

GenAssist: Making Image Generation Accessible

Mina Huh
The University of Texas at Austin
Austin, TX, USA
minahuh@cs.utexas.edu

Yi-Hao Peng
Carnegie Mellon University
Pittsburgh, PA, USA
yihaop@cs.cmu.edu

Amy Pavel
The University of Texas at Austin
Austin, TX, USA
apavel@cs.utexas.edu

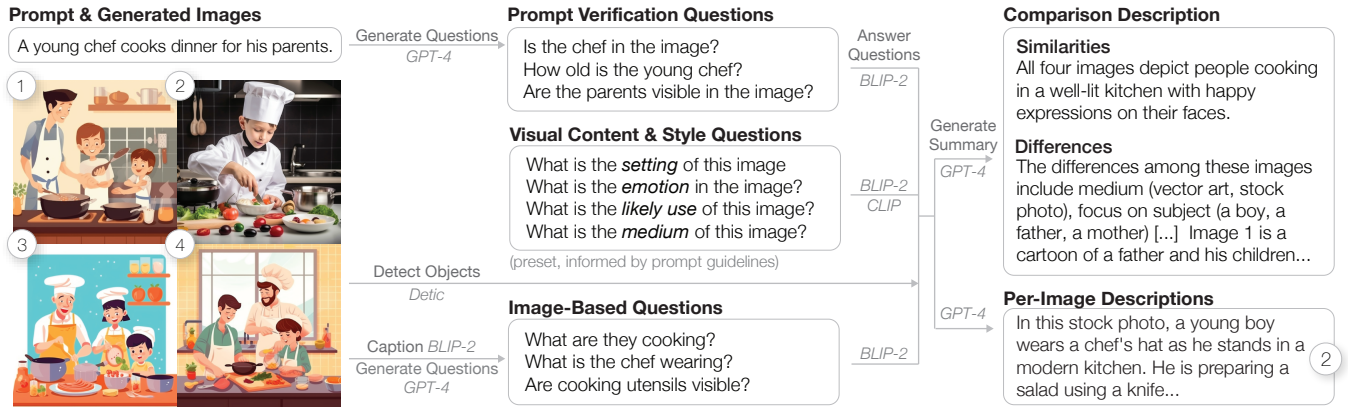


Figure 1: GenAssist makes image generation accessible by providing rich visual descriptions of image generation results. Given a text prompt and set of generated images, GenAssist uses a large language model (GPT-4) to generate *prompt verification questions* from the prompt and *image-based questions* from the image captions. GenAssist then answers the visual questions (BLIP-2) and uses a vision-language model (CLIP) and an object detection model (Detic) to extract additional visual information. GenAssist then uses GPT-4 to summarize all of the information into *comparison descriptions* and *per-image descriptions*.

ABSTRACT

Blind and low vision (BLV) creators use images to communicate with sighted audiences. However, creating or retrieving images is challenging for BLV creators as it is difficult to use authoring tools or assess image search results. Thus, creators limit the types of images they create or recruit sighted collaborators. While text-to-image generation models let creators generate high-fidelity images based on a text description (i.e. prompt), it is difficult to assess the content and quality of generated images. We present GenAssist, a system to make text-to-image generation accessible. Using our interface, creators can verify whether generated image candidates followed the prompt, access additional details in the image not specified in the prompt, and skim a summary of similarities and differences between image candidates. To power the interface, GenAssist uses a large language model to generate visual questions, vision-language models to extract answers, and a large language model to summarize the results. Our study with 12 BLV creators demonstrated that GenAssist enables and simplifies the process of image selection and generation, making visual authoring more accessible to all.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/XXXXXXXXXXXX>

1 INTRODUCTION

BLV creators use images in presentations [52], social media [5], videos [24], and art [8]. To obtain images, creators currently either describe their desired images to the sighted collaborators who then search for or create the image [52, 75], or limit the types of images they create [61]. Large-scale text-to-image generation models, such as DALL-E [58], Stable Diffusion [60], and Midjourney [41], present an opportunity for these creators to generate images directly from text descriptions (i.e., prompts). However, current text-to-image generation tools are inaccessible to BLV creators, as creators must *visually inspect* the content and quality of the generated images to iteratively refine their prompt and select from multiple generated candidate images.

While BLV creators can gain access to images using automated descriptions [34, 40], existing descriptions are intended primarily for image consumption. As a result, the descriptions leave out details that may help authors decide whether or not to use the image (e.g., style, lighting, colors, objects, emotions). Prior work also enables users to gain flexible access to the spatial layout of objects in images [32], but exploring details per image makes it difficult to assess similarities and differences between image options provided during image generation. To make authoring visuals more accessible, prior work has explored describing visuals to help creators author presentations [52] or videos [24]. While such work helps creators identify low-quality visuals (e.g., blurry footage in a video [24]) or graphic design changes (e.g., changing slide layouts [52]), prior work has not yet explored how to improve the accessibility of image generation.

To understand the opportunities and challenges of text-to-image generation, we conducted a formative study with 8 BLV creators who regularly create or search for images. Creators in our study reported their existing strategies for making images themselves (e.g., using SVG editors or code), searching for images, or asking others to search for or create images (similar to prior work [5, 24, 52]). All creators expressed excitement about using image generation to improve their efficiency and expressivity in image authoring. Creators all used image generation for the first time during our study and enjoyed creating high-fidelity images for their own uses (e.g., creating a logo for their website, making a card for their family). While we invited participants to ask the researchers visual questions to gain access to the visual details (e.g. “*What are the differences?*”, “*Is the color calm or aggressive?*”), it remained challenging for participants to: craft a well-specified prompt especially without visual experience, assess how well the generated image followed the prompt, recognize generated details that were not originally specified in the prompt, and understand or remember the similarities and differences between images.

To improve the accessibility of image generation, we present GenAssist, a system that provides access to text-to-image generation results via prompt-guided image descriptions and comparisons (Figure 1). Our system lets creators skim an overview of similarities and differences between images using our comparison descriptions and per image descriptions (Figure 1, right), assess if the images followed their prompt using prompt verification (Figure 1, center), and recognize visual details not in the prompt using our content and style extraction (Figure 1, center). Creators can also interactively ask questions across multiple images to gain additional details. Our interface design enables creators to easily navigate visual information via a screen reader-accessible table format. Our tables let creators selectively gain information about individual images (columns) or visual questions (rows) (Figure 4).

We evaluated GenAssist in a within-subjects study with 12 BLV creators who compared GenAssist with a baseline interface that was designed to encompass practices of accessing images (e.g., automated caption [77], object detection [40], and Visual Question Answering [34]). Participants rated GenAssist as more useful than the baseline interface for understanding similarities and differences between the images, and they reported higher satisfaction with their image generation performance. Participants all expressed excitement about using GenAssist in their own workflows for authoring images and for new uses.

Our work contributes:

- Design opportunities making image generation accessible, derived from a formative study
- GenAssist, a system that provides access to image generation results via prompt-guided summaries and descriptions
- User study that demonstrates how BLV creators use GenAssist to interpret and generate images

2 BACKGROUND

As we aim to enhance the experience of content BLV creators working with AI-powered image-generation tools, our work builds upon prior research that explores: the accessibility of authoring tools and images, and text-to-image generation tools.

2.1 Accessibility of Authoring Tools

Enabling access to authoring tools unlocks new forms of self-expression. Recent research investigated how BLV people take and edit photos and videos [5, 24], compose music [48], draw digital images [8], and make presentations [52, 61, 83]. Such work includes studies of current practices that highlight accessibility concerns of existing authoring tools and the authored visuals. For example, features of current authoring tools remain difficult to access using screen readers [24, 35, 51], and it can be difficult to assess the effect of the visual edits such as color changes [61].

To improve the accessibility of authoring tools, researchers have explored methods for providing feedback to authors as they modify visual elements. For example, prior work has developed tactile devices that assist BLV designers in understanding and adjusting the layout of user interface elements [33, 53]. Tactile feedback has also been used to help developers interpret code structure, such as indentation [15]. Other prior work has used audio notifications to inform users about scene changes when reviewing videos [24, 49], while text descriptions have been used to convey visual details important to authoring such as brightness and layout [24, 52]. Sound and text feedback have also been used to keep blind authors informed about their collaborators’ edits to documents [30]. Similar to prior research, we also aim to make authoring tools accessible by providing in-situ feedback, but we instead provide creation-specific information to facilitate authoring images.

In addition to offering authoring feedback, researchers have developed systems to automate visual authoring. Prior systems recommend 2D layouts for visual elements during graphic design [45] and transform text into visual presentations [29, 63, 79]. To accommodate individual preferences and mitigate the impact of errors produced during generation, these systems typically offer multiple options for users to choose from and allow iterative generation attempts. Iterative generation and selection are not accessible for BLV creators, as it requires visually inspecting the output designs to choose a generated option or revise the input. In this work, we seek to make automated authoring tools, such as image generation, more accessible to BLV creators. Our approach provides a structured format for assessing and comparing generated results, and on-demand access to additional visual details to support creators in selecting a result and revising their input.

2.2 Accessibility of Images

Improving the accessibility of image generation systems involves not only ensuring access to image generation features, but also making the produced images accessible. A primary method for making images more accessible is representing them as text descriptions, such as image captions or alt text (e.g., “*A person walking on the street*”). Early work hired crowd workers to create alt text [6, 72], while recent research has developed machine-learning-based systems that automatically generate image descriptions [34, 71, 81]. Building on auto-generated captions, researchers have developed systems that further improve users’ understanding of images by providing additional information, such as regional descriptions [40, 84], and structuring detailed descriptions into an overview [14, 32, 43]. This approach enables users to review visual information more efficiently and has been found to help blind people better understand

images compared to using captions alone [31]. Our work builds upon this idea by presenting descriptions of image generation results in a hierarchical, easy-to-compare format, and tailoring the descriptions to the task of authoring rather than consuming images.

Automatic descriptions do not always capture all of the important image details. Visual Question Answering (VQA) tools can fill this gap by offering on-demand information to visual questions (e.g., “What is the person walking on the street wearing?”). Previous research has explored what visual questions blind people would like to have answered [9] and provided on-demand visual question answering support using both crowdsourcing [6, 25] and automated methods [17]. While VQA provides control over visual information gathering, it takes effort to ask individual questions. We investigate what types of visual questions BLV creators ask to create images during our formative study (similar to Brady et al. [9]), then use VQA to extract visual information and summarize this information as image descriptions. Thus, we explore how VQA and image descriptions work together as interconnected rather than separate accessibility solutions.

2.3 Text-to-Image Generation Tools

In recent years, significant progress has been made in the field of generative image models, particularly text-to-image models. These models employ pre-trained vision-language models to encode text input into guiding vectors for image generation, allowing users to create images using text prompts. This advancement can be attributed to various factors, including innovations in deep learning architectures (e.g., Variational Autoencoders (VAEs) [26] and Generative Adversarial Networks (GANs) [16]), novel training paradigms like masked modeling for language and vision tasks [10, 12, 13, 70], and the availability of large-scale image-text datasets [62]. With these advancements, recent diffusion-based models like DALL-E 2 [57], Stable Diffusion [60], and Midjourney [41] have successfully demonstrated the ability to synthesize high-quality images in versatile styles, including photorealism. This opens up potential practical applications for the content production industry [37]. However, none of the image generation tools provide text descriptions of the output so they are not accessible to BLV creators. In this work, we chose to use Midjourney due to its popularity among designers and content creators for its high-quality results. Midjourney enables creators to generate 4 candidate images for a single text prompt via a text-based interface hosted on Discord. However, our approach is not limited to any particular model, as we focus on comparing and describing multiple generated results from a single prompt, helping creators select the ideal image from various candidates produced by image generation tools.

With the development of these models, recent works have conducted studies to understand the relationship between content creators and AI generative tools, introducing design guidelines for such systems [28, 36]. These guidelines emphasize the need for more user controllability. Researchers have thus developed various tools to help designers better make use of generative AI, including assistance in exploring and writing better prompts [36, 76], recommending potential illustrations for news articles [37], and supporting collaboration between writers and artists [27]. While these studies offer valuable insights into how designers interact

with generative models, none have focused on creators with disabilities. Given the potential of text-to-image models for BLV creators, our work is the first to explore how to increase inclusivity in the expressiveness of image generation tools and make this emerging authoring approach more broadly accessible.

3 FORMATIVE STUDY

To understand the strategies and challenges of authoring and searching for images, we conducted a formative study with BLV creators. The formative study consisted of a semi-structured interview to investigate current strategies and challenges of obtaining images, and two image generation tasks to explore current strategies and challenges of using text-to-image generation.

3.1 Method

We recruited eight BLV creators who create or use visual assets on a regular basis (P1-P8, Table 4). Participants were recruited using mailing lists and compensated 50 USD for the 1.5-hour remote study conducted via Zoom¹. Participants described their vision as totally blind (6 participants) or legally blind (2 participants) with light and color perception. All participants had previously produced or selected images for their work and they represented a variety of professions: teacher (English, Music), professor (Computer Science, Climate), software engineer, graduate student, and artist. While 7 participants reported prior knowledge of text-to-image generation models, none had previously utilized such tools in their work.

We first conducted a semi-structured interview asking participants how they currently created or used visual assets, and the accessibility barriers that they encountered with their current approaches. We then provided a short tutorial on text-to-image generation and provided participants with Midjourney’s guidelines for creating text prompts [42] and example prompts from a Midjourney dataset [69]. Participants then completed two image generation tasks: a *guided task* in which participants generated a cover image for a news article on healthy eating [44] given the article’s title and full text, and a *freeform task* in which participants generated their own image. We gave participants 20 minutes to complete each image generation task. To limit onboarding time, we asked participants to email us their prompt (text and/or image) instead of using Midjourney’s Discord-based interface, then we shared the four generated candidate images back to the participants. We encouraged participants to ask questions about the four candidate images to decide whether to select one or change the prompt.

We recorded and transcribed the formative studies. To analyze the types of visual questions asked in the image generation task, two of the researchers labeled questions based on their goals and the types of information asked.²

3.2 Findings

Current Practice. Participants reported that they currently use images for a variety of contexts including slides, website images, paintings for commission, cartoons, scientific diagrams, and music album covers (Table 4). Five participants noted that they created

¹This study was approved by our institution’s Institutional Review Board (IRB).

²See Supplemental Material for the full list of prompts, images, and visual questions of the formative study

images on their own using image creation software such as SVG editors, slides, photoshop, and ProCreate (P7, P1, P5, P6), code packages including Python and Latex (P4, P5), or by taking photos (P3). Among them, three participants asked sighted people to review them (P3, P4, P6), and two participants reviewed the images using accessibility tools (e.g., audioScreen, tactile graphs, ZoomText) (P7, P3). Five participants searched for images online (P7, P8, P2, P3, P5), and three participants recruited another person to create or search the images for them (P7, P4, P5).

All participants who searched for images mentioned that they ask sighted people to describe the images for them in addition to reading any available alt text. P7 noted *“Alt text has never been helpful. It’s too short without important details.”* P8 and P5 mentioned that while a few established websites (e.g., New York Times, NASA) have good alt text, Google Image Search returns options other than established websites and *“it is hard to compare the results of the image search”* (P5). Participants also noted barriers to asking others to describe the image search results including finding available people to describe the images and avoiding false perceptions: *“I only ask a handful of people because it might lead to some subconscious bias that I’m not independent, cause it’s a basic task”* (P7).

Generating Prompts. All prompts written by participants specified the content they wanted to appear in the image (e.g., P6 used the prompt *“A person pushing a grocery cart down a produce aisle.”*), and only two participants specified the style of the image (P1 and P7 specified *“a photograph of...”*). Participants mentioned several challenges of creating prompts. First, while prompt guidelines [42] recommend users to specify multiple attributes in their prompt (e.g., style, lighting), participants reported that they were unfamiliar with visual attributes (*“I’m trying not to leave much to system randomness, I want to detail more things. But I don’t know a lot about different styles.”* — P5) and others found it difficult to remember what to mention in the prompt: *“I want the model to behave more like a wizard – asking me a series of questions ‘What do you want to create?’, ‘What style?’ and so on. It is hard to create detailed prompts in one attempt”* (P2). Participants also noticed that it is challenging to create a prompt that AI would be capable of generating: *“If I pin down something really specific or narrow [in the prompt], AI seems to break down”* (P1). P5 mentioned that transparency could inform prompt iteration: *“I want to know how the model works! [...] then I will know how to write a good prompt.”* Finally, while participants easily generated prompts during the free-form task motivated by their own creation goals, they mentioned it was challenging to know what content would effectively convey the article in the guided task: *“I have no experience reading a news article with images, so it’s hard to think of one. What do these images usually contain?”* (P7).

Understanding Image Candidates with Visual Questions. After generating images, participants asked visual questions to understand and select the images. Participants asked a total of 89 questions (47 asked in the guided task, 42 in the freeform task). The goals of the questions asked were to check whether the generated images followed the prompt (51), compare two or more images (34), request clarification of the answer provided by the interviewer (3), or understand a single image (1). The type of visual information asked by participants also varied. Participants asked about medium

(5), settings (6), object presences (18), object types (11), position attributes (11), color/light/perspective (16), and others (22).

Participants typically started by asking general questions, narrowing down to more specific questions as they ruled out images. For example, P4 progressively asked: *“Can you describe the images?”*, *“What are the differences between the four images?”*, *“What are the differences between the [store] isles?”*, *“Is the second image realistic?”*. Alternatively, participants started their questioning by directly checking if the image followed their prompts, such as in P5’s first question: *“Do we actually get the woman sitting at a desk?”* Finally, P1 and P2 started with questions about the style of the images: *“Is it realistic or cartoony?”* (P1) and *“Is the color calm or aggressive?”* (P2). Through asking questions, participants realized differences between their prompt and the generated images: *“it seems like the model generator is filling in details according to the context, even if I didn’t specify some details. I didn’t specify the clothes but in all images, the women are wearing office clothes”* (P5). Participants then asked follow-up questions based on new details. While the visual questions revealed the content and structure of what participants wanted to know about the images, participants reported that asking questions for each image was *“very time-consuming and confusing”* (P4). 5 participants noted that they would prefer to receive descriptions before asking questions, and participants reported that remembering all of the answers was difficult, as P2 summarized: *“I wish there were more description provided in the first place. I don’t know what to ask. Also, it’s hard to remember all the answers for each image.”*

Selecting an Image Candidate. While participants initially asked questions based on their prompt, they ultimately selected the final image considering both prompt-based descriptions and descriptions of extra details produced by the model. P7 suggested that information on whether the prompt is reflected in each image should be presented early so that he can decide whether to explore the image in detail or skip to the next candidate. P8 highlighted the importance of additional details: *“The model has randomness. It showed items I didn’t ask for and didn’t show what I asked for in the prompt. I want much information to be surfaced so that I can make a decision. Whether that unexpected parts can be still used.”* We also observed that similarities between images guided participants in deciding whether to further explore the images or to refine the prompt. For instance, after P3 generated images using a prompt *“A photo looking down on a kitchen table with a plate of pizza, a plate of fried chicken, and a bowl of ice cream on it.”*, he realized that all four images did not display drinks and iterated the prompt to explicitly mention *“fizzy drinks”*. On the other hand, differences between the images ultimately informed the final selection, as participants cited unique backgrounds, objects, and mediums as reasons for selecting the image (e.g., *P3 selected the final image because that was the only image that presented a dog putting his paw on the books.*).

Uses of Image Generation. When participants generated their own images in the free-form task, participants created a variety of images ranging from logos, art, website decorative images, presentations, and music album cover. All participants expressed excitement about using the text-to-image model as part of their image creation process in the future. Participants mentioned with image generation, they can create new types of images they had not created

before. P6 mentioned “With SVG editor, I cannot make realistic images. But now I can!” Also, participants mentioned that the quick creation will lead them to use images more often: “Because it’s so quick, I will use it for communication. Similar to how sighted people draw on a whiteboard during a Zoom meeting, I can quickly generate an image because representing a concept visually is easier for sighted team members.” (P8). P4 also compared the experience of image generation with image search “This simplifies things when I’m looking for things very niche, something that is hard to find online.” Finally, participants also mentioned the benefit of creating images alone. P7 said that because there is no need to ask a sighted person to help search images, it brings more autonomy and privacy. Participants also noted limitations and potential downsides of image generation including potential bias (P8, P4), copyright and training data concerns (P3, P4), wanting to use it only for inspiration (P1), and potential errors (P8). However, P8 expressed that he expected future models to produce fewer errors.

3.3 Reflection

Creators in our formative study currently employ resourceful strategies for creating or searching for images, but all creators expressed excitement to use image generation in their workflow. To improve access to image generation, our formative study reveals design opportunities (D1-D5) to make image generation accessible through technical or social support for:

- D1. Authoring prompts that specify content and style.
- D2. Understanding high-level image similarities and differences.
- D3. Assessing if images followed the prompt.
- D4. Accessing image details not specified by the prompt.
- D5. Organizing responses to visual questions.

These design opportunities address key user tasks in accessible text-to-image generation: generating the prompt (D1), understanding and selecting images (D2, D3, D4, D5), and revising the prompt for iteration (D4). Our work aims to help creators understand their image generation results through prompt-guided descriptions and comparisons (D2-D5). While providing high-quality descriptions may help creators improve their future prompts (D1), future work should explore how to actively support creators in authoring prompts.

4 SYSTEM

We present GenAssist, a system that supports accessible image generation via prompt-guided image descriptions and comparisons (Figure 1). To illustrate GenAssist, we follow Vito, a professional blogger who uses a screen reader to author his articles. Vito recently wrote an article about the benefits of teaching children to cook, and he wants to add an image to the article to engage his sighted readers. He attempts to use image search to find a stock photo of “a young chef” but notices that many of the images are missing detailed captions and alt text, or feature adult chefs instead of children. He decides to create an image using text-to-image generation with the prompt “a young chef is cooking dinner for his parents”. The text-to-image generation model returns four candidates:



To decide whether to use one of these images or change his prompt, Vito enters his prompt and image results into GenAssist.

4.1 Prompt Verification

While the text-to-image model generates output images based on the prompt, the generated image often does not reflect the specifications in the prompt, especially if the prompt is long, complicated, or ambiguous [22]. To help users assess how well their generated images adhered to their prompt, GenAssist provides **prompt verification**. To perform prompt verification, we first use GPT-4 [46] to generate visual questions that verify each part of the prompt. We input the text instruction “Generate visual questions that verify whether each part of the prompt is correct. Number the questions.” followed by the user’s prompt. GPT-4 outputs a series of questions:

| Input | Prompt Verification Questions |
|--|---|
| Generate visual questions that verify whether each part of the prompt is correct. Number the questions. Prompt: A young chef is cooking dinner for his parents. | <ol style="list-style-type: none"> 1. Is there a chef in the image? 2. How old is the young chef? 3. Is the young chef cooking food? 4. Are the parents present in the image? |

We generate answers to the visual prompt verification questions for each of the four generated candidate images using the BLIP-2 model with the ViT-G Flan-T5-XXL setup [34]. For each generated image and prompt verification question, we instruct the BLIP-2 model with the starting sequence “Answer the given question. Don’t imagine any contents that are not in the image.” to reduce hallucinations with non-existent information:

| Prompt Verification Questions | Image Answers (BLIP-2) | | | |
|---------------------------------------|------------------------|-----------|-----------|-----------|
| Is there a chef in the image? | Yes | Yes | Yes | Yes |
| How old is the young chef? | Young kid | Young kid | Young kid | Young man |
| Is the young chef cooking food? | Yes | Yes | Yes | Yes |
| Are the parents present in the image? | Yes | No | Yes | Yes |
| | 1 | 2 | 3 | 4 |

To help users quickly find which images do or do not adhere to the prompt, we use GPT-4 to summarize the responses to each question using the following prompt: “Below are the answers of four similar images to one visual question. Write one sentence summary that captures the similarities and differences of these results. The summary should fit within 250 character limit”. When using GPT-4’s chat completion API, we set the role of the system as “You are a helpful assistant that is describing images for blind and low vision individuals.”. The temperature value was set to 0.8. The summaries either indicate that all images have the same answer (e.g., “All images have a chef in the image”), or they alert users to differences:

| Prompt Verification Questions | Prompt Verification Summary |
|---------------------------------------|---|
| Is there a chef in the image? | Three images depict a young kid, while Image 4 depicts a young man. |
| Are the parents present in the image? | Three images show parents present in the image, while Image 2 does not. |

To enable screen reader users to easily access the answers to each question, we present the prompt verification results as a table including the prompt verification questions (rows, with the question

in column #1), prompt verification summaries (column #2), and per-image prompt verification answers (columns #3-6) (Figure 4).

Using our prompt verification table, Vito reads the answers summaries to check if the images follow his prompt. He notices that the 4th image contains an older chef, so it does not apply to his article about teaching children how to cook. While Vito also realizes the 2nd image does not feature the chef’s parents, he keeps the image in consideration as it may still apply to his article.

4.2 Visual Content & Style Extraction

Generated image candidates often feature similarities or differences that are not present in the original prompt. For example, Vito’s prompt “*A young chef is cooking dinner for his parents*” does not specify the style such that the resulting images include three illustrations and one photo. To enable access to image content and style details that were not specified in the prompt, we extract the **visual content** and **visual style** of the generated image candidates. To surface content and style similarities and differences that are important for improving image generation prompts, we used text-to-image prompt guidelines [20, 42, 47] to inform our approach.

We first created a list of visual questions about the image based on existing prompt guidelines, i.e. *prompt guideline questions*. The prompt guideline questions consist of questions about the content of the image (subjects, setting, objects), the purpose of the image (emotion, likely use), the style of the image (medium, lighting, perspective, color), and an additional question about errors in the image to surface distortions in the generated images such as blurring or unnatural human body features (Table 1).

To answer our prompt guideline questions for each image, we answered 5 questions (setting, subjects, emotion, likely use, colors) using Visual Question Answering with BLIP-2, similar to our prompt verification approach:

| Content & Style Questions | Image Answers (BLIP-2) | | | |
|-------------------------------------|------------------------|---------------------------|----------------------------|------------------------|
| What is the setting of the image? | Kitchen | Kitchen | Kitchen | Kitchen |
| What are the subjects of the image? | Father and children | Chef, kitchen, vegetables | Father, mother and son | Father, mother and son |
| What is the emotion of the image? | Happy | Happy | Happy | Happy |
| Where would this image be used? | On a website | In a cookbook | A children’s cooking class | On a website |
| What are the main colors? | Brown, blue, yellow | Black, white, red, green | Blue and white | Red, yellow, green |
| | 1 | 2 | 3 | 4 |

For our objects question, we used Detic [85], a state-of-the-art object detection model, with an open detection vocabulary and a confidence threshold of 0.3 to enable users to access all objects:

| Content & Style Questions | Image Answers (Detic) | | | |
|------------------------------------|--------------------------------------|--|--|---|
| What are the objects in the image? | Spoon, pot, cup, tub, apron, bowl... | Spoon, sink, tomato, lettuce, hat, bowl... | Spoon, fork, knife, apple, sausage, plate... | Spoon, pot, window, flowerpot, plate, frog... |
| | 1 | 2 | 3 | 4 |

For the remaining questions covering medium, lighting, perspective, and errors, we answer the question for each image candidate by using CLIP [56] to determine the similarity between the image and a limited set of answer choices (similar to CLIP interrogator [19]). To provide answers that could inform future prompts, we curated our answer choices for medium, lighting, and perspective from Midjourney’s list of styles [20] and DALL-E’s prompt book [47]. To

| Category | Name | Question | Model |
|----------------|-------------|--|--------|
| Content | Setting | What is the setting of the image? | BLIP-2 |
| | Subjects | What are the subjects of the image? | BLIP-2 |
| | Objects | What are the objects in this image? | Detic |
| | Emotion | What is the emotion of the image? | BLIP-2 |
| | Usage | Where would this image likely be used? | BLIP-2 |
| Style & Errors | Medium | What is the medium of the image? | CLIP |
| | Lighting | What is the lighting in this image? | CLIP |
| | Perspective | What is the perspective of this image? | CLIP |
| | Colors | What are the main colors used in this image? | BLIP-2 |
| | Errors | What are the errors in this image? | CLIP |

Table 1: Our *prompt guideline questions* including the question category, question name, and question, along with the model we used to answer the question (BLIP-2 [34], CLIP [56], or Detic [85]).

address common image generation errors, we retrieved the answer choices for our errors question from prior work [18, 59]. We include the full list of answer choices in the Supplementary Material. For each question, GenAssist presents the top three answer choices with a similarity score between the answer choice embedding and the image embedding above a threshold of 0.18:

| Content & Style Questions | Image Answers (CLIP) | | | |
|---------------------------------------|----------------------------------|------------------|------------------|----------------------------------|
| What is the medium of the image? | Cartoon, storybook, illustration | A stock photo | Vector art | Cartoon, storybook, illustration |
| What is the lighting of the image? | Natural lighting | Natural lighting | Natural lighting | Natural lighting |
| What is the perspective of the image? | Medium shot | Centered shot | Medium shot | Medium shot |
| What are the errors in this image? | Poorly drawn hands | None | None | None |
| | 1 | 2 | 3 | 4 |

To inform creators about unfamiliar visual style types, GenAssist provides the definition and the usage for each answer choice for visual style questions (Medium, Lighting, Perspective) by generating the description with GPT-4 and the prompt “*Describe the definition and the usage of the following [QUESTION NAME] in one sentence: [STYLE NAME]*”. Similar to the prompt verification table, we present the prompt guideline results in a table format including the prompt guideline questions (rows, with the question in column #1), prompt guideline summaries (column #2), and per-image prompt guideline answers (columns #3-6). We further split the prompt guideline results into two tables to improve ease of navigation: the *visual content table* includes answers to the content and purpose questions, and the *visual style table* includes answers to the style and errors questions. Finally, users can ask their own questions at the bottom of either table and GenAssist adds a row to the table by generating the answer for each image using BLIP-2, and the summary of answers using GPT-4. Using the visual content table, Vito notices from the objects summary that Image 1 has more food items than Images 2-4. As the purpose of the article is partially to introduce children to more ingredients, he decides to remove Image 1 from consideration. Using the visual style table, Vito realizes that Image 2 is a photo, while the other images are illustrations. As Vito was initially searching for a photo, he notes he may want to further refine his prompt to get more photo results. Vito also wants to check if the images will match his blog which is primarily black and white, so he adds a question about the background color:

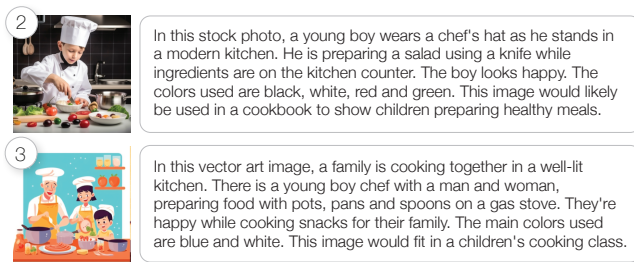


Figure 2: GenAssist's per-image descriptions.

| User Question | Image Answer Summary |
|-------------------------------|--|
| What color is the background? | Image 1 and Image 4 are light brown, Image 2 is black and Image 3 is blue. |

As Image 2 fits his article and includes a black background, he ranks Image 2 as his current top choice.

4.3 Description Summarization

To enable users to quickly assess their image results, we summarize the results from our pipeline to create a per-image description for each image and a summary of image similarities and differences.

To generate **per-image descriptions**, we first obtain the BLIP-2 caption for each image that provides a concise overview of the image content (e.g., “A family preparing food in the kitchen with a window.”). Then, we obtain additional detail about the image by generating questions about the caption with GPT-4 with the prompt: “Given the caption, generate 10 visual questions that are likely to be asked by blind and low vision individuals”. Unlike the other questions in our pipeline that are common across all images, this step enables the GenAssist to ask image-specific questions to add detail (e.g., “What is the view outside the window?” is only asked for Image 4). We generate the answers to these questions using BLIP-2.

We create individual image descriptions by first aggregating all information acquired in our pipeline for each image including the prompt verification, prompt guideline, and caption-detail question-answer pairs for each image. Then, we guide GPT-4 with the aggregated visual information and the prompt “Below is the information of an image. Write a description of this image for the blind and low vision audience. Describe the medium first. Your response should fit within 250 character limit. Do not add additional information that was not provided. Do not describe parts that are not clear or cannot be determined from the given information.” GPT-4 generates rich descriptions for each image (Figure 2).

To generate the **comparison description**, we simply provide all the information extracted from our pipeline to GPT-4 with the prompt “Below is the information for four images. Write one paragraph about the similarities between the four images and one paragraph about the differences between the four images. The summary should be concise.”. GPT-4 briefly summarizes the image similarities and differences (Figure 3). To help users quickly assess whether to revise their prompt or continue exploring, we present the **comparison description** and **per-image description** at the top of the page before the prompt verification and prompt guidelines tables.

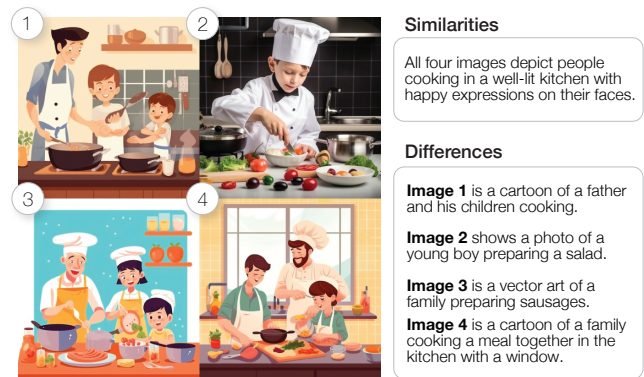


Figure 3: GenAssist's image comparison descriptions.

With the per-image description, Vito can quickly recall the content of Image 2 before making his final selection. With the comparison descriptions, Vito can quickly notice that Image 2 was the only image that contained a photo, then updated his prompt to get additional photos rather than illustrations.

4.4 Implementation

We implemented GenAssist using Gradio [1], an open-source Python library for the front-end web interface. The interface was deployed through Hugging face³ space with an NVIDIA A100 GPU (large, 40GB GPU Memory). Users' interaction logs were saved in the Firebase database. We followed the guidelines of W3C [73] and tested the compatibility of the GenAssist with all three major screen readers: NVDA, JAWS, and VoiceOver. GenAssist's tables follow the recommendations of W3C tables with two headers⁴.

5 PIPELINE EVALUATION

We measured the *coverage* of the descriptions generated by GenAssist and the *accuracy* of the information presented in GenAssist's tables. We compare the coverage of GenAssist-generated caption with the human-generated caption and the caption generated by a state-of-the-art image captioning model BLIP-2 [34].

5.1 Method

We selected 20 image sets (20 prompts x 4 generated images for each prompt = 80 total images) from Midjourney's community feed spanning different prompt lengths, content types, and styles. We recruited two people with experience describing images to provide descriptions for 10 randomly selected image sets each. For each image set, the describers provided descriptions of each individual image, and the similarities and differences between the images. We provided describers with prompt guidelines [42], image description guidelines [2], an example set of descriptions created by GenAssist, and the prompt for each image set to inform their descriptions. Both describers spent 3.5 hours to create descriptions for the 10 sets of images — or around 21 minutes per image set.

³<https://huggingface.co/spaces>

⁴<https://www.w3.org/WAI/tutorials/tables/two-headers/>

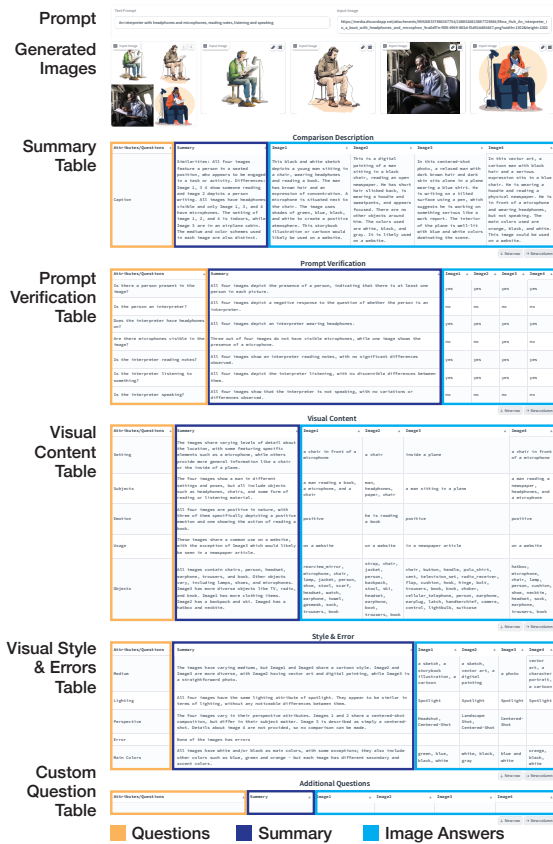


Figure 4: The GenAssist interface consists of screen reader accessible tables that enable users to flexibly gain more information about the content of interest.

We compared the coverage of GenAssist-generated descriptions to those generated by a baseline captioning tool (BLIP-2) and human describers. For comparison, we annotated the similarities and differences descriptions for all 20 sets of images and annotated the individual descriptions for 10 sets of images. We chose the 10 sets with the longest human descriptions to compare GenAssist with the highest quality descriptions. Because BLIP-2 cannot take multiple images as input to extract similarities and differences, we generated captions of the 4 images using BLIP-2, then prompted GPT-4 with the same prompt we used in our system to generate summary descriptions. We tallied whether the descriptions contained details about the image in each of our set of pre-defined visual information categories (Table 1). We counted only the correct information in the descriptions. One of the researchers annotated the descriptions and the other researcher reviewed the annotations. To compute the accuracy of the detailed visual information in GenAssist, one of the researchers examined the 20 sets of images with the three tables generated by the GenAssist (prompt verification table, visual content table, and visual style table) and counted the number of correct and incorrect answers in each table.

| Category | Sub-category | Correct (%) | Correct (#) |
|---------------------|--------------|-------------|-------------|
| Prompt verification | | 92.82 | 418 |
| Content | Setting | 97.53 | 81 |
| | Subjects | 98.60 | 143 |
| | Objects | 82.86 | 1243 |
| | Emotion | 87.5 | 80 |
| | Usage | 97.50 | 80 |
| Style | Medium | 82.76 | 174 |
| | Lighting | 94.33 | 141 |
| | Perspective | 71.83 | 142 |
| | Colors | 99.1 | 221 |
| | Errors | | 60.00 |

Table 2: We report the accuracy (percentage and number of correctly predicted information) of the pipeline results (Prompt verification, Content, Style, and Errors) with 20 sets of images.

5.2 Results

5.2.1 Coverage. We summarize our coverage evaluation results in Table 3. Overall, GenAssist’s comparison descriptions covered more similarities and differences than the human describers’. In the coverage of differences, GenAssist spotted more than twice the number of total differences than the human describers (4.55 vs. 2.25). The coverage of GenAssist’s individual image descriptions was comparable to that of human describers. When compared to human-generated description, GenAssist captured more information about the content and styles but revealed fewer image generation errors. For instance, one human describer specified in the comparison description “...All of the images have some AI generation error with fingers or clothing. ”. While GenAssist and the baseline used the same GPT-4 prompt to extract the similarities and differences, the baseline’s comparison description did not capture many differences.

5.2.2 Accuracy. Table 2 summarizes the results of the accuracy evaluation. Prompt verification, content, and style categories all achieved over 90% accuracy except for medium, perspective and emotion. In the 80 images in the dataset, GenAssist only detected five images as having errors, and detected the correct error types in three of them. The most common errors made in our pipeline were from perspective, medium, and error categories which are all extracted using the CLIP score. For perspective and medium, the majority of the errors were due to CLIP matching images to common style expressions (e.g., natural lighting, centered-shot) which likely reflects prevalence of these expressions in the training data. In the incorrect output of errors, GenAssist detected cartoon or sketch images as ‘poorly drawn faces’ errors. One reason for the relatively low accuracy of object detection results is that we empirically set the output threshold of GenAssist’s object detection (Detic) as 0.3 to present diverse objects to users in addition to information about the main subject extracted by BLIP-2 in our pipeline.

6 USER EVALUATION

We conducted a user study with 12 BLV visual content creators to compare GenAssist with a baseline interface.

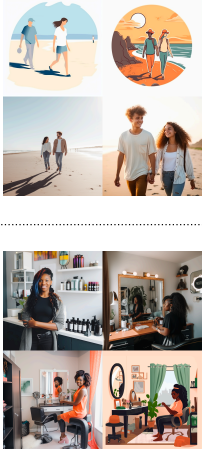
| | Human Description | Baseline Description | GenAssist Description |
|---|---|--|--|
|  | <p>All of the pictures are about a man and woman walking on a sunny beach</p> <p>The styles are different in cartoon and realistic. The haircut of the man and woman are different. The sun's appearance in frame are different.</p> | <p>All four images show people walking on the beach.</p> <p>Image 2 specifies the time of day as sunset. Image 4 emphasizes that the couple is holding hands. Image 3 specifies the age of the couple as young.</p> | <p>All four images depict a young couple walking on a sunny beach</p> <p>They differ in the type of medium, lighting, shot types, and colors used. The first three images mainly use blue and white, while the fourth image has a black and white theme. Image 1 and 3 show the couple walking on dry sand, while image 2 and 4 show the couple walking on the shore. The couple is holding hands in image 2, while in image 1 and 3, they are not. The images have different intended uses, ranging from a website or album cover to a poster for promoting beach trips.</p> |
| | <p>All feature a single black woman who could reasonably be hair stylists given their setting.</p> <p>Images that feature a mirror have varying degrees of realism; the second image may reasonably feature the same person, though the third image does not show the correct side of the subject.</p> | <p>All images involve women and mirrors.</p> <p>Image 1 takes place in a hair salon. Image 2 depicts a woman getting her hair done. Image 3 shows a woman sitting on a chair. Image 4 takes place in a woman's room.</p> | <p>All images feature a black woman with long hair in a positive and happy mood.</p> <p>Differences include setting, color scheme, and activity. Image 1 shows a hair stylist in an empty salon with blue and black hair, while Image 2 features a woman getting her hair styled in a mirror with black and white lighting. Image 3 also features a hair stylist with dark brown hair, but in a room with an orange dress and no tools. Image 4 is a digital illustration of a woman in her room with a plant, wearing a hoodie, and looking into a mirror hanging on a wall.</p> |

Figure 5: Two image sets and the descriptions of the similarities and differences used in the pipeline coverage evaluation (each image set described by a different human describer). GenAssist captured more information in the similarities and differences caption than the human describers.

| (Correct Only) | | Total Content (#) | | | Total Style (#) | | | Total Error (#) | | | Total (#) | | |
|------------------------|----------|-------------------|----------|-------------|-----------------|----------|-------------|-----------------|----------|-----------|-------------|----------|-------------|
| | | Human | Baseline | GenAssist | Human | Baseline | GenAssist | Human | Baseline | GenAssist | Human | Baseline | GenAssist |
| Similarities | μ | 1.5 | 1.65 | 2.45 | 0.70 | 0.00 | 0.80 | 0.10 | 0.00 | 0.00 | 2.35 | 1.65 | 3.25 |
| | σ | 0.61 | 0.59 | 1.10 | 0.80 | 0.00 | 0.83 | 0.31 | 0.00 | 0.00 | 0.83 | 0.85 | 1.29 |
| Differences | μ | 1.50 | 1.95 | 2.35 | 0.65 | 0.35 | 2.20 | 0.05 | 0.00 | 0.00 | 2.25 | 2.30 | 4.55 |
| | σ | 0.69 | 0.39 | 0.49 | 0.75 | 0.49 | 1.01 | 0.22 | 0.00 | 0.00 | 0.84 | 0.93 | 1.26 |
| Per-Image Descriptions | μ | 1.71 | 0.69 | 1.71 | 0.71 | 0.04 | 0.68 | 0.05 | 0.00 | 0.01 | 2.47 | 0.73 | 2.41 |
| | σ | 0.39 | 0.10 | 0.26 | 0.22 | 0.07 | 0.30 | 0.05 | 0.00 | 0.03 | 0.74 | 0.33 | 0.75 |

Table 3: We compared the coverage of GenAssist-generated descriptions to those generated by a baseline captioning tool and human describers. GenAssist captured more similarities and differences than the human describers.

6.1 Method

In a within-subjects study, participants used *GenAssist* and a *baseline* interface to interpret image generation results (*interpretation task*) and to generate images (*generation task*).

Participants. We recruited BLV creators who create or use visual assets on a regular basis using mailing lists (P7-P18, Table 4). Participants described their vision as totally or legally blind and they were students, consultants, software engineers, video creators, and artists. P7 and P8 participated in the formative study.

Baseline. The baseline interface included for each image: the image caption from BLIP-2 [34], a list of objects from Detic [85], and the ability to interactively ask visual questions powered by BLIP-2 [34]. We designed the baseline to encompass commonly used captioning and object detection tools available in commercial devices and applications (e.g., SeeingAI [40]). As such captions tend to be concise, we added visual question answering via BLIP-2 [34] to let participants gain additional information on-demand.

Procedure. We first asked participants demographic and background questions about how they use images in their work. We then gave a 15-minute tutorial on both the *GenAssist* interface and the baseline interface using S0 (Figure 6). Participants then completed two tasks: the interpretation task and the generation task.

In the interpretation task, participants used both interfaces to evaluate pre-generated images (Figure 6). For each set of images, we provided participants with an example scenario (e.g., Select an image for a blog post titled ‘*My grandma still dances!*’). Using *GenAssist* or the baseline interface, participants were asked to identify the similarities and differences in the image candidates and choose a final image. For each interface, users were given one short prompt image set (S1 or S3) and one long prompt image set (S2 or S4). The order of the interfaces and image sets were counterbalanced and randomly assigned to participants. After each interface, we conducted a post-stimulus survey that included the following ratings: Mental Demand, Performance, Effort, Frustration, and Usefulness of the caption in understanding differences between images. All ratings were on a 7-point Likert scale.

In the generation task, we provided participants with the title and first 5 paragraphs of two articles, then asked participants to create a relevant image for the article by coming up with their own prompts. We selected the two articles from the New York Times: ‘*Why Multitasking is Bad for You*’ and ‘*My Kids Want Plastic Toys. I Want to Go Green.*’ [67, 68]. The order of the interfaces and articles was counterbalanced and randomly assigned to participants. After each interface, we asked the participants to choose one image from

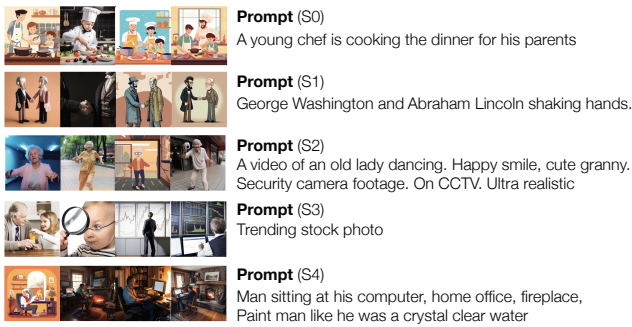


Figure 6: We selected two sets of images from Midjourney’s community feed generated with a short prompt without detailed descriptions of objects or styles (S1, S3) and two sets with a long prompt with detailed descriptions of objects or styles (S2, S4). We selected long and short prompts to explore how users compared images when they are similar (long prompts) vs. dissimilar (short prompts).

the generated images and explain their reasoning. We also conducted a post-stimulus survey that included the following ratings: Mental Demand, Performance, Effort, Frustration, Usefulness of the caption, Satisfaction with the final image, and Confidence in posting the final image. All ratings were on a 7-point Likert scale. At the end of the study, we conducted a semi-structured interview to understand participants’ strategies using GenAssist and the pros and cons of both GenAssist and the baseline.

The study was 1.5 hours long, conducted in a 1:1 session via Zoom, and approved by our institution’s IRB. We compensated participants 50 USD for their time.

Analysis. We recorded the study video, user-generated prompts and images, and the survey responses. We transcribed the exit interviews and participants’ spontaneous comments during the tasks and grouped the transcript according to (1) strategies of using GenAssist and (2) perceived benefits and limitations of our system.

6.2 Results

Overall, all participants stated they would like to use GenAssist rather than the baseline interface to create images in the future. Participants expressed that GenAssist would be immediately useful in their workflows: “This is usable out of the box!” [...] “I need access to this technology” (P14), “I’d even pay for this! I really need this” (P15). In particular, participants rated GenAssist to be significantly more useful for understanding the differences between images in both tasks (interpretation: $\mu=1.50$, $\sigma=1.00$ vs. $\mu=3.58$, $\sigma=4.00$; $Z=-2.31$; $p<0.05$; generation: $\mu=1.92$, $\sigma=2.00$ vs. $\mu=4.33$, $\sigma=5.00$; $Z=-2.77$; $p<0.01$) (Figure 7). For the interpretation task, participants reported significantly better performance ($\mu=1.83$, $\sigma=2.00$ vs. $\mu=3.67$, $\sigma=3.00$; $Z=-2.47$; $p<0.05$), significantly less frustration ($\mu=1.75$, $\sigma=1.00$ vs. $\mu=3.50$, $\sigma=3.50$; $Z=2.46$; $p<0.05$), and effort ($\mu=2.25$, $\sigma=2.00$ vs. $\mu=4.00$, $\sigma=4.00$; $Z=-2.00$; $p<0.05$). For generation tasks, participants rated that they were significantly more satisfied with the final image ($\mu=3.17$, $\sigma=3.00$ vs. $\mu=5.00$, $\sigma=5.50$; $Z=-2.17$; $p<0.05$). Significance was measured with the Wilcoxon Signed Rank test.

Gaining a summary of image content. With GenAssist across both tasks, all participants started by reading the summary table including the comparison description (summary of similarities and differences), as well as the per-image descriptions. Participants all stated that the summary table was helpful for understanding the images they generated, as P6 explained: “I cannot do without the summary. Highlighting the differences was very useful.” (P6). In addition, participants noted that the summary table’s per-image descriptions were valuable for understanding the images. For example, P19 mentioned “This is more like an audio description because I can make a very clear mental image!” and slowed down his screen reader pace to mimic the experience of listening to an audio description. P20 reported “I always thought that AI is not as capable of describing as humans, because usually alt-text generated by AI is short and doesn’t capture much information. But reading this, I am rethinking AI’s capabilities.”. P12 found the detailed descriptions particularly helpful when authoring rather than interpreting images: “The first table (comparison description table) is so comprehensive. When I’m authoring images I need more information than when I’m looking at what others uploaded.” (P12).

Using the baseline, participants all initially read all of the information they had access to (the caption and objects) for each image. all participants mentioned the inconvenience of having concise image captions for gaining an overview, especially when the generated images are similar to each other. For example, after reading the BLIP-2 caption of S4, P18 asked “Are they all same images?”

Selectively accessing additional information. While all participants accessed the summary table first, we observed multiple strategies of using additional information provided by GenAssist to understand the differences between the generated images. First, P9, P7, P16, P18, and P20 checked the information from all tables before making their decision. P20 mentioned “They are equally important but in different ways. If the generated images are different, the summary table would be sufficient. For similar ones, I’d have to go down the tables more.” P16 noted “We never have too much information. All the details provided here matter to me”. After checking all the tables, P18 and P20 revisited the summary table again to remember and organize all information. The other seven participants (P10-P12, P8-P15, P17, P19) checked the tables selectively. Participants’ preferences reflected their prior experiences creating images. For instance, P7 who typically creates images using an SVG editor prioritized the prompt verification table. He said “I detail more things in the prompt and want everything to be in the image, ‘cause I am more used to programming-drawing.” P13 skipped the style and errors table as he was not familiar with the concepts despite the definitions provided: “As a born blind person, most information in the visual attributes is not useful as it’s hard to imagine those.” Participants also mentioned that they liked that GenAssist provided the breakdown of the summary description into multiple tables. P16 described that GenAssist has “So much transparency because it provides access to intermediate tables that constitute the summary table, just like a [programming tool]! I can look at the inside of the models and see what they’re doing.” P10 and P11 both mentioned that they appreciated the order of the tables: “The summary [table] is the bigger picture. Then the tables go into the details. I also like that the prompt questions come first because they’re important.”

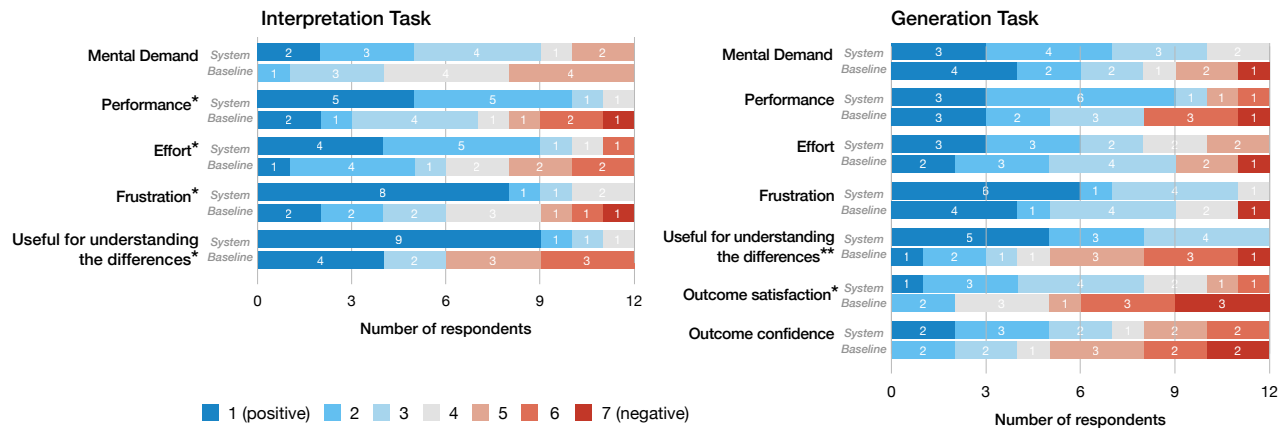


Figure 7: Distribution of the rating scores for GenAssist and the baseline interface (1 = positive, 7 = negative) in the two tasks. Note that a lower value indicates positive feedback and vice versa. The asterisks indicate the statistical significance as a result of Wilcoxon text ($p < 0.05$ is marked with * and $p < 0.01$ is marked with **). In the interpretation task, GenAssist significantly outperformed the baseline interface in performance, effort, frustration, and usefulness for understanding the differences between images. In the generation task, GenAssist was significantly lower in being useful for understanding the differences and in outcome satisfaction.

Participants also employed multiple strategies for navigating within the tables. Participants browsed through questions in the tables to identify questions they found to be important and skipped questions that were less important (e.g., not interested, or already appeared in the summary descriptions). We also identified multiple patterns of navigating within the tables. Participants checked all cells in a row when they found the table to be important. For instance, P11 checked the answers of all four images in the prompt verification table. In other cases, participants first checked the questions, then decided whether to read the row or skip to the next row. Participants skipped rows if the answers to the questions were already mentioned in the summary table, or if they were not interested in the question. For example, P8 skipped the medium, lighting, and perspective row in the visual style & errors table and only attended to the error row. Sometimes, participants only checked the answer cells if the summary column highlighted the differences between the images and skipped to the next row if the summary stated mainly the similarities between the images. Participants stated that GenAssist’s table format was easy to navigate. P19 noted the ease of navigation within the table: “I like having control with the tables. If the question or summary doesn’t seem interesting, I can skip to the next row instead of reading all answers of four images.”

Asking additional information. With the baseline, most participants (12 participants in the interpretation task, 9 participants in the generation task) asked follow-up questions to try to understand the images, while with our system participants rarely asked follow-up questions (1 participant in the interpretation task and none in the generation task). P16 was the only participant who asked additional visual questions with GenAssist after reading the table (‘Is the data showed falling or rising?’ and ‘What is the date of the x-axis?’ for S3 in Figure 6). When asked about the reason for not asking any additional questions, P18 said “Looking at captions I

already had a big picture so I didn’t ask additional questions.” P7 similarly reflected: “I like that [GenAssist] asks questions that I haven’t thought of but are still important. The answers to the questions told me additional stuff about the images.” In contrast, with the baseline interface, participants asked many additional visual questions. Because each image was presented separately, participants often asked the same question for each image to compare the answers. Most of the questions were about the objects detected, especially when the object was not mentioned in the caption or did not seem relevant to the setting (e.g., P11 asked “Where is the beachball in the picture?” after reading the object detection results of an image with the kitchen setting). P10 who experienced the baseline condition after GenAssist reflected that “This one [Baseline] is not simply laid out for me. The previous one [GenAssist] is easy peasy presenting everything for me. And this one is ‘Here you have to figure out.’”

Refining and Iterating Prompt. In the generation task, none of the participants refined the prompt using the baseline and five participants refined the prompt when using GenAssist (P9, P10, P13, P16, P17). Among the remaining 7 participants, 5 participants reported that they did not iterate as they were satisfied with the results, and 2 participants were unsure how to iterate the prompt after realizing that the image generation model did not reflect some parts of the original prompt (P15, P20).

Participants often quickly made the decision to revise the prompt while reading the summary table and before they moved on to other tables. For instance, while generating an image about an article about multitasking, P10 first attempted to generate an image with the following prompt ‘A woman who is holding the iPhone is texting on it while she glances at another device which displayed some funny videos going on. She’s in the kitchen trying to cook. it looks like the food is smoking’ Figure 8. However, she quickly noticed that most of the images generated depicted the woman as smoking instead of

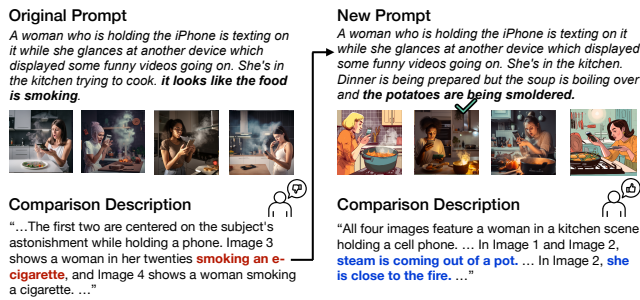


Figure 8: P10 generated the first set of images, noticed that the image generation model has made errors in the image (depicting the woman smoking instead of the food smoking), and corrected her prompt by replacing “smoking” with “being smoldered”.

the food as smoking. She quickly iterated the prompt by replacing the word with ‘smoldering’ to generate a new set of images.

In addition, participants reported that GenAssist informed them about the capabilities of the image generation model and guided them to refine their prompts. P20 mentioned “After reading the tables, it makes me think of what AI is capable of generating and what is not. It can’t exactly reflect what I try to accomplish when the prompt is too complicated, so I will have to adjust my expectation and adjust my prompt.” Participants also noted that GenAssist is helpful for learning how to generate a detailed prompt (P7, P16, P17). P16 stated “Visual [styles & errors] table is helpful for learning new styles.” Similarly, P7 said “If I don’t specify the styles, I think AI is generating [the styles] based on the context and content. So I know which style is good for which.”

Selecting an Image Candidate. To choose the final image from the four image candidates in the generation task, participants using GenAssist often considered whether the image followed the prompt, whether additional details added by the generation model were relevant, and whether the image style or emotion was appropriate to the usage context. P17 said “I choose the third image because it has the information that I described. Also, P7 mentioned “I will not choose the cartoon image because I want to be more serious here.” Some participants changed the choice of image as they moved on to the next tables in the GenAssist. For example, P8 who generated images of multiple plastic containers to portray the pollution problem updated his choice as he read the style and errors table: “Oh so the last image has many colors, I want to change to this one because I want it to be colorful!”

Noticed and unnoticed errors. Participants encountered errors using both interfaces. In the baseline, all participants read the objects following the captions, but objects occasionally contained errors (e.g., labeling as another object that has similar shapes, colors, or textures). When the participants noticed objects irrelevant to the context, they often asked about the object but the questions about non-existent objects often led to further confusion. For instance, P11 asked ‘Where is the television?’ for an image where a television is not present. Because the answer generated by BLIP-2 was ‘There

is no television.’, P11 was more confused and did not consider the image due to uncertainty. Also, P16 asked ‘Where is the lollipop in the image?’ for an image without a lollipop (S1 in Figure 6) and BLIP-2 answered with a hallucination ‘In the man’s mouth.’, misleading P16. While GenAssist features the same list of objects, participants did not experience this issue as they prioritized other information or recognized misinformation by referencing across multiple information sources. While using GenAssist, P10, P7, P16 pointed out that some visual information in the tables conflicted with one another. For instance, in the second image of S2 (Figure 6), the summary table stated that the woman is walking in the street, but when the GenAssist asked ‘Is she dancing?’ for prompt verification, BLIP-2 answered with ‘Yes’, which confused the participants. P16 hypothesized that the caption mentioned walking because the dancing action is hard to capture in one image frame and thus the image is actually showing her dancing. Still, participants did not notice inaccurate information in GenAssist if there was no conflict. For example, a woman was described as looking happy but had a neutral expression (the 4th image for P10’s 2nd prompt in Figure 8). P10 removed the image from consideration as she wanted the woman to look stressed rather than happy.

Future improvements for GenAssist. Participants noted suggestions on how to improve GenAssist’s description in the future. First, P9 and P8 participants noted that the visual information provided by GenAssist was long and difficult to process at once. This reflects users’ subjective ratings on mental demand which is comparable in GenAssist and the baseline in the interpretation task. Participants suggested allowing users to remove image columns and question rows from consideration. P8 mentioned “I want to filter images based on certain answers so that from then on, I won’t consider all four images and it will be easier!” P17 also shared that he wanted GenAssist to learn from his interactions with the cells so that gradually it will present only the rows of interest.

Participants mentioned the difficulties of writing good prompts. P13 said “Even if I read the definitions about the style, it’s hard to feel what effect it will give.” In the generation task, none of the participants specified the medium in the prompt as they were not familiar with it. This often resulted in the image generations having varied styles. In addition, P7 and P16 mentioned that it is difficult to decide on what content to put in the prompt to effectively convey the message. P16 mentioned “I want to give it the whole book and make it generate.” After experiencing that the generation model cannot reflect all the details in the prompt when the prompt is too long and complex, P12 stated “I want [GenAssist] to tell me what [the generation model] can generate and what it can not.”

7 DISCUSSION

In this section, we reflect on our findings from the development and evaluation of GenAssist. We also discuss future opportunities for research exploring accessible media authoring tools.

Scope of GenAssist. GenAssist uses a text-to-image generation model [41] to generate image candidates, vision-language models [34, 56] to extract visual information, and a large language model [46] to synthesize descriptions. The scope of GenAssist reflects the limitations of the models it uses. First, we designed GenAssist to support the images that text-to-image generation

models currently support: content-driven photos or illustrations with simple structures. However, both text-to-image generation and GenAssist do not yet support images that are information-rich or densely structured such as information visualizations [64, 65] or diagrams [3, 66]. As text-to-image generation improves, future research will explore extending GenAssist to complex graphics with text. For example, GenAssist could help creators recognize if their prompt-generated diagram contains the desired text (by integrating Optical Character Recognition), relationships, and perceptual qualities (e.g., legibility, saliency of important information).

Second, the descriptions that GenAssist is capable of providing are also limited by the capabilities of the pre-trained vision-language models [34, 56, 85]. For example, while GenAssist helped creators notice image generation errors such as omitted prompt details [36], distortions to human bodies [78], and objects placed illogically [80], some errors remained undetected. Also, GenAssist occasionally included hallucinations (e.g., missing or non-existent objects) in the descriptions. While these issues may be mitigated with improvements to text-to-image models (e.g., better aligning with human preferences [78]) and vision language models (e.g., better composition reasoning [38], reducing hallucinations [7]), GenAssist could also learn what prompts are prone to generation errors and guide BLV creators in creating strong prompts.

Finally, while GenAssist’s pipeline surfaced large differences between images (e.g., different objects, characters, expressions, or styles), its descriptions often missed smaller differences between images that were less likely to be described in training data captions (e.g., slightly different compositions or makeup styles). Thus, GenAssist is currently useful in the early stages of prompt iteration, where large differences between images remain. In the future, GenAssist could detect detailed changes by adding more detailed or domain-specific content and style questions, or integrating vision models that explicitly compare images [74].

Understanding Multiple Images. Creators in the formative study revealed that it is difficult to understand multiple images at the same time (D2. Understanding high-level image similarities and differences). To tackle this challenge, we designed GenAssist with three strategies: (1) providing the overview of similarities and differences between the generated image candidates, (2) progressively disclosing the information from high-level to low-level to give the user control over the level of detail received [23, 43, 50], and (3) presenting the descriptions in a table format so that users can easily navigate between images to compare them. Participants highlighted that not only these detailed summaries but also the ability to selectively gain information about the underlying questions were helpful in narrowing down their choices. For example, some participants prioritized the prompt verification table to assess if the image followed their instructions (D3. Assessing if images followed the prompt), and other participants used the content and style table to learn how to improve their prompts (D4. Accessing image details not specified by the prompt). In the future, GenAssist could support sorting or filtering images based on visual attributes to limit the number of images they consider at once (e.g., sorting images based on prompt adherence or filtering images that have AI-generated distortions). GenAssist could also read image descriptions with multiple voice styles to help creators distinguish generation candidates.

GenAssist’s ability to attend to multiple similar images and surface differences can be useful in broader contexts. Our study participants expressed interest in using GenAssist for comparing image search results or similar photos in social media. It can also help BLV people in decision-making situations based on visual information (e.g., online shopping, communicating with the design team in the software development, selecting a photo from similar shots).

Implications for Visual Question Answering. Comparing GenAssist to our baseline of typical descriptions with visual question answering (VQA), all participants rated GenAssist as more useful for understanding differences between images and creators asked fewer follow up questions with GenAssist. GenAssist reduced follow-up questions by predicting visual questions based on the formative study and applying the questions to multiple images. Our *predict-ask-summarize* approach also reduced the requirement for reading individual question answers. Future VQA systems intended for real-world environments may benefit from our approach as repetitive questions, “unknown unknowns”, and complex visuals are likely.

Support in Creating Prompts. In the formative study, we distilled the need to support creating prompts (D1. Authoring prompts that specify content and style). While we do not directly support prompt creation, we designed our system to reveal visual content and styles based on prompt guidelines to inform users about details the model filled in. In the user study, participants cited that reading the tables in GenAssist helped inform their prompt iterations and learn about what styles to use. Prior work has explored using structured search for visual concepts for writing prompts [37, 39], and combining our system with such prior work is a promising avenue for future work. We are currently exploring suggesting content and styles for the prompt when the user specifies the context of image use and new ways to help users add specificity to their prompt (e.g. a chatbot, as suggested in the formative study). In addition to text input, we can also consider multimodal input from users in the future such as image prompts [54], sketch prompts [11, 82], or music prompts [55] to create an image for a music album cover, as desired by P6.

Supporting Creators with Different Visual Impairments. BLV creators’ interest in color or style information (e.g., medium, lighting, angle) often depended on their prior experience with visuals and onset of blindness. GenAssist supports creators in selectively accessing description details, but in the future GenAssist will let creators control which details to filter out or prioritize. To support creators without knowledge of visual style, GenAssist could recommend popular styles given the image’s intended use, provide style descriptions, or deliver style in another modality (e.g., sound [21], tactile interfaces). We will also improve GenAssist in the future to support users with remaining vision beyond providing descriptions. For example, GenAssist could provide descriptions based on the current zoom viewing window or support further visual edits to the generated images, as desired by P1.

Implications of GenAssist on Creativity. Text-to-image generation models have sparked conversations about their implications for creativity. For BLV creators, image generation can improve creative agency compared to existing approaches for creating or selecting images. In our formative study, creators wanted to use image generation as it provided fewer limits over content and style

than searching for images online and greater autonomy than asking a sighted person to create the image. GenAssist supports BLV creators in exercising creative control over generated images by letting creators examine image details to revise the prompt or make an informed selection. Compared to sighted artists who use generated images primarily as references [37], BLV creators often intend to use generated images directly. In the future, GenAssist will further creative control by supporting prompt-based editing [4].

Implications of GenAssist on Communication. We designed GenAssist to support communication goals of BLV creators. BLV creators in our formative study aimed to create images to express their ideas to a broad audience and achieve self-expression. Images are particularly useful for capturing visual attention and communicating with sighted people who have difficulty reading text. For example, P4 generated an image of his family to share with his child. BLV creators also wanted to use GenAssist in the workplace and on digital platforms. As GenAssist exists in an ableist environment that prioritizes visual communication, there is a risk that GenAssist may cause sighted people to expect image-based communication from BLV people. Tools like GenAssist must be coupled with research and activism to make digital, workplace, and educational environments accessible — e.g., enabling non-visual communication and providing access to existing visuals. Our work also reveals that generated images themselves should be shared with descriptions in addition to the prompt that might not accurately reflect the image.

Generative AI for Accessible Media Authoring. Advances in large-scale generative models enable people to create new types of content, yet no existing research has explored people with disabilities as the users of these tools [28]. We see opportunities for generative AI models to broaden the type of content that people with disabilities can create. For example, our study participants mentioned that they are interested in using generative models for creating dynamic graphics like cartoons and videos. Similarly, generative models may be useful for people with motor impairments authoring visual media, or people with hearing impairments authoring music.

8 CONCLUSION

We created GenAssist, an accessible text-to-image generation system for BLV creators. Informed by our formative study with 8 BLV creators, our interface enables users to verify the adherence of generated images to their prompts, access additional image details, and quickly assess similarities and differences between image candidates. Our system is powered by large language and vision-language models that generate visual questions, extract answers, and summarize the visual information. Our user study with 12 BLV creators demonstrated the effectiveness of our approach. We hope this research will catalyze future work in supporting people with disabilities to express their creativity.

REFERENCES

- [1] Abubakar Abid, Ali Abdalla, Ali Abid, Dawood Khan, Abdulrahman Alfozan, and James Zou. 2019. Gradio: Hassle-free sharing and testing of ml models in the wild. *arXiv preprint arXiv:1906.02569* (2019).
- [2] AccessiblePublishing.ca. 2023 (accessed Apr 2, 2023). Guide to Image Descriptions. <https://www.accessiblepublishing.ca/a-guide-to-image-description/>
- [3] David Austin and Volker Sorge. 2023. Authoring Web-accessible Mathematical Diagrams. In *Proceedings of the 20th International Web for All Conference*. 148–152.
- [4] Omer Bar-Tal, Dolev Ofri-Amar, Rafail Fridman, Yoni Kasten, and Tali Dekel. 2022. Text2live: Text-driven layered image and video editing. In *European conference on computer vision*. Springer, 707–723.
- [5] Cynthia L Bennett, Jane E, Martez E Mott, Edward Cutrell, and Meredith Ringel Morris. 2018. How teens with visual impairments take, edit, and share photos on social media. In *Proceedings of the 2018 CHI conference on human factors in computing systems*. 1–12.
- [6] Jeffrey P Bigham, Chandrika Jayant, Hanjie Ji, Greg Little, Andrew Miller, Robert C Miller, Robin Miller, Aubrey Tatarowicz, Brandyn White, Samuel White, et al. 2010. Vizviz: nearly real-time answers to visual questions. In *Proceedings of the 23rd annual ACM symposium on User interface software and technology*. 333–342.
- [7] Ali Furkan Biten, Lluís Gómez, and Dimosthenis Karatzas. 2022. Let there be a clock on the beach: Reducing object hallucination in image captioning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 1381–1390.
- [8] Jens Bornschein and Gerhard Weber. 2017. Digital drawing tools for blind users: A state-of-the-art and requirement analysis. In *Proceedings of the 10th International Conference on Pervasive Technologies Related to Assistive Environments*. 21–28.
- [9] Erin Brady, Meredith Ringel Morris, Yu Zhong, Samuel White, and Jeffrey P Bigham. 2013. Visual challenges in the everyday lives of blind people. In *Proceedings of the SIGCHI conference on human factors in computing systems*. 2117–2126.
- [10] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [11] Tao Chen, Ming-Ming Cheng, Ping Tan, Ariel Shamir, and Shi-Min Hu. 2009. Sketch2photo: Internet image montage. *ACM transactions on graphics (TOG)* 28, 5 (2009), 1–10.
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiuhua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
- [14] Facebook. 2021. How Facebook is using AI to improve photo descriptions for people who are blind or visually impaired. <https://ai.facebook.com/blog/how-facebook-is-using-ai-to-improve-photo-descriptions-for-people-who-are-blind-or-visually-impaired/>
- [15] Olutayo Falase, Alexa F Siu, and Sean Follmer. 2019. Tactile code skimmer: A tool to help blind programmers feel the structure of code. In *Proceedings of the 21st International ACM SIGACCESS Conference on Computers and Accessibility*. 536–538.
- [16] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. Generative adversarial networks. *Commun. ACM* 63, 11 (2020), 139–144.
- [17] Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. 2018. Vizviz grand challenge: Answering visual questions from blind people. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3608–3617.
- [18] <https://github.com/mikhail-bot/>. 2023 (accessed Apr 2, 2023). Stable Diffusion Negative Prompts. <https://github.com/mikhail-bot/stable-diffusion-negative-prompts>
- [19] <https://github.com/pharmapsychotic/>. 2023 (accessed Apr 2, 2023). CLIP Interrogator. <https://github.com/pharmapsychotic/clip-interrogator>
- [20] <https://github.com/willwulfken/>. 2023 (accessed Apr 2, 2023). Midjourney Styles and Keywords. <https://github.com/willwulfken/MidJourney-Styles-and-Keywords-Reference>
- [21] <https://huggingface.co/spaces/ffiloni/>. 2023 (accessed Apr 2, 2023). Image to Music. <https://huggingface.co/spaces/ffiloni/img-to-music>
- [22] Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith. 2023. TIFA: Accurate and Interpretable Text-to-Image Faithfulness Evaluation with Question Answering. *arXiv preprint arXiv:2303.11897* (2023).
- [23] Mina Huh, YunJung Lee, Dasom Choi, Haesoo Kim, Uran Oh, and Juho Kim. 2022. Cocomix: Utilizing Comments to Improve Non-Visual Webtoon Accessibility. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–18.
- [24] Mina Huh, Saellyne Yang, Yi-Hao Peng, Xiang'Anthony' Chen, Young-Ho Kim, and Amy Pavel. 2023. AVscript: Accessible Video Editing with Audio-Visual Scripts. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*.
- [25] Jiho Kim, Arjun Srinivasan, Nam Wook Kim, and Yea-Seul Kim. CHI 2023. Exploring Chart Question Answering for Blind and Low Vision Users. (CHI 2023).
- [26] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).

- [27] Hyung-Kwon Ko, Subin An, Gwanmo Park, Seung Kwon Kim, Daesik Kim, Bohyoung Kim, Jaemin Jo, and Jinwook Seo. 2022. We-toon: A Communication Support System between Writers and Artists in Collaborative Webtoon Sketch Revision. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology*. 1–14.
- [28] Hyung-Kwon Ko, Gwanmo Park, Hyeon Jeon, Jaemin Jo, Juho Kim, and Jinwook Seo. 2023. Large-scale text-to-image generation models for visual artists' creative works. In *Proceedings of the 28th International Conference on Intelligent User Interfaces*. 919–933.
- [29] Mackenzie Leake, Hijing Valentina Shin, Joy O Kim, and Maneesh Agrawala. 2020. Generating Audio-Visual Slideshows from Text Articles Using Word Coherence. In *CHI*, Vol. 20. 25–30.
- [30] Cheuk Yin Phipson Lee, Zhuohao Zhang, Jaylin Herskovitz, JooYoung Seo, and Anhong Guo. CHI 2022. CollabAlly: Accessible Collaboration Awareness in Document Editing. (CHI 2022).
- [31] Jaewook Lee, Jaylin Herskovitz, Yi-Hao Peng, and Anhong Guo. 2022. Image-Explorer: Multi-Layered Touch Exploration to Encourage Skepticism Towards Imperfect AI-Generated Image Captions. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–15.
- [32] Jaewook Lee, Yi-Hao Peng, Jaylin Herskovitz, and Anhong Guo. 2021. Image Explorer: Multi-Layered Touch Exploration to Make Images Accessible. In *Proceedings of the 23rd International ACM SIGACCESS Conference on Computers and Accessibility*. 1–4.
- [33] Jingyi Li, Son Kim, Joshua A Miele, Maneesh Agrawala, and Sean Follmer. 2019. Editing spatial layouts through tactile templates for people with visual impairments. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. 1–11.
- [34] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597* (2023).
- [35] Junchen Li, Garreth W. Tigwell, and Kristen Shinohara. 2021. Accessibility of high-fidelity prototyping tools. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–17.
- [36] 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*. 1–23.
- [37] Vivian Liu, Han Qiao, and Lydia Chilton. 2022. Opal: Multimodal Image Generation for News Illustration. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology*. 1–17.
- [38] Zixian Ma, Jerry Hong, Mustafa Omer Gul, Mona Gandhi, Irena Gao, and Ranjay Krishna. 2023. CREPE: Can Vision-Language Foundation Models Reason Compositionally?. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10910–10921.
- [39] Shane McGeehan. 2023 (accessed Apr 2, 2023). Prompter. <https://prompterguide.com/prompter/>
- [40] Microsoft. 2021. Seeing AI. <https://www.microsoft.com/en-us/ai/seeing-ai>
- [41] Midjourney. 2023 (accessed Apr 2, 2023). Midjourney. <https://www.midjourney.com>
- [42] Midjourney. 2023 (accessed Apr 2, 2023). Midjourney Propmt Guidelines. <https://docs.midjourney.com/docs/prompts>
- [43] Meredith Ringel Morris, Jazette Johnson, Cynthia L Bennett, and Edward Cutrell. 2018. Rich representations of visual content for screen reader users. In *Proceedings of the 2018 CHI conference on human factors in computing systems*. 1–11.
- [44] Hospital News. 2016. You are what you eat. <https://hospitalnews.com/you-are-what-you-eat-why-nutrition-matters/>
- [45] Peter O'Donovan, Aseem Agarwala, and Aaron Hertzmann. 2015. Designscape: Design with interactive layout suggestions. In *Proceedings of the 33rd annual ACM conference on human factors in computing systems*. 1221–1224.
- [46] OpenAI. 2023. GPT-4 Technical Report. [arXiv:2303.08774](https://arxiv.org/abs/2303.08774) [cs.CL]
- [47] Guy Parsons. 2023 (accessed Apr 2, 2023). DALL-E2 Propmt Book. <https://dallery.gallery/the-dalle-2-prompt-book/>
- [48] William Christopher Payne, Alex Yixuan Xu, Fabiha Ahmed, Lisa Ye, and Amy Hurst. 2020. How blind and visually impaired composers, producers, and songwriters leverage and adapt music technology. In *Proceedings of the 22nd International ACM SIGACCESS Conference on Computers and Accessibility*. 1–12.
- [49] Yi-Hao Peng, Jeffrey P Bigham, and Amy Pavel. 2021. Slidecho: Flexible Non-Visual Exploration of Presentation Videos. In *The 23rd International ACM SIGACCESS Conference on Computers and Accessibility*. 1–12.
- [50] Yi-Hao Peng, Peggy Chi, Anjali Kannan, Meredith Morris, and Irfan Essa. 2023. Slide Gestalt: Automatic Structure Extraction in Slide Decks for Non-Visual Access. (2023).
- [51] Yi-Hao Peng, JiWoong Jang, Jeffrey P Bigham, and Amy Pavel. 2021. Say It All: Feedback for Improving Non-Visual Presentation Accessibility. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–12.
- [52] Yi-Hao Peng, Jason Wu, Jeffrey Bigham, and Amy Pavel. 2022. Diffscrber: Describing Visual Design Changes to Support Mixed-Ability Collaborative Presentation Authoring. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology*. 1–13.
- [53] Venkatesh Potluri, Liang He, Christine Chen, Jon E Froehlich, and Jennifer Mankoff. 2019. A multi-modal approach for blind and visually impaired developers to edit webpage designs. In *Proceedings of the 21st International ACM SIGACCESS Conference on Computers and Accessibility*. 612–614.
- [54] 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*. 15–28.
- [55] Yue Qiu and Hirokatsu Kataoka. 2018. Image generation associated with music data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. 2510–2513.
- [56] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- [57] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125* (2022).
- [58] 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 *International Conference on Machine Learning*. PMLR, 8821–8831.
- [59] Mr D Murahari Reddy, Mr Sk Masthan Basha, Mr M Chinnaiahgari Hari, and Mr N Penchalaiah. 2021. Dall-e: Creating images from text. *UGC Care Group I Journal* 8, 14 (2021), 71–75.
- [60] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2021. High-Resolution Image Synthesis with Latent Diffusion Models. [arXiv:2112.10752](https://arxiv.org/abs/2112.10752) [cs.CV]
- [61] Anastasia Schaadhardt, Alexis Hiniker, and Jacob O Wobbrock. 2021. Understanding blind screen-reader users' experiences of digital artboards. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–19.
- [62] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv preprint arXiv:2210.08402* (2022).
- [63] Athar Sefid, Prasenjit Mitra, and Lee Giles. 2021. SlideGen: an abstractive section-based slide generator for scholarly documents. In *Proceedings of the 21st ACM Symposium on Document Engineering*. 1–4.
- [64] Ather Sharif, Olivia H Wang, Alida T Muongchan, Katharina Reinecke, and Jacob O Wobbrock. 2022. Voxlens: Making online data visualizations accessible with an interactive javascript plug-in. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–19.
- [65] Ather Sharif, Andrew M Zhang, Katharina Reinecke, and Jacob O Wobbrock. 2023. Understanding and Improving Drilled-Down Information Extraction from Online Data Visualizations for Screen-Reader Users. In *Proceedings of the 20th International Web for All Conference*. 18–31.
- [66] Volker Sorge, Mark Lee, and Sandy Wilkinson. 2015. End-to-end solution for accessible chemical diagrams. In *Proceedings of the 12th International Web for All Conference*. 1–10.
- [67] NY Times. 2023 (accessed Apr 2, 2023). My Kids Want Plastic Toys. I Want to Go Green. <https://time.com/6126981/my-kids-want-plastic-toys-i-want-to-go-green-heres-a-fix/>
- [68] NY Times. 2023 (accessed Apr 2, 2023). Why Multitasking is Bad for You. <https://time.com/4737286/multitasking-mental-health-stress-texting-depression/>
- [69] Iulia Turc and Gaurav Nemade. 2022. Midjourney User Prompts & Generated Images (250k). <https://doi.org/10.34740/KAGGLE/DS/2349267>
- [70] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [71] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3156–3164.
- [72] Luis Von Ahn and Laura Dabbish. 2004. Labeling images with a computer game. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. 319–326.
- [73] W3C Web Accessibility Initiative (WAI). 2022 (accessed Dec 12, 2022). Introduction to web accessibility. <https://www.w3.org/WAI/fundamentals/accessibility-intro/>
- [74] Jiuniu Wang, Wenjia Xu, Qingzhong Wang, and Antoni B Chan. 2020. Compare and reweight: Distinctive image captioning using similar images sets. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*. Springer, 370–386.
- [75] Ruolin Wang, Zixuan Chen, Mingrui Ray Zhang, Zhaoheng Li, Zhixiu Liu, Zihan Dang, Chun Yu, and Xiang'Anthony' Chen. 2021. Revamp: Enhancing Accessible Information Seeking Experience of Online Shopping for Blind or Low Vision Users. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–14.
- [76] Yunlong Wang, Shuyuan Shen, and Brian Y Lim. 2023. RePrompt: Automatic Prompt Editing to Refine AI-Generative Art Towards Precise Expressions. *arXiv preprint arXiv:2302.09466* (2023).

- [77] Shaomei Wu, Jeffrey Wieland, Omid Farivar, and Julie Schiller. 2017. Automatic alt-text: Computer-generated image descriptions for blind users on a social network service. In *proceedings of the 2017 ACM conference on computer supported cooperative work and social computing*. 1180–1192.
- [78] Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. 2023. Better aligning text-to-image models with human preference. *arXiv preprint arXiv:2303.14420* (2023).
- [79] Haijun Xia. 2020. Crosspower: Bridging graphics and linguistics. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology*. 722–734.
- [80] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. 2023. Imagereward: Learning and evaluating human preferences for text-to-image generation. *arXiv preprint arXiv:2304.05977* (2023).
- [81] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*. PMLR, 2048–2057.
- [82] Lvmin Zhang and Maneesh Agrawala. 2023. Adding Conditional Control to Text-to-Image Diffusion Models. *arXiv:2302.05543* [cs.CV]
- [83] Zhuohao Zhang and Jacob O. Wobbrock. CHI 2023. A11yBoard: Making Digital Artboards Accessible to Blind and Low-Vision Users.
- [84] Yu Zhong, Walter S Lasecki, Erin Brady, and Jeffrey P Bigham. 2015. Regionspeak: Quick comprehensive spatial descriptions of complex images for blind users. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*. 2353–2362.
- [85] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. 2022. Detecting Twenty-thousand Classes using Image-level Supervision. In *ECCV*.

A STUDY PARTICIPANTS DEMOGRAPHICS

| PID | Gender | Age | Visual Impairment | Onset | Job | Images Produced |
|-----|------------|-----|-------------------|---------------|--------------------------|------------------------------------|
| P1 | Non-binary | 40 | Legally blind | Congenital | Artist | Paintings, Cartoons |
| P2 | Male | 50 | Totally blind | Congenital | Professor (CS) | Presentations, Scientific figures |
| P3 | Female | 29 | Legally blind | Congenital | Teacher (English) | Presentations, Course website |
| P4 | Male | 28 | Totally blind | Acquired | Teacher (Music) | Website logos |
| P5 | Male | 59 | Totally blind | Congenital | Professor (Climate) | Presentations, Scientific figures |
| P6 | Male | 42 | Totally blind | Acquired | Software engineer | Website images, Music album cover |
| P7 | Male | 32 | Totally blind | Acquired | Software engineer | Website images |
| P8 | Male | 30 | Totally blind | Acquired | Graduate student | Presentations |
| P9 | Female | 41 | Totally blind | Congenital | Graduate student | Presentations, Social media images |
| P10 | Female | 30 | Totally blind | Acquired | Graduate student | Presentations, Website images |
| P11 | Female | 37 | Totally blind | Congenital | Accessibility consultant | Website images |
| P12 | Male | 50 | Legally blind | Totally blind | Finance consultant | Charts, Graphs |
| P13 | Male | 61 | Totally blind | Congenital | YouTuber, Musician | Video thumbnails |
| P14 | Male | 44 | Totally blind | Congenital | Author, Photographer | Book covers |
| P15 | Male | 20 | Totally blind | Congenital | University student | Book covers |
| P16 | Male | 36 | Totally blind | Congenital | Artist | Event flyers |
| P17 | Male | 26 | Totally blind | Congenital | Accessibility consultant | Icons, Video thumbnails |
| P18 | Male | 47 | Legally blind | Acquired | Software engineer | Brochures, Website images |

Table 4: Participant table for formative and comparison study.