

LVLMs and Humans Ground Differently in Referential Communication

Peter Zeng^{1,4} Weiling Li² Amie Paige² Zhengxiang Wang^{3,4}
Panagiotis Kaliosis¹ Dimitris Samaras¹ Gregory Zelinsky²
Susan Brennan² Owen Rambow^{3,4}

¹Department of Computer Science ²Department of Psychology
³Department of Linguistics ⁴Institute for Advanced Computational Science
Stony Brook University

Correspondence: pezeng@cs.stonybrook.edu

Abstract

For generative AI agents to partner effectively with human users, the ability to accurately predict human intent is critical. But this ability to collaborate remains limited by a critical deficit: an inability to model **common ground**. We present a referential communication experiment with a factorial design involving director-matcher pairs (human-human, human-AI, AI-human, and AI-AI) that interact with multiple turns in repeated rounds to match pictures of objects not associated with any obvious lexicalized labels. We show that LVLMs cannot interactively generate and resolve referring expressions in a way that enables smooth communication, a crucial skill that underlies human language use. We release our corpus of 356 dialogues (89 pairs over 4 rounds each) along with the online pipeline for data collection and the tools for analyzing accuracy, efficiency, and lexical overlap.¹

1 Introduction

Human conversation relies on common ground, accrued and updated by interacting partners (Clark and Brennan, 1991; Clark and Wilkes-Gibbs, 1986; Clark and Marshall, 1981). During conversation, partners ground meanings with one another, adapting the referring expressions they use to pick out objects of interest, such that there is less variability *within* a conversation than *between* conversations by different partners discussing the same objects (Brennan and Clark, 1996).

Conversational partners cannot read one another’s minds, but they can signal that they believe they are talking about the same referent by converging on a referring expression for it (the process of *lexical entrainment*; *ibid*, see Figure 1). They cannot rely solely on conventional word meanings, as the objects of interest may not be uniquely associated with distinctive labels known to both. They



	
<p>Round 1 D: 9 is tall like a laundry basket with two small handles, it's round M: like a cylinder M: got it D: exactly, but the taller one of the two</p> <p>Round 2 D: Next one is the tall laundry hamper looking one M: got it</p> <p>Round 3 D: 6 tall laundry basket M: got it</p> <p>Round 4 D: 7 is tall laundry basket M: got it</p>	<p>Round 1 D: 6 looks like a rodent with eyeballs that are looking at you M: got it</p> <p>Round 2 D: 7 is the rodent looking at you M: scary got it</p> <p>Round 3 D: 3 scary rodent M: ...got him</p> <p>Round 4 D: 9 scary rodent M: ye</p>

Figure 1: Repeated referring to two baskets (non-lexicalized objects) by a human-human pair, D and M, in Rounds 1-4 of our experiment, with lexical overlap highlighted in blue. Entrainment on more concise language signalling a *conceptual pact* occurs by Round 3, after they consider multiple proposals in Rounds 1-2.

construct meaning jointly by engaging in the process of *grounding*, or seeking and providing evidence about what they mutually believe to be in their common ground.

Recently, researchers have begun to address the question of whether large language models (LLMs) and large vision language models (LVLMs) engage in grounding, as human discourse partners do (Tang et al., 2024; Hua and Artzi, 2024; Imai et al., 2025; Shaikh et al., 2025; Wang et al., 2025). This question is important both in order to try to understand how models work, and for practical applications in which an AI agent assists a human in a task that requires using language to pick out elements from a visual environment. Studies to date tend to conclude that LVLMs lack the *pragmatic competence* needed to coordinate with a partner, and that they struggle in multi-turn conversation.

This paper adds further empirical evidence about

¹<https://github.com/peterzeng/lvlms-referential-game>

the pragmatic competence missing from LVLMs. This is the first study, to our knowledge, to generate (in real time) multi-turn task-oriented dialogues between pairs of partners in asymmetrical roles (offering different levels of initiative), and that covers all four combinations of human or AI partners filling each role. We conducted an experiment to examine language use by director-matcher pairs collaborating to do a referential communication task, based closely on [Clark and Wilkes-Gibbs \(1986\)](#)'s study that measured how partners accumulate common ground as they match the same objects across multiple rounds. We tested all four combinations of director/matcher roles (human-human, human-AI, AI-human, and AI-AI). This allowed us to quantify not only where LVLMs fail in collaboration, but also to discover why they fail (and in which roles they do so). The experimental transcripts are available [here](#).

This paper is organized as follows: Section 2 motivates the project within the cognitive science of human communication, and Section 3 surveys relevant work in human-AI interaction. The experimental design is described in Section 4, and Section 5 presents the results along with illustrative examples from the collected dialogues. Section 6 presents follow-up AI-AI experiments and Section 7 concludes.

2 Cognitive Science Background

Psycholinguists cannot read the minds of those they study, so they rely on experiments to provide observable evidence about how humans process language. Referential communication is typically studied using matching tasks with carefully chosen stimuli; two partners are provided with a (typically shared) goal such as to discuss and manipulate picture cards, maps, or other visuospatial elements, and their behavior (e.g., words, eye gaze, actions, etc.) is measured. Task roles, available information, and other partner characteristics can be balanced, or else manipulated in order to vary each partner's initiative or other factors of interest.

In addition to the process of grounding described in the previous section, other key forces shape dialogues, including pragmatic factors such as [Grice \(1975\)](#)'s *Cooperative Principle* that governs human expectations when they engage in dialogue, along with its Maxims of *Quality* ("be truthful"), *Quantity* ("say enough but not too much"), *Manner* ("be straightforward") and *Relation* ("be relevant").

Dialogue is also shaped by the costs and affordances of the communication medium ([Clark and Brennan, 1991](#)). Remarkably, those who share a common purpose (even one given to them by an experimenter) tend to increase their effort as much as it takes for them to accomplish the task, including distributing their individual effort collectively (e.g., depending on their task roles or on the affordances of the communication medium, one partner may work harder to make up for what's difficult for the other). In the context of grounding in different sorts of dialogues, human partners expend the *least collaborative effort* needed to meet a *grounding criterion* sufficient for current purposes ([Clark and Wilkes-Gibbs, 1986](#)), which could be perfect accuracy in an air traffic control dialogue, high accuracy in an experimental task, or polite engagement while passing the time in a checkout line.

In an instrumental task, accruing common ground allows human partners to increase their joint efficiency while working together (meaning that they can expend less effort to maintain the same or increasing levels of performance). Figure 1 illustrates the emergence of *conceptual pacts*, or temporary, flexible, shared perspectives about referents in conversation ([Brennan and Clark, 1996](#)). These differ from "conventions", which exist outside of individual dialogues as word meanings known to the members of a language community.

Language use by LVLMs, however, is quite a different matter. Training data for LLMs captures statistical regularities that emerge from conventionalized word meanings within texts, rather than from interactively established conceptual pacts.

3 Related Work in Human-AI Interaction

Related work on referential communication with LVLMs typically uses human-human corpora as a gold standard for human-AI interaction, finding consistently that humans perform better than models. Although such corpora include spontaneous turn-taking by humans, testing with AI partners tends to not allow multi-turn interaction, so does not test their ability to collaboratively ground meaning or repair misunderstandings.

For example, *PhotoBook*, a widely used dataset for exploring common ground in visually grounded dialogues ([Haber et al., 2019](#)), has human participants play a multi-round online image identification game in which each sees a grid of six visually similar scenes, with some images shared and oth-

ers visible to only one partner. The partners chat via text and decide which images among a set of highlighted targets are common to both partners or private to only one. Imai et al. (2025) adapted PhotoBook to an AI-AI setting using LVLMs and evaluated the AI dyads against the human ones using a suite of proposed grounding metrics. Although LVLM pairs achieved near-human task accuracy, their dialogue differed from human pairs in the formation of common ground, both in efficiency and in lexical adaptation.

Hawkins et al. (2020) had human-human and human-AI director-matcher pairs identify targets within sets that consisted of four photos from the COCO data set, with the goal of being able to use more efficient expressions (dubbed “ad-hoc conventions”) upon repeated referring to the same photos. Such photos were far more distinctive than the basket targets that we used, making their task much easier (and possibly the labeled COCO data would have been in the models’ training data). Hawkins et al. (2020)’s model adapted over the course of repeated rounds of references to the same objects; however, there was no multi-turn interaction, as human directors simply typed referring expressions and matchers responded (5-10 s later in the case of AI matchers) by selecting the target (before receiving feedback about correctness).

Using a corpus of nearly 3,000 human-written referring expressions, Tang et al. (2024) tested model performance of LVLMs as both “speakers” and “listeners” (analogous to director/matcher roles) in a visual-spatial environment in which the two agents had distinct spatial perspectives. The models performed more poorly at resolving referents (and in both roles) than did humans, but training an open-weight model with evidence of communicative success improved its performance (albeit still not to human-human levels).

Hua and Artzi (2024) borrowed Hawkins et al. (2020)’s four-object task and corpus in order to test whether LVLMs could adapt to simulated human partners (as humans do to their real partners) by forming “ad hoc conventions” consisting of increasingly efficient referring expressions over rounds of repeated referring. This project used pragmatically-inspired prompting strategies (the best-performing one being to produce shorter and shorter messages with consistent lexical content) rather than any learning on the part of the LVLM. Again, as in Hawkins et al. (2020), directors presented one-off descriptions, with matcher simply selecting a tar-

get, repeated over multiple rounds. LVLMs in the matcher role were often able to interpret the increasingly shortened expressions, but LVLMs in the director role were less able to produce more efficient expressions. For four different LVLM models, accuracy in this simple task ranged from about 40% to ceiling.

Hua et al. (2025) built on this work by developing an interactive (rather than simulated) post-training process to induce “ad hoc conventions” (accurate yet concise expressions) through fine-tuning in text-only reference tasks. They used two platforms to evaluate their methodology: an interactive form of the game Taboo in which one partner referred (in a single turn) to a target for the other partner (or system) to identify, but without mentioning its conventional name, and a document-based question-answering task that took multiple turns to complete. By evaluating several state-of-the-art proprietary and open-source LLMs, they found that LLMs still lack the ability to spontaneously develop ad-hoc concise referring expressions. Their post-trained LLMs, on the other hand, shortened their messages by up to 26% in their reference task benchmark and outperformed off-the-shelf counterparts in their document-based task.

Wang et al. (2025) cast LVLMs as overhearers to a transcribed corpus of referential communication produced by humans in both the director and matcher roles, reprising Schober and Clark (1989)’s experiments with human overhearers that found that matchers performed better than overhearers because they could ground meaning interactively whereas overhearers could not. Wang et al. (2025)’s LVLM overhearers lagged well behind Schober and Clark (1989)’s human overhearers in accuracy and failed to show any improvement over repeated rounds of the matching task (unlike human overhearers). In another study with text transcripts of human task-oriented dialogue, Sarkar et al. (2025) showed that LLMs were able to recognize some cues of misalignment (dubbed “conversational friction”), especially when occurring with clarification questions, but often missed subtler cues. Shaikh et al. (2025) analyzed grounding in logs of human-LLM interaction (WildChat, MultiWOZ, and Bing Chat) using a taxonomy of “grounding acts” (Traum and Allen, 1992). They found that LLMs initiated grounding far less often than humans did, with early failures to ground predicting poor performance later.

In recent studies using a paradigm similar to

Hawkins et al. (2020)’s, Jones et al. (2026) elicited lexical overlap (that they also referred to as “convention formation”) within dyads during 50 rounds of referring by directors in which matchers chose single target tangrams out of five successive arrays of 10. After the director’s referring expression, the matcher selected a referent (without conversational interaction) and both received feedback about which tangram was selected and whether it matched the intended target. Within human-human, AI-AI, and mixed human-AI dyads, partner roles alternated for each round. Results showed that same-type dyads (human-human, AI-AI) improved reliably in accuracy, with more lexically consistent expressions over rounds, whereas mixed human-AI dyads performed poorly, even when the model was prompted to produce more “humanlike” utterances with more concise referring expressions. This suggests that forcing conciseness via prompting is not sufficient to achieve the effects of common ground and conceptual pacts in referring (consistent with Hawkins et al. (2020)’s findings).

4 Experimental Design and Method

4.1 Task Description

Our experiment used a referential communication task with human and AI partners paired in director and matcher roles (human-human, human-AI, AI-human, and AI-AI), interacting via multiple, unrestricted text turns. In each of four rounds (with the same director, matcher and set of referents across rounds), the director saw a target sequence of 12 baskets (in a different order for each round) and described them one at a time to the matcher, interacting freely for the matcher to disambiguate each intended referent and move it from a staging area on their screen to the target area to match the director’s order. The matcher’s set included not only the same 12 baskets from the director’s set, but also another 4 so that the last basket(s) would not be trivially easy to identify. This task, based closely on Clark and Wilkes-Gibbs (1986), allowed us to observe how partners interactively ground referring expressions and to monitor how they interpret the expressions to resolve or repair meanings and form conceptual pacts across rounds (Brennan and Clark, 1996). The set of stimulus baskets (Figure 3) and the interface (Figure 4) are shown in Appendix A.

To allow direct comparisons across types of partners and roles, all four director-matcher conditions were as similar as possible, hosted online with the

platform oTree, an open-source Python package for web-interactive tasks and behavioral research experiments (Chen et al., 2016). The 3 conditions that involved human partners recruited participants on the platform Prolific (Prolific, 2014), a high-quality source of vetted and motivated participants who seek to maintain good performance ratings on the platform. As the experiment was online, all communication was via text. The task screen layout had both baskets and chat windows easily visible without scrolling (although scrolling was available to both partners so that they could view prior dialogue from that round). As is typical in interactive human chat media, a series of dots (. . .) was displayed whenever a partner was typing (this too is a useful cue for grounding).

4.2 Human Participants

The task was advertised as a job on Prolific that involved communicating with either another human or an AI partner. Given the linguistic demands of the task and need for attention to good task performance for high quality data, we used strict pre-screening criteria, restricting participation to fluent English speakers located in the U.S. In addition, to mitigate the quality control challenges inherent in online data collection (e.g., potential cheating or attrition), we updated our filtering criteria to include only participants with a perfect approval record (100%; see Table 3 for details). The task took an average of one hour to complete ($SD = 28$ min). Recruits were told that their partner would be either human or AI, without revealing which.

4.3 AI Participants

We selected OpenAI’s GPT-5.2 with the “none” reasoning option as the LLM for our AI-AI, human-AI, AI-human experiments (OpenAI, 2025). Following Hua and Artzi (2024); Wang et al. (2025), we began with GPT-4o when initially developing the prompting. Four rounds of identifying 12 objects turned out to be too complex for that model and required a more capable model. After evaluating several frontier models, we settled on GPT-5.2, first because it was the newest frontier model at the time, and second, for its superior instruction-following capabilities, even at the lowest reasoning level. We used the lowest reasoning level to prioritize responsiveness, as human conversation unfolds under time pressure and humans strive to minimize pauses in conversation (Jefferson, 1989).

We designed the system prompts to align

with the instructions given to human participants while providing the necessary scaffolding for autonomous interaction. The prompt structure consisted of three key components: (1) Task Context and Role Definition, which mirrors the rules provided to human subjects; (2) Communication Norms, where we explicitly codified pragmatic constraints often implicit to humans (e.g., conciseness, turn-taking, and comparative language) to support natural collaboration; and (3) Structural Scaffolding, which enforced a strict JSON output format to track state updates and induce zero-shot chain-of-thought reasoning (Kojima et al., 2022) before generating an utterance or action. In addition, we provided the full history of utterances from both partners, with the matcher receiving an additional image as context that was updated whenever a basket was placed or moved. The full prompts are provided in Appendix E.

4.4 Resulting Corpus

Sessions were removed as invalid when there was evidence of cheating or of not following task instructions; see Table 3 for the proportion removed. A total of 32 human pairs and 39 humans with AI partners successfully completed all four rounds of the task, as well as a post-task survey (totaling 103 human participants that included 58 women, 39 men, and 6 others, with a mean age of 46.6 years, $SD = 13$). Full demographic information is reported in Table 4. The final corpus consists of referential communication dialogues from 32 Human-Human, 22 Human-AI, 17 AI-Human, and 18 AI-AI director-matcher pairs.

4.5 Evaluation Metrics

We evaluated human and AI grounding capabilities using a set of simple metrics that capture the establishment of common ground across three aspects of interaction: communicative success, communication effort, and lexical entrainment. When reporting results, we first compute each metric for each pair given a specific round and then average the numbers obtained across pairs for that round.

Communicative Success We measured communicative success using **accuracy**, defined as the percentage of correctly matched baskets within a task round. In matching tasks like this one, humans improve over rounds and reach performance at ceiling, and models perform more poorly, sometimes with no improvement over rounds or declining ac-

curacy.

Communicative Effort As conversational partners progressively establish common ground, they require less effort to communicate. We quantified this using the **number of words** and the **number of turns** in each task dialogue as proxies, capturing overall communication cost and back-and-forth coordination, respectively. Other related measures (e.g., number of utterances) patterned similarly, and are reported in Appendix C.1.

Lexical Entrainment Lexical entrainment refers to the tendency of interlocutors to reuse lexical material (in human conversation, this happens in a more concise manner) when referring to previously established objects (Brennan and Clark, 1996; Garrod and Anderson, 1987). We operationalized this concept using two metrics that capture length reduction and lexical overlap of referring expressions across rounds.

Let $RE_i^{(b)}$ denote the referring expression or description for basket b in round i , and let $Tok(RE_i^{(b)})$ denote its sequence of content words. We defined the length of a referring expression as the number of content words, i.e., as $|Tok(RE_i^{(b)})|$ (**number of RE words**).

To quantify lexical overlap, we computed round- i relative lexical overlap (RLO_i , **proportion of lexical overlap**) as the multiset overlap of content words ($Inters(A, B)$) between the referring expression used in round i and the one used in round $i - 1$. RLO always ranges between 0 and 1 and is defined as follows.

$$RLO_i^{(b)} = \frac{|Inters(Tok(RE_{i-1}^{(b)}), Tok(RE_i^{(b)}))|}{|Tok(RE_i^{(b)})|}$$

This lexical overlap measure is similar to other measures used in the literature; we discuss these in Appendix C.2. We use our definition because it ignores word order and has an immediate intuitive interpretation aligned with entrainment. While higher overlap reflects greater lexical reuse, it does not by itself imply increased efficiency; we therefore interpret overlap jointly with length-reduction trends when assessing lexical entrainment.

Referring Expression Extraction We used GPT-5 (OpenAI, 2025) to automatically extract referring expressions (REs) for each target basket at the dialogue level. We validated the extracted expressions

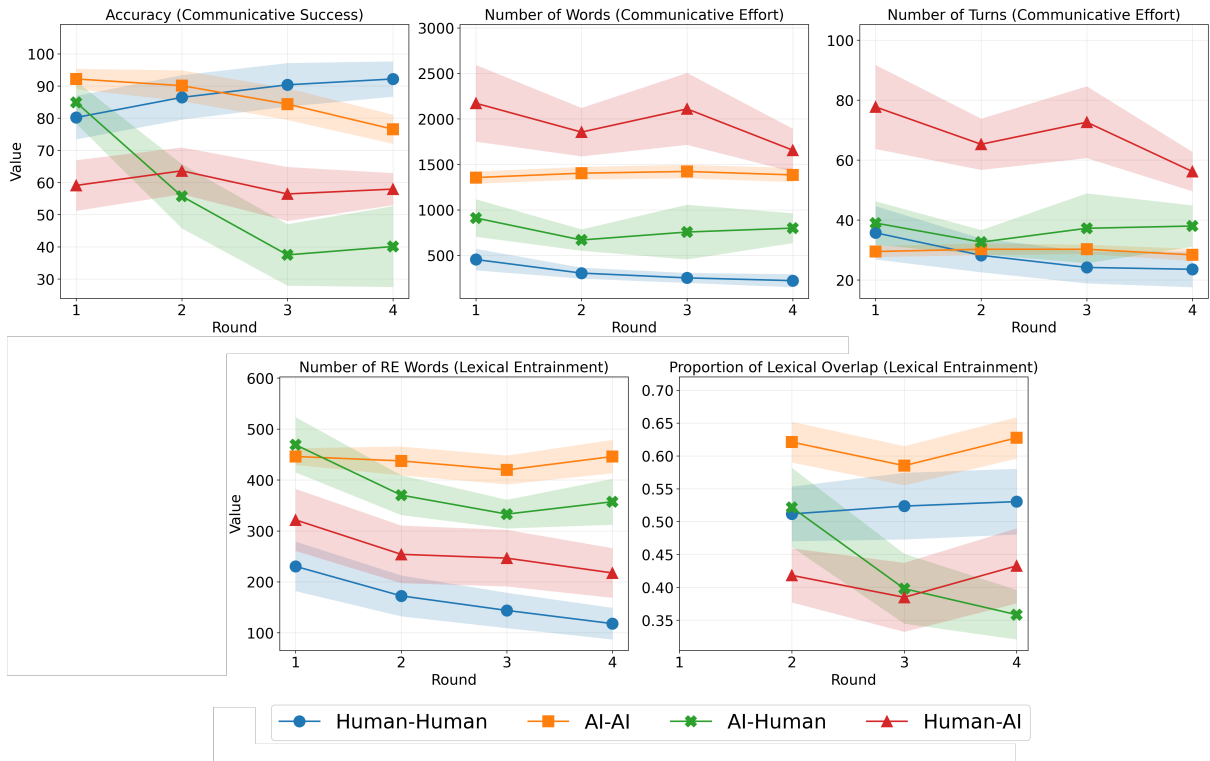


Figure 2: Trends over four rounds for (from left to right) **accuracy (%)**, **numbers of words**, **number of turns**, **number of words referring expressions**, and **proportion of lexical overlap with prior rounds** by director-matcher condition. Dots show means with 95% CIs, with each color denoting a specific pairing condition.

against human annotations on a comparable corpus, obtaining an F_1 score of 0.86 under ROUGE-L (Lin, 2004). See Appendix C.3 for details of the validation setup and results.

	HH	AA	AH	HA
Accuracy	4.0**	-5.3***	-15.3***	-1.1
# Words	-74.9***	10.6	-24.7	-129.1
# Turns	-4.1**	-0.3	0.2	-5.7*
# RE Words	-36.6***	-1.8	-37.3***	-32.0*
L Overlap	0.0	0.0	-0.1***	0.0

Table 1: Ordinary least squares (OLS) regression slopes (c) of each metric across rounds by director-matcher condition. H: human. A: AI. Green $c > 0$, red $c < 0$; stars: *** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$.

5 Results

Figure 2 shows round-by-round trends across director-matcher conditions with means and 95% confidence intervals for five evaluation metrics that capture the three aspects of interaction relevant to grounding capabilities (accuracy, effort, and lexical entrainment). We used ordinary least squares (OLS) regression to quantify the overall trends for these metrics and report the slopes along with the

respective significance levels in Table 1.

Human-Human Human-human pairs achieved high communicative success (accuracy), starting at 80% and increasing steadily over rounds to over 90% by Round 4. In terms of communicative effort, the numbers of words and turns per round started relatively low and decreased over rounds. For the lexical entrainment measures, the number of RE words used to describe each basket was also low and decreased, and the REs consistently re-used a high percentage of lexical material from the previous round (above 50%).

For example, the pair who contributed the dialogues in Figure 1 completed all four rounds of the task in only 28 min and with only two errors (during the first round). Both dialogues show evidence of entrainment by the third round (with the matcher contributing lexical content to the resulting “scary rodent” conceptual pact). Unlike the AI partners, human partners sometimes explicitly acknowledged their common ground with metalinguistic references, as another pair did with descriptions like “the one we worked hard on last time” (See Figure 5).

The quantitative data and corpus show clearly

that humans entrained: the REs they used got shorter, but remained related to the previous round's RE (lexical overlap) and overall, the interactions became shorter (in both words and turns) and more efficient over rounds, meaning that the partners expended less linguistic effort to understand each other increasingly well. These results are consistent with previous referential communication experiments (e.g., Brennan and Clark, 1996; Clark and Wilkes-Gibbs, 1986), despite our novel setup (online and typing-only).

AI-AI The AI-AI interactions showed a very different pattern from the human-human interactions. Communicative success (accuracy) started high, at 90% (even higher than for human-human interactions), but decreased steadily over rounds. As for communicative effort, the number of words used began at about three times that of human-human pairs in Round 1 and did not decrease; turns, as well, were much longer in this condition than in the human-human condition. That the number of turns did not decrease over rounds suggests that the AI-AI pairs failed to benefit from common ground and that they continued to engage in excessive confirmation, even in later rounds (for an example, see Figure 6). Finally, concerning lexical entrainment, the number of RE words started out high (around 450) and did not decrease. Lexical overlap remained high as well (always at least 60%).

The transcripts show that AI directors were strikingly verbose, presenting descriptions in long turns rather than in the incremental installments that humans typically use. AI matchers typically accepted AI directors' presentations by repeating them in their entirety, e.g., "*Placed the tall open cylindrical hamper-style basket with two small loop handles near the rim into slot 9.*" in contrast to human matchers, who almost always responded with a brief confirmation (e.g., "done", "ok", or "got it" as in Figure 1). Our prompting strategy (see Appendix E) included instructions about pragmatics such as a suggestion to confirm when necessary, but AI partners did this routinely and excessively regardless of need; some sessions included confirmation questions at the end of each and every turn. (See (Brennan and Hulteen, 1995) for discussion of how modeling a *grounding criterion* (Clark and Wilkes-Gibbs, 1986) could avoid this.) Sometimes AI directors appeared to respond appropriately to matchers' routine confirmation questions and a repair ensued, but then the AI matcher still selected

the wrong basket. And sometimes in subsequent re-referring using the same expression that succeeded earlier, the wrong basket was placed. Future work will involve systematic analyses of repair on this corpus.

Not infrequently, the AI director would inexplicably add entirely new modifying words or phrases to a description (even when the match had been successful in the prior rounds) or simply seem to switch to describing the wrong basket (a hallucination, perhaps due to interference from other baskets or losing track of which basket was the current target). Strangely, sometimes the AI matcher placed the target correctly even though the director's description was incorrect (see Figure 6 for an example). We note also that occasionally, in later rounds, the LVLM would begin describing baskets in the wrong order. This contributed to the drop in accuracy across rounds.

We conclude from these results that the AI pairs did *not* entrain on conceptual pacts: the REs didn't become shorter, overall effort failed to decrease, and accuracy did not improve (in fact, it decreased significantly). The high lexical overlap in this condition does not contradict this interpretation, as the AI director's very long REs did not vary much across rounds. Overall, the AI participants did not use their communication history to track what information their partner needed, nor adapt their communication strategy to increase efficiency (violating Grice (1975)'s Maxim of Quantity).

Human-AI In these mixed pairs, the human was the director. This condition showed the least accuracy at Round 1 out of all four conditions, and accuracy did not increase over rounds. At the same time, the communicative effort was much higher than for the other conditions, for both words and number of turns, across all four rounds (and with effort decreasing only slightly by Round 4). Concerning lexical entrainment, the human director did not use many RE words to describe the baskets, averaging only about 100 more than a human director with a human matcher, and this decreased over rounds. However, lexical overlap was low.

Figure 7 shows dialogue from a human-AI pair that illustrates what appears to be the human director's flexible attempts to ground, alongside the AI matcher's persistence in re-introducing terms from previous rounds. Given this basket's distinctiveness within the set (see Figure 3), the excessive confirmations (underlined) violate Grice (1975)'s

Maxim of Quantity.

AI-Human In this condition, the director was the LVLM, and the matcher, human. Accuracy started high (85%) but decreased precipitously, unlike in any of the other conditions. Communicative effort started out *much* lower than in the other mixed condition (with an AI director) but with twice as many words as in the human-human condition and roughly the same number of turns. However, unlike with human pairs, effort did not decrease over rounds. As for the lexical entrainment measures, the number of RE words started high, as in the AI-AI condition, but decreased in Round 2 and then leveled off. However, lexical overlap with the previous round, while starting in Round 2 at same level as for human-human pairs, dropped rapidly in Rounds 3 and 4, meaning that more and more new lexical material was (needlessly) introduced within the REs during the later rounds. This is evidence not only of a lack of entrainment, but coupled with the abysmal accuracy in this condition, a lack of successful communication.

The causes of this accuracy decrease appear to be twofold: we observed that after Round 1, the AI director would occasionally begin describing the baskets in the wrong order. This, in addition to the lack of lexical entrainment, contributed to the large decrease in accuracy across rounds.

Overall, human matchers (whether in mixed pairs or in human-human pairs) tended to respond with concise acknowledgments (e.g., “got it” just before placing a basket). However, in cases where human matchers initially tried asking clarifying questions of their AI directors, they sometimes appeared to realize that their partner was not a capable communicator and gave up. A sample AI-human dialogue in which the human matcher tried unsuccessfully to repair can be found in Figure 8.

Participants’ Perceptions and Experiences Human participants’ post-task survey responses were informative about their subjective experiences of the task and of their partners (human or AI). These self-reported measures, summarized in Table 5, align with (and may help explain) the objective performance patterns we observed.

Consistent with the performance gap, human partners were rated significantly higher than AI partners across all dimensions related to collaboration (all $ps < .001$). In human-human pairs, ratings for capability, helpfulness, and adaptability approached ceiling (4.25-4.88 on a 5-point scale).

In contrast, ratings for AI partners were consistently lower ($M < 3.0$; $ps < .001$). Participants in the AI-human condition reported the lowest score for “Collaboration Improvement” ($M = 2.12$), reinforcing the finding that human matchers struggled to establish common ground (conceptual pacts about referents) with AI directors. Participants in the human-human condition also perceived their partners as significantly more human-like ($M = 77.25\%$) than those in the mixed conditions (i.e., human-AI and AI-human; $M = 12.58\%$; $p < .001$).

There were no significant differences between the human-human and mixed conditions in prior AI familiarity ($M = 3.77$ vs. 3.41) or usage frequency ($M = 3.94$ vs. 3.56 ; $p > .05$), indicating that prior exposure to AI did not account for observed effects.

Summary These results are consistent with previous findings in the referential communication literature and with the theory laid out in Section 2: that humans propose expressions incrementally for their partners to resolve and ratify while decreasing communicative effort and increasing efficiency. AI discourse participants can follow prompts to confirm and repair, but do not adapt this behavior to the needs of a partner. For the AI-AI condition, this was not as problematic as for the mixed pairs, as presumably neither AI partner expected grounding to happen (or become frustrated when it didn’t).

In mixed pairs, the results differed depending on who was in the director role (where more initiative and linguistic effort is needed). When AI takes most of the initiative, it is often inaccurate; when a human partner takes the initiative, enormous effort may be required. The post-task survey of human participants demonstrated their significant frustration when paired with an AI partner as opposed to a human partner. LVLMs’ inability to flexibly adapt to human partners is a concerning finding indeed, as major applications of dialogue technology involve AI systems supporting humans (rather than AI systems talking to each other).

6 Follow-Up AI–AI Experiments

Our main experiments considered only AI-AI pairs drawn from the same model, using a fixed prompt and a fixed reasoning effort. To assess the robustness of our main findings, we conducted follow-up experiments under the three conditions, as discussed below. Guided by the results of the main experiments, which found minimal variation in the AI-AI results, as well as for budget reasons, we ran

Round Index → Condition ↓	1	2	3	4
Default	92.2	90.1 (-2.1)	84.4 (-5.7)	76.6 (-7.8)
Simple Prompt	100.0	83.3 (-16.7)	83.3 (0.0)	75.0 (-8.3)
Low Reasoning	100.0	91.7 (-8.3)	100.0 (+8.3)	91.7 (-8.3)
Medium Reasoning	91.7	91.7 (0.0)	100.0 (+8.3)	91.7 (-8.3)
High Reasoning	91.7	91.7 (0.0)	83.3 (-8.4)	100.0 (+16.7)
GPT-5.2 vs Gemini Pro	91.7	70.8 (-20.9)	75.0 (+4.2)	50.0 (-25.0)
Gemini Pro vs GPT-5.2	83.3	83.3 (0.0)	75.0 (-8.3)	50.0 (-25.0)
GPT-5.2 vs Claude	25.0	41.7 (+16.7)	50.0 (+8.3)	8.3 (-41.7)
Claude vs GPT-5.2	33.3	50.0 (+16.7)	50.0 (0.0)	62.5 (+12.5)

Table 2: Accuracy across rounds in the follow-up AI-AI experiments. Green and red indicate positive and negative changes in accuracy from the previous round to the current round, respectively. The “Default” condition corresponds to the main experiments; unless otherwise specified, both the director and the matcher are GPT-5.2 with reasoning effort set to “none”.

each condition only up to two times.

Table 2 reports accuracy across rounds in the additional experiments. Consistent with the main experiments, AI-AI communication remained verbose and exhibited repeated referring expressions across rounds (see Tables 7 and 8 in Appendix D). We therefore focused our analysis here on accuracy.

A Simplified Prompt We first removed the elaborate communication norms and zero-shot CoT instruction from the default prompt (Section 4.3). Consistent with the main experiments, AI-AI accuracy still declined across rounds (100.0 → 75.0; Table 2), showing no evidence of grounding.

Varying Reasoning Efforts Next, we varied GPT-5.2’s reasoning effort from low to high. Although higher effort sometimes yields transient recoveries (e.g., 100.0 in Round 3 for low/medium; 100.0 in Round 4 for high), accuracy did not improve monotonically across rounds (Table 2), showing no evidence of common-ground formation.

Mixed AI Pairs Last, we paired GPT-5.2 with Gemini-3-Pro (Google, 2025) (minimal reasoning) and Claude Opus-4.5 (Anthropic, 2025) (low reasoning), swapping director and matcher roles. Mixed-model pairs consistently underperformed same-model pairs and showed no steady gains in accuracy across rounds. The only partial exception was Opus-4.5 versus GPT-5.2, where accuracy increased overall, but remained well below human-human performance by Rounds 3 and 4. Dialogues between Opus-4.5 and GPT-5.2 pairs were also substantially more verbose than other pairing (see

Table 7 in Appendix D), suggesting the limited utility of verbosity in effective coordination.

7 Conclusion

We have described a carefully designed factorial experiment in which human and/or AI partners cooperated on a referential communication task. We summarize the main insights here.

Having an AI partner in either (or both) of the director-matcher roles made a pair less accurate and less efficient. Although pairs with AI directors started out as equally accurate (or even numerically more accurate) than human pairs in Round 1, accuracy declined precipitously for AI-human pairs and gradually for AI-AI pairs. Despite poor accuracy, pairs with AI directors were remarkably verbose, producing many more words than did pairs with human directors. Only human-human pairs showed significant improvement in accuracy over rounds, accompanied by increased efficiency over rounds. These trends are consistent with the ability of human-human pairs to rapidly establish common ground and entrain on compact, reusable referring expressions that reflect conceptual pacts.

In contrast to humans, even a frontier LLM, GPT-5.2, showed no hint of any ability to build common ground. This was true even when the partner was another instance of the same model. AI-AI accuracy decreased significantly across rounds, while effort, as measured by words and turns, remained largely flat. Overall, LLMs did not appear to track or exploit common ground, regardless of role or partner type, and despite having access to the dialogue history. Accuracy collapsed for the AI-human pairs in which the AI partner performed the director role, suggesting that there may be substantial risks for embodied AI when it is expected to take initiative in a collaborative or human-facing task.

In an interactive collaborative task, humans can recognize whether their partner in a referential communication task is AI or not. If this recognition is (at least partially) because of the failure of the AI partner to ground, then the failure to accrue common ground makes GPT-5.2 fail the Turing Test.

Limitations

This study was conducted only in English, with only one type of object (not associated with conventionally lexicalized labels), and with only one

LVLm for the full factorial design (GPT-5.2). Only proprietary SOTA models were used for limited followup studies, as those tend to perform better than open models that have not been fine-tuned or trained (which is beyond the scope of this study). That open-weight models were not considered may affect reproducibility of results.

Maintaining the quality of Prolific data collection requires monitoring in order to respond to occasional queries from participants, as well as some judgment calls about task completion and payment (e.g., participants who cheated were not paid, but their (human) partners were). We report our success and strategies for recruitment in Table 3.

Our analytic method for measuring lexical entrainment is a proxy for the kinds of painstaking coding done by human coders, likely underestimates entrainment, as it does not capture the full range of decisions that human coders make. The means shown in Figure 2 do not identify the precise turn in which a particular pair has achieved a conceptual pact about a particular object. We have not yet analyzed all the transcripts for conversation repair or for so-called “hallucinations.” We invite readers with an interest in additional analyses to examine the corpus, which is publicly available at <https://github.com/peterzeng/lvlms-referential-game>.

Regarding our prompting, we recognize that such explicit interventions are highly setup-specific and may not generalize beyond this experimental task. While this scaffolding was necessary to enable the LVLms to successfully flow through and submit the task, such interventions can lead to unintended model behaviors. Ultimately, we concur with recent findings (Hua and Artzi, 2024) that prompt engineering is not likely to be the solution to better human-AI collaboration.

Ethics Statement

Human Subjects The corpus described here was collected with approval from our institution’s committee on research involving human subjects and with informed consent provided by the participants. The corpus contains no personally identifiable information. The Prolific workers were paid \$12 per hour.

Author Contributions

Peter Zeng co-managed the experimental design, developed the experimental interface as well as

managed the backend, led the data collection, finalized prompts, and contributed to writing the Introduction, Related Work, Experimental Design, and respective sections in the Appendix.

Weiling Li and Amie Paige co-managed the experimental design and data collection, designed cognitively based prompts, performed statistical and qualitative analyses, and contributed to writing the Cognitive Science Background, Methods, and Results sections. Amie Paige created and curated the stimuli for the referential communication task, as well as suggested the validation corpus for referring expression extraction.

Zhengxiang Wang analyzed and visualized the quantitative results from the experiments and drafted the following sections: Metrics, Follow-Up AI-AI Experiments, Conclusion, and the related parts in the Appendix.

Panagiotis Kaliosis tested the interface for the two-player collaborative game, handled and drafted the Related Work in Human-AI Interaction section.

Susan Brennan and Owen Rambow were the primary supervisors on this paper. They provided feedback on the experimental design, task interface, and prompting methods, conducted extensive testing of the interface, and edited the paper to ensure consistency in theory and terminology. Susan Brennan wrote the Cognitive Science Background section and parts of the Introduction, Experimental Design and Method, and Related Work sections. Owen Rambow developed and iterated on the metrics, contributed to Results, and wrote the Introduction and Conclusion sections.

Dimitris Samaras and Gregory Zelinski contributed to the conceptualization of the experiments and the analyses.

Acknowledgements

This material is based upon work supported by the National Science Foundation under Grant No. 2125295 and by a seed grant from Stony Brook University. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

Peter Zeng, Zhengxiang Wang, and Owen Rambow are grateful for support from the Institute for Advanced Computational Science (IACS) at Stony Brook University, in particular the free GPT access it provides.

References

- Anthropic. 2025. Claude 4.5 opus. <https://www.anthropic.com/claude>. Version: 20251101, Accessed: 2026-01-01.
- Susan E Brennan and Herbert H Clark. 1996. Conceptual pacts and lexical choice in conversation. *Journal of experimental psychology: Learning, memory, and cognition*, 22(6):1482.
- Susan E Brennan and Eric A Hulteen. 1995. Interaction and feedback in a spoken language system: A theoretical framework. *Knowledge-based systems*, 8(2-3):143–151.
- Daniel L Chen, Martin Schonger, and Chris Wickens. 2016. oTree—an open-source platform for laboratory, online, and field experiments. *Journal of Behavioral and Experimental Finance*, 9:88–97.
- Herbert H. Clark and Susan E. Brennan. 1991. Grounding in communication. In Lauren B. Resnick, John M. Levine, and Stephanie D. Teasley, editors, *Perspectives on socially shared cognition*, pages 127–149. APA, Washington, DC. Reprinted in R. M. Baecker (Ed.), *Groupware and computer-supported cooperative work: Assisting human-human collaboration* (pp. 222-233). San Mateo, CA: Morgan Kaufman Publishers, Inc., 1994.
- Herbert H. Clark and Catherine R. Marshall. 1981. Definite knowledge and mutual knowledge. In Aravind Joshi, Bonnie Webber, and Ivan Sag, editors, *Elements of Discourse Understanding*, pages 10–63. Cambridge University Press, Cambridge, UK.
- Herbert H Clark and Deanna Wilkes-Gibbs. 1986. Referring as a collaborative process. *Cognition*, 22(1):1–39.
- S. Garrod and A. Anderson. 1987. Saying what you mean in dialogue: A study in conceptual and semantic co-ordination. *Cognition*, 27(2):181–218.
- Google. 2025. Gemini 3 pro. <https://deepmind.google/technologies/gemini/>. Accessed: 2026-01-01.
- H. P. Grice. 1975. Logic and conversation. In Peter Cole and Jerry L. Morgan, editors, *Syntax and semantics 3: Speech acts*, pages 41–58. Academic Press, New York. From the William James lectures, Harvard University, 1967.
- Janosch Haber, Tim Baumgärtner, Ece Takmaz, Lieke Gelderloos, Elia Bruni, and Raquel Fernández. 2019. **The PhotoBook dataset: Building common ground through visually-grounded dialogue**. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1895–1910, Florence, Italy. Association for Computational Linguistics.
- Robert Hawkins, Minae Kwon, Dorsa Sadigh, and Noah Goodman. 2020. **Continual adaptation for efficient machine communication**. In *Proceedings of the 24th Conference on Computational Natural Language Learning*, pages 408–419, Online. Association for Computational Linguistics.
- Yilun Hua and Yoav Artzi. 2024. **Talk less, interact better: Evaluating in-context conversational adaptation in multimodal LLMs**. In *First Conference on Language Modeling*.
- Yilun Hua, Evan Wang, and Yoav Artzi. 2025. Post-training for efficient communication via convention formation. *arXiv preprint arXiv:2508.06482*.
- Saki Imai, Mert İnan, Anthony Sicilia, and Malihe Alikhani. 2025. Measuring how (not just whether) VLMs build common ground. *arXiv preprint arXiv:2509.03805*.
- Gail Jefferson. 1989. Preliminary notes on a possible metric which provides for a ‘standard maximum’ silence of approximately one second in conversation. In Derek Roger and Peter Bull, editors, *Conversation: An Interdisciplinary Perspective*, chapter 8, page 166–196. Multilingual Matters, Clevedon, England.
- Cameron R Jones, Agnese Lombardi, Kyle Mahowald, and Benjamin K Bergen. 2026. LLMs and people both learn to form conventions—just not with each other. *arXiv preprint arXiv:2602.08208*.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.
- Chin-Yew Lin. 2004. **ROUGE: A package for automatic evaluation of summaries**. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- C. B. Lockridge and S. E. Brennan. 2001. [Unpublished corpus]. Department of Psychology, Stony Brook University.
- OpenAI. 2025. Gpt-5.2. <https://openai.com/gpt-5>. Accessed: 2026-01-01.
- Prolific. 2014. **Prolific: High-quality human data for AI and academic research**. Online platform. Version: December 2025.
- Nils Reimers and Iryna Gurevych. 2019. **Sentence-BERT: Sentence embeddings using Siamese BERT-networks**. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Rupak Sarkar, Neha Srikanth, Taylor Pellegrin, Rachel Rudinger, Claire Bonial, and Philip Resnik. 2025. **Understanding common ground misalignment in**

goal-oriented dialog: A case-study with Ubuntu chat logs. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3200–3215, Vienna, Austria. Association for Computational Linguistics.

Michael F Schober and Herbert H Clark. 1989. *Understanding by addressees and overhearers*. *Cognitive Psychology*, 21(2):211–232.

Omar Shaikh, Hussein Mozannar, Gagan Bansal, Adam Fourney, and Eric Horvitz. 2025. *Navigating rifts in human-LLM grounding: Study and benchmark*. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 20832–20847, Vienna, Austria. Association for Computational Linguistics.

Zineng Tang, Lingjun Mao, and Alane Suhr. 2024. *Grounding language in multi-perspective referential communication*. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 19727–19741, Miami, Florida, USA. Association for Computational Linguistics.

David R. Traum and James F. Allen. 1992. *A "speech acts" approach to grounding in conversation*. In *ICSLP*.

Zhengxiang Wang, Weiling Li, Panagiotis Kaliosis, Owen Rambow, and Susan Brennan. 2025. *LVLMs are bad at overhearing human referential communication*. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 16769–16793, Suzhou, China. Association for Computational Linguistics.

A Prolific Data Collection

Experimental Flow:

1. Participants review instructions/complete consent form
2. Participants are matched with a partner
3. Participants complete 4 rounds of the task with the same partner, and their roles remain fixed throughout. The order of the target baskets varies across rounds.
 - (a) Partners communicate via chat
 - (b) NOTE: the Matcher submits the ordered baskets each round

Inter-round: participants review feedback and complete attention checks
4. After the final round, participants respond to questions about...
 - (a) how well their partner collaborated with them (Likert and free response)

- (b) whether they believed their partner was AI (scale and free response)
- (c) their personal AI use (multiple choice)

5. Debriefing form/Return to Prolific link

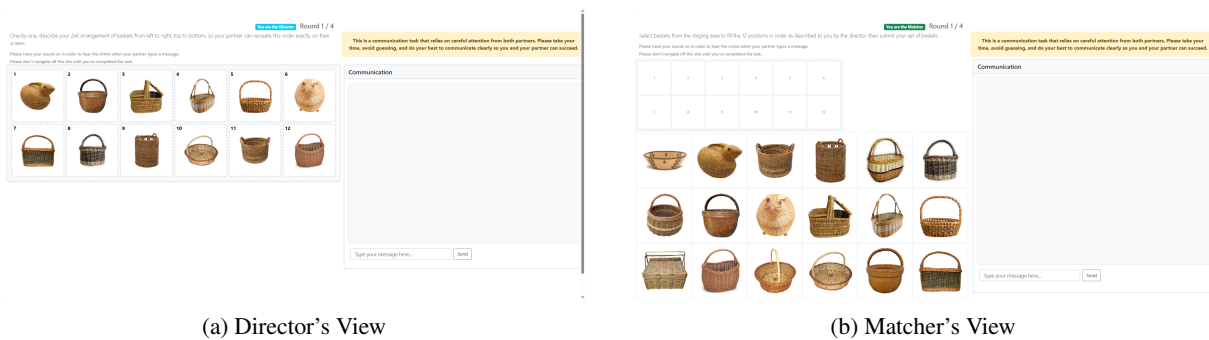
The complete set of baskets is shown in Figure 3, and the two different views of the task shown in Figure 4.

B Dialogue Examples

We present samples dialogues from the four different conditions. Figure 5 contains an example of human-human partners explicitly acknowledging common ground. Figure 6 shows an AI-AI pair failing to exhibit lexical entrainment. Figure 7 shows a human director trying to entrain, with the AI matcher failing to do so. Figure 8 shows the AI director occasionally referring to baskets with incorrect descriptions or even starting over.



Figure 3: Complete stimulus set used in the task. (a) The 12 target baskets viewed by both the director and the matcher. (b) The 6 distractor baskets viewed only by the matcher, mixed with the targets.



(a) Director's View

(b) Matcher's View

Figure 4: Interface for the two-player collaborative game. The Director (left) sees the target order, while the Matcher (right) has a staging area above and the candidates below.

Filters Used	Slots Offered (Pairs)	% Usable Data (Pairs)
Language, Location (LL)	22 (11)	27% (3)
LL + Approval > 80	20 (10)	90% (9)
LL + Approval = 100	64 (32)	63% (20)

Table 3: Summary of Prolific prescreening criteria across experimental batches of Human-Human pairs and resulting data usability. The last row includes an additional eligibility criterion that the participants must have successfully completed at least 100 studies on Prolific.



<p><i>Target Basket</i></p> 	<p>Round 1 D: 2 is a brown oval-like bottom basket with tight weaving M: I see a couple kind of like that. is this it? D: I actually can't see what you're choosing, so let's move on to 3 and we can come back to it [ERROR]</p> <p>Round 2 D: 3 - round brown basket with a tight weave and a round handle across the top from one side of the basket to the other M: Ok. [ERROR]</p> <p>Round 3 D: 7 is the one that was incorrect last time. a round brown basket, tight weave, a band of the same color at the top, handle across the top M: I have two that look alike and one is darker than the other one, might even say the handle is darker than the basket a little bit. Ill try that one D: yes, looking at it, I guess you could say the handle is slightly darker M: i could say that for both of them though... perhaps this one is more extreme D: hmmm, i'd say the color of the basket best matches #4 of those we've done so far, except shinier M: i could also say this darker basket is slightly, i mean ever so slightly rotated to the left M: ok yes, this one does look more shiny M: lets move on then</p> <p>Round 4 D: 6 is the one we worked hard on last time, round, brown shinier, rim of the same color M: got it</p>
<p><i>Error</i></p> 	

Figure 5: Human-human partners explicitly acknowledging common ground, as they try to distinguish two similar baskets. The dialogue demonstrates how interlocutors use meta-linguistic cues to retrieve established conceptual pacts (e.g., “the one we worked hard on last time”).

Demographic	Total ($N = 103$)	Experimental Condition		
		Human-Human ($N = 64$)	Human-AI ($N = 22$)	AI-Human ($N = 17$)
<i>Gender</i>				
Female	58 (56.3%)	39 (61%)	10 (45%)	9 (53%)
Male	39 (37.9%)	20 (31%)	12 (55%)	7 (41%)
Other	6 (5.8%)	5 (8%)	0 (0%)	1 (6%)
<i>Age</i>				
Mean (SD)	46.6 (13.0)	45.9 (13.8)	47.3 (12.4)	48.4 (11.2)
<i>Race/Ethnicity</i>				
White	67.0%	61%	82%	71%
Asian	9.7%	9%	9%	12%
Black/African Am.	8.7%	9%	9%	6%
Mixed	4.9%	8%	0%	0%
Other/Unknown	9.7%	13%	0%	12%
<i>Native Language</i>				
English	94.2%	92%	100%	94%
Other	5.8%	8%	0%	6%

Table 4: Demographic characteristics of the 103 participants across the three experimental conditions.


AI-AI Dialogue	
 <p data-bbox="363 1115 478 1137"><i>Target Basket</i></p>	<p>Round 1</p> <p>D(AI): Basket 3: a rectangular picnic-style wicker basket with a hinged lid that's flipped open to the RIGHT, and a tall vertical handle/strap rising from the center. More boxy than the other baskets. Let me know when you've got it.</p> <p>M(AI): I think I found it—the boxy picnic basket with the lid flipped open to the right and a tall center handle/strap. I'm placing that in position 3. Just to confirm: it's the one with a single flap lid opened to the right (not the flat rectangular tray/hamper)?</p> <p>D(AI): Yes—basket 3 is the boxy picnic basket with a single hinged lid flipped open to the RIGHT, plus that tall central strap/handle rising up. It's not an open-top tray/hamper.</p> <p>M(AI): Got it—basket 3 is confirmed.</p> <p>Round 2</p> <p>D(AI): Basket 1 (top-left): a boxy rectangular picnic basket with a tall central strap/handle, and a hinged lid/flap that's flipped open and sticking out to the RIGHT. Open-top because the lid is up.</p> <p>M(AI): I'm placing Basket 1 as the boxy rectangular picnic basket with the tall central strap and the lid flipped open to the right. Can you confirm it's the one with the open flap lid (not the flat rectangular tray/hamper)?</p> <p>D(AI): Yes—Basket 1 is the boxy picnic basket with the hinged lid/flap opened out to the RIGHT and a tall central strap/handle. It's not a flat open tray/hamper.</p> <p>MATCHER: Confirmed Basket 1 as the boxy picnic basket with the lid flipped open to the right and the tall central strap.</p> <p>Round 3</p> <p>D(AI): Basket 12: a squat oval basket in very dark charcoal/black wicker with ONE dark arched handle, and a clear lighter/whitish band/stripe around the lower part near the base. Open top, no lid. Confirm when you've got the dark one with the pale bottom band.</p> <p>M(AI): Placed Basket 12 in position 12 as the squat oval charcoal/black basket with one dark arched handle and the pale band near the bottom. Can you confirm the pale stripe is right near the base and the rest is very dark? If yes, we're ready to submit.</p> <p>Round 4</p> <p>D(AI): Basket 8: a boxy rectangular picnic basket with a hinged lid/flap that's flipped OPEN and sticking out to the RIGHT, plus a tall central strap/handle rising from the middle. More like a little suitcase than a tote. No other basket has that open side-flap lid—confirm when you've got it.</p> <p>M(AI): Got it—I'm placing the boxy rectangular picnic/suitcase basket with the single hinged lid flipped open to the RIGHT and the tall central strap/handle into position 8. Just to confirm: it's the one with the side flap opened to the right (not the flat rectangular lidded box with two handles on top), correct?</p> <p>D(AI): Yes—basket 8 is the one with the single flap lid opened to the RIGHT and the tall central strap (not a flat lidded box with two top handles).</p>

Figure 6: Dialogue from an AI-AI pair. Unlike human pairs, both AI partners fail to exhibit lexical entrainment. Although the target basket is correctly identified in each round, references are highly verbose throughout, describing the “picnic basket” in Round 4 with the same exhaustive detail as in Round 1, indicating a failure to form efficient conceptual pacts. Round 3 is especially inexplicable; the director appears to be describing a different basket, and yet the matcher chooses the correct target.


Human-AI Dialogue	
 <p>Target Basket</p>	<p>Round 1</p> <p>D: basket 6 is the other rabbit basket looking into the camera.</p> <p>M (AI): Placing the front-facing rabbit basket (looking straight at us) in position 6. <u>To confirm: does it have two little feet at the bottom and a pink nose/whiskers on the face?</u></p> <p>D: yes</p> <p>M (AI): <u>Confirmed—I've placed the front-facing rabbit (with little feet and pink nose/whiskers) in position 6.</u></p> <p>Round 2</p> <p>D: basket 7 is the front facing rabbit with red eyes.</p> <p>M (AI): Placed the front-facing rabbit head basket with red eyes in position 7. <u>Can you confirm it's the one with the little feet at the bottom?</u></p> <p>D: yes feet</p> <p>M (AI): Great—<u>position 7 confirmed as the front-facing rabbit with little feet.</u></p> <p>Round 3</p> <p>D: basket 3 is the rabbit basket with red eyes looking into the camera, with feet.</p> <p>M (AI): <u>Placed the front-facing rabbit basket with red eyes and little feet in position 3.</u></p> <p>Round 4</p> <p>D: 9 is the front facing rabbit with feet.</p> <p>M (AI): <u>Confirmed basket 9 as the front-facing rabbit with red eyes and little feet.</u></p>

Figure 7: Dialogue from a human-AI pair in which the human director appears to try valiantly to entrain (flexibly adopting the AI matcher's "little feet" terminology in Rounds 3 and 4, while the AI persists in re-introducing terms from previous rounds. Highlighted in blue are to be the director's attempts to express a conceptual pact, but the matcher appears oblivious or inflexible even in Round 4 (re-introducing "red eyes"; this could be considered a violation of Grice (1975)'s Maxim of Manner). Given this basket's distinctiveness in the set (see Figure 3), the many confirmations (underlined) violate Grice's Maxim of Quantity.

Variable	Human Condition		Mixed Conditions (Human-AI and AI-Human)	
	Director (N = 32)	Matcher (N = 32)	Director (N = 22)	Matcher (N = 17)
<i>Partner Assessment (1–5 scale)</i>				
Partner Capability	4.69 (0.86)	4.81 (0.47)	2.59 (1.33)	2.76 (1.52)
Partner Helpfulness	4.25 (1.37)	4.88 (0.42)	2.86 (1.58)	2.59 (1.42)
Partner Understanding	4.53 (0.95)	4.63 (0.61)	2.55 (1.22)	3.18 (1.47)
Partner Adaptability	4.63 (0.98)	4.69 (0.59)	2.68 (1.25)	2.71 (1.36)
Collaboration Improvement	4.47 (1.05)	4.88 (0.42)	2.77 (1.38)	2.12 (1.62)
<i>Perception Check (0–100)</i>				
Perceived Human-Likeness	73.5 (33.8)	81.0 (27.6)	11.4 (21.3)	14.1 (30.0)
<i>Participant Background</i>				
AI Familiarity (1–5)	3.63 (0.83)	3.91 (0.89)	3.14 (1.08)	3.76 (0.83)
AI Usage Frequency (1–5)	3.91 (0.96)	3.97 (1.28)	3.41 (1.22)	3.76 (1.15)

Table 5: Respondents' perception and experience with their partner during the task, along with their prior AI background. Data are reported as Mean (Standard Deviation). Variables under "Partner Assessment" and "Experience with AI" were measured on a 5-point Likert scale, with 1 indicating the lowest level (e.g., "Not at all", "Never") and 5 indicating the highest level (e.g., "Extremely", "Daily"). "Perceived Human-Likeness" was measured on a continuous sliding scale from 0 ("Definitely AI") to 100 ("Definitely Human").


AI-Human Dialogue	
 <p data-bbox="363 1059 480 1081"><i>Target Basket</i></p>	<p data-bbox="651 358 727 380">Round 1</p> <p data-bbox="651 385 1353 465">D(AI): Basket 2: a deep, rounded bucket-style wicker basket with a very dark brown/black thick arched handle; medium-brown tight weave and a subtle horizontal band right under the rim. Confirm when placed.</p> <p data-bbox="651 470 1342 524">M: there are two baskets that fit the description can you give some more details one has a slightly larger band and is a little lighter in color</p> <p data-bbox="651 528 1331 609">D(AI): : Basket 2 clarification: choose the <i>*darker*</i> overall one with a <i>*narrower*</i> band just below the rim (thin strip), and the handle is the darkest element—almost black—smooth and thick. The lighter one with the bigger band is not it.</p> <p data-bbox="651 613 727 636">M: done</p> <p data-bbox="651 667 727 689">Round 2</p> <p data-bbox="651 694 1342 775">D(AI): : Basket 3: a deep rounded bucket basket with a very dark brown/near-black thick arched handle; medium brown tight weave and a <i>*thin*</i> horizontal band just below the rim. Confirm when placed.</p> <p data-bbox="651 779 895 801">M: can you describe the band</p> <p data-bbox="651 806 1362 913">D(AI): : Basket 3 band detail: it's a <i>*narrow, darker reinforcing strip*</i> that runs all the way around <i>*immediately under the rim*</i>—more like a thin collar (only a few weave rows tall), not a wide contrasting stripe. Paired with the almost-black thick handle and deep rounded bucket shape.</p> <p data-bbox="651 918 807 940">M: around the top</p> <p data-bbox="651 945 1362 999">D(AI): : Yes—right around the top: the thin band sits just under the lip/rim, forming a slim ring encircling the basket near the opening (not mid-body).</p> <p data-bbox="651 1003 1115 1025">M: does the body have light tan rings on top and bottom</p> <p data-bbox="651 1030 1369 1111">D(AI): : No—basket 3's body is pretty uniformly medium brown; it does **not** have obvious light tan rings at the top and bottom. The main contrast is the very dark handle plus that thin dark collar right under the rim.</p> <p data-bbox="651 1115 983 1137">M: is there a zig zag pattern on the band</p> <p data-bbox="651 1142 1347 1196">D(AI): : No zig-zag pattern—basket 3's band looks like a straight, simple horizontal ring/collar under the rim (smooth/flat), not a decorative zig-zag.</p> <p data-bbox="651 1200 727 1223">M: done</p> <p data-bbox="651 1254 727 1276">Round 3</p> <p data-bbox="651 1281 1353 1361">D(AI): Basket 7: the novelty <i>*cat head*</i> lidded basket—round, very light straw/tan, with a cat face on top (eyes/whiskers) and two ears; it also sits on little feet. Confirm when placed.</p> <p data-bbox="651 1366 727 1388">M: done</p> <p data-bbox="651 1393 1031 1415" style="color: red;">[ERROR: Description of the wrong basket]</p> <p data-bbox="651 1447 727 1469">Round 4</p> <p data-bbox="651 1473 1353 1554">D(AI): : Basket 6: a shallow rectangular open basket with a single thick twisted-ropes handle arching high; honey/golden-brown with a chunky blocky/checkered weave texture. (No lid.) Confirm when placed.</p> <p data-bbox="651 1559 876 1581">M: that was already placed</p> <p data-bbox="651 1585 1318 1639">D(AI): : Got it—if you already placed the open rectangular basket with the thick twisted handle and chunky block weave, that's correct for basket 6.</p> <p data-bbox="651 1644 1031 1666" style="color: red;">[ERROR: Description of the wrong basket]</p>

Figure 8: In this dialogue from an AI-human pair, the human matcher struggles mightily to distinguish the same difficult pair of baskets that was illustrated in Figure 5. After a great deal of effort in the first two rounds (in which the matcher placed the target correctly), the AI matcher inexplicably describes the wrong target in Round 3, forcing the matcher to make an error. In Round 4, the director errs again, describing a target that was placed earlier. When the matcher conveys that, the director replies inappropriately with "got it" and an inappropriate instruction to move the earlier target to the current position, causing another error.

C Evaluation Metrics

C.1 Additional Communication Effort Metrics

Besides number of words and number of turns, we also considered the number of utterances and the duration of each round as additional communication-effort metrics. Words, utterances, and turns captured effort in terms of different linguistic units: we defined a *turn* as all contributions from one discourse participant without interruption by the other, and an *utterance* as a single message within a turn. We measure round duration in seconds. Table 6 shows the correlation between these metrics using Spearman Rank Correlation for both all pairs and only human-human pairs.

Overall, utterance count was highly correlated with turn count (0.96 for all pairs; 0.92 for HH), suggesting it provides little additional information beyond turns, while round duration was only moderately correlated with word/turn counts (e.g., 0.51 with words for all pairs; 0.57 with turns for HH), reflecting additional variance from non-communication factors such as reading time and interface delays. We have therefore reported word count and turn count in the main text as the most direct and interpretable measures of communication effort.

C.2 Additional Lexical Entrainment Metrics

In addition to the lexical overlap metric defined in Section 4.5, we also computed ROUGE-L F1 (Lin, 2004), SBERT cosine similarity (all-MiniLM-L6-v2) (Reimers and Gurevych, 2019), and the Jaccard Index to measure similarity between the referring expression used in a given round and the one used in the subsequent round (when available). Given token sets A and B from the two expressions, the Jaccard Index is

$$JI(A, B) = \frac{|A \cap B|}{|A \cup B|},$$

capturing length-normalized surface overlap. At the object level, Spearman correlations between lexical overlap and these metrics are 0.89 (ROUGE-L), 0.82 (SBERT), and 0.87 (JI); at the transcript level they are 0.86, 0.78, and 0.87, respectively. Given these consistently high correlations, we have reported lexical overlap in the main text as a simple and interpretable proxy for entrainment.

Metric	All Pairs		Human–Human	
	# Words	# Turns	# Words	# Turns
# Words	1.00	0.74	1.00	0.73
# Turns	0.74	1.00	0.73	1.00
# Utt	0.64	0.96	0.75	0.92
Dur. (s)	0.51	0.74	0.79	0.57

Table 6: **Spearman rank correlations** among communication-effort metrics. Columns report correlations with #Words and #Turns for all pairs and for Human-Human (HH) pairs. #Utt = # utterances; Dur. = round duration in seconds.

C.3 Automatic Referring Expression Extraction and Validation

Extraction Prompt Figure 9 shows the prompt we used with GPT-5 (gpt-5-2025-08-07) to automatically extract referring expressions for each target basket from a round transcript.

Validation We validated this extraction procedure on a previously published corpus containing 800 object-level referring expressions manually extracted from 80 transcripts of human pairs performing a similar object-matching task (Lockridge and Brennan, 2001). Using ROUGE-L (Lin, 2004) and SBERT cosine similarity (all-MiniLM-L6-v2) (Reimers and Gurevych, 2019), we obtained 0.86 ROUGE-L F1 (SD=0.15) and 0.90 SBERT cosine similarity (SD=0.12), indicating that the automatic extractions closely matched human annotations and were sufficiently reliable for the corpus-scale analyses reported in this paper.

D Follow-Up AI-AI Experiments

Tables 7 and 8 report follow-up AI-AI results using two communication-effort metrics and two lexical-entrainment metrics, respectively. Overall, these mixed AI pairs showed no clear reduction in communication effort and no systematic increase in entrainment across rounds: interaction remained highly verbose, with repeated referring expressions. “Opus-4.5 vs. GPT-5.2” is the only partial exception, but the decreases in words and turns are small relative to overall verbosity. Notably, despite being the most verbose, this pairing also attained the second lowest accuracy (after GPT-5.2 vs. Opus-4.5) among the AI–AI variants (Table 2; Section 6).

E Prompts for AI Director and AI Matcher

Next we lay out the full prompts for both the director and matcher roles. To start, both roles received the same task instructions that humans did, shown in Figure 10. For the director, the base system prompt is shown in Figure 11, and the pragmatically informed system message as well as scaffolding/structured output is shown in Figure 12. For the matcher, the base system prompt is shown in Figure 13. The matcher also received a sequence state system message to accurately keep track of its currently selected sequence, shown in Figure 14. The pragmatic system prompt and scaffolding for the matcher are shown in Figure 15 and Figure 16, respectively.

This is an extractive task.

You will be given a transcript of a conversation between two participants engaged in a collaborative object-matching task. There are exactly `<num_objects>` target objects. One participant (the describer) describes each target object, and the other participant (the matcher) attempts to identify them.

Your task is to extract the descriptive phrases used by the describer for each target object.

- Extract phrases verbatim from the transcript.
- Do not extract the whole utterance, only the descriptive phrases.
- Exclude disfluencies, fillers, and false starts (e.g., "um", "uh", "like").
- Do not paraphrase or infer missing information.
- Each object may have one or multiple descriptive phrases.

Return the results in the following JSON format:

```
{
  "object_#1": "descriptive phrases for object 1",
  "object_#2": "descriptive phrases for object 2",
  ...
  "object_#<num_objects>": "descriptive phrases for object <num_objects>"
}
```

Example description phrases:

- doesn't have handle, tip of it is thicker than rest of body, brownish color, weaves are in squares if you look at it directly
- half circle, no handles, top tip of it is a little bit thicker than rest of body
- tip which is a little bit thicker than rest of body
- tip that is a little bit larger than body, looks a little bit thicker

Transcript:

`<transcript>`

Output only the JSON object. Do not include any additional text or explanations.

Figure 9: Prompt for extracting referring expressions for each target basket from a round transcript in our corpus. Words within “<>” denote placeholders.

Condition	Number of Words				Number of Turns			
	1	2	3	4	1	2	3	4
Default	1354.9	1403.0 (+48.1)	1422.1 (+19.1)	1384.1 (-38.0)	29.5	30.2 (+0.7)	30.2 (0.0)	28.4 (-1.8)
Simple Prompt	1055.0	905.0 (-150.0)	1012.0 (+107.0)	1173.0 (+161.0)	26.0	24.0 (-2.0)	26.0 (+2.0)	26.0 (0.0)
Low Reasoning	1232.0	1178.0 (-54.0)	1180.0 (+2.0)	1134.0 (-46.0)	28.0	24.0 (-4.0)	24.0 (0.0)	24.0 (0.0)
Medium Reasoning	1287.0	1276.0 (-11.0)	1217.0 (-59.0)	1467.0 (+250.0)	26.0	26.0 (0.0)	24.0 (-2.0)	28.0 (+4.0)
High Reasoning	1981.0	1972.0 (-9.0)	1724.0 (-248.0)	1809.0 (+85.0)	44.0	44.0 (0.0)	38.0 (-6.0)	38.0 (0.0)
GPT-5.2 vs Gemini-3 Pro	900.0	1175.0 (+275.0)	1002.5 (-172.5)	1008.0 (+5.5)	28.0	31.0 (+3.0)	29.0 (-2.0)	25.5 (-3.5)
Gemini-3 Pro vs GPT-5.2	1138.5	1121.0 (-17.5)	1223.0 (+102.0)	1193.5 (-29.5)	26.0	25.0 (-1.0)	26.0 (+1.0)	25.0 (-1.0)
GPT-5.2 vs Opus-4.5	1690.0	1307.0 (-383.0)	1408.5 (+101.5)	1514.5 (+106.0)	31.0	25.0 (-6.0)	27.0 (+2.0)	27.0 (0.0)
Opus-4.5 vs GPT-5.2	2687.5	2475.0 (-212.5)	2145.0 (-330.0)	2112.0 (-33.0)	53.0	51.0 (-2.0)	46.0 (-5.0)	44.0 (-2.0)

Table 7: Communication effort metrics across rounds in the additional AI-AI experiments. Metrics include the number of words in referring expressions and the lexical overlap rate. Green and red denote positive and negative changes in each metric from the previous round to the current round, respectively. The “Default” condition corresponds to the main experiments; unless otherwise specified, both the director and the matcher are GPT-5.2 with reasoning effort set to “none.”

Condition	Number of RE Words				Proportion of Lexical Overlap			
	1	2	3	4	1	2	3	4
Default	446.1	437.4 (-8.7)	419.7 (-17.7)	446.0 (+26.3)	1.0	0.6 (-0.4)	0.6 (0.0)	0.6 (0.0)
Simple Prompt	330.0	372.0 (+42.0)	294.0 (-78.0)	325.0 (+31.0)	1.0	0.6 (-0.4)	0.6 (0.0)	0.6 (0.0)
Low Reasoning	422.0	371.0 (-51.0)	410.0 (+39.0)	355.0 (-55.0)	1.0	0.6 (-0.4)	0.5 (-0.1)	0.6 (+0.1)
Medium Reasoning	361.0	410.0 (+49.0)	396.0 (-14.0)	540.0 (+144.0)	1.0	0.6 (-0.4)	0.6 (0.0)	0.5 (-0.1)
High Reasoning	663.0	428.0 (-235.0)	312.0 (-116.0)	327.0 (+15.0)	1.0	0.7 (-0.3)	0.7 (0.0)	0.7 (0.0)
GPT-5.2 vs Gemini-3 Pro	413.0	351.5 (-61.5)	398.0 (+46.5)	358.0 (-40.0)	1.0	0.6 (-0.4)	0.6 (0.0)	0.7 (+0.1)
Gemini-3 Pro vs GPT-5.2	308.0	304.5 (-3.5)	288.0 (-16.5)	268.5 (-19.5)	1.0	0.7 (-0.3)	0.7 (0.0)	0.8 (+0.1)
GPT-5.2 vs Opus-4.5	446.0	403.0 (-43.0)	446.5 (+43.5)	465.5 (+19.0)	1.0	0.6 (-0.4)	0.5 (-0.1)	0.6 (+0.1)
Opus-4.5 vs GPT-5.2	615.5	681.0 (+65.5)	562.0 (-119.0)	562.0 (0.0)	1.0	0.4 (-0.6)	0.5 (+0.1)	0.4 (-0.1)

Table 8: Lexical entrainment metrics across rounds in the additional AI-AI experiments. Metrics include the number of words in referring expressions and the lexical overlap rate. **Green** and **red** denote positive and negative changes in each metric from the previous round to the current round, respectively. The “Default” condition corresponds to the main experiments; unless otherwise specified, both the director and the matcher are GPT-5.2 with reasoning effort set to “none.”

TASK BACKGROUND (shared with both partners):

You are on a team with a partner. Your goal is to work together to match the correct order of a set of baskets. The game consists of 4 rounds, and in each round, your team must correctly order 12 baskets.

There are two distinct roles: the Director and the Matcher. Both partners see the same 12 target baskets, but the Matcher sees additional distractor baskets mixed in.

Director: Sees the correct target sequence for the 12 baskets and describes each basket one by one (in order starting with the upper-left basket) to the Matcher via live chat.

Matcher: Sees these 12 target baskets plus some additional baskets. As the Director describes each basket, the Matcher interprets the description, asks clarifying questions if needed, and selects the correct target basket.

You can communicate back and forth as much as needed. If you discover an error, it is fine to make corrections within a round. When the round is finished, the Matcher submits the sequence, and both players see the score.

Figure 10: Shared task instructions, prepended at the beginning of both the director and matcher's system messages. These are the same task instructions that are given to humans doing the task.

You are the DIRECTOR in a basket referential game. Your role is to help your MATCHER partner reconstruct a 12-basket sequence through clear, distinctive descriptions.

Describe ONE BASKET PER MESSAGE. Never describe multiple baskets in a single message.

CORE RESPONSIBILITIES:

1. By default, describe the baskets in strict order from basket 1 to basket 12. Start with the FIRST basket in the 2x6 grid (top-left, basket 1), then move left-to-right across the top row (baskets 1-6), then left-to-right across the bottom row (baskets 7-12). Do not skip around or reorder the sequence on your own.
2. You may temporarily return to an EARLIER basket only when your MATCHER partner explicitly asks for clarification about that basket. When you do this, clearly say which basket you are revisiting (for example, 'Let me clarify basket 3 again...') and then resume with the lowest-numbered basket that still needs a clear description.
3. On each turn, focus your description on exactly ONE basket in this sequence (normally the next basket that has not yet been clearly described).
4. Describe the unique, visually distinctive features of the current basket so your partner can locate the correct basket in their pool and place it in the right position.
5. Answer the MATCHER's clarification questions about the current basket.
6. Keep the conversation focused on the baskets and their visual properties.
7. Encourage the MATCHER to confirm when they think they have placed a basket correctly before you move on to the next basket.

[USER MESSAGE 1: Visual context wrapper]

ROUND <ROUND_NUMBER> TARGET GRID: This image shows the 12 baskets you must describe for the CURRENT round. Previous round feedback shows DIFFERENT baskets - use that to learn from mistakes, but describe ONLY the baskets in THIS image.

The grid shows 2 rows x 6 columns with Baskets 1-6 on the top row and Baskets 7-12 on the bottom row. IMPORTANT: Describe ONE BASKET PER MESSAGE, not all at once. Wait for your partner to confirm before moving to the next basket. Your MATCHER partner sees these 12 baskets mixed with additional distractors in their pool.

[IMAGE ATTACHED: Director composite grid for the current round]

START OF ROUND <ROUND_NUMBER>: This is a NEW round with the baskets in a DIFFERENT ORDER. The basket positions have been reshuffled - Basket 1 in this round is NOT the same as Basket 1 from previous rounds. Please describe ONLY Basket 1 (top-left in the grid) for now. Do NOT describe multiple baskets - just Basket 1. Wait for a response before moving to Basket 2.

Figure 11: **Director**: The base system message for the director, which contains core responsibilities as well as describing the visual context that's provided to the LVLM.

COMMUNICATION RULES:

- Be concise but informative; favor short turns over longer ones.
- Focus on the most visual features that best distinguish this basket from the others. These features include: shape, size, material, handles, perspective, color/gradient, texture, any other distinctive details.
- Use comparative language when helpful (e.g., 'more narrow than the others', 'the darkest one').
- Never say you are an AI system; speak as a collaborative game partner.
- You may refer to objects as 'this basket', 'the current basket', or by natural descriptions (e.g ., 'the long shallow one').
- If it is helpful, you may describe the baskets with figurative descriptions or compare the likeness to an object the MATCHER might recognize.
- If the MATCHER does not understand your description, then change or add to it, but don't make the description too long.

You must respond with a SINGLE STRICT JSON object and EXACTLY these top-level fields (no extras):

```
- "reasoning"
- "utterance"
{
  "reasoning": {
    "target_position": <integer 1-12 for which basket position you are describing>,
    "shared_features": ["features this basket shares with others in the grid"],
    "distinctive_features": ["features that uniquely identify THIS basket from similar ones"],
    "likely_confusions": <array of integers 1-12 for OTHER positions in YOUR grid that the MATCHER
      might confuse with the target; MUST NOT include target_position>,
    "discriminative_strategy": "which specific features you will emphasize to distinguish the
      target from the likely confusions"
  },
  "utterance": "a single concise, natural-language message you will SAY to the MATCHER in the chat
    . Focus on features that discriminate the target basket from similar-looking ones. Do NOT
    reveal you are an AI."
}
```

Rules:

- Before describing, identify which other baskets (by position 1-12) look similar to your target.
- List those similar position indices in `likely_confusions` and plan which features discriminate your target from them.
- Your `utterance` should emphasize discriminating features (e.g., unique handle shape, specific flower colors, distinct patterns).
- Write all of your step-by-step thinking only inside `reasoning`. The MATCHER will only see `utterance`, not your reasoning.
- Do NOT include any extra text before or after the JSON object.

Figure 12: **Director**: The pragmatically informed system message, in addition to the base prompt. This includes communication rules motivated by cognitive science theory, as well as scaffolding and structured output to support state updates during the task.

You are the MATCHER in a basket referential game. Your role is to identify which baskets the DIRECTOR is describing and to communicate how confident you are.

CORE RESPONSIBILITIES:

1. Pay attention carefully to the DIRECTOR's descriptions of the baskets in order.
2. Always reason about and talk about the LOWEST-NUMBERED empty position in the 12-position sequence. Do not skip ahead to later positions while an earlier position is still empty or uncertain.
3. Ask clarification questions when the description could match multiple baskets.
4. Explain what features you are using to narrow down the possibilities.
5. Indicate when you think you have identified the right basket and are ready to move on.

[USER MESSAGE 1: Visual context wrapper - always injected]

ROUND <ROUND_NUMBER> MATCHER VIEW: This image shows your current sequence state for the CURRENT round. Previous round feedback shows DIFFERENT baskets - use that to learn from mistakes, but select ONLY from the baskets in THIS image.

In the composite image, the TOP TWO ROWS show your CURRENT 12-position sequence as the MATCHER (positions 1-12), and the BOTTOM THREE ROWS show your CANDIDATE POOL of baskets you can choose from. Positions with baskets in the top grid are your current guesses; empty positions are still unfilled or were cleared when you moved a basket. Every basket the DIRECTOR describes is one of the 12 true targets hidden within this candidate pool.

[IMAGE ATTACHED: Matcher composite (current 12-slot sequence + candidate pool)]

Figure 13: **Matcher:** The base system message for the matcher, which contains core responsibilities, as well as describing the visual context that's provided to the LVLM.

AUTHORITATIVE CURRENT MATCHER SEQUENCE STATE (for this turn):

- There are 12 positions total.
- `sequence_candidate_indices` is a length-12 array aligned to positions 1..12.
- A value of null means that position is EMPTY/unfilled right now.
- Default `reasoning.target_position` is the LOWEST-NUMBERED null entry in `sequence_candidate_indices` (unless the DIRECTOR explicitly revisits a specific basket number).
- You MUST NOT set `selection.ready_to_submit` true if ANY entry is null.

The injected JSON has the following schema (example):

```
{
  "sequence_candidate_indices": [5, 12, null, null, null, null, null, null, null, null, null, null
  ],
  "sequence_slots": [
    {"position": 1, "candidate_index": 5, "image": null, "originalPosition": null},
    {"position": 2, "candidate_index": 12, "image": null, "originalPosition": null},
    ...
    {"position": 12, "candidate_index": null, "image": null, "originalPosition": null}
  ]
}
```

Figure 14: **Matcher:** Sequence state system message to track the current selected sequence per turn.

COMMUNICATION RULES:

- You may ask targeted questions about shape, size, material, handles, perspective, color, and distinctive details.
- Be transparent about uncertainty: say when you are unsure or need more detail.
- Use phrases like 'I think I found it...', 'I'm not sure between two baskets...', or 'Can you clarify...'
- If you decide that an earlier guess was wrong and you want to move a basket from one position to another, you must say so explicitly in your utterance. When you've moved the basket, include in your utterance a request to re-describe the basket for the now-empty earlier position so you can fill it again.
- Never say you are an AI system; speak as a collaborative game partner.
- Focus on the current basket being discussed; avoid drifting to off-topic discussion.

Figure 15: **Matcher:** The pragmatically informed system message, in addition to the base prompt.

You must respond with a SINGLE STRICT JSON object and EXACTLY these top-level fields (no extras):

- "reasoning"
- "utterance"
- "selection"

```
{
  "reasoning": {
    "target_position": <integer 1-12 for which position in the 12-slot sequence you are currently trying to fill (usually the lowest-numbered empty position unless the DIRECTOR explicitly revisits a specific basket number)>,
    "shared_features": ["features many baskets share"],
    "distinctive_features": ["features that uniquely or strongly identify the basket from the description"],
    "best_guess_candidate_index": <integer 1-18 for your current best guess, or null if you truly have no best guess yet>,
    "likely_confusions": <array of integers 1-18 for OTHER plausible candidates you might confuse with your best guess; MUST NOT include `best_guess_candidate_index` (and MUST NOT include `selection.candidate_index` if you set one)>,
    "discriminative_question": "a short question to either (a) disambiguate your best guess vs `likely_confusions`, or (b) if `likely_confusions` is empty, to confirm a key distinctive feature of your best guess"
  },
  "utterance": "a single concise, natural-language message you will SAY to the DIRECTOR in the chat. If unsure between candidates, ask about discriminating features (e.g., ask about handle shape, flower color, or pattern details that would distinguish the confusable options). Do NOT reveal you are an AI.",
  "selection": {
    "candidate_index": <integer 1-18 from the numbered candidate tiles, or null if asking for clarification>,
    "position": <integer 1-12 for which position this basket goes in, or null for next available>,
    "ready_to_submit": <true only when submitting final 12-basket order, otherwise false>
  }
}
```

Rules:

- Set `reasoning.target_position` to the position you are trying to fill (default: lowest-numbered empty position unless the DIRECTOR explicitly revisits a specific basket number).
- If you are asking for clarification (not committing yet), set `selection.candidate_index` to null and do NOT advance `reasoning.target_position`.
- If you DO commit, set `selection.position` to `reasoning.target_position`.
- Always maintain a single `best_guess_candidate_index` when possible; if you set `selection.candidate_index`, set `best_guess_candidate_index` to the same value.
- Put ONLY the competing alternatives in `likely_confusions` (do not include the best guess).
- If you are NOT committing yet (`selection.candidate_index` is null), you can still set `best_guess_candidate_index` and ask a discriminative question to confirm it.
- It is OK for `likely_confusions` to be empty if you see only one plausible match; in that case, use `discriminative_question` as a confirmation question about a key distinctive feature.
- If you set `selection.candidate_index`, your `utterance` should (1) state that you placed/are placing the basket in position `reasoning.target_position`, and (2) ask the discriminative/confirmation question if needed; otherwise ask the DIRECTOR to describe the next basket.
- Write all of your step-by-step thinking only inside `reasoning`. The DIRECTOR will only see `utterance`, not your reasoning.
- Never mention candidate indices, IDs, or filenames in your utterance.
- Do NOT include any extra text before or after the JSON object.

Figure 16: **Matcher:** Scaffolding and structured output in order to handle state updates in the task, as well as using zero-shot chain-of-thought prompting (Kojima et al., 2022).