

1 **Enhancing Art Perception for Individuals with Visual Impairments Using**
2 **Computer Vision, Depth Perception and Multimodal Feedback**
3

4 ADYA AGGARWAL, IIIT-Delhi
5

6 PANKHURI SINGH, IIIT-Delhi
7

8 PRANAV JAIN, IIIT-Delhi
9

10 Art remains largely inaccessible to individuals with visual impairments, limiting their ability to experience paintings. This paper
11 presents a multimodal haptic-audio system that enables users to perceive artwork through touch and sound. A camera tracks a user's
12 finger, marked with a distinct color, as it moves across a canvas, mapping its position to the artwork. A depth estimation algorithm
13 translates image depth into haptic feedback using ERM vibration motors or piezoelectric actuators on the fingertips. Simultaneously,
14 color information is converted into auditory cues via headphones. This dual-feedback mechanism enables vision-impaired users to
15 experience spatial textures and color representations, allowing for a richer and more intuitive engagement with visual art. Additionally,
16 sighted users gain a tactile appreciation of paintings that are typically untouchable, enriching their interaction with art.
17

18 Additional Key Words and Phrases: Multisensory Art Accessibility, Haptic Feedback for Art, Computer Vision, Depth Mapping for
19 Tactile Art, Color Sonification
20

21 **ACM Reference Format:**

22 Adya Aggarwal, Pankhuri Singh, and Pranav Jain. 2018. Enhancing Art Perception for Individuals with Visual Impairments Using
23 Computer Vision, Depth Perception and Multimodal Feedback. In *Proceedings of Make sure to enter the correct conference title from your*
24 *rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 20 pages. <https://doi.org/XXXXXX.XXXXXXXX>
25

26 The complete source code is available on GitHub.
27

28 **1 Introduction**
29

30 Traditional two-dimensional (2D) paintings and artworks rely almost entirely on visual perception, making them
31 inaccessible to individuals with visual impairments and limiting tactile engagement for sighted users. While some
32 accessibility solutions exist, they often fail to provide a comprehensive sensory experience. Bas-relief techniques
33 and tactile graphics attempt to translate paintings into touchable forms, but they frequently lose fine details and
34 spatial accuracy. Audio descriptions, though helpful, often provide subjective interpretations rather than a structured
35 understanding of an artwork's composition.
36

37 A major challenge in making visual art more accessible is the inability to convey depth, texture, and color in a way
38 that is intuitive for non-visual users. Color, in particular, remains difficult to represent through tactile or auditory means,
39 leaving individuals with an incomplete understanding of an artwork's visual elements. Additionally, existing accessibility
40 tools lack interactivity, preventing users from actively exploring paintings at their own pace. These limitations create a
41 gap in how visually impaired individuals and even sighted users interact with and appreciate visual art, highlighting
42 the need for a more effective and inclusive solution.
43

44 To address these limitations, this research proposes an AI-driven, multimodal system that enhances art perception
45 through computer vision, haptic feedback, and auditory cues.
46

47 Authors' Contact Information: Adya Aggarwal, IIIT-Delhi, New Delhi; Pankhuri Singh, IIIT-Delhi, New Delhi; Pranav Jain, IIIT-Delhi, New Delhi.
48

49 2018. Manuscript submitted to ACM
50

51 Manuscript submitted to ACM
52

53 *Problem Statement.* Despite advancements in accessibility tools, there remains a lack of affordable, interactive, and
54 intuitive systems for experiencing visual art without sight. Current solutions are either expensive or cognitively
55 demanding, often sacrificing detail or interactivity. There is a need for a cost-efficient and engaging approach that can
56 simulate depth and color perception while minimizing cognitive overload for the user.
57

58
59 *Research Questions.* To tackle the above challenge, this work explores the following key research questions:
60

- 61** (1) **How can we design a cost-effective haptic feedback system using basic vibration motors to simulate**
62 **texture for visually impaired users?**
- 63** (2) **How can we reduce cognitive overload during art exploration to enhance the learning curve and**
64 **ensure sustained user engagement?**
- 65** (3) **Can Large Language Models (LLMs) be integrated to provide contextual and adaptive learning experi-**
66 **ences for users interacting with artworks?**

70 2 Literature Review

72 Accessibility in visual perception extends beyond the realm of art, influencing various fields such as sports, performance
73 arts, music, and education. Researchers have explored different technologies and methodologies to enhance sensory
74 experiences for individuals with visual impairments or other sensory limitations. Studies in haptic feedback, auditory
75 substitution, computer vision, and multisensory integration have contributed to making information more accessible
76 across diverse domains. The following sections review literature from multiple disciplines, highlighting advancements,
77 limitations, and potential applications relevant to creating a more inclusive and immersive experience.
78

81 2.1 Haptic Assistance in Sports

83 A drone-guided haptic system helps visually impaired runners by sending directional cues via a wearable vibrotactile
84 device. Continuous feedback proved most effective, keeping users on track 93% of the time while increasing their
85 confidence and safety. Future developments aim to support visually impaired athletes in team sports like blind soccer[1].

87 A skiing assistance system uses vibrating bracelets on the forearms to provide haptic guidance. Ski instructors
88 control these vibrations via augmented ski poles, allowing users to navigate without relying solely on auditory cues.
89 This system enhances safety and reduces cognitive load for visually impaired skiers[2].

91 2.2 Multisensory Communication in Performing Arts

93 A vibrating bracelet enhances communication for deaf performers by providing tactile feedback linked to emotional
94 states. Using heart rate monitoring, it helps actors understand and express emotions, improving audience engagement.
95 The study highlights the role of multisensory technology in making the arts more inclusive[3].

98 2.3 Assistive Technologies for Deafblind Individuals

100 Deafblind individuals lack both visual and auditory channels for communication, making interaction difficult. Assistive
101 devices use haptic stimuli to deliver textual information and translate speech and visual input into tactile feedback.
102 These systems expand access to digital platforms and support independent communication[4].

105 **2.4 Music Accessibility Through Multisensory Integration**

106
107 The Auris System enables music experiences for individuals with hearing impairments by converting sound into tactile
108 and visual stimuli. Devices like the Auris Chair and Bracelet translate musical elements into vibrations, activating
109 the auditory cortex through multisensory integration. EEG analysis helps refine these representations, improving
110 accessibility for deaf users[5].
111

112 **2.5 Wearable Safety and Alert Systems**

113
114 The Vibration Alert Bracelet provides emergency notifications for visually and hearing-impaired users through wireless
115 alerts and distress signals. It features a Wi-Fi-enabled bracelet and a mobile app, offering a lightweight and efficient
116 solution. Future upgrades may include BLE for power efficiency, GPS tracking, and enhanced security[6].
117

118 **2.6 Crossmodal Translation Between Sound and Vision**

119
120 Innovative approaches in sensory substitution explore the translation between auditory and visual modalities. One such
121 method involves converting sounds into visual patterns or vice versa to support individuals with sensory impairments.
122 Research highlights systems that use visual displays to represent auditory properties such as pitch, timbre, and rhythm,
123 enabling deaf users to perceive music visually. Conversely, images can be sonified to convey visual scenes to blind users
124 through sound patterns. These techniques leverage the brain's neuroplasticity and underscore its ability to interpret
125 sensory data across modalities. The development of such systems reveals potential for inclusive media, education, and
126 multisensory artistic performances[20]
127

128 **2.7 Color-Based Quantification and Multisensory Translation**

129
130 The integration of color-based representation into analytical systems extends beyond scientific measurement and into
131 broader accessibility. Cantrell et al. introduced the use of the hue (H) parameter from the HSV color space as a stable
132 and precise metric for bitonal optical sensors. Unlike RGB intensity-based approaches, the hue value remains consistent
133 across variations in membrane thickness, indicator concentration, and imaging conditions. This makes it a reliable
134 parameter for portable or disposable sensors used in mobile diagnostics and quality control. The study demonstrates
135 how a traditionally qualitative visual attribute—color—can be quantified and used analytically, bridging the gap between
136 visual perception and data-driven analysis[21]
137

138 **2.8 Low-Cost Haptic Glove Designs**

139
140 Recent advancements have led to the development of affordable haptic gloves aimed at enhancing accessibility and
141 immersion in various applications.
142

143 DOGlove is an open-source glove offering 21 degrees of freedom (DoF) motion capture and 5-DoF force feedback,
144 assembled under \$600, designed for precise teleoperation and dexterous manipulation.
145

146 Pulse Haptic Glove, designed for enthusiasts, uses knuckle-mounted motors to apply pressure via fingertip caps,
147 providing realistic haptic feedback at approximately \$300.
148

149 TactGlove features 10 individually controllable linear resonant actuators (LRAs) on the fingertips, offering enhanced
150 haptic feedback compatible with camera-based hand tracking systems, priced around \$299.
151

152 Affordable haptic gloves beyond the fingertips propose a novel design extending haptic feedback beyond fingertips
153 to include intermediate and proximal phalanges, utilizing a ratchet and pawl mechanism for more realistic feedback.
154

157 **3 Methodology**

158 The methodology is structured into four key phases: data collection, system pipeline design, code development, and
 159 hardware integration.

161 **3.1 Data Collection**

162 Multiple data collection methods were employed to ensure a user-centered design and effective system development.
 163 These included interviews, surveys, and persona building, providing both qualitative and quantitative insights.

167 **3.2 System Pipeline and Workflow**

169 This section elaborates on the technical pipeline driving the system. It integrates computer vision, depth estimation,
 170 color classification, and multisensory feedback to create an immersive experience.

172 **3.2.1 Image and Surface Setup.** A printed artwork is placed on a flat surface. A camera mounted overhead captures
 173 the entire frame including the image and the user's hand. This is done by recognizing the image by a big green rectangle.

175 **3.2.2 Depth Map Generation.** The image is processed via the SculptOK depth estimation API, which returns a
 176 pixel-wise depth map using pre-trained models. This map simulates surface texture for tactile conversion.

178 **3.2.3 Real-Time Hand and Finger Tracking.** Using OpenCV and MediaPipe Hands, the system detects and continu-
 179 ously tracks the index fingertip (Landmark ID 8). The fingertip's (x, y) coordinates are mapped to the image frame and
 180 corresponding depth data.

183 **3.2.4 Vibration Feedback (Depth-to-Haptics).**

$$185 \quad I(d) = \begin{cases} 0 & \text{if } d \leq 0.01 \\ 186 \quad d^{1.5} \cdot (1 - (1 - d)^{1.5}) & \text{if } d > 0.01 \end{cases}$$

188 **Where:**

- 189
 - 190 • $d \in [0, 1]$ is the *normalized depth* value (i.e., `depth_map[y, x] / 255.0`),
 - 191 • $I(d)$ is the *vibration intensity*.

193 This function produces a *bell-shaped curve* peaking at intermediate depth values.

194 **3.2.5 Color Detection and Audio Feedback.** The RGB color at the fingertip location is classified using a trained
 195 K-Nearest Neighbors (KNN) classifier. The following mappings are used:

- 197
 - 198 • **Red:** Guitar
 - 199 • **Black:** Saxophone
 - 200 • **Blue:** Flute
 - 201 • **Green:** Congo
 - 202 • **Yellow:** Piano

204 These are sonified and played using the Python `sounddevice` library. All sounds are defined by 3 important metrics -
 205 Frequency(Hz), Amplitude and Waveform. The waveforms have been modified to simulate different instruments, with a
 206 frequency of 440Hz and equal amplitu.

209 3.2.6 **Integrated Feedback Experience.** As the user moves their finger across the artwork:

- 210
- 211 • Texture (depth) is conveyed via vibration intensity.
 - 212 • Color is conveyed through musical audio cues.

213

214 This creates a synchronized haptic-audio experience for visually impaired users to explore visual art through touch and

215 sound.

216 217 3.3 Hardware Integration

218

219 The system relies on a combination of software and hardware to deliver synchronized feedback. Below is an overview

220 of the hardware connections and purpose:

221 *1. Arduino UNO Acts as the interface between depth-based intensity and the ERM motor. Receives vibration levels

222 via serial input from the main system.

223

224 *2. ERM Vibration Motor Connected to a PWM-enabled digital pin (e.g., D9) on the Arduino. Provides real-time

225 feedback at varying intensities.

226

227 *3. Camera Mounted overhead, continuously streams frames to the Python script for fingertip detection and position

228 mapping.

229

230 *4. Headphones Connected to the system output. Sound cues are generated on the host machine and played back

231 instantly to indicate color recognition.

232

233 *5. Power Supply The Arduino is powered via USB from the main system, and external power may be used for

234 consistent motor strength if required.

235

236 This tight integration allows the user to interact naturally with an image and receive immediate sensory feedback,

237 enabling an intuitive and accessible form of art exploration.

238 3.4 Chatbot Integration

239

240 In the language model pipeline, the process begins with an image captured during the initial phase, accompanied by a

241 natural language question provided by the user. This image-question pair is first processed by BLIP-2, a state-of-the-art

242 Vision-Language model designed for Visual Question Answering (VQA). BLIP-2 interprets the visual content of the

243 image in conjunction with the user's query and generates a relevant textual response based on its understanding of

244 both modalities.

245

246 The output from BLIP-2 is then passed to Mistral-7B, a powerful open-weight large language model known for its

247 instruction-following and text-generation capabilities. Mistral-7B takes the VQA response and contextualizes it further,

248 refining the output into a coherent, informative, and well-articulated sentence. This final response combines the visual

249 insight from BLIP-2 with Mistral-7B's linguistic fluency, enabling a more conversational and context-aware interaction.

250 251 3.5 Modalities

252

253 The system employs a multimodal interaction framework designed to reduce cognitive load and enhance efficiency

254 by distributing information across various sensory channels. This approach aligns with principles of multisensory

255 integration, where combining inputs from different modalities can lead to improved perception and performance.

256

257 In the primary system, users explore digital artworks through tactile and auditory feedback. A camera tracks the

258 user's finger movements across a canvas, mapping these to corresponding elements in the artwork. Depth information

259 is conveyed through haptic feedback using Eccentric Rotating Mass (ERM) vibration motors on the fingertips, allowing

261 users to perceive spatial structures. Simultaneously, color information is translated into distinct auditory cues delivered
 262 via headphones, enabling users to identify colors through sound.
 263

264 Complementing this, a conversational chatbot interface supports both text and speech interactions. Users can input
 265 queries either by typing or speaking, and the chatbot responds with both textual and spoken answers. This dual-mode
 266 communication caters to different user preferences and situational needs, facilitating seamless interaction.
 267

268 By integrating tactile, auditory, and linguistic modalities, the system distributes cognitive demands, allowing users
 269 to process information more naturally and efficiently. This multimodal design not only enhances user engagement but
 270 also supports accessibility, particularly for individuals with visual impairments.
 271

272 4 Evaluation and Explaination

273 The project code is modularized into the following key components:
 274

275 4.1 Explainability

276 To enhance the transparency and accountability of our model's predictions, we applied **SHAP (Shapley Additive**
 277 **exPlanations)**, a leading explainability framework grounded in cooperative game theory. SHAP provides instance-level
 278 attributions that clearly explain how much each input feature (Red, Green, and Blue channel intensities) contributed to
 279 the model's final prediction.
 280

281 The following table displays SHAP values along with the raw RGB inputs and the predicted color label for some
 282 samples:
 283

286 Predicted_Label	287 Index	288 R	289 G	290 B	291 SHAP_R	292 SHAP_G	293 SHAP_B
288	289 0	290 4	291 222	292 10	293 0.077	294 -0.293	295 0.117
289	290 Green						
290	291 1	292 19	293 17	294 234	295 0.111	296 0.083	297 -0.294
291	292 Blue						
292	293 2	294 223	295 35	296 1	297 -0.270	298 0.050	299 0.120
293	294 Red						
294	295 3	296 48	297 243	298 32	299 0.071	300 -0.293	301 0.122
295	296 Green						
296	297 4	298 227	299 34	300 39	301 -0.270	302 0.050	303 0.120
297	298 Red						

298 Table 1. SHAP Empowers Confident Interpretation
 299

300 Key Observations

- 301
- 302 (1) The model consistently predicts correct classes despite complex input combinations. For instance, in Sample
 303 0, even though the Green (G) channel shows a strongly negative SHAP value, the model balances it with
 304 positive contributions from the Red and Blue channels to confidently classify it as "Green". This shows a mature
 305 decision boundary, where the model does not rely on a single dominant feature but understands color nuances
 306 holistically.
 307
 - 308 (2) In Sample 1, the model predicts "Blue" despite the Blue (B) channel having a negative SHAP value (-0.294). This
 309 again proves that the model has learned fine-grained relationships, where high blue intensity alone does not
 310

- 313 blindly drive the decision – instead, it factors in the subtle push from the Red and Green channels for a more
314 context-aware classification.
315
316 (3) The Red class predictions in Samples 2 and 4 reflect that high Red intensity does not necessarily mean a strong
317 SHAP value. Interestingly, both have slightly negative SHAP contributions, meaning the model is learning to differentiate reds from redsw
318 dimensional ways.
319

320 4.2 Model Evaluation: Hand Recognition System Using MediaPipe

321 To evaluate our hand recognition pipeline, we tested it on the HaGRID 30K sample dataset (384p resolution). The model
322 leverages MediaPipe Hands, a lightweight and highly optimized framework by Google that performs real-time hand
323 tracking using a machine learning model trained to detect up to 21 3D landmarks per hand. The detection performance
324 is summarized below:

- 325 (1) **Precision (1.000)** : This indicates that every detected hand was a true positive—there were no false positives.
326 MediaPipe’s confidence-based detection and landmark refinement pipeline ensures that once a hand is identified,
327 it is highly reliable and accurate.
328
329 (2) **Recall (0.674)**: While MediaPipe is precise, it misses about 32.6% of the hands in the dataset. This could be
330 due to: occlusions, multiple hands in frame, low lighting or resolution issues, MediaPipe’s internal thresholding
331 (designed for speed over exhaustive detection).
332
333 (3) In a high-variance dataset like HaGRID, which includes diverse poses and gestures, this trade-off is expected.
334
335 (4) **F1 Score (0.805)**: A harmonic mean of precision and recall, the F1 Score reflects an overall balanced performance.
336 While recall could be improved with more aggressive detection tuning or ensemble methods, the score indicates
337 a strong and stable performance.
338
339 (5) **Accuracy (0.674)**: This metric shows the overall percentage of correct detections (true positives + true negatives).
340 Since this is a detection task, accuracy here is consistent with recall, reaffirming that most misses are due to
341 undetected hands, not incorrect ones.
342
343

344 Why MediaPipe Works Well :

- 345 (1) **Real-time optimized**: Extremely efficient on CPU and mobile.
346
347 (2) **Precision-first**: Prioritizes correctness over exhaustive detection, as reflected in the high precision.
348
349 (3) **Modular**: Can be extended with custom classifiers (e.g., gesture recognition, hand sign classification).

350 Our MediaPipe-based hand recognition system demonstrates excellent precision and reliable performance, particularly
351 suitable for real-time applications where false positives must be minimized. With additional tuning or post-processing,
352 recall can be improved to further enhance detection coverage without sacrificing accuracy.

353 Performance of functions used in contour detection

- 354 (1) **cv2.cvtColor(img, cv2.COLOR_BGR2GRAY)**: Converts a BGR image to grayscale. On CUDA-enabled GPUs,
355 the device execution time is approximately 0.03 ms, with host times ranging from 0.087 ms to 0.28 ms, depending
356 on system load and initialization overhead.
357
358 (2) **cv2.GaussianBlur(gray, (5, 5), 0)**: Applies Gaussian blur to reduce image noise and detail. Ikomia +15. Virtual
359 Edge Solution | IO River +15. BoofCV +15. On a GeForce GTX 750 GPU.
360
361 (3) **cv2.Canny(blur, 50, 150)**: Performs Canny edge detection. The execution time varies with image complexity.
362 For instance, on an 8-threaded system, optimization reduced the Canny algorithm’s runtime from 1.7 seconds
363 to 1.45 seconds.
364

365 4.3 Model Evaluation: Color Identification Using KNN

366 The K-Nearest Neighbors (KNN) classifier used for color identification demonstrates strong predictive performance, as
 367 reflected in the following metrics:

- 368
- 369 (1) **F1 Score (0.9715):** This high F1 score highlights that the model maintains an excellent balance between precision
 370 and recall, meaning it is effective at correctly identifying colors while minimizing both false positives and false
 371 negatives. This is particularly important in color classification, where overlaps between visually similar colors
 372 can occur.
- 373 (2) **Accuracy (0.9720):** An accuracy of 97.2% confirms that the model classifies the vast majority of test samples
 374 correctly. It suggests that the current feature set (RGB values) and KNN's distance-based classification are
 375 highly effective for this task.

376 With both high accuracy and a strong F1 score, the KNN model proves to be a robust and reliable choice for color
 377 identification. Its performance can serve as a solid baseline, with opportunities to further optimize through parameter
 378 adjustments.

384 4.4 Chatbot Evaluation

385 :

387 1. Speech-to-Text (STT) – Google Cloud API

- 388
- 389 (1) **Word Error Rate (WER):** Approximately 8% under favorable conditions.
- 390 (2) **Latency:** Typically 1–2 seconds for short audio inputs.
- 391 (3) **Multilingual Support:** Supports over 125 languages, including English and Hindi.

393 2. Text-to-Speech (TTS) – ElevenLabs

- 395
- 396 (1) **Mean Opinion Score (MOS):** Ranging from 4.2 to 4.5 on a 5-point scale, indicating high naturalness.
- 397 (2) **Latency:** Approximately 5–6 seconds for generating short responses.
- 398 (3) **Features:** Offers emotional tone control, multilingual support, and real-time streaming capabilities.

400 3. Chatbot Evaluation (Visual Question Answering)

- 402 (1) **BLEU Score:** Typically ranging from 0.25 to 0.35, reflecting moderate alignment with reference answers.

404 4.5 Interpretability

406 The Odin Vision system is explainable and interpretable because it uses transparent, human-understandable logic to
 407 convert visual elements like color and depth into sound and vibration, allowing visually impaired users to meaningfully
 408 interact with paintings. Its modular design, reliance on simple models like KNN for color classification, and real-time
 409 visual overlays ensure the decision-making process is easy to understand, debug, and explain—aligning with the
 410 course definitions of explainability and interpretability. By mapping distinct colors to culturally intuitive sounds and
 411 varying vibration based on proximity, the system offers consistent, predictable feedback that fosters trust, usability, and
 412 inclusivity. Additionally, features like voice-based interaction and multilingual support make it accessible to diverse
 413 users, supporting the principles of human-centered AI, including transparency, user control, and cultural sensitivity.

417 **5 User Study and Evaluation**
418
419
420435 Fig. 1. Participant 1
436435 Fig. 2. Participant 2
436439 **5.1 Pre-Study Survey**
440

441 To evaluate the effectiveness of our system designed for visually impaired individuals, we conducted a study involving
442 10 participants who were blindfolded for the duration of the experiment. Before interacting with the system, each
443 participant was asked to fill out a preliminary survey. This survey included questions about any pre-existing visual
444 impairments. The majority of participants reported no such impairments. When asked about their engagement with
445 visual arts, 27 participants stated they engage frequently, 36 occasionally, 27 rarely, and 9 never. On being asked whether
446 they had ever experienced art through non-visual means, 63% responded negatively. Most participants rated their sense
447 of touch for object recognition as a 3 on a scale of 1 to 5.
448

450 **5.2 Study**
451

452 The core system used in this study tracks a user's finger via a camera, mapping its position across a canvas in
453 correspondence with the artwork. A depth estimation algorithm translates image depth into haptic feedback using
454 ERM vibration motors attached to the fingertips. Simultaneously, color information from the artwork is converted into
455 auditory cues delivered through headphones. The task for participants was to experience the painting through this
456 system and then attempt to recreate it by drawing.
457

460 **5.3 Compositional Variation**
461

462 The images that participants drew displayed consistent variation, yet key elements from the original artwork were
463 present. These included a large red circle against a blue background, boxes of various colors in the foreground, and a
464 green triangle set against a yellow background. In terms of shape recognition, most participants were able to identify
465 the basic geometric forms; however, they faced challenges with accurate composition, particularly when visual elements
466 were overlapping.
467

469
470
471
472
473
474
475
476
477

5.4 Colour Variation

478 Color identification showed promising results. Most participants correctly identified the colors used in the artwork.
479 Since each color was paired with a distinct auditory cue, recognition was relatively successful. However, the outline
480 of the shapes remained difficult for participants to reproduce accurately. While the locations where the colors were
481 applied were mostly correct, the shapes themselves were often distorted or imprecise.
482

483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520

5.5 Post-Study Reflection

519 In the post-study reflection, 45% of the participants rated the difficulty of understanding the artwork through touch
520 and sound as 3 out of 5, and another 45% rated it as 4. A majority of participants (63%) found the vibration feedback to
521 be somewhat clear, and 45% reported that the auditory cues were somewhat helpful in identifying colors. While the
522 system showed potential, participants noted areas for improvement. Suggestions included enhancing the language
523 model's responses to give more informative guidance, and improving the comfort of the haptic glove. Notably, 72% of
524 the participants expressed interest in using such a system in real-world settings like museums and galleries.
525

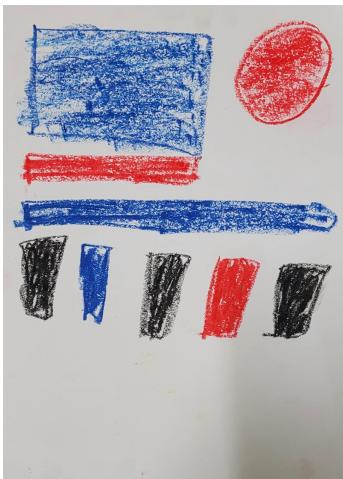


Fig. 3. Participant recreations of the artwork based on multisensory interaction

5.6 Conclusion

519 Overall, the study highlighted both the promise and the limitations of the system. While participants could recognize
520 and reproduce fundamental aspects of the artwork through multisensory cues, challenges remain in achieving precision
521 in shape and composition. Nonetheless, the positive reception suggests strong potential for further development and
522 application of such assistive technologies in inclusive art experiences.
523

6 User Evaluation HCI Principles

Tactile Feedback and Affordance

518 The glove integrates a rubber-dot textured surface to enhance tactile interaction, supporting the principle of *affordance*.
519 The textured dots offer clear interaction cues, enabling users to understand how to interact with the system based
520 Manuscript submitted to ACM

521 on surface feedback. This tactile design aligns with HCI's goal of reducing cognitive load and ensuring intuitive user
522 interaction by providing immediate, localized haptic feedback.
523

524 **Cognitive Load Optimization**

525 Initially, the system's mapping of *Hue*, *Saturation*, and *Value* to *Instrument*, *Frequency*, and *Amplitude*, respectively, was
526 found to be cognitively demanding. This violated the HCI principles of simplicity and error prevention. In response, the
527 system was redesigned to map only **Hue to Instrument**, reducing the cognitive load and supporting easier recognition
528 over recall. This change also adheres to the principle of *minimizing cognitive load*, making the interaction process more
529 intuitive and user-friendly.
530

531 **Depth-Based Object Perception**

532 To enhance the user experience, the system utilizes depth mapping to aid object recognition. By providing haptic
533 feedback based on the depth of objects within a scene, users can perceive the spatial arrangement of these objects. This
534 design follows the HCI principles of *spatial awareness* and *natural interaction*, ensuring users can intuitively navigate
535 and interact with the environment in a way that is both rewarding and engaging.
536

537 **Inclusive Dataset Utilization**

538 The system's use of the HaGRID dataset, which includes diverse skin tones and hand types, demonstrates a commitment
539 to *inclusive design*. By ensuring the model's robustness across varied demographics, the system minimizes bias and
540 provides consistent interaction experiences for all users. This supports the HCI principle of *designing for all*, ensuring
541 the technology is accessible and effective for a wide range of users.
542

543 **Multimodal Interaction via Voice Interfaces**

544 The system employs voice interfaces, utilizing **Google Speech-to-Text** optimized for Indian accents and **ElevenLabs**
545 **Text-to-Speech** using Indian voice models. This integration enhances the system's *accessibility* and *engagement* by
546 providing a familiar, culturally inclusive interface. By facilitating natural voice-based interaction, the system fosters
547 a more sociable and intuitive user experience, aligning with HCI goals of reducing barriers and promoting seamless
548 communication.
549

550 **User Feedback Loop**

551 At the end of each interaction, the system actively gathers user feedback. This feedback loop embodies *user-centered*
552 *design*, where user insights directly inform future design iterations. Encouraging ongoing user involvement aligns with
553 HCI's emphasis on *continuous improvement*, building trust and reinforcing the value of the user's perspective in shaping
554 the system.
555

556 *Error Handling and System Status Visibility.* The system includes real-time error notifications and progress indicators,
557 offering users clear visibility into the system's status. This transparency aligns with the HCI principles of *visibility of*
558 *system status and feedback*, ensuring users are always informed and never left uncertain about the system's operations.
559 Clear error handling and status visibility minimize frustration, fostering a positive user experience and increasing trust
560 in the system.
561

573 7 Discussion

574 7.1 Low-Resource Haptic Glove

575 Our system utilizes a cost-effective haptic glove built with ERM vibration motors controlled via an Arduino UNO.
 576 This glove offers five discrete levels of vibration based on image depth, allowing users to perceive texture variations.
 577 Although more advanced gloves such as DOGlove or TactGlove offer greater precision and multimodal feedback, our
 578 implementation prioritizes accessibility and affordability. It can be easily deployed in resource-constrained environments
 579 like schools or community spaces, enabling a broader impact.
 580

581 7.2 LLM as a Guide: Chatbot for Knowledge Validation

582 We propose the integration of a Large Language Model (LLM) to act as an intelligent guide or chatbot within the
 583 system. Beyond tactile and auditory feedback, the LLM would provide real-time explanations, historical context,
 584 and conversational validation for user understanding. For example, after identifying a region of the artwork, users
 585 could inquire about its symbolism, artist, or era. The LLM could also quiz the user to validate their comprehension,
 586 facilitating an engaging and educational interaction—particularly valuable for visually impaired users who rely heavily
 587 on descriptive interpretation.
 588

589 8 Limitations

590 While the proposed system shows promise, several limitations need to be acknowledged:
 591

- 592 • **Basic Vibration Feedback:** The use of ERM motors limits the system's ability to convey subtle tactile differences,
 593 affecting the perception of intricate textures and gradients.
- 594 • **Resource Constraints:** The focus on affordability restricts the use of higher-end sensors, actuators, and
 595 computing resources that could enhance the fidelity of the experience.
- 596 • **Limited Participation from Visually Impaired Individuals:** Due to access constraints, most evaluations
 597 were conducted with sighted users simulating visual impairment. This may not reflect true usability for the
 598 intended demographic.
- 599 • **Color-to-Sound Mapping Simplification:** The mapping of six colors to six instruments simplifies interpreta-
 600 tion but fails to represent the full diversity of hues and tonal variations found in art.
- 601 • **Environmental Sensitivity:** The system's color recognition is sensitive to lighting changes, which can impact
 602 the accuracy of auditory feedback.

603 9 Future Work

604 Several extensions and improvements are envisioned for future iterations of the system:
 605

- 606 • **Higher-Quality Haptics:** Incorporation of Linear Resonant Actuators (LRAs) or piezoelectric feedback mech-
 607 anisms could offer more realistic and localized haptic responses.
- 608 • **Complex Image Interpretation:** Future versions will aim to handle high-resolution artworks, segment artistic
 609 elements, and identify stylistic features using advanced computer vision.
- 610 • **Color Gradation to Soundscape:** Instead of mapping discrete colors to instruments, continuous color gradients
 611 could be translated into ambient soundscapes using variations in pitch, tone, and reverb.
- 612 • **Emotion-Based Sonification:** The system could identify the emotional tone of artwork and generate corre-
 613 sponding music to enhance interpretive immersion.

- **Expanded Multimodal Interfaces:** Adding voice narrations, spatial audio, or even tactile displays could improve inclusivity and enrich the user experience.
- **Collaborative Exploration:** Introducing a multi-user mode could facilitate group learning and enhance social engagement with the artwork.

631 References

- [1] Boreas Technologies. *Guidelines of Haptic UX Design*. Retrieved from <https://pages.boreas.ca/blog/piezo-haptics/guidelines-of-haptic-ux-design>, accessed April 2025.
- [1] . Liao, J. Salazar, and Y. Hirata, "Robotic Guidance System for Visually Impaired Users Running Outdoors Using Haptic Feedback," 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 8325–8331, Sep. 2021, doi: <https://doi.org/10.1109/iros51168.2021.9636567>.
- [2] . Aggravi, G. Salvietti, and D. Prattichizzo, "(PDF) Haptic Assistive Bracelets for Blind Skier Guidance," [www.researchgate.net](https://www.researchgate.net/publication/292681848_Haptic_Assistive_Bracelets_for_Blind_Skier_Guidance), Jan. 2016. https://www.researchgate.net/publication/292681848_Haptic_Assistive_Bracelets_for_Blind_Skier_Guidance (accessed Mar. 22, 2022).
- [3] . H. Huang and Y. Li, "The Development and Application of the Wearable Device for the Deaf Performers," ISSN: 2189-101X – The Asian Conference on Education International Development 2023 Official Conference Proceedings, pp. 579–591, May 2023, Accessed: Feb. 20, 2025. [Online]. Available: <https://papers.iafor.org/submission68872/>
- [4] . Ranasinghe et al., "EnPower: Haptic Interfaces for Deafblind Individuals to Interact, Communicate, and Entertain," Advances in Intelligent Systems and Computing, pp. 740–756, Nov. 2020, doi: https://doi.org/10.1007/978-3-030-63089-8_49.
- [5] . Alves Araujo, F. Lima Brasil, A. Candido Lima Santos, L. de Sousa Batista Junior, S. Pereira Fonseca Dutra, and C. Eduardo Coelho Freire Batista, "Auris System: Providing Vibrotactile Feedback for Hearing Impaired Population," BioMed Research International, vol. 2017, pp. 1–9, 2017, doi: <https://doi.org/10.1155/2017/2181380>.
- [6] . Conley et al., "Vibration Alert Bracelet for Notification of the Visually and Hearing Impaired," Journal of Open Hardware, Oct. 2019, Available: https://epublications.marquette.edu/electric_jac/622/
- [8] Q. Zeng, R. R. Martin, L. Wang, J. A. Quinn, Y. Sun, and C. Tu, "Region-Based Bas-Relief Generation from a Single Image," *Graphical Models*, vol. 76, no. 3, pp. 140–151, 2014.
- [9] X. Sun, P. L. Rosin, R. R. Martin, and F. C. Langbein, "Bas-Relief Generation Using Adaptive Histogram Equalization," *IEEE Transactions on Visualization and Computer Graphics*, vol. 15, no. 4, pp. 642–653, 2009.
- [10] J. Kerber, A. Tevs, A. Belyaev, R. Zayer, and H.-P. Seidel, "Feature-Sensitive Bas-Relief Generation," *IEEE International Conference on Shape Modeling and Applications (SMI)*, 2009.
- [11] X. Wu, Y. Zhao, J. Luo, M. Zhang, and W. Yang, "Bas-Relief Modeling from RGB Monocular Images with Regional Division Characteristics," *Scientific Reports*, vol. 12, 2022.
- [12] N. Suciati, M. T. Baihaqi, H. Tjandrasa, A. Z. Arifin, et al., "Converting Image into Bas Reliefs Using Image Processing Techniques," *Journal of Physics: Conference Series*, vol. 1196, 2019.
- [13] Y. Fu, H. Yu, C.-K. Yeh, J. Zhang, and T.-Y. Lee, "High Relief from Brush Painting," *IEEE Transactions on Visualization and Computer Graphics*, 2018.
- [14] Z. Yang, B. Chen, Y. Zheng, X. Chen, and K. Zhou, "Human Bas-Relief Generation from a Single Photograph," *IEEE Transactions*, 2021.
- [15] P. Sobociński, M. Maik, and K. Walczak, "Multimodal Presentation of 3D Relief Sculptures in Virtual Reality," *International Conference on Cyberworlds (CW)*, 2022.
- [16] S. Kratz and T. Dunnigan, "ThermoTouch: Design of a High Dynamic Temperature Range Thermal Haptic Display," *CHI Extended Abstracts*, 2016.
- [17] Thulin, S., "ThermoTouch: Sound maps matter: expanding cartophony. *Social & Cultural Geography*, 19(2), 192–210. <https://doi.org/10.1080/14649365.2016.1266028>
- [18] H. Ö. Sertlek, H. Slabbeekorn, C. ten Cate, and M. A. Ainslie, "Source specific sound mapping: Spatial, temporal and spectral distribution of sound in the Dutch North Sea," *Environmental Pollution*.
- [19] W. Lin, "The hearing, the mapping, and the Web: Investigating emerging online sound mapping practices," *Landscape and Urban Planning*, vol. 142, 2015.
- [20] Cantrell, K., Erenas, M. M., de Orbe-Payá, I., & Capitán-Vallvey, L. F. (2010). Use of the hue parameter of the hue, saturation, value color space as a quantitative analytical parameter for bitonal optical sensors. *Analytical Chemistry*, 82(2), 531–542. <https://doi.org/10.1021/ac901753c>
- [21] Yakhontova, E. N. (2016). Translating sounds into visual images, and vice versa. *Psychology and Psychotechnics*, (6), 103–108. <https://cyberleninka.ru/article/n/translating-sounds-into-visual-images-and-vice-versa>

670 10 Appendix

671 10.1 Low-Fidelity and System Design

672 A low-fidelity prototype was developed using basic materials to simulate system flow. The user interface consists of:

- A static image placed on a table.

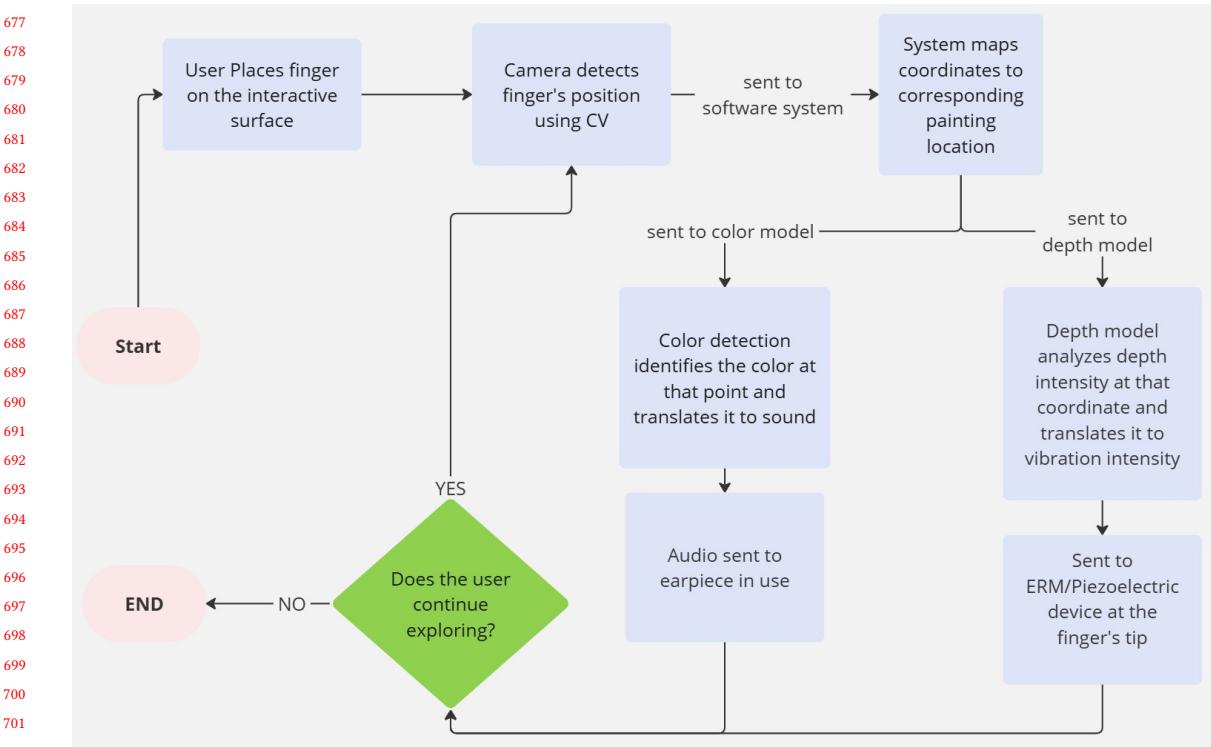


Fig. 4. User Taskflow

- A webcam overhead to track hand movement.
- Colorful fingertip markers to improve tracking reliability.
- Vibrating motors and a headphone setup to deliver multisensory feedback.

The system design follows a camera–processing–feedback loop, where all tracking and classification occur in real-time. This helps in validating the interaction flow and user response before progressing to final hardware implementation.

10.2 Low Fidelity Prototype

In Figure 1, the cursor, controlled by the fingertip, provides both haptic and auditory feedback based on its position within the visual field. When the cursor moves over darker regions, the vibration intensity remains low, and the corresponding sound frequency in the earphones is also lower. As the cursor moves into brighter, white regions, the vibration intensity increases, and the sound frequency rises accordingly. This combination of tactile and auditory feedback allows for an enhanced perception of different colors in the virtual environment

10.2.1 User Taskflow. See Figure 2.

10.2.2 Overview. Figure 3. illustrates an interactive system where a camera captures the user's hand movements, allowing them to explore a virtual interface using their fingertip as a cursor. As the user moves their finger, they

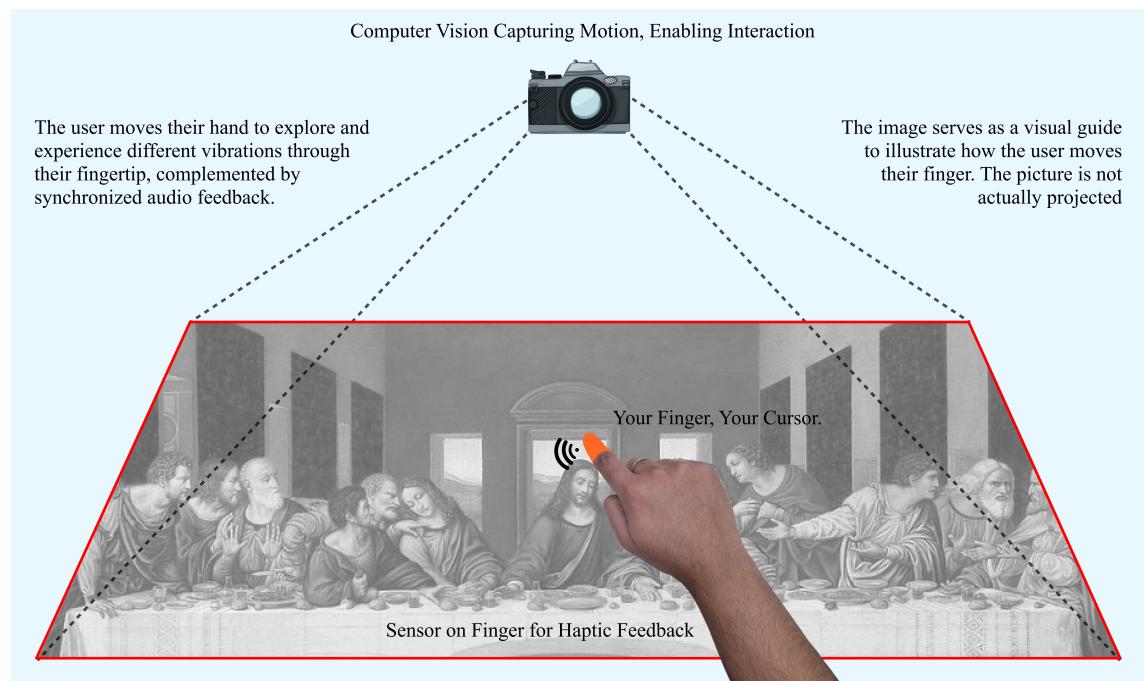


Fig. 5. Low Fidelity

experience different vibrations through a sensor attached to their fingertip, providing haptic feedback synchronized with audio cues. The displayed image serves only as a visual guide and is not actually projected.

10.3 Problems Solved

- **Limited Accessibility to Art for the Visually Impaired**: Traditional paintings are entirely visual, making them inaccessible to blind and visually impaired individuals. Existing solutions, like audio descriptions, fail to provide a fully immersive and interactive experience.
- **Lack of Tactile Engagement for Sighted Users**: Museums and galleries prohibit touching artwork, preventing visitors from experiencing the texture and fine details of paintings, limiting engagement and appreciation.
- **Insufficient Detail in Bas-Relief Generation**: Current 2D-to-bas-relief conversion methods often lack precision in preserving fine textures, focusing only on general depth rather than intricate details.
- **Absence of a Unified Multisensory Approach**: Previous research has addressed depth mapping, color sonification, and haptic feedback separately, but no comprehensive system integrates all three for a richer, more immersive artistic experience.

10.4 Mid Fi Design

This section presents the mid-fidelity prototype of the system, illustrating the functional flow and interactions between different components.

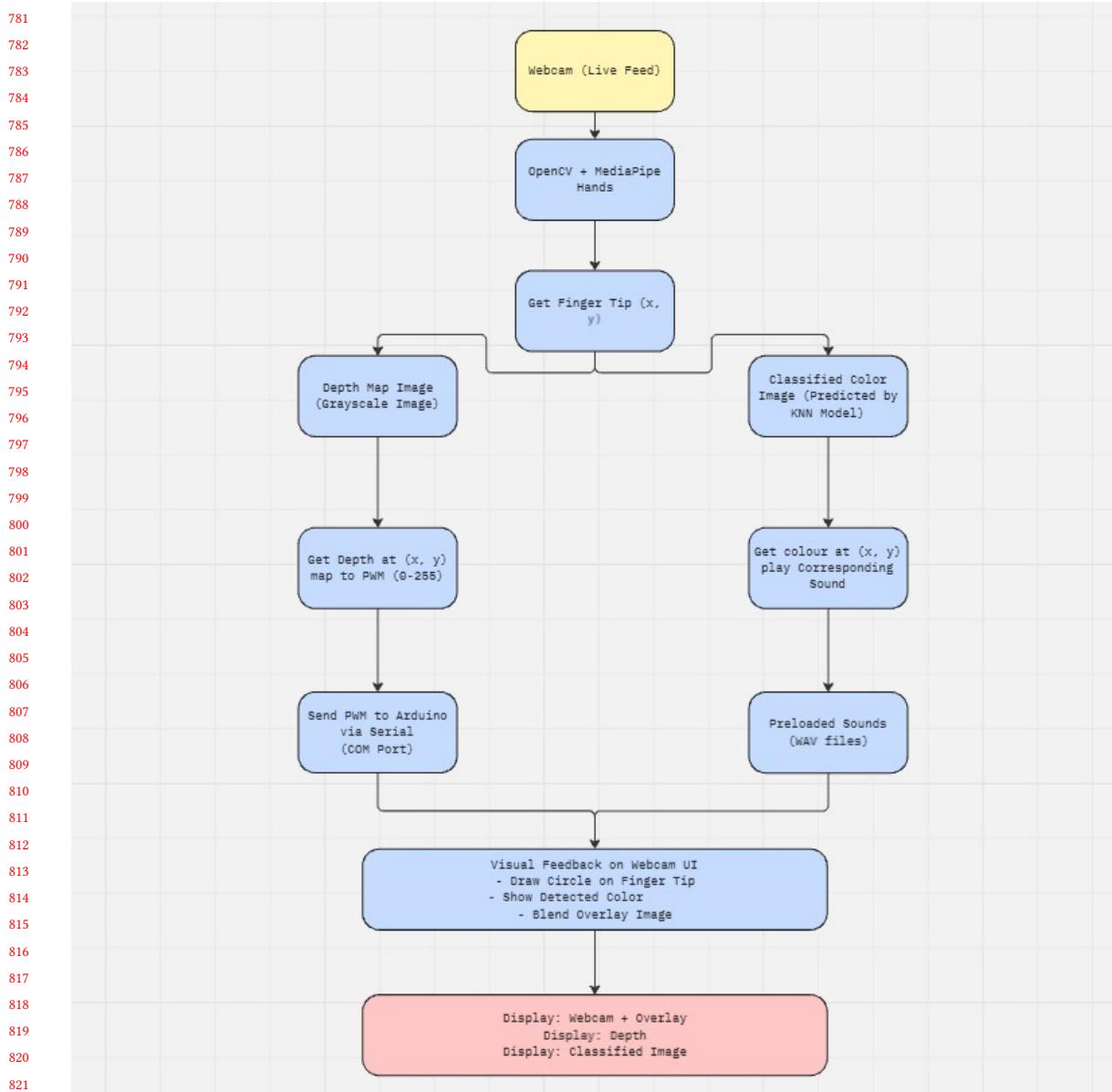


Fig. 6. TaskFlow

10.4.1 *Flow Diagram.* The flow diagram (see Fig 4)provides a structured representation of how the system processes input from the user, including depth mapping, color classification, and sensory feedback mechanisms.

10.4.2 *Circuit Diagram.* This circuit diagram (see Fig 5)illustrates the connection between an Arduino Uno (U1) and a ERM motor (M1). The motor is controlled via digital pin D8, which serves as the signal output to drive the motor. The ground (GND) of the motor is connected to the GND of the Arduino, completing the circuit. This setup enables the Manuscript submitted to ACM

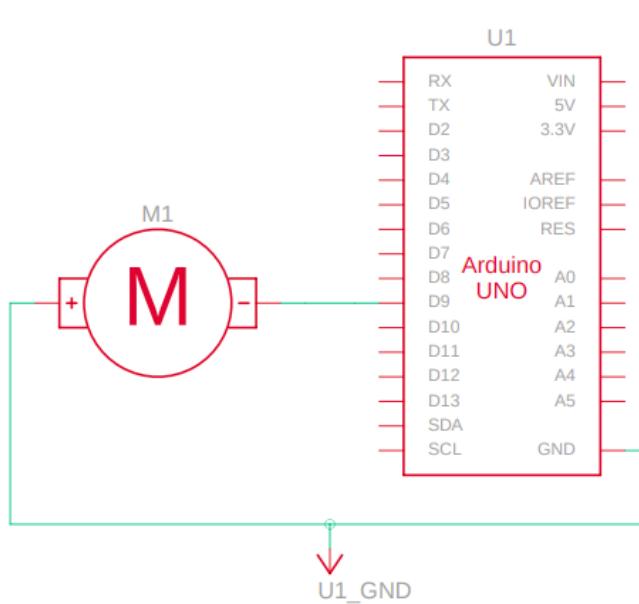


Fig. 7. Circuit Diagram

Arduino to control the motor's activation based on input signals, likely corresponding to vibration feedback in the haptic system.

10.4.3 Connections. Fig 6 showcases the real-world implementation of the system. The laptop runs the depth mapping algorithm, machine learning model, and computer vision algorithm to track finger movements and detect colors. A camera mounted on a tripod captures the interaction. The sheet serves as a reference surface where visually abled participants move their fingers, triggering different sensory feedback.

10.5 Functionalities and Baseline Implementation

10.5.1 Components/Technology Used: Software

- OpenCV (cv2)
- MediaPipe Hands
- NumPy, Pandas, Scikit-Learn (KNN Classifier)
- SoundDevice (sd.play)

Hardware

- Arduino UNO
- Bluetooth module
- ERM Motor
- Camera
- Breadboard

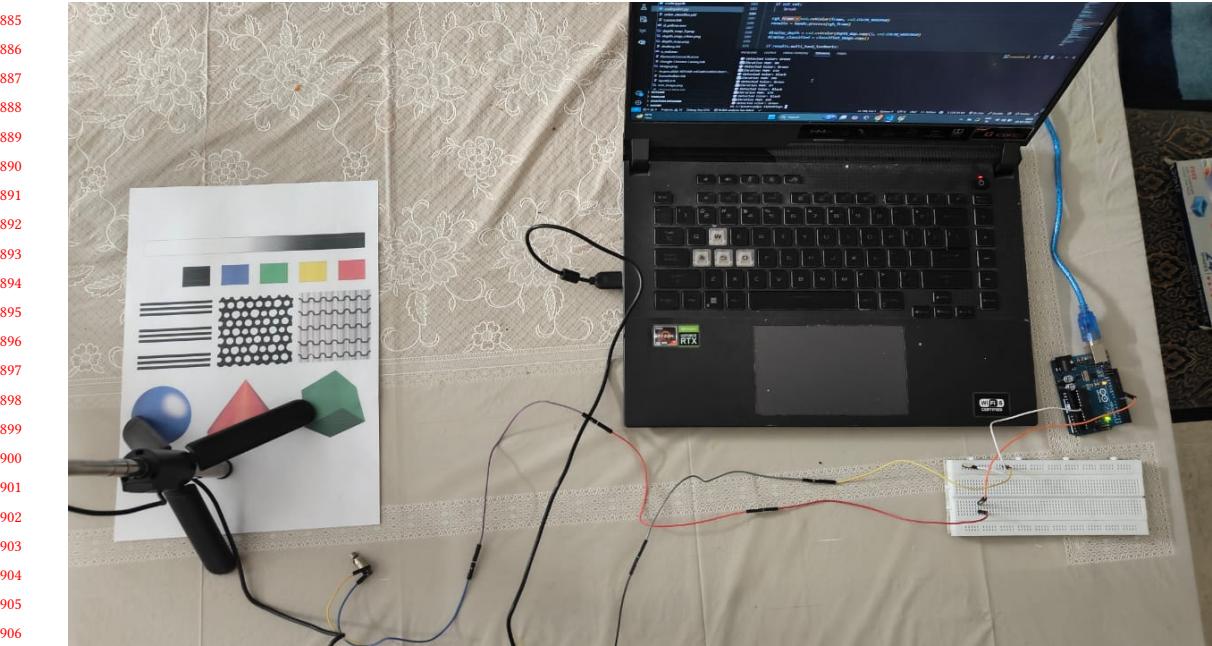


Fig. 8. Component connections

- Jumper wires

10.5.2 *Haptic Feedback: Depth to Vibration.* **Finger Tracking** OpenCV (cv2) is used along with MediaPipe Hands to detect and track the position of fingers in real-time. The HandTracker extracts landmark coordinates (x, y) of key points, particularly the index fingertip (landmark ID: 8). Once the hand is detected, the (x, y) position of the index fingertip is extracted from the image frame. These coordinates are normalized within the frame dimensions to ensure they fit within the required range for mapping.

Depth Mapping Algorithm Using the Depth Mapping algorithm, we generate a depth-map image where pixel intensity represents object depth. Instead of implementing a custom depth estimation model, we utilize Sculpt OK, an online depth estimation tool, as it provides a pre-trained, high-performing model. This ensures better accuracy and efficiency compared to a custom-built solution. The extracted depth data is then processed to map depth values into a scaled range suitable for haptic feedback generation.

Vibration Motor Control The processed depth data is transmitted to an Arduino microcontroller, which controls a ERM vibration motor to provide tactile feedback. The intensity of vibration is inversely proportional to the depth value—closer objects trigger stronger vibrations, while distant objects result in weaker vibrations.

The depth-map image is analyzed to determine the depth value at a specific region (e.g., the center or a tracked point). The motor vibrates with varying intensity based on that point's depth, allowing users to perceive depth through touch.

10.5.3 *Auditory Feedback: Colour to Sounds.* **Training Machine Learning Algorithm** The machine learning model used for color classification is a Scikit-Learn K-Nearest Neighbors (KNN) classifier, trained on a dataset containing RGB values of six distinct colors: Red, Blue, Green, Yellow, Black, and White. Each pixel's RGB triplet served as a feature,

Manuscript submitted to ACM

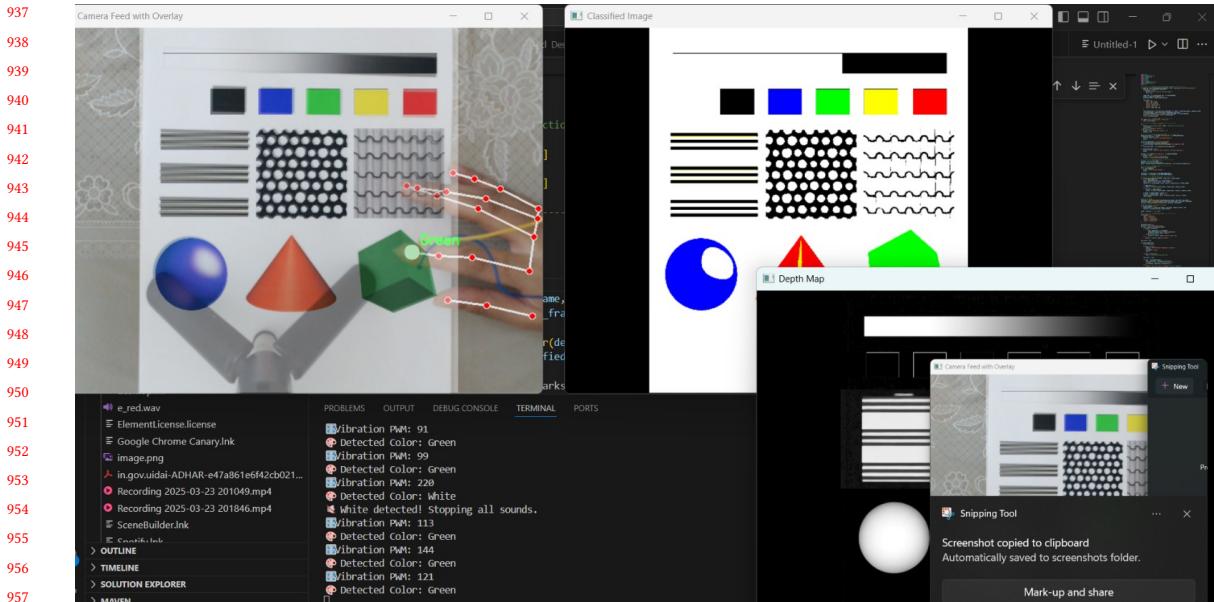


Fig. 9. Working Prototype

while the corresponding color label was the target. Then, the trained model was saved as a pickle file. This was done to ensure that complex colours can be classified into our broader colour base.

Identification of colours and Auditory mapping Once a frame is captured, it is processed using the trained KNN model, which predicts the color of each pixel. The output is then mapped to predefined audio cues that correspond to the detected colors. Each color is assigned a unique frequency pattern as described below :

- White : No Sound
- Red : E note - 659.26 Hz
- Yellow : D note - 587.33 Hz
- Green : C note - 523.25 Hz
- Blue : B note - 493.88 Hz
- Black : A note - 440 Hz

Real-Time Audio Feedback Implementation To provide continuous real-time auditory feedback, the system processes live video frames. Wherever our fingers trace the image, pixel colour values are translated to a colour label. The detected colours are converted into corresponding sound waves, which are played using Python's SoundDevice library.

10.6 Preliminary Results

Utilizing the aforementioned techniques and components, we have successfully developed a functional system. A demonstration of the system in action can be viewed in the following video: [LINK](#) (Also see Figure 7)

989 10.7 Evaluation Metrics

- 990
991 (1) Virtual Texture Perception: By leveraging haptic feedback with varying vibration intensities, the system enables
992 users to perceive texture variations in an image without physically touching it. Different depth levels correspond
993 to different vibration strengths, allowing users to sense surface details, edges, and contours, effectively translating
994 visual texture into a tactile experience.
995
996 (2) Responsiveness : The system ensures seamless real-time interaction—fingers are tracked instantly, the vibration
997 motor responds immediately to depth changes, and sound feedback is generated without delay, creating a
998 dynamic and fluid user experience.
999
1000 (3) F1 Score and Accuracy : The KNN model achieves an impressive F1 score of 0.9715, indicating a strong
1001 balance between precision and recall. We also have an accuracy of 0.9720. To further enhance performance, we
1002 can explore hyperparameter tuning. This could help refine the model's decision boundaries and improve its
1003 classification accuracy.
1004
1005 (4) Learnability : To improve learnability, we employ a structured mapping of sound frequencies that follows a
1006 descending tonal progression—moving from warmer to cooler tones. This gradual descent aligns with natural
1007 perceptual tendencies, making it easier for users to intuitively associate colors with their respective sounds.
1008
1009 (5) LLM Integration : In the future, we plan to integrate large language models (LLMs) to enhance the overall user
1010 experience and accuracy. LLMs can provide contextual assistance by dynamically adjusting sound mappings
1011 based on user preferences, offering natural language feedback, and improving object recognition through
1012 multimodal learning. This integration will make the system more adaptive, intuitive, and personalized for users.
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039