scores, as depicted in Figure 8, indicate a consistent representation across subjects, although *cats* show slightly higher scores, hinting at a less consistent inference process. For the image hash metric, only 30 to 40 instances for each subject show scores below the success threshold, indicating that lower scores, which signify better accuracy, are not frequently achieved.
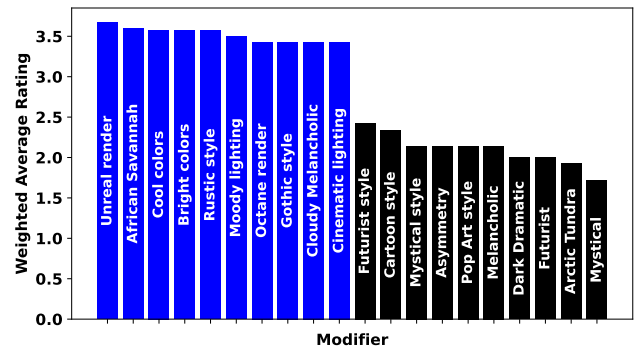
However, the performance in terms of perceptual similarity is notably poorer across all subjects, with fewer than 12 instances falling below the success threshold.

According to surveyed similarity ratings (Figure 9), images of *cats* and *astronauts* predominantly received "Very Similar" ratings (67 and 61, respectively), denoting high inference accuracy. *Robots* also achieved several "Very Similar" ratings (57) but with a noticeable amount of "Slightly Similar" ratings (45) as well, indicating some variability. Images depicting *man* are mostly rated as "Somewhat Similar," (70) which suggests moderate inference accuracy, with a distribution of ratings across different similarity levels, similar to images of *woman*.

Despite certain subjects such as *cats* and *robots* aligning moderately well in certain metrics, the overall analysis across subjects underscores that human prompt inference generally does not meet the success thresholds, emphasizing the inherent challenge of the task. These results suggest that while there are instances of alignment, the overall capability of humans to accurately infer prompts remains below par.

**Modifiers.** We next explored how modifiers present in the image generation prompts could affect the participants' prompt inference accuracy. To better quantify this, we calcu-

late a *Weighted Average Rating* to assess the ranking of the modifiers through the images shown in Part IV. This metric is determined by assigning a normalized value to each familiarity rating, where 1 represents "Not at All Similar" and 4 represents "Very Similar". The process involves counting the frequency of each rating, multiplying it by its normalized value, and then calculating the average of these products. Formally, the Weighted Average Rating, $W$, is expressed as:

$$W = \frac{\sum_{i=1}^{n} f_i \cdot v_i}{\sum_{i=1}^{n} f_i}$$

where $f_i$ is the frequency of the $i$-th rating, $v_i$ is the normalized value of the $i$-th rating, and $n$ is the number of rating types. Figure 10 shows 20 modifiers from our list of 121, the top 10 modifiers with the highest weighted average, in blue; and the bottom 10 modifiers, with the lowest weighted average, in black. Notably, some of the top modifiers such as "Unreal render," and "Cloudy Melancholic" are potentially challenging to be inferred directly by participants. This observation suggests two possibilities: either the top modifiers exert substantial influence, making them easily identifiable and thus straightforward to infer by participants, or, despite their complexity, these modifiers do not significantly impact the image generation process, thereby not detrimentally affecting the participants' ability to recognize the intended imagery.

**Art Background vs. Other Backgrounds.** Figure 24 in Appendix F shows that participants with an art background tend to introduce more variability and slightly lower scores in semantic similarity and CLIP scores across both controlled (where the median of semantic similarity from B32 is 0.465% lower on average and median CLIP score of B32 is 0.172% lower on average) and uncontrolled (where the median of semantic similarity from B32 is 0.292% lower on average, CLIP score B32 is 0.057% lower on average) datasets. This suggests that their refined ability to discern certain image characteristics may lead to more critical or complex interpretations of the AI-generated images, potentially diverging from the target images more than those without an art background.

Participants for both groups have their CLIP score below the success thresholds (green and red dashed lines in Figure 24 indicate the thresholds $\theta_{B32}$ and $\theta_{L14}$ respectively.). In the controlled data set, only 1 out of 134 participants in the Non-Art were able to have an average CLIP B32 score above the threshold; while the inverse is true for the uncontrolled dataset; where, notably, there were only 1 participant with Art background (out of 36) and 13 Non-Art participants (out of 134) were able to achieved scores above the thresholds.

Perceptual similarity and image hash scores (Figure 25) show negligible differences between the art and non-art groups (e.g. perceptual similarity median only shows a 0.02% difference between the group in the uncontrolled dataset, and 0.326% difference in the controlled dataset). Surveyed similarity ratings (Figure 23) also does not provide any significant difference between two participant groups. In summary, the involvement of individuals with certain expertise in the field, such as those with an art background, does not appear to significantly enhance the accuracy of human prompt inference.

**MTurk Workers vs. University Population.** Figure 26 in Appendix F indicates that university participants generally achieved higher median scores across all categories compared to MTurk workers (semantic similarity B32 is 3.168% higher and CLIP score B32 is 4.311% higher for uncontrolled dataset on average, while the controlled dataset sees the semantic similarity B32 to be 2.2% higher, and CLIP score B32 to be be 4.446% higher on average), suggesting a more accurate inference of prompts. This difference is particularly notable in the controlled dataset, where the tasks may be inherently more straightforward due to the structured nature of the prompts. However, a significant observation is the broader range of scores among MTurk participants, pointing to a greater variability in their prompt interpretations. In terms of success thresholds, participants for both groups have their CLIP scores below the success threshold. In the uncontrolled dataset, there were 7 participants from either group who achieved a CLIP B32 scores above the success threshold. While in the controlled dataset, there were only 1 out of 59 participants in the university group with their score above the threshold.

Perceptual similarity and image hash scores (Figure 27 in Appendix F) further supports this observation, showing that median scores are closely aligned between university and MTurk participants. In summary, while university participants tend to perform better on average, indicating potentially more precise prompt inferences, the variability among MTurk workers highlights a diverse range of interpretations and evaluations.

**AI Tool Familiarity.** Participants with limited familiarity with generative AI tools exhibited higher inference accuracy across metrics, such as higher B32 and L14 CLIP scores (Figures 28a and 28b in Appendix F), as well as lower image hash and perceptual similarity scores (Figures 28c and 28d) in the controlled dataset. This outcome may suggest that participants without extensive AI tool exposure are either applying broader interpretations that align well with AI-generated images or that their fresh intuition-based prompt inference aligns better with the `txt2img` models' generative behavior. The uncontrolled dataset reaffirms this trend, where participants with minimal AI tool familiarity again demonstrated slightly higher prompt inference accuracy (Figures 29a to 29d). Being more familiar with generative AI tools does not seem to affect how close the participant's score is to the success thresholds and overall all familiarity categories still fall short of the success thresholds across all the metrics.

The cumulative findings thus far indicate that human performance in prompt inference is generally below par, regardless of the image generation models, subjects or modifiers involved, or the demographic backgrounds of the participants. Yet, before drawing any definitive conclusions, we will explore the potential for enhancing prompt inference accuracy through a synergistic approach that combines AI and human capabilities.
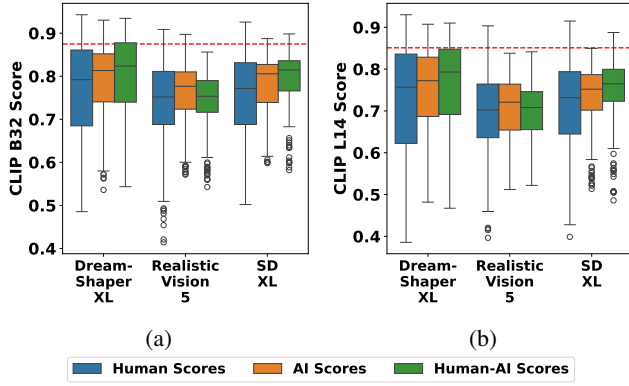
Figure 11: B32 (a) and L14 (b) CLIP scores for different models from prompts generated by human, AI, and human-AI combination.
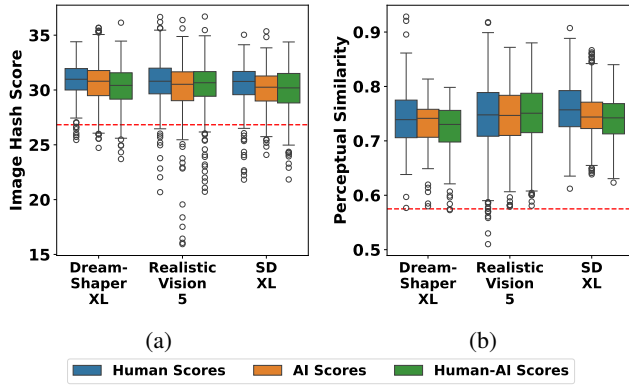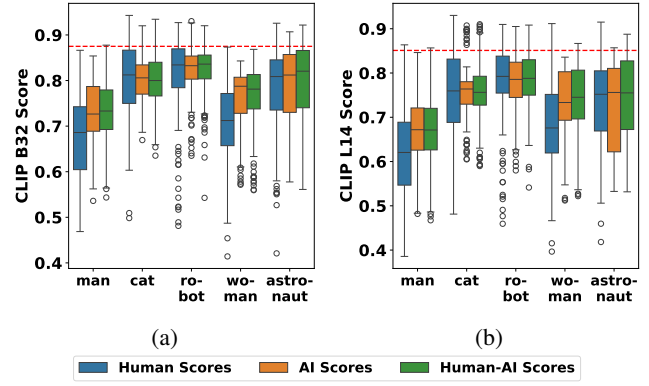


Figure 13: B32 (a) and L14 (b) CLIP scores for different subjects and modifiers from prompts generated by human, AI, and human-AI combination.



Figure 12: Image hash score (a) and perceptual similarity (b) for different models from prompts generated by human, AI, and human-AI combination.
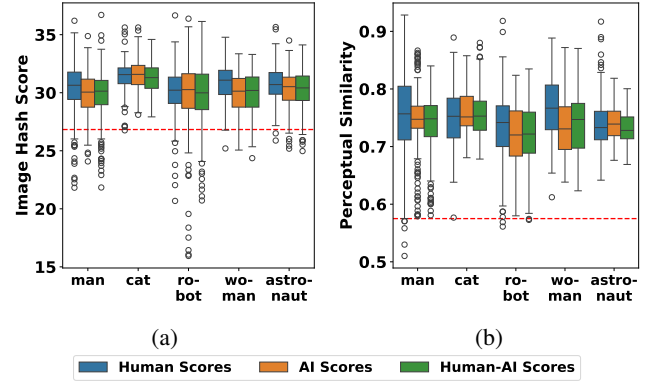


Figure 14: Image hash score (a) and perceptual similarity score (b) for different subjects and modifiers from prompts generated by human, AI, and human-AI combination.

## 5.3 Human-AI Combined Inference

We next examine the effectiveness of combining human inferences with AI inferences, particularly using CLIP interrogator for inference and then merging them using GPT-4 (as outlined in Section 4.3), towards accurately inferring prompts for AI-generated images. CLIP score analysis shows that combination of humans and AI inferred prompts generally results in higher image CLIP scores (Figure 11). Specifically, the B32 scores were found to be 3.126% higher on average in a combined setting compared to individual efforts. However, when compared to the images generated by prompts inferred by CLIP Interrogator, the B32 score saw a 0.23% decrease on average, due to the Realistic Vision 5-generated images performing worse; although Midjourney-generated images performed 1.265% better, and SDXL images with a 1.099% increase in performance (Figure 11a). With that being said, the overall improvement highlights how combining human inference with AI inference can lead to slightly better accuracy in matching prompts with

generated images. The overlap in scores across different `txt2img` models suggests that while human-AI combined prompts enhances accuracy, the degree of improvement may not be drastically distinct among different AI models.

Perceptual similarity and image hash score (Figure 12) also indicate that images generated through human-AI combined prompts align more closely with target images, as evident from the slightly lower average hash scores (by 1.417%) and perceptual similarity scores (by 0.905%). Figures 13 and 14 further explores this combined prompt inference, showcasing that human-AI combined prompts not only consistently improves inference across various subjects but also indicates a nuanced enhancement in how specific subjects and modifiers are interpreted and matched (hash score median is on average 1.438% lower, while perceptual similarity median is on average 1.431% lower).

For all 4 types of metrics, combining AI inferred prompts through CLIP interrogator with human inferred prompt has a positive effect on the overall accuracy. While the human-AI combined prompts still did not reach the respective metrics'

success thresholds, its performance is slightly better than that of purely human inferred prompts. This suggests that accurate prompt inference by a human would benefit from combination with AI inference.

The enhanced accuracy in prompt inference through the combination of human and AI efforts can be attributed to the complementary skills and knowledge each brings. Humans excel in creative and contextual understanding, while AI provides extensive data analysis and pattern recognition capabilities. This combination allows for error reduction, iterative refinement, and a broader knowledge base, leading to more accurate prompts. Despite these benefits, the fact that performance did not reach the success thresholds suggests that further optimizations in human-AI interaction are required.

In summary, the empirical evidence gathered from the evaluation, shows that while humans alone may struggle with prompt inference, their efforts, when combined with AI, demonstrate a modest improvement. The findings suggest there's room for growth in human-AI collaborations, which could eventually enhance the efficacy of prompt reconstruction. Until such advancements are realized, prompt marketplaces maintain their practicality as a business model in the realm of AI-generated art.

## Conclusion

The study presented in this paper explores the intellectual property considerations of prompts in the realm of AI-generated art, particularly focusing on the ability of humans to infer these prompts by solely examining the resulting artworks. The below par results from our human only prompt inference experiments suggest that the performance of human participants is influenced by the complexities in the prompt and the generated artwork. However, the combination of prompts from AI models such as CLIP interrogator has a positive effect on improving the inference efforts of the human participants. Collectively, these observations point towards the feasibility of generative AI prompt marketplaces as viable business models, where the uniqueness and creativity of prompts can be valued and traded while maintaining their secure intellectual property. Moreover, the exploration of combining human insights with AI in improving prompt inference opens new avenues for research.

## References

Brade, S.; Wang, B.; Sousa, M.; Oore, S.; and Grossman, T. 2023. Promptify: Text-to-image generation through interactive prompt exploration with large language models. In *ACM UIST*.

CLIP Interrogator. (web). https://huggingface.co/spaces/pharmapsychotic/CLIP-Interrogator.

Hugging Face. (web)a. runwayml/stable-diffusion-v1-5. https://huggingface.co/runwayml/stable-diffusion-v1-5.

Hugging Face. (web)b. SG161222/Realistic-Vision-V5.1-noVAE. https://huggingface.co/SG161222/Realistic-Vision-V5.1-noVAE.

Hugging Face. (web)c. stabilityai/stable-diffusion-xl-base-1.0. https://huggingface.co/stabilityai/stable-diffusion-xl-base-1.0.

Lanz, J. A. (web). Inside the Lucrative New Business of Selling AI Prompts. https://decrypt.co/137689/lucrative-new-business-selling-ai-prompts.

Li, J.; Li, D.; Xiong, C.; and Hoi, S. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*.

Lyu, Y.; Wang, X.; Lin, R.; and Wu, J. 2022. Communication in human–AI co-creation: Perceptual analysis of paintings generated by text-to-image system. *Applied Sciences*, 12(22).

MidJourney. (web)a. methexis-inc / img2prompt. https://replicate.com/methexis-inc/img2prompt.

MidJourney. (web)b. Midjourney Model Versions. https://docs.midjourney.com/docs/model-versions.

Nichol, A. Q.; Dhariwal, P.; Ramesh, A.; Shyam, P.; Mishkin, P.; Mcgrew, B.; Sutskever, I.; and Chen, M. 2022. GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models. In *ICML*.

OpenAI, R. 2023. GPT-4 technical report. *ArXiv*, 2303.

Podell, D.; English, Z.; Lacey, K.; Blattmann, A.; Dockhorn, T.; Müller, J.; Penna, J.; and Rombach, R. 2023. SDXL: Improving latent diffusion models for high-resolution image synthesis. *ArXiv*.

PromptBase. (web). https://promptbase.com/tandcs.

Promptrr.io. (web). Terms of Service — Promptrr.io. https://promptrr.io/terms-of-service/.

Ramesh, A.; Pavlov, M.; Goh, G.; Gray, S.; Voss, C.; Radford, A.; Chen, M.; and Sutskever, I. 2021. Zero-shot text-to-image generation. In *ICML*.

Reimers, N.; and Gurevych, I. 2019. Sentence-BERT: Sentence embeddings using siamese bert-networks. *ArXiv*.

Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *IEEE/CVF CVPR*.

Shen, X.; Qu, Y.; Backes, M.; and Zhang, Y. 2023. Prompt Stealing Attacks Against Text-to-Image Generation Models. *ArXiv*.

Struppek, L.; Hintersdorf, D.; and Kersting, K. 2022. Rickrolling the artist: Injecting invisible backdoors into text-guided image generation models. *ArXiv*.

Timothy, M. (web). Are Premium AI Prompts Worth the Money? https://www.makeuseof.com/should-you-buy-ai-prompts/.

Wang, J.; Chan, K. C.; and Loy, C. C. 2023. Exploring clip for assessing the look and feel of images. In *AAAI Conference on Artificial Intelligence*, volume 37.

Wang, J.; Liu, Z.; Zhao, L.; Wu, Z.; Ma, C.; Yu, S.; Dai, H.; Yang, Q.; Liu, Y.; Zhang, S.; et al. 2023. Review of large vision models and visual prompt engineering. *Meta-Radiology*, 100047.

Wu, Y.; Yu, N.; Li, Z.; Backes, M.; and Zhang, Y. 2022. Membership inference attacks against text-to-image generation models. *ArXiv*.

Wulfken, W. (web). MidJourney Styles and Keywords Reference. https://github.com/willwulfken/MidJourney-Styles-and-Keywords-Reference/tree/main/Pages/MJ_V4/Style_Pages/Just_The_Style.

Zhai, S.; Dong, Y.; Shen, Q.; Pu, S.; Fang, Y.; and Su, H. 2023. Text-to-image diffusion models can be easily backdoored through multimodal data poisoning. *ArXiv*.

Zhang, R.; Isola, P.; Efros, A. A.; Shechtman, E.; and Wang, O. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *IEEE/CVF CVPR*.

## Paper Checklist

1. For most authors...

   (a) Would answering this research question advance science without violating social contracts, such as violating privacy norms, perpetuating unfair profiling, exacerbating the socio-economic divide, or implying disrespect to societies or cultures? Yes, this survey should not have any of the negative impact mentioned above.

   (b) Do your main claims in the abstract and introduction accurately reflect the paper's contributions and scope? Yes

   (c) Do you clarify how the proposed methodological approach is appropriate for the claims made? NA

   (d) Do you clarify what are possible artifacts in the data used, given population-specific distributions? NA

   (e) Did you describe the limitations of your work? NA

   (f) Did you discuss any potential negative societal impacts of your work? No

   (g) Did you discuss any potential misuse of your work? NA

   (h) Did you describe steps taken to prevent or mitigate potential negative outcomes of the research, such as data and model documentation, data anonymization, responsible release, access control, and the reproducibility of findings? NA

   (i) Have you read the ethics review guidelines and ensured that your paper conforms to them? Yes

2. Additionally, if your study involves hypotheses testing...

   (a) Did you clearly state the assumptions underlying all theoretical results? NA

   (b) Have you provided justifications for all theoretical results? NA

   (c) Did you discuss competing hypotheses or theories that might challenge or complement your theoretical results? NA

   (d) Have you considered alternative mechanisms or explanations that might account for the same outcomes observed in your study? NA

   (e) Did you address potential biases or limitations in your theoretical framework? NA

   (f) Have you related your theoretical results to the existing literature in social science? NA

   (g) Did you discuss the implications of your theoretical results for policy, practice, or further research in the social science domain? NA

3. Additionally, if you are including theoretical proofs...

   (a) Did you state the full set of assumptions of all theoretical results? NA

   (b) Did you include complete proofs of all theoretical results? NA

4. Additionally, if you ran machine learning experiments...

   (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? NA

   (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? NA

   (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? NA

   (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? NA

   (e) Do you justify how the proposed evaluation is sufficient and appropriate to the claims made? NA

   (f) Do you discuss what is "the cost" of misclassification and fault (in)tolerance? NA

5. Additionally, if you are using existing assets (e.g., code, data, models) or curating/releasing new assets, **without compromising anonymity**...

   (a) If your work uses existing assets, did you cite the creators? NA

   (b) Did you mention the license of the assets? NA

   (c) Did you include any new assets in the supplemental material or as a URL? NA

   (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? NA

   (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? NA

   (f) If you are curating or releasing new datasets, did you discuss how you intend to make your datasets FAIR? NA

   (g) If you are curating or releasing new datasets, did you create a Datasheet for the Dataset? NA

6. Additionally, if you used crowdsourcing or conducted research with human subjects, **without compromising anonymity**...

   (a) Did you include the full text of instructions given to participants and screenshots? Yes, details the Appendix.

   (b) Did you describe any potential participant risks, with mentions of Institutional Review Board (IRB) approvals? No, the study was granted an Exempt status, and posed no risks to participants.

(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? NA

(d) Did you discuss how data is stored, shared, and deidentified? NA

# A    Prompts and Marketplaces

## A.1    Prompt Marketplaces

The integration of `txt2img` models in digital art creation has led to the development of AI prompt marketplaces, such as PromptHero[6], Promptrr.io[7], Prompti AI[8], PromptBase[9], and CivitAI[10]. These platforms facilitate the buying, selling, or sharing of prompts designed for various `txt2img` models (Figure 15). Prompt marketplaces operate on a business model in which users can create and submit their own prompts or purchase those created by others. This has turned prompt creation into a profitable activity, as effective prompts greatly enhance the quality of AI-generated art (Lanz (web)). Users selling prompts on these platforms can earn income, while buyers gain access to a diverse range of ready-to-use prompts, enhancing their productivity and creativity in AI art generation. These marketplaces typically use a commission-based revenue model, where prompt authors pay the platform a part of their revenue whenever the platform sells their prompt to a customer (Timothy (web)). These platforms also generally treat prompts as intellectual property (PromptBase; Promptrr.io), recognizing and protecting the creative effort involved in their creation.
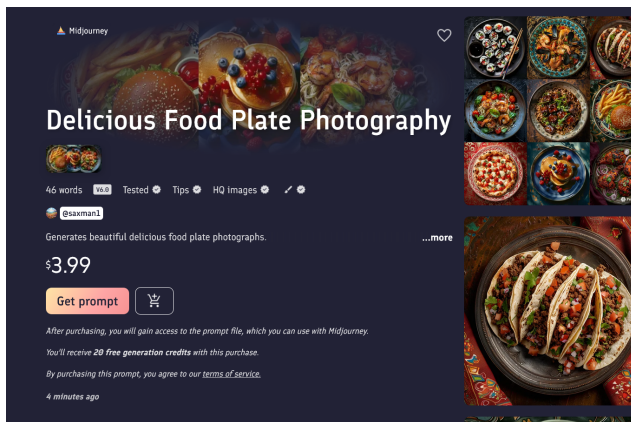


Figure 15: A PromptBase listing selling the prompt to a specific style of artwork.

## A.2    Subjects and Modifiers

The use of *subjects* and *modifiers* in prompts plays an important role in guiding a `txt2img` model to produce images that align with the creator's vision. The subject serves as the core theme or focus of the image, while modifiers provide the specifications that refine and shape the final output. Subjects can range from concrete objects, like "cat," "forest" or "robot," to more abstract concepts, such as "solitude" or "chaos." On the other hand, the modifiers adjust the images's aesthetic by specifying attributes such as color, texture, time of day, or emotional ambiance, effectively guiding the AI towards a more targeted and refined artistic rendition. Figure 1 demonstrates an example of how adding different modifiers alongside a subject (cat) produces varied image generations. Prompts for `txt2img` models are typically limited by 300 to 400 input characters (depending on the model) on how many subjects and modifiers can be specified for each image generation.

# B    Analyzed `txt2img` Model Details

- **MidJourney v5.0** (MidJourney (web) from Midjourney Inc. is a closed-source `txt2img` model which is recognized for generating images that often exhibit a unique artistic and abstract quality. Due to its closed-source nature, not much is known about the architecture and parameters of its image generator model or the dataset used to train it. MidJourney v5.0 can natively generate images of 1024×1024 pixels, using a Discord bot with `/imagine [Prompt]` command.

- **Stable Diffusion XL** (SDXL) (Podell et al. 2023; Rombach et al. 2022), created by Stability AI, employs a variational autoencoder (VAE) in conjunction with a cross-attention transformer-based architecture as the generator. Stable Diffusion XL can also natively generate images of 1024×1024 pixels, using direct interface with the publicly available model (Hugging Face c).

- **DreamShaper XL** was trained and fine-tuned over Stable Diffusion XL (Podell et al. 2023), with a focus on generating photo-realistic fantasy images. Similar to Stable Diffusion XL, DreamShaper XL can also natively generate images of 1024×1024 pixels.

- **Realistic Vision v5** (Hugging Face b) was trained and fine-tuned over Stable Diffusion 1.5 (Hugging Face a), with a focus on generating photo-realistic images. However, because Stable Diffusion 1.5 was trained using 512×512 images, Realistic Vision 5 output is also limited to 512×512 pixels. Generating images at a resolution of 1024×1024 pixels with a model originally trained for 512×512 pixels often leads to the emergence of undesirable artifacts, as the diffusion models tend to replicate the quantity of subjects. But for an equitable comparison with other images in our study, we upscaled Realistic Vision v5 images with *Highres.fix* and ScuNET-PSNR upscaler to 1024×1024 pixels. *Highres.fix* is a generation technique where the first half of sampling iterations are done at native resolution, and then the later half at the targeted upscaled resolution. This can add more accurate details to the upscaled images while the generation prompt is still in context, without the aforementioned artifacts being introduced.

---

[6]https://prompthero.com
[7]https://promptrr.io
[8]https://prompti.ai
[9]https://promptbase.com
[10]https://civitai.com

## C   Pre-Survey Details and Results

---

**Age:**

☐ 18-24
☐ 25-34
☐ 35-44
☐ 45-54
☐ 55-64
☐ 65-74
☐ 75+

**Gender:**

☐ Male
☐ Female
☐ Other

**[For University Participants] Major/Department:**
----------------

**[For MTurk Participants] Occupation:**

☐ Professional/Managerial
☐ Skilled Trades
☐ Service Industry
☐ Sales and Marketing
☐ Healthcare and Medical
☐ Arts and Entertainment
☐ Education and Academia
☐ Manufacturing and Production
☐ Other

**Country:**
----------------

**Familiarity with AI generation tools:**
**Text Generation Tools:**

☐ Not at All Familiar
☐ Slightly Familiar
☐ Somewhat Familiar
☐ Very Familiar

**Image Generation Tools:**

☐ Not at All Familiar
☐ Slightly Familiar
☐ Somewhat Familiar
☐ Very Familiar

**Audio Generation Tools:**

☐ Not at All Familiar
☐ Slightly Familiar
☐ Somewhat Familiar
☐ Very Familiar

**Video Generation Tools:**

☐ Not at All Familiar
☐ Slightly Familiar
☐ Somewhat Familiar
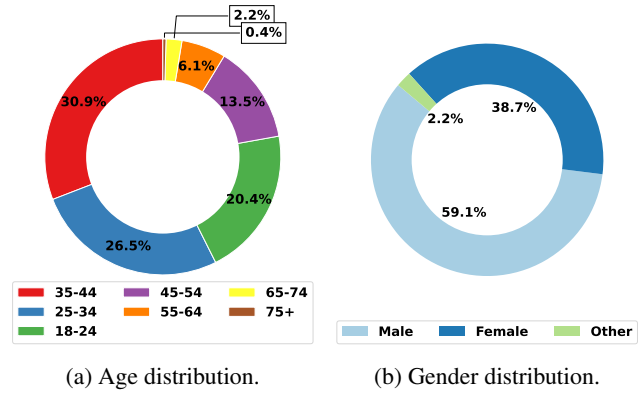☐ Very Familiar

Figure 16: Demographic survey form.

(a) Age distribution.

(b) Gender distribution.

Figure 17: Demographic distribution of survey participants.

## D   Details of the Metrics

### D.1   MSQ Scores (Survey Part I)

The scoring mechanism for evaluating multiple select questions (MSQs) in Part I with one or more correct answers on subject and modifiers pertaining to each displayed image is formalized through the following formula. The formula aims to give an MSQ score from 0 to 2 for each question, and is commonly used by learning management systems such as Canvas[11]. This metric serves to give an initial baseline for human prompt inference, as it only requires user to select the correct subject and modifiers from a set of given options.

$$MSQScore = \max\Big\{0, \min\Big\{T, \big(T \cdot \mathbf{1}_{\{S_c=C \text{ and } I=0\}}\big) + \frac{-0.1 \times I}{C}\Big\}\Big\}$$

where:

- $T = 2$ represents the total available points for a question, 1 point for correct subject selection, 1 point for correct modifier(s) selection.
- $S_c$ is the number of correctly selected options by the participant.
- $C$ is the total number of correct options for a question.
- $I$ is the number of incorrect options selected by the participant.
- $\mathbf{1}_{\{S_c=C \text{ and } I=0\}}$ is an indicator function that equals 1 if the participant selects all correct options without selecting any incorrect ones ($S_c = C$ and $I = 0$), and 0 otherwise.

This formula integrates the conditions for awarding full points and applying penalties for incorrect selections into a singular expression. The use of the *max* and *min* functions ensures that the final score remains within the acceptable range of 0 to $T = 2$ points. The indicator function facilitates the condition under which full points are awarded, while the penalty for incorrect selections is universally applied but is effectively neutralized when full points are granted due to correct selection criteria being met. This scoring framework was devised to offer a balanced assessment of participant knowledge and decision-making skills, reflecting both the breadth of correct understanding and the penalties for inaccuracies.

---

[11]https://www.instructure.com/canvas

## D.2 Image Hash

The `imagehash` library in Python[12] offers a compelling approach for measuring prompt inference accuracy, specifically by comparing the original AI-generated art with an art generated from an inferred prompt. `imagehash` generates perceptual hashes of images, where images that are similar result in hashes that are closely aligned. The similarity between these hashes is measured using the Hamming distance, a measure for comparing two binary data strings. By calculating the Hamming distance between the hashes of two images, we can quantitatively assess their similarity. A lower image hash score (difference) implies the two images are similar, and vice versa. For example, images in Figure 18a and Figure 18b have an image hash score of 26, where they were generated from the same prompt. In contrast, Figure 18b and Figure 18c have a slightly higher image hash score of 28 when generated from a slightly different prompt, and Figure 18b and Figure 18d have a significantly higher image hash score of 34 as they were generated from very different prompts.

## D.3 Perceptual Similarity

Perceptual similarity metric (Zhang et al. 2018) captures the degree to which two images are perceived to be similar by human observers. This approach goes beyond traditional pixel-based comparisons, which often fail to capture the nuances of human visual perception. Instead, perceptual similarity considers factors like texture, color distribution, structural elements, and contextual information, which are more aligned with how humans process visual information. This method can provide a more accurate and meaningful measure of similarity between the original AI-generated art with an art generated from an inferred prompt, particularly in applications where the visual impression of an image is more important than its exact pixel composition. Similar to image hash, a lower perceptual similarity implies the two images are similar, and vice versa. For example, images in Figure 18a and Figure 18b have perceptual similarity score of 0.526, where they were generated from the same prompt. In contrast, Figure 18a and Figure 18c have a slightly higher perceptual similarity score of 0.585 when generated from a slightly different prompt, and Figure 18a and Figure 18d
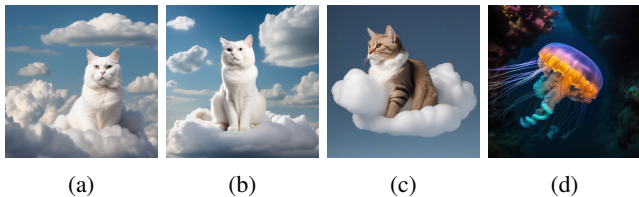
---

[12]https://pypi.org/project/ImageHash/



|       |       |       |       |
| :---: | :---: | :---: | :---: |
| (a)   | (b)   | (c)   | (d)   |

Figure 18: SDXL generations using the same prompt "*photo of a huge white cat sitting on a cloud in the sky*" for images (a) and (b), a similar but less specific prompt "*cat on cloud*" for image (c), and an entirely different prompt "*a glowing jellyfish underwater, breathtaking*" for image (d).

have a significantly higher perceptual similarity score of 0.793 as they were generated from very different prompts.

## D.4 Image Embedding Similarity (CLIP Score)

The CLIP score between an image and a prompt (text) measures the cosine similarity between their embeddings (Wang, Chan, and Loy 2023). A higher score indicates greater relevance or similarity, as perceived by the CLIP model, between the text and the image. This score is particularly useful in our objective to determine the accuracy or relevance of an image in relation to the textual input that was inferred to be used in its creation. We employ two state-of-the-art CLIP models in our evaluation, OpenAI's ViT-L/14 Transformer[13] (L14) and ViT-B/32 Transformer[14] (B32). More specifically, we utilize SentenceTransformers (Reimers and Gurevych 2019) wrapped L14 and B32 models for simultaneously calculating both the CLIP score and semantic similarity when the embeddings are generated, which is described next. Unlike image hash and perceptual similarity, a higher CLIP score implies the two images are similar, and vice versa. For example, images in Figure 18a and Figure 18b have a CLIP score of 0.970, where they were generated from the same prompt. In contrast, Figure 18a and Figure 18c have a lower CLIP score of 0.923 when generated from a slightly different prompt, and Figure 18a and Figure 18d have a significantly lower CLIP score of 0.629 as they were generated from very different prompts.

## D.5 Text Embedding Similarity (Semantic Similarity)

In `txt2img` generation, where prompts are used to guide the creation of visual content, semantic text similarity is also an important metric in evaluating how accurately an inference deduces or reconstructed the original text prompt from the generated images. We use SentenceTransformers (Reimers and Gurevych 2019), an extension of the BERT model for efficiently producing sentence embeddings, for calculating the cosine similarity between the original prompt and participants' response prompts. Similar to CLIP score, a higher semantic similarity implies the two prompts are similar, and vice versa. For example, the prompts "*photo of a huge white cat sitting on a cloud in the sky*" and "*cat on cloud*" used in Figure 18 have a high semantic similarity of 0.745, as measured by the L14 model. In contrast, the prompts "*photo of a huge white cat sitting on a cloud in the sky*" and "*a glowing jellyfish underwater, breathtaking*" have a semantic similarity of only 0.374.

**Surveyed Similarity Rating (from Survey Part IV).** As outlined in Section 3.6, participants in Part IV of the study evaluated pairs consisting of an original image and a newly generated image created from a prompt provided by a previous participant. They rated these pairs on a Likert scale consisting of "Not at All Similar", "Slightly Similar", "Somewhat Similar", and "Very Similar." This provides an additional measurement of image similarity, and thus prompt inference success, based on human interpretations.

---

[13]https://huggingface.co/openai/clip-vit-large-patch14
[14]https://huggingface.co/openai/clip-vit-base-patch32

# E   Survey Details

## E.1   Survey Part I Question Example



**Question 1: Select the best fitting subject:**
- ○ man
- ○ astronaut
- ○ robot
- ○ cat

**Question 2: Select the best modifiers (check all that apply):**
- ☐ Pop Art
- ☐ High Saturation
- ☐ Australian Outback
- ☐ Tone mapped
- ☐ Vibrant colors

Next

Figure 19: An example of a question for Part I. The participant must select the best fitting subject and between 1 to 5 best fitting modifiers.

## E.2   Survey Part II Question Example



**Question: Enter the best fitting prompt**

Enter the prompt

Next

Figure 20: An example of a question for Part II. The participant must type in a prompt for the given image.

## E.3   Survey Part III Question Example



**Question: Enter the best fitting prompt**

Enter the prompt

Next

Figure 21: An example of a question for Part III. The participant must type in a prompt for the given image.

## E.4   Survey Part IV Question Example



**How similar are the images?**
- ○ Not at All Similar
- ○ Slightly Similar
- ○ Somewhat Similar
- ○ Very Similar

Next

Figure 22: Example of a Part IV question, where participants rated the similarity between pairs of AI-generated images.

## E.5   121 Controlled Dataset Modifiers and Frequency of Use

abstract style (5), acrylic (3), aerial (2), african savanna (1), african savannah (1), american wild west (2), ancient chinese dynasty (4), ancient mayan (1), angry (2), arctic tundra (1), art deco (1), art deco style (1), art nouveau style (1), asymmetry (1), australian outback (2), backlit lighting (3), baroque (1), baroque style (1), bokeh lighting (3), bollywood inspired (4), bright colors (1), caribbean island (1), cartoon style (1), cel shading (4), charcoal (1), cinema render (2), cinematic lighting (1), close-up (5), cloudy lighting (2), cloudy melancholic (1), comic book (2), comic book

style (3), cool colors (3), dark (3), dark colors (5), dark dramatic (1), digital art (4), digital camera (4), dramatic (4), dramatic acrylic (1), dramatic lighting (5), dreamy lighting (5), dusk lighting (5), dystopian style (2), elegant (2), engraving (2), excited (4), flat shading (1), futurist (1), futurist style (1), golden hour lighting (3), gooch shading (2), gothic european castle (6), gothic style (3), gouraud shading (2), high contrast (5), high saturation (4), high saturation sad (1), hyper real (3), japanese garden (1), joyful (1), juxtaposed (2), light-hearted (2), low contrast (4), low saturation (1), maasai village (2), macro lens (5), magical (3), mediterranean seaside (1), melancholic (1), metallic colors (7), metallic noir (1), middle eastern bazaar (2), minimalist (1), minimalist style (2), moody (1), moody lighting (1), moonlight lighting (3), moonlit (1), mysterious (1), mystical (1), mystical style (1), natural light lighting (2), octane render (1), pastel colors (2), peaceful (2), pen and ink (2), pencil sketch (5), phong shading (3), photorealistic (3), pixel art (4), polynesian island (1), pop art (1), pop art style (1), ray traced (4), realistic (5), renaissance italy (2), renaissance style (2), rustic (1), rustic style (1), sad (6), serene (3), shallow depth of field (1), sketch (2), soft lighting (1), steampunk style (3), studio (1), studio lighting (3), sunrise lighting (3), surrealist style (2), symmetry (4), telephoto lens (3), tilt-shift lens (1), tone mapped (2), unreal render (1), vibrant colors (6), victorian england (2), warm colors (7), watercolor (4), whimsical (2), wide angle (1).

## E.6  100 Controlled Dataset Prompts

- astronaut, cool colors, digital art, serene
- astronaut, art deco
- astronaut, ancient chinese dynasty, high contrast, pencil sketch, dramatic
- astronaut, realistic, gothic european castle, symmetry, rustic
- astronaut, dreamy lighting, surrealist style, peaceful, metallic colors
- astronaut, dusk lighting, minimalist style, close-up, ancient chinese dynasty
- astronaut, comic book style, symmetry, hyper real, sunrise lighting
- astronaut, steampunk style, angry
- astronaut, phong shading, dark colors
- astronaut, dramatic lighting, renaissance style, warm colors, pen and ink
- astronaut, excited, phong shading, high saturation, abstract style
- astronaut, abstract style, symmetry, high saturation, photorealistic, arctic tundra
- astronaut, gothic style, telephoto lens, dusk lighting
- astronaut, watercolor, elegant, cool colors, natural light lighting
- astronaut, low contrast, ray traced, backlit lighting, close-up
- astronaut, vibrant colors, sunrise lighting, ancient chinese dynasty

- astronaut, bollywood inspired, dusk lighting, metallic colors
- astronaut, gothic style, dark, high contrast, studio lighting
- astronaut, cinematic lighting, watercolor
- cat, dreamy lighting
- cat, dark dramatic, futurist
- cat, cel shading, australian outback
- cat, metallic noir, realistic, gothic european castle
- cat, pixel art, dark colors
- cat, dramatic lighting
- cat, digital camera, melancholic, moonlight lighting, comic book style, ray traced
- cat, metallic colors, golden hour lighting
- cat, moody lighting, minimalist style, acrylic, macro lens, cel shading
- cat, digital camera, serene, macro lens, cool colors, african savannah
- cat, ancient chinese dynasty, bokeh lighting, gooch shading, elegant
- cat, charcoal, tone mapped, australian outback, light-hearted, art deco style
- cat, dreamy lighting
- cat, pencil sketch, middle eastern bazaar, vibrant colors
- cat, dystopian style
- cat, magical, vibrant colors
- cat, warm colors, dramatic lighting
- cat, warm colors, tone mapped, engraving
- cat, octane render, soft lighting
- cat, telephoto lens, abstract style
- man, baroque, comic book, minimalist
- man, flat shading, maasai village, peaceful
- man, bollywood inspired, acrylic, comic book, dark
- man, dramatic, pastel colors, gothic european castle, photorealistic
- man, realistic, dramatic lighting, excited
- man, steampunk style, cinema render, digital art, metallic colors
- man, digital camera, excited, renaissance italy
- man, vibrant colors
- man, metallic colors, whimsical, digital camera, golden hour lighting
- man, serene
- man, studio lighting, american wild west, unreal render
- man, cartoon style, gouraud shading, natural light lighting, dark colors, digital art
- man, ray traced
- man, pop art style, phong shading
- man, magical, watercolor, bokeh lighting, middle eastern bazaar
- man, macro lens, high saturation sad

- man, dark colors, caribbean island, dark, bokeh lighting, aerial
- man, sunrise lighting
- man, warm colors
- man, dramatic, cinema render, warm colors, backlit lighting
- robot, american wild west, dramatic, metallic colors
- robot, gothic style
- robot, moody, victorian england, sad
- robot, pop art, high saturation
- robot, high contrast, angry
- robot, vibrant colors
- robot, low contrast, realistic
- robot, cel shading
- robot, symmetry, warm colors, ray traced, sad, dusk lighting
- robot, low contrast, photorealistic, magical, art nouveau style, juxtaposed
- robot, studio lighting, surrealist style
- robot, dystopian style
- robot, dark colors
- robot, dreamy lighting, pixel art, close-up, bollywood inspired, baroque style
- robot, asymmetry, mystical style, cloudy lighting
- robot, watercolor, futurist style, maasai village, golden hour lighting, whimsical
- robot, sad, pixel art
- robot, sad, low contrast, gothic european castle, sketch
- robot, renaissance italy
- robot, wide angle, vibrant colors, cel shading, sad
- robot, close-up, pastel colors, ancient mayan, abstract style
- woman, hyper real, pencil sketch, aerial
- woman, high contrast, pencil sketch, realistic, japanese garden
- woman, victorian england, juxtaposed, cloudy melancholic
- woman, light-hearted, mystical
- woman, macro lens, studio, moonlit, hyper real, gothic european castle
- woman, rustic style, acrylic, cloudy lighting, bright colors
- woman, dramatic lighting, telephoto lens
- woman, steampunk style, high saturation, close-up, mediterranean seaside
- woman, renaissance style, moonlight lighting, warm colors, dramatic acrylic
- woman, engraving, backlit lighting, joyful
- woman, dusk lighting, pixel art, low saturation
- woman, gooch shading, moonlight lighting, tilt-shift lens, abstract style, mysterious
- woman, bollywood inspired, comic book style

- woman, polynesian island
- woman, shallow depth of field, pen and ink, sad, metallic colors
- woman, pencil sketch, gouraud shading
- woman, excited, african savanna, macro lens
- woman, gothic european castle, dreamy lighting, digital art
- woman, sketch
- woman, high contrast

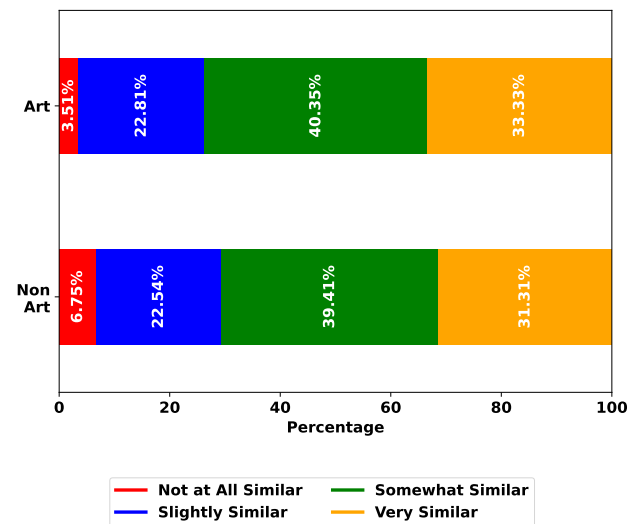## F  Additional Results



Figure 23: Part IV similarity rating percentage from participants with and without art background.
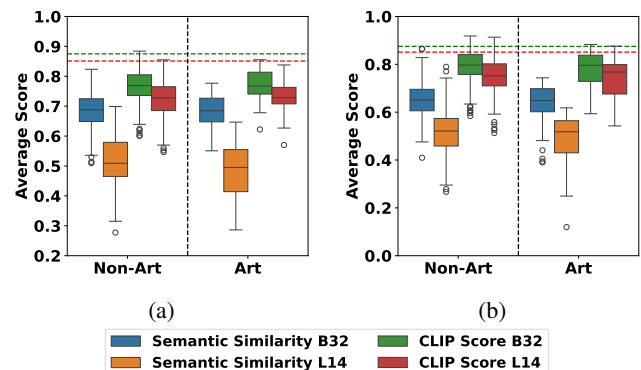


Figure 24: Semantic similarity and CLIP score for participants with and without art background, in (a) controlled dataset, and (b) uncontrolled dataset.
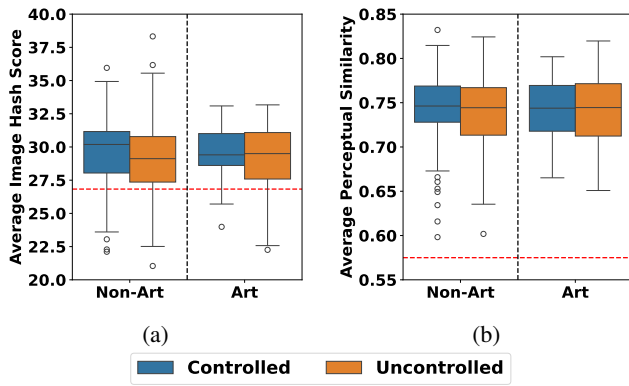
Figure 25: Image hash score and perceptual similarity for participants with and without art background, in (a) controlled dataset, and (b) uncontrolled dataset.
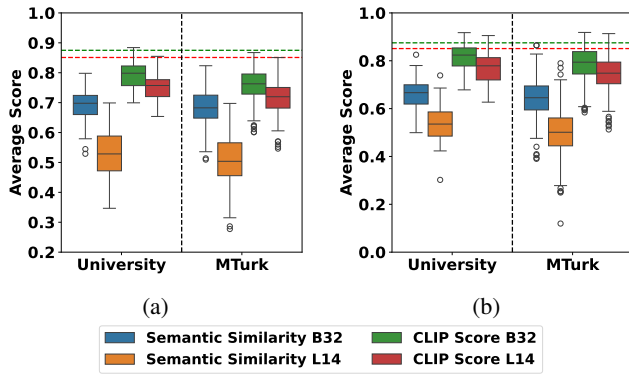


Figure 26: Semantic similarity and CLIP score for participants recruited on campus vs. on MTurk, in (a) controlled dataset, and (b) uncontrolled dataset.
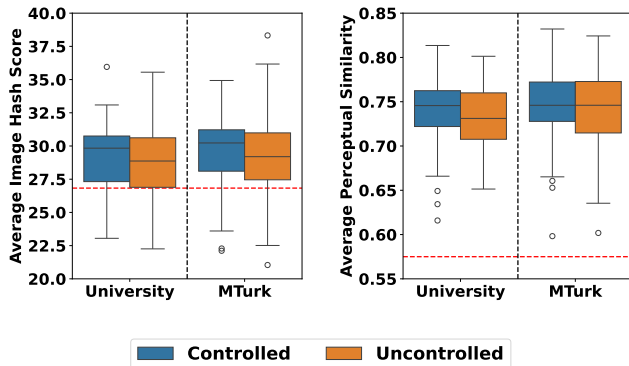


Figure 27: Image hash score and perceptual similarity for participants recruited on campus vs. on MTurk, in (a) controlled dataset, and (b) uncontrolled dataset.
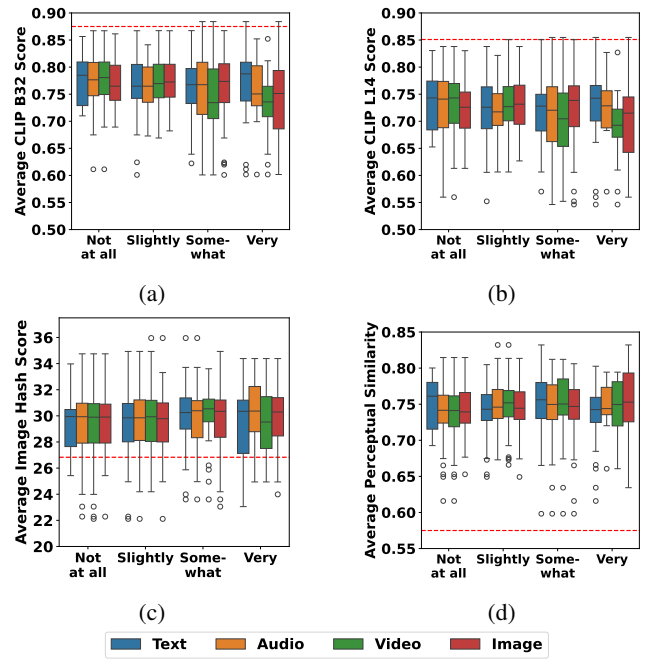


Figure 28: Prompt inference accuracy measured for participants with different levels of AI tools familiarity, using (a,b) B32 and L14 CLIP scores, respectively, (c) image hash, and (d) perceptual similarity, in the controlled dataset.
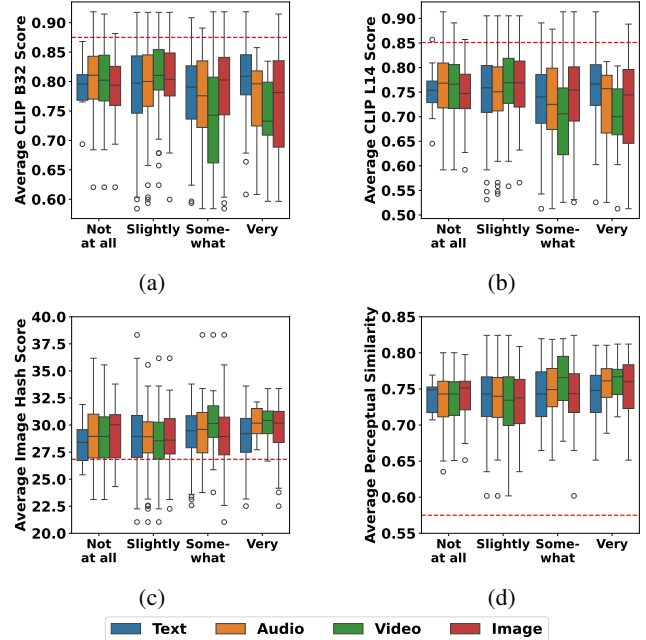


Figure 29: Prompt inference accuracy measured for participants with different levels of AI tools familiarity, using (a,b) B32 and L14 CLIP scores, respectively, (c) image hash, and (d) perceptual similarity, in the uncontrolled dataset.