# DanmuA11y: Making Time-Synced On-Screen Video Comments (Danmu) Accessible to Blind and Low Vision Users via Multi-Viewer Audio Discussions

Shuchang Xu
Hong Kong University of Science and Technology
Hong Kong, China
sxuby@connect.ust.hk

Xiaofu Jin
Hong Kong University of Science and Technology
Hong Kong, China
xjinao@connect.ust.hk

Huamin Qu
Hong Kong University of Science and Technology
Hong Kong, China
huamin@cse.ust.hk

Yukang Yan*
University of Rochester
New York, United States
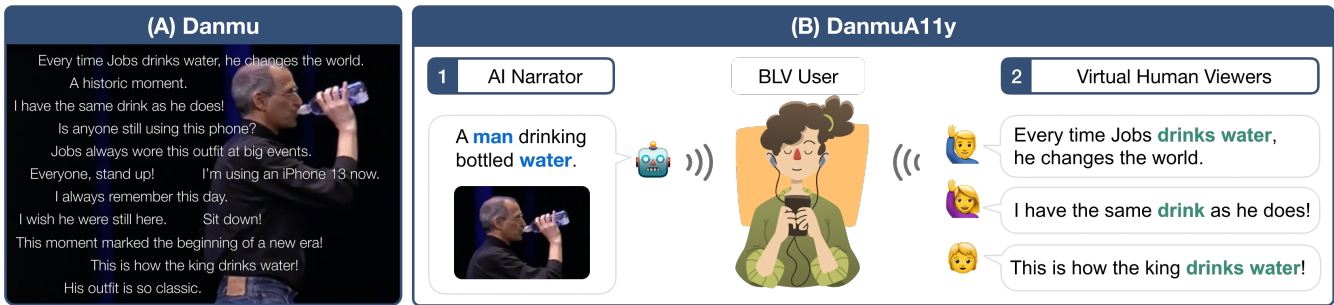yukang.yan@rochester.edu

Figure 1: (A) Danmu is a video-commenting feature that overlays time-synced user comments onto videos. It creates a co-watching experience for online viewers. However, its visual-centric design poses significant challenges for blind and low vision (BLV) viewers. (B) DanmuA11y makes Danmu accessible to BLV viewers by transforming it into multi-viewer audio discussions. It includes two types of virtual viewers: (1) an AI narrator offering visual descriptions, and (2) virtual human viewers verbalizing curated Danmu comments. These audio discussions are seamlessly integrated into the video, creating an enjoyable and socially engaging experience for BLV viewers. (All Danmu comments in this paper were originally in Mandarin and have been translated into English.)

## ABSTRACT

By overlaying time-synced user comments on videos, Danmu creates a co-watching experience for online viewers. However, its visual-centric design poses significant challenges for blind and low vision (BLV) viewers. Our formative study identified three primary challenges that hinder BLV viewers' engagement with Danmu: the lack of visual context, the speech interference between comments and videos, and the disorganization of comments. To address these challenges, we present DanmuA11y, a system that makes Danmu accessible by transforming it into multi-viewer audio discussions. DanmuA11y incorporates three core features: (1) *Augmenting Danmu with visual context*, (2) *Seamlessly integrating Danmu into videos*, and (3) *Presenting Danmu via multi-viewer discussions*. Evaluation with twelve BLV viewers demonstrated that DanmuA11y significantly improved Danmu comprehension, provided smooth viewing experiences, and fostered social connections among viewers. We further highlight implications for enhancing commentary accessibility in video-based social media and live-streaming platforms.

*This is the corresponding author.

## CCS CONCEPTS

• **Human-centered computing → Accessibility systems and tools**.

## KEYWORDS

Visual Impairment, Blind, Low Vision, Video, Social Media, Danmaku, Danmu, Bullet Comment

## 1 INTRODUCTION

Online video platforms have become important channels for viewers to exchange thoughts, emotions, and knowledge related to videos [56, 88]. Traditionally, these interactions occur in a separate comments section, detached from the video itself, as exemplified by platforms like YouTube[1]. This structure often fails to capture viewers' immediate reactions to specific video moments, leading to a disconnect between the video and the commentary [60]. To address this limitation, several Asian video platforms have introduced a video-commenting feature known as Danmu. It enables viewers to post comments that are synchronized with the video timeline. These *time-synced* comments are *visually overlaid* on the video screen (see Figure 1 (A)), allowing viewers to see others' instant reactions as they watch the video. Moreover, viewers can interact with each other by responding to previous comments [60]. Consequently, Danmu creates a socially engaging experience that closely simulates the sensation of co-watching videos with other viewers [15, 60, 75].

However, Danmu is primarily designed for *visual consumption*. Danmu comments are visually overlaid onto videos, often in overwhelming volumes [7, 13, 60]. These comments are scattered across the screen, requiring viewers to visually process multiple simultaneous comments and discern social interactions, such as identifying which comments are replies to others. These visual-centric characteristics pose significant challenges for blind and low vision (BLV) viewers, making it difficult for them to interpret Danmu comments and socially connect with other viewers.

To understand the practices and challenges of BLV viewers in accessing Danmu, we conducted formative interviews and co-watching sessions with eight BLV participants who regularly watched videos on Danmu-enabled platforms. Participants accessed Danmu using a design probe modeled after the most effective tool they currently use: an auto-scrolling list of Danmu comments. Our research revealed three significant challenges that hinder their engagement with Danmu: First, Danmu comments often discuss visual elements that are not accessible to BLV viewers. This **lack of visual context** makes it difficult for them to understand the discussion topics. Second, while BLV viewers use screen readers to access Danmu, the audio from screen readers often overlaps with the video's speech. This **speech interference** prevents them from enjoying both content at the same time. Third, while Danmu utilizes screen space to display a large number of comments simultaneously, they become disorganized when accessed sequentially via screen readers. This **disorganization of comments** makes it tedious for BLV viewers to follow audience discussions and socially engage with other viewers.

To address these challenges, we present DanmuA11y, a system that makes Danmu accessible by transforming Danmu into multi-viewer audio discussions. DanmuA11y encompasses three core features: (1) ***Augmenting Danmu with Visual Context***: DanmuA11y supplements Danmu comments with descriptions of the visual context, allowing BLV viewers to easily grasp the discussion topics; (2) ***Seamlessly Integrating Danmu into Videos***: By optimizing the insertion timing of Danmu comments in the video, DanmuA11y enables viewers to enjoy both content without speech overlap; and (3) ***Presenting Danmu via Multi-Viewer Discussions***: To create a co-watching experience, DanmuA11y organizes Danmu comments into dialogues and uses spatial audio to simulate the sensation of other viewers conversing around the user. Powered by these features, DanmuA11y offers an enjoyable and socially engaging experience for BLV video viewers.

To evaluate DanmuA11y, we conducted a within-subject study with 12 BLV participants, who compared DanmuA11y to a baseline system that simulated their current practices. Each participant watched two similar groups of videos using two systems, respectively. The results showed that DanmuA11y significantly improved Danmu comprehension ($p < .01$), provided a smoother viewing experience ($p < .01$), and enhanced the sense of co-watching with others ($p < .01$) compared to the baseline. Participants reported that DanmuA11y made video watching more enjoyable and socially engaging, with a strong willingness to use it in the future ($\mu = 6.92$ on a seven-point Likert scale). Based on our findings, we discussed directions for personalizing DanmuA11y, and summarized implications for improving commentary accessibility in broader contexts, such as video-based social media and live-streaming platforms.

Our contributions are threefold:

- We identify the motivations, needs, current practices, and challenges of BLV viewers in accessing Danmu, derived from a formative study;
- We present DanmuA11y, a system that makes Danmu accessible to BLV viewers through three core features: (1) augmenting Danmu with visual context, (2) seamlessly integrating Danmu into videos, and (3) presenting Danmu via multi-viewer audio discussions;
- We contribute an evaluation study that demonstrates how BLV viewers engage with Danmu using DanmuA11y, and derive design implications to improve commentary accessibility in video and live-streaming platforms.

## 2 BACKGROUND AND RELATED WORK

Our work extends prior research in three areas: (1) Danmu and its accessibility, (2) social media accessibility, and (3) video accessibility for BLV viewers.

### 2.1 Danmu and its Accessibility

Danmu, also known as Danmaku, is a video-commenting feature that displays time-synced user comments as moving subtitles overlaid on the video screen. Originating in Japan, this feature has gained widespread popularity on Asian video platforms. For example, in China, nearly all major online video platforms support

---

Danmu [15]. As of 2024, China's first Danmu-enabled video platform, Bilibili[2], has accumulated over 20 billion Danmu comments.
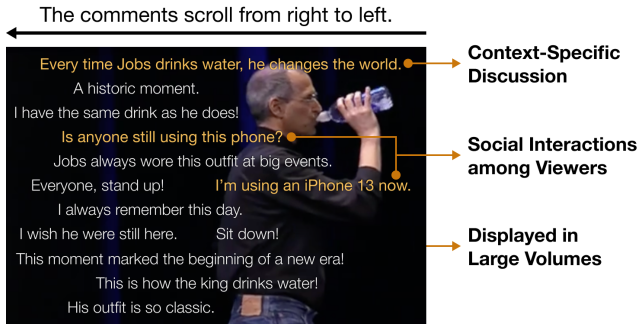


**Figure 2: The key characteristics of Danmu comments.**

Compared to traditional comments, Danmu creates a closer connection between the video and the commentary [15, 60]. Danmu comments are synchronized with the video timeline and visually overlaid on the screen, resulting in several unique features (see Figure 2). First, Danmu comments appear exactly at the video moment when they are posted. This allows viewers to exchange real-time, **context-specific** information (e.g., identifying an actor in a scene) rather than post-viewing reflections [15]. Second, Danmu comments are immediately displayed on the screen, visible to both the poster and other viewers. This enables more direct **social interaction among users**, creating a pseudo co-watching experience [15, 60]. Third, Danmu comments often come in **large volumes** simultaneously [59, 88]. This concurrent display enables sighted viewers to skim through the information efficiently. Fourth, Danmu comments are anonymous [13, 60]—usernames are hidden, and comments are not structured into threads. This anonymity leads viewers to address each other through semantic references [31, 60]. For instance, users might reply with phrases like "*The one saying 'Why not buy an airplane', are you serious?*" or repeat the same comment to express shared opinions.

The HCI community has emphasized Danmu's role in enhancing social connection: it creates a co-watching experience [15, 60, 87], facilitates information exchange [12, 30, 49], and fosters a sense of belonging [60, 90]. However, the lack of accessibility support for Danmu may prevent certain individuals from participating in this social interaction. For example, people who are hard of hearing may find it challenging to engage with text-based Danmu comments [12], as they often prefer sign language over text for information access. More significantly, Danmu's inherently visual design presents major challenges for BLV viewers. Despite these barriers, there has been limited understanding on how BLV viewers engage with Danmu. To fill this gap, we conducted a study where we invited BLV viewers to participate in co-watching sessions. Our work identified three key challenges that hinder BLV viewers' engagement with Danmu: the lack of visual context, the speech interference between comments and videos, and the disorganization of comments. In response, we designed DanmuA11y to address these challenges and create a socially engaging video viewing experience for BLV viewers.

[2]https://www.bilibili.com/

## 2.2 Social Media Accessibility

Addressing the accessibility of social media has been a consistent focus in HCI research [22, 23]. The emergence of new types of visual content continues to challenge BLV users in fully participating on social media platforms [72, 73, 80]. In response, researchers have explored methods to help BLV users access various forms of visual content, including images [26, 89, 91], memes [25], GIFs [24, 100], emojis [77, 99], comics [35], and videos [56, 78, 84]. When engaging with social media content, BLV users seek more than just factual descriptions (e.g., identifying objects in an image); they also aim to establish social and emotional connections with other users [25]. For instance, BLV users hope to understand the humor in memes [25], interpret comics through others' perspectives [35], and share emotional resonance with other video viewers [8, 43].

To help BLV viewers engage with others' perspectives, providing accessible user comments play a crucial role [35, 82]. Tools like Cocomix [35] utilizes user comments to offer additional context to comics, enabling BLV users to enjoy humor and engage with the community. However, unlike traditional user comments, Danmu comments are time-synced and visually overlaid onto videos, presenting unique challenges for BLV users. Our research addresses these challenges by integrating Danmu comments into videos as multi-viewer audio discussions, enabling BLV users to engage more deeply with other audiences. By studying time-synced video comments, our research offers implications for broader contexts, such as enhancing the accessibility of real-time commentary on live-streaming platforms [70].

## 2.3 Video Accessibility for BLV Viewers

To enhance video accessibility for BLV viewers, prior research has explored various methods, including audio descriptions and spatial audio representations. In this section, we review these techniques and discuss how our work extends existing literature.

*2.3.1 Audio Descriptions.* Audio descriptions make videos accessible to BLV viewers by narrating visual elements in sync with the video content [5]. Prior works have summarized guidelines for creating effective audio descriptions [66]: (1) Describe important visual elements that are essential for understanding the video, (2) Do not overlap audio descriptions with the video's dialogue, and (3) Use an objective style and avoid subjective interpretation. Unlike audio descriptions that provide objective visual information, Danmu comments emphasize social discussions among viewers and often contain subjective opinions (e.g., "*His outfit is so classic.*") and emotional reactions (e.g., "*A historic moment!*") [60], leading to distinct information needs. In response, our work identifies BLV viewers' needs for Danmu access through interviews and co-watching exercises with end users.

Prior research has explored various methods to facilitate audio description creation, including manual [46], semi-automatic [57, 66], and automatic approaches [84]. For instance, CrossA11y [57] offers semi-automatic support for authors by identifying accessibility issues in videos through the detection of visual-auditory inconsistencies. Similarly, Rescribe [66] utilizes dynamic programming to optimize the placement of audio descriptions within videos. Tiresias [84] presents a fully automatic pipeline that detects visual-auditory discrepancies, generates audio descriptions, and refines them into

a coherent output. CustomAD [63] further enables BLV viewers to personalize the content and style of audio descriptions. Compared to audio descriptions, Danmu comments are abundant and disorganized, presenting new challenges in content curation and integration. Our work addresses these challenges by proposing an automatic pipeline for curating and integrating Danmu comments, ensuring a seamless video viewing experience.

*2.3.2 Spatial Audio Representations.* Besides audio descriptions, prior works [11, 20, 40, 64] have explored using spatial audio to represent spatial information in videos. For instance, Front Row [40] employs spatial audio to encode the positions and movements of players in sports broadcasts, allowing BLV viewers to directly perceive the action in real time. Similarly, SPICA [64] leverages spatial sound to help BLV viewers identify the location of objects within video key frames. Following this, Dang et al. [20] utilizes spatial audio descriptions to enhance BLV viewers' spatial understanding of music performances in virtual reality. Recent research [38, 39] has demonstrated that spatial audio can enhance social presence in multi-person remote communication. Building on this insight, our work explores transforming Danmu comments into multi-viewer spatial audio discussions to enhance BLV users' social presence. Evaluation results showed that this approach effectively strengthened BLV viewers' sense of co-watching videos with others. Based on our findings, we derive insights for improving BLV users' social engagement in broader contexts.

## 3 FORMATIVE STUDY

We conducted a formative study with eight BLV viewers to understand their needs, practices, and challenges in accessing Danmu. Through interviews and co-watching exercises, we addressed the following questions:

(1) **Motivations**: Why do BLV viewers want to access Danmu?

(2) **Practices**: How do BLV viewers currently access Danmu?

(3) **Challenges**: What challenges do BLV viewers encounter in accessing Danmu?

### 3.1 Methods

*3.1.1 Participants.* We recruited eight BLV viewers (P1-P8; four male, four female; Table 1 lists their demographics) who frequently watched videos on Danmu-enabled platforms. These participants were recruited from an online support community, with ages ranging from 26 to 38 (mean = 31.2, SD = 4.4). Four participants were totally blind and four participants were legally blind with light perception. All participants had prior experience accessing Danmu on platforms such as Bilibili, Douyin[3], and Youku[4]. They mainly used these platforms on mobile phones and their viewing frequency ranged from several times a week to daily. All participants were native Mandarin speakers.

*3.1.2 Design probe.* Prior to the study, seven out of eight participants identified an mobile application (Etong)[5] as the **best currently available tool** for accessing Danmu. This tool presents

Danmu comments in an auto-scrolling list but is limited to livestreaming. Since no Danmu-enabled video platforms offer similar functionality, we replicated Etong's interface in our design probe.

As shown in Figure 3, the design probe displays Danmu comments in a list, which is accessible via screen readers [36, 41]. The list auto-scrolls as the video plays, allowing users to easily view the latest comments. Additionally, the probe supports standard video playback controls, including a play/pause button and a slider.

The design probe contained six videos, each representing a genre commonly watched by BLV viewers [43]: an educational video, a comedic video, a tutorial video, a news video, a music video, and a film clip. The video lengths ranged from 33 seconds to 7 minutes (mean = 3.5 minutes, SD = 2.2 minutes). Table 2 lists their details. These videos were randomly chosen from videos with over 1,000 Danmu comments on Bilibili. We downloaded the latest 1,000 Danmu comments, the default maximum provided by the platform. We imported videos in MP4 format, Danmu comments in XML format, and creator-written titles and subtitles as texts into the design probe. The probe was implemented as an iOS app on an iPhone 15 and has been tested for compatibility with VoiceOver[6].



**Figure 3: The design probe presents Danmu comments in an auto-scrolling list. Participants used the probe in Mandarin.**

*3.1.3 Procedure.* The study consisted of two successive phases: a 45-minute interview followed by a 45-minute co-watching exercise, with a 5-minute break in between.

**Phase 1: Semi-structured Interview.** We asked participants about their motivations, practices, and challenges in accessing Danmu. To elicit details, we asked participants to show their practices on their preferred platforms. When interesting points arose, we followed up with questions to probe further details.

---

[3]https://www.douyin.com/

[4]https://www.youku.tv/?lang=en_US

[5]https://apps.apple.com/app/id6471562648

[6]VoiceOver is the built-in screen reader on iOS devices.

**Phase 2: Co-watching Exercise.** In this session, participants first received a brief tutorial on the design probe and then used it to view the six pre-selected videos in a random order. They were encouraged to think aloud about any confusion while viewing videos. After the session, we conducted an exit interview to understand the challenges of using the design probe. The entire study was conducted one-on-one, in-person. Participants were compensated approximately 17 USD in local currency for their time.

*3.1.4 **Analysis***. We recorded the study audio, transcribed it[7], and then used an open-coding approach [19] to analyze participants' motivations, current practices, and challenges in accessing Danmu. Two authors independently reviewed the transcripts, developing and applying codes through an iterative process until reaching a consensus. The codes were then grouped into clusters representing the emerging themes from the study data. After a thorough discussion and review of all potential themes, the final themes were reported as the study findings.

## 3.2 Findings

*3.2.1 **Motivations for Accessing Danmu***. Participants primarily viewed Danmu as a *real-time* communication channel, where they sought comments that are "*closely related to the current video moment*" (P1). During the interview, they mainly reported three motivations for accessing Danmu:

**(1) To Enjoy Creative Comments**: Seven participants sought joy from Danmu comments that were "*creative*" (P1), "*unexpected*" (P3), or "*hilarious*" (P7). For example, P2 found a comment in a bakery video particularly creative: "*Use soy sauce if you don't have dark chocolate.*" Similarly, P5 described some Danmu comments as "*non-visual bloopers*", citing the unexpected comment, "*There is a bottle of mineral water on the table!*" in a historical drama. These comments often added an element of surprise and enhanced the enjoyment of their viewing experience.

**(2) To Gain Supplementary Information**: Six participants valued Danmu for providing additional information not present in the original video. For instance, P5 recalled many informative Danmu comments when Steve Jobs released the first iPhone, such as "*Jobs was really nervous because there were so many issues with this prototype iPhone.*" Such information "*enriched the video*" (P5) and "*deepened their understanding of the moment*" (P8).

**(3) To Engage with Diverse Opinions**: Six participants appreciated the diverse perspectives in Danmu comments. For example, P2 noted that in a video discussing whether to buy a smart speaker, viewers shared contrasting opinions like, "*Don't buy it! Its sound quality is bad.*" and "*You can't miss it! The voice assistant is really useful.*" Such debates "*offered new perspectives*" (P5), "*heated up the discussion*" (P7), and "*made the conversation more engaging*" (P2).

*3.2.2 **Current Practices***. Despite the benefits of Danmu, all participants noted that accessibility support for Danmu on mainstream video platforms is quite limited. These platforms either offer only a single Danmu comment (near the playback time) or make all comments inaccessible to screen reader users, forcing participants to adopt resourceful **workarounds**.

Seven participants attempted to use screen recognition tools (e.g., VoiceOver Recognition[8]) to access Danmu as on-screen text but found them ineffective due to issues like "*having to pause the video*" (P2) and "*recognition errors from cluttered text*" (P5). In addition to technical workarounds, three participants mentioned asking friends to read out interesting Danmu comments when watching videos together, while two occasionally encountered videos that curated Danmu comments in other videos. However, all participants agreed that "*the opportunity to get human support is limited*" (P6). Overall, the lack of accessible support for Danmu excludes participants from fully enjoying Danmu and engaging in the social discussions.

*3.2.3 **Three Key Challenges***. During the co-watching exercise, participants primarily read Danmu while watching videos because "*it feels like a live discussion happening alongside the video*" (P3). All participants described the design probe as "*better than any existing tools*" (P1) because "*it allows us to access what is otherwise mostly hidden*" (P7). However, they still faced three significant challenges in fully engaging with Danmu:

**First, Danmu comments often lack descriptions of the visual context, making it difficult for BLV viewers to fully understand the discussions.** For example, four participants were confused by the comment, "*Is she living in summer or winter?*", which humorously referred to an actress wearing both a winter scarf and a summer T-shirt. However, BLV viewers found it confusing because the visual context was inaccessible. This lack of visual context hindered their ability to grasp the topics and "*achieve the same level of understanding as sighted viewers*" (P1).

**Second, the speech interference between Danmu comments and videos prevents BLV viewers from enjoying both content simultaneously.** Although participants wanted to access Danmu while watching videos, the comments often overlapped with the original speech, making it "*difficult to grasp both at once*" (P2). Furthermore, Danmu comments often appeared incoherent with the surrounding video speech. For example, P1 noted, "*Many comments seemed irrelevant to the original video, making them more like a noise.*" This lack of coherence resulted in an experience that was "*obtrusive*" (P4) and "*distracting*" (P7).

**Third, Danmu comments are not structured for sequential access, making it tedious for BLV viewers to follow audience discussions.** Danmu comments typically come in large volumes and appear on screen simultaneously. However, when presented sequentially by screen readers, they become disorganized and shift rapidly between topics (e.g., "*Any one still using this phone?; The one saying 'Kidney', don't go!; Look at the manual on the table!; ...*"). This disorganization of comments made it challenging for participants to follow audience discussions, leading to an experience that felt "*more like a reading test rather than an enjoyable chit-chat*" (P8), ultimately preventing them from connecting with other viewers.

Overall, while Danmu serves as a valuable source of social engagement, participants encountered significant challenges due to the lack of visual context, the speech interference between comments and videos, and the disorganization of comments. To address these challenges, we designed DanmuA11y.

---

[7]Participants' original speech was translated from Mandarin to English.

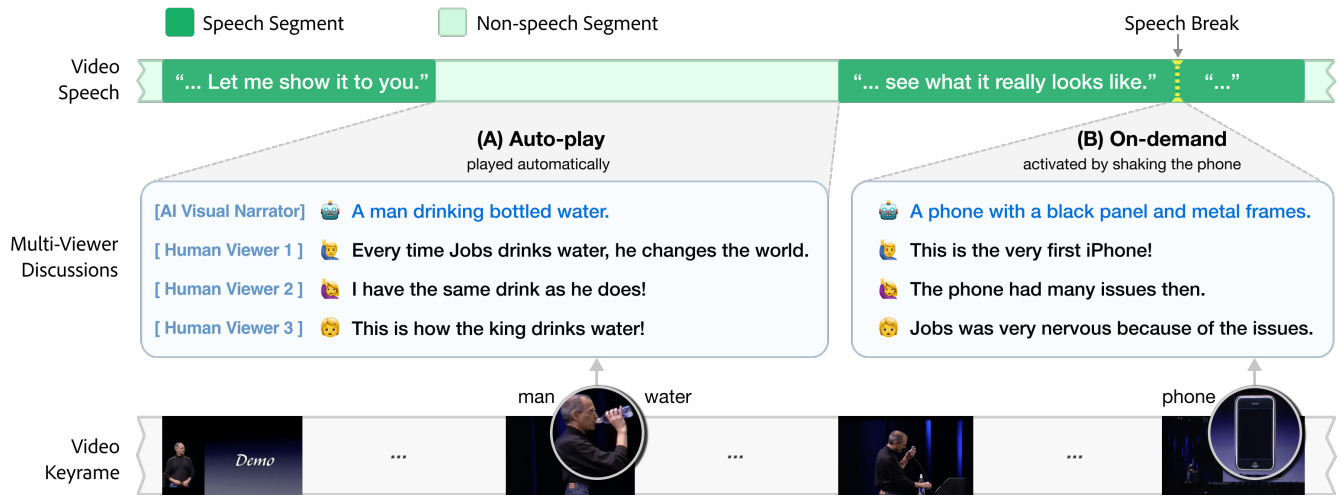[8]https://support.apple.com/en-us/111799

Figure 4: An example walk-through of DanmuA11y: (A) The system automatically plays multi-viewer discussions during non-speech segments, using spatial audio and varied tones to create a co-watching experience. (B) The system provides on-demand access at speech breaks, notifying users with a sound and allowing them to access Danmu comments by shaking their phones.

## 4 SYSTEM

DanmuA11y makes Danmu accessible by transforming Danmu into multi-viewer discussions. It incorporates three core features:

(1) *Augmenting Danmu with Visual Context*: To help BLV viewers understand the visual context related to audience discussions, DanmuA11y supplements Danmu comments with descriptions of the visual context.

(2) *Seamlessly Integrating Danmu into Videos*: To help viewers enjoy Danmu while watching videos, DanmuA11y curates Danmu topics and optimizes their insertion timing within videos to avoid Danmu-speech overlaps.

(3) *Presenting Danmu via Multi-Viewer Discussions*: To create a co-watching experience, DanmuA11y organizes Danmu comments into dialogues and uses spatial audio to simulate the sensation of other viewers conversing around the user.

Powered by these features, DanmuA11y offers an enjoyable and socially engaging experience for BLV viewers. In the following sections, we first provide an example walk-through in Section 4.1. Next, we describe the system pipeline in Sections 4.2 - 4.6. Finally, we report the implementation details in Section 4.7.

### 4.1 User Walk-through

To demonstrate DanmuA11y, we follow Jason, who is watching a video of Steve Jobs launching the first iPhone (see Figure 4). He uses DanmuA11y to engage with the audience discussions in Danmu comments.

DanmuA11y **seamlessly integrates** Danmu comments with the video, ensuring they do not overlap with the video's speech. For instance, during Steve Jobs' speech, the system does not insert Danmu comments, allowing Jason to focus on the speech. When the system detects a non-speech segment in the video, such as after Steve says, "*... Let me show it to you*", it **automatically plays** a sequence of curated Danmu comments: "[Visual Description] *A man*

*is drinking bottled water.* [Viewer 1] *Every time Jobs drinks water, he changes the world.* [Viewer 2] *I have the same drink as he does! ...*" (see Figure 4 (A)). These comments are curated for their creativity and informativeness, organized to reflect interactions among multiple viewers, and preceded by a visual description to help Jason grasp the visual context. Hearing these comments, Jason feels like he is co-watching with other viewers, because the comments are delivered using spatial audio and varied human tones, creating the sensation of live discussions happening around him.

The system also provides **on-demand Danmu access** at speech breaks. To minimize disruption, the system plays a lightweight notification sound at detected speech breaks. For example, after Steve says, "*... You get to see what it really looks like*", Jason hears a notification sound and becomes curious about what others are discussing. To access these discussions, Jason simply shakes his phone to pause the video and listen to the comments: "[Visual Description] *A phone with a black front panel and metal frames.* [Viewer 1] *This is the very first iPhone!* [Viewer 2] *The phone had many issues then ...*" (see Figure 4 (B)). These comments help Jason gain a deeper understanding of this video moment. Once the comments conclude, the video auto-rewinds to the previous speech break and resumes at the start of a sentence, allowing Jason to refocus on the video easily.

### 4.2 Pipeline Overview

As shown in Figure 5, DanmuA11y incorporates three features. The first feature supplements Danmu topics with descriptions of the visual context. The second feature integrates Danmu topics into the video, employing an optimization method to maximize both topic quality and insertion quality. The third feature transforms Danmu topics into multi-viewer audio discussions. In the following, we first describe the video segmentation method that underpins these features, followed by a detailed introduction to each feature.
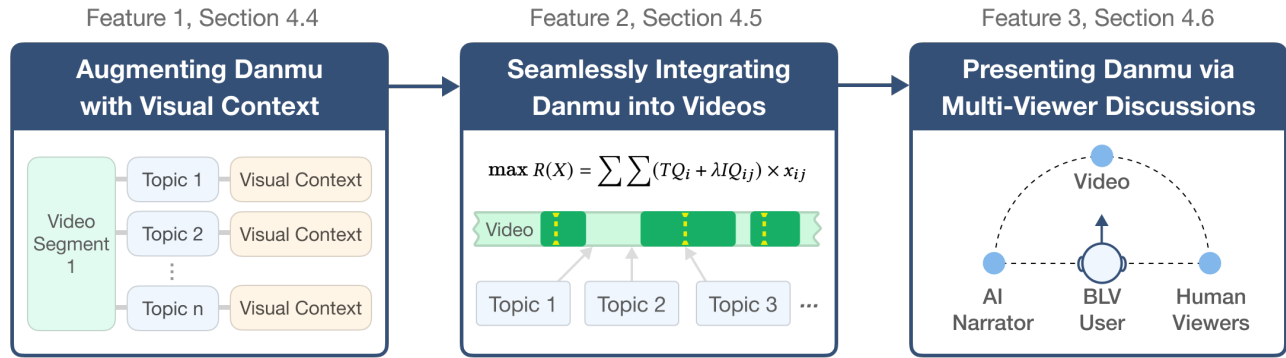
**Augmenting Danmu with Visual Context**

Video Segment 1

Topic 1 — Visual Context

Topic 2 — Visual Context

Topic n — Visual Context

**Seamlessly Integrating Danmu into Videos**

$$\max R(X) = \sum \sum (TQ_i + \lambda IQ_{ij}) \times x_{ij}$$

Video

Topic 1   Topic 2   Topic 3   ...

**Presenting Danmu via Multi-Viewer Discussions**

Video

AI Narrator    BLV User    Human Viewers

**Figure 5: The pipeline of DanmuA11y.**

## 4.3 Video Segmentation

DanmuA11y divides the video into two types of segments: non-speech and speech segments. This speech-based segmentation allows the system to identify proper insertion points for Danmu.

*4.3.1 Non-Speech Segment.* A non-speech segment is a continuous period free of human speech, providing a space for Danmu insertions without overlapping with the spoken content. These segments are identified using time-aligned transcripts generated by the Volcengine Auto-Subtitle API[9], which provides start and end times for each sentence with millisecond-level precision (e.g., "*00:00:01,010 –> 00:00:04,100 Hello, welcome to the video!*"). Following prior research [57], we classify any gap between sentences that lasts longer than two seconds as a non-speech segment. Furthermore, to prevent the insertion of Danmu in non-speech areas with high volumes (e.g., loud background music), we apply a one-second sliding window to identify regions where the root-mean-square volume exceeds a threshold of 0.8 using librosa[10]. These high-volume areas are marked as inappropriate for insertion.

*4.3.2 Speech Segment.* After identifying non-speech segments, we obtain the initial speech segments. These speech segments are further divided into smaller segments by speech breaks. A speech break is a specific point that separates continuous speech into relatively independent paragraphs. This allows the video to be paused without significantly disrupting the original speech flow. We identify these breaks using the natural language processing capabilities of GPT-4o [1], inputting consecutive sentences (i.e., without any non-speech segments in between) and prompting the model with "*Split this text into several segments, ensuring that each segment is relatively independent.*" This process typically identifies break points where the topic or idea changes (e.g., "*... That's our new product.* [Break point] *Now let's move on to the demo...*"). To avoid excessive pauses, each speech segment is constrained to have at least 20 words. The end of each speech segment, unless it is followed by a non-speech segment, is then marked as a speech break. The identified video segments and speech breaks are subsequently utilized by the three core modules.

## 4.4 Feature 1: Augmenting Danmu with Visual Context

Danmu comments often discuss specific visual elements that are inaccessible to BLV viewers, making it challenging for them to understand the discussion topics. To address this issue, DanmuA11y groups the comments into topics and supplements these topics with descriptions of the visual context.

*4.4.1 Grouping Danmu into Topics.* To identify the main topics in each video segment, we utilize a prompt-based topic modeling method [67]. Specifically, we feed GPT-4o all the Danmu comments[11] from a segment and prompt[12] it to "*Summarize the main topics discussed by viewers and categorize the comments under each topic.*" This approach generates a list of topics, each accompanied by a brief summary and their associated comments. Subsequently, each topic is re-input into GPT-4o with instructions to "*Remove redundant comments and rearrange the remaining ones to simulate a dialogue among multiple viewers.*" After this step, each video segment contains several topics, with each topic including a reordered list of comments that simulate a conversation among viewers.

*4.4.2 Generating Visual Descriptions.* To help BLV viewers understand the discussion topics, DanmuA11y generates visual descriptions for topics using GPT-4o [1, 33]. Specifically, the GPT-4o model is provided with the topic summaries and key frames for each video segment, and is prompted with the following instruction: "*Determine whether each topic discusses a visual object in the key frames. If so, describe the object's appearance in 15 words or fewer. If not, output 'None'.*" The key frames are pre-extracted from the video using SceneDetect[13], a content-aware scene detection model that divides the video into shots by comparing adjacent frames in the HSV color space. For each shot, the middle frame is selected as the representative key frame. These key frames are assigned to their respective video segments based on their timestamps. After this step, each Danmu topic includes a visual description (which may be empty) and a list of Danmu comments. Figure 7 (B) illustrates an example of the content within a topic.

---

[9]https://www.volcengine.com/
[10]librosa.org

[11]We use the video time associated with Danmu comments to assign them to their respective segments.
[12]The full prompts are provided in the supplementary materials.
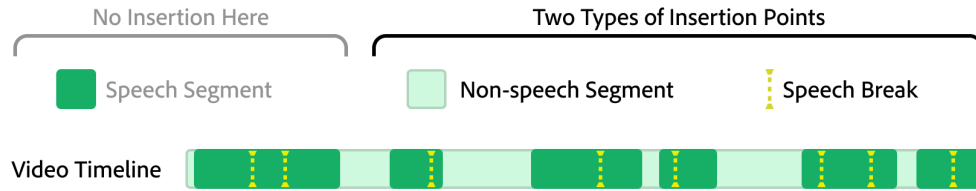[13]https://www.scenedetect.com/

Figure 6: DanmuA11y divides the video into non-speech and speech segments. The points between adjacent speech segments are called speech breaks.
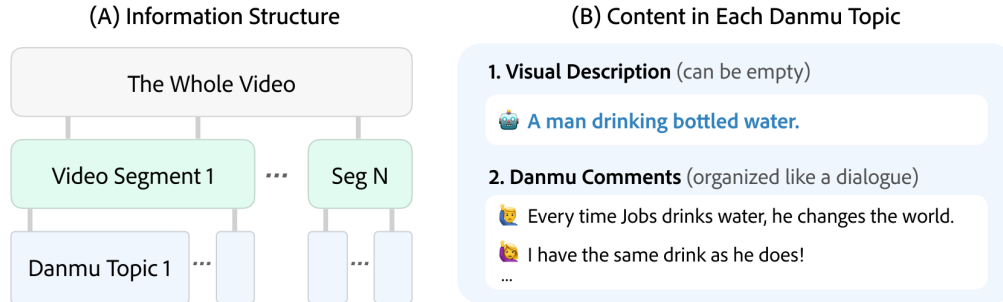


Figure 7: (A) Danmu topics belong to each video segment. (B) Each topic includes a visual description and a list of comments.

## 4.5 Feature 2: Seamlessly Integrating Danmu into Videos

To ensure viewers can enjoy Danmu while watching videos, DanmuA11y curates Danmu topics and determines their insertion points using an optimization algorithm. The algorithm optimizes two aspects: (1) **Topic Quality**: the quality of Danmu topics, and (2) **Insertion Quality**: the suitability of inserting Danmu topics at specific points in the video.

*4.5.1 Topic Quality*. DanmuA11y curates Danmu topics based on the three metrics derived from the formative study (as shown in Section 3.2.1): informativeness, creativity, and opinion diversity. Each metric is measured using established methods from prior works [53, 65, 84].

**Informativeness** assesses how much new information a topic adds to the original speech. Following [84], we quantify informativeness by measuring the semantic correlation between a Danmu topic and the transcript of the corresponding video segment. Specifically, we generate their vector embeddings using the *Universal Sentence Encoder* [9] and compute the cosine similarity *sim* between the two encoded vectors. The informativeness score $S_{\text{info}}$ is then defined as $1 - sim$, where a low semantic correlation results in a high informativeness score.

**Creativity** measures how creative a topic is. Given the subjective nature of creativity, we follow [65] by using GPT-4o for subjective ratings. The model is provided with the topic's content and instructed to "*Rate the creativity on a scale of 1 to 10, where 1 indicates completely mundane and 10 indicates exceptionally creative.*" This results in a rating of 8 for unusual descriptions like "*The snack looks like feet.*" and a rating of 2 for greetings like "*I'm here for a*

*second time!*". These ratings are then normalized to a 0-1 scale by dividing by 10, yielding the creativity score $S_{\text{creativity}}$.

**Opinion Diversity** measures the number of different opinions within a topic, as our formative study shows that topics with different opinions could provide a sense of heated discussion. To quantify opinion diversity [53], we apply sentiment analysis using the *XLM-T*[14] model [4], which classifies each comment as "*positive*", "*neutral*", or "*negative*". The diversity score $S_{\text{diversity}}$ is calculated by dividing the number of distinct sentiment labels by the total number of possible labels (i.e., three). A score of 1.0 indicates all three sentiments are present (e.g., "[positive] *My favorite snack;* [negative] *It is too sour;* [neutral] *I may give it a try; ...*").

**Topic Quality** *TQ* is defined as the weighted sum of the three metrics: $TQ = \lambda_i S_{\text{info}} + \lambda_c S_{\text{creativity}} + \lambda_d S_{\text{diversity}}$. Each metric is scaled from 0 to 1. We set $\lambda_i = \lambda_c = \lambda_d = 1$ to assign equal weight to each metric. The weights could be further personalized based on user preferences, which we discuss in Section 7.1.

*4.5.2 Insertion Quality*. Since Danmu topics are often related to nearby video moments, we restrict each topic's candidate set of insertion points to two locations before and two locations after the topic, including both non-speech segments and speech breaks. The quality of inserting a Danmu topic at each candidate point is measured based on two metrics: coherence and pause.

**Coherence.** Ideally, the inserted Danmu topic should be coherent with the original speech. We measure coherence using the method from [66], employing a pre-trained language model, *davinci-002*[15], to estimate the likelihood of word sequences. We combine the inserted topic $t_i$ with the surrounding speech at the insertion

---

[14]https://huggingface.co/cardiffnlp/twitter-xlm-roberta-base-sentiment
[15]*davinci-002* is a GPT base model. https://platform.openai.com/docs/models/gpt-base

## (A) Danmu-Video Integration

## (B) Optimization Formula

$$\textbf{Maximize } R(X) = \sum_{i=1}^{m} \sum_{j=1}^{n} (TQ_i + \lambda IQ_{ij}) \times x_{ij}$$

$$\textbf{Topic Quality } TQ = \lambda_i S_{\text{info}} + \lambda_c S_{\text{creativity}} + \lambda_d S_{\text{diversity}}$$

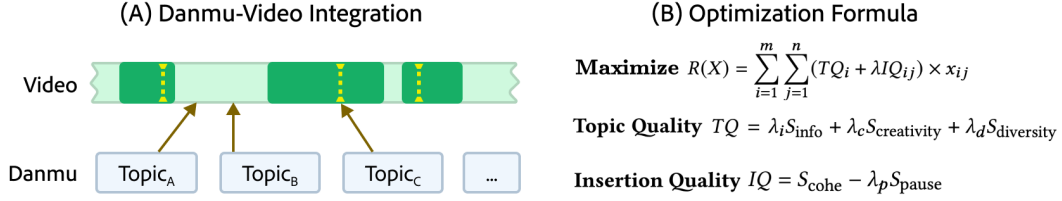$$\textbf{Insertion Quality } IQ = S_{\text{cohe}} - \lambda_p S_{\text{pause}}$$

**Figure 8: (A) The system places Danmu topics into insertion points. (B) It maximizes both topic quality and insertion quality.**

point $p_j$, input the combined text into the model, and calculate the sum log probability of $t_i$ at $p_j$. We apply min-max normalization to the log probabilities for each topic, and use the resulting normalized values as the coherence score $S_{\text{cohe}}$.

**Pause Penalty.** Long pauses at speech breaks can distract users from the original video. To minimize unnecessary pauses, we apply a penalty based on pause duration. The duration is estimated by dividing the word count in a topic by three, representing a typical speech rate of three Chinese words per second. This value is then normalized by the maximum pause duration at speech breaks (set to ten seconds in this work), resulting in the pause score $S_{\text{pause}}$.

**Insertion Quality $IQ$** is defined as: $IQ = S_{\text{cohe}} - \lambda_p S_{\text{pause}}$. We set $\lambda_p = 0.25$ to allow for pauses when necessary.

*4.5.3 Insertion Optimization.* When inserting Danmu topics into the video, our objective is to maximize both the topic quality $TQ$ and the insertion quality $IQ$, subject to the length constraints of insertion points. We formulate this as an optimization problem. We have a set of topics $\{t_i\}$, where $i = 1, \ldots, m$, each having a length $l_i$, and a list of insertion points $\{p_j\}$, where $j = 1, \ldots, n$, each having a maximum length $L_j$. Each topic has a topic quality score $TQ_i$ and there is an insertion quality score $IQ_{ij}$ for placing topic $t_i$ into point $p_j$. Our goal is to maximize the overall reward:

$$R(X) = \sum_{i=1}^{m} \sum_{j=1}^{n} (TQ_i + \lambda IQ_{ij}) \times x_{ij}$$

This optimization is achieved by determining the set of decision variables $X = \{x_{ij}\}$, where $x_{ij}$ represents whether topic $t_i$ is inserted into point $p_j$. The constraints are: (1) the total length of topics inserted into each point must not exceed the point's maximum length: $\sum_{i=1}^{m} l_i x_{ij} \leq L_j, \forall j = 1, 2, \ldots, n$; (2) each topic can be inserted at most once: $\sum_{j=1}^{n} x_{ij} \leq 1, \forall i = 1, 2, \ldots, m$; and (3) the decision variables $x_{ij}$ are binary: $x_{ij} \in \{0, 1\}, \forall i, j$. The optimal solution can be found by solving $X^* = \arg\max R(X)$ using integer linear programming. After obtaining the optimal solution $X^*$, we sort all the topics inserted at the same insertion point by their topic quality score in descending order.

In this work, the maximum length of each speech break is constrained to ten seconds to prevent overly long pauses. The length of each topic is dividing the word count by three, which reflects a typical speech rate of three Chinese words per second. The weight $\lambda$ is empirically set to 10 to prioritize insertion coherence.

Overall, this approach prioritizes inserting high-quality topics into non-speech segments first, followed by speech breaks, and ultimately discards low-quality topics if they do not fit within the available length constraints.

## 4.6 Feature 3: Presenting Danmu via Multi-Viewer Discussions

After inserting Danmu topics into the video, DanmuA11y transforms it into multi-viewer discussions using spatial audio and varied tones, aiming to create a socially engaging experience.

*4.6.1 Multi-Viewer Audio Discussions.* DanmuA11y transforms Danmu into audio discussions from two types of virtual viewers: (1) an *AI narrator* who delivers visual descriptions and (2) multiple *human viewers* who verbalize Danmu comments. These discussions are played via spatial audio, aiming to simulate people conversing around the user.

**Spatial Layout.** As shown in Figure 9 (A), we position the AI narrator on the left side and the human viewers on the right side, with the video placed in the front. This setup ensures that each sound source is distinct from the others. All three sound sources are set to equal volume. To avoid concurrent speech, only one virtual viewer speaks at a time.

**Alternating Tones.** We alternate the tones to simulate multi-person conversations. For each comment from human viewers, we randomly select a tone from three options and ensure that adjacent comments have different tones. A fourth tone is assigned to the AI narrator. The tones are synthesized using the Volcengine Text-to-Speech API[16].

**Auto-play Strategy.** During non-speech segments, the audio discussions are auto-played alongside the video. At speech breaks, the system offers on-demand access to the discussions using the following design.

*4.6.2 Easy On-demand Access.* When audio discussions are inserted at speech breaks, the system plays a notification sound to signal their presence, allowing users to access them selectively using a simple shake gesture.

**Notification Sound.** We use a half-second bubble sound[17] to indicate available discussions. The sound plays from the left if the discussions include AI visual descriptions or from the right if they only contain viewer comments. The video continues to play (i.e., it is not paused) while the sound is played.

**Shake gesture.** On hearing the notification sound, users can shake their phone within five seconds, which will pause the video and play the discussions. These discussions use the same spatial layout and tones as described earlier. When the comments end, the

---

[16]https://www.volcengine.com/product/tts. The tone IDs are BV002 for the AI narrator and BV011, BV064, and BV411 for human viewers.
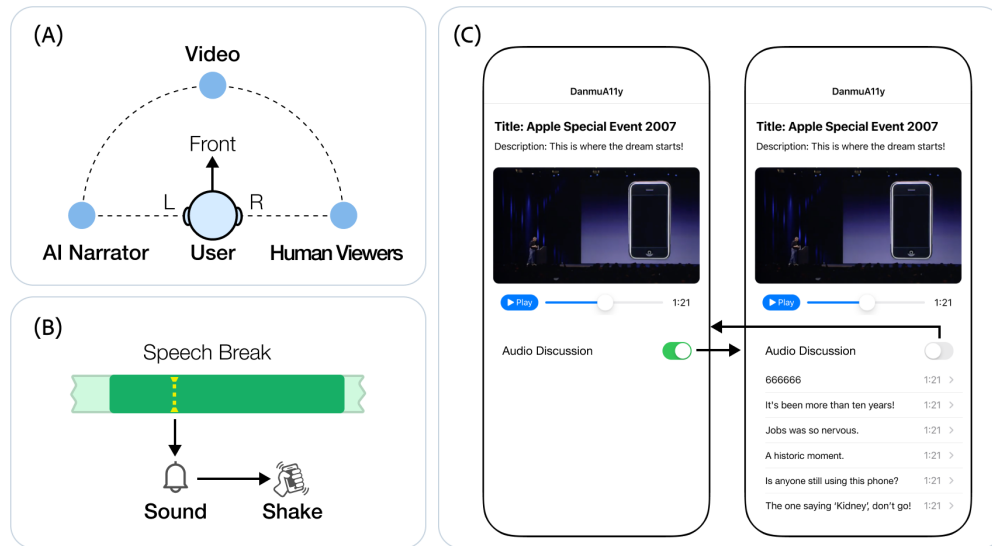[17]The sound is provided in the supplementary materials.

**Figure 9: DanmuA11y presents Danmu via multi-viewer discussions. (A) The system positions the AI narrator on the left side and the human viewers on the right side, with the video placed in the front. (B) The system plays a notification sound at speech breaks, allowing users to shake their phones to access Danmu on demand. (C) The Audio-Discussion toggle is enabled by default, allowing users to engage with the multi-viewer discussions. If the toggle is turned off, the system displays the original Danmu list.**

video auto-rewinds to the last speech break, resuming at the start of a sentence.

**Custom Controls.** Besides the shake gesture, the system provides video playback controls, including a play/pause button and a slider. There is an Audio-Discussion toggle that enables or disables audio discussions (both auto-played and on-demand). The toggle is on by default, allowing users to enjoy Danmu alongside the video. Users can turn it off to watch only the original video, such as when focusing on instrumental music. When the toggle is turned off, the system shows the original Danmu list (see Figure 9 (C)), which is accessible via screen readers.

## 4.7 Implementation

*4.7.1* **System**. The system is implemented as an iOS app for iPhone, tested on iPhone 15 with AirPods Pro (second generation). The phone-shake gesture is triggered when the device acceleration exceeds $6m/s^2$ (detected by CoreMotion API, excluding gravity). The audio playback is handled using AVFoundation.

*4.7.2* **Pipeline**. The video processing pipeline is implemented in Python. For modules utilizing GPT-4o, the temperature is set to 0.8. To solve the integer programming problem, we use the Gurobi Optimizer[18] via a Python interface, which completes within five seconds for 50 topics and 30 insertion points. The audio discussions are rendered using the Volcengine Text-to-Speech API at a speech rate between 1.1 and 1.2, adjusted to fit the duration estimated by word count. Spatial audio is generated via Pydub and exported in MP3 format.

---

[18]https://www.gurobi.com

*4.7.3* **Video Dataset**. We randomly selected 18 videos from Bilibili to represent six video types frequently watched by BLV viewers [43]: educational, comedic, tutorial, news, music, and film clips, with three videos for each type. None of the 18 videos were included in the formative study. Table 3 provides the video details. All videos had over 1,000 Danmu comments. We downloaded the latest 1,000 Danmu comments, the default maximum provided by the platform. The selected videos ranged from 24 seconds to seven minutes in length (mean = 3.8 minutes, SD = 1.7 minutes). The speech ratio (i.e., ratio of speech segments) in each video varied from 0% to 100% (mean = 59%, SD = 38%).

*4.7.4* **Technical Performance**. We ran the pipeline on the above 18-video dataset. None of the videos had been tested during the system's development. After processing the 18 videos, the pipeline yielded 706 Danmu topics. Based on the data, we report the technical performance regarding (1) visual description accuracy and (2) creativity scoring.

**(1) Visual Description Accuracy.** To evaluate the visual description accuracy, we randomly sampled 200 topics with visual descriptions, reviewed the topics and their associated key frames, and labelled whether the visual descriptions matched the visual elements being discussed. One researcher labeled all the data, and a second researcher reviewed the labels to ensure reliability. The accuracy rate was calculated by dividing the number of correct descriptions by the total number of sampled descriptions, yielding an accuracy rate of 92.5%. The errors mainly arose from difficulties in describing rapid visual changes. For example, in V11, a fast-paced human fight was mistakenly described as "*A person is performing CPR*" based on the comment "*CPR*". Additionally, incorrect associations between comments and visual elements were noted, such

as linking the comment "*He's the CEO!*" to the wrong person in the video. We discuss potential methods to mitigate these errors in Section 7.2.

**(2) Creativity Rating.** To assess how well GPT-4o rated creativity, we analyzed the distribution of creativity scores across all 706 topics. The scores ranged from 0.1 to 0.8 ($\mu = 0.408$, $\sigma = 0.139$). Specifically, 32.1% of the scores fell between 0.1 and 0.3, typically for greetings (e.g., "*I'm here!*") or simple emotions (e.g., "*Haha! So funny.*"). The majority, 54.3%, fell between 0.4 and 0.5, representing commonly seen comments (e.g., "*This cake looks tasty.*"). The remaining 13.6% scored between 0.6 and 0.8, mainly for unexpected descriptions (e.g., "*The snack looks like feet.*") and bloopers (e.g., "*Wait! He didn't even pay the bill!*"). This distribution aligns with the general observation that creative comments are not commonly found. One issue identified was that the model assigned scores lower than five for certain memes or puns (e.g., using common words to denote a humorous fact), which could potentially be improved by fine-tuning the model with specialized datasets [74].

The performance of other modules can be found in prior works [53, 66, 67, 84]. Additionally, we evaluated the overall Danmu curation quality and integration quality in the user evaluation study.

## 5 USER EVALUATION

We conducted a within-subject study with 12 BLV viewers to evaluate the effectiveness of DanmuA11y compared to a baseline. Specifically, we aim to answer the following research questions:

**(RQ1) System Usability**: How do BLV viewers perceive the usability of DanmuA11y?

**(RQ2) Danmu Comprehension**: How does DanmuA11y impact BLV viewers' Danmu comprehension?

**(RQ3) Video Viewing Experience**: How effectively does DanmuA11y provide a smooth viewing experience?

**(RQ4) Social Connection**: How does DanmuA11y affect BLV viewers' sense of co-watching with other viewers?

### 5.1 Participants and Materials

*5.1.1 Participants.* We recruited 12 BLV viewers (P9-P20; seven male, five female, Table 1 lists their demographics) who regularly watched videos on Danmu-enabled platforms. These participants were recruited from an online support community, with ages ranging from 22 to 41 (mean = 30.8, SD = 5.3). Eight participants were totally blind and four participants were legally blind with light perception. All participants had prior experience accessing Danmu on video platforms, with a viewing frequency ranging from two or three times a week to daily. None of the participants took part in the formative study. Additionally, all participants were iOS VoiceOver users, had experience listening to spatial audio via headphones, and had normal hearing.

*5.1.2 Apparatus.* Each participant used two systems: DanmuA11y and a baseline. The baseline was the design probe used in the formative study, which simulated their current practices. Both systems were run on an iPhone 15. While using each system,

participants wore AirPods Pro (second generation) with Personalized Spatial Audio[19] enabled. Each participant adjusted VoiceOver's speech rate and audio ducking according to their preferences.

*5.1.3 Videos.* To ensure a fair comparison between the two systems, we selected six videos (V3-V8) from our 18-video dataset (Table 3) and split them into two groups with comparable speech ratios, length, and content. Each group included three videos: an educational video (speech ratio > 95%, length ≈ 2 minutes, featuring people explaining the scientific principles behind everyday tricks), a comedic video (speech ratio ≈ 80%, length ≈ 5 minutes, featuring people humorously trying snacks), and a music video (speech ratio = 0%, length ≈ 5 minutes, featuring instrumental music). Additionally, participants watched another two videos (V1-V2) during the tutorial phase and ten videos (V9-V18) after the comparison study.

### 5.2 Design and Procedure

*5.2.1 Procedure.* We first collected participants' demographics and asked about their prior experiences with Danmu. They then received a 10-minute tutorial on both systems, using the same two videos (V1-V2).

After the tutorial, participants proceeded with the comparison study. Each participant used both systems to consume Danmu-embedded videos. To ensure study control, the order of systems and video groups was counterbalanced by creating four combinations, which were evenly assigned to the twelve participants. The videos were presented in random order. Participants were instructed to watch videos as they would in daily life (e.g., fast-forward, rewind, or pause at any time), to engage with Danmu as long as it did not disrupt their viewing experience, to **report any confusion** about Danmu, and to complete each video unless the system was frustrating. After watching each video, participants wrote a short **video summary**. Upon completing both systems, participants completed a questionnaire, with questions shown in Figure 10. A semi-structured interview was then conducted to gather participants' feedback.

After the comparison study, participants watched the remaining videos (V9-V18) using only DanmuA11y and were asked to think aloud about their viewing experiences. At the end of the study, participants were invited to share any suggestions or concerns they had regarding the system. The entire study lasted three hours and was conducted one-on-one, in-person. Participants were compensated approximately 34 USD in local currency for their time.

*5.2.2 Metrics.* We assessed both systems using comprehension metrics and subjective ratings.

**Comprehension Metrics.** To assess Danmu comprehension, we followed prior work [84], tallying the *instances of confusion* about Danmu reported by participants while viewing the video. For video comprehension, we adopted the method in [78], analyzing the *number of errors* present in participants' video summaries.

**Subjective Ratings.** We adapted established questionnaires [3, 29, 50, 61] to investigate the four research questions. For usability (RQ1), we assessed participants' willingness to use the system in the future, ease of use, and ease of learning using relevant questions

---

[19]This feature personalizes the spatial audio to each participant's head-related transfer function (HRTF). https://support.apple.com/en-us/102596

from the System Usability Scale [50]. For Danmu comprehension (RQ2), participants rated how easy it was to comprehend three key aspects of Danmu: the discussion topics, the visual context of those topics, and the interactions among viewers. To assess the video viewing experience (RQ3), we examined whether the Danmu-video integration caused minimal disruption, following the ambient display heuristics developed by Mankoff et al. [61]. For social connection (RQ4), we measured participants' sense of closeness with other viewers using the Inclusion of Other in the Self (IOS) Scale [3] and whether they felt as if they were co-watching videos with others [34]. All questions used a 7-point Likert scale.

*5.2.3 Analysis.* We recorded the study's audio, tracked participants' interaction histories, and collected their questionnaire responses. We used the Wilcoxon signed-rank test [86] to analyze the significance of the questionnaire responses. We tallied instances where participants expressed confusion about Danmu during video viewing. To analyze the video summaries written by participants, one researcher first shuffled the summaries for each video to obscure participant and system information. Another researcher, who was unaware of these details, identified and labeled any factual errors in the summaries. The interview audio was transcribed and categorized according to the research questions. Our findings are reported based on this analysis.

## 6 EVALUATION RESULTS

During the comparison study, a total of 36 trials (12 participants × 3 videos) were conducted for each system. Participants completed all 36 trials using DanmuA11y. However, seven participants did not complete a total of nine trials[20] when using the baseline due to frustration caused by the system. In the following, we report the results regarding the four RQs: system usability, Danmu comprehension, video viewing experience, and social connection.

## 6.1 System Usability (RQ1)

**All participants preferred using DanmuA11y over the baseline for accessing Danmu while watching videos.** Participants reported a significantly higher willingness to use DanmuA11y in the future compared to the baseline ($Z = -2.82$, $p < .01$). They found DanmuA11y significantly easier to use ($Z = -2.70$, $p < .01$), with reduced mental effort ($Z = -2.54$, $p < .05$) and physical effort ($Z = -2.83$, $p < .01$). All participants rated DanmuA11y the highest score for ease of learning, because the system incorporated familiar design elements they were accustomed to, such as audio discussions mimicking "*group calls*" (P11, P20) and notification sounds resembling "*new message alerts*" (P9, P17). Additionally, participants rated DanmuA11y on how easy it was to distinguish between three sound sources—the original video, the AI narrator, and the human viewers—on a scale from 1 (very hard) to 7 (very easy). The average rating of 6.92 indicates that DanmuA11y was effective in helping participants differentiate between these sources using spatial audio and varied tones.

**DanmuA11y was particularly praised for its minimal cognitive and physical demands.** Participants felt "*totally relaxed*"

(P16) when using DanmuA11y. They highlighted the Danmu insertion as "*seamless*" (P12) and praised the on-demand access to Danmu via the shake gesture as "*effortless*" (P19), "*convenient*" (P14), and "*flexible*" (P15), as it eliminated the need for "*multi-step, repetitive operations required by screen readers*" (P13). In contrast, the baseline made participants feel like they were "*multi-tasking*" (P20), with P12 explaining, "*I had to constantly tap the screen, digest disorganized information, manage concurrent speeches, and shift my focus between the video and comments. It took the enjoyment out of watching videos.*"

## 6.2 Danmu Comprehension (RQ2)

In terms of Danmu comprehension, DanmuA11y outperformed the baseline in both comprehension metrics and subjective ratings. The results are as follows.

*6.2.1 Comprehension Metrics.* **DanmuA11y significantly reduced confusion about Danmu reported by participants** compared to the baseline ($Z = -2.94$, $p < .01$). With the baseline, the twelve participants reported a total of 52 instances of confusion (1.44 times per video). These instances were caused by missing visual context (31 instances), unfamiliar Internet slang (6 instances, e.g., the name of video games), near-homophones (3 instances), and incomplete dialogues (12 instances). For example, P11 expressed confusion over the comment, "*The one mentioning 'meat', are you a devil?*" because preceding comments that provided necessary context were absent. In contrast, when using DanmuA11y, participants only reported nine instances of confusion (0.25 times per video) . These cases were caused by near-homophones (5 instances, e.g., P10 misheard "*hospital, 'yee-yuan'* " as its Chinese near-homophone "*music, 'yin-yweh'* ") and unfamiliar Internet slang (4 instances). These results demonstrate that DanmuA11y effectively reduced confusion by providing visual descriptions and organizing comments into dialogues among viewers.

**DanmuA11y improved participants' video comprehension.** After using DanmuA11y, participant-written video summaries contained significantly more words than the baseline ($t_{11} = 6.25$, $p < .01$)[21]. On average, participants wrote 41.8 words per video with DanmuA11y, compared to 28.1 words with the baseline. Moreover, participants made significantly fewer factual errors in the video summaries with DanmuA11y ($Z = -2.27$, $p < .05$). When using the baseline, they made a total of 14 factual errors (0.33 times per video), of which 11 were visual errors resulting from "*guessing visuals using comments as clues*" (P9). In contrast, with DanmuA11y, the total number of errors were reduced to just two instances (0.06 times per video). Participants attributed their improved video comprehension to the inclusion of visual descriptions (mentioned 11 times) and their ability to focus more attentively on the videos (seven times).

*6.2.2 Subjective Ratings.* **Participants reported that DanmuA11y significantly enhanced Danmu comprehension** across three key aspects: discussion topics ($Z = -2.69$, $p < .01$), visual context ($Z = -2.83$, $p < .01$), and interactions among viewers ($Z = -2.70$, $p < .01$). Participants noted that organizing comments

---

[20]The uncompleted trials are listed in Table 5.

[21]We used a paired t-test for significance analysis, as the data satisfied the assumptions of normality and homogeneity of variances according to the Shapiro-Wilk Test and Levene's Test, respectively.
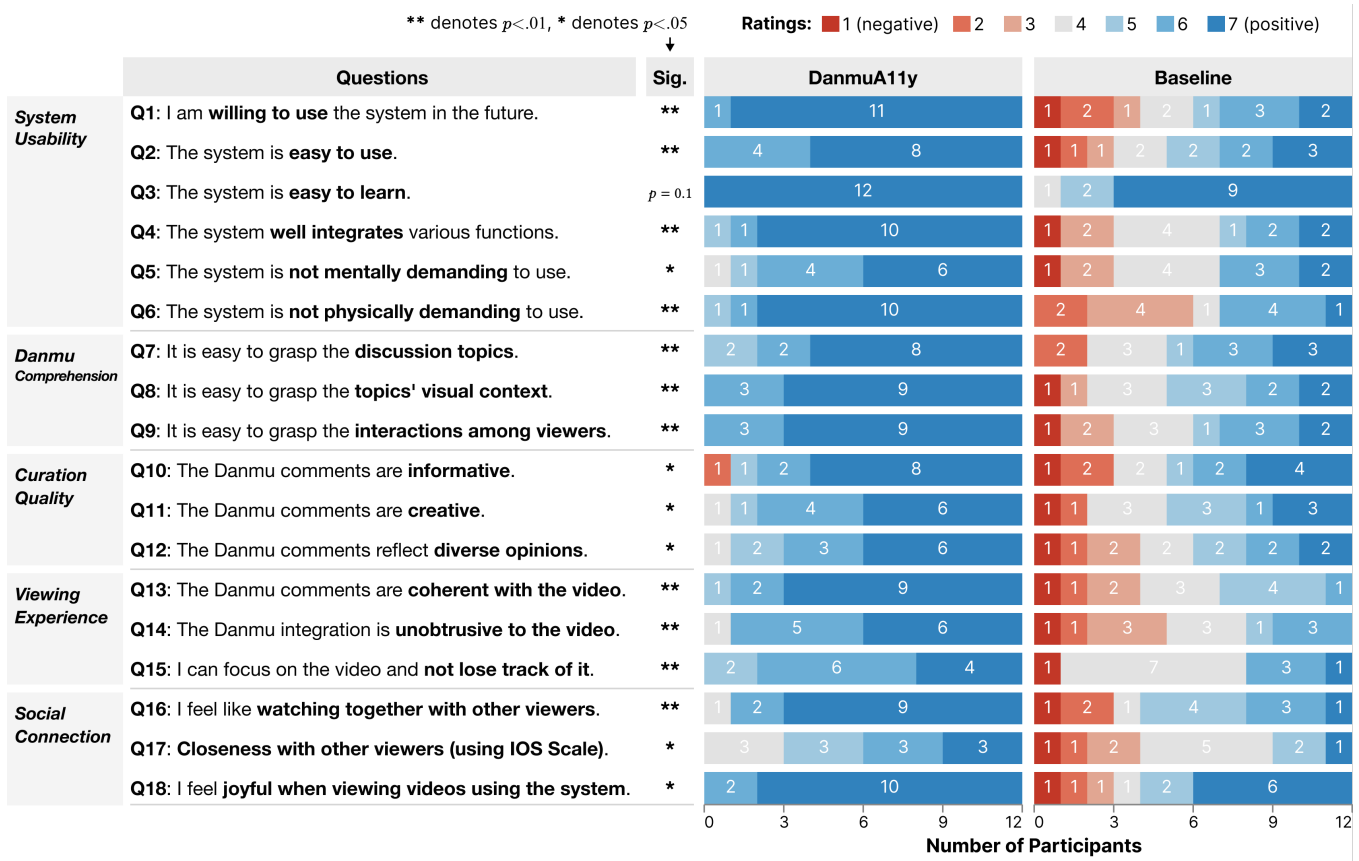
Figure 10: Distributions of the ratings for DanmuA11y and the baseline (1=strongly negative, 7=strongly positive). The asterisks indicate the statistical significance as a result of the Wilcoxon signed-rank test. Detailed statistics are provided in Table 4.
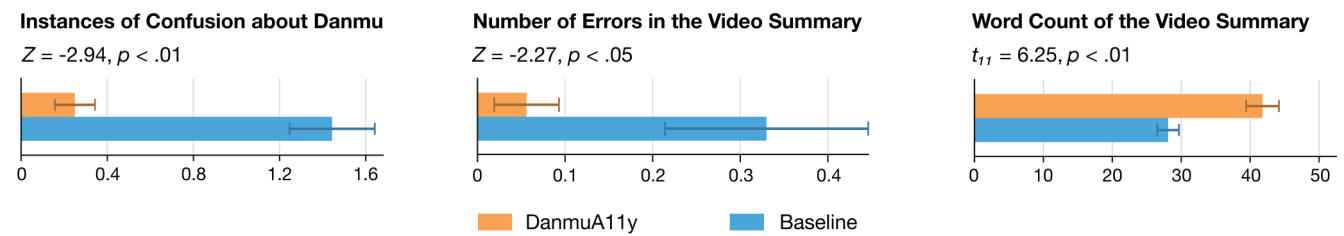


Figure 11: Comprehension metrics in the evaluation study. These metrics were calculated as averages for each participant per video. Error bars indicate standard errors.

into dialogues "*made the comments much easier to understand*" (P13) and helped them "*follow the conversations*" (P12).

**Participants highlighted that visual descriptions were essential for deeply engaging with Danmu.** For example, after watching V8, P10 remarked, "*Without the description, 'a cat lying on the musical instrument', I would've been quite confused. Why did people mention cats in a music video? Thanks to that description, I didn't feel excluded from the discussion.*" Participants appreciated how AI-generated descriptions complemented viewer comments, offering a more comprehensive understanding. As P12 observed,

"*When others said 'It's like eating concrete', the AI description promptly addressed my curiosity: 'a plate of brown bricks'. These pieces came together into a full picture.*" Additionally, participants noticed errors in visual descriptions when they conflicted with the video's audio, such as when the description "*A man dancing*" contradicted the audio "*When you walk outdoors*" in V6. However, they did not notice errors when there was no such conflict, which aligns with issues reported in previous research [37, 78]. To address this issue, future work should focus on enhancing visual description accuracy or providing descriptions from multiple sources [48].

**DanmuA11y presented comments of significantly higher quality compared to the baseline.** Participants found the comments in DanmuA11y to be significantly more informative ($Z = -2.21$, $p < .05$), creative ($Z = -2.56$, $p < .05$), and reflective of diverse opinions ($Z = -2.51$, $p < .05$). For example, P10 found the comment "*(the candy bar is) like my dog's chew toy*" to be both creative and helpful in visualizing the snack. Similarly, P11 noted, "*It's so funny when someone says, 'Coach, I want to learn', and another person replies, 'Coach: Why don't you learn everything?' It (DanmuA11y) really helped me discover these creative interactions.*" In contrast, participants described the comments in the baseline as "*chaotic*" (P15) and "*messy*" (P17), which made it challenging to "*get key points and find highlights*" (P18). These results indicate that DanmuA11y effectively curated content according to the three criteria. Additionally, participants suggested incorporating more curation criteria to better personalize the content based on their preferences, which we discuss in Section 7.1.

## 6.3 Video Viewing Experience (RQ3)

**DanmuA11y offers a smooth video viewing experience.** Participants rated the Danmu comments in DanmuA11y as significantly more coherent with the video ($Z = -3.08$, $p < .01$) and unobtrusive ($Z = -2.82$, $p < .01$) compared to the baseline. They described the viewing experience as "*smooth*" (P16) and noted that "*the seamless blend of Danmu and the videos creates a cohesive piece of work*" (P10). Participants also appreciated that the notification sound during speech breaks was "*unobtrusive*" (P11) and "*natural*" (P9), as "*the sound didn't disrupt the speech flow*" (P13, P18).

**The Danmu-video integration in DanmuA11y resulted in minimal split attention.** Participants reported that DanmuA11y allowed them to stay more focused on the video and not lose track of it ($Z = -2.75$, $p < .01$), thanks to the absence of overlapping speech (twelve times), the well-timed comments (nine times), and the auto-rewinding feature (ten times, described in Section 4.6.2). As P10 explained, "*The video automatically rewound to the point where the bubble sound occurred, not where I shook the phone. This is ingenious because it allowed me to refocus on the video effortlessly, without the trouble of manually adjusting the progress bar.*"

**DanmuA11y enabled participants to easily access high-quality comments while watching videos.** On average, participants accessed more comments using DanmuA11y (51 comments per video) compared to the baseline (33 comments per video). This improvement was attributed to DanmuA11y's simple interactions (P12) and high-quality content (P17), which fostered user engagement with comments. In contrast, the baseline required repetitive and tiring operations, preventing participants from accessing comments. For example, participants had to repeatedly tap the screen every few seconds to "*keep up with the comment list*" (P14). These repetitive actions caused fatigue within one to two minutes, leading participants to either stop reading the comments or end the trial early. P12 succinctly expressed this frustration: "*The frequent screen reader operations made me tired. Plus, the comments were so chaotic that I couldn't enjoy the video.*" These findings highlight the importance of enabling easy access to high-quality comments to provide an enjoyable viewing experience.

**Participants exhibited varied usage patterns with DanmuA11y, indicating directions for personalization.** DanmuA11y offers user control through several features: on-demand Danmu access via shake gestures, an Audio-Discussion toggle, and video playback controls. We observed personalized usage patterns for these features:

**(1) On-demand Danmu Access**: On average, participants responded to 87.4% (SD=9.2%) of Danmu notifications by performing shake gestures, indicating that they accessed most on-demand comments. Participants typically ignored new comments when previous ones did not match their preferences. For instance, P15 skipped the last two notifications in V3, explaining, "*I found previous comments in this video less informative than I wanted, so I skipped the last ones.*" This finding highlights the opportunity to refine comment curation according to individual preferences.

**(2) Audio-Discussion Toggle**: Participants primarily used the audio-discussion toggle to control their engagement with audience discussions. For example, nine participants disabled the toggle during the chorus of an instrumental music video (V7) to "*focus on the best part of the music*" (P9) and re-enabled it afterward to "*see what other people think of it*" (P14). This behavior suggests the potential for tailoring Danmu presentation based on video types—such as minimizing comments during key moments in music videos or increasing discussions for controversial news videos.

**(3) Playback Control**: Participants used playback controls according to their daily habits. Eight participants did not manually rewind videos, while the other four participants rewound an average of 1.42 times per video. When participants rewound, their goals included "*revisiting funny comments*" (P16) and "*double-checking video details such as food names*" (P20). These varied usage patterns highlight the importance of providing flexible controls to accommodate individual needs, which we further discuss in Section 7.1.

## 6.4 Social Connection (RQ4)

In terms of social connection, DanmuA11y significantly enhanced participants' sense of social presence, connection, and enjoyment compared to the baseline. The results are as follows.

**DanmuA11y significantly enhanced the sense of watching together with other viewers** as indicated by higher ratings of social presence compared to the baseline ($Z = -2.83$, $p < .01$). As P11 noted, "*It really felt like there's a group of people watching with me.*" Participants attributed this co-watching experience to several features of DanmuA11y, particularly spatial audio and alternating tones. P18 remarked, "*The alternating tones on my left and right immediately gave me the impression of several people conversing nearby.*" This finding aligns with prior research [39], which indicates that spatial audio improves social presence in remote communication. Moreover, the dialogue-like organization of the comments further contributed to the shared viewing experience. P11 described the comment organization as "*a group of people sharing their funny thoughts about videos*", which made the experience "*truly feel like co-watching with friends*" (P14).

**This co-watching experience also increased participants' sense of connection with other viewers**, with higher ratings of closeness compared to the baseline ($Z = -2.55$, $p < .05$). Participants reported feeling deeply connected with other viewers,

especially when the comments reflected their own thoughts or interests. For instance, P12 noted, "*It's a pleasant surprise when someone said exactly what I was thinking.*" Similarly, P18 mentioned, "*I'm really happy when I hear others mention my favorite movies—It's like we're emotionally in sync.*" Such emotional resonance not only made participants feel as if they were "*meeting someone who understands me well*" (P11), but also "*added extra joy to the viewing experience*" (P17).

**Consequently, participants reported greater enjoyment when using DanmuA11y** compared to the baseline ($Z = -2.23$, $p < .05$). In addition to "*the emotional resonance with other viewers*" (P11), participants highlighted sources of joy, such as "*the creative comments*" (P9) and "*the co-watching atmosphere*" (P16: "*Watching with others makes me happy. I'm no longer alone; I'm accompanied.*"). Moreover, DanmuA11y enhanced user enjoyment by "*filling the awkward silence*" (P16). As P16 explained, "*Previously, when I didn't get a joke in videos, it felt boring and awkward. But now, other people actually fill that awkward silence and make me laugh.*" Overall, the results indicate that by transforming Danmu into multi-viewer discussions, DanmuA11y effectively created an enjoyable and socially engaging experience for BLV video viewers.

## 7 DISCUSSION

Our formative study revealed three key challenges that hinder BLV viewers' engagement with Danmu: the lack of visual context, the speech interference between comments and videos, and the disorganization of comments. To address these challenges, we designed DanmuA11y with three core features: augmenting Danmu with visual context, seamlessly integrating Danmu into videos, and presenting Danmu via multi-viewer discussions. User evaluations demonstrated that DanmuA11y effectively improved Danmu comprehension, offered smooth video viewing experiences, and fostered social connections for BLV viewers. In the following sections, we discuss: (1) directions for personalizing DanmuA11y, (2) generalizability of DanmuA11y across videos, (3) integration of DanmuA11y with existing accessibility techniques, and (4) implications for enhancing commentary accessibility in broader contexts.

### 7.1 Personalization of DanumA11y

Based on the evaluation study, we identify several directions for tailoring DanmuA11y to different user contexts:

**Personalize the weights of curation criteria.** While DanmuA11y assigns equal importance to its three curation criteria, participants expressed individualized preferences. For example, P15 favored informative comments, while P16 preferred humorous or creative comments. To better align with personal preferences, future systems could allow users to adjust the weight of each criterion, or adapt to user preferences by learning from users' interaction histories [32].

**Tailor curation to different video types.** Participants' preferences for comments also varied depending on the video type. For instance, humorous comments were favored for comedic content (V3-V4), diverse opinions were valued for news videos (V9-V10), and informative comments were preferred for action movies with complex visuals (V11-V12). This suggests the potential to fine-tune comment curation based on the video types.

**Incorporate user-defined content filters.** Participants proposed adding more filtering options, such as avoiding spoilers (P20), highlighting replies to the video creators' questions (P14), prioritizing comments posted by creators (P19), and emphasizing comments with many likes (P12). Future systems could better address diverse user needs by supporting user-defined filters [42] or question-answering interactions [2, 47].

**Enable flexible control in daily scenarios.** Currently, DanmuA11y relies on phone-shake gestures for on-demand access. To improve usability in everyday situations, future systems could support alternative input methods, such as head gestures using earphones [94, 95], or hand gestures with hand-worn devices [27, 93].

**Support acoustic feature customization.** Participants suggested customizable acoustic features, such as adjusting the speech rate, tone, comment density, and adding reverberation to simulate different environments (e.g., a theater). Additionally, future work could explore more diverse auditory representations [28, 62], such as utilizing dynamic sound sources to simulate the presence of moving viewers.

### 7.2 System Generalizability Across Videos

DanmuA11y was evaluated using 18 videos with diverse visual styles and speech ratios. Based on the evaluation results, we reflect on the generalizability of DanmuA11y and identify potential directions for improvement.

#### 7.2.1 *AI Visual Description.* The 18 videos covered a wide range of visual styles, from nearly static imagery (e.g., V18) to rapidly changing scenes (e.g., V11). When generating visual descriptions, the pipeline achieved an accuracy of 92.5%, suggesting that the AI descriptions were accurate for most cases. However, there are several areas for improvement.

**The first area for improvement is describing complex visual changes**, such as the fight scenes in V11 and V12. Because the current pipeline uses key frames for generating descriptions, it cannot capture movement over time and can lead to incorrect descriptions of human actions. To address these issues, future work could incorporate recent advances in video-to-language models [17, 44], utilize vision-language pre-training methods [83, 97], or explore specialized pipelines, such as reconstructing 3D meshes of rapidly changing objects and describing their movements [21].

**The second area for improvement is filtering out similar descriptions across topics**. Our pipeline generates different descriptions when Danmu topics discuss different visual elements (e.g., hair, clothing, or facial expressions in V18). However, when the topics center on the same visual element, the descriptions become similar. Future systems could address this by removing similar descriptions through methods like semantic comparison [69].

**Additionally, the coverage of AI descriptions could be enhanced.** The current AI descriptions aim to provide visual context for comments but may not capture the most salient visual objects. They are also not as comprehensive as audio descriptions. To enhance the coverage of AI descriptions, future work could use saliency detection methods [55] or integrate DanmuA11y with existing audio descriptions, which we discuss in Section 7.3.

*7.2.2* **Danmu-Video Integration.** The current pipeline restricts Danmu placement to non-speech segments or speech breaks. For videos with few speech segments (e.g., V9, with 100% speech), most comments were added at speech breaks, which could result in frequent notifications. To address this, future work could explore extending the video timeline with AI-generated background music [18] to create more spaces for Danmu insertion.

*7.2.3* **Audio Presentation.** DanmuA11y utilizes spatial audio and varied human tones to present Danmu as multi-viewer discussions. However, some videos may have used similar spatial-audio layouts or speech tones, which may lead to confusion. Future work could explore detecting similar audio [52] and adaptively adjusting the audio presentation (e.g., using other tones) to ensure clear differentiation between Danmu and videos.

## 7.3 Integration with Existing Techniques

Based on observations and user feedback from the evaluation study, we discuss how DanmuA11y could be further integrated with existing accessibility techniques.

*7.3.1* **Audio Descriptions.** Participants in our study valued the AI visual descriptions for providing additional visual context. However, some participants expressed a preference for professional audio descriptions (AD) due to their higher accuracy and better alignment with the video content. To address this, DanmuA11y could be enhanced by integrating ADs as part of the video input. Based on user feedback and video accessibility guidelines [66], we suggest three considerations for this integration: (1) *Avoiding redundant descriptions*: the AI visual descriptions should provide provide new, complementary information rather than repeating the ADs. (2) *Minimizing speech interference*: Comments should not overlap with speech from the video or the ADs. This can be achieved by placing comments at speech breaks or by extending the video timeline with looping background music, a method accepted by BLV viewers [66]. (3) *Distinguishing different content*: To help users differentiate between the video, ADs, and Danmu comments, these elements could be optionally presented using distinct spatial positions or varying audio tones.

*7.3.2* **Screen Readers.** In our formative study, participants reported that using screen readers to access the fast-scrolling list of comments was tedious, as it required repetitive touch input to retrieve the latest comments. To address this, DanmuA11y introduced a shaking feature, enabling users to access curated comments more easily. User evaluations showed that all participants found this feature more effortless and convenient compared to screen reader access. However, one participant (P17) noted that the feature lacked the flexibility to navigate between comments, which is particularly useful for "*skipping or revisiting comments*". To provide users with flexible access to the curated comment list, future systems could consider also presenting the list through screen readers. This approach would allow viewers to easily access on-demand comments by shaking their phones while maintaining the agency to navigate between curated comments via screen reader controls.

## 7.4 Broader Implications for Commentary Accessibility

By addressing Danmu accessibility, our research offers insights for improving accessibility in broader contexts, especially in scenarios involving real-time or time-synced comments.

For **real-time commentary in live-streaming** [45, 51], prior research indicates that BLV users face challenges in keeping up with rapidly updating comments and understanding the discussions among other viewers [45]. Our design insights suggest potential directions for enhancing commentary accessibility in live-streaming: (1) supplementing comments with visual context, such as descriptions of relevant visual objects in the stream; (2) curating comments based on users' information preferences; (3) avoiding overlaps with the streamer's speech, for instance, by pausing comments when the streamer is speaking; and (4) using spatial audio and varied tones to create a co-watching experience.

To implement the proposed design, future work needs to address both **technical and human-factor challenges**. The key technical challenge lies in achieving real-time processing, which involves several aspects: generating visual descriptions (e.g., via real-time computer vision models [10, 16]), curating comments (e.g., with pre-trained sentence encoders [9, 69]), identifying non-speech periods (e.g., using human speech detection algorithms [68, 76]), and optimizing comment placements (e.g., sorting comments in descending order of quality). Regarding the human-factor challenge, live-streaming often involves two-way communications between streamers and their audiences [58]. This interaction might create unique information needs, such as prioritizing comments that respond to live-streamers' questions, filtering out abusive remarks [71], or removing promotional messages that disrupt the stream's flow [14]. Further research is necessary to identify these specific information needs.

In **video-based social media** platforms without Danmu, our multi-viewer discussion design could be adapted by using time-referenced comments [98] (e.g., "*2:50 this kid is soooo funny*") or by aligning traditional comments with the video timeline [85] to create a co-watching experience for BLV viewers. Moreover, DanmuA11y's key design insight—facilitating easy access to high-quality content—can be applied to broader contexts like audio descriptions. To balance limited video time with the need for detailed descriptions [64, 78], future audio descriptions could insert essential visual information into non-speech segments and provide supplementary explanations during speech breaks. Users could access these additional details through simple interactions like shake gestures, allowing for a seamless viewing experience while retaining the option for more content.

## 7.5 Limitations and Future Work

DanmuA11y is designed to make Danmu consumption accessible to BLV viewers. Future research should explore methods to support BLV users in composing and sharing their own Danmu comments, which may lead to new needs for Danmu consumption, such as curating comments that resemble their own posts. Our work focuses on improving video viewing experiences via auditory feedback. Future works could explore multi-sensory stimuli—such as vibrations [54, 92] or lights [96]—to further enhance user engagement. While

our work focuses on using Danmu for entertainment and social interaction, future studies could explore the accessibility of Danmu for other applications, such as enhancing online learning outcomes [49] or assisting BLV content creators in gathering viewer feedback [13]. In addition to using automatic approaches, future work could leverage social support to further enhance Danmu accessibility. For instance, providing automatic description suggestions to sighted viewers could encourage them to include visual descriptions in their comments, thereby benefiting BLV users.

The current pipeline is designed to process existing comments. To keep the system updated with new Danmu comments, it could classify new comments into existing topics [67] or periodically re-run the pipeline. The optimization algorithm does not account for interactions between topics, such as re-ordering them to form a coherent sequence, which could be considered in future work. Our use of GPT-4o for topic modeling [67], image captioning [33], and creativity rating [65] occasionally results in hallucinations and errors. Future work could address these issues by incorporating specialized models [6, 74, 79] or model advancements. Beyond generating visual descriptions to enhance Danmu accessibility, future research could explore how to leverage Danmu comments to improve video accessibility [84], such as using comments to generate creative audio descriptions [81]. Additionally, exploring the long-term impact of DanmuA11y on the social engagement of BLV viewers presents a promising direction for future research.

## 8 CONCLUSION

DanmuA11y is a system designed to improve the accessibility of Danmu for BLV viewers. It addresses three primary challenges: the lack of visual context, the speech interference between comments and videos, and the disorganization of comments. Through user evaluations, we demonstrated that DanmuA11y effectively improved Danmu comprehension, provided a smooth viewing experience, and fostered social connections among viewers. We also identified future directions for personalizing DanmuA11y to accommodate diverse user needs and distilled implications for improving commentary accessibility in video and live-streaming platforms. We hope this work offers valuable insights for enhancing social media accessibility and inspires researchers to develop tools that facilitate social engagement for BLV users.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).

[2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*. 2425–2433.

[3] Arthur Aron, Elaine N Aron, and Danny Smollan. 1992. Inclusion of other in the self scale and the structure of interpersonal closeness. *Journal of personality and social psychology* 63, 4 (1992), 596.

[4] Francesco Barbieri, Luis Espinosa Anke, and Jose Camacho-Collados. 2021. XLM-T: Multilingual language models in Twitter for sentiment analysis and beyond. *arXiv preprint arXiv:2104.12250* (2021).

[5] Virginia P Campos, Tiago MU de Araújo, Guido L de Souza Filho, and Luiz MG Gonçalves. 2020. CineAD: a system for automated audio description script generation for the visually impaired. *Universal Access in the Information Society* 19, 1 (2020), 99–111.

[6] Virgínia P Campos, Luiz MG Gonçalves, Wesnydy L Ribeiro, Tiago MU Araújo, Thaís G Do Rego, Pedro HV Figueiredo, Suanny FS Vieira, Thiago FS Costa, Caio C Moraes, Alexandre CS Cruz, et al. 2023. Machine generation of audio description for blind and visually impaired people. *ACM Transactions on Accessible Computing* 16, 2 (2023), 1–28.

[7] Shuxian Cao, Dongliang Guo, Lina Cao, Shuo Li, Junlan Nie, Amit Kumar Singh, and Haibin Lv. 2023. VisDmk: visual analysis of massive emotional danmaku in online videos. *The Visual Computer* 39, 12 (2023), 6553–6570.

[8] Marina Ramos Caro. 2016. Testing audio narration: the emotional impact of language in audio description. *Perspectives* 24, 4 (2016), 606–634.

[9] D Cer. 2018. Universal sentence encoder. *arXiv preprint arXiv:1803.11175* (2018).

[10] Ruei-Che Chang, Yuxuan Liu, and Anhong Guo. 2024. WorldScribe: Towards Context-Aware Live Visual Descriptions. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*. 1–18.

[11] Ruei-Che Chang, Chao-Hsien Ting, Chia-Sheng Hung, Wan-Chen Lee, Liang-Jin Chen, Yu-Tzu Chao, Bing-Yu Chen, and Anhong Guo. 2022. Omniscribe: Authoring immersive audio descriptions for 360 videos. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology*. 1–14.

[12] Si Chen, Haocong Cheng, Jason Situ, Desirée Kirst, Suzy Su, Saumya Malhotra, Lawrence Angrave, Qi Wang, and Yun Huang. 2024. Towards Inclusive Video Commenting: Introducing Signmaku for the Deaf and Hard-of-Hearing. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–18.

[13] Shuai Chen, Sihang Li, Yanda Li, Junlin Zhu, Juanjuan Long, Siming Chen, Jiawan Zhang, and Xiaoru Yuan. 2022. DanmuVis: Visualizing danmu content dynamics and associated viewer behaviors in online videos. In *Computer Graphics Forum*, Vol. 41. Wiley Online Library, 429–440.

[14] Xinyue Chen, Si Chen, Xu Wang, and Yun Huang. 2021. " I was afraid, but now I enjoy being a streamer!" Understanding the Challenges and Prospects of Using Live Streaming for Online Education. *Proceedings of the ACM on Human-Computer Interaction* 4, CSCW3 (2021), 1–32.

[15] Yue Chen, Qin Gao, and Pei-Luen Patrick Rau. 2017. Watching a movie alone yet together: understanding reasons for watching Danmaku videos. *International Journal of Human–Computer Interaction* 33, 9 (2017), 731–743.

[16] Tianheng Cheng, Lin Song, Yixiao Ge, Wenyu Liu, Xinggang Wang, and Ying Shan. 2024. Yolo-world: Real-time open-vocabulary object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16901–16911.

[17] Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, et al. 2024. VideoLLaMA 2: Advancing Spatial-Temporal Modeling and Audio Understanding in Video-LLMs. *arXiv preprint arXiv:2406.07476* (2024).

[18] Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre Défossez. 2024. Simple and controllable music generation. *Advances in Neural Information Processing Systems* 36 (2024).

[19] Juliet Corbin and Anselm Strauss. 2015. *Basics of qualitative research*. Vol. 14. sage.

[20] Khang Dang and Sooyeon Lee. 2024. Musical Performances in Virtual Reality with Spatial and View-Dependent Audio Descriptions for Blind and Low-Vision Users. In *The 26th International ACM SIGACCESS Conference on Computers and Accessibility*. 1–5.

[21] Ginger Delmas, Philippe Weinzaepfel, Thomas Lucas, Francesc Moreno-Noguer, and Grégory Rogez. 2022. Posescript: 3d human poses from natural language. In *European Conference on Computer Vision*. Springer, 346–362.

[22] Cole Gleason, Patrick Carrington, Lydia B Chilton, Benjamin Gorman, Hernisa Kacorri, Andrés Monroy-Hernández, Meredith Ringel Morris, Garreth Tigwell, and Shaomei Wu. 2020. Future research directions for accessible social media. *ACM SIGACCESS Accessibility and Computing* 127 (2020), 1–12.

[23] Cole Gleason, Patrick Carrington, Lydia B Chilton, Benjamin M Gorman, Hernisa Kacorri, Andrés Monroy-Hernández, Meredith Ringel Morris, Garreth W Tigwell, and Shaomei Wu. 2019. Addressing the accessibility of social media. In *Companion Publication of the 2019 Conference on Computer Supported Cooperative Work and Social Computing*. 474–479.

[24] Cole Gleason, Amy Pavel, Himalini Gururaj, Kris Kitani, and Jeffrey Bigham. 2020. Making GIFs Accessible. In *Proceedings of the 22nd International ACM

*SIGACCESS Conference on Computers and Accessibility.* 1–10.

[25] Cole Gleason, Amy Pavel, Xingyu Liu, Patrick Carrington, Lydia B Chilton, and Jeffrey P Bigham. 2019. Making memes accessible. In *Proceedings of the 21st International ACM SIGACCESS Conference on Computers and Accessibility.* 367–376.

[26] Cole Gleason, Amy Pavel, Emma McCamey, Christina Low, Patrick Carrington, Kris M Kitani, and Jeffrey P Bigham. 2020. Twitter A11y: A browser extension to make Twitter images accessible. In *Proceedings of the 2020 chi conference on human factors in computing systems.* 1–12.

[27] Yizheng Gu, Chun Yu, Zhipeng Li, Weiqi Li, Shuchang Xu, Xiaoying Wei, and Yuanchun Shi. 2019. Accurate and low-latency sensing of touch contact on any surface with finger-worn IMU sensor. In *Proceedings of the 32nd annual ACM symposium on user interface software and technology.* 1059–1070.

[28] João Guerreiro, Yujin Kim, Rodrigo Nogueira, SeungA Chung, André Rodrigues, and Uran Oh. 2023. The design space of the auditory representation of objects and their behaviours in virtual reality for blind people. *IEEE Transactions on Visualization and Computer Graphics* 29, 5 (2023), 2763–2773.

[29] Sandra G Hart. 2006. NASA-task load index (NASA-TLX); 20 years later. In *Proceedings of the human factors and ergonomics society annual meeting*, Vol. 50. Sage publications Sage CA: Los Angeles, CA, 904–908.

[30] Changyang He, Lu He, Tun Lu, and Bo Li. 2021. Beyond Entertainment: Unpacking Danmaku and Comments' Role of Information Sharing and Sentiment Expression in Online Crisis Videos. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW2 (2021), 1–27.

[31] Ming He, Yong Ge, Enhong Chen, Qi Liu, and Xuesong Wang. 2017. Exploring the emerging type of comment for online videos: Danmu. *ACM Transactions on the Web (TWEB)* 12, 1 (2017), 1–33.

[32] Steven CH Hoi, Doyen Sahoo, Jing Lu, and Peilin Zhao. 2021. Online learning: A comprehensive survey. *Neurocomputing* 459 (2021), 249–289.

[33] Yushi Hu, Hang Hua, Zhengyuan Yang, Weijia Shi, Noah A Smith, and Jiebo Luo. 2023. Promptcap: Prompt-guided image captioning for vqa with gpt-3. In *Proceedings of the IEEE/CVF International Conference on Computer Vision.* 2963–2975.

[34] Zeyu Huang, Xinyi Cao, Yuanhao Zhang, and Xiaojuan Ma. 2024. Sharing Frissons among Online Video Viewers: Exploring the Design of Affective Communication for Aesthetic Chills. In *Proceedings of the CHI Conference on Human Factors in Computing Systems.* 1–19.

[35] 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.

[36] Mina Huh, Yi-Hao Peng, and Amy Pavel. 2023. GenAssist: Making image generation accessible. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology.* 1–17.

[37] Mina Huh, Saelyne 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.* 1–17.

[38] Felix Immohr, Gareth Rendle, Christian Kehling, Anton Lammert, Steve Göring, Bernd Froehlich, and Alexander Raake. 2024. Subjective Evaluation of the Impact of Spatial Audio on Triadic Communication in Virtual Reality. In *2024 16th International Conference on Quality of Multimedia Experience (QoMEX).* IEEE, 262–265.

[39] Felix Immohr, Gareth Rendle, Annika Neidhardt, Steve Göring, Rakesh Rao Ramachandra Rao, Stephanie Arevalo Arboleda, Bernd Froehlich, and Alexander Raake. 2023. Proof-of-concept study to evaluate the impact of spatial audio on social presence and user behavior in multi-modal VR communication. In *Proceedings of the 2023 ACM International Conference on Interactive Media Experiences.* 209–215.

[40] Gaurav Jain, Basel Hindi, Connor Courtien, Xin Yi Therese Xu, Conrad Wyrick, Michael Malcolm, and Brian A Smith. 2023. Front Row: Automatically Generating Immersive Audio Representations of Tennis Broadcasts for Blind Viewers. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology.* 1–17.

[41] Mohit Jain, Nirmalendu Diwakar, and Manohar Swaminathan. 2021. Smartphone usage by expert blind users. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems.* 1–15.

[42] Shagun Jhaver, Quan Ze Chen, Detlef Knauss, and Amy X Zhang. 2022. Designing word filter tools for creator-led comment moderation. In *Proceedings of the 2022 CHI conference on human factors in computing systems.* 1–21.

[43] Lucy Jiang, Crescentia Jung, Mahika Phutane, Abigale Stangl, and Shiri Azenkot. 2024. "It's Kind of Context Dependent": Understanding Blind and Low Vision People's Video Accessibility Preferences Across Viewing Scenarios. In *Proceedings of the CHI Conference on Human Factors in Computing Systems.* 1–20.

[44] Peng Jin, Ryuichi Takanobu, Wancai Zhang, Xiaochun Cao, and Li Yuan. 2024. Chat-univi: Unified visual representation empowers large language models with image and video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 13700–13710.

[45] Joonyoung Jun, Woosuk Seo, Jihyeon Park, Subin Park, and Hyunggu Jung. 2021. Exploring the experiences of streamers with visual impairments. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW2 (2021), 1–23.

[46] Daniel Killough and Amy Pavel. 2023. Exploring Community-Driven Descriptions for Making Livestreams Accessible. In *Proceedings of the 25th International ACM SIGACCESS Conference on Computers and Accessibility.* 1–13.

[47] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics* 7 (2019), 453–466.

[48] 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.

[49] Yi-Chieh Lee, Wen-Chieh Lin, Fu-Yin Cherng, Hao-Chuan Wang, Ching-Ying Sung, and Jung-Tai King. 2015. Using Time-Anchored Peer Comments to Enhance Social Interaction in Online Educational Videos. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems* (Seoul, Republic of Korea) *(CHI '15).* Association for Computing Machinery, New York, NY, USA, 689–698. https://doi.org/10.1145/2702123.2702349

[50] James R Lewis. 2018. The system usability scale: past, present, and future. *International Journal of Human–Computer Interaction* 34, 7 (2018), 577–590.

[51] Franklin Mingzhe Li, Franchesca Spektor, Meng Xia, Mina Huh, Peter Cederberg, Yuqi Gong, Kristen Shinohara, and Patrick Carrington. 2022. "It Feels Like Taking a Gamble": Exploring Perceptions, Practices, and Challenges of Using Makeup and Cosmetics for People with Visual Impairments. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems.* 1–15.

[52] Juncheng Li, Wei Dai, Florian Metze, Shuhui Qu, and Samarjit Das. 2017. A comparison of deep learning methods for environmental sound detection. In *2017 IEEE International conference on acoustics, speech and signal processing (ICASSP).* IEEE, 126–130.

[53] Chengzhong Liu, Shixu Zhou, Dingdong Liu, Junze Li, Zeyu Huang, and Xiaojuan Ma. 2023. CoArgue: Fostering Lurkers' Contribution to Collective Arguments in Community-based QA Platforms. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems.* 1–17.

[54] Guanhong Liu, Tianyu Yu, Chun Yu, Haiqing Xu, Shuchang Xu, Ciyuan Yang, Feng Wang, Haipeng Mi, and Yuanchun Shi. 2021. Tactile Compass: Enabling Visually Impaired People to Follow a Path with Continuous Directional Feedback. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) *(CHI '21).* Association for Computing Machinery, New York, NY, USA, Article 28, 13 pages. https://doi.org/10.1145/3411764.3445644

[55] Nian Liu, Ni Zhang, Kaiyuan Wan, Ling Shao, and Junwei Han. 2021. Visual saliency transformer. In *Proceedings of the IEEE/CVF international conference on computer vision.* 4722–4732.

[56] Xingyu Liu, Patrick Carrington, Xiang'Anthony' Chen, and Amy Pavel. 2021. What makes videos accessible to blind and visually impaired people?. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems.* 1–14.

[57] Xingyu" Bruce" Liu, Ruolin Wang, Dingzeyu Li, Xiang Anthony Chen, and Amy Pavel. 2022. CrossA11y: Identifying Video Accessibility Issues via Cross-modal Grounding. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology.* 1–14.

[58] Zhicong Lu, Haijun Xia, Seongkook Heo, and Daniel Wigdor. 2018. You watch, you give, and you engage: a study of live streaming practices in China. In *Proceedings of the 2018 CHI conference on human factors in computing systems.* 1–13.

[59] Guangyi Lv, Kun Zhang, Le Wu, Enhong Chen, Tong Xu, Qi Liu, and Weidong He. 2019. Understanding the users and videos by mining a novel danmu dataset. *IEEE Transactions on Big Data* 8, 2 (2019), 535–551.

[60] Xiaojuan Ma and Nan Cao. 2017. Video-based evanescent, anonymous, asynchronous social interaction: Motivation and adaption to medium. In *Proceedings of the 2017 ACM conference on computer supported cooperative work and social computing.* 770–782.

[61] Jennifer Mankoff, Anind K Dey, Gary Hsieh, Julie Kientz, Scott Lederer, and Morgan Ames. 2003. Heuristic evaluation of ambient displays. In *Proceedings of the SIGCHI conference on Human factors in computing systems.* 169–176.

[62] Keenan R May, Brianna J Tomlinson, Xiaomeng Ma, Phillip Roberts, and Bruce N Walker. 2020. Spotlights and soundscapes: On the design of mixed reality auditory environments for persons with visual impairment. *ACM Transactions on Accessible Computing (TACCESS)* 13, 2 (2020), 1–47.

[63] Rosiana Natalie, Ruei-Che Chang, Smitha Sheshadri, Anhong Guo, and Kotaro Hara. 2024. Audio description customization. *arXiv preprint arXiv:2408.11406* (2024).

[64] Zheng Ning, Brianna L Wimer, Kaiwen Jiang, Keyi Chen, Jerrick Ban, Yapeng Tian, Yuhang Zhao, and Toby Jia-Jun Li. 2024. SPICA: Interactive Video Content Exploration through Augmented Audio Descriptions for Blind or Low-Vision Viewers. In *Proceedings of the CHI Conference on Human Factors in Computing Systems.* 1–18.

[65] Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology.* 1–22.

[66] Amy Pavel, Gabriel Reyes, and Jeffrey P Bigham. 2020. Rescribe: Authoring and automatically editing audio descriptions. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology.* 747–759.

[67] Chau Minh Pham, Alexander Hoyle, Simeng Sun, and Mohit Iyyer. 2023. TopicGPT: A prompt-based topic modeling framework. *arXiv preprint arXiv:2311.01449* (2023).

[68] Javier Ramırez, José C Segura, Carmen Benıtez, Angel De La Torre, and Antonio Rubio. 2004. Efficient voice activity detection algorithms using long-term speech information. *Speech communication* 42, 3-4 (2004), 271–287.

[69] N Reimers. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. *arXiv preprint arXiv:1908.10084* (2019).

[70] Ethan Z Rong, Mo Morgana Zhou, Zhicong Lu, and Mingming Fan. 2022. "It Feels Like Being Locked in A Cage": Understanding Blind or Low Vision Streamers' Perceptions of Content Curation Algorithms. In *Proceedings of the 2022 ACM Designing Interactive Systems Conference.* 571–585.

[71] Sanjiban Sekhar Roy, Akash Roy, Pijush Samui, Mostafa Gandomi, and Amir H Gandomi. 2023. Hateful sentiment detection in real-time tweets: An LSTM-based comparative approach. *IEEE Transactions on Computational Social Systems* (2023).

[72] Woosuk Seo and Hyunggu Jung. 2021. Understanding the community of blind or visually impaired vloggers on YouTube. *Universal Access in the Information Society* 20 (2021), 31–44.

[73] Woosuk Seo and Hyunggu Jung. 2022. Challenges and opportunities to improve the accessibility of YouTube for people with visual impairments as content creators. *Universal Access in the Information Society* 21, 3 (2022), 767–770.

[74] Luning Sun, Hongyi Gu, Rebecca Myers, and Zheng Yuan. 2023. A New Dataset and Method for Creativity Assessment Using the Alternate Uses Task. In *Bench-Council International Symposium on Intelligent Computers, Algorithms, and Applications.* Springer, 125–138.

[75] Zhida Sun, Mingfei Sun, Nan Cao, and Xiaojuan Ma. 2016. VideoForest: interactive visual summarization of video streams based on danmu data. In *SIGGRAPH ASIA 2016 symposium on visualization.* 1–8.

[76] S Gökhun Tanyer and Hamza Ozer. 2000. Voice activity detection in nonstationary noise. *IEEE Transactions on speech and audio processing* 8, 4 (2000), 478–482.

[77] Garreth W Tigwell, Benjamin M Gorman, and Rachel Menzies. 2020. Emoji accessibility for visually impaired people. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems.* 1–14.

[78] Tess Van Daele, Akhil Iyer, Yuning Zhang, Jalyn C Derry, Mina Huh, and Amy Pavel. 2024. Making Short-Form Videos Accessible with Hierarchical Video Summaries. In *Proceedings of the CHI Conference on Human Factors in Computing Systems.* 1–17.

[79] Ike Vayansky and Sathish AP Kumar. 2020. A review of topic modeling methods. *Information Systems* 94 (2020), 101582.

[80] Violeta Voykinska, Shiri Azenkot, Shaomei Wu, and Gilly Leshed. 2016. How blind people interact with visual content on social networking services. In *Proceedings of the 19th acm conference on computer-supported cooperative work & social computing.* 1584–1595.

[81] Agnieszka Walczak and Louise Fryer. 2017. Creative description: The impact of audio description style on presence in visually impaired audiences. *British Journal of Visual Impairment* 35, 1 (2017), 6–17.

[82] 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.

[83] Xiao Wang, Guangyao Chen, Guangwu Qian, Pengcheng Gao, Xiao-Yong Wei, Yaowei Wang, Yonghong Tian, and Wen Gao. 2023. Large-scale multi-modal pre-trained models: A comprehensive survey. *Machine Intelligence Research* 20, 4 (2023), 447–482.

[84] Yujia Wang, Wei Liang, Haikun Huang, Yongqi Zhang, Dingzeyu Li, and Lap-Fai Yu. 2021. Toward automatic audio description generation for accessible videos. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems.* 1–12.

[85] Jonatas Wehrmann, Camila Kolling, and Rodrigo C Barros. 2020. Adaptive cross-modal embeddings for image-text alignment. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 34. 12313–12320.

[86] Frank Wilcoxon, S Katti, Roberta A Wilcox, et al. 1970. Critical values and probability levels for the Wilcoxon rank sum test and the Wilcoxon signed rank test. *Selected tables in mathematical statistics* 1 (1970), 171–259.

[87] Qunfang Wu, Yisi Sang, and Yun Huang. 2019. Danmaku: A new paradigm of social interaction via online videos. *ACM Transactions on Social Computing* 2, 2 (2019), 1–24.

[88] Qunfang Wu, Yisi Sang, Shan Zhang, and Yun Huang. 2018. Danmaku vs. forum comments: understanding user participation and knowledge sharing in online videos. In *Proceedings of the 2018 ACM International Conference on Supporting Group Work.* 209–218.

[89] 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.

[90] Ying Xiang and Seong Wook Chae. 2022. Influence of perceived interactivity on continuous use intentions on the danmaku video sharing platform: Belonging-ness perspective. *International Journal of Human–Computer Interaction* 38, 6 (2022), 573–593.

[91] Shuchang Xu, Chang Chen, Zichen Liu, Xiaofu Jin, Lin-Ping Yuan, Yukang Yan, and Huamin Qu. 2024. Memory Reviver: Supporting Photo-Collection Reminiscence for People with Visual Impairment via a Proactive Chatbot. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology* (Pittsburgh, PA, USA) *(UIST '24).* Association for Computing Machinery, New York, NY, USA, Article 88, 17 pages. https://doi.org/10.1145/3654777.3676336

[92] Shuchang Xu, Ciyuan Yang, Wenhao Ge, Chun Yu, and Yuanchun Shi. 2020. Virtual Paving: Rendering a Smooth Path for People with Visual Impairment through Vibrotactile and Audio Feedback. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 4, 3, Article 99 (Sept. 2020), 25 pages. https://doi.org/10.1145/3411814

[93] Xuhai Xu, Jun Gong, Carolina Brum, Lilian Liang, Bongsoo Suh, Shivam Kumar Gupta, Yash Agarwal, Laurence Lindsey, Runchang Kang, Behrooz Shahsavari, et al. 2022. Enabling hand gesture customization on wrist-worn devices. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems.* 1–19.

[94] Yukang Yan, Yingtian Shi, Chun Yu, and Yuanchun Shi. 2020. Headcross: Exploring head-based crossing selection on head-mounted displays. *Proceedings of the ACM on interactive, mobile, wearable and ubiquitous technologies* 4, 1 (2020), 1–22.

[95] Yukang Yan, Chun Yu, Xin Yi, and Yuanchun Shi. 2018. Headgesture: Hands-free input approach leveraging head movements for hmd devices. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 2, 4 (2018), 1–23.

[96] Ciyuan Yang, Shuchang Xu, Tianyu Yu, Guanhong Liu, Chun Yu, and Yuanchun Shi. 2021. LightGuide: Directing Visually Impaired People along a Path Using Light Cues. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 5, 2, Article 84 (June 2021), 27 pages. https://doi.org/10.1145/3463524

[97] Sibei Yang, Guanbin Li, and Yizhou Yu. 2019. Cross-modal relationship inference for grounding referring expressions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition.* 4145–4154.

[98] Matin Yarmand, Dongwook Yoon, Samuel Dodson, Ido Roll, and Sidney S Fels. 2019. " Can you believe [1: 21]?!" Content and Time-Based Reference Patterns in Video Comments. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems.* 1–12.

[99] Mingrui Ray Zhang, Ruolin Wang, Xuhai Xu, Qisheng Li, Ather Sharif, and Jacob O Wobbrock. 2021. Voicemoji: Emoji entry using voice for visually impaired people. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems.* 1–18.

[100] Mingrui Ray Zhang, Mingyuan Zhong, and Jacob O Wobbrock. 2022. Ga11y: An automated gif annotation system for visually impaired users. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems.* 1–16.

# A SUPPLEMENTARY TABLES AND FIGURES

This section provides the following tables and figures.

- Table 1 shows participants' demographics.
- Table 2 lists the six videos used in the formative study.
- Table 3 and Figure 12 provide descriptions and screenshots of the 18 videos used in the evaluation study.
- Table 4 presents detailed statistics on the subjective ratings from the evaluation study.
- Table 5 outlines the uncompleted trials in the controlled comparison.

**Table 1: Participants' demographics. P1-P8 participated in the formative study, and P9-P20 participated the evaluation study.**

| PID | Age | Gender | Visual Condition | Danmu-enabled Video Platforms | Experience | Usage Frequency |
|-----|-----|--------|------------------|-------------------------------|------------|-----------------|
| P1  | 27  | F | Legally blind  | Bibibili, Douyin, Youku         | 4 years | Daily |
| P2  | 34  | M | Totally blind  | Bibibili, Youku                 | 1 year  | Daily |
| P3  | 26  | M | Legally blind  | Bibibili, Douyin                | 3 years | 2-3 times per week |
| P4  | 27  | M | Totally blind  | Douyin, Youku                   | 3 years | 4-5 times per week |
| P5  | 38  | F | Totally blind  | Bibibili, Douyin                | <1 year | Daily |
| P6  | 35  | F | Legally blind  | Bibibili                        | 4 years | 1 time per week |
| P7  | 30  | F | Legally blind  | Bibibili, Tencent               | 1 year  | 2-3 times per week |
| P8  | 33  | M | Totally blind  | Douyin                          | 3 years | Daily |
| P9  | 22  | M | Totally blind  | Douyin, Youku                   | 2 years | Daily |
| P10 | 30  | F | Totally blind  | Bibibili, Douyin                | 1 year  | 2-3 times per week |
| P11 | 32  | F | Totally blind  | Bibibili, Douyin, iQiyi, Tencent| 3 years | Daily |
| P12 | 41  | M | Legally blind  | Tencent, Douyin                 | 4 years | Daily |
| P13 | 32  | M | Totally blind  | Bibibili, Douyin, Youku         | 5 years | 4-5 times per week |
| P14 | 36  | M | Totally blind  | Bibibili, Youku                 | 2 years | 2-3 times per week |
| P15 | 28  | F | Legally blind  | Bibibili, Douyin                | 4 years | 2-3 times per week |
| P16 | 23  | M | Totally blind  | Bibibili, Douyin, iQiyi         | 1 year  | Daily |
| P17 | 36  | M | Totally blind  | Tencent, iQiyi                  | 4 years | Daily |
| P18 | 30  | F | Legally blind  | Bibibili, Douyin                | 2 years | 4-5 times per week |
| P19 | 31  | M | Totally blind  | Bibibili, Douyin, iQiyi         | 2 years | Daily |
| P20 | 29  | F | Legally blind  | Bibibili, Youku, iQiyi, Tencent | 3 years | Daily |

**Table 2: The six videos used in the formative study.**

| Type | Video Content | Length |
|------|---------------|--------|
| Educational | A person explains how fast humans can really run. | 01:04 |
| Comedic | A collection of funny bloopers. | 03:17 |
| Tutorial | A person shows how to prepare "egg dishes" for a meal. | 04:28 |
| News | A news clip of a three-year-old girl caught by a kite. | 00:33 |
| Music | A person sings the music "Little Love Song". | 04:33 |
| Film Clip | Dialogues from the dinner party scene in "Eat Drink Man Woman". | 06:56 |

**Table 3: The 18-video dataset used in the evaluation study. V3-V8 are used for the controlled comparison.**

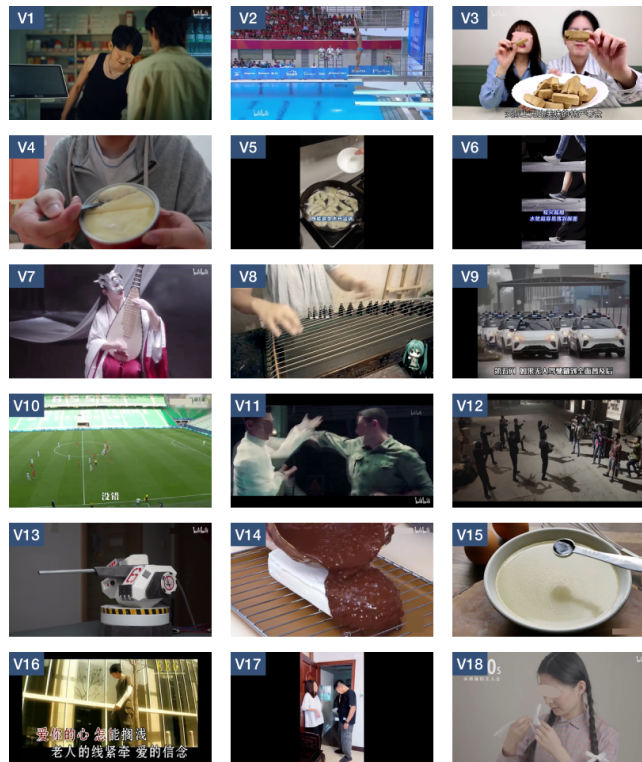| ID | Type | Video Content | Length | Speech Ratio |
|----|------|---------------|--------|--------------|
| V1 | Film Clip | Dialogues between two people in a grocery store. | 2:04 | 48% |
| V2 | News | A sports commentary on a diving event. | 1:22 | 80% |
| **V3** | **Comedic** | **Two people humorously try trending internet snacks.** | **5:08** | **83%** |
| **V4** | **Comedic** | **One person humorously tries celebrity-recommended food.** | **4:47** | **83%** |
| **V5** | **Educational** | **A person explains the proper way to cook dumplings.** | **1:45** | **97%** |
| **V6** | **Educational** | **A person explains why shoes get wet on rainy days.** | **1:48** | **100%** |
| **V7** | **Music** | **Instrumental music "Battle for Prince of Lan Ling".** | **4:59** | **0%** |
| **V8** | **Music** | **Instrumental music "Senbonzakura".** | **4:56** | **0%** |
| V9 | News | A report on the controversy surrounding self-driving cars. | 2:48 | 100% |
| V10 | News | A report on an event during the Paris Olympics. | 3:43 | 100% |
| V11 | Film Clip | A group of people fight using traditional Chinese kung fu. | 4:19 | 17% |
| V12 | Film Clip | A group of people fight with each other using cameras. | 5:33 | 9% |
| V13 | Tutorial | A person demonstrates how to make a mosquito killer device. | 5:58 | 66% |
| V14 | Tutorial | A person shows how to make a large ice-cream bar. | 6:49 | 86% |
| V15 | Tutorial | A person demonstrates how to make pudding. | 2:47 | 96% |
| V16 | Music | Two people sing the music "The Old Man and the Sea". | 4:23 | 64% |
| V17 | Comedic | A person tries to escape home with a funny excuse. | 0:24 | 12% |
| V18 | Educational | The evolution of Chinese women's hairstyles over the last century. | 3:55 | 12% |



**Figure 12: Screenshots from the 18-video dataset. Video creators' IDs and human faces are obscured for privacy.**

**Table 4: Detailed statistics of subjective ratings from the evaluation study (1 = strongly negative, 7 = strongly positive). SD denotes the standard deviation. Significance was analyzed using the Wilcoxon signed-rank test.**

| Aspects | Questions | DanmuA11y | | Baseline | | Significance |
|---|---|---|---|---|---|---|
| | | Mean | SD | Mean | SD | |
| System Usability | Q1: I am **willing to use** the system in the future. | 6.92 | 0.29 | 4.42 | 2.07 | $Z = -2.82,\ p < .01$ |
| | Q2: The system is **easy to use**. | 6.67 | 0.49 | 4.75 | 2.01 | $Z = -2.70,\ p < .01$ |
| | Q3: The system is **easy to learn**. | 7.00 | 0.00 | 6.42 | 1.08 | $Z = -1.63,\ p = 0.1$ |
| | Q4: The system **well integrates** various functions. | 6.75 | 0.62 | 4.50 | 1.78 | $Z = -2.82,\ p < .01$ |
| | Q5: The system is **not mentally demanding** to use. | 6.25 | 0.97 | 4.58 | 1.83 | $Z = -2.54,\ p < .05$ |
| | Q6: The system is **not physically demanding** to use. | 6.75 | 0.62 | 4.25 | 1.82 | $Z = -2.83,\ p < .01$ |
| Danmu Comprehension | Q7: It is easy to grasp the **discussion topics**. | 6.50 | 0.80 | 5.00 | 1.81 | $Z = -2.69,\ p < .01$ |
| | Q8: It is easy to grasp the **topics' visual context**. | 6.75 | 0.45 | 4.75 | 1.71 | $Z = -2.83,\ p < .01$ |
| | Q9: It is easy to grasp the **interactions among viewers**. | 6.75 | 0.45 | 4.67 | 1.83 | $Z = -2.70,\ p < .01$ |
| Curation Quality | Q10: The Danmu comments are **informative**. | 6.25 | 1.48 | 4.83 | 2.21 | $Z = -2.21,\ p < .05$ |
| | Q11: The Danmu comments are **creative**. | 6.25 | 0.97 | 4.75 | 1.91 | $Z = -2.56,\ p < .05$ |
| | Q12: The Danmu comments reflect **diverse opinions**. | 6.17 | 1.03 | 4.42 | 1.93 | $Z = -2.51,\ p < .05$ |
| Video Viewing Experience | Q13: The Danmu comments are **coherent with the video**. | 6.67 | 0.65 | 3.92 | 1.44 | $Z = -3.08,\ p < .01$ |
| | Q14: The Danmu integration is **unobtrusive to the video**. | 6.33 | 0.89 | 3.92 | 1.62 | $Z = -2.82,\ p < .01$ |
| | Q15: I can focus on the video and **not lose track of it**. | 6.17 | 0.72 | 4.50 | 1.57 | $Z = -2.75,\ p < .01$ |
| Social Connection | Q16: I feel like **watching together** with other viewers. | 6.58 | 0.90 | 4.50 | 1.88 | $Z = -2.83,\ p < .01$ |
| | Q17: **Closeness with other viewers** (using IOS Scale [3]). | 5.50 | 1.17 | 3.83 | 1.53 | $Z = -2.55,\ p < .05$ |
| | Q18: I feel **joyful** when viewing videos using the system. | 6.83 | 0.39 | 5.17 | 2.21 | $Z = -2.23,\ p < .05$ |

**Table 5: Details of the nine uncompleted trials with the baseline. (All 36 trials were completed with DanmuA11y.)**

| Uncompleted Videos | V3 | V4 | V5 | V6 | V7 | V8 |
|---|---|---|---|---|---|---|
| Participants | P15 | P9, P17 | P10, P15, P18 | P12, P17 | - | P20 |