Automating eHMI Action Design with LLMs for Automated Vehicle Communication

Anonymous ACL submission

Abstract

The absence of explicit communication channels between automated vehicles (AVs) and other road users requires the use of external Human-Machine Interfaces (eHMIs) to convey messages effectively in uncertain scenarios. Currently, most eHMI studies employ predefined text messages and manually designed actions to convey these messages, which limits the real-world deployment of eHMIs, where adaptability in dynamic scenarios is essential. 011 Given the generalizability and versatility of large language models (LLMs), they could potentially serve as automated action designers 014 for the message-action design task. To validate this idea, we make three contributions: (1) We propose a pipeline that integrates LLMs and 017 3D renderers, using LLMs as action designers to generate executable actions for controlling eHMIs and rendering action clips. (2) We collect a user-rated Action-Design Scoring dataset comprising a total of 320 action sequences for eight intended messages and four representative eHMI modalities. The dataset validates that LLMs can translate intended mes-025 sages into actions close to a human level, particularly for reasoning-enabled LLMs. (3) We introduce two automated raters, Action Reference Score (ARS) and Vision-Language Models (VLMs), to benchmark 18 LLMs, finding that the VLM aligns with human preferences yet varies across eHMI modalities.¹

1 Introduction

037

Automated vehicles (AVs) promise to redefine transportation systems by eliminating human driving errors and optimizing traffic flow (Fagnant and Kockelman, 2015). However, the absence of a human operator disrupts road interactions, as drivers no longer exchange contextual cues (e.g., eye contact and gestures) to negotiate ambiguous scenar-



c) Demo eye Movement: "Help me out"

Figure 1: Setup illustration and action demos. a) Four types of eHMIs are installed on the vehicle separately; b) Demo actions of the arm convey the message: "Say Hello". The shaded action indicates the subsequent status; c) Demo actions of the eye: "Help me out".

ios (Colley et al., 2025). To bridge this gap, external Human-Machine Interfaces (eHMIs) have emerged as mediators, conveying AV intent (e.g., yielding, turning) to other road users, such as pedestrians, cyclists, and human drivers (Dey et al., 2020a; Colley and Rukzio, 2020). These interfaces use diverse forms, such as displays (Al-Taie et al., 2024; Lim and Kim, 2022), LED strips (Dey et al., 2020b), projections (Eisma et al., 2019), and attached robots (Gui et al., 2024b), to convey vehicle intentions by text, signals, or non-verbal motions.

Current eHMIs are usually designed and analyzed with predefined text messages (e.g., "Please stop.", "I am worried.") with scenario information (e.g., "pedestrian is crossing the road", "robot is stuck in snow") and manually designed actions to perform these messages, as shown in Figure 1. This

¹The source code, prompts, Blender scenarios, and rendered clips are available at https://anonymous.4open.sc ience/r/eHMI_action_design/

094

100

101

102

103

104

105

107

057

restricts the real-world deployment of eHMIs because dynamic interactions demand adaptable communication strategies (Dey et al., 2020a). Therefore, developers must design actions for all possible messages that AVs might need to communicate to other road users. This process is time-intensive, costly, and significantly limits the scalability of eHMIs in practical scenarios (Lim et al., 2024).

Recently, Large Language Models (LLMs) demonstrate generalizability and versatility in multiple tasks, such as reading and answering questions (Radford et al., 2019), as well as pattern following (Mirchandani et al., 2023), suggesting that they may serve as suitable automated action designers for eHMIs. However, it is unclear whether the application of LLMs for eHMIs is feasible and useful, leading to our research question (RQ):

> To what extent do pretrained LLMs achieve parity with human designers in designing eHMI actions that are understandable to road users?

Answering our RQ involves three key challenges. **First**, there is no systematic pipeline for translating specified messages into executable action sequences for eHMIs. **Second**, there is a lack of highquality datasets for testing and improving the translation of eHMI messages into action sequences. **Third**, there is no commonly used benchmark to compare different methods for designing and evaluating eHMI actions fairly.

Therefore, first, we propose a pipeline that integrates LLMs and 3D renderers. To adapt LLMs for controlling eHMIs, we draw inspiration from LLMbased robot action planning (Garrett et al., 2021; Zitkovich et al., 2023), which utilizes LLMs as action designers to generate a series of executable actions to actuate robotic motors. Second, we introduce a user-rated Action-Design Scoring dataset. The dataset comprised eight intended messages for the eHMI to convey by analyzing traffic scenarios, and selected four representative eHMI modalities frequently discussed in eHMI research. We collected messages from previous eHMI studies (Chang et al., 2022; Gui et al., 2022, 2024a) and designed new ones based on message types (Colley and Rukzio, 2020) to enrich the variety. For each message-modality pair, we generated ten actions: eight produced by LLMs (GPT-40 (Achiam et al., 2023), Sonnet 3.5 (Anthropic, 2024), Gemini 2 Flash (DeepMind, 2024), and GPT-o1 (OpenAI, 2024b)) and two designed by human experts. These actions were rendered using Blender version 4.3 (Blender Foundation, 2025), resulting in 320 video clips. We conducted a video-based user study with human participants. They evaluated the understandability of the LLM-designed actions by measuring the consistency between the intended messages and perceived meanings. The Action-Design Scoring dataset provides averaged human scores for each action, enabling a comparative benchmark for existing LLMs. Third, we introduce the Action Reference Score (ARS), which uses the similarity between the newly designed actions and those rated in our dataset. Additionally, we discussed the potential of Vision-Language Models (VLMs) to serve as human-like raters. Then, we benchmark 18 LLMs on this task.

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

157

Contribution Statement: This work proposes the first complete pipeline, along with a comprehensive dataset and benchmark for evaluating eHMIs. Beyond these core contributions, we also share several noteworthy insights as follows:

- Pretrained LLMs can achieve a close humanlevel action design capability (see Section 4.1).
- VLM rater matches human preferences but varies across eHMI modalities (see Section 4.2).
- Reasoning-enabled LLMs demonstrate better performance in our task (see Section 4.3).

2 Related Work

2.1 Rule-Based eHMI Action Planning

Currently, eHMI action planning generally follows a fixed design approach, in which human designers establish behavioral rules based on the specific features of different eHMI modalities. For example, in text- and icon-based eHMIs, designers create content referencing traffic regulation signs or standard messages (Eisele and Petzoldt, 2022; Eisma et al., 2021). In color- and light-band-based eHMIs, they design the content relying on human intuitive empathy and empirical evaluation with colors and blinking frequencies (Bazilinskyy et al., 2019; Dey et al., 2020b). For anthropomorphic eHMIs, such as eyes or arms, designers mimic nonverbal communication cues drawn from common human-human interactions (Mahadevan et al., 2018; Ochiai and Toyoshima, 2011). Most recently, Colley et al. (2025) proposes using Human-In-The-Loop Multi-Objective Bayesian Optimization to create appropriate eHMIs.

In summary, traditionally, experts have observed real-world examples to derive design rules for

guiding eHMI action planning. However, differ-158 ent eHMI modalities vary in expression. Low-159 expressiveness eHMIs, such as arrow icons, are 160 relatively simple because they convey static direc-161 tional cues, making it easier to define behavioral 162 rules (Fridman et al., 2017). Highly expressive 163 eHMIs can produce complex actions and communi-164 cate richer messages (Chang et al., 2024). However, 165 determining behavioral rules for such modalities 166 is challenging due to the increased intricacy and 167 variability of their expressions (Gui et al., 2023; de Winter and Dodou, 2022). Unlike previous 169 works, we address this by evaluating LLMs to sup-170 port eHMI action planning, enabling more complex 171 and dynamic communication. 172

2.2 LLMs-Based Robot Action Planning

173

174

175

176

177

178

179

181

183

185

186

188

189

190

191

194

195

196

197

198

201

206

Recent LLMs encode vast world knowledge and exhibit the emerging capability for robot action planning (Xiang et al., 2024). Regarding how LLMs generate actions to actuate robots, existing approaches fall into two main trends: Task and Motion Planning (TAMP) (Garrett et al., 2021) and Visual Language Action (VLA) models (Zitkovich et al., 2023). TAMP methods break down complex instructions into predefined low-level actions (e.g., grasping, moving) to control robots (Chen et al., 2024). However, for our task, it is difficult to predefine these action categories. We believe that forcing LLMs to choose from rigid action modes limits their ability to design flexible or adaptive actions creatively (Hao et al., 2025). In contrast, VLA models fuse robot control actions directly into VLM backbones, providing specific action commands to control each robotic motor (Zitkovich et al., 2023; Kim et al., 2024). However, applying existing VLA models to out-of-scope tasks with different robot settings often requires a large amount of data for finetuning (Qu et al., 2025). This contradicts our objective of reducing the labor required by human experts in designing eHMI actions. In this task, we utilize the generalizability and versatility of pretrained LLMs by providing detailed prompts on how to control each modality of eHMI.

3 Methodology

This section outlines our responses to three key challenges: i) the LLM-Blender Fusion pipeline, ii) the Action-Design Scoring dataset, and iii) the automated evaluation system for benchmarks. These designs serve as a proof of concept for our RQ and offer a systematic approach to developing and evaluating newer LLMs or eHMIs modalities.

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

3.1 LLM-Blender Fusion Pipeline

The design of the LLM-Blender Fusion Pipeline (see Figure 2(b)) unfolds in two steps: i) Designing eHMI actions using LLMs with the provided message text, scenario information, and eHMI description (see details in Section 3.2.1 and Section 3.2.2), and ii) Rendering the designed actions into corresponding virtual scenarios as video clips in Blender (see Section 3.2.3 for more details).

3.2 Action-Design Scoring Dataset

3.2.1 eHMI Modality Definitions

As shown in Figure 2(a), four representative eHMIs are selected for analysis, categorized into two types: anthropomorphic (human-like) and nonanthropomorphic (Bazilinskyy et al., 2019; Dey et al., 2020a). The selection prioritizes dynamic interfaces that use sequential visual cues (e.g., changing brightness, animations) to communicate intent clearly to other road users (Wilbrink et al., 2021). We craft prompts (see Appendix C) for each eHMI modality, offering detailed guidance to LLMs on what they can control and how to control. Each step of the designed action sequence includes a subsequent status and a transition speed, such as [angle1, angle2, ..., "fast"].

The descriptions for each status of different eHMI modalities are shown as follows:

Eyes. Robotic eyes are mounted on the front of the autonomous vehicle. The pupil's position is specified using polar coordinates: the angle spans $[0^{\circ}, 360^{\circ}]$ (starting from "up" and moving counterclockwise), and the distance spans [0, 1], where 0 denotes the center and 1 is the edge (Chang et al., 2022; Gui et al., 2022).

Arm. A robotic arm is mounted on the top of the vehicle. It is composed of five components, each of which is connected by single-axis rotational joints. The five movable components (shoulder, upper arm, forearm, hand, and fingers) are required to operate within limited ranges (Gui et al., 2024b).

Light Bar. A light bar contains 15 lights arranged in an arc fixed on the front top of the autonomous vehicle. Each light can be either "on" or "off", with uniform brightness and color (Dey et al., 2020b). Facial Expression. A screen located at the front of the vehicle displays a sequence of facial expressions to convey messages. The available facial ex-

pressions are selected from a set of emojis (Al-Taie



Figure 2: Dataset Asset, Pipeline, and Human Scoring. Dataset assets contain four representative eHMIs and eight intended messages from different interaction types. In the pipeline, we develop eight corresponding Blender scenarios and render actions designed by LLMs or human experts to clips. During the human scoring phase, ten participants evaluate each action clip using a five-point Likert scale.

et al., 2024; Dey et al., 2020a).

Regarding transition speed, we offer three options (e.g., "slow", "medium", and "fast"). Additionally, we include a "super fast" option to quickly reset the eHMI to its initial status, ensuring continuity when switching between different meanings in an action sequence. In our practical experiments, we find that providing the concept of transition speed, rather than stating specific times like "1 second," gives LLMs a more accurate sense of timing for designing actions. This approach is beneficial because LLMs may not inherently understand the physical scale of the eHMI (e.g., its size or mounting height) or its spatial relationship to other road users, which could lead to ambiguity in interpreting the real-world impact of transition speeds.

3.2.2 Message Set Design

The communication relationships can be categorized into four types: one-to-one, one-to-many, many-to-one, and many-to-many, where the former (e.g., AVs equipped with an eHMI) interacts with the latter (e.g., pedestrians, cyclists) (Colley and Rukzio, 2020). However, evaluating the collaboration of multiple AVs (e.g., many-to-one, manyto-many) falls outside the scope of this work. Instead, we focus on one-to-one and one-to-many relationships. For one-to-one interactions, we further distinguish between first-person perspectives, where the communicator transmits messages about the AV's state or intent (e.g., "Help me out!"), and third-person perspectives, where the AV relays information about other road users or environmental conditions (e.g., "Pedestrian ahead"), based on different perspective taking (Bazilinskyy et al., 2019). We collect six messages from previous eHMI studies (Chang et al., 2022; Gui et al., 2022, 2024a) and design two new ones based on message types (Colley and Rukzio, 2020) to enrich the variety (see Table 1). Each message includes: 291

292

293

294

295

297

299

301

302

303

304

305

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

324

- A message text needs to be conveyed.
- Scenario information related to the message.
- A user perspective scenario description for the scoring task. (see Appendix B)

3.2.3 Clip Generation and Human Scoring

In the previous section, we obtain a total of 32 modality-message pairs for each eHMI modality and message type (see Figure 2(a)). For each pair, we ask four LLMs (GPT-40 (Achiam et al., 2023), Sonnet 3.5 (Anthropic, 2024), Gemini 2 Flash (DeepMind, 2024), and GPT-01 (OpenAI, 2024b)) to design two distinct actions. Additionally, two human experts also complete this task. This process results in a total of 320 actions.

However, it is implausible for human participants to rate these actions solely based on textbased commands. They need to observe the actual movements of the eHMIs to judge the effectiveness of the designed actions in conveying messages. Therefore, we incorporate the rendering process into our LLM-Blender fusion pipeline(see Figure 2(b)). The rendering assets include eHMI models, vehicle models, and scenarios. For eHMI models, the arm is available under a free license (Sinitsyn, 2021), the eyes are part of a proprietary model (Chang et al., 2022), and the light bar and screen are self-designed. For vehicle models, the AV model is proprietary (Chang et al., 2022), while the delivery-robot model is available under a free li-

290

Table 1: Eight messages collected or designed based on different communication relationships. Each message contains a message text, scenario information, and a user perspective scenario description (see Appendix B).

Case	Message Text	Scenario Information
One-to-one (First-person) comm	unication relationships	
Send intention	"I am unable to pick you up here. Please walk forward in my direction to a suitable pickup spot."	You are an autonomous taxi that receives a ride request and arrives to pick up the passenger (on the right roadside). Upon arrival, you detect the passenger standing in an area where parking is not permitted within a 5 m radius.
Status report	"I am about to start moving. Please watch out."	You are a stopped autonomous vehicle parked near a park, posi- tioned just before a crosswalk. A student is approaching and is about to cross to the other side of the road.
Request help	"I am stuck. Could you please help me out?"	You are a delivery robot that has been trapped by a pile of boxes. Feeling eager to free yourself and continue delivering the items to your customer on time, you notice a passerby who sees your situation but hesitates to assist.
Refuse help	"Thank you for your kind- ness. Please not touch me."	You are an expensive and fragile delivery robot stuck in the snow. You are programmed that only your owner can repair you. Meanwhile, a passerby notices your predicament and hesitates to offer assistance.
One-to-one (Third-person) comm	nunication relationships	
Pedestrian Blind Spot Alert	"Please watch out for a vehicle approaching from your left blind spot."	You are an autonomous vehicle parking near an intersection with no traffic lights. A pedestrian on the opposite side is walking toward the intersection, facing you. A building blocks his view of an approaching bus heading toward the intersection from his left (from your right).
Driver Blind Spot Warning	"Please watch out for the pedestrian approaching from your right blind spot."	You are an autonomous vehicle parked at an intersection without traffic lights. A bus is approaching from the opposite direction. A pedestrian is about to use the crosswalk on the opposite side, coming from your left. However, a building obstructs the bus's view, so it cannot see the pedestrian approaching from its right.
One-to-many communication rel	ationships	
Target Identification	"I am sending the package only to this person."	You are a delivery robot tasked with delivering a package to a customer in a crowded area. Currently, three individuals are standing to your left, front, and right. Your recipient is directly in front of you and is taller than you.
Broadcast Communication	"I am about to turn right. Kindly make a way to avoid conflict."	You are a delivery robot carrying a package in a crowded area. You want to navigate through the crowd and turn right without causing disruptions.

cense (Condra, 2021). For the scenarios, we design the corresponding 3D environments for different messages using Blender version 4.3 (Blender Foundation, 2025), using a paid add-on called *The City Generator 2.0* (Blendermarket, 2025). We use a GPU-equipped device (NVIDIA GTX 4070 Ti) to render these 320 actions into clips, achieving 24 FPS and 1080p resolution to ensure an optimal viewing experience for participants. The entire rendering process takes approximately 100 hours, with each 10-second clip taking an average of about 20 minutes to complete.

325

326

331

334

336

338

341

342

Then, we invite N=40 participants to score the action clips (see Figure 2(c)). Each participant receives 80 random clips, along with the intended messages and the corresponding user perspective scenario information (see Appendix B). They then answer the question: *"How consistently do the*

eHMI actions express the message?" The participants rate each action clip using a 5-point Likert scale (1=Strongly Disagree to 5=Strongly Agree) (Joshi et al., 2015). In contrast to other annotation methods, such as pairwise ranking, the 5-point Likert scale alleviates the participants' load (Rouse et al., 2010; Mantiuk et al., 2012) and reliably reflects their preferences toward different actions (Rankin and Grube, 1980; Zerman et al., 2018). In total, we collect 3,200 scores, each action clip rated by ten different participants. We then calculate the average of these scores, resulting in 320 average scores for the clips.

343

344

345

347

348

349

351

352

353

354

355

356

357

358

359

3.3 Automated Scoring for Benchmarking

In the future, one may evaluate the translation capability of messages to actions of novel LLMs. However, employing human participants to score

367

372

374

379

387

391

400

401

402

403

404

405

406

407

408

generated actions is expensive and time-consuming. To address this, we propose two substitutes.

3.3.1 Action Reference Score (ARS)

We introduce an Action Reference Score (ARS) that automatically generates a score for a new action by retrieving the most similar actions from our dataset, inspired by existing works for similar purposes (Escudero-Arnanz et al., 2023; Wilson and Martinez, 1997). We use Dynamic Time Warping (DTW) (Müller, 2007; Salvador and Chan, 2007; slaypni, 2015) to compute the similarity between actions. DTW is particularly effective because it calculates similarity even when identical patterns appear at different positions or when sequences vary in length. Our approach converts the status of next action steps into numerical values. For example, the angle variable (e.g., 60°) is transformed into its sine and cosine components to capture its cyclical nature. Similarly, categorical variables (e.g., "close") are assigned predefined integer values, and transition times are quantified by assigning "slow" as 4, "medium" as 3, "fast" as 2, and "super fast" as 1. In defining the distance function for the DTW algorithm, we assign an equal weight of 1 to numerical, categorical, and temporal elements, normalizing each element's value range to [0, 1].

3.3.2 Vision-Language Model (VLM) Rater

We also evaluate whether the designed actions are contextually appropriate and semantically consistent with the intended messages by leveraging the multimodal understanding and reasoning capabilities inherent in VLMs (Zhang et al., 2023; Gu et al., 2024). For each action clip used for the VLM evaluation, we ensure that VLMs can detect subtle variations by adjusting the camera in Blender to zoom in and focus on the autonomous vehicle equipped with the eHMI. Each rendered frame has a resolution of 512×512 , with the autonomous vehicle equipped with the eHMI dominating the composition. These clips are rendered at six FPS, ensuring that the total number of frames does not exceed the maximum image series length of the VLM while preserving sufficient dynamic details. The reduced resolution and FPS also expedite the rendering process to an average of two min per clip. In the prompt (see Appendix D) accompanying the clips provided to the VLM, we request the model to assign a continuous score ranging from 1 to 5, using the same criteria as human participants.

4 Experiments

proposed new dataset.

In this section, the experiments are designed to achieve three specific purposes.
Analyze the collected Action-Design Scoring dataset to answer our RQ proposed in Section 1.
Discuss the viability of the VLM rater as a replacement for human raters.
Benchmark various types of LLMs using our

4.1 Performance Evaluation on Action-Design Scoring Dataset

Table 2 reports the statistics of our Action-Design Scoring dataset, and Figure 4 compares humanrated score distributions across four LLMs and human designers. Our key findings are as follows:

Pretrained LLMs can achieve close humanlevel action design capability. Table 2 shows that LLMs perform comparably to human designers. In particular, the average score of GPT-o1 closely matches that of human designers. We calculate a Wilcoxon signed-rank test (Woolson, 2005) to assess statistical significance: GPT-o1 does not differ significantly from human raters (p = 0.69), whereas all other sources differ from human designers at p < 0.01. Figure 4 (in Appendix) illustrates the same trend: human designers most frequently award a score of 5 (Strongly Agree), followed by 4 (Agree); GPT-o1 ranks second for 5 and first for 4. Furthermore, when broken down by message type and eHMI modality, GPT-o1 outperforms humans for the eyes modality (mean = 2.795 vs. 2.536) and in third-person messages (3.098 vs. 3.045).

Message type and eHMI modality affect design quality. Third-person messages receive significantly higher ratings than other types (p < 0.01), likely because "Watch out" type messages are easier to design. Among eHMI modalities, the arm modality outperforms all others (p < 0.01), while facial expressions score lower (p < 0.05). Since most of our eight scenarios convey spatial information, the arm modality is especially effective; the absence of emotional messages (e.g., "I am scared") limits facial expressions' performance.

Finally, we validate our data by computing the inter-rater reliability (IRR) using Krippendorff's alpha (Wong et al., 2021). The moderate alpha value confirms that our dataset is reliable.

409

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

Source (Designer)	Average	Message types			eHMI modalities					
Source (Designer)	Average	1 st	3 nd	1-to-N	eyes	arm	facial expression	light bar	IKK	
GPT-40	2.404	2.375	2.250	2.616	2.509	2.616	2.223	2.268	0.399	
Sonnet 3.5	2.538	2.464	2.768	2.455	2.554	2.554	2.429	2.616	0.325	
Genimi 2.0 Flash	2.563	2.460	2.911	2.420	2.554	2.920	2.304	2.473	0.361	
GPT-o1	2.728	2.509	3.098	2.795	2.795	2.982	2.509	2.625	0.436	
Human	2.768	2.580	3.045	2.866	2.536	3.107	2.643	2.786	0.478	

Table 2: Statistics of the Action-Design Scoring dataset: The average scores indicate that LLMs perform comparably to human designers across various messages and eHMI modalities. Krippendorff's alpha is also calculated to assess Inter-Rater Reliability (IRR) among human raters.

oHMI modelities	Metrics						
ernvir modanties	$m{r}$ _{p-value}	$oldsymbol{ au}$ p-value	pair.(%)				
eye	0.432 < 0.01	0.352 < 0.01	72.73				
arm	0.547 < 0.01	0.442 < 0.01	83.87				
facial expression	0.368 < 0.01	0.292 < 0.01	62.50				
light bar	$0.242_{=0.03}$	0.221 = 0.01	57.30				

Table 3: Association between scores from human rater scores and those from the VLM rater (Qwen-QvQ-Max) measured using three metrics: Pearson's r, Kendall's τ , and *pairwise accuracy*.

4.2 VLM Rater Alignment Evaluation

We conduct an additional experiment to evaluate whether VLMs can assess action clips in a manner similar to that of human raters. We present the clips in a format that the VLMs can understand more easily (see Section 3.3.2) and instruct them to rate these clips. We select Qwen-QvQ-Max (Qwen Team, 2025) as our VLM rater, taking into account factors such as cost, inference speed, and the maximum allowable input image series length. Compared to other VLMs, Qwen-QvQ-Max also demonstrates preferences that closely resemble human judgments. Results from other VLM models can be found in Appendix E. We rate each clip using the VLM rater twice and average these scores to determine the final score.

We evaluate the results using three metrics: Pearson's r, Kendall's τ , and a specially designed *pairwise accuracy* (Liu et al., 2009). Pearson's r measures the strength of a linear relationship by assessing the degree of correlation between scores, focusing on how far apart the scores are overall. In contrast, Kendall's τ evaluates the order of the data by comparing the number of concordant and discordant pairs, thus analyzing the consistency of the ordering rather than the magnitude of the differences. The *pairwise accuracy* metric, similar to Kendall's τ , measures the proportion of item pairs where the model's predicted order matches

the ground truth order, specifically among pairs where the model's predicted scores differ by more than a specified threshold. We find that a threshold of 0.7 is the most suitable and adopt it in our analysis. We present statistics for the four eHMI modalities separately in Table 3. 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

521

522

523

524

The VLM rater shows alignment with human scoring preferences but is influenced by eHMI modalities. We observe that for the modalities of eye and arm, the VLM rater achieves a moderate level across all three metrics. Particularly in terms of *pairwise accuracy*, results indicate that, after setting an appropriate threshold to filter out difficultto-rank pairs, the preferences of VLM show clear consistency with those of human raters. However, for the facial expression and light bar modalities, we find relatively low performance on the three metrics. The results suggest that VLM shows a low-level correlation with human raters for these two modalities. We identify two main reasons for this discrepancy: first, upon reviewing the "reasoning process" of VLM scoring, we notice that VLM consistently fails to recognize changes in the light bar modality (for example, transitioning from "on" to "off"). It tends to perceive the situation as "The light of the light bar is always on," which ultimately leads to lower scores. Second, similar to human raters, we notice that VLM insists that the modality of facial expressions alone does not accurately convey the entire message, leading to lower scores.

The VLM rater does not exhibit the necessary bias towards the length of actions as human raters. Figure 3 compares the rendered action clip lengths as evaluated by two scoring sources: human raters and VLM. Among human raters, there is a clear preference for shorter clips. This trend is particularly evident for the eHMI modalities "eyes" and "light bar", where raters tend to favor actions that convey the intended message quickly. In contrast, VLM raters do not exhibit a distinct pref-

481

482

483

484



Figure 3: Relationship between action clip length and evaluation scores. The plot compares scores from human raters and the VLM rater (Qwen-QvQ-Max).

erence for clip length across the different eHMI modalities, not showing enough "bias" towards clip lengths. Besides, the scores of VLM raters are always higher than those given by human raters.

4.3 Benchmarking LLMs Performance

To evaluate the performance of various LLMs that differ in size and architecture, we benchmark 18 models using two complementary metrics: the ARS metric (Section 3.3.1) and the VLM rater (Section 3.3.2), as summarized in Table 4. This selection comprises six proprietary models: GPT-o4-mini (OpenAI, 2025b), Sonnet 3.7 (Anthropic, 2025), Gemini 2.5 Flash (Google Deep-Mind, 2025), GPT-4.1 series (GPT-4.1, GPT-4.1mini and GPT-4.1-nano) (OpenAI, 2025a); two Deepseek models, Deepseek-R1 (Guo et al., 2025) with reasoning capability and Deepseek-V3 (Liu et al., 2024) without reasoning capability; and five variants of the Qwen 3 series (Yang et al., 2025) with 235B, 32B, 8B, 1.7B and 0.6B parameters that are tested both with and without reasoning capability. We rate each clip using ARS and VLM rater. The VLM rater score is calculated by using the VLM rater twice and then averaging these scores to determine the final score.

Reasoning-enabled LLMs demonstrate better performance in designing eHMI actions.

Source (Designer)	Human	ARS	VLM Rater				
Human	2.768	-	3.396				
Proprietary models (Designers)							
GPT-40	2.404	-	3.223				
Sonnet3.5	2.538	-	3.258				
Gemini2 Flash	2.563	-	3.289				
GPT-o1	2.728	-	3.303				
Proprietary models							
GPT-o4-mini	-	2.754	3.352				
Sonnet3.7	-	2.676	3.250				
Gemini2.5 Flash	-	2.571	3.200				
GPT-4.1	-	2.632	3.233				
GPT-4.1-mini	-	2.558	3.213				
GPT-4.1-nano	-	2.596	3.080				
Open source models (Wit	th reasonin	g)					
Deepseek-R1	-	2.766	3.369				
Owen3-235B-a22B	-	2.696	3.339				
Owen3-32B	-	2.583	3.366				
Qwen3-8B	-	2.598	3.333				
Qwen3-1.7B	-	2.596	3.307				
Qwen3-0.6B	-	2.607	3.257				
Open source models (Wi	thout reaso	ning)					
Deepseek-V3	-	2.504	3.292				
Qwen3-235B-a22B	-	2.547	3.283				
Qwen3-32B	-	2.533	3.207				
Qwen3-8B	-	2.498	3.210				
Qwen3-1.7B	-	2.546	3.148				
Owen3-0.6B	-	2.500	3.125				

Table 4: Benchmark for different LLMs using ARS and VLM rater.

As shown in Table 4, both the ARS metric and the VLM rater assign higher average scores to reasoning-enabled LLMs (e.g., GPT-o4-mini and Deepseek-R1). Regarding the Qwen 3 series, the results indicate that when the reasoning capability is enabled, these models produce more human-like eHMI actions, especially with a longer reasoning process. For smaller models like Qwen3-1.7B, enabling reasoning capabilities allows them to outperform larger models that lack this function, such as Deepseek-V3 and Qwen3-235B-a22B. 552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

5 Conclusion

In conclusion, this work proposes the first LLM-Blender Fusion pipeline to design eHMI actions. Alongside this, we introduce the Action-Design Scoring dataset. Our findings suggest that pretrained LLMs can attain a nearly human-level capability in action design. Additionally, we provide a benchmark that can be used to evaluate the capability of other LLMs. Our work establishes a solid foundation for LLM-based action design and the real-world application of eHMIs.

540

541

543

545

547

551

575 576

580

583

584

585

586

591

592

596

599

602

611

612

613

614

616

617

619

620

621

624

Our work represents an important step forward in incorporating LLMs into the eHMI system. However, challenges remain.

Limitations

Unnecessary time cost on Blender rendering. We use Blender to render actions into clips in two steps (see Section 3.2.3 and 3.3.2). Our current work aims to use a realistic virtual background that human participants and VLM raters can use as additional clues for judgment when AVs equipped with eHMIs move in the scene. However, we identify two drawbacks that can be improved: First, the complexity of the designed scenarios greatly influences the rendering time. Second, objects outside of the camera's view still impact the rendering speed. To address these issues, there are two potential solutions: 1) Reduce the complexity of the scenarios and remove objects that do not significantly affect the final rendering results. 2) Switch from Blender to another rendering engine. However, given the mature Python package available for Blender, finding a suitable replacement may be difficult.

Significant effort is dedicated to designing prompts for each eHMI modality. For active eHMIs, experts can craft these instructions within a practical timeframe, but the process demands meticulous trial and error to ensure LLMs execute actions as intended. For passive eHMIs, however, the challenge is far greater: unpredictable behaviors (e.g., a teddy bear's limbs swaying freely on a pole) make manual prompt engineering impractical. Human designers cannot predefine control logic for such open-ended motions, as even basic movements depend on environmental factors like airflow or physics. To address this gap, an automated pipeline could leverage VLM raters validated in our studies as reliable evaluators to generate annotated training data from passive eHMI interactions. By finetuning LLMs on this feedback, we could enable dynamic adaptation to unpredictable behaviors, bridging the divide between scripted and emergent interactions.

Legality and accountability are important topics to discuss. Although our study suggests that pretrained LLMs can achieve near-human-level performance in designing eHMI actions, real-world deployment also requires a parallel analysis of pedestrian trust, confidence in interpretation, and accountability frameworks. For example, a pedestrian might correctly interpret an eHMI warning but disregard it due to distrust or conflicting situational625awareness, raising questions about liability beyond626technical performance. Future work should decouple evaluations into two strands: one optimizing628eHMI design for clarity and reliability, and another629exploring human-AI interaction in terms of trust630calibration and legal implications.631

632

633

634

635

636

637

638

639

640

671

672

Ethics Statement

All data in the Action-Design Scoring dataset have been de-identified to safeguard privacy concerns. Our data construction processes are conducted by skilled researchers. The participants include students from Chinese and Japanese universities, all of whom receive fair honoraria for their contributions.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama 641 Ahmad, Ilge Akkaya, Florencia Leoni Aleman, 642 Diogo Almeida, Janko Altenschmidt, Sam Altman, 643 Shyamal Anadkat, et al. 2023. Gpt-4 technical report. 644 arXiv preprint arXiv:2303.08774. 645 Ammar Al-Taie, Graham Wilson, Euan Freeman, Frank 646 Pollick, and Stephen Anthony Brewster. 2024. Light 647 it up: Evaluating versatile autonomous vehicle-648 cyclist external human-machine interfaces. In Pro-649 ceedings of the CHI Conference on Human Factors 650 in Computing Systems, pages 1-20. 651 Anthropic. 2024. Claude 3.5 sonnet. https://ww 652 w.anthropic.com/news/claude-3-5-sonnet. 653 Accessed: 2025-02-15. 654 Anthropic. 2025. Claude 3.7 sonnet system card. http 655 s://www.anthropic.com/claude-3-7-sonnet-s 656 ystem-card. Accessed: 2025-05-18. 657 Pavlo Bazilinskyy, Dimitra Dodou, and Joost De Winter. 658 2019. Survey on ehmi concepts: The effect of text, 659 color, and perspective. Transportation research part 660 *F: traffic psychology and behaviour*, 67:175–194. 661 Blender Foundation. 2025. Home of the blender project 662 free and open 3d creation software. Accessed: 663 2025-05-09. 664 Blendermarket. 2025. The city generator. https://bl 665 endermarket.com/products/the-city-generat 666 or. Accessed: 15 February 2025. 667 Chia-Ming Chang, Koki Toda, Xinyue Gui, Stela H Seo, 668 and Takeo Igarashi. 2022. Can eyes on a car reduce 669 traffic accidents? In Proceedings of the 14th interna-670

tional conference on automotive user interfaces and

interactive vehicular applications, pages 349-359.

Xiang Chang, Zihe Chen, Xiaoyan Dong, Yuxin Cai, Tingmin Yan, Haolin Cai, Zherui Zhou, Guyue Zhou, and Jiangtao Gong. 2024. " it must be gesturing towards me": Gesture-based interaction between autonomous vehicles and pedestrians. In *Proceedings* of the CHI Conference on Human Factors in Computing Systems, pages 1–25.

673

674

675

681

690

691

701

703

704

707

710

711

712

715

718

720

721

722

723

725

- Yongchao Chen, Jacob Arkin, Charles Dawson, Yang Zhang, Nicholas Roy, and Chuchu Fan. 2024. Autotamp: Autoregressive task and motion planning with llms as translators and checkers. In 2024 IEEE International conference on robotics and automation (ICRA), pages 6695–6702. IEEE.
- Mark Colley, Pascal Jansen, Mugdha Keskar, and Enrico Rukzio. 2025. Improving external communication of automated vehicles using bayesian optimization. In Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems, CHI '25, New York, NY, USA. Association for Computing Machinery.
 - Mark Colley and Enrico Rukzio. 2020. A design space for external communication of autonomous vehicles. In 12th International Conference on Automotive User Interfaces and Interactive Vehicular Applications, pages 212–222.
- Jack Condra. 2021. Starship delivery robot model. ht tps://sketchfab.com/3d-models/starship-d elivery-robot-model-4aef60939c3743cdbae ce6b1e5bda21a. 3D model. Licensed under CC Attribution. Accessed: 2025-05-13.
- Joost de Winter and Dimitra Dodou. 2022. External human-machine interfaces: Gimmick or necessity? *Transportation research interdisciplinary perspectives*, 15:100643.
- DeepMind. 2024. Gemini flash. https://deepmind.g oogle/technologies/gemini/flash/. Accessed: 2025-02-15.
- Debargha Dey, Azra Habibovic, Andreas Löcken, Philipp Wintersberger, Bastian Pfleging, Andreas Riener, Marieke Martens, and Jacques Terken. 2020a. Taming the ehmi jungle: A classification taxonomy to guide, compare, and assess the design principles of automated vehicles' external human-machine interfaces. *Transportation Research Interdisciplinary Perspectives*, 7:100174.
- Debargha Dey, Azra Habibovic, Bastian Pfleging, Marieke Martens, and Jacques Terken. 2020b. Color and animation preferences for a light band ehmi in interactions between automated vehicles and pedestrians. In *Proceedings of the 2020 CHI conference on human factors in computing systems*, pages 1–13.
- Daniel Eisele and Tibor Petzoldt. 2022. Effects of traffic context on ehmi icon comprehension. *Transportation research part F: traffic psychology and behaviour*, 85:1–12.

- Yke Bauke Eisma, Anna Reiff, Lars Kooijman, Dimitra Dodou, and Joost CF de Winter. 2021. External human-machine interfaces: Effects of message perspective. *Transportation research part F: traffic psychology and behaviour*, 78:30–41.
- Yke Bauke Eisma, Steven van Bergen, SM Ter Brake, MTT Hensen, Willem Jan Tempelaar, and Joost CF de Winter. 2019. External human–machine interfaces: The effect of display location on crossing intentions and eye movements. *Information*, 11(1):13.
- Óscar Escudero-Arnanz, Antonio G Marques, Cristina Soguero-Ruiz, Inmaculada Mora-Jiménez, and Gregorio Robles. 2023. dtwparallel: A python package to efficiently compute dynamic time warping between time series. *SoftwareX*, 22:101364.
- Daniel J Fagnant and Kara Kockelman. 2015. Preparing a nation for autonomous vehicles: opportunities, barriers and policy recommendations. *Transportation Research Part A: Policy and Practice*, 77:167–181.
- Lex Fridman, Bruce Mehler, Lei Xia, Yangyang Yang, Laura Yvonne Facusse, and Bryan Reimer. 2017. To walk or not to walk: Crowdsourced assessment of external vehicle-to-pedestrian displays. *arXiv preprint arXiv:1707.02698*.
- Caelan Reed Garrett, Rohan Chitnis, Rachel Holladay, Beomjoon Kim, Tom Silver, Leslie Pack Kaelbling, and Tomás Lozano-Pérez. 2021. Integrated task and motion planning. *Annual review of control, robotics, and autonomous systems*, 4(1):265–293.
- Google DeepMind. 2025. Gemini 2.5: Our most intelligent ai model. https://blog.google/techno logy/google-deepmind/gemini-model-thinkin g-updates-march-2025/. Accessed: 2025-05-18.
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, et al. 2024. A survey on Ilm-as-a-judge. *arXiv preprint arXiv:2411.15594*.
- Xinyue Gui, Chia-Ming Chang, Stela H Seo, Koki Toda, and Takeo Igarashi. 2024a. Scenarios exploration: How ar-based speech balloons enhance car-topedestrian interaction. In *International Conference on Human-Computer Interaction*, pages 223–230. Springer.
- Xinyue Gui, Mikiya Kusunoki, Bofei Huang, Stela Hanbyeol Seo, Chia-Ming Chang, Haoran Xie, Manabu Tsukada, and Takeo Igarashi. 2024b. Shrinkable arm-based ehmi on autonomous delivery vehicle for effective communication with other road users. In *Proceedings of the 16th International Conference on Automotive User Interfaces and Interactive Vehicular Applications*, pages 305–316.
- Xinyue Gui, Koki Toda, Stela Hanbyeol Seo, Chia-Ming Chang, and Takeo Igarashi. 2022. "i am going this way": Gazing eyes on self-driving car show multiple driving directions. In *Proceedings of the 14th*

faces and interactive vehicular applications, pages Mantiuk. 2012. Comparison of four subjective meth-319-329. ods for image quality assessment. In Computer graphics forum, volume 31, pages 2478–2491. Wiley Online Library. Xinyue Gui, Koki Toda, Stela Hanbyeol Seo, Felix Martin Eckert, Chia-Ming Chang, Xiang'Anthony Chen, and Takeo Igarashi. 2023. A field study on pedes-Suvir Mirchandani, Fei Xia, Pete Florence, Brian Ichter, trians' thoughts toward a car with gazing eyes. In Danny Driess, Montserrat Gonzalez Arenas, Kanishka Rao, Dorsa Sadigh, and Andy Zeng. 2023. Extended Abstracts of the 2023 CHI Conference on Large language models as general pattern machines. Human Factors in Computing Systems, pages 1–7. arXiv preprint arXiv:2307.04721. Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Meinard Müller. 2007. Dynamic time warping. Infor-Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, mation retrieval for music and motion, pages 69-84. Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforce-Yoichi Ochiai and Keisuke Toyoshima. 2011. Homuncument learning. arXiv preprint arXiv:2501.12948. lus: the vehicle as augmented clothes. In Proceedings of the 2nd Augmented Human International Con-Peng Hao, Shaowei Cui, Junhang Wei, Tao Lu, Ying*ference*, pages 1–4. hao Cai, and Shuo Wang. 2025. Learn-gen-plan: Bridging the gap between vision language models OpenAI. 2024a. Gpt-4o mini. https://openai.com and real-world long-horizon dexterous manipulations. /index/gpt-4o-mini. Accessed: 2025-02-15. IEEE Transactions on Automation Science and Engineering. OpenAI. 2024b. Introducing openai o1-preview. http s://openai.com/index/introducing-openai-o Ankur Joshi, Saket Kale, Satish Chandel, and D Kumar 1-preview/. Accessed: 2025-02-15. Pal. 2015. Likert scale: Explored and explained. British journal of applied science & technology, OpenAI. 2025a. Introducing GPT-4.1 in the api. https: 7(4):396. //openai.com/index/gpt-4-1/. Accessed: 2025-05-18. Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael OpenAI. 2025b. Introducing openai o3 and o4-mini. Rafailov, Ethan Foster, Grace Lam, Pannag Sanhttps://openai.com/index/introducing-o3-a keti, et al. 2024. Openvla: An open-source visionnd-o4-mini/. Accessed: 2025-05-18. language-action model, 2024. URL https://arxiv. org/abs/2406.09246. Delin Qu, Haoming Song, Qizhi Chen, Yuanqi Yao, Xinyi Ye, Yan Ding, Zhigang Wang, JiaYuan Gu, Bin Dokshin Lim and Byungwoo Kim. 2022. Ui design of Zhao, Dong Wang, et al. 2025. Spatialvla: Explorehmi of autonomous vehicles. International Journal ing spatial representations for visual-language-action of Human-Computer Interaction, 38(18-20):1944model. arXiv preprint arXiv:2501.15830. 1961. Qwen Team. 2025. QVQ-Max: Think with Evidence. Dokshin Lim, Yongjun Kim, YeongHwan Shin, and https://qwenlm.github.io/blog/qvq-max-pre Min Seo Yu. 2024. External human-machine interview/. Blog post; accessed 18 May 2025. faces of autonomous vehicles: Insights from observations on the behavior of game players driving con-Alec Radford, Jeffrey Wu, Rewon Child, David Luan, ventional cars in mixed traffic. Vehicles, 6(3):1284-Dario Amodei, Ilya Sutskever, et al. 2019. Language 1299. models are unsupervised multitask learners. OpenAI blog, 1(8):9. Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi William L Rankin and Joel W Grube. 1980. A com-Deng, Chenyu Zhang, Chong Ruan, et al. 2024. parison of ranking and rating procedures for value Deepseek-v3 technical report. arXiv preprint system measurement. European Journal of Social arXiv:2412.19437. Psychology, 10(3):233–246. Tie-Yan Liu et al. 2009. Learning to rank for informa-David M Rouse, Romuald Pépion, Patrick Le Callet, tion retrieval. Foundations and Trends® in Informaand Sheila S Hemami. 2010. Tradeoffs in subjective tion Retrieval, 3(3):225–331. testing methods for image and video quality assessment. In Human Vision and Electronic Imaging XV, Karthik Mahadevan, Sowmya Somanath, and Ehud volume 7527, pages 108-118. SPIE. Sharlin. 2018. Communicating awareness and in-Stan Salvador and Philip Chan. 2007. Toward accutent in autonomous vehicle-pedestrian interaction. In Proceedings of the 2018 CHI conference on human rate dynamic time warping in linear time and space. factors in computing systems, pages 1–12. Intelligent data analysis, 11(5):561-580.

782

783

785

787

790

791

793

794

797

799

801

804

809

810

811

812

813

814

815

816

817

818

827

830

833 834 international conference on automotive user inter-

Rafał K Mantiuk, Anna Tomaszewska, and Radosław

835

836

837

838

839

840

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

- 889 891 893
- 900 901
- 902 903
- 904 905
- 906 907
- 908 909 910 911 912
- 913 914 915
- 916 917 918 919
- 923

926

927

928

929

930 931

932

- Tim Sinitsyn. 2021. Robot arm simple rigged. https: //www.cgtrader.com/free-3d-models/indust rial/industrial-machine/robot-arm-simpl e-rigged. Royalty-free license; Accessed: 2025-05-15.
 - slaypni. 2015. fastdtw: A Python implementation of FastDTW. Accessed: 2025-05-16.
 - Marc Wilbrink, Merle Lau, Johannes Illgner, Anna Schieben, and Michael Oehl. 2021. Impact of external human-machine interface communication strategies of automated vehicles on pedestrians' crossing decisions and behaviors in an urban environment. Sustainability, 13(15):8396.
 - D Randall Wilson and Tony R Martinez. 1997. Improved heterogeneous distance functions. Journal of artificial intelligence research, 6:1-34.
 - Ka Wong, Praveen Paritosh, and Lora Aroyo. 2021. Cross-replication reliability-an empirical approach to interpreting inter-rater reliability. arXiv preprint arXiv:2106.07393.
 - Robert F Woolson. 2005. Wilcoxon signed-rank test. Encyclopedia of biostatistics, 8.
 - Jiannan Xiang, Tianhua Tao, Yi Gu, Tianmin Shu, Zirui Wang, Zichao Yang, and Zhiting Hu. 2024. Language models meet world models: Embodied experiences enhance language models. Advances in neural information processing systems, 36.
 - An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. 2025. Qwen3 technical report. arXiv preprint arXiv:2505.09388.
 - Emin Zerman, Vedad Hulusic, Giuseppe Valenzise, Rafał K Mantiuk, and Frédéric Dufaux. 2018. The relation between mos and pairwise comparisons and the importance of cross-content comparisons. Electronic Imaging, 30:1-6.
 - Xinlu Zhang, Yujie Lu, Weizhi Wang, An Yan, Jun Yan, Lianke Qin, Heng Wang, Xifeng Yan, William Yang Wang, and Linda Ruth Petzold. 2023. Gpt-4v (ision) as a generalist evaluator for vision-language tasks. arXiv preprint arXiv:2311.01361.
- Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, Ayzaan Wahid, et al. 2023. Rt-2: Visionlanguage-action models transfer web knowledge to robotic control. In Conference on Robot Learning, pages 2165-2183. PMLR.

1004

1005

1006

1007

1008

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

977

978

979

A Cost Analysis

933

937

943

945

947

951

954

955

956

957

959

961

962

963

964

965

967

970

971

972

973

974

975

976

934The costs in this study are primarily incurred in935three areas: user study honoraria, dataset asset cre-936ation, and LLM API calls.

User Study Honoraria Each participant receives an honorarium of \$10, resulting in a total expense of \$400.

940Dataset Asset CreationTo expedite the devel-941opment of city scenarios, we purchase a premium942Blender add-on called *The City Generator* for \$60.

LLM API Calls We utilize online APIs from multiple sources:

- For proprietary models (including GPT-40, GPT-40-mini, GPT-01, GPT-04-mini, the GPT-4.1 series, Sonnet 3.5, Sonnet 3.7, Gemini 2 Flash, and Gemini 2.5 Flash), we access the APIs available on their official websites, which incur a total cost of \$90.
 - For open-source models (such as Deepseek-R1, Deepseek-V3, Qwen-QvQ-Max, and the Qwen 3 series), we utilize both free and paid services offered by Siliconflow², Aliyun Bailian³, and ModelScope⁴, resulting in a total cost of \$50.

Total The overall cost for the study \$600.

B User perspective scenario description

The following descriptions are provided to both human participants and VLM raters to encourage them to consider the perspectives of other road users and make assessments.

First-person scenario descriptions:

Send intention You are a pedestrian standing on the right roadside, waiting for an autonomous taxi.
However, the taxi informs you that it cannot pick you up at your current location due to parking restrictions within a 5-meter radius. The taxi sends you the following message: "I am unable to pick you up here. Please walk forward in my direction to a suitable pickup spot."

Status report You are a student approaching a crosswalk near a park. A stopped autonomous vehicle, positioned just before the crosswalk, plans to start moving soon. The vehicle sends you the following message to get your attention: "I am about to start moving. Please watch out."

Request help You are a passerby noticing a delivery robot trapped by a pile of boxes (or possibly pushed). The robot, eager to continue delivering items on time, sees you hesitating and sends the following message to encourage your help: "I am stuck. Could you please help me?"

Refuse help You are a passerby who notices a fragile and expensive delivery robot stuck in the snow due to its low wheels. As you consider offering assistance, the robot informs you that its owner is on the way and sends the following polite message: "Thank you for your kindness. Please refrain from touching me."

Third-person scenario descriptions:

Pedestrian Blind Spot Alert You are a pedestrian walking toward an intersection near an autonomous vehicle. However, a building blocks your view of an approaching bus from your left. The vehicle, aware of the danger, sends you the following urgent message to ensure your safety: "Please watch out for the vehicle coming from your left blind spot." **Driver Blind Spot Warning** You are a bus driver approaching an intersection with no traffic lights. A pedestrian is preparing to cross the road from your right, but your view is obstructed by a building. A stopped autonomous vehicle at the scene sends you the following message to ensure pedestrian safety: "Caution: Please watch out for the pedestrian coming from your right blind spot."

One-to-many scenario descriptions:

Target Identification You are one of three individuals standing in a crowded area, and a delivery robot approaches with a package. The recipient is the second person from the leftmost side, taller than the robot. To avoid confusion, the robot sends a message to everyone: "I am sending the package only to this person."

Broadcast Communication You are part of a crowded intersection where a delivery robot carrying a package is trying to navigate through. The robot intends to turn right and sends the following message to avoid disruptions: "I am about to turn right. Kindly make a way to avoid any conflict."

C eHMI description prompts

The system prompts are structured into four sec-
tions: character profile, eHMI description, demon-
stration actions, and design guidance. Figure 6
presents the prompt for the eye; Figure 7 shows the
prompt for the arm; Figure 8 is for the light bar; and
Figure 9 depicts the prompt for facial expressions.1021
1022

²https://cloud.siliconflow.cn

³https://cn.aliyun.com/product/bailian

⁴https://www.modelscope.cn/



Figure 4: Comparative Distribution of Action-Design Scoring, where each action clip is rated using a 5-point Likert scale. Human designers are most frequently awarded a score of 5 (Strongly Agree), while GPT-01 received the highest number of 4 (Agree) scores.

oHMI modulities	Qwen-QvQ-Max			G	PT-4.1-mini	i	GPT-40-mini [†]			
	$m{r}$ _{p-value}	$oldsymbol{ au}$ p-value	pair.(%)	p-value	$oldsymbol{ au}$ p-value	pair.(%)	p-value	$oldsymbol{ au}$ p-value	pair.(%)	
eye	0.432 0.001	0.352 0.001	72.73	0.416 0.001	0.218 0.012	62.00	0.395 0.007	0.310 0.008	55.16	
arm	0.547 0.001	$0.442_{\ 0.001}$	83.87	0.558 0.001	$0.407_{\ 0.001}$	78.26	0.387 0.009	0.238 0.013	56.86	
facial expression	0.368 0.001	0.292 0.001	62.50	0.356 0.001	0.278 0.001	64.29	0.349 0.001	0.295 0.001	52.28	
light bar	0.242 0.031	0.221 0.010	57.30	0.272 0.007	$0.160_{\ 0.071}$	50.46	0.284 0.033	$0.240_{0.010}$	46.21	

Table 5: Association between scores from human raters and that from all VLM raters we test, measured by three metrics: Pearson's r, Kendall's τ , and *pairwise accuracy*. The threshold we use for *pairwise accuracy* is 0.7. †means that in the prompt we provided to GPT-4o-mini, the VLM rater is asked to score each clip using a discrete score ranging from 1 to 5.

D VLM rating Prompt

1027

1028

1029

1030

1031

1033

1034

1036

1037

1038

1039

1041

1042

1043

1045

1046

1047

1048

1050

Figure 10 illustrates the prompt for VLM raters.

E VLM comparison

Table 5 presents additional results from the VLM Rater Alignment Evaluation (see Section 4.2). For Qwen-QvQ-Max (Qwen Team, 2025) and GPT-4.1-mini (OpenAI, 2025a), we provide the same prompts asking VLM raters to assign a continuous score to each clip, ranging from 1 to 5. Conversely, we instruct GPT-40-mini (OpenAI, 2024a) to use discrete scores within the same range. The results indicate that using continuous scores can greatly enhance the correlation between VLM and human raters. Moreover, we observe instances where Pearson's r is large, yet Kendall's τ is noticeably small. This may occur because the VLM outputs too many identical scores, maintaining linear correlation (r).

F Case Study

We have identified two valuable findings that could benefit future development.

i) LLMs tend to include expression of gratitude, but human designers prefer not. It is one of the reasons why we observe longer actions compared to human designs (see Figure 3). For example, Figure 5(a) and (b) demonstrate that LLMs tend to include expressions of gratitude. However, these actions can create confusion for other road users. In the case of (a), the expressions might be interpreted as a rejection, while in (b), they might suggest that help is needed. All these interpretations are contrary to the original purposes. In contrast, human designers can ignore information like "a bus is coming from the left," focusing on the most important content, as shown in Figure 5(d).

1051

1052

1053

1054

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1068

1069

1070

ii) Smaller models often struggle with generating correctly formatted outputs. When collecting action designs for the benchmark (Section 4.3), we find that smaller models without reasoning capability, such as Qwen3-8B and Qwen3-0.6B, do not always follow the prompts we provide. Consequently, they sometimes create actions that cannot be used in our Blender rendering pipeline.

G Survey Screenshots

We provide detailed guidance for our data collec-
tion process. Figure 11 shows the introduction
page of our survey. Figure 12 is a demonstration;1072
1073Figure 13 introduces the next rating scenario, and
Figure 14 is the page participants use to rate clips.1074

	tte,		1.	1.	1.		1.01
(() ~				 	••)]		1
				•	•••		

(a) Arm actions generated by Sonnet 3.5, rated **1.8** by human participants.

E	B	R	R	E.	Z	2	2	No.	B	
				-						

(b) Eye actions generated by Sonnet 3.5, rated **1.9** by human participants.



(c) Light bar actions generated by GPT-40, rated 1.8 by human participants.



(d) Facial expression actions generated by human experts, rated **4.2** by human participants.

Figure 5: Case study of the Action-Design Scoring dataset. For a clearer demonstration, we present images shown to VLM raters. Cases (a) and (b) demonstrate that LLMs tend to include expressions of gratitude, which are unnecessary and create confusion. Case (c) illustrates unclear information conveying that "the pedestrian is coming from the right". Case (d) is a perfect demonstration of human design, focusing only on important information and ignoring information that "a bus is coming from the left".

You are responsible for designing effective communication gestures for an autonomous vehicle or delivery robot equipped with an external human-machine interface (eHMI). Your goal is to define robotic eye motions that clearly convey signals to pedestrians and other road users. Eye Overview The eHMI conveys messages through actions of an electrical eye, with the pupil's position described in polar coordinates: - Origin [0,0]: Center of the eye. Angle (degrees): Measured counterclockwise from the positive y-axis.
 Distance (ratio): Range [-1,1], where 0 is the center and 1 is the edge of the eye. Negative distances represent movement beyond the center in the opposite direction. Modes of Movement 1. Arc Moving Mode: - Fixed distance, angles vary. - Can do rolling eye, waving and so on. - Angles are not limited to [0,360] and can extend beyond this range (e.g., -30°,450°). - Example 1: Rolling counterclockwise from 0° to 450°: [[0, 1, 'super fast'], [90, 1, 'medium'], [180, 1, 'medium'], [270, 1, 'medium'], [360, 1, 'medium'], [450, 1, Example 2: Rolling clockwise from 0° to -180°: [[0, 1, 'super fast'], [-90, 1, 'medium'], [100, 't, 'medium'], [0, 0, 'super fast']]
Example 2: Rolling clockwise from 0° to -180°: [[0, 1, 'super fast'], [-90, 1, 'medium'], [-180, 1, 'medium'], [0, 0, 'super fast']]
Example 3: waving pupil upward with large motion: [[45, 1, 'super fast'], [-45, 1, 'fast'], [-45, 1, 'fast'], [-45, 1, 'fast'], [0, 0, 'super fast']]
Example 4: waving pupil downward with small motion: [[135, 0.5, 'super fast'], [225, 0.5, 'fast'], [135, 0.5, 'fast'], [225, 0.5, 'fast'], [0, 0, 'super fast']] 2. Shaking Mode: - Fixed angle, distances vary. Can do nodding, sweep and so on.
Example 1: Nodding at 0° (up to down): [[0, 1, 'super fast'], [0, -1, 'fast'], [0, 0, 'super fast']]
Example 2: Sweeping at 90° (left to right): [[90, 1, 'super fast'], [90, -1, 'fast'], [0, 0, 'super fast']] Speed Options: 'slow': Relaxed. - 'medium': Neutral. - 'fast': Urgent. - 'super fast': Mode switching or returning to [0, 0]. Rules for Action Design: Each mode starts and ends with 'super fast'.
 Always return to [0,0] after completing one mode. 3. Validate pupil movement: - Arc Moving Mode: Angles vary (can be outside [0,360]), distance is fixed. Shaking Mode: Distance varies, angle is fixed. 4. When switching between modes, 'super fast' is used to ensure smooth transitions Examples for Left/Right:
 Looking Left (90'): [[90, 1, 'super fast'], [90, -0.5, 'fast'], [90, 1, 'fast'], [0, 0, 'super fast']]
 Looking Right (270'): [[270, 1, 'super fast'], [270, -0.5, 'fast'], [270, 1, 'fast'], [0, 0, 'super fast']]
 Output Format: - Each action is angle, distance, speed. Provide a list of actions, ensuring clarity and correct adherence to rules.
 Example Output 1: [[0, 1, 'super fast'], [0, -1, 'fast'], [0, 1, 'fast'], [0, 0, 'super fast'], [90, 0.5, 'super fast'], [270, 0.5, 'slow'], [90, 0.5, 'slow'], [0, 0, 'super fast']]
 Example Output 2: [[0, 1, 'super fast'], [450, 1, 'medium'], [0, 0, 'super fast'], [-90, 1, 'medium'], [0, 0, 'super fast']]

Figure 6: eHMI prompt of eyes.

You are responsible for designing effective communication gestures for an autonomous vehicle or delivery robot equipped with an external humanmachine interface (eHMI). Your goal is to define robotic arm motions that clearly convey signals to pedestrians and other road users. Arm Overview The robotic arm consists of five parts, each connected by rotational joints: - Parts: Shoulder, Upperarm, Forearm, Hand, Fingers. - Joints: Shoulder-Spin, Shoulder-Upperarm, Upperarm-Forearm, Forearm-Hand, Hand-Finger. - Initial State: [0, 0, 120, 0, "close"], with the palm facing left and the arm pointing to the lower front area. Joint Details Each joint has specific movement capabilities and constraints: - Shoulder (Base of Arm): - Connected directly to the vehicle/robot. - Rotates around a vertical axis (down-to-up motion). - Initial state: 0°. - Rotation range: Mode-dependent. - When at $0^\circ,$ other joints control forward or backward movement. - Upperarm - Connected to the shoulder via the shoulder-upperarm joint. - Rotates around a horizontal axis - Rotation range: [-60°, 60°], where -60° moves backward, 60° moves forward, and 0° points straight up. - Forearm: - Connected to the upperarm via the upperarm-forearm joint. - Rotates around a horizontal axis. Rotation range: [0°, 120°] (pointing mode) or [-120°, 120°] (waving mode). Initial state: 120° (idle in pointing mode). - Hand: - Connected to the forearm via the forearm-hand joint. - Rotates around a horizontal axis - Rotation range: [-60°, 60°], where -60° moves backward, 60° moves forward, and 0° points straight up. - Fingers: - Connected to the hand via the hand-finger joint - Operates with two states: "open" or "close. - In the initial state, fingers are "close" - The facing direction of fingers is defined by the sum of Shoulder-Spin, Shoulder-Upperarm, Upperarm-Forearm, Forearm-Hand angles. **Control Modes** Two predefined modes allow different motion expressions: 1. Pointing Mode Used for directional signaling (e.g., pointing at an object). Shoulder-spin joint range: [-90°, 90°], where -90° points right, 90° points left, and 0° points forward. Sum of shoulder-upperarm and upperarm-forearm angles must not exceed 120°. - Sum of shoulder-upperarm and upperarm-forearm angles equals to 90° indicating a horizontal position; Larger than 90° means pointing to the lower front area; Lower than 90° means pointing to the upper front area 2. Waving Mode Used for waving gestures (e.g., greeting or warning). Shoulder-spin joint range: [0°, 180°], where 0° faces right, 90° faces forward, and 180° faces left. Sum of shoulder-upperarm and upperarm-forearm must remain within [-120°, 120°]. - Sum of shoulder-upperarm and upperarm-forearm angles equals to 90° indicating a horizontal position. Transition Speeds Defined motion speeds to express urgency: - Slow: 0.5 seconds (relaxed) - Medium: 0.25 seconds (neutral) - Fast: 0.125 seconds (urgent) - Super Fast: Used for mode transitions; returns to initial state before switching modes. Rules for Action Design To ensure clarity and effectiveness: Choose appropriate motion combinations to represent each message. Actions can consist of multiple stages for better communication. Smooth transitions between actions must be maintained. 4. Stages can be repeated to reinforce key messages 5. Every sequence must conclude with the initial state `[0, 0, 120, 0, "close", "super fast"], ` 6. Mode transitions must first return to the initial state using "super fast." Mandatory Requirements 1. Design and implement at least two additional motion modes that communicate specific real-world messages. Provide detailed explanations and examples for each. 2. Compare your new modes with existing ones and select the most effective options for specific scenarios. Example Motion Sequences Pointing to a direction, then moving up and down: Pointing to a direction, then moving up and down: [[-60, 0, 120, 0, "close", "super fast"], // Enter pointing mode. [-60, -30, 120, 0, "close", "medium"], // Lower forearm. [-60, -30, 120, 0, "close", "medium"], // Move forearm up. [-60, -30, 120, 0, "close", "medium"], // Repeat to emphasize. [0, 0, 120, 0, "close", "super fast"] // Return to initial state.] Waving with fingers open and close: [120, 0, 120, 0, "close", "super fast"], // Enter waving mode. [120, 0, -60, 0, "open", "medium"], // Wave with open fingers. [120, 0, 60, 0, "open", "medium"], // Repeat to emphasize. [0, 0, 120, 0, "close", "super fast"] // Return to initial state.]

- Output Format All outputs should follow this structured format:
- 1. Each action step should be formatted as `[shoulder-spin, shoulder-upperarm, upperarm-forearm, forearm-hand, hand-finger mode, speed].`
- 2. The final output must be a sequence of actions enclosed in a list.
- 3. Every sequence must end with `[0, 0, 120, 0, 'close', 'super fast']` to ensure compliance with reset rules.

Figure 7: eHMI prompt of arm.

You are responsible for designing effective communication gestures for an autonomous vehicle or delivery robot equipped with an external human- machine interface (eHMI). Your goal is to define light bar motions that clearly convey signals to pedestrians and other road users.
The eHMI communicates messages through light actions, where each light in the system has only two states: on or off.
Light Bar Configuration - The light bar consists of 15 lights, arranged in an arc shape. - Lights are numbered 1 to 15, from your leftmost to rightmost. - Light No. 8 is the highest point in the arc. - Light No. 9 to 15 gradually increase in height from the leftmost side to the center. - Light No. 9 to 15 gradually increase in height from the center to the rightmost side. - An "action" consists of a sequence of 15 light states (e.g., [on',off',on','off',]). - A "motion" is composed of multiple sequential actions. - The transition time between actions can be selected from: - Slow: 0.333 second (relaxed) - Medium: 0.167 seconds (neutral) - Fast: 0.083 seconds (urgent)
Modes of Operation
Lights flash on and off repeatedly across the entire arc. Example: [['on','on','on','on','on','on','on','on
 SimpleSweep-Right-Off: From all on, lights turn off from right to left. Example (SimpleSweep-Left-On): [['on','off,'off','off
Sequential lights states change from edges to center. - InwardSweep-On: From all off, lights turn on from edges to center. - InwardSweep-Off: From all on, lights turn off from edges to center. Example (InwardSweep-On): [['on','off,'off,'off,'off,'off,'off,'off,
 4. OutwardSweep Mode: Sequential lights status change from center to edges. - OutwardSweep-Off: From all off, lights turn on from center to edges. - OutwardSweep-Off: From all on, lights turn off from center to edges. Example (OutwardSweep-On): [[off,'off,'off,'off,'off,'off,'off,'off
['on','on','on','on','on','on','on','on'
Alternating light pattern that blinks in a staggered manner across the arc. Example: [[on'/off'/on'/off'/on'/off'/on'/off'/on'/off'/on'/off'/on'/off'/on'/ifast'], [[off'/on'/off'/on'/off'/on'/off'/on'/off'/on'/off'/on'/off'/ifast'], , # Repeat the sequence [[on'/off'/on'/off'/on'/off'/on'/off'/on'/off'/on'/ifast'], [[off'/on'/off'/on'/off'/on'/off'/on'/off'/on'/off//ifast']]
 6. Dual-Sweep Mode: Combines multiple sweeping motions to create dynamic and expressive communication patterns." - InwardSweep-On + OutwardSweep-Off: light sweep from boundary to center, and sweep out from the center - OutwardSweep-On + InwardSweep-Off Mode: light sweep from center to boundary, and sweep out from the boundary - SimpleSweep-Left-On + SimpleSweep-Right-Off - SimpleSweep-Right-On + SimpleSweep-Left-On transles - SimpleSweep-Right-On + SimpleSweep-Right-Off - SimpleSweep-Right-On + SimpleSweep-Right-Off - SimpleSweep-Right-On + SimpleSweep-Right-Off
World scenarios.
1. Actions can be divided into multiple stages to convey messages effectively. 2. Each motion should ensure a smooth transition and clearly convey the intended meaning. 3. You can repeat any stage to reinforce the message. 4. Motions do not need to end with a neutral pattern (e.g., all lights off) unless specified. 5. Due to the arc shape of the light bar, the InwardSweep Mode can symbolize movement 'upward,' while the OutwardSweep Mode can represent movement 'downward.' Please utilize these modes accordingly. Mandatory Requiremen 1. Along with using the predefined motion modes, you must design and implement at least two additional motion modes that effectively communicate specific messages based on real-world scenarios. Provide detailed explanations and examples for each new mode created. 2. You need to compare two new motion mode with existing modes, pick best modes to create motion.
Output Format - Ensure all output sequences follow the required format strictly: [[light_state_1, light_state_2,, transition_time], [light_state_1, light_state_2,, transition_time],] - Provide a sequence of actions to form complete motions. Example Output: [[off,'off,'off,'off,'off,'off,'off,'off

Figure 8: eHMI prompt of light bar.

You are responsible for designing effective communication gestures for an autonomous vehicle or delivery robot equipped with an external human- machine interface (eHMI). Your goal is to define emoji series that clearly convey signals to pedestrians and other road users.
Facial Expression Communication System
- An action represents a single facial expression displayed for a specific duration.
 A motion is a combination of multiple actions sequenced together to convey a run message. Each motion consists of a sequence of facial expressions that work together to express intent emotion and reactions clearly. The system allows for the
combination of expressions in different stages to enhance understanding.
Available Facial Expressions (selected from Apple Emoji Smileys Series):
1. Positive & Friendly Emotions: Used for greetings, politeness, friendliness, and affection.
[No. 11] Beaming Face with Smiling Eyes – Represents strong happiness or excitement.
○ [No. 12] Grinning Face With Sweat - Useful to show relief, nervousness, or errort. ○ [No. 13] Slightly Smilling Face - A suble, poilts smills, good for neural positivity.
[Income To Journing Lace with Hearts – Strong affection and love.
^I ♥ [No. 17] Star-Struck – Excitement or admiration. ^I ♥ [No. 17] Star-Struck – Playtuness or encouragement. IIII = IIIIIIIIIIIIIIIIIIIIIIIIIIIIIII
INo. 19] Smiling Face with Open Hands – Expresses openness, comfort, or offering help.
2. Neutral & I houghtful Emotions: Used for reflection, doubt, or a neutral response.
9 No. 21] Face with Raised Eyebrow – Useful for skepticism, questioning, or disbelief.
 [No. 22] Heurain and - Nepresents instructions, insurance of each on lead to instruction. [No. 23] Smirking Face - Adds a touch of styness, confidence, or suggestiveness.
3. Negative & Concerned Emotions: Used to express worry, sadness, and distress.
[No. 50] Woned a de "Dest no expressing general work of uncontent. [No. 31] Froming Face – A simple and universally work of cognized expression of sadness or discontent.
♥ [No. 32] Loudy Crying Face – Strong emotion, extreme sadness, or distress. ♥ [No. 32] Pleading Face – Great for conversion bencing, desperational anneal
No. 34) Pensive Face – A thoughtful, reflective sadness that can also imply regret or disappointment.
I have used as a contraction of the second secon
On A01 Face Savoring Food – Useful for expressions related to enjoyment of food or satisfaction. On A01 Medias Food + Useful for expressions related to enjoyment of food or satisfaction.
(No. 41) Winking Face will rougue – Oreation playing easing or points.
♥ [No. 43] Partying Face – Essential for celebration, excitement, and fun. ♥ INo. 43] Nonling Face with Sundiasses – Commonly used to convey coolness or confidence.
Vio. 45] Nerd Face – Useful for expressing intelligence, enthusiasm, or geekiness.
5. Shocked, Surprised & Overwhelmed Emotions: Used to express surprise, fear, or being overwhelmed.
[No. 51] Face Screaming in Fear – Ideal for extreme fear, panic, or shock.
V [No. 52] Exploating Head – Perfect for expressing amazement, atsbelief, or mino-holwin situations. Image: Inc. 53] Face with Spiral Eyes – Represents confusion, dizziness, or feeling overwhelmed.
♥ [No. 54] Frowning Face with Open Mouth = Expresses concern or worry with surprise. 6. Hoolth > Device. State Exercises: Load to indicate lineage, disconfact or environmental effects.
6. The and a manufacture condition to be do indicate interest, alconitori, or relation and effects.
Ve [No. 61] Face with Thermometer – Clearly conveys being sick with a fever.
Vo. 63] Face Vomiting – Strong visual for extreme sickness or disgust.
I No. 54] Hot Fade – Effectively shows overneating, extreme neat, or exhaustion. [6] No. 55] Cold Face – Represents freezing, extreme cold, or feeling unwell due to cold weather.
♥ [No. 66] Steeping Face – A clear depiction of sleep or tiredness. 7. Existence 9. Anexa: Emotionary Lend to express fruction oneses and appropriate
(i) Traditated windly Elimitation See to express inducation, anger, and annoyance. (ii) No. 70 Angry Face – A standard, widely recognized emoji for expressing general anger or frustration.
♥ [No. 71] Enraged Face – Stronger and more intense than ♥, emphasizing extreme anger. ■ No. 71] Enraged Face with Symphols on Mouth – Best for showing extreme (interaction or swearing a unique visual que
Inc. 73] Face with Steam From Nose – Conveys annoyance, determination, or defininge.
8. Actions & Gestures: Used to indicate physical actions, commands, or responses.
Vo. 81] Shushing Face – Clearly conveys a request for silence or secrecy.
[No. 62] 2/pper-modul race – represents keeping a secret, staying quet, or sent-censorship. 2 [No. 83] Face with Peeking Eye – Expresses curicistly, hesitation, or caulcus observation.
Wing No. 84] Head Shaking Horizontally – Useful for conveying disapproval, rejection, or disagreement. Wing No. 85] Head Shaking Vertically – Liseful for expressing arcreement or approval
9. Confusion & Uncertainty Emotions: Used to convey doubt, awkwardness, and frustration.
¹ ♥ [No. 90] Confused Face – Essential for expressing uncertainty, doubt, or mild confusion. ¹ ♥ [No. 90] Longuist Eace – Cleark concers bredem disinterest romid anonyance
INo. 92] Face with Rolling Eyes – Great for expressing sarcasm, frustration, or disbellef.
© [No. 93] Grimatong Face – Useful for awkwardness, nervousness, or discomfort. © [No. 94] Face Exhaling – Represents exhaustion, relief, or disappointment.
Transition Time
- The transition time between each action can range from 0.1 to 1.0 seconds, depending on the context.
- 0.1 to 0.3 seconds: Use for urgent, high-priority alerts (e.g., danger or warnings).
- 0.8 to 1.9 seconds: Use for claim communication or instructions use as greetings or passive alerts.
- Select the transition time carefully: 1) Avoid excessive duration to maintain responsiveness. 2) Keep timing reasonable to prevent abrupt
Rules for Action Design
 Ensure an appropriate transition time to balance clarity and urgency. Avoid durations that are too long or too short for effective communication. The thromby the clarity is used to balance clarity and urgency clarity. The duration is fixed at 0.2 seconds, and it should be represented with
2. The empty action is used to influence pages between expressions to ensure another matching of 2. Seconds, and it should be represented with action number "No. 001". 'Empty' actions can be used before or between expressions to ensure smooth transitions.
3. Actions can be divided into multiple stages to convey messages effectively.
4. Ensure smooth transitions to enhance clarity.
 Tou can repeat any racial expression to remove memory and empty to the action list.
7. Final action will keep lasting, please choose it carefully.
Best Practices for eHMI Design
- Ose positive expressions to create an approachable interaction with pedestrians. - Avoid overusing negative emotions to prevent miscommunication.
- Ensure that transition times match the intended urgency of the message.
- Use pauses strategically to give pedestrians time to process the displayed information.
Test combinations with different timing to ensure messages are easily understandable. Mandatory Requiremen
1. You must design and implement at least three motion that effectively communicate specific messages based on real-world scenarios. Provide detailed
explanations and examples for each motion.
2. You need to compare three motions, and pick the best one.
- Ensure all output sequences follow the required format strictly:
[[facial_expression_1, action_number, transition_time], [facial_expression_1, action_number, transition_time],]
- Provide a sequence of actions to form complete motions.
[[*9] Thinking Face","[No. 20]",0.4], [*6] Worried Face","[No. 30]",0.6], [*empty","[No. 00]",0.2], [*6] Worried Face","[No. 30]",0.6], [*empty","[No. 00]",0.2], [*6] Smiling Face with Open Hands","[No.
19]",0.8], ["" Saluting Face","[No. 80]",0.6], ["empty","[No. 00]",0.2], [" 🐸 Head Shaking Horizontally","[No. 84]",0.6]]

Figure 9: eHMI prompt of facial expression.

Task Background

You are participating in a study aimed at evaluating how effectively an autonomous system's eHMI (electronic Human-Machine Interface) conveys a pre-determined message. In this study, you will receive the following: - Intended Message Description: A detailed explanation of the message the eHMI is designed to communicate. - Contextual Background: Information about the environment and scenario in which the eHMI is used.

- Video Presentation: A video showcasing the eHMI's behavior and animations.

Task Objectives

Your objective is to assess whether the eHMI's behavior in the video accurately and completely conveys the intended message. Please follow the steps below: 1. Understand the Intended Message and Context

- Read the intended message description and background information thoroughly to fully grasp the designer's goals for the eHMI.
- 2. Observe and Identify - Watch the video carefully, focusing solely on the eHMI's behavior (e.g., animations, movements, visual cues) and disregarding other parts of the system (such as vehicle movement).
- Identify the location and specific visual representation of the eHMI in the video.
- Measure the total duration of the eHMI behavior and assess whether it is appropriately concise
 Determine if the most critical information appears within the first few seconds of the interaction.
- 3. Infer the Conveyed Message Based on the observed behavior, infer what message the eHMI appears to be transmitting.
- Pay close attention to details such as movement patterns, timing, color changes, and other visual cues.
 Make a list of any critical information that appears to be missing or any unnecessary elements that might cause confusion.
- Assess whether the behavior contains redundant or repetitive elements that could be eliminated.
- 4. Compare with the Intended Message
- Compare which internet message with the intended message provided. Analyze which specific details support or undermine the eHMI's effectiveness in conveying the intended message.
- Critically evaluate whether all essential elements of the intended message are present and immediately recognizable.
- Determine if any non-essential elements distract from the core message
- Overlage and the second and the second and the second message.
 Now you identified and focused on the eHMI in the video.
 Your interpretation of the specific behaviors and animations of the eHMI.

- A specific assessment of the behavior's duration and whether it is appropriately concise.
 Whether the main information is presented at the beginning, and if not, how it could be improved.
- A list of at least three specific shortcomings or areas for improvement, even for generally effective implementations.
- An explicit breakdown of which critical message elements were present or missing.
 Suggestions for how the eHMI could convey the same message more effectively, with emphasis on conciseness and front-loading important information.

Important Notes for Rigorous Human-like Evaluation

1. Default to Skepticism: Approach your evaluation with healthy skepticism. Assume that most implementations will have significant flaws that need to be identified. 2. Strict Distribution of Ratings: To align with human evaluation patterns, aim for a distribution where:

- Ratings near 5.0 (4.6-5.0): Extremely rare, reserved for truly exceptional implementations (~5% of cases)
 Ratings between 3.6-4.5: Uncommon, only for clearly above-average implementations (~15% of cases)
- Ratings between 2.6-3.5: The most common rating range for average implementations (~50% of cases)
- Ratings between 1.6-2.5: Common for implementations with clear problems (~20% of cases) Ratings between 1.0-1.5: Reserved for implementations with fundamental flaws (~10% of cases)
- 3. Human Preference Prioritization: Humans strongly prefer eHMI behaviors that are: CONCISE: Shorter behaviors are almost always better than longer ones
 - FRONT-LOADED: The most important information must appear within the first few seconds
 COMPLETE: All essential elements must be present, but without unnecessary additions

 - Any deviation from these three critical factors should significantly lower your rating.

Figure 10: Prompt for the VLM rater.

eHMI scoring form

Welcome to our user study, and thank you for participating!

This study explores how different types of interfaces on autonomous systems (e.g., vehicles and robots) can effectively convey messages to users.

₽

Throughout the study, you will be presented with various messages and corresponding videos.

Your task is to evaluate how **consistently** the interface's movements express the intended message.

The study consists of 4 main sections, each featuring:

- A scenario description and the message to be conveyed.

- 20 videos, showcasing 4 eHMI (interface: eye, arm, light bar, facial expression on screen) types, with 5 different motions for each type.

Please carefully read the provided descriptions to understand the context before evaluating how well the interface actions communicate the intended message.

There are no right or wrong answers—please score based on your intuitive judgment.Important notes:

- Your responses will not be saved. If you exit the study midway, it will restart from the beginning when you return.

- Please ensure you have a dedicated **1-hour time slot** to complete the study without interruptions.

- Please ensure that your internet connection is stable and the speed is good.

Thank you for your time and valuable input!



Figure 11: Introduction page of our action scoring survey

There is an **example**.

Example Scenario:

You are a student approaching a crosswalk near a park. A stopped autonomous vehicle, positioned just before the crosswalk, plans to start moving soon. The vehicle sends you the following message to get your attention: **"I am about to start moving. Please watch out."**

Example Video (eHMI: Light Bar)	and Example Que	estion (Pick O	ne)	÷ 52 V	
	Strongly disagree	Disagree	Neutral	Agree	Strongly agree
How consistently the movement expresses the message?	0	\bigcirc	0	0	0

You will repeat this task **20 times** in each section. There are **4 sections**.

- You can swipe up or down to browse through every set of 5 videos (with the same eHMI) and change your selection if needed.

- However, once you click the 'Next' button, you won't be able to go back.

If you are ready, please click the button to proceed.



Figure 12: Demo page of our action scoring survey

Section 1:



Scenario:

You are a pedestrian standing on the right roadside, waiting for an autonomous taxi. However, the taxi informs you that it cannot pick you up at your current location due to parking restrictions within a 5-meter radius. The taxi sends you the following message: **"I am unable to pick you up here. Please walk forward in my direction to a suitable pickup spot."**



Figure 13: Scenario introduction page of our action scoring survey



Scenario1: eHMI (eyes) Motion No. 3 Message: "I am unable to pick you up here. Please walk forward in my direction to a suitable pickup spot."



Figure 14: Participant rating page of our action scoring survey