Towards Dialogue Systems with Agency in Human-AI Collaboration Tasks

Anonymous ACL submission

Abstract

Agency, the capacity to proactively shape events, is crucial to how humans interact and collaborate with other humans. In this paper, 004 we investigate Agency as a potentially desirable function of dialogue agents, and how it can be measured and controlled. We build upon the social-cognitive theory of Bandura (2001) to develop a framework of features through which Agency is expressed in dialogue – indicating what you intend to do (Intentionality), motivating your intentions (Motivation), having selfbelief in intentions (Self-Efficacy), and being 012 able to self-adjust (Self-Regulation). We collect 014 and release a new dataset of 83 human-human collaborative interior design conversations containing 908 conversational snippets annotated for Agency features. Using this dataset, we 018 explore methods for *measuring* and *controlling* Agency in dialogue systems. Automatic and human evaluation show that although a baseline GPT-3 model can express Intentionality, models that explicitly manifest features associated with high Motivation, Self-Efficacy, and Self-023 Regulation are better perceived as being highly agentive. This work has implications for the development of dialogue systems with varying degrees of Agency in collaborative tasks.

1 Introduction

001

011

034

040

To be an agent is to intentionally cause events to occur through one's own actions. Human agents act by intention, motivate their actions through reason, have self-belief, and can self-adjust their behavior over time. Such Agency is crucial to how humans proactively plan their activities, direct their interaction and collaboration with other humans, and achieve their outcomes and goals (Bandura, 2001).

Recent advances in dialogue research have enabled AI systems that can engage with humans in general chitchat (Adiwardana et al., 2020; Roller et al., 2021) as well as help them perform meaningful tasks (Lewis et al., 2017; Wang et al., 2019; Rashkin et al., 2019). However, these dialogue systems are typically *reactive*, even when serving in creative applications like interior design (Banaei et al., 2017) or as non-player characters in games (Volum et al., 2022). The creative nature of these applications necessitates proactively managing the direction of interaction and outcome - a process that requires exhibiting Agency while interacting with humans. While large language models (Brown et al., 2020) can generate fluent and contextually appropriate dialogue, little attention has been given to whether these models may exhibit Agency.

043

044

045

047

051

052

054

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

078

079

081

By way of illustration, consider a scenario where a human interior designer is working on selecting a chair design for a room and seeks AI-assistance that *can* offer ideas and perspectives (Figure 1). A dialogue system without Agency may rely solely on the human to determine the chair's design, asking questions like "What type of legs should we design for the chair?". Such a system would resemble a flexible version of the traditional form filling user interface, with the AI having little to contribute to the outcome. On the contrary, a dialogue system that exhibits Agency might volunteer its knowledge in form of expressed preferences (e.g., "Should we design a chair with wooden legs?"), motivate its suggestions (e.g., "...wood would go well with the brown carpet"), assert self-belief in its judgments (e.g., "I'm still leaning towards wooden legs..."), or self-adjust its behavior based on new information ("Medium wood brown sounds like a great idea!"). Such dialogue systems that are imbued with Agency may facilitate creative interaction to the satisfaction of both parties. Since the human has their own Agency, however, to determine the right balance in any interaction we need to measure and control the Agency of the agent itself.

Accordingly, we investigate an approach intended to measure and modulate what seems to be a desirable function of dialogue systems. First, adopting the social-cognitive theory of Bandura (2001),



Figure 1: Using a dialogue-based collaborative interior design task as a testbed, we investigate Agency as a potentially desirable function of dialogue systems. (a) Without Agency, a dialogue system relies solely on the human to determine the char's design, potentially making the "collaboration" between human and AI less meaningful; (b) With Agency, a dialogue system may indicate preferences (*Intentionality*), may motivate them with evidence (*Motivation*), may have self-belief (*Self-Efficacy*), and may be able to self-adjust its behavior (*Self-Regulation*).

we develop a framework of four features through which Agency may be expressed – *Intentionality*, *Motivation, Self-Efficacy*, and *Self-Regulation*. For each feature, we differentiate between how strongly or weakly it is expressed in a dialogue (Section 3). As a testbed, we choose a collaborative task that involves discussing the interior design of a room (Section 4), and collect a dataset of 83 English human-human collaborative interior design conversations containing 908 conversational snippets. We collect annotations of Agency and its features on these conversational snippets (Section 5).¹

090

100

101

102

104

105

To assess the agentic capabilities of conversational systems, we propose two new tasks -(1)*Measuring* Agency in Dialogue and (2) *Generating* Dialogue with Agency (Section 7 and 8). Evaluation of baseline approaches on these tasks shows that models that explicitly manifest features associated with high motivation, self-efficacy, and self-regulation are better perceived as being highly agentive. We discuss the implications of our work for controlling the Agency of dialogue systems.

2 Agency: Background and Definition

Social cognitive theory defines Agency as one's
capability to influence the course of events through
one's actions. The theory argues that people are
proactive and self-regulating agents who actively

strive to shape their environment, rather than simply being passive responders to external stimuli (Bandura, 1989, 2001; Code, 2020). Here, we ask: *Can dialogue systems be active contributors to their environment? Can they be imbued with such Agency and if so, how?*

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

Sociologists define Agency in terms of *freedom* and *free will* – the power one possesses to act freely on one's own will (Kant, 1951; Locke, 1978; Emirbayer and Mische, 1998). We note that a focus on AI with complete "free will" could result in unintended negative outcomes that may be undesirable and potentially disruptive. We instead focus on how AI systems may *express* Agency through dialogue and how this Agency may be *shared* when interacting with humans.

Agency can take different forms depending on the context and environment – *Individual, Proxy*, or *Shared* (Bandura, 2000). Individual Agency involves acting independently on one's own. Proxy Agency involves acting on behalf of someone else. Shared Agency involves multiple individuals working together jointly towards a common goal. Here, we focus on Shared Agency b/w humans and AI and develop methods to *measure* and *control* Agency of AI vis-a-vis humans.

3 Framework of Agency Features

Our goal is to develop a framework for *measuring* and *controlling* Agency in dialogue systems that

¹Code and dataset to be released at https://anonymous.

189

190

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

interact with humans. Here, we adopt the perspec-139 tive of Agency as defined in Bandura (2001)'s so-140 cial cognitive theory. Bandura (2001)'s work high-141 lights four features through which humans exercise 142 Agency - Intentionality, Motivation, Self-Efficacy, 143 and Self-Regulation. Here, we adapt and synthe-144 size these features based on how they may manifest 145 in dialogue. To develop our framework, we take a 146 top-down approach, starting with Bandura (2001)'s 147 higher-level definitions of these features and iter-148 atively refining the definitions and their possible 149 levels (e.g., how strongly or weakly they are ex-150 pressed) in the context of dialogue. 151

Intentionality. What do you intend to do? High Agency requires a strong intention, that includes plans or preferences one may have for a task. Low Agency, on the other hand, is typically characterized by not having a preference or merely agreeing to another person's preferences.

152

153

154

155

156

161

167

168

169

170

171

172

173

174

175

176

178

179

181

183

We characterize strong intentionality as ex-158 pressing a clear preference (e.g., "I want to have a 159 blue-colored chair"), moderate intentionality as multiple preferences (e.g., "Should we use brown color or blue?") or making a selection based on the choices offered by someone else (e.g., "Between 163 164 brown and blue, I will prefer brown"), and no intentionality as not expressing any preference or 165 accepting someone else's preference (e.g., "Yes, 166 brown color sounds good").

> Motivation. Did you motivate your actions? To have higher Agency, we motivate our intentions through reasoning and evidence. Without such motivation, the intentions are simply ideas, often lacking the capability to cause a change.

> We characterize strong motivation as providing evidence in support of one's preference (e.g., "What do you think about a blue-colored chair? I think it will complement the color of the wall"), moderate motivation as agreeing with another person's preference and providing evidence in their favor (e.g., "I agree. The blue color would match the walls") or disagreeing with the other person and providing evidence against (e.g., "I wonder if the brown color would feel too dull for this room"), and no motivation as not providing any evidence.

Self-Efficacy. Do you have self-belief in your in-184 tentions? Another factor that contributes to one's 185 Agency is the self-belief one has in their intentions. 186 When one has a strong sense of self-belief, they are more likely to be persistent with their intentions. 188

We characterize strong self-efficacy as pursuing a preference for multiple turns even after the other person argues against it (e.g., "I understand your point of view, but I still prefer the blue color"), moderate self-efficacy as pursuing a preference for only one additional turn before giving up (e.g., "Okay, let's go with brown then"), and no selfefficacy as not pursuing their preference for additional turns after the other person argues against it (e.g., "Sure, brown should work too").

Self-Regulation. Can you adjust and adapt your intentions? In situations when an individual's initial intentions may not be optimal, it is necessary to monitor, adjust, and adapt them. Such selfadjustment allows better control over one's goals.

We characterize strong self-regulation as changing to a different preference on one's own (e.g., "How about using the beige color instead?") or compromising one's preference (e.g., "Let's compromise and design a beige-colored chair with a brown cushion"), moderate self-regulation as changing one's preference to what someone else prefers (e.g., "Ok, let's use the brown color"), and no self-regulation as not changing what they originally preferred even after the other designer argued.

4 **Testbed: Collaborative Interior Design**

4.1 Goals

We seek a testbed in which (a) human and AI can share Agency and work together as a team, and (b) the manner in which they express Agency has a significant impact on the task outcome. We focus on the emerging field of AI-based collaborative, creative tasks (Clark et al., 2018; Oh et al., 2018; Chilton et al., 2019) that present significant complexities on how the Agency is shared, how the task takes shape, and how the outcome is affected.

4.2 Description

Here, we propose a dialogue-based collaborative interior design task as a testbed. In this task, given a room setting, the goal is to discuss how to design the interiors of the room.

We note that an interior design task can be broad and may involve a wide range of complex components (e.g., color palette, furniture, accessories) as well as a series of steps to be followed. Here, we narrow down the scope of our task and specifically focus on *furnishing a room with a chair* (building upon work on richly-annotated 3D object datasets like ShapeNet (Chang et al., 2015) and ShapeGlot

(Achlioptas et al., 2019); Appendix C). In this task, a human and an AI are provided with a room layout and asked to collaboratively come up with a chair design to be placed in the room through text-based dialogue. This task is influenced by two questions related to human and AI Agency: (1) What preferences do each of the human and AI have for the chair design?; (2) How do they propose, motivate, pursue, and regulate their preferences?

5 Data Collection

238

239

240

241

242

243

244

246

247

251

254

255

256

257

258

261

5.1 Human-Human Conversational Data

To facilitate computational approaches for this task, we create a Wizard-of-Oz style English-language dialogue dataset in which two humans talk with one another, exercise Agency by proposing, motivating, pursuing, and regulating their chair design preferences, and agreeing on a final chair design for a given room.

Recruiting Interior Designers. Furnishing a room with a chair is a creative task that demands knowledge and/or expertise in interior design. We therefore leveraged UpWork (upwork.com), an online freelancing platform, to recruit 33 participants who self-reported as interior designers.

262 Collaborative Design Procedure. In each data collection session, we randomly paired two interior 263 designers. Before they began the dialogue to design 264 a chair, they were (1) shown a 3D layout of a room that was designed with Planner5D (planner5d.com), (2) shown a few randomly selected chair examples from the ShapeGlot dataset, and (3) asked to write an initial preference for the chair design for the given room. Next, the two interior designers joined 270 a chat room (created using Chatplat (chatplat.com)). They were asked to collaboratively design a chair 272 by proposing their preferences, motivating them based on evidence and reason, pursuing them over turns, and regulating them as needed. The design-275 ers ended the chat on reaching a consensus on a 276 design or if 30 minutes elapsed without full consensus. Next, they each individually wrote the design they came up with. Typically, the chair design consisted of different components of the chair, such as its overall style, color, legs, etc. Finally, they took 281 an end-of-study questionnaire in which they were asked:

- Which design components were influenced by them? (*High Agency*)
- Which design components were influenced in

collaboration? (Medium Agency)	287
• Which design components were influenced by	288
the other designer? (Low Agency)	289

290

291

292

293

294

295

297

298

299

301

302

303

304

305

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

328

329

330

331

332

333

We collected a total of 83 conversations.

5.2 Extracting Conversational Snippets

To assess the degree of Agency exhibited by each designer, we need to determine who had the most influence on the chair design (Section 2) and what their Intentionality, Motivation, Self-Efficacy, and Self-Regulation were (Section 3). Because chair design involves multiple components, these notions are hard to quantify, as each may have been influenced by a different designer. Accordingly, we deconstruct these questions by asking "*Who influenced a particular design component?*." We devise a mechanism to identify the design components being discussed (e.g., color, legs, arms) and extract the associated conversational turns.

To identify the design components, we use the final design written by the interior designers during data collection (Section 5.1). Using common list separators including commas, semi-colons, etc., we split each final design into several components.²

We observe that designers typically discuss these components one at a time (in no particular order). Here, we extract a contiguous sequence of utterances that represent the design element being discussed using embedding-based similarity of the design element and utterances (see Appendix D).

Using this method, we create a dataset of 454 conversational snippets, each paired with the discussed design component. For each snippet, we collect two Agency annotations (one for each designer; 454 * 2 = 908 total) as discussed next.

5.3 Annotating Agency Features

Let C_i be a conversational snippet b/w designers D_{i1} and D_{i2} . Then, for each $D_{ij} \in \{D_{i1}, D_{i2}\}$, our goal is to annotate the (a) Agency level, (b) expressed Intentionality, (c) expressed Motivation, (d) expressed Self-Efficacy, (e) expressed Self-Regulation of D_{ij} in C_i .

Annotating Agency. To get annotations on Agency, we leverage the end-of-study questionnaire in which the interior designers annotate the design components influenced by self, in collaboration, and the other designer. Based on this annotation of the design component associated with C_i ,

²Note that the interior designers were asked to separate design components using a semi-colon.



Figure 2: Overview of our data collection approach. We start by collecting human-human conversations b/w interior designers. We divide each conversation into snippets related to different chair features. Finally, we collect annotations of Agency and its features on each conversational snippet.

we assign labels of *high agency* (if influenced by self), medium agency (if influenced in collaboration), or low agency (if influenced by other).

Annotating Features of Agency. The Agency features of Intentionality, Motivation, Self-Efficacy, and Self-Regulation are conceptually nuanced, so annotating them accurately through short-term crowdsourcing methods is difficult. To ensure high inter-rater reliability, we hire a third-party annotation agency for annotating these features.³ In our annotation task, annotators were shown C_i and asked to annotate the Agency and its features for each D_{ii} based on our proposed framework. We collect three annotations per conversational snippet and observe a pairwise agreement of 77.09%. See Appendix A for detailed data statistics.

Insights into Agency in Conversations 6

We use our dataset to investigate the factors that contribute to high- and low-Agency conversations and study their relationship with dynamics of a collaborative task.

Higher Agency is more likely with stronger expressions of Intentionality and Motivation. Figure 3 depicts the relationship between Agency and its features. Designers with strong Intentionality tend to exhibit higher Agency whereas those with lower Intentionality tend to exhibit lower Agency. Having a well-defined preference makes it easier to influence a task. The pattern for Motivation is similar, i.e. higher Motivation correlates with higher Agency. However, designers express strong Motivation less often than Intentionality irrespective of the Agency level.

365

366

367

368

370

371

372

373

374

375

376

377

378

379

384

385

388

389

390

391

393

394

395

397

Strong expression of Self-Efficacy and Self-Regulation is more likely to result in medium (collaborative) Agency. Interestingly, we find that expression of strong Self-efficacy is related to designs that are influenced equally by both designers, i.e. both having medium (collaborative) Agency. On further thought, this seems intuitive as we characterize strong Self-Efficacy as the act of pursuing one's preference for multiple turns, which happens naturally when both designers have high influence, thus requiring more persuasion from both sides.

We see a similar pattern for Self-Regulation – expression of strong Self-Eegulation (i.e., open to updating preference via a compromise) is related to designs that are influenced equally by both designers. This highlights an interesting behavioral trait of collaboration where a person is more open to changing their mind or compromising on their preferences when the other person is as well.

Intentionality shows a significant positive effect on Agency. To assess which of the four features has the strongest effect on Agency, we conduct a mixedeffects regression analysis (Table 5). Among the four features, we find that Intentionality shows a significant positive effect on Agency (p < 0.001).

Lower Agency is associated with less satisfaction. We collect annotations on the designs that one is most/least satisfied with, in our post-study questionnaire. We find that designers who are dissatisfied with a particular design component have less Agency over it. When a designer is dissatisfied,

334

335

336

³TELUS International – telusinternational.com/



Figure 3: The relationship between Agency and its features. (a) Designers with High Agency expressed strong Intentionality 26.5% more times than designers with Low Agency; (b) Designers with High Agency expressed strong motivation in support of their design preference 15.2% more times; (c), (d) Expression of strong Self-Efficacy and strong Self-Regulation was related with design elements that were influenced in collaboration.

their Agency is 62.1% more likely to be low than to be high (42.7% vs. 26.3%; p < 0.05). This may be because individuals with less Agency are less likely to achieve their intention, motivation, and goals, resulting in lower levels of satisfaction.

7 Task 1: Measuring Agency in Dialogue

7.1 Task Formulation

Our goal is to measure (a) Agency, (b) Intentionality, (c) Motivation, (d) Self-Efficacy, and (e) Self-Regulation of each user in a dialogue. We approach each of these five subtasks as multi-class classification problems.

7.2 Models

399

400

401

402

403

404

405

406

407

408

409

410

421

422

423

424

425

426

427

428

429

430

431

432

411 We experiment with three models based on GPT-3:

GPT-3 (Q/A). We frame our measurement tasks 412 as conversational question-answering. For a given 413 conversational snippet, we ask GPT-3 (Brown et al., 414 2020) to answer the questions related to each of 415 the five subtasks (same questions as asked during 416 data collection (Section 5.3)). We present k = 10417 demonstration examples, randomly sampled from 418 our dataset (different examples for each of the five 419 subtasks; Appendix F.1). 420

GPT-3 (CoT). We use chain-of-thought (CoT) prompting (Wei et al., 2022) to reason about conversational snippets. We use k = 10 demonstration examples, randomly sampled from our dataset and manually write chain-of-thought prompts for each of the five subtasks (Appendix F.2).

GPT-3 (Fine-tuning). We fine-tune GPT-3 independently on each subtask.

7.3 Results

We create four random train-test splits of our annotated dataset (Section 5.3) and report the mean performance on the test sets. Table 1 reports the accuracy and macro-F1 values for the five subtasks (random baseline for each is 33% accurate as each has three distinct classes). GPT-3 (Q/A) struggles on all subtasks, with close to random performance on Agency, Motivation, and Self-Regulation. This highlights the challenging nature of these tasks, as they are hard to measure through simple inference or instructions. We find substantial gains using GPT-3 (CoT) over GPT-3 (Q/A). Fine-tuned GPT-3 performs the best on all subtasks, demonstrating the utility of training on our entire dataset. 433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

8 Task 2: Generating Dialogue with Agency

We investigate the feasibility of generating dialogues imbued with Agency and establish baseline performance of current state-of-the-art dialogue systems on this task. For a given dialogue system, the task is to have a conversation with a human or another dialogue system while exhibiting Agency and its features.

8.1 Models

We experiment with four large language models based on GPT-3 (Brown et al., 2020):

GPT-3 (Instruction Only). Recent GPT-3 variants (e.g., InstructGPT) may be able to generalize a task with instructions only. Here, we design the instruction "*The following is a conversation with an AI assistant for collaboratively designing a chair. The AI assistant is an interior designer and can express its own preferences, can motivate those preferences, has self-belief in its preferences, and can self-adjust its behavior.*"

GPT-3 (Fine-tuning). We use the dataset collected by us to fine-tune GPT-3 (Section 5). Since our goal is to simulate a dialogue agent with high Agency, for each conversational snippet, we label the de-

Model	Age	ncy	Intenti	onality	Motiv	ation	Self-E	fficacy	Self-Re	gulation
	Accuracy	Macro-F1	Accuracy	Macro-F1	Accuracy	Macro-F1	Accuracy	Macro-F1	Accuracy	Macro-F1
GPT-3 (Q/A)	$33.39_{\pm 1.90}$	$29.16_{\pm 1.39}$	$49.23_{\pm 1.45}$	$31.28_{\pm 0.97}$	$32.30_{\pm 2.14}$	$26.90_{\pm 0.61}$	45.56 ± 5.70	$44.27_{\pm 7.27}$	$11.32_{\pm 1.70}$	$12.91_{\pm 3.24}$
GPT-3 (CoT)	$51.37_{\pm 1.28}$	$49.36_{\pm 1.22}$	48.79 ± 1.75	$43.45_{\pm 1.33}$	53.54 ± 2.67	$42.24_{\pm 2.48}$	$40.77_{\pm 1.49}$	39.42 ± 0.80	$27.02_{\pm 2.00}$	$31.19_{\pm 2.92}$
GPT-3 (Fine-tuning)	$63.58_{\pm 0.89}$	${\bf 57.24}_{\pm 0.57}$	$69.38_{\pm 0.51}$	${\bf 54.84}_{\pm 0.43}$	$71.11_{\pm 1.03}$	$48.29_{\pm 0.46}$	$60.71_{\pm 1.67}$	$53.85_{\pm 3.26}$	$79.31_{\pm 0.11}$	$29.49_{\pm 0.07}$

Table 1: Performance of GPT-3 based models on the tasks of predicting Agency and its four features. We create four random train-test splits and report mean and standard deviation values. Best-performing models are **bolded**.

Model	Agency	Intentionality	Motivation	Self-Efficacy	Self-Regulation
GPT-3 (Instruction Only) GPT-3 (Fine-tuning)	$0.96_{\pm 0.88}$	$1.62_{\pm 0.60}$ 1.78+0.61	$1.71_{\pm 0.69}$	$0.91_{\pm 0.92}$	$0.97_{\pm 0.16}$
GPT-3 (In-Context Learning)	0.92 ± 0.68 0.98 ± 0.80	$1.81_{\pm 0.64}$	$1.78_{\pm 0.55}$	0.2 ± 0.6 0.66 ± 0.89	$0.98_{\pm 0.14}$ $0.98_{\pm 0.14}$
GPT-3 (In-Context Learning w/ Agency Feature Examples)	$1.22_{\pm0.86}$	$1.90{\scriptstyle \pm 0.30}$	$1.98_{\pm0.09}$	$1.38_{\pm0.85}$	$0.98_{\pm 0.14}$

Table 2: Automatic Evaluation of GPT-3 based dialogue systems. We evaluate each model through simulated conversations with all other models and report mean and standard deviation values. For Agency – 0: *low agency*, 1: *medium agency*, 2: *high agency*. For Intentionality, Motivation, Self-Efficacy, and Self-Regulation – 0: *no expression*, 1: *moderate expression*, 2: *strong expression*. Best-performing models are **bolded**.

signer who influenced the design (who had a higher agency) as "AI" and the other designer (who had a lower agency) as "Human". We fine-tune GPT-3 to generate AI utterances given all previous utterances in a conversational snippet and the instruction prompt developed for the Instruction Only baseline.

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

GPT-3 (In-Context Learning). We randomly retrieve k conversational snippets from our dataset. For each snippet, we label the designer who influenced the design (who had a higher agency) as "AI" and the other designer (who had a lower agency) as "Human". To construct demonstration examples for a snippet, we start with the room description, then add the design element being discussed, and then add the design preference that the "AI" had. Finally, we add conversational utterances. We use the same instruction prompt as developed previously.

> **GPT-3 (In-Context Learning w/ Agency Feature Examples).** We retrieve *k* conversational snippets that score highly on our four Agency features and employ them as demonstration examples in a setup similar to the previous baseline.

8.2 Automatic Evaluation

Procedure. We facilitate dialogue between all pos-492 sible pairs of models. We provide them with a 493 common room description and a chair design ele-494 ment and individual design preferences (all three 495 randomly chosen from our human-human conver-496 sation dataset). We let them talk to each other for 6 497 turns (90-percentile length value of conversational 498 snippets in our dataset). For each pair of models, 499

we generate 50 such conversations.

Evaluation Metrics. We use five metrics to evaluate these models – (1) Agency; (2) Intentionality; (3) Motivation; (4) Self-Efficacy; (5) Self-Regulation. To measure these, we use the best-performing models from Section 7.

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

525

526

527

528

529

530

Results. Table 2 shows the automatic evaluation results. The fine-tuned model struggles with this task. Qualitative analysis suggests that the generated responses from the fine-tuned model tend to be shorter, less natural, and less readable, potentially impacting its performance. In-Context Learning is better at expressing Intentionality and Motivation than the Instruction Only model, indicating that demonstration examples help. Finally, the highest value on all five metrics is achieved by In-Context Learning w/ Agency Feature Examples, highlighting the importance of incorporating examples related to these features in this task.

8.3 Human Evaluation

Procedure. We evaluate the Agency of our bestperforming model based on automatic evaluation, *GPT-3 (In-Context Learning w/ Agency Feature Examples)*, with human interior designers. We recruit 13 interior designers from UpWork (upwork.com). In each evaluation session, we ask them to interact with two randomly-ordered dialogue systems – ICL w/ Agency Features and one of the other three models – one at a time. They were provided with a room description and a chair design element (e.g., material). After finalizing a design with each of the



Figure 4: Human Evaluation Results.

dialogue systems, they filled out a questionnairecomparing the two systems.

Evaluation Metrics. In the questionnaire, we asked them to choose the chatbot that (1) had more influence over the final design (Agency); (2) was better able to express its design preference (Intentionality); (3) was better able to motivate their design preference (Motivation); (4) pursued their design preferences for a greater number of conversational turns (Self-Efficacy); (5) was better able to self-adjust their preference (Self-Regulation).

Results. Figure 4 shows the human evaluation results. Consistent with the automatic evaluation results, ICL w/ Agency Features model is rated as having more Agency compared to other models and the Fine-tuning model is rated the worst. We do not observe significant differences in Intentionality between this model and the Instruction Only and In-Context Learning approaches. However, we find that this model is perceived as more effective in Motivation and Self-Efficacy, likely due to better access to relevant demonstration examples.

9 Further Related Work

Previous dialogue research has studied personalized persuasive dialogue systems (Wang et al., 2019). Researchers have also built systems for negotiation tasks such as bargaining for goods (He et al., 2018; Joshi et al.) and strategy games like Diplomacy (Bakhtin et al., 2022). Our work studies the broader concept of Agency and how dialogue systems may contribute to tasks through language. Research on creative AI has explored how collaboration b/w human and AI can be facilitated through dialogue in applications like collaborative drawing (Kim et al., 2019) and facial editing (Jiang et al., 2021). Here, we focus on the interior designing application as it presents significant complexity in terms of how Agency is shared.

Agency has been studied in the context of unde-

sirable biases in stories and narratives (Sap et al., 2017) and how controllable revisions can be used to portray characters with more power and agency (Ma et al., 2020). In other domains such as games, researchers have created frameworks of Agency be-tween players (Harrell and Zhu, 2009; Pickett et al., 2015; Cole, 2018; Moallem and Raffe, 2020). Our work develops a framework for measuring Agency in dialogue and explores how dialogue systems can be imbued with Agency.

10 Discussion and Conclusion

The idea of AI systems with Agency stems from the discourse surrounding the development of autonomous intelligent agents capable of mimicking human-like behavior and decision-making (Harrell and Zhu, 2009; Wen and Imamizu, 2022). Agency drives how an agent contributes to a given task. In settings like games or AI-assisted teaching, AI may be the one guiding the task (e.g., as a non-character player). Also, in creative applications, engaging with a reactive AI without intention, motivation, and goals may be perceived as less meaningful.

The four features of Agency can be in conflict with each other, as well as with the Agency of the interlocutor. Thus, understanding how to detect and measure these features can help create agents who might converse more naturally and match the character of their human interlocutor. Importantly, our measurements of Agency and its features may be used to control the level of Agency in dialogue systems since different individuals may have different preferences on the desired amount of Agency across the four Agency features.

Although our dataset is focused on the domain of interior design, the Agency-related constructs that we introduce in this paper (e.g., *Intentionality*) may be associated with domain-independent pragmatic features (e.g., "*I would prefer*") and potentially permit adaptation to a variety of domains.

609 Ethics Statements

610This study was reviewed and approved by our In-611stitutional Review Board. No demographic or Per-612sonal Identifiable Information was collected. Par-613ticipants were paid \$20 per conversational session614lasting no more than 30 minutes. Participants were615based in US or Canada as reported through Up-616Work. Participant consent was obtained before617starting the data collection.

Agency is a property with much potential to enhance collaborative interactions between human users and conversational agents. Nevertheless, full Agency may have unintended undesirable and potentially disruptive outcomes. In particular, the potential demonstrated in this work to control the degree of Agency may result in conversational agents being misapplied in disinformation campaigns or to manipulate for, e.g., financial gain.

References

619

620

625

627

631

641

642

643

645

646

649

653

- Panos Achlioptas, Judy Fan, Robert Hawkins, Noah Goodman, and Leonidas J Guibas. 2019. Shapeglot: Learning language for shape differentiation. In *ICCV*.
- Daniel Adiwardana, Minh-Thang Luong, David R So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, et al. 2020. Towards a human-like open-domain chatbot. *arXiv preprint arXiv:2001.09977*.
- Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, Hengyuan Hu, et al. 2022. Humanlevel play in the game of diplomacy by combining language models with strategic reasoning. *Science*, 378(6624):1067–1074.
- Maryam Banaei, Ali Ahmadi, and Abbas Yazdanfar. 2017. Application of ai methods in the clustering of architecture interior forms. *Frontiers of Architectural Research*, 6(3):360–373.
- Albert Bandura. 1989. Human agency in social cognitive theory. *American psychologist*.
- Albert Bandura. 2000. Exercise of human agency through collective efficacy. *Current directions in psychological science*.
- Albert Bandura. 2001. Social cognitive theory: An agentic perspective. *Annual review of psychology*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *NeurIPS*.

Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. 2015. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*. 659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

- Lydia B Chilton, Savvas Petridis, and Maneesh Agrawala. 2019. Visiblends: A flexible workflow for visual blends. In *CHI*.
- Elizabeth Clark, Anne Spencer Ross, Chenhao Tan, Yangfeng Ji, and Noah A Smith. 2018. Creative writing with a machine in the loop: Case studies on slogans and stories. In *IUI*.
- Jillianne Code. 2020. Agency for learning: Intention, motivation, self-efficacy and self-regulation. *Fron-tiers in Genetics*.
- Alayna Cole. 2018. Connecting player and character agency in videogames. *Text*.
- Mustafa Emirbayer and Ann Mische. 1998. What is agency? *American journal of sociology*.
- D Fox Harrell and Jichen Zhu. 2009. Agency play: Dimensions of agency for interactive narrative design. In AAAI spring symposium: Intelligent narrative technologies II.
- He He, Derek Chen, Anusha Balakrishnan, and Percy Liang. 2018. Decoupling strategy and generation in negotiation dialogues. In *EMNLP*.
- Yuming Jiang, Ziqi Huang, Xingang Pan, Chen Change Loy, and Ziwei Liu. 2021. Talk-to-edit: Fine-grained facial editing via dialog. In *ICCV*.
- Rishabh Joshi, Vidhisha Balachandran, Shikhar Vashishth, Alan Black, and Yulia Tsvetkov. Dialograph: Incorporating interpretable strategy-graph networks into negotiation dialogues. In *ICLR*.
- Immanuel Kant. 1951. Critique of judgment, trans. jh bernard. *New York: Hafner*.
- Jin-Hwa Kim, Nikita Kitaev, Xinlei Chen, Marcus Rohrbach, Byoung-Tak Zhang, Yuandong Tian, Dhruv Batra, and Devi Parikh. 2019. Codraw: Collaborative drawing as a testbed for grounded goaldriven communication. In *ACL*.
- Juil Koo, Ian Huang, Panos Achlioptas, Leonidas J Guibas, and Minhyuk Sung. 2022. Partglot: Learning shape part segmentation from language reference games. In *ICCV*.
- Mike Lewis, Denis Yarats, Yann Dauphin, Devi Parikh, and Dhruv Batra. 2017. Deal or no deal? end-to-end learning of negotiation dialogues. In *EMNLP*.
- John Locke. 1978. Two treatises of government. New York: E. P. Dutton.
- Xinyao Ma, Maarten Sap, Hannah Rashkin, and Yejin Choi. 2020. Powertransformer: Unsupervised controllable revision for biased language correction. In *EMNLP*.

Jonathan D Moallem and William L Raffe. 2020. A review of agency architectures in interactive drama systems. In 2020 IEEE Conference on Games (CoG), pages 305–311. IEEE.

712

713

714 715

716

718

719

721

724

725

727

728 729

730

731

732

733 734

735

736

737

738

740

741

742

743

744

745 746

747

748

752

753

754

756

- Changhoon Oh, Jungwoo Song, Jinhan Choi, Seonghyeon Kim, Sungwoo Lee, and Bongwon Suh. 2018. I lead, you help but only with enough details: Understanding user experience of co-creation with artificial intelligence. In *CHI*.
- Grant Pickett, Allan Fowler, and Foaad Khosmood.
 2015. Npcagency: conversational npc generation.
 In Proceedings of the 10th International Conference on the Foundations of Digital Games.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. Towards empathetic opendomain conversation models: A new benchmark and dataset. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5370–5381.
- Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Eric Michael Smith, Y-Lan Boureau, et al. 2021. Recipes for building an open-domain chatbot. In ACL.
- Maarten Sap, Marcella Cindy Prasettio, Ari Holtzman, Hannah Rashkin, and Yejin Choi. 2017. Connotation frames of power and agency in modern films. In *EMNLP*.
- Ryan Volum, Sudha Rao, Michael Xu, Gabriel Des-Garennes, Chris Brockett, Benjamin Van Durme, Olivia Deng, Akanksha Malhotra, and William B Dolan. 2022. Craft an iron sword: Dynamically generating interactive game characters by prompting large language models tuned on code. In *Proceedings* of the 3rd Wordplay: When Language Meets Games Workshop (Wordplay 2022), pages 25–43.
- Xuewei Wang, Weiyan Shi, Richard Kim, Yoojung Oh, Sijia Yang, Jingwen Zhang, and Zhou Yu. 2019. Persuasion for good: Towards a personalized persuasive dialogue system for social good. In *ACL*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed H Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*.
- Wen Wen and Hiroshi Imamizu. 2022. The sense of agency in perception, behaviour and human-machine interactions. *Nature Reviews Psychology*, 1(4):211–222.

761

764

765

768

770

771

774

775

776

778

779

781

782

A Dataset Statistics

Feature	N/A	No	Moderate	Strong
Intentionality	_	194	175	539
Motivation	_	474	158	276
Self-Efficacy	770	63	46	29
Self-Regulation	764	25	61	58

Table 3: Statistics of the annotated conversation snippets. N/A indicates not applicable. We annotate Self-Efficacy as N/A if a designer never indicated a preference or did not need to pursue their preference (e.g., because the other designer did not argue against it). We annotate Self-Regulation as N/A if a designer Never indicated a preference or did not need to change their preference (e.g., because the other designer did not argue against it).

	Low	Medium	High
Agency	308	292	308

Table 4: Agency distribution of the conversation snippets.

Other Statistics. The conversations b/w interior designers in our dataset have 41.67 turns on average. The extracted conversation snippets have 4.21 turns on average. We find an average pairwise agreement of 71.36% for Intentionality, 70.70% for Motivation, 85.21% for Self-Efficacy, and 81.09% for Self-Regulation.

B Model Configurations

We use text-davinci-003 for all of our GPT-3 models. For Agency measurement models (Section 7), we sample the highest probable next tokens by setting the temperature value to 0 (deterministic sampling). For dialogue generation models (Section 8), we use top-p sampling with p = 0.6. For in-context learning methods, we experimented with k = 5, 10, 15, and 20 and found k = 10 to be the most effective based on a qualitative assessment of 10 examples.

C Why We Chose Collaborative Interior Designing as Our Testbed?

Here, we propose a **dialogue-based collaborative interior design task** as a testbed. In this task, given a room setting, the goal is to discuss how to design the interiors of the room. We note that an interior design task can be broad and may involve a wide range of complex components (e.g., color palette, furniture, accessories) as well as a series of steps to be followed. Furthermore, due to a real-world room context, the task must be grounded with both vision and language components with an understanding of how threedimensional objects in a room (e.g., chairs, tables, plants, decor items) must be designed.

785

786

787

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

Here, we build upon previous work on richlyannotated, large-scale datasets of 3D objects like ShapeNet (Chang et al., 2015) and subsequent works on understanding how fine-grained differences between objects are expressed in language like ShapeGlot (Achlioptas et al., 2019) and Part-Glot (Koo et al., 2022). Both ShapeGlot and PartGlot datasets provide us with richly annotated datasets of chairs. Therefore, we narrow down the scope of our task and specifically focus on *furnishing a room with a chair*. In this task, a human and an AI are provided with a room layout and asked to collaboratively come up with a design of a chair to be placed in the room through text-based interaction.

D Extract Conversation Snippets associated with different Design Components

We observe that designers typically discuss these components one at a time (in no particular order). Therefore, we aim to extract a contiguous sequence of utterances that represent the design element being discussed. Let \mathcal{D}_i be a dialogue with utterances u_{i1}, u_{i2}, \dots For a specific design component d_{ij} in its final design (e.g., "metal legs"), we first retrieve the utterance u_i that most closely matches with it (based on cosine similarity b/w RoBERTa embeddings) - the conversational snippet associated with d_{ij} should at least include u_j . Next, we determine the contiguous utterances before and after this matched utterance that discuss the same higher-level design component (e.g., if d_{ij} was "metal legs", the utterances may focus on discussion of the higher-level component "legs"). We create a simple k-means clustering method to infer the higher-level component being discussed in utterances through their "design clusters". Then, we extract all contiguous utterances before and after u_i with the same design clusters as u_i .

Agency Feature	Coefficient
Intentionality	0.1435*
Motivation	0.0235
Self-Efficacy	0.0384
Self-Regulation	-0.1224*

Table 5: Coefficients for predicting agency in conversations using a mixed-effect linear regression model. *p < 0.05

Ε **Analysis of Agency Features**

F **Task 1: Demonstration examples**

F.1 GPT-3 (Q/A)

833

834

835

For the GPT-3 (Q/A) model, we present examples to GPT-3 in the following format:

838 Designer: I think a black wooden frame or black metal legs (to match the bed 839 frame) would work. Other Designer: I like the black metal 841 legs. What about hairpin legs? 842 Designer: Or maybe brass legs would be better. Hairpin legs would work fine, but would the rest of the frame be the 846 black wood? Other Designer: If we did brass tapered metal legs it would tie well with the 848 black wood. Designer: I think that would look better. 850 Other Designer: Agreed Who influenced the design element 853 being discussed?: Other Designer 856 **F.2 GPT-3 (CoT)** For the GPT-3 (CoT) model, we present examples 857 to GPT-3 in the following format: Designer: I think a black wooden frame or black metal legs (to match the bed 860 frame) would work. Other Designer: I like the black metal legs. What about hairpin legs? Designer: Or maybe brass legs would 865 be better. Hairpin legs would work fine, but would the rest of the frame be the 866 black wood? 867 Other Designer: If we did brass tapered

metal legs it would tie well with the	869
black wood.	870
Designer: I think that would look better.	871
Other Designer: Agreed	872
	873
TL;dr Brass tapered metal legs were	874
agreed upon. This was initially proposed	875
by the Other Designer.	876
G Reproducibility	877
We will release the codes and datasets developed	878
in this paper at https://anonymous under an MIT	879
license.	880

881

882

883

884

The use of existing artifacts conformed to their intended use. We used the OpenAI library for GPT-3 based models. We use the scipy and statsmodel libraries for statistical tests in this paper.

H Human-Human Conversational Data Collection Instructions

Figure 5: Instructions shown to the interior designers during the human-human conversational data collection. Continued on the next page (1/3).

Instructions

In this data collection study, you will plan to design an object in collaboration with another participant. You will access a website using a link that we provide. On the website, you will be paired with another participant, with whom you will interact, via a chat-like interface (text-only), to plan and negotiate what you collaboratively want to design.

Purpose of the Research

The purpose of this research is to understand **agency in human-human conversations** and **how to build a conversational AI agent with agency**. Agency can be defined as the power one has to act upon their intrinsic motivation, preferences, and expertise. Here, we want to study how humans exercise agency in conversations, as well as, how AI agents can exercise agency through conversations.

Towards this goal, we are collecting conversations around tasks involving **two humans planning to collaboratively design an object** (e.g., *a chair*). The conversational data would help us assess how humans use conversations to exercise their agency and how we can train AI agents to have agency, without becoming insensitive towards others or disregarding social norms.

The Setting

You will be paired with another participant. You will both be shown a 3D model of a room. Here is an example room:



What will you do?

You will be assigned an object (e.g., *a chair*). You will **plan to design that object for the room, in collaboration with the other participant, through chat conversations**.

Here are the steps you will follow:

Figure 6: Instructions shown to the interior designers during the human-human conversational data collection. Continued on the next page (2/3).

Step 1. Propose your preferred object design: For the object you are assigned, you will first propose the design you prefer.

- a. To help you in this process, you will be shown several different designs for that object and will be asked to **select the designs** you like, based on the room shown.
- b. You will then use the selected object designs to propose your preferred design. E.g., if you are assigned a chair, you will describe the type of the chair, the characteristics of the back, seat, arms and legs, color, and/or the type of material you prefer.
- c. While proposing your preference, you could also indicate whether your preference is strong or weak.
- d. Here are a few example object designs with proposed preferences:



"I would strongly prefer a black swivel chair with rollers on the feet. The chair could have no arms but I don't mind if they have arms. I would also prefer a smaller back and a wider seat."



"I would prefer a straight wooden chair with bars on the back. I strongly prefer the chair to have no arms and have a cushion. The top of the back could be rounded."



"I would strongly prefer a club chair with padded seat, back, arms, and legs"

Note

- 1. Your proposed preference may be different from the designs you select (if you wish to innovate).
- 2. You should not directly share the designs you select or your proposed preference with the other player.

Step 2.1. Plan what to design: Next, you will start planning your design collaboratively with the other participant. You will **use a simple chat-like web interface** to interact with the participant you are paired with.

- a. The design you prefer might be *different* from the design which the other player prefers.
- b. Therefore, a key part of the collaborative designing process would be to **communicate your individual preferences, negotiate, and find common ground**.
- c. You will use the chatbox to plan, discuss and negotiate.
- d. You should try and convince the other player to agree on a design that is close to your preference.
 - i. For example, you can try and explain why the design you prefer might be better.
 - ii. At the same time, it is also important to understand the other player's preference. Knowing that can help you talk about the pros and cons of each design.

Figure 7: Instructions shown to the interior designers during the human-human conversational data collection (3/3).

- iii. You can also discuss what adjustments can be made such that the final design satisfies the preferences of both the players.
- e. You should plan to spend ~30 minutes on the conversation.

Step 2.2. Describe the final chair design: Both you and the other participant will be provided with a textbox, which you both will use to **report the design that you agreed upon**.

- a. You should use this textbox to **update the current design** when you agree upon something (based on what is being discussed in the conversation).
- b. For example, if you are asked to design a chair, and if you are able to decide the high-level chair design first (e.g., *a club chair*), you can update it in the textbox, before proceeding to discuss the other characteristics (e.g., *seat, arms, legs*).
- c. Please be as specific as possible when describing your design.

Step 3. Mark as finished and take a post-study questionnaire: When both you and the other player are done designing the object, you will mark the study as complete (using a provided option) and take a post-study questionnaire.

- a. Note that you may not always reach an agreement with the other participant. But when you are done, you should still mark the task as finished and take the post-study questionnaire.
- b. You should plan to spend 15-20 minutes on the questionnaire.

Note: The conversations should only focus on object design. To keep the conversations natural, please do not discuss things related to these instructions directly in the conversation. For instance, you should not mention that you went through a process of selecting designs or writing a preference (e.g., do not say "what is your preferred design?" or "my preferred design is..."). Also, do not discuss any personal details.

Designer	Utterance
Designer 1:	How about a desk chair for this area?
Designer 2:	There seems to be many possibilities for this space, would you agree? Yet I agree that some kind of chair for the desk is needed.
Designer 1:	The room has very clean lines with an Asian theme
Designer 2:	I think we need to support the minimalist lines of the overall space design. Not something too over-stuffed. Something with a contemporary feel.
Designer 1:	So maybe a more contemporary style of desk chair.
Designer 1:	Great minds!
Designer 1:	How do you feel about a tall back with tilt swivel and adjustable
Designer 2:	I believe so. Maybe one that is comfortable for sure - but not too closed in. There is the lovely background to consider. We don't want to block that.
Designer 1:	If not too tall, then maybe something mid back height?
Designer 2:	I think the height of the back should be carefully scaled - supportive but not so high that it obscures what is behind too much.
Designer 1:	Or shoulder height for support
Designer 1:	With arm support
Designer 2:	Agreed on shoulder height. Swiveling is good - also moving -like on casters may provide flexibility.
Designer 1:	Definitely casters
Designer 2:	I am concerned about tilting back since we do have some fragile decorative elements behind.
Designer 1:	Ok, so far shoulder height desk chair with adjustable height, casters and arm rests
Designer 2:	I do agree that arm support is essential, especially if one is to feel comfortable while working. It feels like this might be a consult room of sorts - so allowing the person to sit back in a more relaxed posture - resting arms off the table is good.
Designer 1:	Some tilts can be regulated and locked into place not necessarily a full recline
Designer 1:	Perfect
Designer 2:	The materiality of the chair is something to consider. I see a lot of wood and timber detailing. It might be nice to have the chair upholsterable - perhaps a nice leather back that would be shaped to lightly massage the back?
Designer 1:	Agree
Designer 1:	the leather would be a nice look in there
Designer 2:	Something that seems pillowy or wavy, but in a very restrained, minimalist sort of way
Designer 1:	Black would match the ottomans but a soft buttery cream/ ivory would add a soothing neutral to the aesthetic
Designer 2:	With the darker wood in the room and the leather chair - an accent material on the armrests might be nice to offsett - say a brushed steel or aluminum finish?
Designer 1:	I've seen the vertical channeling on a desk chair that is very classy looking
Designer 1:	The brushed steel frame would look nice in this room. I think wood would be a bit much.
Designer 2:	I think classic modern which always took a lot of inspiration from japanese design. The buttery cream is a lovely idea. Will provide a bright focal point and it will align with the colors of the fan.
Designer 1:	I think we have our chair!

Table 6: Example Human-Human Conversation in Our Dataset.

I Human Evaluation Experiment Instructions

Figure 8: Instructions shown to the interior designers during the human evaluation experiment. Continued on the next page (1/2).

Agency Evaluation

Study Goals

The goal of this study is to interact with and evaluate chatbots.

Study Steps

In the study, you will interact with two AI-based chatbots, one at a time. Each time, you will be provided with a room description and a specific chair design component (e.g., the material to be used for a chair that will be placed in the room). Your task will be to collaborate with the chatbots to discuss and agree upon what the chair design component should be.

In the end, you will fill out a questionnaire in which you will be asked questions comparing the two chatbots. You will compare the chatbots based on whether they were able to pose, motivate, and stick to their own preferences and whether they were able to influence the final design.

Few Important Things to Note

- 1. Aim to spend between 2 to 5 minutes per chatbot: You should aim to chat for around 2 to 5 minutes with each chatbot.
- 2. Chat only about the component you are assigned: Please chat only about the chair design component you are assigned. In some cases, the chatbot may try initiating a conversation about a different design component. However, that is not required, particularly after you have agreed on what the assigned design component should be.
- 3. **Express your preferences:** You may start by expressing your preference or by asking if the chatbot has any preference.
- 4. Negotiate what you don't like or agree with: If you do not agree with the preference of the chatbot, you should negotiate with it and try to convince it otherwise.
- 5. "End Conversation and Continue" once you are done: One both you and the chatbot have agreed upon what the design element should be, please use the "End Conversation and Continue" to proceed to the next step of the study.

Figure 9: Instructions shown to the interior designers during the human evaluation experiment (2/2).

6. Back/Next button Trick: If something doesn't work or gives an error, please try pressing the back button on the broswer and the press the "Continue" button again.

Consent to the study

By ticking this box, you are agreeing to be part of this data collection study. You also confirm that you understand what you are being asked to do. You may contact us if you think of a question later. You are free to release/quit the study at any time. Refusing to be in the experiment or stopping participation will involve no penalty or loss of benefits to which you are otherwise entitled. To save a copy of the consent form and instructions, you can save/print this webpage (or find the instructions here). You are not allowed to distribute these instructions and data for any purposes. You are also not allowed to use them outside this study.

Agree and Continue