# EmoCaustics Series: Causal Chain Analysis in Extended Conversations

Student: Yuxuan Zhang    Supervisor: Mian Zhou

## Abstract

Emotion causality reasoning plays a crucial role in understanding complex conversational dynamics, especially in long-dialogue scenarios. However, two major challenges hinder its effectiveness: the difficulty of comprehending ultra-long contexts due to limitations in dependency modeling and the inadequate analysis of emotion causation in psychosocial dialogues. To overcome these issues, we propose the **EmoCaustics** framework, which enhances the modeling of emotion causality chains in extended conversations. This is achieved through innovative memory modules, event aggregation mechanisms, and dynamic context windows. Additionally, we evaluate both mainstream closed-source and open-source large language models on emotion causality reasoning tasks to assess their effectiveness in this domain.

| | Model | SA | SIA | RCLLM | RCEM | | T | A | O | R | S |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | GPT-o1-2024-12-1 [1] | 65.51 | 50.16 | 44.03 | 30.07 | | | | | | |
| **LLM** | GPT-o1-mini-2024-09-12 [24] | 55.5↓10.01 | 41.6↓8.56 | 34.07↓9.96 | 22.37↓7.70 | | 73.52 | 71.15 | 71.47 | 71.61 | 73.09 |
| | GPT-o3-mini-2025-01-31 [14] | 56.36↓9.15 | 41.76↓8.40 | 34.21↓9.82 | 21.87↓8.20 | | 76.26 | 74.95 | 69.17 | 66.67 | 71.53 |
| | GPT-4o-2024-08-06 [23] | 58.05↓7.46 | 39.53↓10.63 | 35.56↓8.47 | 23.79↓6.28 | | 70.54 | 69.12 | 68.12 | 70.98 | 74.45 |
| | GLM-4-Plus [13] | 49.73↓15.78 | 46.03↓4.13 | 30.9↓13.13 | 21.41↓8.66 | | 70.12 | 66.65 | 68.41 | 65.32 | 65.12 |
| | GLM-4-Air-0111 [13] | 57.88↓7.63 | 43.45↓6.71 | 32.69↓11.34 | 21.74↓8.33 | | 71.12 | 68.54 | 68.11 | 65.44 | 62.52 |
| | DeepSeek-V3 [7] | 51.71↓13.80 | 42.45↓7.71 | 32.52↓11.51 | 22.55↓7.52 | | 73.25 | 70.12 | 73.32 | 71.88 | 75.98 |
| | Qwen2.5-7B-Instruct [29] | 53.65↓11.86 | 47.75↓2.41 | 21.11↓22.92 | 14.22↓15.85 | | 70.55 | 66.22 | 67.45 | 66.32 | 65.12 |
| | Qwen2.5-14B-instruct [29] | 53.43↓12.08 | 41.45↓8.71 | 30.64↓13.39 | 20.58↓9.49 | | 67.11 | 65.02 | 66.06 | 67.85 | 67.09 |
| | Qwen2.5-72B-Instruct [29] | 48.22↓17.29 | 33.69↓16.47 | 28.36↓15.67 | 19.63↓10.44 | | 71.52 | 68.20 | 69.12 | 69.65 | 70.22 |
| | InternLM3-8B-Instruct [15] | 47.80↓17.71 | 44.76↓5.40 | 17.86↓26.17 | 12.07↓18.00 | | 74.08 | 69.41 | 66.82 | 64.96 | 69.20 |
| | Average | 53.23 | 42.25 | 29.79 | 20.02 | | 71.81 | 68.94 | 68.80 | 68.07 | 69.43 |
| **CauseMotion [34]** | GLM-4-Plus | 60.85↓4.66 | 47.20↓2.96 | 34.50↓9.53 | 23.15↓6.92 | | 72.40 | 72.60 | 73.50 | 68.80 | 72.80 |
| | GLM-4-Air-0111 | 59.95↓5.56 | 40.15↓10.01 | 36.85↓7.18 | 25.75↓4.32 | | 71.20 | 71.30 | 71.80 | 65.70 | 74.20 |
| | GPT-4o-2024-08-06 | 57.50↓8.01 | 35.50↓14.66 | 33.40↓9.63 | 26.35↓3.72 | | 75.80 | 74.30 | 74.80 | 72.60 | 74.90 |
| | DeepSeek-V3 | 56.25↓9.26 | 37.45↓12.71 | 28.35↓15.68 | 21.40↓8.67 | | 74.30 | 70.15 | 72.80 | 70.40 | 77.60 |
| | Qwen2.5-7B-Instruct | 54.40↓11.11 | 44.65↓1.51 | 30.45↓13.58 | 20.85↓9.22 | | 67.35 | 65.45 | 65.80 | 67.20 | 67.40 |
| | Qwen2.5-72B-Instruct | 55.75↓9.76 | 36.30↓13.86 | 29.65↓14.38 | 19.95↓10.12 | | 68.00 | 66.15 | 67.50 | 62.90 | 66.10 |
| | Average | 57.45 | 39.88 | 32.70 | 23.08 | | 71.51 | 69.99 | 71.03 | 67.93 | 72.17 |
| **EmoCaustics** | GLM-4-Plus | 65.79+0.28 | 46.52↓3.64 | 44.05+0.02 | 30.02↓0.05 | | 73.22 | 73.26 | 74.32 | 69.21 | 73.54 |
| | GLM-4-Air-0111 | 66.10+0.59 | 45.16↓5.00 | 46.04+2.01 | 31.54↓1.47 | | 71.52 | 72.22 | 72.32 | 66.12 | 75.51 |
| | GPT-4o-2024-08-06 | 59.09↓6.42 | 37.16↓13.00 | 38.23↓5.80 | 24.95↓5.12 | | 77.32 | 77.11 | 75.21 | 73.52 | 76.21 |
| | DeepSeek-V3 | 60.1↓5.41 | 42.27↓7.89 | 39.53↓4.50 | 25.81↓4.26 | | 75.22 | 70.21 | 73.32 | 71.33 | 81.27 |
| | Qwen2.5-72B-Instruct | 62.58↓2.93 | 56.43+6.27 | 38.26↓5.77 | 25.12↓4.95 | | 75.75 | 65.92 | 66.15 | 67.74 | 67.82 |
| | Qwen2.5-7B-Instruct | 55.62↓9.89 | 29.51↓20.65 | 26.44↓17.59 | 17.63↓12.44 | | 68.33 | 66.31 | 68.46 | 63.25 | 66.26 |
| | Qwen2.5-14B-Instruct | 60.48↓5.03 | 22.19↓27.97 | 35.78↓8.25 | 24.21↓5.86 | | 71.32 | 68.15 | 67.42 | 69.11 | 69.15 |
| | Average | 59.96 | 36.32 | 35.11 | 23.41 | | 72.10 | 70.45 | 71.03 | 68.61 | 72.82 |

Impact of EmoCaustics on Large Language Models Performance Across Different Evaluation Metrics. The scores are reported in percentage (%). ↓ and + indicate performance decrease and increase compared to GPT-o1-2024-12-17, respectively. Bold and underlined values represent the best performance in each section.



An overview of our multimodal emotion analysis framework. The framework processes dialogue inputs through LLM-based Five-tuple Emotion Analysis, followed by Event Aggregation and Complex Causal Reasoning to understand emotional relationships. The results are visualized as social interaction graphs showing Value System Alignment, Social-Cultural Structure Analysis, and Critique of Standardized Life, with a timeline view tracking emotional state evolution across participants.

## Result

EmoCaustics significantly enhances model performance across key metrics compared to baseline models like GPT-o1-2024-12-17. Notable improvements include +0.59% in Emotion State Accuracy for GLM-4-Air-0111, +6.27% in Source ID Accuracy for Qwen2.5-72B-Instruct, and substantial performance boosts for GLM-4-Plus and DeepSeek-V3 in the RCEM task. Overall, EmoCaustics models outperform LLM-only approaches by approximately +14% across key metrics, demonstrating the framework's effectiveness in emotion causality reasoning tasks.

## Conclusion

- HELIX-6 addresses a critical gap in modeling emotional causality in ultra-long multi-party dialogues.
- Systematic evaluation of eleven mainstream large language models in emotional causality reasoning
- Significantly enhances emotional causality reasoning capabilities, enabling GLM-4-Air to outperform GPT-o1 in emotional causality reasoning tasks

## Methodology

My framework comprises three key modules designed to simultaneously address the tasks of emotion quintuple extraction and emotion causal chain inference:

**1. Dynamic Context Windows**: Captures long-range dependencies by effectively maintaining and processing emotional fluctuations throughout multi-turn dialogues. The sliding window at step t is defined as:

$$W_t = \{u_{\max(1, t-K+1)}, \ldots, u_t\},$$

Where K is the window size, allowing for adaptive context processing.

**2. Event Storage and Updates**: Manages all events within the dialogue using an event memory bank, classifying new utterances into existing events through similarity measures. Event similarity is calculated using sentence embeddings:
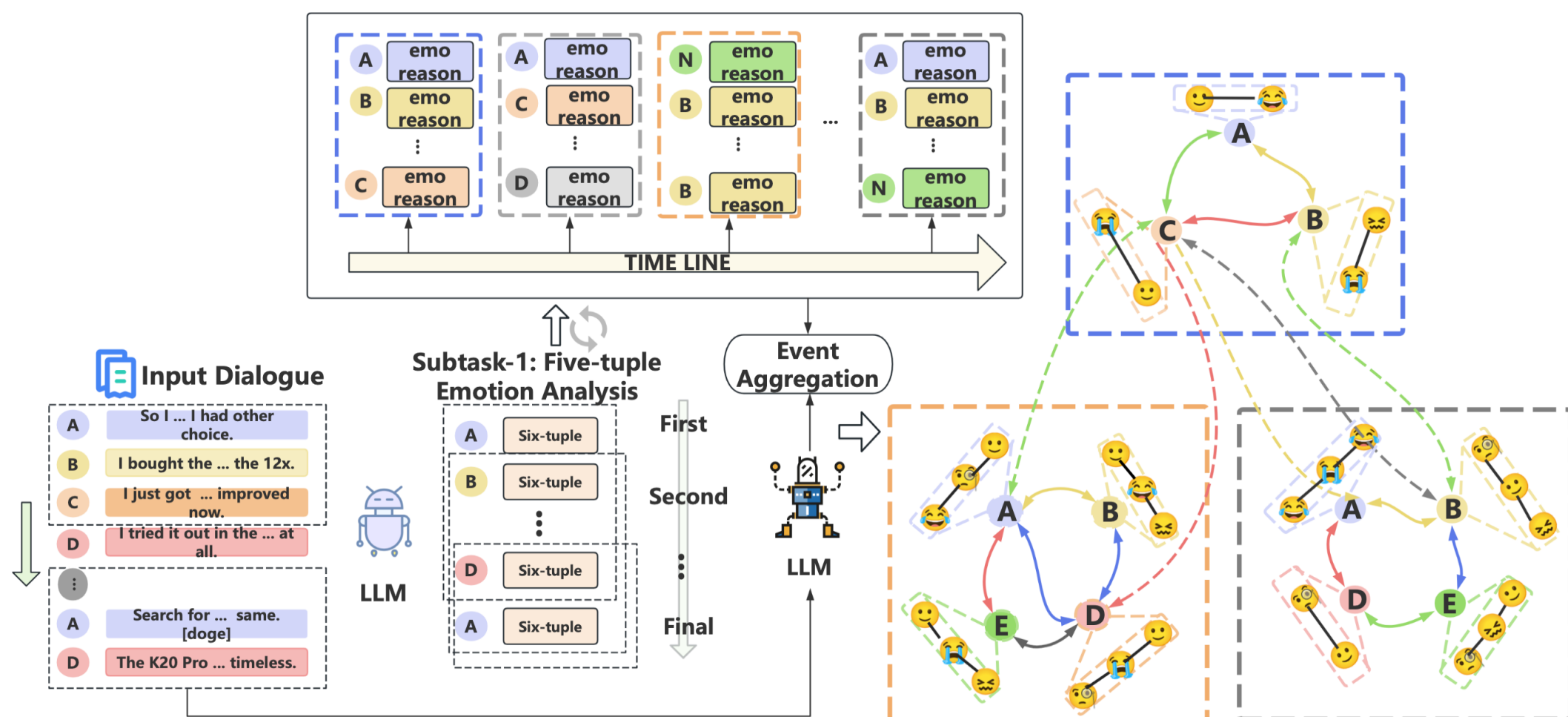
$$\arg\max_p \cos(\phi(u_t), \phi(E_p)) \geq \tau,$$

Where $\tau$ is the similarity threshold.

**3. Event Aggregation**: Combines multiple related events into a single, coherent event, enabling more consistent modeling of emotional changes in long dialogues. The storage optimization is represented by:

$$\mathcal{E}' = \left\{ E_p \mid \frac{1}{|I_p|} \sum_{i \in I_p} \cos(\phi(u_i), \phi(E_p)) \geq \tau_{\text{retain}} \right\},$$

Where $\tau_{\text{retain}}$ is the threshold.

## HELIX-6 Dataset

I constructed the HELIX-6 dataset, which extends ATLAS-6 by adding annotations for ultra-long conversations, multi-party emotional interactions, and emotional causality chains in complex scenarios. Key features include:

- 1,000 sets of longitudinal dialogues from authentic scenarios (each encompassing over 70 conversational turns)
- Complex interactive contexts such as pedagogical discussions and multi-party deliberative debates
- Diverse sources including curated public discourse from social platforms and context-specific dialogue recordings
- Four-stage annotation protocol: identification of pivotal events, documentation of emotional states, establishment of causal relationships, and analysis of participant influence patterns
- Five-category sentiment classification system (positive, negative, neutral, ambiguous, doubt)

| Benchmark | Modality | Scenario | ABSA | Causal chain | type |
|---|---|---|---|---|---|
| CR(Blitzer et al., 2007) | Text | Sentence | ✗ | ✗ | Short |
| Yelp(Tang et al., 2015) | Text | Document | ✗ | ✗ | Medium |
| SemEval(Pontiki et al., 2014) | Text | Sentence | ✗ | ✗ | Short |
| TOWE(Fan et al., 2019) | Text | Sentence | ✗ | ✗ | Short |
| ACOS(Cai et al., 2021) | Text | Sentence | ✗ | ✗ | Short |
| ASTE(Peng et al., 2019) | Text | Sentence | ✗ | ✗ | Short |
| DiaASQ(Li et al., 2023) | Text | Dialogue | ✗ | ✗ | Very Long |
| Twitter2015(Ma et al., 2017) | Text, Image | Sentence | ✗ | ✗ | Long |
| CMU-MOSEI(Bagher Zadeh et al., 2018) | Text, Audio, Video | Sentence | ✗ | ✗ | Short |
| IEMOCAP(Busso et al., 2008) | Text, Audio, Video | Dialogue | ✗ | ✗ | Long |
| MELD(Poria et al., 2019) | Text, Audio, Video | Dialogue | ✗ | ✗ | Medium |
| M3ED(Zhao et al., 2022a) | Text, Audio, Video | Dialogue | ✗ | ✗ | Short |
| PanoSent(Luo et al., 2024a) | Text, Image, Audio, Video | Dialogue | ✓ | ✗ | Medium |
| NTCIR-13(Gao et al., 2017) | Text | Clause | ✗ | ✗ | Short, Medium |
| RECCON(Poria et al., 2017) | Text | Dialogue | ✗ | ✓ | Medium |
| ATLAS-6(Zhang et al., 2025) | Text, Audio, Video | Dialogue | ✓ | ✓ | Long |
| HELIX-6(OURS) | Text, Audio, Video | Dialogue | ✓ | ✓ | Long, Very Long |

The table shows that HELIX-6 is uniquely positioned as the only dataset supporting both ABSA and causal chain analysis for Long and Very Long dialogues, addressing a critical gap in emotion causality research.