

# Are they lovers or friends? Evaluating LLMs’ Social Reasoning in English and Korean Dialogues

Anonymous ACL submission

## Abstract

As LLMs are increasingly deployed in real-world interactions, their social reasoning in interpersonal situations becomes critical. To explore their capabilities, we introduce **SCRIPTS**, a 1k-dialogue dataset in English and Korean, sourced from movie scripts. We furthermore propose a social reasoning task based on **SCRIPTS** that evaluates the capacity of LLMs to infer the social relationships (e.g., friends, sisters, lovers) between speakers in each dialogue. Among nine models’ evaluation results, current proprietary LLMs achieve around 75–80% on the English dataset and 58–69% in Korean. Strikingly, models predict relationships labeled as **UNLIKELY** by humans in 10–25% of responses in both languages. Furthermore, we find that thinking models and chain-of-thought prompting provide minimal benefits for social reasoning and occasionally amplify social biases. In sum, there are significant limitations in current LLMs’ social reasoning capabilities especially for Korean, highlighting the need for efforts to develop socially-aware LLMs.<sup>1</sup>

## 1 Introduction

As LLM-based agents become more prevalent, we also expect frequent interactions among multiple LLM agents and users (Cai et al., 2024; Liu et al., 2024). This trend is already reflected in practice (e.g., Group chats in ChatGPT (OpenAI, 2025)), and for more natural and smooth communication, LLMs should be able to recognize the relationship between the participants (Sehl et al., 2023). We term this ability to recognize and identify the relevant social relationships (e.g., lovers, friends, family members) as *social relationship reasoning*. When LLMs fail in this reasoning, they risk pro-

ducing responses that violate social norms or cause safety issues, as illustrated in Figure 1.

Although earlier studies have made important progress in evaluating LLMs’ ability to infer social relationships, they often use simplified settings that may not fully capture real-world complexity. For instance, some work frames the task as multiple-choice classification (Jia et al., 2021a; Li et al., 2023), considers a limited taxonomy of relationship types (Jia et al., 2021b; Tiginova et al., 2021), or focuses on relatively simple dyadic conversations (Jurgens et al., 2023). Moreover, social relationship inference is often inherently uncertain and context-dependent, so a single “correct” label may be difficult to justify in many cases (Hilton, 1995). For example, the remark “You never listen to me” could express a serious complaint between romantic partners or playful banter between friends depending on the broader conversational context and the language(s) used.

To address this shortcoming of previous studies, we introduce **SCRIPTS**, a novel benchmark for evaluating LLMs’ social relationship reasoning abilities, featuring an answer schema that incorporates inherent uncertainty (Figure 1). It contains 1k dialogues (English 580, Korean 567), derived from U.S. and Korean movie scripts, ensuring conversations are realistic and culturally grounded. We adopt soft labeling, where each relationship type is annotated with likelihood-based categories: **HIGHLY LIKELY**, **LESS LIKELY**, and **UNLIKELY**. By distinguishing relationship labels by likelihood, our dataset supports fine-grained evaluation that differentiates nonsensical responses from plausible but contextually less prominent ones.

We evaluate nine state-of-the-art models and find that current LLMs achieve moderate accuracy but still make frequent socially implausible inferences in both languages. Even the best-performing model (GPT-4o) still responds with relationships that humans annotated as **UNLIKELY** in 10.9%

<sup>1</sup>We include the dataset as a ZIP file with the submission, and we will publicize the dataset when the paper is accepted.

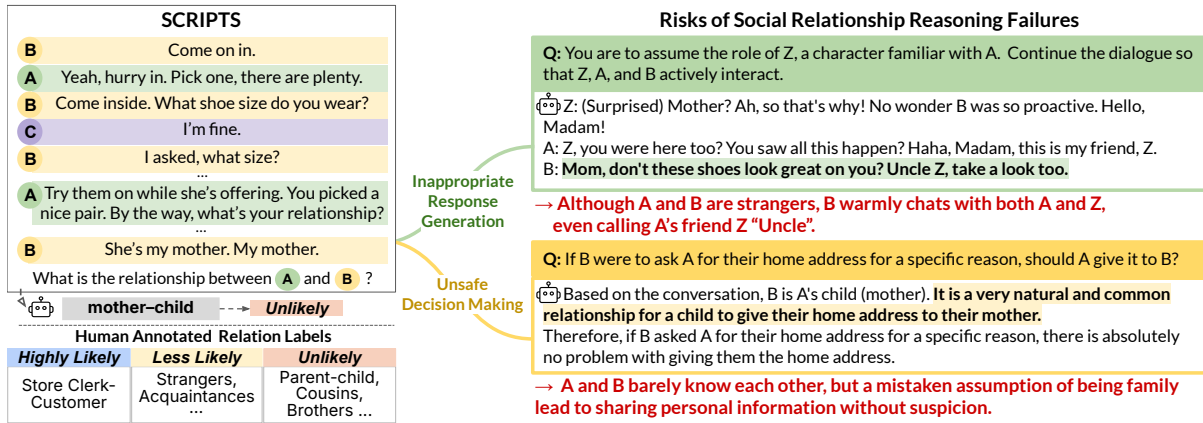


Figure 1: **Illustration of SCRIPTS and potential risks from failures in Social Relationship Reasoning.** **SCRIPTS** provides three human-annotated relationship labels: **HIGHLY LIKELY**, **LESS LIKELY**, and **UNLIKELY**. In this example, Gemini-2.5-Flash incorrectly predicts a store clerk-customer relationship as a mother-child. Such misleading relational reasoning can lead to inappropriate responses and unsafe decisions, such as privacy leakage. Examples are from Korean dialogues, translated into English.

of cases, settles on **LESS LIKELY** responses in 10.0%, and predicts **HIGHLY LIKELY** relationships in 79.1% of cases. Lower-performing models, such as Llama-3.1-8B-Instruct, make **HIGHLY LIKELY** predictions in only 41.3% of the English dataset. Surprisingly, chain-of-thought (CoT) prompting and thinking models, which are effective for logical reasoning, provide minimal benefits for social reasoning and occasionally amplify social biases.

We qualitatively analyze LLM failure cases by examining instances where models assign **UNLIKELY** relationship labels and identify four failure modes in both English and Korean. We further investigate whether the socio-demographic backgrounds (e.g., age, gender, relationship), formality, hierarchy and intimacy between human conversational partners enhance social relationship reasoning and find that these factors reduce **UNLIKELY** outputs and mitigate nonsensical behavior for LLMs in general.

In summary, our contributions are as follows:

- We introduce **SCRIPTS**, an English-Korean benchmark for evaluating LLMs' social relationship reasoning, comprising 1K dialogues with uncertainty-aware relationship labels.
- Evaluating nine LLMs, we find limited social reasoning ability: models often predict **UNLIKELY** relationships and vary substantially in identifying **HIGHLY LIKELY** relationships.
- We find that neither CoT prompting nor thinking models yields consistent gains, whereas adding relational information tends to shift predictions toward more likely inferences.

## 2 Related Work

**Evaluating Social Relationship Reasoning in LLMs** Existing research on computational models for social relationship reasoning is often constrained by simplified tasks and datasets. Many studies frame the problem as a classification or multiple-choice task, which makes it difficult to capture nuanced reasoning process (Li et al., 2023; Jia et al., 2021a; Chen et al., 2020). Some datasets are built from single utterances rather than multi-turn conversations, limiting contextual variation (Jurgens et al., 2023). Other studies using conversational data have methodological limitations. For instance, PRIDE (Tigunova et al., 2021) gathered annotations from movie summaries instead of dialogues, and Rashid and Blanco (2018) used static labels which fail to reflect what can be inferred from an interaction between humans. In contrast, our benchmark uses multi-turn dialogues annotated with multiple human-inferred labels, capturing relationships as they are expressed in dialogue. Also, we evaluate our task with open-ended generation rather than fixed-choice classification, allowing a wider range of possible relationships and better reflecting the social complexity.

**Cultural Dependency in Social Relationship Reasoning** Although most studies focus on English, social relationship reasoning depends on linguistic and cultural context: Korean, for example, relies on *terms of address* and *honorifics* to encode relational information in dialogue (Chung, 2010; Hwang, 1991).

Terms of address (ToA) are expressions used to directly refer to another person and carry discourse and social meaning (Hwang, 1991). In English, ToA rely largely on personal names, whereas in Korean they commonly include kinship terms, titles, and pronouns. For instance, instead of using an adult’s name, speakers may say “child’s name’s dad,” reflecting norms that discourage direct name use in certain contexts. While many languages encode politeness, Korean has a highly grammaticalized honorific system (Harada, 1976; Brown et al., 2014) that conveys roles, status, and formality through verbal morphology (Fukada and Asato, 2004; Brown, 2015; Pizziconi, 2011), unlike English which lacks an equivalent system. These cultural and linguistic differences motivate evaluating social relationship reasoning cross-linguistically and cross-culturally.

### 3 SCRIPTS: Evaluating LLMs’ Interpersonal Social Reasoning

We introduce **SCRIPTS**, a benchmark for social relationship reasoning in multi-turn dialogues in English and Korean. In this section, we outline the motivation (§ 3.1), design (§ 3.2), and construction (§ 3.3) of **SCRIPTS**.

#### 3.1 The Importance of Social Relationship Reasoning in LLMs

To participate in natural social conversations, LLMs must produce utterances that are appropriate for the underlying relationship and context. The right side of Figure 1 illustrates how misinterpreting relationships can lead to undesirable outcomes (e.g., social harm). In the example, an LLM (Gemini-2.5-flash) mistakes a store clerk–customer interaction for a family relationship, resulting in an inappropriate next response and potentially encouraging oversharing of personal information.

#### 3.2 Dataset Design

To capture the complexity and diversity of real-world social dynamics, we leverage movie scripts that contain natural human interactions spanning a wide range of relationships (Table 2). **SCRIPTS** makes two key contributions: (1) it models relationships as they are expressed within the dialogue, rather than relying on static role labels, and (2) it adopts a three-tier probabilistic scheme—**HIGHLY LIKELY**, **LESS LIKELY**, and **UNLIKELY**—for more fine-grained evaluation of social relationship reasoning.

Type	English	Korean	Total
Movies	28	32	60
Dialogues	580	567	1,147
3-Person Dialogues	223	256	479
Unique Highly-Likely Relationships	230	617	–
†Turns	10.21	9.89	10.05
†Highly-Likely Relationships	3.62	3.72	3.67
†Unlikely Relationships	18.50	23.13	20.79

Table 1: **Statistics of SCRIPTS**. **SCRIPTS** contains 1K English–Korean dialogues from 60 movies, annotated with 230+ unique relationship types. (†denotes the average per dialogue.)

**Dialogue-Level Labeling** A key design choice of **SCRIPTS** is to label relationships as they appear in specific conversational contexts. Prior benchmarks often assign fixed character roles from movie metadata (e.g., mother–son) (Jia et al., 2021b). However, social relationships are dynamic, as speakers can shift roles across contexts and may hold multiple relationships simultaneously.

We highlight the value of dialogue-level labeling by comparing our annotations with static, movie-level labels. We find that (1) 19% of movie-level labels are judged irrelevant for a given dialogue, suggesting that a global label can be misleading; and (2) even when a movie-level label applies, our annotations include more than three **HIGHLY LIKELY** relationships per dialogue on average. Together, these show that a single conversation often reflects multiple social facets, motivating our context-aware, dialogue-level labeling approach.

**Probabilistic Labeling** As social situations are inherently ambiguous, one dialogue may suggest multiple relationships with varying plausibility (Figure 1). We adopt a three-tier probabilistic scheme—**HIGHLY LIKELY**, **LESS LIKELY**, and **UNLIKELY**. This design offers two key benefits: (1) **granularity**, allowing metrics to reward models for identifying the most salient relation (**HIGHLY LIKELY**) rather than just plausible ones (**LESS LIKELY**); and (2) a **nonsense penalty**, which penalizes contextually inappropriate predictions (**UNLIKELY**)—a critical failure in social relationship reasoning.

#### 3.3 Dataset Construction

We collect 60 movie scripts: 28 English scripts crawled from IMSDb and 32 Korean scripts obtained via an onsite visit to the Korean Film Archive and crawling an open-access Korean script

Category	Specific Relationships
Family	Parent-Children, Brothers/Sisters, Grandparent-Grandchildren, Cousins, Uncle/Aunt-Niece
Social	Friends, Acquaintances, Neighbors, Strangers
Romance	Romantic Interest, Dating, Married, Engaged, Friends with benefits, Affair, Ex-relationship
Organizational	Coworkers, Professional colleagues, Supervisor-Subordinate relationship
Role-based	Mentor-Mentee, Teacher-Student, Lawyer-Client, Doctor-Patient, Landlord-Tenant
Antagonist	Competitive relationship, Rivalry, Arch-enemies

Table 2: **Initial relationship categories and specific relationships used for UNLIKELY annotation (50 items).**

community.<sup>2</sup> The full movie list and metadata are provided in Appendix Table 14.

We filter scenes to those with at least three turns ( $\geq 3$ ) and two or three participants. Among the remaining scenes, we prioritize those with diverse speaker combinations (avoiding repeated exchanges among the same few participants). From 23k initial scenes, this yields 1,322 high-quality dialogues (698 English; 624 Korean) for human annotation. Additional collection and filtering details are in Appendix A. While prior work primarily studies dyadic interactions, **SCRIPTS** includes three-speaker dialogues (41.8%) to capture more complex social settings. For these scenes, the task remains dyadic relationship inference between two interlocutors; we randomly select the speaker pair.

### 3.3.1 Collecting Human Annotations

Each dialogue is annotated by three annotators who are native or near-native speakers with extensive cultural familiarity (e.g., 10+ years of residence in U.S. and South Korea each). We recruit 17 annotators for English and 14 annotators for Korean (see Appendix A.5 for details).

For human annotation, we employ a four-phase annotation process (See Appendix D for details).

**Phase 1: Labeling UNLIKELY Relationships.** Given a predefined set of 50 relationships compiled from prior work (Tigunova et al., 2021; Jurgens

<sup>2</sup>imsdb.com; kmdb.or.kr; filmmakers.co.kr.

et al., 2023) (Table 2), annotators select all relationships that are clearly contradicted by the dialogue. A label is marked UNLIKELY only when a majority of annotators agree, yielding high-confidence negatives.

**Phase 2: Open-ended Labeling for HIGHLY LIKELY Relationships.** Annotators then provide open-ended text labels for the relationship(s) most strongly supported by the dialogue. After normalization (e.g., merging “mother-child” and “mom-child”), we take the union of these labels as the HIGHLY LIKELY set.

**Phase 3: Deriving LESS LIKELY Relationships.** All remaining relationships from the predefined list that are neither HIGHLY LIKELY nor UNLIKELY are automatically categorized as LESS LIKELY.

**Phase 4: Annotating Auxiliary Labels** Annotators label auxiliary socio-demographic attributes (e.g., age, gender) and relational dimensions (e.g., formality, hierarchy, intimacy). We use these labels for finer-grained analyses of factors shaping relationship inferences, focusing on cues humans are known to rely on for social relationship reasoning (Wish et al., 1981; Nguyen et al., 2016).

**Quality Control** To ensure annotation quality, we run training sessions and a pilot study. We exclude dialogues where the three annotators’ HIGHLY LIKELY labels have no overlap (i.e., are mutually exclusive), indicating low reliability. This yields 1,147 dialogues (580 English; 567 Korean), removing 13.2% of the initial dataset.

Table 1 shows the detailed statistics of **SCRIPTS**. Also, see Appendix A.6 for comparisons of the types and frequencies of relationships in English and Korean dialogues.

## 4 Evaluating LLMs with SCRIPTS

We evaluate 9 LLMs: 3 proprietary models (GPT-4o, o3, Gemini-2.5-flash), 3 widely used open-source models (Qwen-3-{8B/14b}, Llama-3.1-8B-instruct), and 3 open-source multilingual models specialized for Korean (A.X-4.0-light-7B, Kanana-1.5-8B, and Exaone-4.0-30B).<sup>3</sup>

**Evaluation and Metrics** Considering the probabilistic nature of the task, we run each model five times per dialogue and take the model’s majority response among them. We then compute (1) the

<sup>3</sup>See Appendix B.1 for model configurations.

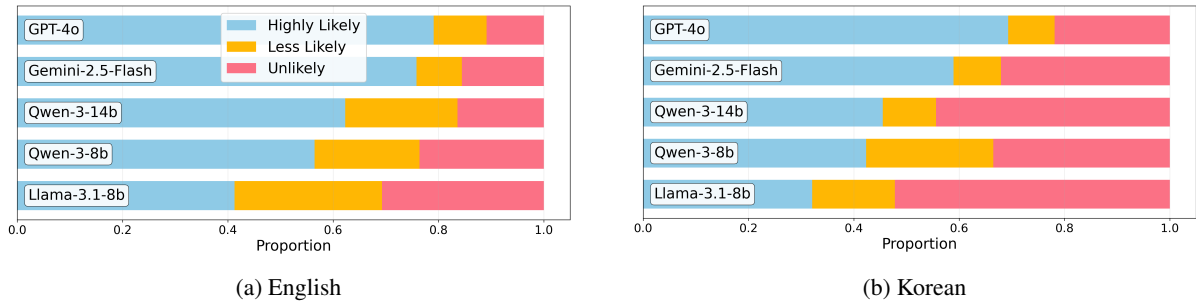


Figure 2: **Comparison of model performance in English and Korean datasets.** **HIGHLY LIKELY** represents the accuracy of the model’s majority response being a highly likely response, while **UNLIKELY** indicates the error rate where the model generate an unlikely response. **LESS LIKELY** indicates the proportion of cases in which the model generates neither a **HIGHLY LIKELY** nor an **UNLIKELY** prediction.

Model	Thinking	En		Ko	
		HIGHLY LIKELY (↑)	UNLIKELY (↓)	HIGHLY LIKELY (↑)	UNLIKELY (↓)
OpenAI/GPT-4o	×	0.767	0.116	0.642	0.215
OpenAI/o3	✓	<b>0.807</b>	<b>0.086</b>	<b>0.742</b>	<b>0.152</b>
Gemini-2.5-flash	×	0.759	0.154	<b>0.582</b>	0.318
Gemini-2.5-flash	✓	<b>0.776</b>	<b>0.138</b>	0.538	<b>0.239</b>
Qwen-3-14b	×	0.623	0.164	0.455	0.444
Qwen-3-14b	✓	<b>0.673</b>	<b>0.107</b>	<b>0.467</b>	<b>0.443</b>

Table 3: **Model comparison with and without Thinking mode across English (En) and Korean (Ko).**

305 proportion of samples whose majority response  
306 falls into the **HIGHLY LIKELY** relation set and (2)  
307 the proportion of which majority response falls  
308 into the **UNLIKELY** relation set. We use GPT-4o to  
309 evaluate each model’s short-form answers. GPT-4o  
310 compares model outputs with ground-truth labels,  
311 a common practice in prior LLM evaluation work.  
312 We report 92.0% human-validated accuracy for  
313 GPT-4o as an evaluator (Appendix B.3).

#### 314 4.1 Overall Performance

315 Figure 2 shows the performance of five multi-  
316 lingual models (i.e., GPT-4o, Gemini-2.5-Flash,  
317 Qwen-3-{8/14}B, and Llama-3.1-8B-Instruct). We  
318 find that GPT-4o achieves the best performance  
319 with **HIGHLY LIKELY** rate of 79% and 69% in  
320 English and Korean, respectively. The models in-  
321 correctly infer an **UNLIKELY** relationship in 10.8–  
322 31.9% of their responses and this tendency is ampli-  
323 fied in the Korean dataset with a rate increasing by  
324 an additional 7.2–16.5%p. Table 7 in Appendix B  
325 provides the exact numerical values. We provide  
326 a case study of these behaviors with the frequent  
327 failure modes of the models in § 5.

#### 328 4.2 Does Thinking Help?

329 We analyze the performance of models when CoT  
330 prompting or internal thinking processes are in-  
331 corporated. These methods have been shown to  
332 be effective for improving reasoning on math and  
333 STEM-related tasks (Wei et al., 2022).

#### 334 Effectiveness of Chain of Thought Prompting

335 We study CoT prompting on four multilingual mod-  
336 els (one per family): GPT-4o, Gemini 2.5 Flash,  
337 Qwen-3-8B, and Llama-3.1-8B-Instruct. As shown  
338 in Table 8, CoT does not consistently help: Gemini  
339 2.5 Flash shows a 1.7%p drop in **HIGHLY LIKELY**  
340 responses in English, and Llama-3.1-8B-Instruct  
341 shows a 3.1%p rise in **UNLIKELY** responses in Ko-  
342 rean. This contrasts with general reasoning tasks,  
343 where CoT often helps, suggesting that social rea-  
344 soning may limit CoT’s benefits.

#### 345 Effectiveness of Thinking Process

346 We evaluate three reasoning models: o3, Gemini-2.5-Flash,  
347 and Qwen-3-14B, comparing the latter two with/  
348 without an internal thinking process. We compare  
349 performance in a single-run setting (due to budget  
350 constraints). As shown in Table 3, the results vary  
351 significantly per language. In English, enabling  
352 thinking yields better performance. The impact of  
353 the thinking process in Korean is negligible (e.g.,

Rank	En	Ko
1	* <b>A.X-4.0-Light (0.589)</b>	* <b>A.X-4.0-Light (0.467)</b>
2	† <b>Qwen-3 (0.565)</b>	† <b>Qwen-3 (0.423)</b>
3	‡ <b>Llama-3.1 (0.413)</b>	Exaone-4.0 (0.409)
4	Kanana-1.5 (0.406)	Kanana-1.5 (0.328)
5	Exaone-4.0 (0.318)	‡ <b>Llama-3.1 (0.321)</b>

Table 4: **Model ranking of Korean-specialized and open-source models in English and Korean**, based on the **HIGHLY LIKELY** response rate (numbers in parentheses indicate the corresponding values). \*/†/‡ denote ranks 1/2/3 in English.

Qwen-3-14b) and enabling the thinking process (in Gemini-2.5-flash) even leads to a 4.4%p decrease in the rate of **HIGHLY LIKELY** responses.

### 4.3 Do Korean-specialized models perform better in the Korean context?

We evaluate three Korean-specialized models: A.X-4.0-light (7B), Kanana-1.5-8B, and Exaone-4.0-32B. A.X-4.0<sup>4</sup> is further trained on Korean data on top of Qwen. The training data for Kanana-1.5 and Exaone-4.0 are undisclosed, but their technical reports indicate Korean performance comparable to or better than English (Kanana LLM Team et al., 2025; LG AI Research, 2025).

Table 4 compares the three models with similarly sized open-source multilingual models (Qwen-3-8B, Llama-3.1-8B-Instruct). The results show that A.X-4.0-Light and Qwen-3-8B achieve the best and second-best performance in both languages, but the 3rd–5th rankings differ. In English, Llama-3.1-Instruct-8B ranks 3rd, while in Korean, Exaone-4.0-32B and Kanana-1.5-8B, take 3rd and 4th, with Llama-3.1-Instruct-8B ranking last. This implies that a language and culture-specific approach to model training may be beneficial for our task. Full results are in Appendix B.4 (Table 9).

## 5 Reasons Behind Failure of LLMs

Based on qualitative analyses of the reasoning processes of LLMs in CoT experiments (§ 4.2), we identify four types of failures.

**Failure to Distinguish Terms of Address and References** Models often misinterpret a term of reference (i.e., a word used to refer to someone) as a term of address (i.e., a word used to call someone directly), leading to a fundamental misunderstanding of the social context.

<sup>4</sup><https://github.com/SKT-AI/A.X-4.0/>

#### Dialogue 1 (English):

[A]: Hi Officer, can I help you?  
 [B]: Yes, I’m hoping you can. An elderly gentleman went missing from the nursing home down the street. Staff seems to think he came here.  
 (...)  
 [A]: (Pause, then) Oh...that’s my Dad. He can’t talk. Had a major stroke a few years back. But he’s doing well. Ain’t ya Pop?  
 [B]: OK, well, thanks for your time. Here’s my number in case you hear of anything. Sorry to bother you.  
 (...)

**Ground Truth:** Police officer–Civilian, Strangers

**Prediction:** Parent–Children / Father–son (Llama, Gemini, GPT)

In Dialogue 1, speaker [A] says, “*that’s my Dad*”, using “*Dad*” as a reference to identify a third person for the police officer [B]. However, the models latch onto this keyword and misinterpret it as a term of address from [A] to [B], leading to the incorrect inference of a Parent-Child relationship. This failure leads the models to ignore clear cues (e.g., [B] being called “*Officer*”) that contradict this interpretation.

In Korean, this error is more pronounced because speakers often use address terms to refer to themselves. For example, a teacher may tell a student, “The teacher (I) told you to do this,” where “the teacher” is a self-reference. However, models may misread “the teacher” as a third party.

**Failure to Aggregate Multiple Cues** Social relationship reasoning requires the ability to identify multiple cues within the context and integrate them to arrive at a conclusion. However, models fail to integrate cues, especially when their combination is atypical.

#### Dialogue 2 (Translated from Korean):

(...)  
 [B]: (Salutes) Hey.  
 [C]: Hey.  
 [B]: Hey... What’s this? A drowned body? Doesn’t even look that deep to me.  
 [A]: Doesn’t seem like he drowned.  
 [B]: Then what, a dumped body?  
 [A]: Nah... doesn’t look dumped either. Go take a closer look. Go on.  
 [B]: Then what the hell is it?  
 [A]: Hey B, you know, brace yourself before you look.  
 [B]: You kidding me? Damn it... shit...

**Ground Truth:** Friend/Coworker

**Prediction:** Supervisor-Subordinate (Qwen, Gemini, GPT)

In Dialogue 2, GPT-4o identifies three cues: (1) A and B converse casually without honorifics, (2) the topic concerns work, and (3) B complies with A’s instruction. The model places greatest weight on the last cue, interpreting the interaction as hierarchical and labels them as supervisor–subordinate. For Korean speakers, this is implausible, since subordinates are expected to use honorifics when addressing a superior. The absence of honorifics indicates the relationship is not hierarchical, but rather that of coworkers or friends. This shows that even

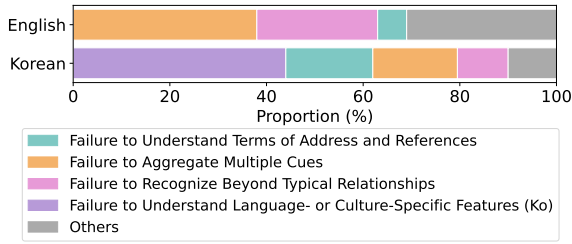


Figure 3: Distribution of GPT-4o’s 30 failure cases by error type in English and Korean.

when the models detect relevant cues, they are unable to prioritize and integrate them within the social context. We provide original Korean scripts for Dialogue 2 in Appendix C.1.

### Failure to Recognize Atypical Relationships

Models frequently struggle to recognize relationships that deviate from conventional or stereotypical patterns, such as non-hierarchical (equal) exchanges between parents and children or hierarchical conversation between married couple.

#### Dialogue 3 (English):

[A]: So you’re seeing Mom tomorrow, huh? At my parent-teacher thing?  
 [B]: Yeah.  
 [A]: First time in a while.  
 [B]: Yeah, but no biggie.  
 [B]: Hey, what’s with the moping?  
 [A]: Nothing. It’s just... there’s this girl.  
 [B]: Oh yeah? You like her?  
 [A]: I like [C]. This girl’s my soulmate. I’m like crazy, stupid, in love with her. And she wants someone else.  
 [B]: But she’s your soulmate?  
 [A]: Yeah.

Ground Truth: Parent-Children

Prediction: Siblings (Llama, GPT, Gemini)

In Dialogue 3, all human annotators agree on the parent–child relationship, yet the models reject it: “less likely since the conversation feels more peer-like rather than hierarchical or guiding” (GPT-4o), “the casual tone and discussion about a crush imply a more peer-like relationship” (Qwen-3-8b), and “if [B] is the parent, they might not discuss the girl in such a casual way” (Llama-3.1-8b-instruct). Gemini-2.5-flash likewise dismisses the parent–child relationship, reasoning that B is a bachelor and A’s parents are deceased, concluding it is not a traditional parent–child relationship, revealing a stereotyped conception of family roles.

### Failure to Understand Language- or Culture-Specific Features (Ko)

In Korean, this error type accounts for the largest share of failures. Most confusions stem from the misinterpretation of terms of address and honorifics. For example, unlike English, *eomeoni*, literally “mother,” can also refer to

a friend’s mother or an older woman, yet the model often predicts a parent–child relationship whenever it appears. This issue is especially pronounced in dialogues containing culturally specific terms such as kinship expressions. For instance, Qwen misinterprets *Hyungsoo* (older brother’s wife) as “older brother”, and *Hyungrim* as “father”, resulting in a complete failure. Honorifics are also frequently misinterpreted—for example, equal relationships (e.g., friendships) predicted as hierarchical, and vice versa.

Additionally, we manually examine 30 failure cases of GPT-4o (best performing model). Figure 3 presents their distribution across error types in English and Korean. In English, the majority of errors arise from Failure to Aggregate Multiple Cues (36.7%). In contrast, Korean errors are predominantly caused by difficulties in handling Language and/or Culture-Specific Features (46%), with smaller proportions attributed to the other categories. This disparity highlights the model’s difficulty in identifying relational cues embedded in Korean-specific cultural and linguistic markers, such as terms of address and honorific systems, thereby revealing the culturally dependent nature of social relationship reasoning.

## 6 Does Providing Additional Social Information Help?

Humans interpret social relationships using demographic cues and relational dimensions (Nguyen et al., 2016). We investigate whether such social information can similarly enhance the social relationship reasoning abilities of LLMs using four models—GPT-4o, Gemini-2.5-flash, Qwen-3-8B, and Llama-3.1-8B-instruct—each representing a different model family while excluding thinking-enabled and Korean-specialized models.<sup>5</sup> Specifically, we examine how providing social information influences performance in inferring social relationships, considering two types of information: demographic cues (age, gender) and relational dimensions (intimacy, hierarchy, formality).

**Experimental Setting** We design six experimental settings with two variables. First, we vary the type of social information: (i) age/gender only, (ii) relational dimensions only, or (iii) both. Second, we vary the source of social information: (a) human-annotated gold data, available in the dataset

<sup>5</sup>We find that models often rely on relational dimensions when inferring relationships; see Appendix C.2 for examples.

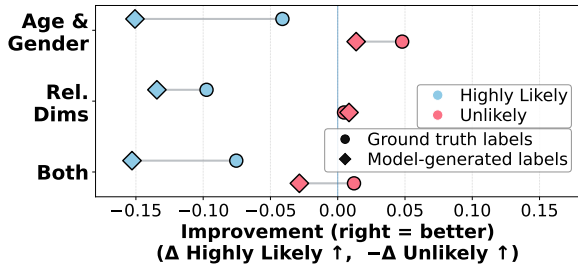


Figure 4: **Impact of Relational Information on GPT-4o’s Performance.** Positive values indicate improvement, while negative values indicate deterioration after adding relational information.

as metadata (see § 3.3.1 for illustration), or (b) model-generated predictions, where the model infers each type of social information and incorporates these predictions into the social relationship reasoning. The accuracy of these predictions is reported in Table 12 in the appendix. Detailed experimental settings, including the prompt, appear in Appendix B.2.3.

**Results With Ground Truth Labels** Figure 4 shows GPT-4o results across six settings on the English dataset. Providing human gold information yields no substantial or consistent performance gains, but it reduces the proportion of UNLIKELY predictions. This suggests that while such information may not directly guide identification of HIGHLY LIKELY relationships, it helps models avoid UNLIKELY ones. The tendency holds across models, except for Qwen-3-8B. For instance, when a model’s initial inference of an UNLIKELY relationship is refined into a more plausible one, the label “Intimate” often shifts toward labels that are typically more intimate: Strangers → Romantic Interest (3.3%). Similarly, when given the label “No hierarchy”, the most common change is also Parent–Children → Friends (2.9%). Thus, dimension labels provide additional cues about relationships, enabling the model to incorporate them and reduce implausible predictions. However, these changes do not always yield correct reasoning. Sometimes models over-rely on dimensional labels rather than context. For instance, in an atypical “close” superior–subordinate relationship, GPT-4o misinterprets the interaction due to the intimate tone, even when clear terms of address are present.

**Results With Model-Generated Labels** With model-generated information, auxiliary labels do not consistently help because they are often inac-

Social. Info.	Improved (Unlikely→Likely)	Deteriorated (Likely→Unlikely)
Age & gender	72.8	65.5
Rel. dims.	53.3	50.7

Table 5: **Accuracy comparison for social information inference between improved and deteriorated cases of social relationship reasoning.**

curate. For example, GPT-4o achieves under 60% accuracy on age and gender and below 75% on relational dimensions (see Table 12 for accuracy across four models). Consistent with this, inferred social information is more accurate in cases where it improves relationship reasoning than in cases where it degrades performance (Table 5).

These results suggest that demographic cues and relational dimensions, which humans naturally rely on, can facilitate social relationship reasoning. However, current LLMs are limited in their ability to infer these dimensions. Therefore, instructing LLMs to infer these factors before identifying the social relationships is ineffective. Results for other models are provided in Table 10-11 in the Appendix D. In table 13 of Appendix D, we examine the link between social-information inference and relationship reasoning using separate logistic regressions for each factor.

## 7 Conclusion

We introduce a bilingual dataset **SCRIPTS** to investigate the limitations of current LLMs in social relationship reasoning. Our experiments show that most models perform suboptimally across English and Korean, and often assign UNLIKELY relationships. Our analyses (§4) reveal that current reasoning techniques such as CoT do not consistently benefit social reasoning. Furthermore, we provide an analysis on where LLMs fail, especially focusing on cases where models respond with “UNLIKELY” RELATIONSHIPS. Our cross-linguistic results, including improved performance of Korean-specific models on Korean dataset relative to that in English, demonstrates the importance of language and culture-specific approaches to advance LLMs’ social reasoning abilities. We hope that **SCRIPTS** serves as a starting point for exploring how to improve LLMs’ social relationship reasoning in diverse contexts.

## 577 Limitations

578 While our dataset includes dialogues in both En- 611  
579 glish (U.S.) and Korean (South Korea), thereby of- 612  
580 fering greater linguistic diversity than many prior 613  
581 studies, the scope of our analysis remains limited 614  
582 to these two languages. As a result, our findings 615  
583 may not fully generalize to interpersonal interac- 616  
584 tions in other languages or cultural contexts. So- 617  
585 cial norms, communication styles, and relational 618  
586 cues can vary significantly across cultures, and fu- 619  
587 ture work should extend this research to a broader 620  
588 range of languages in order to more comprehen- 621  
589 sively evaluate LLMs’ interpersonal social reason- 622  
590 ing abilities across various linguistic and cultural 623  
591 settings. Additionally, as discussed in § 5, we ana- 624  
592 lyze CoT traces to analyze characterize model be- 625  
593 havior. However, we acknowledge that these traces 626  
594 may reflect post-hoc rationalizations rather than 627  
595 the mechanisms that produced the final answer. 628

## 596 Ethics Statement

597 This study involves human annotation based on pre- 631  
598 existing movie scripts, which may contain harmful 632  
599 or offensive content due to the nature of the source 633  
600 material. The study was reviewed and approved by 634  
601 an IRB, and informed consent was obtained from 635  
602 all participants prior to their involvement. Annota- 636  
603 tors were recruited via an institutional participant 637  
604 portal and compensated at hourly rates of KRW 638  
605 15,000 (Korea) and USD 20 (U.S.), approximately 639  
606 1.5× the local minimum wage. 640

## 607 Acknowledgements

608 We used AI assistants, including ChatGPT<sup>6</sup> to sup- 641  
609 port the writing, and Cursor<sup>7</sup> to support coding. 642

<sup>6</sup><https://chatgpt.com/>

<sup>7</sup><https://cursor.com/>

## References

- Lucien Brown. 2015. Honorifics and politeness. *The handbook of Korean linguistics*, pages 303–319. 611 612
- Lucien Brown, Bodo Winter, Kaori Idemaru, and Sven Grawunder. 2014. Phonetics and politeness: Perceiving korean honorific and non-honorific speech through phonetic cues. *Journal of Pragmatics*, 66:45–60. 613 614 615 616 617
- Zhenyao Cai, Seehee Park, Nia Nixon, and Shayan Doroudi. 2024. Advancing knowledge together: integrating large language model-based conversational ai in small group collaborative learning. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*, pages 1–9. 618 619 620 621 622 623
- Yi-Ting Chen, Hen-Hsen Huang, and Hsin-Hsi Chen. 2020. [MPDD: A multi-party dialogue dataset for analysis of emotions and interpersonal relationships](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 610–614, Marseille, France. European Language Resources Association. 624 625 626 627 628 629 630
- Kyung-Sook Chung. 2010. Korean evidentials and assertion. *Lingua*, 120(4):932–952. 631 632
- Atsushi Fukada and Noriko Asato. 2004. Universal politeness theory: application to the use of japanese honorifics. *Journal of pragmatics*, 36(11):1991–2002. 633 634 635
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*. 636 637 638 639 640
- Shin-Ichi Harada. 1976. Honorifics. In *Japanese generative grammar*, pages 499–561. Brill. 641 642
- Denis J. Hilton. 1995. [The social context of reasoning: Conversational inference and rational judgment](#). *Psychological Bulletin*, 118(2):248–271. 643 644 645
- Shin Ja J Hwang. 1991. [Terms of address in korean and american cultures](#). *Intercultural Communication Studies*, 1(2):117–136. 646 647 648
- Qi Jia, Hongru Huang, and Kenny Q. Zhu. 2021a. [Ddrel: A new dataset for interpersonal relation classification in dyadic dialogues](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(14):13125–13133. 649 650 651 652 653
- Qi Jia, Hongru Huang, and Kenny Q Zhu. 2021b. [Ddrel: A new dataset for interpersonal relation classification in dyadic dialogues](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 13125–13133. 654 655 656 657 658
- David Jurgens, Agrima Seth, Jackson Sargent, Athena Aghighi, and Michael Geraci. 2023. [Your spouse needs professional help: Determining the contextual](#) 659 660 661

662	appropriateness of messages through modeling social relationships. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10994–11013, Toronto, Canada. Association for Computational Linguistics.	719
663		720
664		721
665		722
666		723
667		
668	Kanana LLM Team, Yunju Bak, Hojin Lee, Minho Ryu, Jiyeon Ham, Seungjae Jung, Daniel Wontae Nam, Taegyeong Eo, Donghun Lee, Doohae Jung, Boseop Kim, Nayeon Kim, Jaesun Park, Hyunho Kim, Hyunwoong Ko, Changmin Lee, Kyoung-Woon On, Seulye Baeg, Junrae Cho, Sunghye Jung, Jieun Kang, EungGyun Kim, Eunhwa Kim, Byeongil Ko, Daniel Lee, Minchul Lee, Miok Lee, Shinbok Lee, and Gaeun Seo. 2025. <b>Kanana: Compute-efficient bilingual language models.</b>	724
669		725
670		726
671		727
672		
673		
674		
675		
676		
677		
678	LG AI Research. 2025. <b>Exaone 4.0: Unified large language models integrating non-reasoning and reasoning modes.</b>	
679		
680		
681	Hengli Li, Song-Chun Zhu, and Zilong Zheng. 2023. Diplomat: A dialogue dataset for situated pragmatic reasoning. <i>Advances in Neural Information Processing Systems</i> , 36:46856–46884.	
682		
683		
684		
685	Jiawen Liu, Yuanyuan Yao, Pengcheng An, and Qi Wang. 2024. Peergpt: Probing the roles of llm-based peer agents as team moderators and participants in children’s collaborative learning. In <i>Extended abstracts of the CHI conference on human factors in computing systems</i> , pages 1–6.	
686		
687		
688		
689		
690		
691	Dong Nguyen, A. Seza Dođruöz, Carolyn P. Rosé, and Franciska de Jong. 2016. <b>Computational sociolinguistics: A Survey.</b> <i>Computational Linguistics</i> , 42(3):537–593.	
692		
693		
694		
695	OpenAI. 2025. Introducing group chats in chatgpt. <a href="https://openai.com/index/group-chats-in-chatgpt/">https://openai.com/index/group-chats-in-chatgpt/</a> . Accessed: 2025-12-31.	
696		
697		
698	Barbara Pizziconi. 2011. <b>Honorifics: The cultural specificity of a universal mechanism in japanese.</b> In Dániel Z. Kádár and Sara Mills, editors, <i>Politeness in East Asia</i> , pages 45–70. Cambridge University Press.	
699		
700		
701		
702	Farzana Rashid and Eduardo Blanco. 2018. <b>Characterizing interactions and relationships between people.</b> In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 4395–4404, Brussels, Belgium. Association for Computational Linguistics.	
703		
704		
705		
706		
707		
708	Claudia G. Sehl, Ori Friedman, and Stephanie Denison. 2023. <b>The social network: How people infer relationships from mutual connections.</b> <i>Journal of Experimental Psychology: General</i> , 152(4):925–934.	
709		
710		
711		
712	Anna Tigunova, Paramita Mirza, Andrew Yates, and Gerhard Weikum. 2021. <b>PRIDE: Predicting Relationships in Conversations.</b> In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 4636–4650, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
713		
714		
715		
716		
717		
718		
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	719
		720
		721
		722
		723
	Myron Wish, Morton Deutsch, and Susan J Kaplan. 1981. Perceived dimensions of interpersonal relations. In <i>The Psychology of Social Situations</i> , pages 113–129. Elsevier.	724
		725
		726
		727
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	728
		729
		730
		731

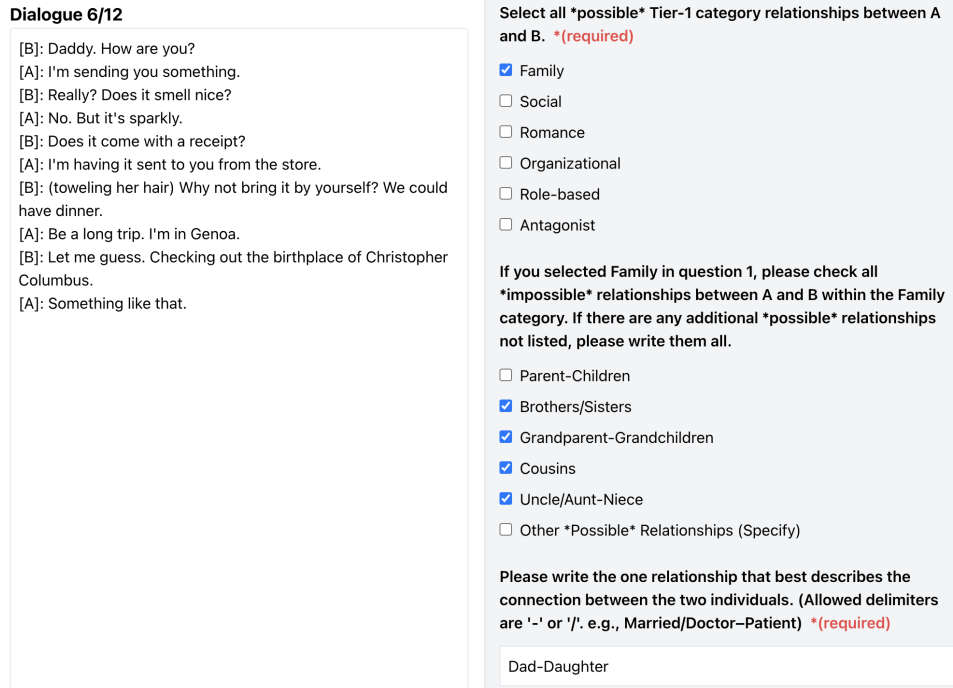


Figure 5: Annotation Platform. The annotators can read the dialogue on the left side and annotate the relationships and relational dimensions on the right side.

## A Dataset

We source dialogues from movie scripts. The collection process consist of the following steps: (1) movie selection, (2) raw scene collection, (3) OCR processing and human verification, (4) scene filtering (5) anonymization, and (6) relationship annotation.

### A.1 Movie Selection

We select the movies based on the following criteria to capture daily real-life interactions. Only modern-day movies released after 2000 are included, excluding medieval fantasy or alien sci-fi. We consider age limits, only including movies up to PG-13 for English movies and up to 15세 관람가 (suitable for audiences aged 15 and above) for Korean movies.

We collect English movie scripts from IMSDb.<sup>8</sup> For Korean movies, due to the limited online dialogue resources, we visited the Korean Film Archive (KOFA)<sup>9</sup> to collect physical copies of movie scripts. From KOFA, we only collect one-third of each movie script to adhere to the data use policy. Also, we collect additional movie scripts from Filmmakers Online Community.<sup>10</sup> In total,

<sup>8</sup><https://imsdb.com>

<sup>9</sup><http://www.kmdb.or.kr/>

<sup>10</sup><https://www.filmmakers.co.kr/>

we collect 60 movies (28 English, 32 Korean) across various genres.

### A.2 Raw Scene Collection and Processing

For physical copies of Korean movies from KOFA, we use NAVER CLOVA OCR API to extract dialogue texts.<sup>11</sup> After OCR, we use GPT-4o to further process and clean the text into structured format. Human annotators then verify and modify the outputs with the original PDF files. As a result, we obtain 16k English and 7k Korean scenes.

### A.3 Scene Filtering

Scenes are filtered to include at least four utterances, and involve two to three speakers. To maximize speaker diversity, we prioritize scenes with unique set of speakers. Using these criteria, we select 1,322 scenes from an initial pool of 23k scenes, comprising 698 English and 624 Korean evaluation sets. The selected movies and the number of scenes per movie are listed in Table 14.

### A.4 Anonymization

To mitigate potential data contamination (e.g., identifying the source movie) and reduce bias (e.g., gender inference), all character names are automatically replaced with placeholders such as [A] and

<sup>11</sup><https://www.ncloud.com/product/aiService/ocr>

[B]. Any names not covered by this process are further verified and anonymized manually.

## A.5 Human Annotation

We construct a golden answer set with human annotators who have over ten years of experience in the target language and culture. Annotators include undergraduate and graduate students in South Korea and the United States, compensated at 1.5 times the minimum hourly wage in their country. Before annotation, they attend an introductory session covering data usage policies and guidelines.

For social relationships, annotators are given an initial set of possible relationships (Table 2), adapted from Tigunova et al. (2021), and asked to mark UNLIKELY ones. To reduce workload, annotators first choose LIKELY relationship categories for each dialogue, then select the UNLIKELY relationships from the list of specific relationships in those categories. They also provide up to five open-ended answers describing relationships that best characterize the interaction. These serve as candidates for LIKELY relationships.

For relational dimensions, we provide annotators with definitions of each dimension and ask them to rate dialogues on a 5-point scale: intimacy (from strongly intimate to strongly unintimate), formality (from strongly formal to strongly informal), and hierarchy (from A»B to A«B). When constructing the golden labels, we collapse the ratings into a 3-point scale (e.g., intimate, neutral, unintimate) and assign the majority-voted label. The inter-annotator agreement is reported in Table 6.

### Recruitment and Management of Human Annotators

Figure 5 shows the human annotation platform: the anonymized scene appears on the left, and annotation questions on the right. All annotations were conducted with IRB approval. We recruited 17 English annotators (7 male, 10 female) and 14 Korean annotators (5 male, 9 female). We explained our study and how the data would be used during a Zoom session with the annotators, and obtained their informed consent. All annotators were undergraduate or graduate students enrolled at universities in the United States or Korea. The Korean annotators were all native speakers, while the English annotators were either U.S. citizens or individuals who had lived in the United States for over ten years. For quality control, applicants were asked to complete the task on three sample items during recruitment, and their responses were

reviewed by the authors to select the final annotators. After selection, annotators participated in an orientation session and a training phase designed to support them in performing the task as effectively as possible.

We provide the participant recruitment announcement below.

### Participant Recruitment (English)

<b>Overview</b>	We are recruiting participants for a research experiment that evaluates the conversational understanding abilities of language models. This study builds a dataset to assess models' social reasoning in dialogue. Participants will read short dialogues from movie scripts and label the social relationships between speakers.
<b>What you will do</b>	<ul style="list-style-type: none"><li>• Identify <b>social role-based relationships</b> (e.g., parent–child, romantic partners, mentor–mentee).</li><li>• Label <b>relationship aspects</b> (e.g., intimacy/closeness, hierarchy/power, purpose: work-oriented vs. casual).</li><li>• Infer <b>speaker attributes</b> (e.g., gender and approximate age).</li></ul>
<b>Eligibility</b>	<ul style="list-style-type: none"><li>• Comfortable using web-based interfaces for research participation.</li><li>• Fluent in English and highly familiar with U.S. culture (e.g., lived in the U.S. for 10+ years).</li><li>• Age 18 or older.</li><li>• Not offended by dialogues that may include profanity, offensive language, or depictions of violence.</li><li>• Registered (or able to register) as a participant on Amazon Mechanical Turk.</li></ul>
<b>IRB Safety Notes</b>	In accordance with Institutional Review Board (IRB) guidelines, we cannot recruit individuals directly supervised by the research lead, nor undergraduate students under the age of 18. Movie scripts may contain profanity, offensive language, and morally questionable situations. Participation is voluntary, and you may withdraw at any time without penalty.

## A.6 Diversity of SCRIPTS

Figure 6 presents the ten most common relationships in each language. Both the English and Korean datasets frequently include social (e.g., friends, acquaintances), organizational (e.g., coworkers, supervisor-subordinate), and familial (e.g., parent-child, siblings) relationships.

Beyond these shared categories, we also examine which relationship types appear exclusively in Korean or English. Each dataset contains culturally specific relationships that reflect distinct social

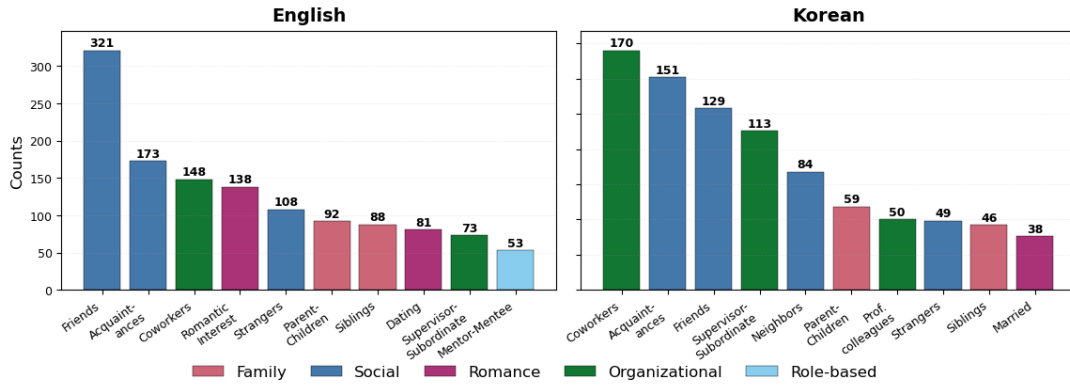


Figure 6: Top 10 Relationships in each Language dataset.

roles and lexicalizations.

**Korean-only relations** include *North Korean soldier-citizen* (1), *shaman-client* (2), *shaman-assistant* (1), *private tutor-student* (1), and *student’s family acquaintance-tutor* (1). In addition, kinship terms are more fine-grained in Korean: for example, distinctions such as *older brother-younger brother* (1) and *older brother-sister-in-law* (1), whereas in English these are typically generalized under a single “siblings” category.

**English-only relations** include roles such as *father figure-child*, *mother figure-child*, *co-parents*, and *babysitter-child*, reflecting cultural and social roles that are more explicitly lexicalized in English.

These findings highlight how cultural context shapes the granularity and salience of social relationships represented in dialogue datasets.

Figure 7 illustrates that our dataset can capture diverse interpersonal dynamics by labeling relational dimensions. The typicality of certain relationships is often defined by their levels of intimacy, formality, and hierarchy (Wish et al., 1981). For instance, friendship is generally characterized as intimate, non-hierarchical (equal), and informal. Yet, our dataset also includes atypical relationships. For instance, over 40% of friend relationships in our dataset deviate from these typical dimensions.

	Hierarchy		Formality		Intimacy	
	All	2>	All	2>	All	2>
EN	0.333	0.416	0.408	0.513	0.314	0.426
KO	0.462	0.550	0.469	0.562	0.375	0.458

Table 6: Inter-annotator agreement (Fleiss’  $\kappa$ ) by relational dimension. “2>” indicates samples with at least two annotators in agreement. We used these filtered samples in our experiments.

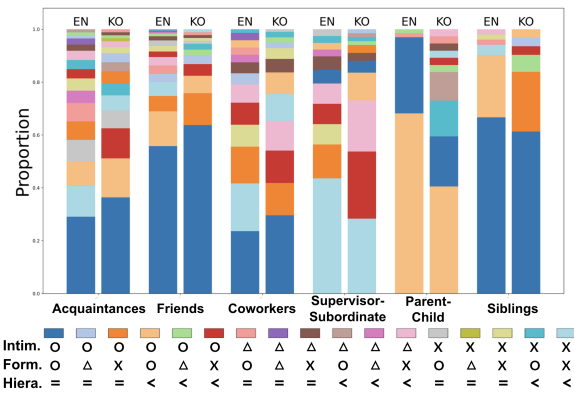


Figure 7: Comparative Analysis of Relational Dimension Distributions for Six Relationship Types Present in Both English and Korean Top-10 Relations. Legend labels denote intimacy (O: intimate, X: unintimate,  $\Delta$ : neutral), formality (O: formal, X: informal,  $\Delta$ : neutral), and hierarchy (<: hierarchical, =: equal).

## B Experiment

### B.1 Model Configuration

We use total nine models, GPT-4o<sup>12</sup>, o3<sup>13</sup>, Gemini-2.5-flash, Qwen-3-{8B/14b} (Yang et al., 2025), Llama-3.1-8B-instruct (Grattafiori et al., 2024), A.X-4.0-light-7B<sup>14</sup>, Kanana-8B<sup>15</sup>, and Exaone-4.0-30B (LG AI Research, 2025). For GPT-4o, we use the gpt-4o-2024-05-13 version via the OpenAI API. For Gemini-2.5-Flash, we use OpenRouter (openrouter.ai). We set the temperature of the models to 1.

<sup>12</sup><https://platform.openai.com/docs/models/gpt-4o>

<sup>13</sup><https://platform.openai.com/docs/models/o3>

<sup>14</sup><https://github.com/SKT-AI/A.X-4.0>

<sup>15</sup><https://github.com/kakao/kanana>

## B.2 Prompt

### B.2.1 Prompt for Evaluation

#### Prompts for Evaluation

You are a judge that evaluates the correctness of relation classification outputs from the model. Your task is to check each relation in the model's output and determine whether it belongs to the \*Possible\* or \*Impossible Relation Sets\* provided. For each relation:

- If it is included in the Possible Relation Sets, mark "is\_in\_possible": 1, otherwise mark it as 0.
- Similarly, if it is included in the Impossible Relation Sets, mark "is\_in\_impossible": 1, otherwise mark it as 0.

Answer **\*\*only\*\*** in the following JSON format (no explanations):

**Possible Relation Sets:** {gt\_list}

**Impossible Relation Sets:**

{impossible\_gt\_list}

**Model output:** {model\_outputs}

**Output in JSON format:**

```
{
  relation1 : {
    "is_in_possible": 1,
    "is_in_impossible": 0
  },
  relation2 : {
    "is_in_possible": 0,
    "is_in_impossible": 1
  },
  ...
}
```

**Output:**

### B.2.2 Prompt for Inference

#### Prompt for Inference(EN)

Read the following conversation and guess the relationship of the participants [A] and [B]. When guessing the relationship, refer to the following examples of relationships:

{example\_relations}

If the relationship matches one of the examples above, use it as is, but if the relationship does not fit any of the examples, describe the relationship yourself.

Your answer about the relationship must be in JSON format:

```
{
  "relation": ""
}
```

{Additional Information}

**Conversation:** {dialogue}

**Output (JSON):**

#### Prompt for Inference(KO)

다음 대화를 읽고, 등장인물 A와 B의 관계를 추측하십시오.

관계를 추측할 때는 다음의 관계 예시를 참고하십시오:

{example\_relations}

만약 위의 예시에 해당하는 관계라면 그대로 사용하고, 예시에 없는 관계라고 판단되면 해당 관계를 직접 서술하십시오.

관계에 대한 최종 답은 반드시 JSON 형식으로 답변하십시오:

```
{
  "relation": ""
}
```

{Additional Information}

**대화:** {dialogue}

**Output (JSON):**

#### Example Relations

##### English Relations

```
[
  "Parent-Children",
  "Brothers/Sisters",
  "Grandparent-Grandchildren",
  "Cousins",
  "Uncle/Aunt-Niece",
  "Friends",
  "Acquaintances",
  "Neighbors",
  "Strangers",
  "Romantic Interest",
  "Dating",
  "Married",
  "Engaged",
  "Friends with benefits",
  "Affair",
  "Ex-relationship",
  "Coworkers",
  "Professional colleagues",
  "Supervisor-Subordinate relationship",
  "Mentor-Mentee",
  "Teacher-Student",
  "Lawyer-Client",
  "Doctor-Patient",
  "Landlord-Tenant",
  "Competitive relationship",
  "Rivalry",
  "Arch-enemies"
]
```

##### Korean Relations

```
[
  "부모-자식",
  "형제/자매/남매",
  "조부모-손주",
  "사촌",
  "삼촌/이모/고모-조카",
  "단짝 친구",
  "친구",
  "지인",
  "이웃",
  "모르는 사이",
  "쌤",
  "연애",
  "부부",

```

"약혼관계",  
 "Friends with benefits",  
 "불륜관계",  
 "전애인 관계",  
 "동료",  
 "직장 동료",  
 "상관-부하직원 관계",  
 "멘토-멘티",  
 "선생-제자",  
 "변호사-고객",  
 "의사-환자",  
 "집주인-세입자",  
 "경쟁관계",  
 "라이벌 관계",  
 "숙적"

]

**CoT setting:** We append “*Think step by step*” at the end of the prompt to encourage chain-of-thought reasoning.

### B.2.3 Prompts Used in § 6

For the experiment in § 6, we add additional information to the prompt. The additional information consists of two types: Age & Gender and Relational Dimensions. The prompts for each type are as follows.

In the Ground Truth Labels setting, we fill the placeholders {age\_gender\_info} and {relational\_dimensions\_info} with human-annotated gold labels. In the Model-Generated Labels setting, the model is first asked to separately infer each type of information, and the infer results are then inserted back into the corresponding placeholders.

#### Additional Information(EN)

**Base setting:** None

**Age & Gender:** The age and gender information of the participants [A] and [B] are as follows. Please refer to them when inferring the nature of their relationship. {age\_gender\_info}

**Relational Dimensions:** The Intimacy level, Pleasure level, and Hierarchy level between A and B in the conversation are as follows. Please refer to them when inferring the nature of their relationship. {Relational Dimensions\_info}

#### Additional Information(Ko)

**Base setting:** None

**Age & Gender:** 등장인물 A와 B의 나이와 성별은 다음과 같다. 그들의 관계의 성격을 추론할 때 참고하라. {age\_gender\_info}

**Relational Dimensions:** 대화에서 A와 B 사이의 친밀감(Intimacy) 수준, 격식(Formality) 수준, 그리고 위계(Hierarchy) 수준은 다음과 같다. 그들의 관계의 성격을 추론할 때 참고하라.

{Relational Dimensions\_info}

### B.3 Validating LLM-as-a-Judge

To validate the accuracy of GPT-4o as an evaluator, we sample 100 question-answer pairs for each language, and two authors independently verify the results. As a result, GPT-4o has accurately evaluated in 97.85% of the responses in English and 86.2% in Korean, with an inter-annotator agreement of 96% (Cohen’s  $\kappa = 0.58$ ) between the two authors.

### B.4 Results

See Table 7 for the complete results, Table 8 for results with CoT prompting, and Table 9 for results from Korean-specialized models.

## C Qualitative Analysis - Cues

### C.1 Original Korean Dialogue

Dialogue 2 (Korean):

(...)  
 [B]: (경례) 어이.  
 [C]: 왔냐?  
 [B]: 야. 뭐냐? 이 익사야? 별로 깊어 보이지도 않는데.  
 [A]: 익사는 아닌 것 같고.  
 [B]: 그럼 뭐 유기?  
 [A]: 아하....유기도 아닌 것 같은데. 너가 가서 한번 봐봐. 한번.  
 [B]: 그럼 뭐야?  
 [A]: 야. B야. 그 마음 단단히 먹고 봐.  
 [B]: 장난하나. 에이씨. 자.. 에이씨.

### C.2 What cues do LLMs rely on in social reasoning?

To understand how models use and integrate cues to infer social relationships, we conduct a qualitative analysis on their CoT reasoning.

**Terms of Address and Reference** LLMs frequently leverage terms of address and references as explicit cues to infer social relationships. For instance, when a speaker use terms like “*Daddy*” or “*Professor [B]*”, the models infer family-based or professional relationship. Self-reference also provide valuable information. For example, a speaker introducing themselves as “*Doctor [A]*” signals their professional identity as a medical practitioner, leading to LLMs suggesting relationships such as Doctor-Patient or Doctor-Doctor. Furthermore, LLMs analyze how individuals refer to third parties to understand the relationship between the referring individuals themselves. For example, if both A and B refer to a third person as “*Sergeant [C]*”, the LLM can infer that A and B are likely colleagues within a military context, and that their shared use

Model	English		Korean	
	HIGHLY LIKELY (↑)	UNLIKELY (↓)	HIGHLY LIKELY (↑)	UNLIKELY (↓)
GPT-4o	0.7910	0.1088	0.6931	0.2187
Gemini-2.5-flash	0.7582	0.1554	0.5894	0.3204
Qwen-3-8b	0.5648	0.2360	0.4233	0.3351
Llama-3.1-8b	0.4128	0.3074	0.3210	0.5220

Table 7: Comparison of model performance in English (En) and Korean (Ko) datasets. HIGHLY LIKELY represents the accuracy of the model’s majority response being a highly likely response, while UNLIKELY indicates the error rate where the model generate an unlikely response.

Model	En		Ko	
	HIGHLY LIKELY (↑)	UNLIKELY (↓)	HIGHLY LIKELY (↑)	UNLIKELY (↓)
GPT-4o	0.8024 (0.0114)	0.0971 (-0.0117)	0.6953 (0.0022)	0.2025 (-0.0162)
Gemini-2.5-flash	0.7412 (-0.0170)	0.1265 (-0.0289)	0.6031 (0.0137)	0.2010 (-0.1194)
Qwen-3-8b	0.6882 (0.1332)	0.1513 (-0.1261)	0.4974 (0.023)	0.3651 (-0.0259)
Llama-3.1-8b	0.5410 (0.0868)	0.2667 (-0.0410)	0.2996 (-0.0578)	0.5769 (0.0306)

Table 8: Comparison of model performance with Chain of Thought Prompting across English (En) and Korean (Ko) with deltas in parentheses.

Model	En		Ko	
	HIGHLY LIKELY (↑)	UNLIKELY (↓)	HIGHLY LIKELY (↑)	UNLIKELY (↓)
ax-4.0-light	0.5889	0.1934	0.4674	0.4127
exaone-4.0-32b	0.3178	0.3074	0.4092	0.4674
kanana-1.5-8b	0.4059	0.2884	0.328	0.3739

Table 9: Performance of Korean Specialized models in English and Korean.

953 of a formal title suggests a potentially task-oriented  
954 conversation.

955 **Conversation context and background** LLMs  
956 also take into account the context of the conver-  
957 sation (e.g., *a school, church, workplace, home*).  
958 They then utilize this background information to  
959 infer the social relationship or level of intimacy  
960 between the individuals involved in the dialogue.

961 **Tone or Atmosphere** LLMs also assess the emo-  
962 tional tone of individuals in a dialogue to judge  
963 their social dimensions, particularly intimacy and  
964 formality, utilizing that information to infer their  
965 social relationship. The models often associate  
966 *casual, friendly, teasing, empathetic, or support-*  
967 *ive* tones with more intimate relationships while  
968 *aggressive, frustrated, or angry* expressions are  
969 linked to less intimate or strained relationships.  
970 Similarly, emotional expressions, whether friendly  
971 or hostile, are often connected to informal relation-  
972 ships while the models associate *serious, indiffer-*  
973 *ent, dismissive, or emotionally neutral* expressions  
974 with formal relationships.

975 **Relational Dimensions** When inferring social  
976 relationships, models often consider relational di-  
977 mensions (*intimacy, hierarchy, formality*) in their  
978 rationale. For instance, in a dialogue where A play-  
979 fully jokes with B while B shares personal con-  
980 cerns, the model infers strong intimacy and sug-  
981 gests a close tie such as friendship or siblinghood.

982 However, it is important to note that while us-  
983 ing social dimensions as cue, particularly hierar-  
984 chy, LLMs often reveal social stereotypes, defining  
985 “*typical*” relational dimensions to certain relation-  
986 ships. For example, models assume that a parent-  
987 child inherently shares a *hierarchical* relationship  
988 while a married couple would generally have a  
989 *non-hierarchical(equal)* relationship. This lead to  
990 failures when the social interaction deviate from  
991 these norms, real-life atypical relationships.

## 992 D Does Providing Additional Social 993 Information Help?

994 This section provides supplementary material for  
995 Section 6. Tables 10–11 present the results across  
996 models, while Table 12 report the accuracy of in-  
997 ferring social information.

Model	$\Delta$ Highly Likely ( $\uparrow$ )	$\Delta$ UnLikely ( $\downarrow$ )	Model	$\Delta$ Highly Likely ( $\uparrow$ )	$\Delta$ UnLikely ( $\downarrow$ )
<b>GPT-4o</b>			<b>GPT</b>		
Age & Gender	-0.0411	-0.0479	Age & Gender	-0.1507	-0.0137
Rel. Dims	-0.0975	-0.0047	Sub Dims	-0.1344	-0.0084
Both	-0.0754	-0.0121	Both	-0.1529	0.0285
<b>Gemini-2.5-flash</b>			<b>Gemini</b>		
Age & Gender	-0.0411	-0.0343	Age & Gender	-0.1027	-0.0137
Rel. Dims	-0.1563	-0.0570	Sub Dims	-0.1194	-0.0460
Both	-0.1452	-0.0534	Both	-0.1711	-0.0017
<b>Qwen-3-8b</b>			<b>Qwen</b>		
Age & Gender	0.0127	0.0392	Age & Gender	-0.0395	0.0321
Rel. Dims	-0.1580	0.1209	Sub Dims	-0.0237	0.0442
Both	-0.1138	0.0398	Both	-0.0916	0.0435
<b>Llama-3.1-8b</b>			<b>Llama</b>		
Age & Gender	0.0205	-0.1027	Age & Gender	0.0548	-0.0137
Rel. Dims	-0.0123	-0.0735	Sub Dims	0.0099	-0.0514
Both	0.0505	-0.0994	Both	-0.0492	-0.0145

(a) With Ground Truth Labels

(b) With Model-Generated Labels

Table 10: Impact of Relational Information on Model Performance (English).

### Associations Between Social Information and Social Relationship Reasoning Performance

To further examine the link between social information inference and relationship reasoning, we run separate logistic regressions for each factor. Table 13 shows that most factors are positively associated with social relationship reasoning. This suggests that models performing well on age, gender, and relational dimension inferences also tend to perform better on overall social relationship reasoning, highlighting the interconnections among these dimensions.



Language	Movie Title	Movie Title (Ko)	Genre	Year	# of Scenes
EN	Amelia	-	Adventure, Biography, Drama	2009	17
	Autumn in New York	-	Drama, Romance	2000	26
	Big Fish	-	Adventure, Epic, Drama	2003	27
	Bruce Almighty	-	Comedy, Fantasy	2003	32
	Crazy Love	-	Documentary, Romance	2007	21
	Crazy, Stupid, Love.	-	Romance, Comedy, Drama	2011	39
	Date Night	-	Romance, Comedy, Crime	2010	24
	Easy A	-	Comedy, Drama, Romance	2010	45
	He's Just Not That Into You	-	Romance, Comedy, Drama	2009	16
	Larry Crowne	-	Comedy, Drama, Romance	2011	15
	Monte Carlo	-	Adventure, Comedy, Family	2011	6
	Moonrise Kingdom	-	Romance, Adventure, Comedy	2012	6
	New York Minute	-	Comedy, Adventure, Crime	2004	59
	Something's Gotta Give	-	Comedy, Drama, Romance	2003	21
	Speed Racer	-	Action, Adventure, Comedy	2008	7
	The Blind Side	-	Drama, Biography, Sport	2009	16
	The Bounty Hunter	-	Comedy, Action, Romance	2010	21
	The Brothers Bloom	-	Comedy, Action, Adventure	2008	15
	The Curious Case of Benjamin Button	-	Drama, Fantasy, Romance	2008	22
	The Fault in Our Stars	-	Drama, Romance	2014	17
	The Italian Job	-	Action, Crime, Thriller	2003	22
	The Invention of Lying	-	Comedy, Fantasy, Romance	2009	20
	The Next Three Days	-	Thriller, Action, Drama	2010	12
	The Pacifier	-	Action, Comedy, Drama	2005	16
	The Secret Life of Walter Mitty	-	Adventure, Comedy, Romance	2013	12
	The Theory of Everything	-	Drama, Biography, Romance	2014	13
Water for Elephants	-	Drama, Romance	2011	15	
Wild Hogs	-	Action, Adventure, Comedy	2007	18	
KO	200 Pounds Beauty	미녀는 괴로워	Comedy, Drama, Music	2006	28
	A Violent Prosecutor	검사외전	Action, Comedy, Crime	2016	8
	Battle for Incheon: Operation Chromite	인천상륙작전	Action, Drama, History	2016	9
	Cold Eyes	감시자들	Action, Crime, Thriller	2013	13
	Deranged	연가시	Drama, Sci-Fi, Thriller	2012	15
	Exit	엑시트	Comedy	2019	19
	Extreme Job	극한직업	Comedy, Crime	2019	1
	Hide and Seek	숨바꼭질	Horror, Mystery, Thriller	2013	3
	Jeon Woochi	전우치	Action, Adventure, Comedy	2009	3
	Marathon	말아톤	Biography, Drama, Sport	2005	20
	May 18	화려한 휴가	Drama, History	2007	13
	Miss Granny	수상한 그녀	Comedy, Fantasy, Music	2014	15
	My Tutor Friend	동감내기	Action, Comedy, Romance	2003	36
	Northern Limit Line	연평해전	Drama, War	2015	11
	Ode to My Father	국제시장	Drama, War	2014	11
	Pandora	판도라	Disaster, Action, Drama	2016	29
	Punch	완득이	Comedy, Drama, Sport	2011	25
	Secret Reunion	의형제	Action, Drama, Thriller	2010	24
	Secretly, Greatly	은밀하게 위대하게	Drama, Action, Comedy	2013	13
	Silmido	실미도	Action, Drama	2003	24
	Sunny	써니	Comedy, Drama	2011	31
	Take Off	국가대표	Comedy, Drama, Sport	2009	31
	The Attorney	변호인	Crime, Drama, History	2013	7
	The Berlin File	베를린	Spy, Action, Thriller	2013	18
	The Himalayas	히말라야	Adventure, Biography, Drama	2015	5
	The Neighbors	이웃사람	Thriller, Mystery	2012	36
	The Priests	검은 사제들	Horror, Mystery, Thriller	2015	9
	The Roundup	범죄도시2	Action, Crime, Thriller	2022	21
The Thieves	도둑들	Action, Comedy, Crime	2012	31	
Tidal Wave	해운대	Action, Drama, Sci-Fi	2009	17	
Tunnel	터널	Disaster, Drama	2016	35	
Veteran	베테랑	Action, Comedy, Crime	2015	6	

Table 14: List of movies in **SCRIPTS**. The genre (top three) and release year are sourced from IMDb. The dataset contains 60 movies (English 28 / Korean 32) spanning various genres.