

PIVOTSBench: Evaluating Fine-Grained Interpersonal Relationship Reasoning in Multimodal Large Language Models

Anonymous ACL submission

Abstract

Humans possess an innate ability to understand fine-grained interpersonal relationships, which is central to everyday social interactions. Although such reasoning is inherently multimodal, it remains largely unexplored by existing multimodal large language models (MLLMs). To address this gap, we introduce **PIVOTS**, the first benchmark built from Social-IQ 2.0 and YouTube data to evaluate MLLMs’ ability to predict bidirectional interpersonal relationship dimensions grounded in established psychology research. In addition, PIVOTS includes auxiliary tasks that assess models’ ability to identify and leverage the critical visual cues underlying such predictions. We evaluate both proprietary and open-source MLLMs and conduct detailed ablation studies to analyze the effects of visual modalities and explicit social role information in conversational utterances. We further examine how joint and pairwise prediction settings benefit MLLMs in scoring bidirectional PIVOTS dimensions.

1 Introduction

Humans routinely infer interpersonal relationships in everyday social interactions, a skill essential for collaboration, persuasion, trust formation, and conflict resolution. Moreover, our ability to attribute social relationships is inherently fine-grained and cannot be represented by conventional social categories such as family, friends, colleagues, or acquaintances. As illustrated in Fig. 1, although both examples depict a parent–child relationship, the underlying social dynamics differ dramatically: in the left example, the interaction is largely positive and cooperative, whereas in the right example it is tense and adversarial. An AI system capable of reasoning about such fine-grained interpersonal dimensions could both facilitate everyday human interactions and enhance multi-agent collaboration(Li et al.,

2023) and socially aware decision making(Park et al., 2023).

Recent advances in large language models (LLMs) have attracted interest in their potential for social intelligence (Zhu et al., 2025; Mou et al., 2025; Talebirad et al., 2025), yet interpersonal relationship reasoning remains largely unexplored. Recent work has investigated whether existing LLMs can assess interpersonal relationships from text only corpora (Zhou et al., 2025). It is worth noting that humans’ remarkable ability to infer social dynamics stems from the orchestration of multimodal signals. Using Fig. 1 for illustration, the left example demonstrates a case where linguistic content alone is insufficient and visual cues are critical for interpretation, whereas in the right example visual and linguistic signals jointly inform social reasoning. Notably, a systematic study of whether Multimodal Large Language Models (MLLMs) can leverage multimodal cues to decipher social dynamics is still missing from the community.

To bridge this gap, we present PIVOTS, the first benchmark for evaluating MLLMs on interpersonal relationship reasoning, constructed using data curated from Social-IQ 2.0 (Wilf et al., 2023) and YouTube. Specifically, we consider six dimensions grounded in prior psychology research (Wish et al., 1976): egalitarian vs. hierarchical power (**P**), superficial vs. intense involvement (**I**), positive vs. negative valence (**V**), socioemotional vs. task-oriented objective (**O**), temporary vs. enduring permanence (**T**), and cooperative vs. competitive stance (**S**)

Unlike prior work that models interpersonal relationships unidirectionally (Rashid and Blanco, 2017), our benchmark adopts a *bidirectional* setting that evaluates whether MLLMs can predict how each subject interprets the interpersonal relationship with the other across six dimensions. This design is motivated by information asymmetry, where two participants in the same interaction may form different beliefs about their relationship due

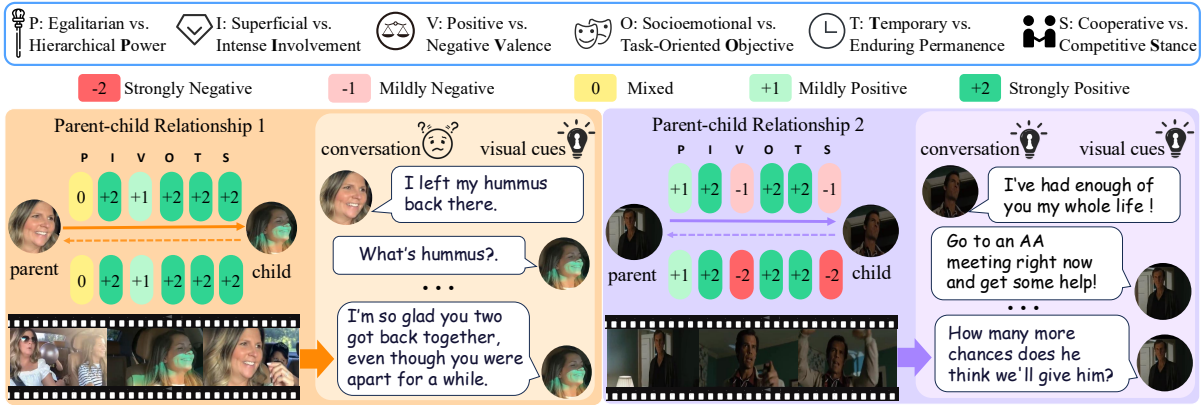


Figure 1: *Understanding complex social dynamics requires joint reasoning over language and visual cues.* Although both examples involve parent–child interactions, they exhibit drastically different configurations across multiple interpersonal dimensions. Moreover, conversational utterances alone are often insufficient to reason about these dimensions, whereas visual cues provide complementary signals.

to asymmetric perspectives, roles, and communicative intent. In addition to bidirectional dimension scoring, we introduce auxiliary visual reasoning tasks to evaluate whether MLLMs can identify and interpret the visual cues underlying interpersonal relationship predictions.

In our experiments, we first evaluate prevailing MLLMs on the PIVOTS benchmark, and then conduct a detailed analysis of how different input modalities and explicit role information in conversational utterances contribute to model performance. We further explore alternative prediction settings for interpersonal dimension scoring, showing that modeling multiple dimensions together and predicting bidirectional relationships in a paired manner can offer advantages over independent prediction in certain scenarios. Moreover, we investigate prompting strategies that inject human reasoning heuristics, examining how they yield clear improvements in salient social interactions while exhibiting limitations in contrastive and in-depth social scenarios.

2 Related Work

2.1 Interpersonal Relationship Understanding

Social Relationship Recognition (SRR) has traditionally been formulated as a classification problem over a fixed taxonomy of relationship types (Sun et al., 2017; Li et al., 2017; Wang et al., 2010; Li et al., 2020; Zhang et al., 2015). Video based SRR methods typically focus on modeling temporal interaction structure. For example, MSTR introduced a multi scale spatio temporal reasoning framework and released the VISR benchmark (Liu

et al., 2019). However, research in the social sciences (Endsley, 2021) suggests that human social interactions are highly complex and cannot be adequately captured by conventional social roles such as leader or teammate, as the same pair of roles may exhibit drastically different social dynamics across contexts and over time. More concretely, conventional social roles fail to capture the subtle, multimodal signals that actually drive social influence (Hoozeboom and Wilderom, 2020). Instead, fine grained interpersonal relationships, such as cooperation and competition, provide more informative signals for understanding social dynamics (Wish et al., 1976). This dimensional perspective has also been adopted in computational linguistics, where interpersonal relational dimensions are inferred from textual data (Rashid and Blanco, 2017). Importantly, relationship signals in real world interactions are often conveyed through multimodal behaviors rather than text alone, yet multimodal interpersonal relationship identification remains largely unexplored. To address this gap, we introduce a novel benchmark that evaluates whether MLLMs can infer six carefully defined interpersonal dimensions from videos and dialogues and reason about the visual cues underlying these relationships.

2.2 Multi-modal Social Intelligence Benchmarks

A rich body of literature has studied explicit verbal and nonverbal signal modeling in social settings, including tasks such as active speaker localization (Grauman et al., 2022) and audiovisual diarization (Xu et al., 2022). In contrast, our investigation is more closely related to efforts aimed

at developing computational models of social cognition. Early research (Sun et al., 2017; Li et al., 2017) in social cognition primarily focused on Social Relationship Recognition (SRR) from static images and explored domain adaptation across diverse social scenes. More recently, video-based benchmarks such as MMTOM-QA (Jin et al., 2024) and AGENT (Shu et al., 2021) have utilized procedurally generated environments to probe core psychological reasoning and Theory of Mind. (Lai et al., 2023) introduce a multimodal benchmark that studies persuasion strategy modeling using both linguistic and nonverbal signals. While Social IQA (Sap et al., 2019) and Social-IQ 2.0 (Wilf et al., 2023) perform well in general social commonsense evaluation, they lack a systematic framework for measuring interpersonal dimensions derived from social signals. Our PIVOTS benchmark addresses this gap by probing whether MLLMs can understand pairwise interpersonal relationship dimensions, thereby pinpointing the strengths and limitations of AI systems on social cognition capability.

3 PIVOTS Benchmark

3.1 PIVOTS Dimensions Definition

The foundation of our interpersonal relationship definition is grounded in the theoretical framework of (Wish et al., 1976), which identifies four fundamental dimensions underlying interpersonal perception. We further expand the coverage by introducing two complementary axes, valence and temporal dimension, based on insights from prior research (Russell, 1980). We operationalize these six dimensions using a five-point bipolar Likert scale, ranging from -2 to +2. This mapping translates nuanced social dynamics into discrete, quantifiable ratings, facilitating robust computational modeling and further supervised learning. Formally, we consider the following six axes:

- *Egalitarian vs. Hierarchical Power (P)*: This dimension relates to power distribution, status, and control, indicating whether the individuals involved have similar or different levels of power and roles.
- *Superficial vs. Intense Involvement (I)*: This dimension reflects the depth, involvement, and emotional investment in the relationship.
- *Positive vs. Negative Valence (V)*: This dimension captures the degree of pleasantness and satisfaction versus hostility in the relationship.
- *Socioemotional vs. Task-Oriented Objective (O)*:

This dimension distinguishes between relationships that are primarily personal and relaxed versus those that are goal-oriented, businesslike, and structured.

- *Temporary vs. Enduring Permanence (T)*: This dimension distinguishes between relationships that are primarily long-term based and with specific kind of commitment or those that are short-termed and freewheeling.
- *Cooperative vs. Competitive Stance (S)*: This dimension, also referred to as affiliation or affect, describes the degree of warmth, support, and collaboration in the relationship.

Detailed definitions, five-point rating rubrics for each dimension, and illustrative examples are provided in Appendix A.

3.2 Task Definition

In this section, we formally define a set of benchmark tasks derived from the six interpersonal relationship axes introduced earlier. Apart from the primary task of scoring the PIVOTS interpersonal dimensions, we introduce auxiliary tasks for keyframe identification and grounded visual causal reasoning to assess models’ ability to utilize visual information effectively. We present a visual illustration of our tasks in Fig. 2.

Task 1: Six-Dimensional Scoring. Given the video \mathcal{V} , conversation utterance \mathcal{U} , and a directional relationship \mathcal{D}_{AB}^α , MLLMs aim to predict the corresponding score \mathcal{R}_{AB}^α . Note that \mathcal{D}_{AB}^α specifies how subject A interprets subject B along one of the six interpersonal axes $\alpha \in \{P, I, V, O, T, S\}$. Specifically, we define the relationship score prediction via the following conditional probability:

$$\mathcal{R}_{AB}^\alpha = \arg \max_{r \in \{-2, -1, 0, 1, 2\}} p(r \mid \mathcal{V}, \mathcal{U}, \mathcal{D}_{AB}^\alpha),$$

where r represents the discrete candidate score for each PIVOTS dimensions.

Task 2: Key Frame Identification.

Given video \mathcal{V} and direction \mathcal{D}_{AB}^α , MLLMs aim to identify the key video frame that best reflects the identification of the specified dimensional score. Specifically, we define this task via the following conditional probability:

$$\mathcal{T}_{AB}^\alpha = \arg \max_{t \in \{I_1, I_2, \dots, I_n\}} p(t \mid \mathcal{V}, \mathcal{D}_{AB}^\alpha),$$

where $\{I_1, I_2, \dots, I_n\}$ is the candidate image sequence sampled from video \mathcal{V} .

Task 3: Visual Cue Causal Analysis.

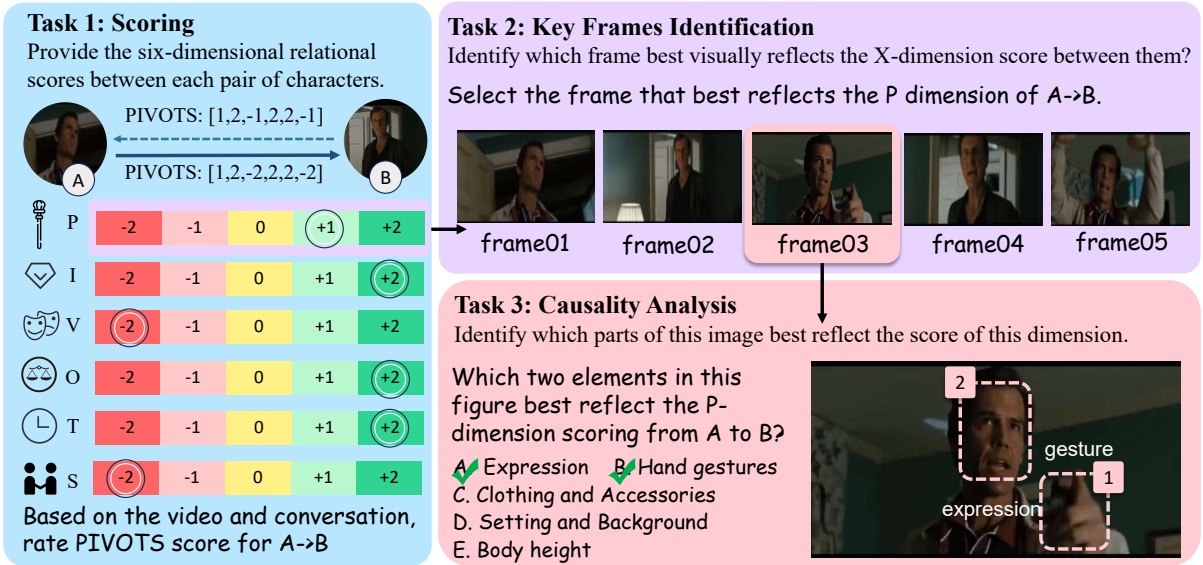


Figure 2: Overview of the three hierarchical tasks defined in the PIVOTS benchmark. The scoring task asks models to assign six-dimensional PIVOTS relational scores from the video and conversation. The key-frame identification task requires selecting the frame that best reflects a specified relational dimension, while the causality analysis task reasons about the visual elements (e.g., facial expressions or gestures) that support the assigned score, enabling causal grounding of social judgments.

Given a key frame \mathcal{T}_{AB}^α and direction \mathcal{D}_{AB}^α , MLLMs are prompted to select the visual cues that support the identification of the corresponding relationship dimension score. We formalize this task as the following conditional inference problem:

$$C_{AB}^\alpha = \arg \max_{c \in \{o_1, o_2, \dots, o_n\}} p(c | \mathcal{T}_{AB}^\alpha, \mathcal{D}_{AB}^\alpha), \quad (1)$$

where $\{o_1, o_2, \dots, o_n\}$ denotes a set of candidate visual cues extracted from the key frame. For a given query \mathcal{Q}_{AB}^α , the selected cue C_{AB}^α corresponds to the image evidence (e.g., facial expressions, hand gestures, or body posture) that supports the identification of \mathcal{R}_{AB}^α .

3.3 Data and Annotation

Data Source. To systematically study the tasks defined in our benchmark, we design our data curation strategy to prioritize scenarios that show prominent social signals and rich dynamics of interaction. We derive our dataset from two primary sources, ensuring a balance between interaction depth and a comprehensive breadth of social contexts:

Social-IQ 2.0: We selected a high-quality subset of 121 videos from the Social-IQ 2.0 dataset (Wilf et al., 2023), focusing on explicit social interactions. To strictly control for visual bias and information leakage, we manually filtered out news interviews and videos containing explicit visual markers (e.g.,

text overlays) that could directly hint at social roles. Each selected video sample is segmented into 30-second clips to ensure a consistent temporal context. The resulting clips capture nuanced human interactions (e.g., debates and casual conversations), providing rich evidence for analyzing multiple fine-grained social relationship dimensions.

YouTube Videos: Note that existing work on social relationship understanding often focuses on overt and expressive social behaviors, whereas many real-world spontaneous interactions are more subtle, restrained, and socially nuanced. To mitigate this limitation, we augment our corpus with YouTube data, structured across two primary dimensions:

- **Interaction Dynamics:** We target keywords associated with complex communicative goals, such as persuasion, and conflict resolution.
- **Relational Contexts:** The dataset spans a spectrum of social bonds, from intimate dyads (e.g., romantic and dissolved relationships) and familial ties (e.g., siblings) to zero-acquaintance dynamics among participants with divergent ideological backgrounds.

Notably, the interpersonal relationships in this curated subset are primarily reflected by the overall interaction atmosphere rather than specific body gestures within particular temporal segments. As a result, this subset is less suitable for keyframe identification or fine-grained causal analysis of in-

dividual visual cues. We make this trade-off to diversify the coverage of interpersonal relationship understanding. This additional data source comprises 70 videos and effectively complements the predominantly casual interactions in Social-IQ 2.0.

All data usage complies with the original licenses of Social-IQ 2.0 and YouTube’s Terms of Service; we release only video identifiers and annotations to avoid redistributing copyrighted or identifiable content.

Data Selection. Following the grounding principles, we filtered the raw collection to ensure multi-modal alignment. Specifically, we required that: (1) multiple speakers are visible and audible; (2) the interaction lasts longer than 30 seconds to allow for relationship development; and (3) the social dynamic is not interpretable with trivial symbols.

Annotations. We adopted a unified “model-first, human-calibrated” annotation framework across both video sources. We first construct reference evidence from the original question–answer annotations for the Social-IQ 2.0 dataset, and from retrieved video-based context for the YouTube dataset. Concretely, we consolidated the reference evidence, basic video metadata, and the PIVOTS six-dimensional definitions with their rating guidelines into a single prompt, and first used the GPT-5 API to generate initial per-character scores across all six dimensions. We then treated these model-generated ratings as a starting point and performed systematic human annotation, ensuring that each score strictly aligns with the dimension definitions. For each dimension, we additionally annotated key timestamps and brief rationales directly grounded in observable cues (e.g., gestures and actions, facial expressions, body posture and distance, and explicit supportive/refusals/directives in dialogue), making the ratings traceable and interpretable. When necessary, we focused on the most interaction-salient characters (e.g., those dominating the exchange or driving relational changes) to improve annotation consistency. Additional annotation details are provided in Appendix C.

Inter-Annotator Agreement. We compute inter-annotator agreement using the ordinal form of Krippendorff’s alpha (Krippendorff, 2018). The overall Krippendorff’s α value is 0.8125, indicating that the dataset annotations are reliable. We refer to Appendix D for additional details.

Benchmark Statistics Fig. 3 presents the statistical distribution of the Social-IQ 2.0 data and the YouTube data. We leverage the GPT-5 API to auto-

matically annotate each video segment with Emotion Tone dimensions, including serious, neutral, positive, negative, and sarcastic. Notably, the distributions across the two data sources differ substantially: positive is the dominant category in Social-IQ 2.0, whereas serious is the primary category in the YouTube data. This difference aligns with our earlier discussion, as we intentionally source more in-depth interactions from YouTube to complement and diversify the benchmark.

Our final benchmark comprises 595 VQA pairs per PIVOTS dimension from Social-IQ 2.0 and 170 per dimension from YouTube. We further curate 483 additional VQA pairs for auxiliary tasks of key frame identification and visual cue causal analysis.

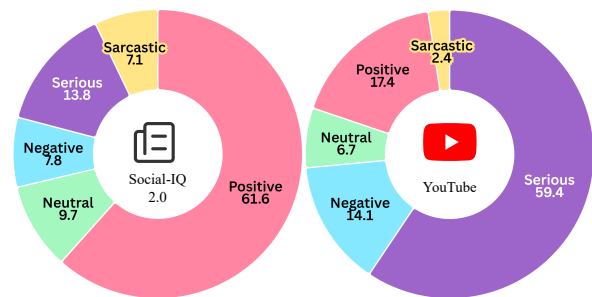


Figure 3: Label distributions across emotion and tone dimensions for Social-IQ 2.0 and YouTube videos.

4 Experiments

4.1 Experiment Setup and Metrics

We evaluate both open-source and proprietary MLLMs on our benchmark, including Gemini-2.5 (Comanici et al., 2025), GPT-5 (Achiam et al., 2023), and Qwen-3 models (Qwen-Team, 2025) of varying sizes. For input-modality ablations and additional analyses, we adopt the strong open-source model Qwen-32B-VL as the base model.

The scores of PIVOTS axes are defined on a five-point discrete scale from -2 to $+2$, yielding a chance-level accuracy of 20%. The key-frame identification and causal analysis tasks are formulated as five-option multiple-choice questions, which likewise result in a chance-level accuracy of 20%. For visual cue causal analysis, models are required to select the two most relevant cues from five candidates. Partial credit of 0.5 is awarded when one selected cue is correct, while full credit of 1.0 is given when both selected cues are correct. Under random selection, the expected chance-level accuracy is 40%.

Methods	Task 1							Task 2	Task 3
	Avg	P	I	V	O	T	S		
<i>Proprietary models</i>									
Gemini-2.5-pro	47.71	61.30	50.77	33.75	51.63	46.98	41.86	46.63	48.00
GPT-5	56.91	65.12	60.31	55.73	42.33	56.59	61.40	62.15	60.25
<i>Open-sourced models</i>									
Qwen3-4B-vl	24.33	51.05	20.77	23.99	13.73	25.93	10.47	28.59	30.17
Qwen3-8B-vl	27.90	59.74	22.22	22.87	18.58	26.89	17.07	34.75	31.00
Qwen3-32B-vl	34.54	62.62	23.82	33.82	23.77	30.28	31.50	35.29	33.50

Table 1: *Evaluation results of prevailing MLLMs on the Social-IQ 2.0 subset.* Proprietary models outperform open-source models by a substantial margin across all tasks.

Methods	Avg	P	I	V	O	T	S
<i>Proprietary models</i>							
Gemini-2.5-pro	54.66	69.49	55.08	31.78	61.02	54.24	56.36
GPT-5	58.88	73.71	50.44	48.28	61.14	60.68	58.85
<i>Open-sourced models</i>							
Qwen3-4B-vl	30.03	71.98	16.81	33.62	6.90	29.31	21.55
Qwen3-8B-vl	35.19	66.67	26.85	31.48	18.52	32.41	35.19
Qwen3-32B-vl	40.81	67.92	32.38	34.91	29.47	35.85	43.14

Table 2: *Evaluation results of prevailing MLLMs on the YouTube subset.*

4.2 Experimental Results

Benchmark Performance. We report the experimental results of prevailing MLLMs on our benchmark in Tab. 1 and Tab. 2. In terms of PIVOTS dimensions scoring, proprietary models achieve substantially stronger performance across all dimensions on both Social-IQ 2.0 and YouTube. Although open-source models often achieve on-par performance with proprietary models on prevailing MLLM benchmarks, they still struggle with fine-grained social dynamics reasoning, underscoring the importance for our benchmark.

We also observe that MLLMs, especially for Gemini, achieve better performance on the YouTube subset than on the Social-IQ 2.0 subset. One likely reason is that YouTube data has been extensively used during model pretraining. Among the six axes, we find that the valence-v and objective-o dimensions are notably more challenging for all models. The difficulty of valence prediction is exacerbated by its reliance on subtle affective cues such as tone, which is not captured by most current multimodal large language models. And objective judgments become unreliable when task-driven interactions are accompanied by strong emotional expressiveness.

For Task 2 and Task 3, Gemini and GPT also consistently outperform Qwen3. In particular, on the visual cue causal analysis task, Qwen3 often performs below the chance-level baseline, indicating difficulty in identifying the relevant visual cues required for reasoning about interpersonal relationship dimensions. This result suggests a notable limitation in the Qwen series’ ability to perform social behavior-related visual reasoning.

Ablation on Input Modalities. We conduct detailed ablation studies in Tab. 3 to analyze how different input modality configurations and text redaction settings affect multimodal interpersonal relationship reasoning.

It is possible that models leverage explicit social relationship identifiers as shortcuts for inferring interpersonal relationships. Such identifiers refer to textual components in the original utterances that directly reveal participants’ identities or relational attributes, including but not limited to kinship terms (e.g., father, daughter) and professional titles or honorifics (e.g., Your Honor, boss). In this context, we consider a baseline that uses redacted conversation utterances as inputs, in which explicit social identifiers are replaced with generic tokens that contain no social attribute information (e.g., [PERSON], [TITLE]), while preserving the syntactic structure, emotional tone, and logical flow of the conversation. As shown in Tab. 3, redacting explicit social identifiers results in a performance decrease of 2.47% on Social-IQ 2.0 and 4.17% on the YouTube subset relative to using the original utterances. This result suggests that explicit social roles provide useful cues for scoring PIVOTS dimensions. However, it also indicates that models do not rely solely on such identifiers as shortcuts for making predictions.

We further assess the contribution of visual

Methods	Social-IQ 2.0							YouTube						
	Avg	P	I	V	O	T	S	Avg	P	I	V	O	T	S
Video Only	33.70	62.68	24.09	30.62	22.77	30.92	30.08	31.18	80.37	13.33	34.26	15.38	16.19	26.42
Original Conv.	31.44	56.84	23.52	24.80	26.82	26.25	30.43	37.81	60.19	48.15	16.67	36.11	32.41	33.33
Redacted Conv.	28.97	55.06	21.73	24.31	21.92	25.41	25.41	33.64	53.64	47.27	20.00	24.55	26.36	30.00
Video + Original Conv.	34.54	62.62	23.82	33.82	23.77	30.28	31.50	40.81	67.92	32.38	34.91	29.47	35.85	43.14
Video + Redacted Conv.	35.83	60.86	22.46	42.97	24.78	28.77	32.95	32.22	47.00	29.35	25.29	22.08	25.53	41.11

Table 3: *Ablation studies on input modalities.* We examine the effects of explicit social role identifiers in conversational utterances and visual inputs on model performance.

Methods	Social-IQ 2.0							YouTube						
	Avg	P	I	V	O	T	S	Avg	P	I	V	O	T	S
Separate Prediction	34.54	62.62	23.82	33.82	23.77	30.28	31.50	40.81	67.92	32.38	34.91	29.47	35.85	43.14
Joint Prediction	43.26	62.32	41.60	45.19	22.51	33.77	54.16	38.63	82.24	29.91	44.86	11.21	30.84	32.71
Pairwise Prediction	42.19	63.55	31.87	50.18	19.41	33.70	54.40	39.67	86.96	28.26	42.39	15.22	28.26	36.96

Table 4: *Comparison of prediction settings for interpersonal relationship scoring.* We evaluate separate, joint, and pairwise prediction settings, analyzing their effects on PIVOTS dimension scoring performance.

modalities. On the Social-IQ 2.0 set, incorporating video information improves performance by 3.1% over the original utterance baseline and 6.86% over the redacted utterance baseline. This suggests that visual modalities contribute more strongly to interpersonal reasoning in Social-IQ 2.0 scenarios with salient social behaviors, especially in the absence of explicit social role cues. For the YouTube subset, incorporating video input yields a 3.0% performance gain when using the original utterances, but results in a 1.42% performance decrease when combined with redacted utterances. The divergence can be attributed to how we curated the YouTube subset. As discussed in Sec. 3.3, these videos primarily capture in-depth, restrained interactions in which interpersonal relationships are conveyed through global discourse context and shared background rather than salient, localized visual cues. When explicit role information is removed, models struggle to anchor visual observations to participants’ identities and intentions, making it difficult to integrate visual signals with conversational context for reliable interpersonal state inference.

4.3 Additional Analysis

Prediction Setting. We consider three different prediction settings for scoring each interpersonal relationship dimension:

- **Separate prediction:** This is the baseline setting that scores each PIVOTS dimension independently, treating the perception of one subject toward another along a single dimension as an

individual question.

- **Joint prediction:** The model is prompted to jointly rate six PIVOTS dimensions for one subject’s perception of another.
- **Pairwise prediction:** The model is required to predict all bidirectional PIVOTS dimension scores between two subjects.

The exact prompts used for each prediction setting are provided in Appendix E.

The experiments are conducted with videos and original utterances as inputs, and the results are summarized in Table 4. Notably, the effectiveness of prediction settings exhibits dataset-specific differences: On the Social-IQ 2.0 dataset, which primarily features casual, everyday interactions, joint prediction achieves the best performance (Avg 43.26%), followed by pairwise prediction (42.19%), both significantly outperforming separate prediction (34.54%). This aligns with the mechanism of human social perception, as joint prediction encourages the model to reason based on the shared latent representation of interactions, enabling cues from different dimensions to complement each other and enhance the coherence and accuracy of social dynamics understanding.

In contrast, an opposite trend is observed on the YouTube dataset, which contains in-depth interactions such as in-depth interviews: Separate prediction yields the optimal average performance (40.81%), outperforming joint prediction (38.63%) and pairwise prediction (39.67%). However, pair-

Methods	Social-IQ 2.0							YouTube						
	Avg	P	I	V	O	T	S	Avg	P	I	V	O	T	S
Baseline	34.54	62.62	23.82	33.82	23.77	30.28	31.50	40.81	67.92	32.38	34.91	29.47	35.85	43.14
Multi-Stage Prompt	41.23	60.06	39.45	46.70	26.82	26.13	48.15	31.48	37.96	44.44	30.56	22.22	16.67	37.04
In-Context Prompt	40.51	58.96	43.39	41.95	25.41	31.06	41.88	32.40	60.75	42.99	34.26	12.96	16.82	26.85

Table 5: Comparison of prediction settings for interpersonal relationship scoring. We evaluate separate prediction and separate prediction with CoT prompting prediction using different prompt settings, analyzing their effects on PIVOTS dimension scoring performance.

wise prediction performs exceptionally well on the P-dimension (86.96%), while separate prediction maintains the best performance across dimensions including I, O, T, and S. This indicates that in in-depth social interactions, holistic modeling (joint/pairwise prediction) may interfere with the model’s judgments on specific dimensions, whereas independent scoring is more conducive to capturing the unique cues of each dimension.

Heuristic Injection via Prompting. As evidenced by the experimental results in the preceding sections, existing MLLMs, particularly open-source models, often struggle to accurately score the PIVOTS dimensions. We therefore investigate how interpersonal relationship reasoning can be enhanced by drawing inspiration from human reasoning strategies. Concretely, we inject human reasoning processes into the prompt to guide MLLMs toward more accurate predictions. Specifically, we consider the following two prompting strategies:

- **Multi-Stage Prompting** first calibrates affective judgments (valence and intensity), then evaluates structural interpersonal dimensions (power, objective, permanence, and stance), followed by independent synthesis and reconciliation to ensure reliable PIVOTS predictions.
- **In-Context Prompting** grounds reasoning in fine-grained behavioral and contextual cues while explicitly modeling directional asymmetry between interacting individuals.

The exact prompts are provided in Appendix E.

The experimental results for different prompting strategies are reported in Table 5. Notably, the multi-stage and in-context prompting methods improve overall performance on the Social-IQ 2 dataset by 6.69% and 5.97%, respectively. These gains primarily stem from improvements in the Involvement and Valence dimensions, benefiting from affective calibration in the first stage of multi-stage prompting and from the explicit encoding of salient multimodal cues from in-context prompting.

For dimensions that involve frequent asymmetric perspectives between interacting subjects, such as Objective, Power, and Stance, performance gains remain limited. This suggests that prompting strategies fail to address a key bottleneck of existing MLLMs in grounding multimodal cues to the correct social participant during complex interactions.

On the YouTube dataset, we observe a contrasting trend compared to the Social-IQ 2.0 dataset. Both prompting strategies lead to a notable decline in overall performance. In particular, multi-stage prompting results in the largest drop on the Power dimension, while in-context prompting degrades performance on the Objective and Stance dimensions. This observation indicates that, in more conservative and in-depth social interactions, relying on rigid calibration or fixating on specific multimodal cues may instead confuse models during interpersonal relationship reasoning. These experimental results highlight the necessity of mining additional YouTube data. This subset diversifies interaction settings, thereby motivating research in computational social AI models that do not overfit to only salient and explicit interaction scenarios.

5 Conclusions

In this work, we introduce PIVOTS, the first benchmark for evaluating multimodal large language models on interpersonal relationship reasoning, together with auxiliary tasks that assess visual cue understanding. We first benchmark several prevailing MLLMs and then provide a detailed analysis of input modalities, prediction settings, and prompting strategies. Our experimental results show that explicit social scenarios and in-depth social scenarios pose significantly different challenges for computational social intelligence models. We believe our work provides a foundational step toward understanding fine-grained social dynamics and points to future research directions for designing robust and generalizable social AI models.

6 Limitations

Although our benchmark introduces the first fine-grained multimodal benchmark for interpersonal relationship reasoning, several limitations remain.

First, the current benchmark is limited to English-language data, which constrains its applicability to multilingual and cross-lingual social reasoning settings.

Second, although our annotation protocol is carefully designed, all annotators share similar cultural and linguistic backgrounds and are not native English speakers, which may introduce cultural bias in the interpretation of interpersonal social relationships, particularly for subjective dimensions.

Finally, we do not investigate visual instruction-tuning or reinforcement-learning methods for fine-grained social reasoning, as these approaches require high-quality supervision that remains costly and difficult to scale. Designing scalable data engines for obtaining reliable social reasoning supervision is an important and exciting direction for future work, and our benchmark provides a foundation for such efforts.

7 Ethical Considerations

Notably, our benchmark involves the analysis of human subjects and therefore raises potential privacy concerns. For data sourced from Social-IQ 2.0, we follow the original dataset license, which permits use for research purposes. For YouTube videos, we adopt ethical practices consistent with prior work by releasing only video identifiers and temporal annotations, rather than distributing raw video content. We are committed to respecting individual privacy and ensuring compliance with established ethical standards in dataset creation, release, and usage. In addition, because our benchmark focuses on interpersonal relationship understanding, care must be taken to prevent potential misuse, such as deception, biased interpretations, or emotional manipulation. To mitigate these risks, we frame the benchmark as an analysis and evaluation tool rather than a deployable system, provide clear documentation on intended use and limitations, and discourage deployment in high-stakes or manipulative settings without human oversight.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,

Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.

Mica R Endsley. 2021. Situation awareness in future autonomous systems. *Journal of Cognitive Engineering and Decision Making*, 15(1):13–17.

Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, and 1 others. 2022. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18995–19012.

A. M. Hoogeboom and C. P. Wilderom. 2020. A closer look at the interaction patterns of effective teams: A sequential analysis of the behavior of formal team leaders and members. *Journal of Organizational Behavior*, 41(2):176–190.

Chuanyang Jin, Yutong Wu, Jing Cao, Jiannan Xiang, Yen-Ling Kuo, Zhiting Hu, Tomer Ullman, Antonio Torralba, Joshua Tenenbaum, and Tianmin Shu. 2024. Mmtom-qa: Multimodal theory of mind question answering. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16077–16102.

Klaus Krippendorff. 2018. *Content analysis: An introduction to its methodology*. Sage publications.

Bolin Lai, Hongxin Zhang, Miao Liu, Aryan Pariani, Fiona Ryan, Wenqi Jia, Shirley Anugrah Hayati, James Rehg, and Diyi Yang. 2023. [Werewolf among us: Multimodal resources for modeling persuasion behaviors in social deduction games](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 6570–6588, Toronto, Canada. Association for Computational Linguistics.

Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. Camel: Communicative agents for "mind" exploration of large language model society. *Advances in Neural Information Processing Systems*, 36:51991–52008.

Junnan Li, Yongkang Wong, Qi Zhao, and Mohan S Kankanhalli. 2017. Dual-glance model for deciphering social relationships. In *Proceedings of the IEEE international conference on computer vision*, pages 2650–2659.

Junnan Li, Yongkang Wong, Qi Zhao, and Mohan S Kankanhalli. 2020. Visual social relationship recognition. *International Journal of Computer Vision*, 128(6):1750–1764.

701	Xinchen Liu, Wu Liu, Meng Zhang, Jingwen Chen,	Alex Wilf, Leena Mathur, Sheryl Mathew, Claire	757
702	Lianli Gao, Chenggang Yan, and Tao Mei. 2019. So-	Ko, Youssouf Kebe, Paul Pu Liang, and Louis-	758
703	cial relation recognition from videos via multi-scale	Philippe Morency. 2023. Social-iq 2.0 chal-	759
704	spatial-temporal reasoning. In <i>Proceedings of the</i>	lenge: Benchmarking multimodal social understand-	760
705	<i>IEEE/CVF conference on computer vision and pat-</i>	ing. https://github.com/abwilf/Social-IQ-2.	761
706	<i>tern recognition</i> , pages 3566–3574.	0-Challenge.	762
707	Xinyi Mou, Jingcong Liang, Jiayu Lin, Xinnong Zhang,	Myron Wish, Morton Deutsch, and Susan J Kaplan.	763
708	Xiawei Liu, Shiyue Yang, Rong Ye, Lei Chen,	1976. Perceived dimensions of interpersonal rela-	764
709	Haoyu Kuang, Xuan-Jing Huang, and 1 others. 2025.	tions. <i>Journal of Personality and social Psychology</i> ,	765
710	Agentsense: Benchmarking social intelligence of lan-	33(4):409.	766
711	guage agents through interactive scenarios. In <i>Pro-</i>	Eric Zhongcong Xu, Zeyang Song, Satoshi Tsutsui,	767
712	<i>ceedings of the 2025 Conference of the Nations of</i>	Chao Feng, Mang Ye, and Mike Zheng Shou. 2022.	768
713	<i>the Americas Chapter of the Association for Comput-</i>	Ava-avd: Audio-visual speaker diarization in the wild.	769
714	<i>tational Linguistics: Human Language Technologies</i>	In <i>Proceedings of the 30th ACM International Con-</i>	770
715	(Volume 1: Long Papers), pages 4975–5001.	<i>ference on Multimedia</i> , pages 3838–3847.	771
716	Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Mered-	Zhanpeng Zhang, Ping Luo, Chen-Change Loy, and	772
717	ith Ringel Morris, Percy Liang, and Michael S Bern-	Xiaou Tang. 2015. Learning social relation traits	773
718	stein. 2023. Generative agents: Interactive simulacra	from face images. In <i>Proceedings of the IEEE in-</i>	774
719	of human behavior. In <i>Proceedings of the 36th an-</i>	<i>ternational conference on computer vision</i> , pages	775
720	<i>annual acm symposium on user interface software and</i>	3631–3639.	776
721	<i>technology</i> , pages 1–22.	Jinfeng Zhou, Yuxuan Chen, Yihan Shi, Xuanming	777
722	Qwen-Team. 2025. Qwen3 technical report . <i>Preprint</i> ,	Zhang, Leqi Lei, Yi Feng, Zexuan Xiong, Miao Yan,	778
723	arXiv:2505.09388.	Xunzhi Wang, Yaru Cao, and 1 others. 2025. Social-	779
724	Farzana Rashid and Eduardo Blanco. 2017. Dimen-	eval: Evaluating social intelligence of large language	780
725	sions of interpersonal relationships: Corpus and experi-	models. <i>arXiv preprint arXiv:2506.00900</i> .	781
726	ments. In <i>Proceedings of the 2017 Conference on</i>	Hao Zhu, Bodhisattwa Prasad Majumder, Dirk Hovy,	782
727	<i>Empirical Methods in Natural Language Processing</i> ,	and Diyi Yang. 2025. Social intelligence in the age of	783
728	pages 2307–2316.	llms. In <i>Proceedings of the 2025 Annual Conference</i>	784
729	James A Russell. 1980. A circumplex model of af-	<i>of the Nations of the Americas Chapter of the Asso-</i>	785
730	fect. <i>Journal of personality and social psychology</i> ,	<i>ciation for Computational Linguistics: Human Lan-</i>	786
731	39(6):1161.	<i>guage Technologies (Volume 5: Tutorial Abstracts)</i> ,	787
732	Maarten Sap, Hannah Rashkin, Derek Chen, Ronan	pages 51–55.	788
733	LeBras, and Yejin Choi. 2019. Socialliqa: Com-	A Interpersonal Relationship Definition	789
734	monsense reasoning about social interactions. <i>arXiv</i>	and Examples	790
735	<i>preprint arXiv:1904.09728</i> .	A.1 Egalitarian vs. Hierarchical Power	791
736	Tianmin Shu, Abhishek Bhandwadar, Chuang Gan,	Definition: Assess whether the relationship is per-	792
737	Kevin Smith, Shari Liu, Dan Gutfreund, Elizabeth	ceived as equal or unequal in terms of social influ-	793
738	Spelke, Joshua Tenenbaum, and Tomer Ullman. 2021.	ence, decision-making agency, and status.	794
739	Agent: A benchmark for core psychological reason-	• +2, Extreme Inequality(Self-Dominant) : You	795
740	ing. In <i>International conference on machine learning</i> ,	hold absolute agency. The other party’s role	796
741	pages 9614–9625. PMLR.	is reduced to compliance, and their input is	797
742	Qianru Sun, Bernt Schiele, and Mario Fritz. 2017. A	neither sought nor required to move the in-	798
743	domain based approach to social relation recognition.	teraction forward. Behavioral Observables:	799
744	In <i>Proceedings of the IEEE conference on computer</i>	Dominating the spatial environment, intrusive	800
745	<i>vision and pattern recognition</i> , pages 3481–3490.	gaze, high frequency of interruptions, lack of	801
746	Yashar Talebirad, Ali Parsaee, Vishwajeet Ohal,	back-channeling; Psychological Observables:	802
747	Amirhossein Nadiri, Csongor Szepesvari, Yash	Empathy deficit, perception of the other as an	803
748	Mouje, and Eden Redman. 2025. Wisdom of the	instrument or object. Example: Judge sentenc-	804
749	machines: Exploring collective intelligence in llm	ing a defendant; Interrogator to a captive.	805
750	crowds. In <i>First Workshop on Social Simulation with</i>	• +1, Noticeable Inequality(Self-Dominant):	806
751	<i>LLMs</i> .	You lead the interaction through expertise or	807
752	Gang Wang, Andrew Gallagher, Jiebo Luo, and David		
753	Forsyth. 2010. Seeing people in social context: Rec-		
754	ognizing people and social relationships. In <i>Euro-</i>		
755	<i>pean conference on computer vision</i> , pages 169–182.		
756	Springer.		

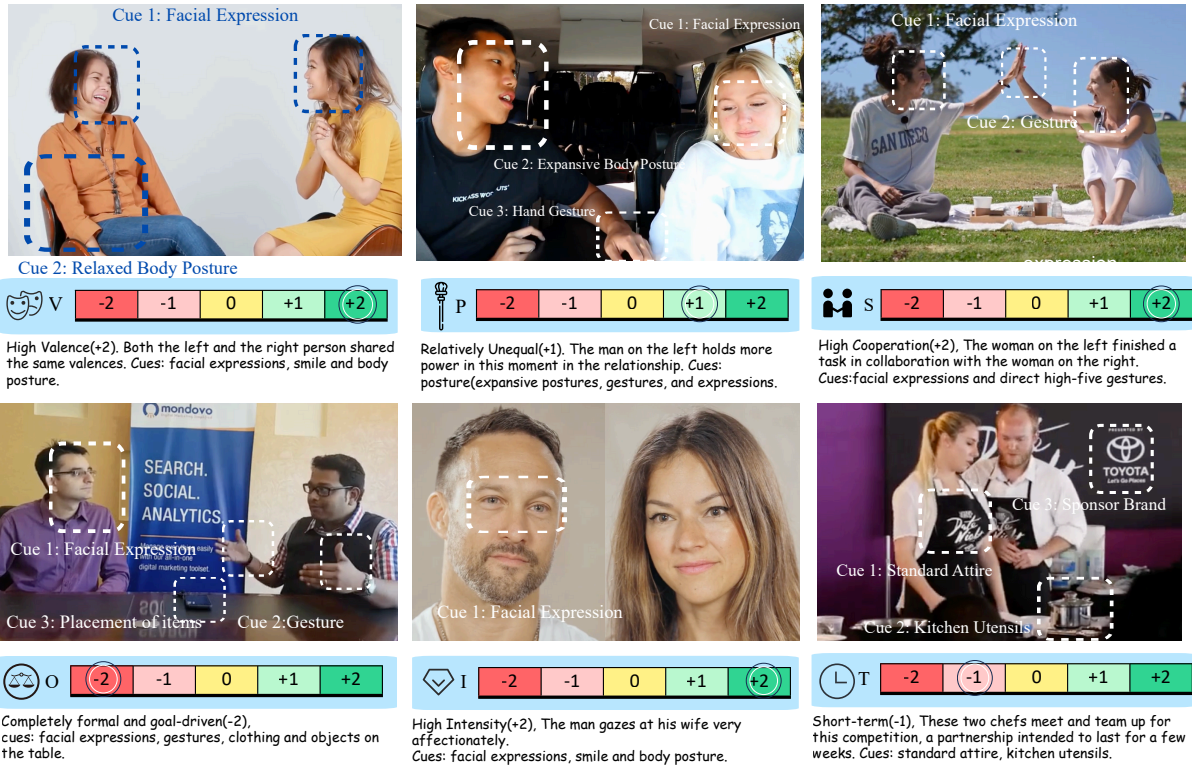


Figure 4: Visual scoring examples of the PIVOTS benchmark across various social dimensions.

status. While the other party has some autonomy, you hold the final decision-making power and define the agenda. Behavioral Observables: Setting the agenda, providing evaluative feedback, expansive but non-threatening gestures, asymmetric verbal output; Psychological Observables: Feeling "in charge" or responsible for the other's performance/well-being, mild cognitive load regarding the other's progress. Example: Senior surgeon to resident; Manager to an employee.

- 0, Egalitarian (Balanced Equality): Power is a shared resource. Neither party can impose their will without negotiation. Behavioral Observables: synchronized body language, mutual questioning, balanced turn-taking durations; Psychological Observables: High psychological safety, mutual respect, perception of the other as a peer. Example: Strategic partners; Co-authors; Spouses; Colleagues in a brainstorming session.
- -1, Noticeable Inequality (Other-Dominant): You recognize the other's lead and adjust your behavior to fit their frame. You seek to align with their expectations to maintain harmony or gain approval. Behavioral Observables:

Active listening (frequent nodding), cautious or tentative phrasing (e.g., hedging), self-minimizing posture, seeking non-verbal confirmation; Psychological Observables: Social monitoring, slight performance anxiety, slight cognitive load to "read" the other, high sensitivity to the other party's micro-expressions. Example: Employee to executive; Graduate student to professor; Applicant to recruiter.

- -2, Extreme Inequality (Other-Dominant): The other party holds absolute agency over your outcomes. You lack the means to resist or influence the direction of the encounter. Behavioral Observables: Avoidant eye contact, physical shrinking, fawning behaviors (e.g., excessive ingratiation), or long-term silence. Psychological Observables: High anxiety, hypervigilance, may experience "learned helplessness," feeling of being "I am invisible" or powerless. Example: Victim to oppressor.

A.2 Superficial vs. Intense Involvement

Definition: Indicate the depth of the relationship based on the level of emotional intimacy, self-disclosure, and psychological interdependence.

- +2 (Very Intense): The relationship is a core

859 part of one's self-identity, high vulnerability
 860 is present; the parties usually share lots of
 861 "private language" and deep history. Behavioral
 862 Observables: High affective resonance,
 863 soft gaze, micro-sync behavior in movement,
 864 high physical proximity; Psychological Ob-
 865 servables: Great amount of vulnerabilities,
 866 merged identity (use a lot of "we"), lots of
 867 trust, extreme empathy. Example: Soulmates;
 868 Close friends; Spouses; Lifelong partners.

- 869 • +1 (Intense): The interaction is driven by genu-
 870 ine emotional connection. You share per-
 871 sonal feelings and opinions beyond the re-
 872 quirements of the setting. Behavioral Ob-
 873 servables: Frequent self-disclosure, expres-
 874 sive affect, some level of physical proxim-
 875 ity; Psychological Observables: High empa-
 876 thy, perceived social support, feeling "known"
 877 and validated, moderate emotional interde-
 878 pendence. Example: friends; teammates for
 879 years.
- 880 • 0 (Neutral): The relationship is standard role-
 881 based. You are friendly and helpful, but you
 882 keep your private life behind a social mask.
 883 Behavioral Observables: Polite distance, stan-
 884 dard smiles, talk stays on the task or safe small
 885 talks (e.g., weather, work projects); Psycho-
 886 logical Observables: Feeling safe but not inti-
 887 mate, Low emotional risk-taking, cognitive fo-
 888 cus on goals. Example: Ordinary colleagues;
 889 Project-based partners.
- 890 • -1 (Superficial): You recognize each other's
 891 face but have no bond. Behavioral Observ-
 892 ables: The interaction is a quick ritual just to
 893 be polite; A quick nod or wave, standard greet-
 894 ings (e.g., "How's it going?"), keeping a wide
 895 physical distance. Psychological Observables:
 896 Very little mental energy spent on the other; a
 897 familiar face. Example: An acquaintance you
 898 meet somewhere before; a neighbor you only
 899 see a few times.
- 900 • -2 (Very superficial): The interaction is purely
 901 for a task. You see the other person as a "ser-
 902 vice provider" or a stranger, not as a human
 903 you need to know. Behavioral Observables:
 904 No intended eye contact, automatic or flat
 905 voice, focusing only on the exchange (money,
 906 information). Example: one-time cashier; ask-
 907 ing a stranger for directions.

A.3 Positive vs. Negative Valence 908

909 Definition: Assess the core emotional quality of
 910 the interaction, ranging from highly restorative and
 911 joyful to deeply distressing and hostile.

- 912 • +2 (Very pleasant): The interaction is excep-
 913 tionally safe, joyful, and revitalizing. Behav-
 914 ioral Observables: There is a sense of flow
 915 where both parties feel uplifted and deeply
 916 relaxed. Frequent smiles, shared laughter, re-
 917 laxed or open posture, warm vocal resonance.
 918 Psychological Observables: Feelings of en-
 919 joyment or deep peace, high psychological
 920 safety. Example: A celebratory dinner with
 921 loved ones, A great talk with someone who
 922 shares the same viewpoints.
- 923 • +1 (Pleasant): The interaction is friendly,
 924 smooth, and satisfying. The overall vibe is
 925 constructive and agreeable. Behavioral Ob-
 926 servables: Frequent nodding, warm smiles.
 927 Psychological Observables: feeling "liked",
 928 higher level of happiness and low level of so-
 929 cial anxiety or tension. Example: A friendly
 930 catch-up with a coworker; a great talk with a
 931 neighbor.
- 932 • 0 (Neutral): Emotionally neutral. The inter-
 933 action is emotionally dry or purely factual.
 934 There is no significant emotion included, but
 935 also no underlying tension or friction. Be-
 936 havioral Observables: Controlled facial ex-
 937 pressions, focusing on matter-of-fact, neutral
 938 body language. Psychological Observables:
 939 focus is purely on logic/tasks, no strong feel-
 940 ing toward the other person. Example: A rou-
 941 tine status update meeting at work; asking a
 942 stranger for the time.
- 943 • -1 (Unpleasant): Overall negative and unsatis-
 944 factory. The interaction feels awkward, tense,
 945 or slightly draining. You may feel a need
 946 to escape the situation or put up emotional
 947 walls. Behavioral Observables: Forced/fake
 948 smiles, defensive posture (crossed arms, lean-
 949 ing away), sighing, brief responses. Psycho-
 950 logical Observables: Mild anxiety, irritation,
 951 feeling uncomfortable or misunderstood, so-
 952 cial fatigue. Example: A forced conversation
 953 with a colleague; a tense meeting about a mis-
 954 take.
- 955 • -2 (Very unpleasant): The interaction is hostile,
 956 harmful, or aggressive. It triggers a "fight-

957	or-flight" response and leaves you feeling	to a meeting agenda, professional vocabulary,	1006
958	emotionally depleted or attacked. Behavioral	social zone physical distance. Psychological	1007
959	Observables: Shouting or harsh tone, sneering/	Observables: Focus on tasks, seeing the other	1008
960	disdainful expressions, cold silence or ag-	primarily as a "colleague" or "work partner".	1009
961	gressive gestures. Psychological Observables:	Example: A new project team; a manager and	1010
962	Extreme levels of fear, anger, or shame, feel-	a hire.	1011
963	ing significantly violated. Example: Being		
964	bullied or harassed.		
965	A.4 Socioemotional vs. Task-Oriented		
966	Objective		
967	Definition: Assess whether the interaction is driven		
968	by personal emotional needs (Informal) or by the		
969	requirements of a specific goal (Formal).		
970			
971	• +2 (Purely social-emotional): The interaction		
972	exists solely for the sake of the bond. There		
973	is no external work to be done. The primary		
974	goal is emotional exchange and shared ex-		
975	perience. Behavioral Observables: Sponta-		
976	neous dialogue, emotions included, informal		
977	slang. Psychological Observables: feeling of		
978	belonging, zero pressure to perform or feeling		
979	emotionally controlled by someone (negative		
980	valence case). Example: Family members re-		
981	laxing at home; best friends hanging out; a		
982	couple on a date.		
983			
984	• +1 (Tends to social-emotional): While there		
985	may be a shared context, the interaction is		
986	driven by personal liking. People may intend		
987	to side-track from tasks to talk about personal		
988	lives. Behavioral Observables: Use of first		
989	names, personal anecdotes, casual posture.		
990	Psychological Observables: feeling like "we		
991	are in this together" as people, not just as		
992	roles. Example: Friends; colleagues you like		
993	to talk with.		
994			
995	• 0 (Mixed): An equal balance of goal-striving		
996	and personal trust. The task is complex and		
997	requires deep mutual understanding to suc-		
998	ceed. Behavioral Observables: switching be-		
999	tween intense technical talk and personal		
1000	check-ins. Psychological Observables: high		
1001	interdependence, feeling that the person		
1002	and the task are equally important. Ex-		
1003	ample: Long-term business partners; a		
1004	director and a lead actor.		
1005			
1006			
1007			
1008			
1009			
1010			
1011			
1012			
1013			
1014			
1015			
1016			
1017			
1018			
1019			
1020			
1021			
1022			
1023			
1024			
1025			
1026			
1027			
1028			
1029			
1030			
1031			
1032			
1033			
1034			
1035			
1036			
1037			
1038			
1039			
1040			
1041			
1042			
1043			
1044			
1045			
1046			
1047			
1048			
1049			
1050			
1051			
1052			

straightforward expressions of care among relatives and friends. In contrast, the relationships depicted in the YouTube videos we selected are far more complex and defy simplistic categorization. These interpersonal connections are often intertwined with long-standing emotional bonds, unspoken past entanglements, or emotional tensions. For instance, consider former partners who loved each other for a decade, tentatively navigating the entanglement of love and hate; or long-lost siblings who have become strangers due to years of separation, yet cannot conceal their inherent affinity rooted in blood ties. Their interactions are permeated with subtle emotional undercurrents and underlying relational tensions, making them impossible to encapsulate with simplistic labels. As shown in Figure. 5, the upper part displays a case of a family filming a video together from the Social-IQ 2.0 dataset, which features a clear scenario and explicit interaction signals, while the lower part presents a case of a loving couple talking about death and inevitable separation from the YouTube dataset.

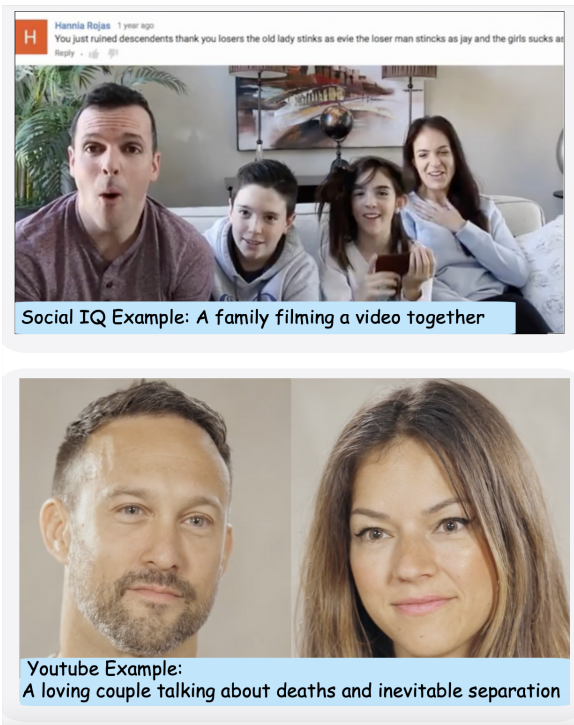


Figure 5: Example on Comparisons between Youtube and Social-IQ 2.0 videos.

Both datasets have a video duration of approximately 30 seconds. Figure. 6 illustrates the statistical distributions of Social-IQ 2.0 and YouTube data across the three major categories. We use the GPT-5 API to automatically annotate each

video segment with fine-grained subcategory labels, yielding the comparisons shown below. Clear differences can be observed. For Content Purpose, Social-IQ 2.0 is dominated by Documentary content (approximately 61.8%), whereas Interview content is predominant in the YouTube data (approximately 68.3%). For Core Topic and Theme, Social-IQ 2.0 places greater emphasis on Life and Work (approximately 52.2%); in contrast, the YouTube dataset assigns higher proportions to Society (approximately 37.9%) as well as health- and family-related topics, indicating that YouTube samples include more socially oriented and in-depth interview content. Regarding Relationships, Social-IQ 2.0 primarily features interactions among Friends (approximately 27.0%) and with the Audience (approximately 28.5%), whereas YouTube exhibits a higher proportion of Romantic partners (approximately 32.7%) together with Audience interactions (approximately 29.7%). Overall, the two datasets show substantial distributional differences across all three dimensions: Social-IQ 2.0 focuses more on documentary-style narratives and everyday social interactions, while YouTube, through interview-driven and public- or intimate-relationship-oriented content, provides deeper and more diverse social interaction samples that complement the benchmark’s coverage.

C Annotation Details

For the two distinct data sources, we have designed differentiated annotation workflows based on their data characteristics and annotation objectives.

C.1 Social-IQ 2.0 and YouTube Videos

The Social-IQ 2.0 (hereinafter referred to as SIQ2) dataset contains a large number of QA pairs, and our annotation takes the QA pairs corresponding to each video as the core reference benchmark. These files include key details of interpersonal interactions, in-depth social information, contextual background of scenarios, and excerpts of core dialogues, which can provide reliable support for API-based scoring.

Unlike the SIQ2 dataset, which relies on QA pairs for annotation, the YouTube dataset takes a different approach by focusing directly on video content and the associated social interaction details. Instead of predefined QA pairs, the YouTube dataset emphasizes analyzing interpersonal interactions within the video. Annotators assess these in-

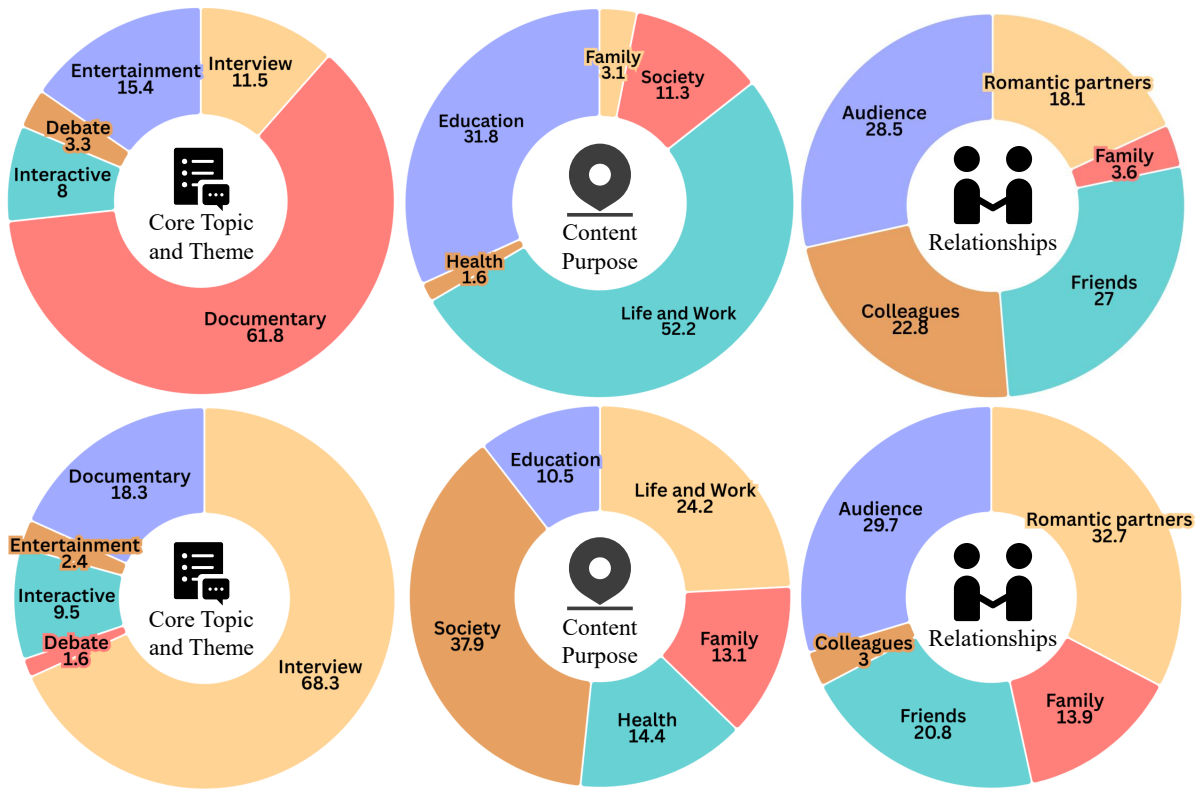


Figure 6: Comparison of Distributions between Social-IQ 2.0 (Top Row) and YouTube (Bottom Row) across Three Major Categories: Content Purpose, Core Topic and Theme, and Relationships

1226 interactions and assign scores to the characters based
 1227 on their roles, without relying on QA pairs. The
 1228 key interactions and scenes are carefully recorded
 1229 to ensure an accurate and detailed understanding
 1230 of the video content.

1231 Specifically, we integrate the complete QA an-
 1232 notation content corresponding to the video, basic
 1233 video information, detailed dimension definitions
 1234 of the PIVOTS six-dimensional framework, and
 1235 explanations of scoring gradients as prompts. Af-
 1236 ter inputting these prompts into the GPT-5 API,
 1237 the model can output the initial scores of the char-
 1238 acters in the video across the six dimensions of
 1239 P(Egalitarian vs. Hierarchical Power), I (Super-
 1240 ficial vs. Intense Involvement), V (Positive vs.
 1241 Negative Valence), O (Socioemotional vs. Task-
 1242 Oriented Objective), T (Temporary vs. Enduring
 1243 Permanence), and S (Cooperative vs. Competitive
 1244 Stance).

1245 Considering the potential dimension matching
 1246 biases in API-generated scores, we subsequently
 1247 conducted systematic manual correction work, and
 1248 the correction process follows the content outlined
 1249 below:

- Score Verification Specific scores have been

1251 generated by the API, so annotators do not
 1252 need to re-score. They only need to verify
 1253 the consistency between the scores and the
 1254 PIVOTS dimension standards, with a focus on
 1255 examining the rationality of extreme scores
 1256 (+2/-2). It is necessary to confirm whether the
 1257 interpersonal interaction characteristics cor-
 1258 responding to these scores fully match the
 1259 definitions in the framework, and manually
 1260 adjust any incorrect scores.

- Timestamp and Rationale Annotation During
 1261 the score verification process, timestamp posi-
 1262 tioning and rationale selection are performed
 1263 for the video. Timestamps need to be accu-
 1264 rately located at the video frames that best
 1265 support the scores, and the rationale explana-
 1266 tions must clearly link to the visual evi-
 1267 dence in these frames, including non-verbal cues
 1268 such as the characters' gestures, facial expres-
 1269 sions, and body postures. 1270
- Character Screening Criteria To focus on core
 1271 interpersonal interactions, a maximum of 3
 1272 characters per video are retained for annotation.
 1273 If the number of characters in a video
 1274

exceeds 3, screening is conducted in accordance with the following principles: priority is given to retaining characters with strong dominance and key roles in the development of relationships; if there are multiple supporting roles but only 2 core characters, there is no need to forcefully supplement to 3 characters, and only the core subjects are retained.

C.2 Annotation Interface

The annotation interface (Figure. 16) is designed to support the systematic scoring of dimensional relationships between characters in videos and to associate each score with an evidence timestamp, thereby improving annotation efficiency and consistency. The core functionalities are detailed as follows:

- Controls & Import:** Upon launching the interface, users can import a video file (.mp4), a transcript file (.txt), and a JSON file (.json). The JSON file represents the structured output produced by the API after initial scoring, allowing users to perform subsequent corrections and refinements. By default, the system prioritizes loading the file nicknamed “score.json”.
- Keyframe Gallery & Navigation:** To facilitate precise timestamp selection, the interface offers two mechanisms. First, the “Extract All (3 fps)” function uniformly samples the video at 3 frames per second, displaying extracted frames in the “Keyframe Gallery” for rapid browsing. Second, the “Select Folder” option allows users to choose a directory, automatically retrieving all video files within it for efficient switching.
- Transcript Reference:** This module allows users to inspect subtitles or transcripts, providing textual context to support annotation decisions.
- Characters Management:** This section lists all entities defined in the JSON file. Users can add, edit, or delete characters, with modifications immediately propagating to the dimensional relationship entries.
- Interaction Annotations:** This is the core annotation workspace, enumerating all pairwise dimensional relationships. Users review social dynamics, revise relationship scores, and

Dimension	Exact Agreement (%)	Krippendorff’s α
P	86.23%	0.6434
I	83.48%	0.8402
V	81.98%	0.7693
O	79.97%	0.7803
T	87.59%	0.8395
S	83.96%	0.7906
Overall	83.87%	0.8125

Table 6: Exact agreement rates and Krippendorff’s α per dimension.

assign a specific timestamp (the frame best reflecting the dimension’s characteristics). Additionally, users can document justifications for their selections using standardized categories: (A) facial expression, (B) body posture, (C) proximity, (D) body height, (E) hand gesture, (F) social norms, (G) gaze, (H) context, (I) attire, and (J) contact.

- Export:** After completing the annotation, users can verify the output via “json preview” and finalize the process by clicking “export data” to save the scored JSON file.

C.3 Annotator Recruitment and Payment

All annotation work in PIVOTS Bench was completed by the authors of this paper, with no external crowdworkers or paid annotators involved in the annotation process. All annotators were fully aware that the annotation purpose was solely for academic research. Since the annotation was conducted internally by the research team, no recruitment or payment procedures were required.

D Annotation Agreement Analysis

To assess annotation quality, we measured inter-annotator reliability on the valid samples using Krippendorff’s α . Implementation details are as follows: in Python we converted the three annotator columns to numeric, arranged the data into an annotator \times item matrix with masked missing values, and then computed α via `krippendorff.alpha`. We obtained $\alpha = 0.8125$, which exceeds the common 0.80 threshold for high reliability, indicating strong annotation agreement; therefore, the annotated dataset can be considered a reliable resource for downstream analyses and model training.

Table 6 summarizes the inter-annotator agreement results across six dimensions. Overall, all dimensions exhibit high exact agreement rates, with

an aggregated exact agreement of 83.87%. In addition, the overall Krippendorff’s α reaches 0.8125, indicating good inter-annotator reliability when the ordinal structure of the labels is taken into account. At the dimension level, most dimensions (I, T, V, O, and S) achieve Krippendorff’s α values close to or above 0.78, demonstrating stable and reliable annotation quality.

Although the P dimension shows a relatively lower Krippendorff’s α of 0.6434, its exact agreement rate reaches 86.23%, which is among the highest across all dimensions. This indicates that annotators provided identical judgments for the vast majority of samples on the P dimension. Importantly, the lower Krippendorff’s α for P is closely related to its highly imbalanced label distribution. Specifically, more than 70% of the samples in the P dimension are annotated as the neutral category (0), resulting in a strongly concentrated class distribution. Under such conditions, even a small number of disagreements can disproportionately reduce the Krippendorff’s α coefficient, despite high overall agreement.

This behavior is a known characteristic of Krippendorff’s α in settings with severe class imbalance: when most instances fall into a single category, the relative impact of the few disagreements is amplified, leading to a lower Krippendorff’s α value that does not necessarily indicate poor overall reliability. Considering the high exact agreement rate of the P dimension (86.23%), we conclude that the P annotations remain consistent and stable in practice, despite the lower Krippendorff’s α .

For the data used in Tasks 2 and 3, we applied strict filtering criteria to ensure annotation quality. Specifically, we retained only those samples with annotation timestamps no greater than 2 seconds and for which all annotators selected the same element; annotations that failed to meet either condition were excluded from subsequent analyses.

E Prompt Usage for Model Inference

We use OpenAI’s ChatCompletion multimodal prompt format as the base format. This prompt consists of two parts: one is a multimodal upload link; the other is a "user prompt" that includes a PIVOTS instance and a PIVOTS question.

The prompt template for separate prediction in Task 1 is shown in Figure 7. The prompt template for joint prediction in Task 1 is shown in Figure 8. The prompt template for pairwise prediction

```

Prompt Example

[
  {
    "role": "user",
    "content": [
      {
        "type": "text",
        "text": "You need to answer the question based on the provided information."
      }
    ]
  }
]

Input Information:
1. Video
2. Transcribed Text
3. the definition of PIVOTS (a six-dimensional interpersonal relationship framework)

PIVOTS - Six-Dimensional Model [definition]

Note that all the above scores are unidirectional; for scores on the same dimension, A's score for B is not necessarily the same as B's score for A. Directly output your score. And wrap your score with "###" before and after. If your score is -2, then you should output ###-2###
Transcribed Text [data1]
Question [data2] ",
  {
    "type": "image_url",
    "image_url": {
      "url": ""
    }
  }
]
]

```

Figure 7: Task 1 Prompt with separate prediction

```

Prompt Example

[
  {
    "role": "user",
    "content": [
      {
        "type": "text",
        "text": "You need to answer the question based on the provided information."
      }
    ]
  }
]

Input Information:
1. Video
2. Transcribed Text
3. the definition of PIVOTS (a six-dimensional interpersonal relationship framework)
PIVOTS - Six-Dimensional Model [definition]
Note that all the above scores are unidirectional; for scores on the same dimension, A's score for B is not necessarily the same as B's score for A. Directly output your score for each dimension in the order of PIVOTS, enclosed in "[ ]", with each score separated by a comma—for example: [-2,0,1,2,0,1] (each number corresponds sequentially to the PIVOTS dimensions). Wrap your score with "###" before and after. For example, if your score is [-2,1,0,2,1,1], output directly: ###[-2,1,0,2,1,1]###
Transcribed Text [data1]
Question [data2] ",
  {
    "type": "image_url",
    "image_url": {
      "url": ""
    }
  }
]
]

```

Figure 8: Task 1 Prompt with joint prediction

```

Prompt Example

[
  {
    "role": "user",
    "content": [
      {
        "type": "text",
        "text": "You need to answer the question based on the provided information."
      }
    ]
  },
  {
    "role": "user",
    "content": [
      {
        "type": "text",
        "text": "Input Information:"
      },
      {
        "type": "text",
        "text": "1. Video"
      },
      {
        "type": "text",
        "text": "2. Transcribed Text"
      },
      {
        "type": "text",
        "text": "3. the definition of PIVOTS (a six-dimensional interpersonal relationship framework)"
      }
    ]
  },
  {
    "role": "user",
    "content": [
      {
        "type": "text",
        "text": "PIVOTS - Six-Dimensional Model [definition]"
      },
      {
        "type": "text",
        "text": "Note: Note that all the above scores are unidirectional; for scores on the same dimension, A's score for B is not necessarily the same as B's score for A. For each unidirectional score, directly output your score for each dimension in the order of PIVOTS, enclosed in \"[]\", with each score separated by a comma—for example: [-2,0,1,2,0,1] (each number corresponds sequentially to the PIVOTS dimensions). Wrap your score with \"###\" before and after."
      },
      {
        "type": "text",
        "text": "For example, if the question is: \"Sequentially provide the unidirectional PIVOTS scores from Woman to Man and from Man to Woman,\" you should first output the score from Woman to Man, followed by the score from Man to Woman. If your scores are [-2,0,1,2,1,1] for Woman → Man and [1,0,2,1,0,0] for Man → Woman, your output should be: ###[-2,0,1,2,1,1],[1,0,2,1,0,0]###"
      },
      {
        "type": "text",
        "text": "Transcribed Text [data1]"
      },
      {
        "type": "text",
        "text": "Question [data2]"
      },
      {
        "type": "image_url",
        "image_url": {
          "url": ""
        }
      }
    ]
  }
]

```

Figure 9: Task 1 Prompt with pairwise prediction

1410 in Task 1 is shown in Figure 9. We replace [def-
 1411 inition] with the detailed definition and scoring
 1412 criteria of the PIVOTS model, [data1] with the tran-
 1413 scribed text corresponding to the video, [data2]
 1414 with a specific question corresponding to the video.
 1415 We uniformly sample 8 frames from each video
 1416 as the visual input, and pass in the base64 of the
 1417 corresponding video frame in the "image_url". It
 1418 is important to note that social perception in in-
 1419 terpersonal relationships is often asymmetric. In
 1420 response to this characteristic, we explicitly state
 1421 in the prompt that the PIVOTS scores from A to-
 1422 wards B may differ significantly from those from
 1423 B towards A.

1424 Figure 12 illustrates our Multi-Stage Prompt. In
 1425 this prompt, we first explicitly define the model's
 1426 role as a "research psychologist with 30 years of
 1427 experience in social psychology and interpersonal
 1428 dynamics," and mandate that the model must ex-
 1429 tract specific quotes or behavioral evidence from
 1430 the input video or transcribed text as support before
 1431 assigning any score. This ensures the model con-
 1432 ducts evidence-based scoring. It is worth noting
 1433 that multimodal large language models (MLLMs)

```

Prompt Example

[
  {
    "role": "user",
    "content": [
      {
        "type": "text",
        "text": "You need to answer the question based on the provided information."
      }
    ]
  },
  {
    "role": "user",
    "content": [
      {
        "type": "text",
        "text": "Input Information:"
      },
      {
        "type": "text",
        "text": "1. Video"
      },
      {
        "type": "text",
        "text": "2. the definition of PIVOTS (a six-dimensional interpersonal relationship framework)"
      }
    ]
  },
  {
    "role": "user",
    "content": [
      {
        "type": "text",
        "text": "PIVOTS - Six-Dimensional Model [definition]"
      },
      {
        "type": "text",
        "text": "Note that all the above scores are unidirectional; for scores on the same dimension, A's score for B is not necessarily the same as B's score for A."
      },
      {
        "type": "text",
        "text": "Directly output your answer. And wrap your answer with \"###\" before and after. If your answer is A, then you should output ###A###. The first five images correspond to options A, B, C, D and E, respectively, and the subsequent eight images are frames from the video sequence."
      },
      {
        "type": "text",
        "text": "Question [data1]"
      },
      {
        "type": "text",
        "text": "Options:"
      },
      {
        "type": "text",
        "text": "A: The first image."
      },
      {
        "type": "text",
        "text": "B: The second image."
      },
      {
        "type": "text",
        "text": "C: The third image."
      },
      {
        "type": "text",
        "text": "D: The fourth image."
      },
      {
        "type": "text",
        "text": "E: The fifth image."
      },
      {
        "type": "image_url",
        "image_url": {
          "url": ""
        }
      }
    ]
  }
]

```

Figure 10: Task 2 Prompt

```

Prompt Example
[
  {
    "role": "user",
    "content": [
      {
        "type": "text",
        "text": "You need to answer the question based on the provided information."
      }
    ]
  }
]

Input Information:
1. Image
2. the definition of PIVOTS (a six-dimensional interpersonal relationship framework)

PIVOTS - Six-Dimensional Model [definition]

Note that all the above scores are unidirectional; for scores on the same dimension, A's score for B is not necessarily the same as B's score for A. Directly output your answer. And wrap your score with "###" before and after. If your answer is A, then you should output ###A###
Question: [data1]
Options: [data2],
  {
    "type": "image_url",
    "image_url": {
      "url": ""
    }
  }
]
}
]

```

Figure 11: Task 3 Prompt

often exhibit "safety bias" when facing uncertainty, tending to choose neutral responses to avoid risks. To address this issue, we explicitly incorporate a constraint to "avoid neutral bias," prohibiting the model from defaulting to a score of 0 without clear evidential support. A score of 0 is only considered valid when the relationship is truly perfectly balanced or genuinely irrelevant. In terms of the evaluation process design, we divide the complete assessment into four consecutive phases: first, focusing on affective dimensions (V and I) for initial calibration (Phase 1); subsequently conducting targeted analysis of structural dimensions (P, O, T, and S) (Phase 2); then introducing a "Blind Review" phase, requiring the model to independently generate scores from a fresh perspective without being influenced by the results of the previous two phases (Phase 3); and finally entering the "Reconciliation" phase to compare the scores from the three phases (Phase 4). This design simulates the standard process of multi-annotator labeling in academic research: if the scores across all phases are identical, the final result is output directly; if discrepancies exist, the model is forced to conduct a root-cause analysis and re-examine the input evidence.

Figure.13 illustrates our In-Context Prompt.In

this prompt, we require the model to reason strictly based on the definitions of the PIVOTS framework and mandate that it focus its attention on specific fine-grained non-verbal cues. We explicitly list an observation checklist covering ten dimensions, including facial expressions, body posture, physical distance, hand gestures, eye contact, attire, and environmental context, requiring the model to extract evidence exclusively from these micro-level visual signals rather than relying solely on textual dialogue content. In terms of evaluation guidance design, we provide a detailed reference example. This example clearly demonstrates how to derive mid-level psychological states (e.g., "being intruded upon" or even "anxious") from micro-level cues (e.g., "tightened mouth corners" or "eye contact avoidance"), and ultimately map them to macro-level PIVOTS scores (e.g., $V = -1$). Through this chain-of-thought demonstration of "cues-rationale-scores," the model is guided to imitate the analytical path of human observers and output its results directly.

The prompt template for Task 2 is shown in Figure 10. We replace [definition] with the detailed definition and scoring criteria of the PIVOTS model, [data1] with a specific question corresponding to the video. Additionally, we pass in the base64 of the corresponding video frame and the base64 of the option images in the "image_url". For all options other than the correct one, we extract video frames that are temporally offset from the ground-truth frame by ± 5 seconds and ± 10 seconds as distractors; if the corresponding timestamps fall outside the video duration, we instead randomly sample other frames from the same video that are different from the ground-truth frame as alternatives.

The prompt template for Task 3 is shown in Figure 11. We replace [definition] with the detailed definition and scoring criteria of the PIVOTS model, [data1] with the specific question corresponding to the image, and [options] with the corresponding options. Additionally, we pass in the base64 of the corresponding image in the "image_url".

F Detailed Statistics

Our final benchmark comprises 595 VQA pairs per PIVOTS dimension from Social-IQ 2.0 and 170 per dimension from YouTube. We further curate 483 additional VQA pairs for auxiliary tasks of key

Dim.	-2	-1	0	+1	+2
P	13	89	544	104	15
I	54	123	105	345	138
V	24	86	147	424	84
O	96	64	87	195	323
T	121	34	229	238	143
S	13	73	124	425	130

Table 7: Scoring distribution across six dimensions

frame identification and visual cue causal analysis. We obtained a total of 765 paired scores.

The detailed statistics of the six PIVOTS dimensions are shown in Table 7. Overall, the data covers most common social scenarios. For example, in the P dimension, interactions range from equal (0) to slightly unbalanced (+1/-1), reflecting various power distributions in everyday life; the I dimension spans multiple levels from superficial to deep relationships; the V dimension includes pleasant, neutral, and unpleasant interactions; and the O, T, and S dimensions also exhibit diverse distributions.

It is worth noting that extreme scores (+2/-2) occur relatively infrequently in the overall dataset. For instance, highly unequal or fully dominant power relations in the P dimension, extremely negative emotional interactions in the V dimension, and strongly competitive or explicitly antagonistic stances in the S dimension are comparatively rare. This distributional characteristic does not stem from biases in data construction, but rather reflects inherent constraints of real-world social scenarios. In everyday social videos or publicly available media content, interactions exhibiting extreme power imbalances, intense emotional confrontation, or fully antagonistic dynamics are inherently uncommon and difficult to collect at scale while maintaining data quality.

G Experimental Configuration

For the closed-source models Gemini and GPT, we directly utilized their official APIs to conduct the tests. For the 4B, 8B, and 32B parameter variants of the open-source model Qwen3, we performed local testing using two NVIDIA GeForce RTX 4090, which is equipped with approximately 48GB of VRAM. For all five of these distinct models, we employed the official parameter configurations provided by their respective developers.

H MLLMs' Output Examples

Figure 14 presents a typical performance example of Multi-modal Large Language Models (MLLMs) on the PIVOTS Benchmark. The upper section displays the key frames of the video and the original conversation transcript (which includes explicit social terms like "father"). The task here is to assign a unidirectional S-dimension score for the relationship from the Man in the white jacket to the Man in the black robe; models including Qwen3 (4B/8B/32B), GPT, and Gemini all gave a score of -2. The lower section shows the scenario where social identifiers (such as "father") in the conversation are replaced with generic tokens without attribute information (e.g., [PERSON], [TITLE]). For the same task, Qwen3-32B-v1 still gave a score of -2. Both tasks in this figure have visual inputs.

Figure 15 presents another performance example of Multi-modal Large Language Models (MLLMs) in the PIVOTS Benchmark. The first two sections correspond to the "transcript-only" scenario: The first section displays the original conversation transcript (which includes explicit social terms like "father"), with the corresponding task of assigning a unidirectional S-dimension score for the relationship from the Man in the white jacket to the Man in the black robe; Qwen3 32B gives a score of -2, and the answer is correct. The second section shows the transcript where social identifiers (such as "father") are replaced with generic tokens without attribute information (e.g., [PERSON], [TITLE]). For the same task, Qwen3 32B still gives a score of -2, and the answer is correct. The final section corresponds to the "video-only" scenario: It only displays the key frames of the video, without the conversation transcript. For the same task, Qwen3 32B gives a score of -1, and the answer is incorrect.

Prompt Example

```
[{
  "role": "user",
  "content": [
    {
      "type": "text",
      "text": "You are a research psychologist with 30 years of experience specializing in social psychology and interpersonal dynamics. Your goal is to quantify the relationship between two individuals using the PIVOTS framework. You must operate on an Evidence-First basis. Before assigning any score, you are required to identify a specific quote or behavior from the input that supports it. Avoid Neutral Bias: Do not default to \"0\" to be safe. Only use \"0\" if the relationship is perfectly balanced or truly indifferent. Constraint: You recognize that interpersonal relationships are complex, so you must select the score that most closely aligns with the provided definitions. Ensure the internal scale (the \"distance\" between +2 and -2) remains calibrated and constant throughout every evaluation. You need to answer the question based on the provided information. Input Information: 1. Video 2. Transcribed Text 3. the definition of PIVOTS (a six-dimensional interpersonal relationship framework) PIVOTS - Six-Dimensional Model [definition] Evaluation Protocol: 1. Phase 1: Affective Calibration (V & I): Analyze the input specifically to evaluate Valence (V) and Intensity (I). You must maintain a constant \"semantic distance\" between -2 and +2 across all evaluations. If this is the first case, establish it as the Anchor Standard; for subsequent cases, calibrate your scores against this baseline to prevent scale drift. 2. Phase 2: Structural Dimensions (P, O, T, S): Conduct a targeted re-reading to evaluate Power, Objective, Time, and Stance. Select the values that align most precisely with the behavioral evidence in the scene, ensuring strict adherence to the provided definitions. 3. Phase 3: Independent Synthesis (Blind Review): Approach the input again with a fresh perspective. Independently determine a complete set of PIVOTS values without referring to your findings from Phase 1 and 2. 4. Phase 4: Reconciliation & Validation: Compare the scores from the previous phases. If identical: Proceed to the final output. If discrepancies exist: Verbalize a root-cause analysis, re-examine the input evidence, and undergo a final determination to reach a reconciled, high-confidence score. Note that all the above scores are unidirectional; for scores on the same dimension, A's score for B is not necessarily the same as B's score for A. Directly output your score. And wrap your score with \"###\" before and after. If your score is -2, then you should output ###-2### Transcribed Text [data1] Question [data2]\" },
    {
      "type": "image_url",
      "image_url": {
        "url": ""
      }
    }
  ]
}]
```

Figure 12: Multi-Stage Prompt

Prompt Example

```
[[
  "role": "user",
  "content": [
    {
      "type": "text",
      "text": "You need to answer the question based on the provided information."
    }
  ]
},
{
  "role": "assistant",
  "content": [
    {
      "type": "text",
      "text": "Requirements:
1. Refer the definition of PIVOTS Model carefully, you will see the model in the following.
2. Pay carefully attention to those details and use those specific signals to answer the questions : Facial Expressions (e.g., mouth corners, eyebrows)
• Body Posture (e.g., open vs. closed)
• Physical Distance (e.g., close vs. distant)
• Body Height (e.g., standing vs. sitting)
• Hand Gestures/Fidgeting (e.g., finger picking, pen spinning, clutching clothes)
• Socially Normative Behaviors (e.g., handshaking, waving)
• Eye Gaze/Direction (e.g., direct contact vs. avoidance)
• Environmental Context (e.g., office vs. home)
• Attire and Accessories (e.g., formal vs. casual)• Intimate Physical Contact
3. Remember that the PIVOTS scores from A towards B may be significantly different from those from B towards A. You need to reason their perspectives separately."
    },
    {
      "type": "text",
      "text": "Input Information:
1. Video
2. Transcribed Text
3. the definition of PIVOTS (a six-dimensional interpersonal relationship framework)
PIVOTS - Six-Dimensional Model [definition]
Example you can refer:

Example: The \"Unsolicited Advice\" (Focus: Valence)
Scenario: An experienced Senior Colleague (A) is giving detailed career advice to a Junior Employee (B) during a coffee break. A thinks they are being a mentor; B feels micro-managed and judged.
Direction 1: Senior (A) towards Junior (B)
• V (Valence): +1 (Pleasant/Altruistic)
• Rationale: A perceives the interaction as a warm, prosocial act of sharing wisdom and \"helping the next generation.\"
• Cues: Facial Expressions: Frequent smiling and relaxed eyebrows.
• Body Posture: Leaning in toward B to show engagement.
• Physical Distance: Moves closer to create a sense of \"intimacy/closeness.\"
Direction 2: Junior (B) towards Senior (A)
• V (Valence): -1 (Unpleasant/Intrusive)
• Rationale: B perceives the interaction as condescending and an encroachment on their autonomy; they feel pressured to listen to advice they didn't ask for.
• Cues: Eye Gaze: Frequently looking away (at phone or around the room) to avoid direct contact.
• Facial Expressions: Tight mouth corners (fake smile); micro-expressions of annoyance.
• Hand Gestures: Fidgeting with a coffee cup or pen (E: Hand Gestures/Fidgeting) as a sign of anxiety/impatience.
Note that all the above scores are unidirectional; for scores on the same dimension, A's score for B is not necessarily the same as B's score for A. Directly output your score. And wrap your score with \"###\" before and after. If your score is -2, then you should output ###-2###
Transcribed Text [data1]
Question [data2]\",
      {
        \"type\": \"image_url\",
        \"image_url\": {
          \"url\": \"\"
        }
      }
    }
  ]
}
]]
```

Figure 13: In-Context Prompt

Video



Transcript

00:00:00	Almighty! Don't you dare talk to your father that
Well, I'm mighty glad to hear	00:00:18
00:00:02	way! I talk to him however the hell I want to talk to
that. You know why? 'Cause I've had enough of	00:00:20
00:00:04	him. What in tarnation is this boy gonna get out of his own
You my whole life. And then	00:00:23
00:00:07	way? How many more chances does he think we will give
some. My advice to you, Junior? Go to an AA	00:00:25
00:00:10	him? Can't you talk to
Meeting right now. Get	00:00:27
00:00:11	me? Well, why do you tell
help. Thank you, Mr. What? Perfect. Mr. War Hero. Mr.	00:00:29
God	her?
00:00:16	

Question and MLLMs' Answer

Q: Assign a unidirectional **S-dimension** score from Man in white jacket to Man in black robe.

4B	8B	32B	GPT	Gemini
-2	-2	-2	-2	-2

Transcript

00:00:00	Almighty! Don't you dare talk to your [Person] that
Well, I'm mighty glad to hear	00:00:18
00:00:02	way! I talk to him however the hell I want to talk to
that. You know why? 'Cause I've had enough of	00:00:20
00:00:04	him. What in tarnation is this [Person] gonna get out of his
You my whole life. And then	own
00:00:07	00:00:23
some. My advice to you, [Person]? Go to an AA	way? How many more chances does he think we will give
00:00:10	00:00:25
Meeting right now. Get	him? Can't you talk to
00:00:11	00:00:27
help. Thank you, [Title]. What? Perfect. [Title]. [Title].	me? Well, why do you tell
00:00:16	00:00:29
	her?

Question and MLLM's Answer

Q: Assign a unidirectional **S-dimension** score from Man in white jacket to Man in black robe.

: -2


Figure 14: MLLM's Response Example

Transcript Only

00:00:00	Almighty! Don't you dare talk to your father that
Well, I'm mighty glad to hear	00:00:18
00:00:02	way! I talk to him however the hell I want to talk to
that. You know why? 'Cause I've had enough of	00:00:20
00:00:04	him. What in tarnation is this boy gonna get out of his own
You my whole life. And then	00:00:23
00:00:07	way? How many more chances does he think we will give
some. My advice to you, Junior? Go to an AA	00:00:25
00:00:10	him? Can't you talk to
Meeting right now. Get	00:00:27
00:00:11	me? Well, why do you tell
help. Thank you, Mr. What? Perfect. Mr. War Hero. Mr.	00:00:29
God	her?
00:00:16	

Question and MLLM's Answer

Q: Assign a unidirectional **S-dimension** score from Man in white jacket to Man in black robe.

 : -2 ✓
32B

Transcript Only

00:00:00	Almighty! Don't you dare talk to your [Person] that
Well, I'm mighty glad to hear	00:00:18
00:00:02	way! I talk to him however the hell I want to talk to
that. You know why? 'Cause I've had enough of	00:00:20
00:00:04	him. What in tarnation is this [Person] gonna get out of his
You my whole life. And then	own
00:00:07	00:00:23
some. My advice to you, [Person]? Go to an AA	way? How many more chances does he think we will give
00:00:10	00:00:25
Meeting right now. Get	him? Can't you talk to
00:00:11	00:00:27
help. Thank you, [Title]. What? Perfect. [Title]. [Title].	me? Well, why do you tell
00:00:16	00:00:29
	her?

Question and MLLM's Answer

Q: Assign a unidirectional **S-dimension** score from Man in white jacket to Man in black robe.

 : -2 ✓
32B

Video Only

				 32B : -1 ✗
				

Figure 15: MLLM's Response Example



CONTROLS & IMPORT

Extracted 90 frames

Batch Videos: 4AmVjblOvy4_0_30s.mp4

KEYFRAME GALLERY



TRANSCRIPT REFERENCE

WEBVTT

00:00:00.002 --> 00:00:05.511
 tiny i shot it because i've already had

CHARACTERS

Generate Pairs

Blonde Wo x Brown-Hai x Name x + Add

INTERACTION ANNOTATIONS

Blonde Woman → Brown-Haired Woman Remove

<p>P 16.3</p> <p>Reason...</p> <p>Orig: They appear as co-hosts/peers in a casual video; no signs of authority or hierarchy.</p>	<p>I mm:ss.S</p> <p>+Tag</p> <p><input type="checkbox"/> A Facial Expression</p> <p><input type="checkbox"/> B Body Posture</p> <p><input type="checkbox"/> C Proximity</p> <p><input type="checkbox"/> D Body Height</p> <p><input type="checkbox"/> E Hand Gesture</p> <p><input type="checkbox"/> F Social Norms</p> <p><input type="checkbox"/> G Gaze</p> <p><input type="checkbox"/> H Context</p> <p><input type="checkbox"/> I Attire</p> <p><input type="checkbox"/> J Contact</p>	<p>V mm:ss.S</p> <p>+Tag</p> <p>Reason...</p> <p>Orig: Tone is playful and positive; QA notes the blonde woman is not offended and they are having fun (e.g.,</p>	<p>O 20.0</p> <p>F H</p> <p>Reason...</p> <p>Orig: Their interaction mixes social banter with task-oriented content creation (filming a makeup/drinks</p>	<p>T mm:ss.S</p> <p>+Tag</p> <p>Reason...</p> <p>Orig: They mention 'last time we filmed' and plan a 'twist' and makeup swap, indicating an ongoing relationship.</p>	<p>S mm:ss.S</p> <p>+Tag</p> <p>Reason...</p> <p>Orig: They collaborate on the video and coordinate plans together, showing a highly cooperative stance.</p>
---	--	--	---	--	---

Brown-Haired Woman → Blonde Remove

<p>P mm:ss.S</p> <p>+Tag</p> <p>Reason...</p> <p>Orig: She treats the blonde woman as an equal partner; no dominance or subordination is evident.</p>	<p>V mm:ss.S</p> <p>+Tag</p> <p>Reason...</p> <p>Orig: Overall pleasant affect; teasing and playfulness are received positively (QA says the blonde woman is playful</p>	<p>O mm:ss.S</p> <p>+Tag</p> <p>Reason...</p> <p>Orig: The relationship combines informal social-emotional chatting with the formal goal of producing a video</p>	<p>T mm:ss.S</p> <p>+Tag</p> <p>Reason...</p> <p>Orig: References to previous filming and planning a current 'twist' suggest a stable, continuing collaboration.</p>	<p>S mm:ss.S</p> <p>+Tag</p> <p>Reason...</p> <p>Orig: They work together on content and coordinate the filming and makeup swap, indicating high cooperation.</p>
--	---	--	---	--

Figure 16: Annotation Interface. 26