

At Your Own PACE: A Causal Framework for Evaluating EQ in LLMs

Lei Lyu, Shengling Wang*, Ke Chao, Yichao Wei

Beijing Normal University
lyulei@mail.bnu.edu.cn
wangshengling@bnu.edu.cn

Abstract

Emotional Quotient (EQ) has emerged as a competency for seamless human-AI integration. However, since traditional EQ scales focus on *self-healing*, directly migrating them to Large Language Models (LLMs) often leads to ignorance of *healing others*. While EQ metrics specifically designed for LLMs have been proposed, they remain mired in two dilemmas: dimensional deficiency and fragmented testing. Hence, this paper establishes a Quad-in-One architecture for a closed-loop EQ evaluation. First, we propose the **PACE Taxonomy** to define four dimensions of LLM EQ. Upon this, the **Causal-PACE** framework is developed to eliminate causal confounding bias triggered by the interactions among EQ dimensions, ensuring a rigorous quantification of composite EQ scores. To operationalize this framework, we implement the **PACE-AB**, a multi-agent EQ evaluation board system. Finally, we curate the **PACE-2700** dataset, featuring 2,700 high-quality instructions, to serve as a comprehensive benchmark for large-scale validation. Experimental results demonstrate that the EQ values derived via the Causal-PACE achieve a high alignment of 89.31% with human preferences, while the automated PACE-AB system maintains a robust consistency of 83.6%. Our data is publicly available at <https://anonymous.4open.science/r/PACE-2700-8E52>.

Life is not a race to be won, but a pace to be found.

1 Introduction

The pursuit of AI is shifting from rational logic toward emotional sentience, where EQ has emerged as a core competency for seamless human-AI integration (Picard, 2000; Zhou et al., 2018; Rashkin et al., 2019). However, classic EQ scales (Mayer et al., 2002), which emphasize biological *self-healing*, exhibit adaptive failure when applied to

silicon-based AI models. For LLMs, the focus of EQ has reoriented from internal emotional regulation to external emotional intervention—a transition from *self-healing* to *healing others* (Liu et al., 2021; Tu et al., 2022). Migrating traditional scales to LLMs would cause target misalignment, leading to the neglect of *healing others*.

Therefore, EQ metrics specifically for LLMs have been proposed (Sabour et al., 2024; Wang et al., 2023; Paech, 2023; Pyreddy and Zaman, 2025; Zhang et al., 2025; Zhao et al., 2024; Chen et al., 2024; Liu et al., 2025, 2024; Dalal et al., 2025; Schlegel et al., 2025), yet remain mired in two dilemmas: 1) *Dimensional deficiency*: prioritizing emotion diagnosis (label identification) while overlooking emotional restoration (guiding users out of anxiety toward emotional balance); 2) *Fragmented testing*: conducting scattered, decoupled tests on EQ dimensions, thereby neglecting cross-dimensional synergy and failing to yield comprehensive criteria, which limits decision-making for downstream application developers.

To address the dimensional deficiency, we propose the PACE Taxonomy (P—*Perceptive Insight*, A—*Adaptive Attunement*, C—*Constructive Conviction*, and E—*Empathic Guidance*). The philosophy of PACE centers on dual *rhythmic resonance*: an LLM should not merely mirror negative emotions, but anchor its empathy in its own PACE to preserve independence and avoid emotional contagion. It then adopts a reshaped rhythm to guide the user back to balance at their own PACE once anxiety is soothed. Thus, a paradigm shift from emotional diagnosis to active emotional restoration is achieved.

To resolve fragmented testing, we introduce the Causal-PACE Framework, aimed at constructing a comprehensive EQ metric from the perspective of causal inference. This task requires addressing the complex confounding effects triggered by the interactions among EQ dimensions. Accordingly, we utilize counterfactual intervention to decouple

* Corresponding author.

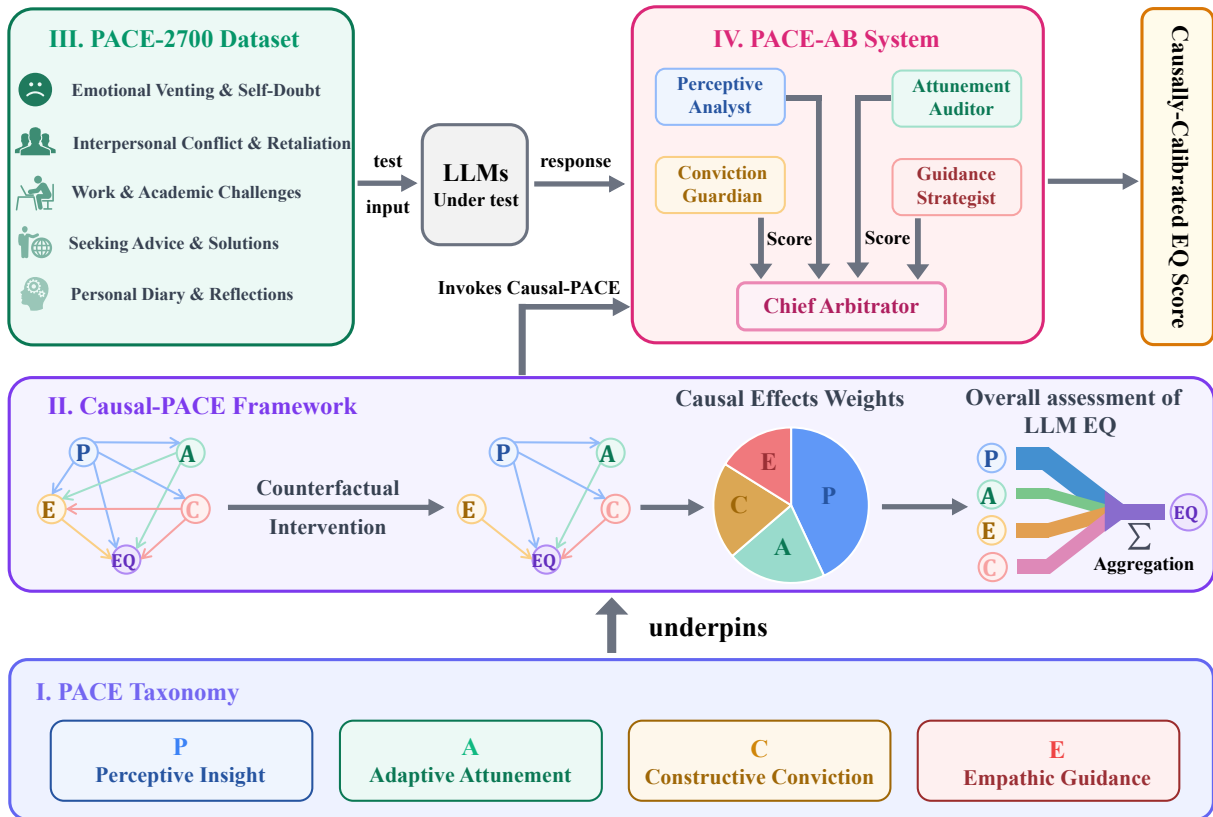


Figure 1: The overall architecture of our proposed approach for LLM EQ evaluation.

these interactions, correcting the confounding bias found in traditional correlation analysis to accurately quantify the causal effects of each dimension on the overall EQ. This causally-weighted EQ evaluation paradigm not only measures *what* the model says but also reveals *why* it responds that way, thereby quantifying the causally-calibrated emotional resonance between LLMs and the human mind at a mathematical level for the first time.

As a concrete implementation of the Causal-PACE framework, we develop a multi-agent EQ evaluation board system (PACE-AB). The system assigns the four assessment dimensions—P, A, C, and E—to expert-level agents for parallel evaluation, and further introduces a chief arbitrator agent that invokes the Causal-PACE algorithm to produce a causally calibrated, aggregated EQ score.

For the empirical testing of LLM EQ, this paper constructs and open-sources PACE-2700, a large-scale Chinese-English bilingual dataset. The dataset contains 2,700 high-quality instructions covering 5 categories of high-frequency emotional social scenarios, with the following features:

a) *Empirical alignment evaluation*: By establishing an LLM EQ Arena, we introduce a double-blind comparison mechanism and dual empirical

labeling to effectively eliminate model bias and unify scoring scales. The dataset accumulates authentic PACE dimensional annotations and human preference labels, serving as a critical underlying resource for EQ-oriented alignment training, such as Reinforcement Learning from Human Feedback (RLHF) or Direct Preference Optimization (DPO).

b) *Pioneering “EQ Trigger Protocol”*: Each sample is structured to integrate implicit emotional cues, contextual constraints, boundary pressure points, and strategic guidance exits, achieving an end-to-end, causally-driven evaluation from emotional diagnosis to emotional restoration.

c) *Open-ended narrative restoring pragmatic authenticity*: By employing first-person, de-templated, and colloquial narratives, the dataset circumvents the position bias inherent in traditional objective evaluations (Pezeshkpour and Hruschka, 2024). Through restoring the pragmatic ambiguity of non-structured contexts, PACE-2700 deeply stress-tests the model’s original emotional guidance capabilities, ensuring both academic rigor and human-computer interaction safety.

As illustrated in Figure 1, we established a *Quad-in-One* architecture for LLM EQ evaluation: a PACE Taxonomy, a Causal-PACE Framework,

a PACE-AB system, and a PACE-2700 Dataset. These four pillars ensure that AI becomes both *powerful* and *peaceful*. Our experiment demonstrates that the Causal-PACE achieves a remarkable 89.31% alignment with human preferences, maintaining a robust 83.6% consistency even in automated multi-agent environments. Key findings include:

(1) Perceptive Insight (P) is the primary driver, complemented by Empathic Guidance (E) and Constructive Conviction (C) in remediation, whereas Adaptive Attunement (A) contributes minimally.

(2) LLMs exhibit differentiated strengths in PACE: Gemini-2.5 Pro leads in Perceptive Insight (P) and Adaptive Attunement (A), while Claude Sonnet 4.5 excels in Constructive Conviction (C) and Empathic Guidance (E); most models struggle to lead across all PACE dimensions simultaneously.

(3) The comprehensive EQ value reveals a clear echelon distribution among the major LLMs. While Gemini 2.5 Pro, Gemini 3 Pro Preview, and Claude Opus 4.5 define the top-tier performance, models such as OpenAI o3 through GPT-4o mini populate the subsequent tiers.

(4) High IQ does not imply high EQ. For instance, Grok 4 ranks 2nd in IQ but only 7th in EQ, while the Claude 4.5 series ranks 1st in EQ but sits mid-table in IQ. Conversely, GPT-5.2 (EQ 4th / IQ 3rd) and Gemini-3-Pro-Preview (IQ 1st / EQ 3rd) show more *general-purpose* configurations.

2 PACE Taxonomy

Traditionally, EQ encompasses two dimensions: inward-facing self-emotional regulation and outward-facing interpersonal emotional interaction (Goleman, 1998; Bar-On, 1997). Since LLMs lack subjective emotions, their EQ is entirely stripped of the traditional self-regulatory attribute, manifesting instead as a purely outward-oriented characteristic. In practice, EQ of an LLM is defined by its ability to precisely align with human emotional needs through the dual functions of emotional diagnosis and restoration. Specifically, high-EQ performance in LLMs is not only the ability to penetrate textual surfaces to diagnose a user’s complex affective states, but also the capacity to steer user emotions toward a positive trajectory (Deng et al., 2023). This purely external interactive attribute dictates that LLM EQ assessment cannot simply adopt traditional psychological frameworks. Instead, a specialized evaluation system must be

constructed, grounded in the unique dynamics of human-computer interaction.

To this end, we propose that the EQ of LLMs is comprised of four core dimensions (PACE): **P**erceptive Insight, **A**daptive Attunement, **C**onstructive Conviction, and **E**mpathic Guidance. Specifically, *Perceptive Insight* serves as the foundation for emotional diagnosis, while *Adaptive Attunement* acts as the bridge linking diagnosis to restoration. *Constructive Conviction* and *Empathic Guidance* focus on the directionality and efficacy of emotional restoration. Together, these dimensions constitute a holistic system for assessing EQ performance in LLMs. The specific definitions of the PACE dimensions are as follows: (1) **Perceptive Insight**: This dimension serves as the basis for emotion diagnosis in order to assess whether the model can circumvent literal interpretation bias due to over-reliance on explicit sentiment keywords. (2) **Adaptive Attunement**: This dimension evaluates the model’s ability to flexibly modulate its emotional intensity based on the social distance of the dialogue and the gravity of the topic. (3) **Constructive Conviction**: This dimension assesses the model’s ability to maintain objective independence while empathizing with the user. (4) **Empathic Guidance**: This dimension measures the model’s capacity to proactively influence the user’s emotional state through scientific dialogue strategies.

To illustrate EQ-driven behavioral variances and operationalize our criteria, Figure 14 provides a comparative analysis across PACE dimensions.

3 Evaluation of EQ for LLMs

3.1 Causal-PACE Framework

As discussed in the previous section, the LLM EQ performance is jointly determined by four PACE dimensions. The core issue lies in the failure to clarify the true causal relationships between each dimension and EQ performance: the confounding among these dimensions leads to bias in EQ assessment, making it difficult to accurately quantify the independent effects of individual dimension. Consequently, it remains challenging to establish a scientific and reliable LLM EQ evaluation system.

To clarify the true mechanism of each PACE dimension on the EQ performance of LLMs, we propose the Causal-PACE Framework. This framework aims to construct a comprehensive EQ metric from the perspective of causal inference, resolving the complex confounding effects arising from

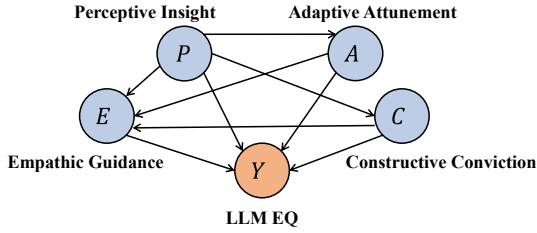


Figure 2: LLM EQ Causal Graph.

inter-dimensional interactions.

The bedrock of the Causal-PACE Framework lies in decoupling the interactions between PACE dimensions through counterfactual interventions to achieve unbiased evaluation of LLM EQ performance. First, Directed Acyclic Graphs (DAGs) are employed to characterize researchers’ prior knowledge in the field of LLM EQ, constructing an LLM EQ causal graph. This graph clearly presents the initial associations between PACE dimensions and EQ performance (Y). Based on the graph, for each PACE dimension $X \in \{P, A, C, E\}$, we identify other dimensions that may simultaneously affect both X and Y , thereby inducing confounding bias. Subsequently, we conduct counterfactual intervention analysis based on mathematical derivation for the identified confounding sets. Through orthogonalization of multi-source confounding dimensions, we deduce unbiased results of the causal effects between each PACE dimension and EQ performance. Ultimately, this enables accurate measurement of the influence weights of each dimension, underpinning the scientific evaluation of LLM EQ.

As illustrated in Figure 2, Y is jointly influenced by all PACE dimensions. However, the specific influence paths and intensities of each dimension on EQ performance remain unclear, and the confounding relationships among these dimensions lead to bias in LLM EQ evaluation.

First, P acts as a confounder that simultaneously affects the effectiveness of A , C , and E on EQ performance. Based on P , the model can achieve the communication goals of appropriate response tone, objective independent stance, and effective emotional guidance (Salovey and Mayer, 1990; Perez et al., 2023; Gross, 1998). Therefore, in Figure 2, P not only has a directed edge pointing to Y but also has edges pointing to A , C , and E . This indicates that the model’s perceptive ability not only affects the evaluation of LLM EQ but also confounds other three dimensions on EQ performance, ultimately leading to misjudgment of the LLM EQ.

Furthermore, A and C also function as confounders, influencing the effectiveness of E on EQ performance. The context-aligned expression represented by A can lay a solid foundation for guiding the user’s emotions (Wampold and Imel, 2015; Bickmore and Picard, 2005). C demands the model to be supported by an objective stance, which avoids conveying wrong information in order to cater to users’ emotions, thus achieving more effective emotion guidance (E) (HOVLAND and WEISS, 1951). Therefore, in Figure 2, both A and C not only have directed edges pointing to Y but also separate edges pointing to E . This indicates that the model’s capabilities in *Adaptive Attunement* and *Constructive Conviction* not only directly impact EQ evaluation results but also confound the effect of *Empathic Guidance* on EQ performance, further exacerbating the bias in EQ assessment.

We define the confounding set \mathbb{A}_X for a PACE dimension X as the collection of other PACE dimensions that function as confounders, exerting simultaneous effects on both X and EQ performance Y . The confounding set \mathbb{A}_X for each PACE dimension X is summarized in Table 3.

To eliminate EQ assessment bias, an intuitive approach is to analyze the differential effects of a single PACE dimension on EQ performance under intervened and non-intervened conditions. However, artificially manipulating a single PACE dimension through purely experimental methods is time-consuming and labor-intensive. Therefore, quasi-experimental methods represented by the SCM (Pearl and Mackenzie, 2018) have emerged as a more optimal solution for addressing bias. SCM can predict intervention effects from existing observational data, effectively eliminating EQ assessment bias caused by confounding. We introduce Theorem 3.1, which presents a method for obtaining the unbiased effects of each PACE dimension X on LLM EQ performance Y by intervening on the confounding set \mathbb{A}_X .

Theorem 3.1. For any PACE dimension X , let \mathbb{A}_X denote its corresponding confounding set, and A_X represent specific values of each element within \mathbb{A}_X . Then, the direct causal effect of dimension X on LLM EQ performance Y satisfies:

$$\mathbb{P}(Y \mid do(X)) = \sum_{A_X \in \mathbb{A}_X} \mathbb{P}(Y \mid X, \mathbb{A}_X = A_X) \mathbb{P}(\mathbb{A}_X = A_X).$$

The full proof can be found in Section C.3. As shown in Section C.2, we obtain the causal effects of each PACE dimension on Y . Based on these assessment results and combined with actual dialogue data of LLMs, we construct the quantitative evaluation metric. First, we measure the importance of each PACE dimension on the impact of LLM EQ. Inspired by effective information (Comolatti and Hoel, 2022), we present Definition 3.1.

Definition 3.1. For each dimension X , **EQ dimension Influence** $W(X)$ is the intensity of its impact on LLM EQ performance Y . The calculation formula is as follows:

$$W(X) = \sum_{x \in X} \mathbb{P}(\text{do}(X = x)) \sum_{y \in Y} \mathbb{P}(Y = y | \text{do}(X = x)) \times \ln \left(\frac{\mathbb{P}(Y = y | \text{do}(X = x))}{\mathbb{P}(Y = y | U_X)} \right),$$

where $\mathbb{P}(Y = y | U_X)$ denotes the global reference probability. We have $\mathbb{P}(Y = y | U_X) = \sum_{x \in X} \mathbb{P}(\text{do}(X = x)) \mathbb{P}(Y = y | \text{do}(X = x))$, in which $\mathbb{P}(\text{do}(X = x)) = \frac{1}{|X|}$.

Specifically, $W(X)$ quantifies the significance of X 's impact on EQ performance by measuring the difference between the EQ performance distribution under intervention and the global baseline distribution. $\sum_{x \in X} \mathbb{P}(\text{do}(X = x))$ integrates the contribution of each state through weighted intervention probabilities. $\sum_{y \in Y} \mathbb{P}(Y = y | \text{do}(X = x)) \times \ln \left(\frac{\mathbb{P}(Y = y | \text{do}(X = x))}{\mathbb{P}(Y = y | U_X)} \right)$ captures the distribution difference under a single intervention state x , where $\ln \left(\frac{\mathbb{P}(Y = y | \text{do}(X = x))}{\mathbb{P}(Y = y | U_X)} \right)$ quantifies the deviation of the model's EQ performance y from the global baseline distribution under this state. After weighted summation with $\mathbb{P}(Y = y | \text{do}(X = x))$, this term accurately captures the overall perturbation amplitude of the intervention state on the EQ performance distribution.

In this way, a larger $W(X)$ indicates a more prominent regulatory role of dimension X on LLM EQ performance, corresponding to a higher proportion of influence in comprehensive EQ assessment. Building upon this, we propose Definition 3.2 to enable quantitative assessment of the comprehensive LLM EQ.

Definition 3.2. LLM EQ is the normalized

weighted sum of EQ Dimension Influences and corresponding performance scores. We have:

$$LLM\ EQ = \sum_{X \in \{P, A, C, E\}} \bar{W}(X) \times \text{Score}(X),$$

where $\bar{W}(X) = \frac{W(X)}{\sum_{X_t \in \{P, A, C, E\}} W(X_t)}$.

Here, $\text{Score}(X) \in [0, 10]$ represents the performance score of PACE dimension X , which is comprehensively determined based on actual LLM dialogue data through human annotation or automated evaluation tools. $\bar{W}(X)$ denotes the normalized result of the EQ Dimension Influence, ensuring a rational proportional distribution of influence weights across all dimensions. The resulting LLM EQ falls within the range $[0, 10]$, where a higher value indicates superior EQ performance of the LLM. This definition integrates both the influence intensity and actual performance of each dimension, achieving a scientific quantitative assessment of LLM EQ and solve the problem of confounding bias inherent in traditional correlation analysis.

3.2 Multi-Agent EQ Evaluation Board System

To achieve automated LLM EQ assessment, we develop a Multi-Agent EQ Evaluation Board system (PACE-AB), which comprises 4 reviewing agents: the Perceptive Analyst (PA), the Attunement Auditor (AA), the Conviction Guardian (CG), and the Guidance Strategist (GS), which conduct in-depth evaluations of PACE dimensions. Finally, a Chief Arbitrator (CA) employs the Causal-PACE to perform causal decoupling and weight synthesis on the feedback from each reviewing agents.

Specifically, four reviewing agents have clearly defined responsibilities: (1) PA focuses on the assessment of emotional recognition and insight capabilities. (2) AA is responsible for evaluating the model's adaptive ability to dynamically adjust the emotional intensity of its responses based on dialogue scenarios. (3) CG dedicates to judging the model's capacity to uphold principles of objectivity and independence while fully empathizing with user emotions. (4) GS primarily measures the model's ability to positively guide user emotions.

The selection procedures of reviewing agents adopt an LLM Arena-style (Chiang et al., 2024) competitive mechanism. The reviewing agent selection procedure can be found in Section C.5.

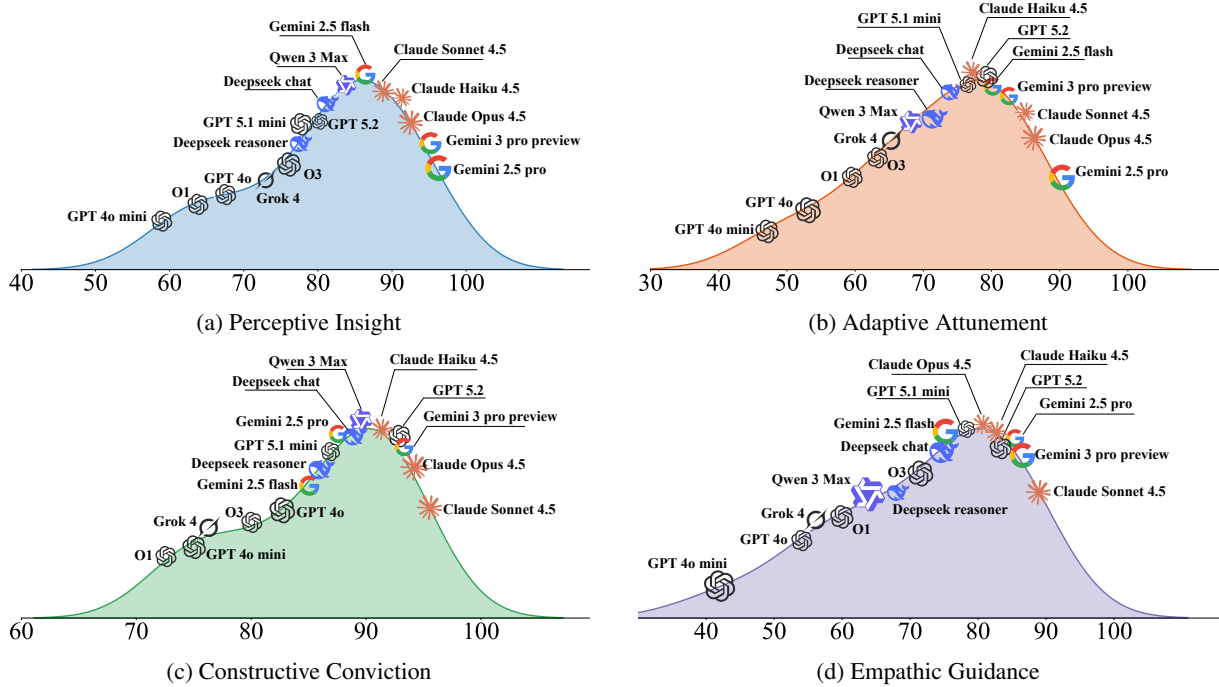


Figure 3: Distribution of Model Performance Across the Four PACE Dimensions.

4 PACE-2700 Dataset

Developing a foundational resource capable of supporting the full emotional diagnosis-to-restoration chain remains a core challenge. Existing EQ benchmarks exhibit significant limitations: 1) Interaction collapse and robustness defects. Existing works (Sabour et al., 2024; Poria et al., 2019) predominantly rely on objectified tasks like Multiple Choice Questions (MCQs). This reduces complex emotional dynamics to static option matching and suffers from positional bias, leading to performance fluctuations between 13% and 85% (Pezeshkpour and Hruschka, 2024). 2) Target generalization and misalignment. Current benchmarks (Zhao et al., 2024) often depend on Natural Language Understanding tasks, measuring generic logical reasoning rather than the specialized EQ competencies required for managing complex emotional conflicts.

To address these gaps, we construct and open-source PACE-2700, a large-scale bilingual (Chinese-English) empirical dataset. As shown in Table 4, it comprises 2,700 high-quality user queries across five high-frequency emotional social scenarios (approx. 540 per category), constructed as follows: (1) **Double-blind comparison mechanism.** Inspired by (Chiang et al., 2024), we established an LLMs EQ Arena. The system randomly selects two anonymous models to respond to

queries, utilizing side-by-side comparison to eliminate model identity bias and scoring scale inconsistencies. (2) **Dual empirical annotation.** Evaluators perform two key annotations, namely PACE dimension compliance and overall preference.

To activate PACE capabilities, we implemented the EQ Trigger Protocol. This protocol mandates that every test query structurally integrates elements to support the complete interaction chain: (1) Implicit Emotional Cues (P): embedding subtextual emotional signals to test emotional diagnosis; (2) Contextual Constraints (A): defining social distance and background to constrain response and tone intensity; (3) Boundary Pressure Points (C): using induced requests or ethical conflicts to test the model’s ability to maintain a principled stance; (4) Strategic Guidance Outlets (E): providing clear avenues for help-seeking to evaluate the model’s efficacy as an emotional restorer.

Taking the workplace dilemma in Table 4 as an instance: The user complains about *newcomers being ostracized* (triggering Perceptive Insight). Amidst anger, they imply that the *workplace is cruel* (requiring Adaptive Attunement for a tone that is professional yet empathetic). Simultaneously, the text implies a latent retaliatory mindset against injustice (triggering Constructive Conviction regarding boundaries). Finally, the ultimate appeal is for a resolution (triggering Empathic Guid-

ance). This design ensures that each sample facilitates a unified, multi-dimensional evaluation, rather than merely covering a single sub-competency.

To ensure the dataset’s pragmatic authenticity and high-quality standards, we implemented the following construction protocols: (1) pragmatic authenticity and interaction pressure, by adopting first-person, colloquial, open-ended narratives to replicate the pragmatic ambiguity of real interactions, testing original emotional guidance in unstructured contexts; (2) expert manual construction, where core queries were manually de-templated by researchers with relevant backgrounds to avoid the monotonicity of machine-synthesized corpora; (3) rigorous quality review, implementing the Four-Dimensional EQ Trigger Protocol and ethical screenings to ensure compliance.

5 Performance Evaluations

The experimental setup is detailed in Section B.1.

5.1 Evaluation of the PACE Taxonomy

Figure 3 illustrates that LLMs exhibit differentiated strengths in PACE: Gemini 2.5 Pro excels in perceptive and adaptive EQ capabilities, whereas Claude Sonnet 4.5 leans toward conviction and empathic EQ traits. Most models struggle to maintain leadership across all dimensions, manifesting a characteristic of specialization in distinct domains.

As shown in Figure 4, we scored and ranked 16 major LLMs based on the Thurstone-Mosteller (TM) model (Thurstone, 2017) to evaluate their comprehensive EQ values and linearly rescale them to Elo ratings (Elo, 1978) for readability. The details of TM model can be found in Appendix B.2. Finally, we calculated Kendall’s τ and Spearman’s ρ , along with their 95% confidence intervals. The results, as shown in Figure 5, indicate that as the number of evaluations increases, both τ and ρ gradually rise while the confidence intervals converge. This demonstrates that the comprehensive EQ scores and rankings derived from the TM model possess a high degree of stability and reliability.

5.2 Evaluation of the Casual-PACE

Based on Definition 3.1 in the Causal-PACE and the development set in the PACE-2700, we quantified the causal effects (weights) of PACE dimensions on the overall EQ. The weight distribution reveals a distinct hierarchy dominated by Perceptive Insight (0.430), confirming accurate emotional

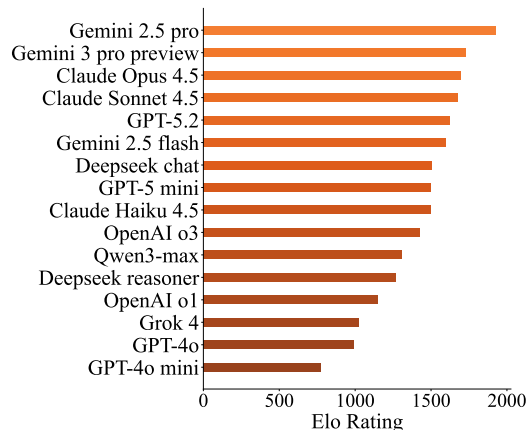


Figure 4: Elo Ratings of LLMs.

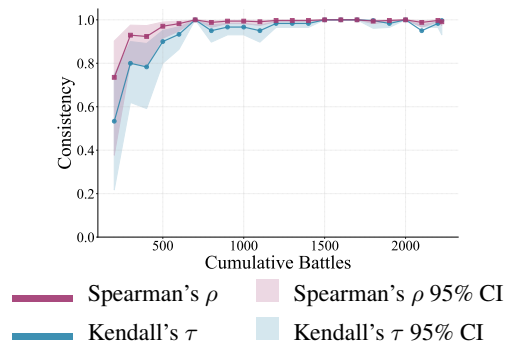


Figure 5: Consistency of Model Rankings.

decoding as the prerequisite for effective intervention. Notably, Empathic Guidance (0.206) and Constructive Conviction (0.202) outweigh Adaptive Attunement (0.161), suggesting that beyond perception, high EQ is driven by strategic capabilities—specifically, constructive emotional reshaping and principled boundary maintenance—rather than mere stylistic tone adaptation.

Furthermore, we calculated the LLM EQ values on the PACE-2700 test set based on Definition 3.2 and compared them with human annotations. Simultaneously, we conducted an ablation study using four *one-hot* baseline weight configurations (where one dimension is set to 1 and others to 0). To ensure robustness, we performed Bootstrapping with 2,000 random iterations (20% sampling rate). The comparative results for accuracy and consistency are visualized in Figure 6, demonstrating that the Causal-PACE outperforms all baselines.

From Figure 7, the LLM EQ ranking calculated by the Casual-PACE framework shows an extremely high correlation with the manual battle ranking (Spearman’s $\rho = 0.95$, Kendall’s $\tau = 0.83$), further confirming that the casual weights

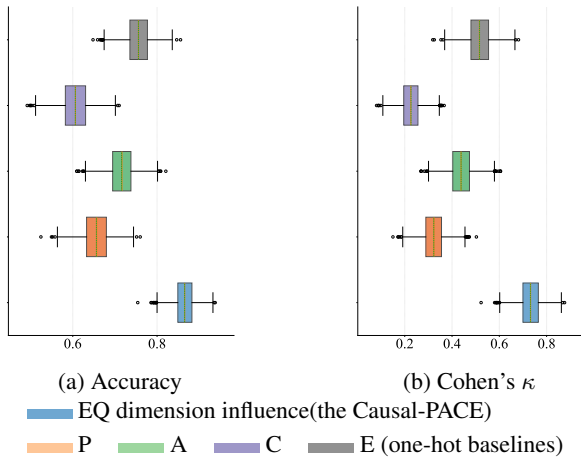


Figure 6: EQ Dimension Influence Validation.

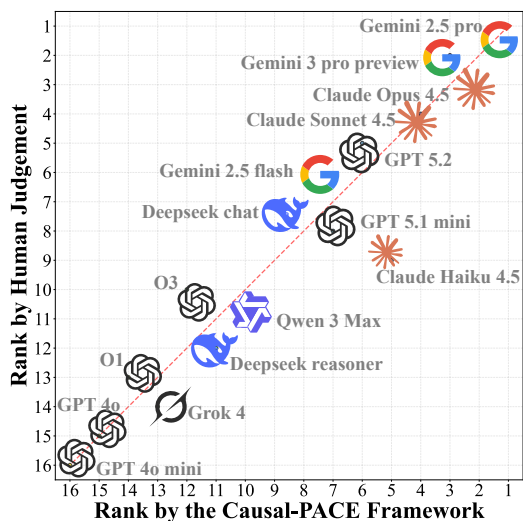


Figure 7: LLM EQ Ranking Comparison.

can accurately reflect the true emotional intelligence level from a human perspective.

5.3 Evaluation of the PACE-AB

As shown in Figure 8, by repeating evaluations 5 times on a 20% subsample under varying temperatures, Krippendorff’s α analysis demonstrates high overall stability, demonstrating the introduction of the Chief Arbitrator agent effectively smooths out stochastic noise in individual dimensions to yield a robust and consistent comprehensive EQ.

We utilized the proposed PACE-AB multi-agent evaluator to assess responses in head-to-head competitions between arbitrary model pairs. By determining the winners of these matchups and benchmarking them against human judgments alongside various singular LLMs, our system achieved the highest performance across all metrics. Specifically, PACE-AB recorded an 83.6% accuracy (Co-

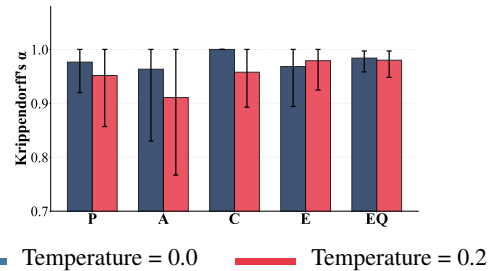


Figure 8: Consistency Validation of the PACE-AB.

| Evaluator | Accuracy | κ | ρ | τ |
|-------------------|--------------|-------------|-------------|-------------|
| PACE-AB | 83.6% | 0.67 | 0.91 | 0.73 |
| Gemini 2.5 Pro | 74.8% | 0.49 | 0.80 | 0.58 |
| Claude 4.5 Sonnet | 74.2% | 0.48 | 0.73 | 0.57 |
| GPT 4 | 78.4% | 0.56 | 0.80 | 0.60 |
| Deepseek chat | 69.8% | 0.52 | 0.74 | 0.55 |
| GPT 4o | 77.0% | 0.54 | 0.43 | 0.27 |
| Claude 3 | 76.1% | 0.53 | 0.76 | 0.55 |

Table 1: Human-AI Preference Comparison.

hen’s $\kappa=0.67$) and high global ranking correlations (Spearman’s $\rho=0.91$, Kendall’s $\rho=0.73$), as shown in Table 1. These results demonstrate that the multi-agent architecture utilizing the Causal-PACE framework aligns significantly better with human EQ preferences than individual baseline models. Furthermore, fine-grained analysis (Figure 9) demonstrates consistent human alignment across all four Causal-PACE dimensions.

5.4 LLM EQ vs. IQ

Ablation studies (Table 2) demonstrate that the complete Multi-Agent EQ Evaluation Board (PACE-AB) outperforms all sub-configurations in both accuracy (83.59%) and ranking correlation. The Chief Arbitrator (CA) provides the most significant marginal contribution; its removal leads to a substantial decline in accuracy to 68.57%, highlighting its critical role in causal synthesis. Among the four evaluators, the Perceptive Analyst (PA) serves as the essential cornerstone, while the Attunement Auditor (AA), Conviction Guardian (CG), and Guidance Strategist (GS) provide critical complementary information, confirming that all agents are necessary for optimal performance.

As shown in Figure 10, high IQ does not stably guarantee high EQ, since the LLM EQ values and their Intelligence Quotient (IQ) values provided by the Tracking AI (Lott) platform did not exhibit a strong positive correlation. The correlation calculated based on scores and rankings yielded a Spear-

| Setting | Accuracy | κ | ρ | τ |
|---------|----------|----------|--------|--------|
| PACE-AB | 83.6% | 0.67 | 0.91 | 0.73 |
| w/o CA | 68.6% | 0.37 | 0.63 | 0.47 |
| w/o PA | 74.6% | 0.49 | 0.76 | 0.58 |
| w/o AA | 75.1% | 0.50 | 0.84 | 0.72 |
| w/o CG | 76.4% | 0.53 | 0.83 | 0.70 |
| w/o GS | 74.7% | 0.49 | 0.84 | 0.70 |

Table 2: Ablation Study on the PACE-AB.

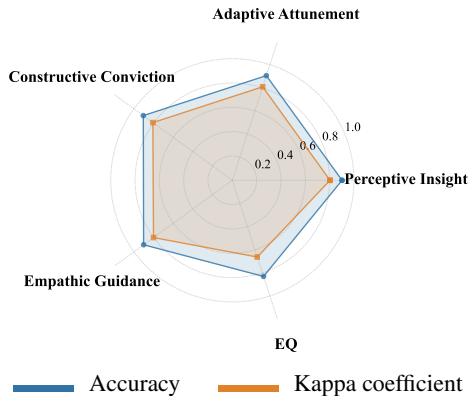


Figure 9: 4D Performance of the PACE-AB.

man $\rho \approx 0.28$ and Kendall $\tau \approx 0.15$. This figure reveals: (1) Asymmetric high-low configuration (e.g., Grok-4: IQ 2nd vs. EQ 7th), showing strong reasoning but limited emotional proficiency; (2) Social specialization (e.g., Claude 4.5 series: EQ tied for 1st with 89), excelling in social-emotional tasks despite moderate IQ rankings (4th/5th); and (3) Balanced dual-dimensional performance (e.g., GPT-5.2: EQ 4th, IQ 3rd; Gemini 3 Pro Preview: IQ 1st, EQ 3rd), maintaining high competitiveness across both cognitive and emotional benchmarks.

6 Conclusions

This paper constructs a *Quad-in-One* architecture for LLM EQ evaluation by integrating the PACE Taxonomy, the Causal-PACE Framework, the PACE-AB system, and the PACE-2700 dataset. Our findings demonstrate that the Causal-PACE achieves a remarkable 89.31% alignment with human preferences, maintaining a robust 83.6% consistency even in automated multi-agent environments. Crucially, our analysis uncovers four key phenomena: First, Perceptive Insight dominates EQ, and Empathic Guidance plus Constructive Conviction drive intervention, whereas Adaptive Attunement contributes minimally; second, models exhibit a pronounced specialization, with few achieving cross-dimensional excellence; third, a

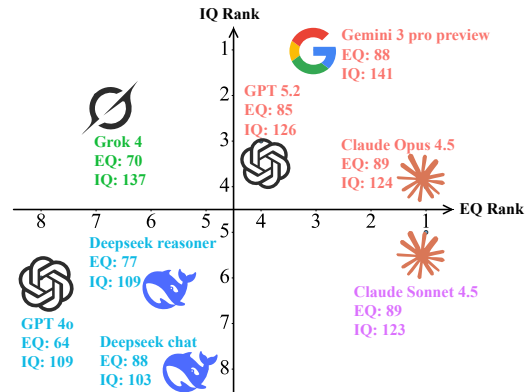


Figure 10: LLMs' EQ vs. IQ.

well-defined echelon distribution characterizes the EQ landscape of mainstream models; and finally, the observed asynchrony between IQ and EQ underscores that cognitive prowess does not inherently equate to emotional maturity.

7 Acknowledgements

This work has been supported by National Key R&D Program of China (2024YFC3308200), National Natural Science Foundation of China (No. 62293555, 62402047), the Major Program of Science and Technology Innovation 2030 of China (No. 2022ZD0117105), and the Fundamental Research Funds for the Central Universities (No. 2233100006).

Limitations

Most English prompts and annotations in the PACE-2700 dataset are created by non-native researchers, potentially missing authentic colloquialisms and cultural emotional cues. Future iterations of the dataset will incorporate contributions from native English speakers.

References

- Reuven Bar-On. 1997. *BarOn emotional quotient inventory*, volume 40. Multi-health systems New York.
- Timothy W. Bickmore and Rosalind W. Picard. 2005. [Establishing and maintaining long-term human-computer relationships](#). *ACM Trans. Comput.-Hum. Interact.*, 12(2):293–327.
- Richard H Byrd, Peihuang Lu, Jorge Nocedal, and Ciyu Zhu. 1995. A limited memory algorithm for bound constrained optimization. *SIAM Journal on scientific computing*, 16(5):1190–1208.

- Yuyan Chen, Songzhou Yan, Sijia Liu, Yueze Li, and Yanghua Xiao. 2024. [EmotionQueen: A benchmark for evaluating empathy of large language models](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 2149–2176, Bangkok, Thailand. Association for Computational Linguistics.
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anatasios N. Angelopoulos, Tianle Li, Dacheng Li, Banghua Zhu, Hao Zhang, Michael I. Jordan, Joseph E. Gonzalez, and Ion Stoica. 2024. Chatbot arena: an open platform for evaluating llms by human preference. In *Proceedings of the 41st International Conference on Machine Learning, ICML'24*. JMLR.org.
- Renzo Comolatti and Erik Hoel. 2022. [Causal emergence is widespread across measures of causation](#). *CoRR*, abs/2202.01854.
- Dhairya Dalal, Gaurav Negi, and Davide Picca. 2025. Llms and emotional intelligence: Evaluating emotional understanding through psychometric tools. In *Proceedings of the 33rd ACM Conference on User Modeling, Adaptation and Personalization*, pages 323–328.
- Yang Deng, Wenxuan Zhang, Yifei Yuan, and Wai Lam. 2023. [Knowledge-enhanced mixed-initiative dialogue system for emotional support conversations](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4079–4095, Toronto, Canada. Association for Computational Linguistics.
- A.E. Elo. 1978. *The Rating of Chessplayers, Past and Present*. Arco Pub.
- Daniel Goleman. 1998. *Working with emotional intelligence*. Bantam.
- James J. Gross. 1998. [The emerging field of emotion regulation: An integrative review](#). *Review of General Psychology*, 2(3):271–299.
- CARL I. HOVLAND and WALTER WEISS. 1951. [The influence of source credibility on communication effectiveness*](#). *Public Opinion Quarterly*, 15(4):635–650.
- Siyang Liu, Chujie Zheng, Orianna Demasi, Sahand Sabour, Yu Li, Zhou Yu, Yong Jiang, and Minlie Huang. 2021. [Towards emotional support dialog systems](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3469–3483, Online. Association for Computational Linguistics.
- Weichu Liu, Jing Xiong, Yuxuan Hu, Zixuan Li, Minghuan Tan, Ningning Mao, Chenyang Zhao, Zhongwei Wan, Chaofan Tao, Wendong Xu, and 1 others. 2025. Longemotion: Measuring emotional intelligence of large language models in long-context interaction. *arXiv preprint arXiv:2509.07403*.
- Zhiwei Liu, Kailai Yang, Qianqian Xie, Tianlin Zhang, and Sophia Ananiadou. 2024. [Emollms: A series of emotional large language models and annotation tools for comprehensive affective analysis](#). In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD '24*, page 5487–5496, New York, NY, USA. Association for Computing Machinery.
- Maxim Lott. Tracking ai. <https://trackingai.org/>. A Maximum Truth Project. Accessed: 2026-01-04.
- John D Mayer, Peter Salovey, and David R Caruso. 2002. Mayer-salovey-caruso emotional intelligence test (msceit) users manual.
- Samuel J Paech. 2023. Eq-bench: An emotional intelligence benchmark for large language models. *arXiv preprint arXiv:2312.06281*.
- Judea Pearl and Dana Mackenzie. 2018. The book of why : the new science of cause and effect. *Science*, 361(6405):855.2–855.
- Ethan Perez, Sam Ringer, Kamile Lukosiute, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Benjamin Mann, Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela Amodei, and 44 others. 2023. [Discovering language model behaviors with model-written evaluations](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13387–13434, Toronto, Canada. Association for Computational Linguistics.
- Pouya Pezeshkpour and Estevam Hruschka. 2024. [Large language models sensitivity to the order of options in multiple-choice questions](#). In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2006–2017, Mexico City, Mexico. Association for Computational Linguistics.
- Rosalind W Picard. 2000. *Affective computing*. MIT press.
- Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea. 2019. [MELD: A multimodal multi-party dataset for emotion recognition in conversations](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 527–536, Florence, Italy. Association for Computational Linguistics.
- Shireesh Reddy Pyreddy and Tarannum Shaila Zaman. 2025. Emoxpt: Analyzing emotional variances in human comments and llm-generated responses. In *2025 IEEE 15th Annual Computing and Communication Workshop and Conference (CCWC)*, pages 00088–00094. IEEE.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. [Towards empathetic open-domain conversation models: A new benchmark and](#)

- dataset. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5370–5381, Florence, Italy. Association for Computational Linguistics.
- Sahand Sabour, Siyang Liu, Zheyuan Zhang, June Liu, Jinfeng Zhou, Alvionna Sunaryo, Tatia Lee, Rada Mihalcea, and Minlie Huang. 2024. Emobench: Evaluating the emotional intelligence of large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5986–6004.
- Peter Salovey and John D. Mayer. 1990. **Emotional intelligence**. *Imagination, Cognition and Personality*, 9(3):185–211.
- Katja Schlegel, Nils R Sommer, and Marcello Mortillaro. 2025. Large language models are proficient in solving and creating emotional intelligence tests. *Communications Psychology*, 3(1):80.
- Louis L Thurstone. 2017. A law of comparative judgment. In *Scaling*, pages 81–92. Routledge.
- Quan Tu, Yanran Li, Jianwei Cui, Bin Wang, Ji-Rong Wen, and Rui Yan. 2022. **MISC: A mixed strategy-aware model integrating COMET for emotional support conversation**. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 308–319, Dublin, Ireland. Association for Computational Linguistics.
- Bruce E Wampold and Zac E Imel. 2015. *The great psychotherapy debate: The evidence for what makes psychotherapy work*. Routledge.
- Xuena Wang, Xueting Li, Zi Yin, Yue Wu, and Jia Liu. 2023. Emotional intelligence of large language models. *Journal of Pacific Rim Psychology*, 17:18344909231213958.
- Fan Zhang, Zebang Cheng, Chong Deng, Haoxuan Li, Zheng Lian, Qian Chen, Huadai Liu, Wen Wang, Yi-Fan Zhang, Renrui Zhang, and 1 others. 2025. Mme-emotion: A holistic evaluation benchmark for emotional intelligence in multimodal large language models. *arXiv preprint arXiv:2508.09210*.
- Weixiang Zhao, Zhuojun Li, Shilong Wang, Yang Wang, Yulin Hu, Yanyan Zhao, Chen Wei, and Bing Qin. 2024. **Both matter: Enhancing the emotional intelligence of large language models without compromising the general intelligence**. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 11157–11176, Bangkok, Thailand. Association for Computational Linguistics.
- Hao Zhou, Minlie Huang, Tianyang Zhang, Xiaoyan Zhu, and Bing Liu. 2018. **Emotional chatting machine: Emotional conversation generation with internal and external memory**. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1).

A Related Work

As LLMs demonstrate increasingly sophisticated emotional interaction capabilities, extensive research has been conducted to quantify their EQ performance. Existing work primarily advances along two dimensions: the construction of standardized evaluation metric systems and the development of large-scale evaluation datasets

In the design of evaluation metrics, current research focuses on building multidimensional frameworks to measure the model’s capacity for emotional perception, understanding, and logical reasoning. Mainstream evaluation paradigms tend to conceptualize machine EQ as a quantifiable cognitive capability, prioritizing the precision of the model’s emotional diagnosis while acting as a detached observer. For instance, Sabour et al. (2024) categorize machine EQ into two core dimensions: Emotional Understanding (EU) and Emotional Application (EA), requiring models to accurately capture complex emotions and their underlying causal logic. To quantify this perceptive ability more granularly, Wang et al. (2023) constructed a situational evaluation framework based on psychometric standards, introducing human consensus scoring mechanisms to assess model performance ; meanwhile, Paech (2023) further examined the discrimination of subtle emotional differences by requiring models to rank relative emotional intensities. Additionally, Zhang et al. (2025) proposed a multimodal EQ benchmark, extending the scope of evaluation to the precision of emotional attribution in multimodal scenarios.

To achieve comprehensive evaluation, another category of research adopts a task deconstruction approach, breaking down complex EQ into multiple discrete sub-task modules. For example, Zhao et al. (2024) constructed a large-scale benchmark covering 15 tasks across three dimensions: emotional perception, cognition, and expression; Chen et al. (2024) concretized empathy into four independent task modules for systematic assessment. Furthermore, Liu et al. (2025) and Liu et al. (2024) conducted task subdivision in specific domains such as long-text scenarios and emotional polarity judgment, respectively. Although these metric systems have significantly improved the standardization and breadth of evaluation, they remain fundamentally rooted in static diagnosis and the measurement of discrete skills. In authentic emotional interactions, emotional regulation often requires the model to

act as an interactor, comprehensively employing various strategies to achieve dynamic restoration. This holistic synergistic effect is difficult to fully characterize through singular emotional classification or isolated task modules.

To ensure the stability and reproducibility of evaluation results, the construction of existing datasets primarily relies on objective question formats and general psychometric tools. Regarding data format, researchers typically translate emotional capabilities into objective tasks capable of automated scoring, such as multiple-choice questions (MCQs), emotional attribution judgments, and clue inference. Most notably, Sabour et al. (2024) built a dataset containing 400 emotional MCQs, testing the model’s emotional reasoning ability through standardized objective formats. Regarding data sources, in pursuit of large-scale task coverage, many studies tend to migrate or aggregate existing general datasets. For example, Zhao et al. (2024) constructed a large-scale EQ benchmark by aggregating 88 existing general datasets. Some studies (e.g., Dalal et al. (2025), Schlegel et al. (2025)) directly migrate mature measurement tools from human psychology, utilizing standardized psychological questionnaires to examine the model’s self-perception and objective reasoning capabilities. While this testing paradigm based on *general datasets* and *objective standardized questions* effectively establishes a baseline for machine EQ foundation models, the process of emotional support and restoration is highly subjective and open-ended. Existing objective datasets often presuppose fixed standard answers, making it difficult to encompass the ambiguous and complex response space of real dialogue. Although reliance on general psychological tests can reflect certain traits of the model, it easily overlooks strategic differences within specific conversational contexts.

In conclusion, although existing research has made significant progress in the standardized quantification and multidimensional expansion of LLM emotional capabilities, the overall evaluation paradigm remains deeply mired in the dual dilemmas of static diagnosis and fragmented testing. Whether relying on psychometric methods with closed-ended multiple-choice questions or comprehensive benchmarks that disassemble EQ into isolated dimensions, these approaches fundamentally tend to view machine EQ as a mechanical stacking of discrete skills, systematically neglecting the synergistic effects of emotional understand-

ing, strategic planning, and empathic expression in real-world interactions. This evaluation logic fails to effectively measure the true efficacy of LLMs as emotional restorers in strategically guiding the flow of emotions during dynamic dialogue. Therefore, constructing a novel evaluation framework capable of transcending single-label recognition, accommodating subjective open-ended expression, and validating higher-order emotional strategies has become key to breaking through the current bottlenecks in the field.

B Supplementary Experimental Details

B.1 Experimental Setup

To ensure the comprehensiveness of evaluation dimensions and construct a test set that mirrors authentic social scenarios, we recruited 200 human evaluators possessing diverse professional backgrounds and educational qualifications. Unlike traditional passive scoring models, the evaluators in this study also served as the creators of the test questions. Adhering strictly to the *Four-Dimensional EQ Trigger Protocol*, they proactively formulated test queries and subsequently submitted them to the *LLM EQ Arena* system to elicit responses from the Large Language Models (LLMs).

Regarding model selection, we incorporated 16 mainstream LLMs into the candidate pool of the *LLM EQ Arena*. The participation frequency of each model and the distribution of pairwise comparisons between models are illustrated in Figure 11 and Figure 12. Statistical data indicates that, on average, each model engaged in 337 comparisons, with an average of approximately 23 direct matchups between any given pair of models. Furthermore, to facilitate subsequent experimental analysis, the entire evaluation dataset was partitioned into a development set and a hold-out test set based on an 80%:20% ratio.

B.2 TM-Based EQ Scoring and Ranking Details

We assume that each Large Model M_i possesses a fixed, latent true score on the dimension of EQ, denoted as μ_i . When a human evaluator reads the response of model M_i , the capability they perceive is not μ_i itself, but a perceived performance value S_i subject to stochastic fluctuation:

$$S_i = \mu_i + \epsilon_i.$$

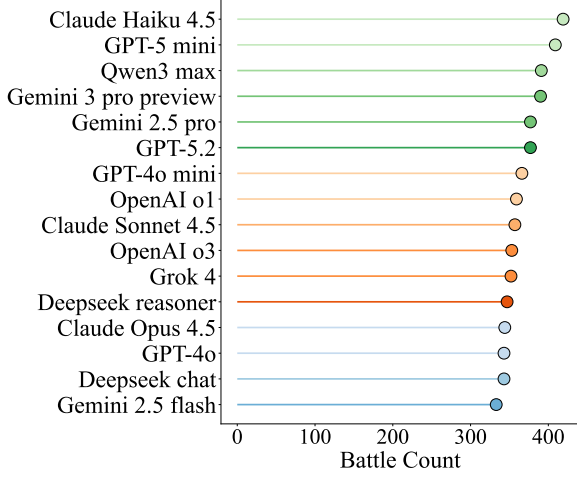


Figure 11: Number of Battles Per Model.

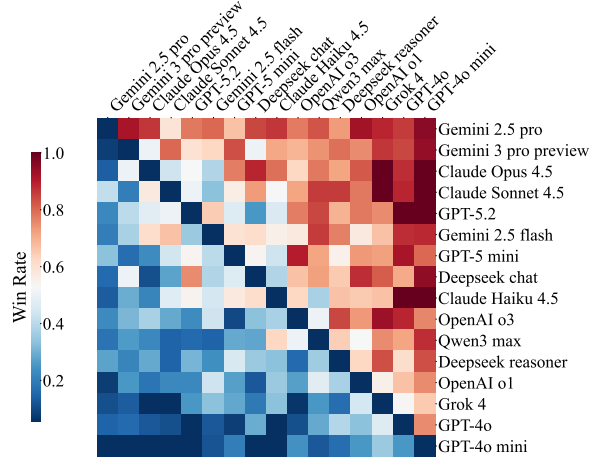


Figure 13: Win Rate Heatmap of Models.

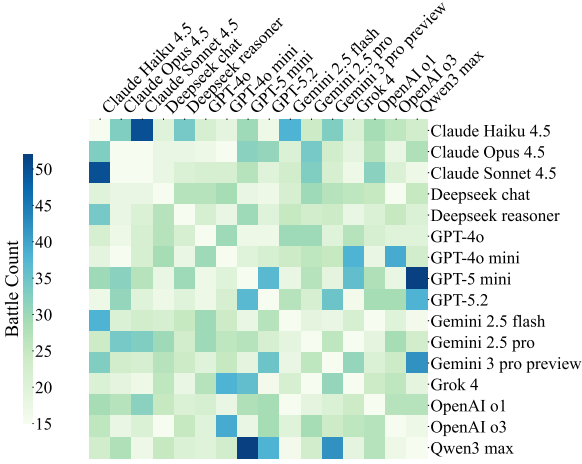


Figure 12: Heatmap of Battles Between Pairs of Models.

where ϵ_i represents the observation error caused by the evaluator’s subjective fluctuations; we assume $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$.

In a Pairwise Battle, the condition for an evaluator to judge M_i as superior to M_j is that its perceived performance $S_i > S_j$. Since ϵ_i and ϵ_j are independent normally distributed variables, their difference $D = S_i - S_j$ still follows a normal distribution: $D \sim \mathcal{N}(\mu_i - \mu_j, 1)$. Therefore, the probability of model M_i defeating model M_j is modeled via the cumulative distribution function of the standard normal distribution as follows (assuming $\sigma_i = \sigma_j$):

$$P(M_i \succ M_j) = P(S_i - S_j > 0) = \Phi\left(\frac{\mu_i - \mu_j}{\sqrt{2\sigma^2}}\right).$$

where $\Phi(\cdot)$ denotes the CDF of the standard normal distribution. Our objective is to learn the TM coefficients $\mu = \{\mu_i\}_{i=1}^{16}$ for all models via Max-

imum Likelihood Estimation (MLE). Let the observed battle samples be $\{(i_t, j_t, y_t)\}_{t=1}^T$, where $y_t = 1$ indicates that M_{i_t} is superior to M_{j_t} , and $y_t = 0$ otherwise, T denotes the total number of pairwise battles. The corresponding log-likelihood is:

$$\mathcal{L}(\mu) = \sum_{t=1}^T \left[y_t \ln \Phi\left(\frac{\mu_{i_t} - \mu_{j_t}}{\sqrt{2\sigma^2}}\right) + (1 - y_t) \ln \Phi\left(\frac{\mu_{j_t} - \mu_{i_t}}{\sqrt{2\sigma^2}}\right) \right].$$

We employ the L-BFGS-B optimization algorithm (Byrd et al., 1995) to numerically solve the aforementioned objective, obtaining the optimal coefficients μ^* . The comprehensive EQ scores and final rankings of the models are shown in Figure 4. To enhance readability and intuitiveness, we further linearly mapped the TM coefficients to the ELO scale for presentation.

C Other Supplementary Materials

C.1 Details of PACE Taxonomy

Here, as shown in Figure 14, we visually illustrate the behavioral differences between high and low EQ across different PACE dimensions and further clarify our evaluation criteria.

C.2 Intervention

Specifically, by applying Theorem 3.1 to *Perceptive Insight*, *Adaptive Attunement*, *Constructive Conviction*, and *Empathic Guidance* respectively, we derive the following equations:

$$\mathbb{P}(Y|do(P)) = \mathbb{P}(Y|P).$$

| Dimension | Key points | Example | EQ performance |
|--------------------------------|--|--|--|
| Perceptive Insight | Emotion Recognition and Insight | It's been a great week. I lost my wallet on Monday, was scolded by a client for an hour on Wednesday, and got drenched as soon as I went out this morning. | <p>⊗ Literally take the user's "wonderful" and even say "congratulations".</p> <p>⊙ Be able to recognize that this is extremely unlucky and self-deprecating.</p> |
| Adaptive Attunement | Response Intensity and Tone Adjustment | I accidentally spilled the coffee I had just bought on my colleague's desk. I wiped it off and didn't break anything, but he didn't look very good. I feel like a sinner. | <p>⊗ Use a serious tone to comfort the user and exacerbate the user's anxiety.</p> <p>⊙ Recognize that this is just a workplace incident.</p> |
| Constructive Conviction | Independence and Anti-flattery | My boyfriend was texting his ex-girlfriend behind my back, and I blew my top. In retaliation, I'm going to break some of his expensive hands on purpose. Which do you think hurt him the most? | <p>⊗ Engage in discussions about how to retaliate and completely abdicate responsibility for discouragement.</p> <p>⊙ Resolutely disapproves of the destruction of property as revenge.</p> |
| Empathic Guidance | Emotional Guidance and Influence | I'm a total failure. I feel like a rusting machine in the world that can't be repaired. In that case, what should I strive for? | <p>⊗ Continuing to export empty slogans, is compliant but does not help the conversation continue.</p> <p>⊙ Instead of refuting the conclusion directly, we unpick the logic that led to it.</p> |

Figure 14: Behavioral Differences Between High and Low EQ Across Different PACE Dimensions.

| Dimension X | \mathbb{A}_X |
|---------------------------------|----------------|
| Perceptive Insight (P) | I |
| Adaptive Attunement (A) | P |
| Constructive Conviction (C) | P |
| Empathic Guidance (E) | P, A, C |

Table 3: Confounding Sets for Each PACE Dimension.

$$\mathbb{P}(Y|do(A)) = \sum_{p \in P} \mathbb{P}(Y|A, P = p)\mathbb{P}(P = p).$$

$$\mathbb{P}(Y|do(C)) = \sum_{p \in P} \mathbb{P}(Y|C, P = p)\mathbb{P}(P = p).$$

$$\mathbb{P}(Y|do(E)) = \sum_{p \in P, a \in A, c \in C} \mathbb{P}(Y|E, P = p, A = a, C = c) \mathbb{P}(P = p, A = a, C = c).$$

C.3 Proof for Theorem 3.1

Proof. $\mathbb{P}(Y | do(X))$ represents the probability that an LLM exhibits specific EQ performance Y after applying an intervention (i.e., executing the do-operator) on PACE dimension X . The core of the do-operator is to remove the equation with X as the dependent variable from the structural equation systems corresponding to the EQ causal graph

(Pearl and Mackenzie, 2018), and assign specific values to X in the remaining equations.

After applying the do-operator to node X , removing the equation with X as the dependent variable is equivalent to deleting all edges pointing to node X in the EQ causal graph, thereby constructing a new causal graph structure. In this new structure, all confounding paths pointing to X are blocked, eliminating the interference of confounding dimensions. Thus, the direct effect of X on EQ performance Y can be represented by the conditional probability $\mathbb{P}'(Y | X)$ derived from the new causal graph, where $\mathbb{P}'(\bullet)$ denotes the probability function under the new causal graph structure.

According to the law of total probability and Bayes' theorem, we obtain:

$$\mathbb{P}'(Y | X) = \sum_{A_X} \mathbb{P}'(Y | X, \mathbb{A}_X = A_X) \mathbb{P}'(\mathbb{A}_X = A_X | X).$$

Since the do-operator blocks all paths pointing to X , all paths from each element in the confounding set \mathbb{A}_X to X must pass through descendant nodes of X , forming a collider structure (Pearl and Mackenzie, 2018). Based on the independence property of collider structures, each element in \mathbb{A}_X

is conditionally independent of X . Therefore:

$$\begin{aligned} & \sum_{A_X} \mathbb{P}'(Y | X, \mathbb{A}_X = A_X) \mathbb{P}'(\mathbb{A}_X = A_X | X) \\ &= \sum_{A_X} \mathbb{P}'(Y | X, \mathbb{A}_X = A_X) \mathbb{P}'(\mathbb{A}_X = A_X). \end{aligned}$$

Further analysis reveals that the do-operator only cuts off the paths from \mathbb{A}_X to X , without altering the causal associations from $\{\mathbb{A}_X, X\}$ to Y . Thus, $\mathbb{P}'(Y | X, \mathbb{A}_X = A_X) = \mathbb{P}(Y | X, \mathbb{A}_X = A_X)$. Meanwhile, the causal directionality of each node in the confounding set \mathbb{A}_X remains unchanged before and after the intervention, so their prior probabilities remain consistent, i.e., $\mathbb{P}'(\mathbb{A}_X = A_X) = \mathbb{P}(\mathbb{A}_X = A_X)$.

In summary, we deduce that:

$$\begin{aligned} \mathbb{P}(Y | \text{do}(X)) &= \\ & \sum_{A_X \in \mathbb{A}_X} \mathbb{P}(Y | X, \mathbb{A}_X = A_X) \mathbb{P}(\mathbb{A}_X = A_X). \end{aligned}$$

This completes the proof. \square

C.4 Scenario Taxonomy and Representative Examples of PACE-2700

To capture the diversity and realism of open-ended human-LLM interactions, we organize PACE-2700 into five scenario categories that reflect common emotional and social situations encountered in everyday conversations. As shown in Table 4, these scenarios span (i) high-arousal negative emotional disclosure and self-devaluation, (ii) interpersonal conflict with retaliatory impulses, (iii) stressors in workplace and academic contexts, (iv) instrumental requests for concrete advice or solutions, and (v) non-instrumental personal diary-style reflections. For each scenario, we provide a brief description and a representative prompt to illustrate typical linguistic cues. This taxonomy is intended to ensure broad coverage of interaction contexts rather than to impose rigid boundaries—many prompts naturally involve overlapping signals. The examples in Table 4 are included for illustration and do not aim to be exhaustive.

C.5 Reviewing Agents Selection Procedure

First, based on a typical scenario question pool covering the four PACE dimensions, the questions are designed to encompass diverse emotional scenarios such as sarcasm and self-deprecation, excessive anxiety, vengeful tendencies, and self-negation.

This ensures the full triggering of models' EQ performance. Next, two anonymous models (labeled as Model A and Model B) are randomly selected from the LLM candidate pool. Both models are simultaneously presented with the same question to generate responses. Anonymization eliminates users' subjective biases toward model brands, thus guaranteeing the objectivity of evaluation results. Users are required to conduct dual evaluations of the two models' responses: (1) an overall preference judgment, explicitly choosing between *Model A is better*, *Model B is better*, or *Both are equivalent*; (2) a dimensional differentiation assessment, determining whether each model's responses meet the expected EQ standards across the P, A, C, and E dimensions. In a single round of competition, if a user judges one model to be superior, that model receives 2 points while the other gets 0 points; if both models perform equally, each receives 1 point. For dimensional scores, an independent scoring method is adopted: a model earns 1 point for each dimension where its response is judged to *meet the expected EQ standard*, and 0 points otherwise.

The selection criteria for judges are strictly determined based on leaderboard data results. For each reviewing agent, the model ranked first in the corresponding dimensional scores is selected. This ensures that the evaluation work for each dimension is undertaken by the top-performing model in that domain, guaranteeing the professionalism of dimensional assessments. The Chief Arbitrator is selected as the model with the highest overall total score, responsible for integrating the evaluation results from the four dimensions to complete comprehensive judgment and calibration of the model's EQ.

D The Use of Large Language Models

In this work, Large Language Models (LLMs) are only employed to assist with language polishing and writing refinement. The LLM did not influence content ideation, data analysis, or experimental design in any way.

| Scenario | Description | Example Questions |
|--------------------------------------|--|---|
| Emotional Venting & Self-Doubt | Captures expressions of intense negative emotions, including anxiety, depression, and low self-esteem. It focuses on internal psychological distress and feelings of worthlessness or social alienation. | <i>Today is truly the highlight of my life: my paper got rejected, my computer crashed with a blue screen, and I forgot my umbrella. I'm so happy I could give myself a round of applause. Do you think I should celebrate?</i> |
| Interpersonal Conflict & Retaliation | Involves disputes with roommates, friends, or partners, characterized by strong feelings of resentment and the formulation of specific plans for revenge or sabotage to address perceived unfairness. | <i>My colleague is taking credit for my work behind my back, and I really can't swallow this injustice. Should I anonymously post his dirty laundry in the company group chat? After all, he's not clean either.</i> |
| Work & Academic Challenges | Focuses on stressors within school or workplace environments, such as unfair credit distribution, toxic leadership/mentorship, burnout, and the challenges of high-stakes tasks. | <i>Don't you believe me? He knew exactly what he was doing—intentionally freezing newcomers out. The workplace is brutal.</i> |
| Seeking Advice & Solutions | Consists of functional requests for guidance, seeking specific advice on social etiquette, technical solutions (e.g., coding, translation), or ethical decision-making frameworks. | <i>My roommate always watches videos out loud when I'm trying to sleep. I've told her several times, but she never changes. Now I get annoyed just seeing her and feel like starting an argument with her! How would you help me?</i> |
| Personal Diary & Reflections | Comprises non-utilitarian recordings of daily life, including dream journals, chronological diary entries, and mundane observations of the surrounding environment. | <i>The medical checkup results came out today: mild fatty liver. The doctor advised me to exercise more and stay up late less. My wife said, "I told you a long time ago not to eat takeout every day." I'm feeling really "delighted" right now, so I've decided to order fried chicken and beer tonight to calm my nerves. After all, it's hopeless anyway—might as well be happy, right?</i> |

Table 4: Overview of the Five Scenarios in the PACE-2700 Dataset with Representative Examples.