

# MIRAI: EVALUATING LLM AGENTS FOR INTERNATIONAL EVENT FORECASTING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We present MIRAI, a **benchmark** designed to systematically evaluate LLM agents as temporal forecasters to predict international events. Our benchmark features an agentic environment with tools for accessing an extensive database of historical, structured events and textual news articles. We refine the GDELT<sup>1</sup> event database with careful cleaning and parsing to curate a series of relational prediction tasks with varying forecasting horizons, assessing LLM agents’ abilities from short-term to long-term forecasting. We further implement APIs to enable LLM agents to utilize different tools via a code-based interface. Notably, MIRAI features a dynamic data construction pipeline that supports periodically downloading recent news and events, and automatically generates the most recent test split. This allows us to evaluate any newly released model in a contamination-free manner as we can always construct a test split later than its knowledge cutoff date. MIRAI comprehensively evaluates the agents’ capabilities in three dimensions: 1) autonomously source and integrate critical information from large global databases; 2) write codes with both domain-specific APIs and libraries for tool-use; and 3) jointly reason over historical knowledge from diverse formats and timespan to accurately predict future events. We establish a benchmark for evaluating LLM agents’ international event forecasting abilities, advancing the development of more reliable models for analyzing international relations.<sup>2</sup>

## 1 INTRODUCTION

Accurate forecasting of international events is crucial (Brown & Lee, 2018), as understanding the evolution of geopolitical developments enables stakeholders to make well-informed decisions, mitigate risks, and seize opportunities in the interconnected world. Traditionally, researchers in international relations rely on domain expertise (Smith & Doe, 2020; Johnson & Roberts, 2019). They conduct detailed analyses of the complex interplay among nations, considering alliances, trade agreements, ideological affinities, and historical rivalries to forecast events such as conflicts, collaborations, or alliance shifts (Davis & Nguyen, 2017). With the rapid development of deep learning techniques, forecasting through data-driven neural networks becomes an attractive alternative. Despite their success, current methods rely on single types of information—either structured knowledge graphs (Mahdizoltani et al., 2015; Jin et al., 2020b; Li et al., 2021c) or textual datasets (Zou et al., 2022; Reddy et al., 2023). Knowledge graphs, although organized, can suffer from incompleteness (Huang et al., 2023; Galárraga et al., 2017) or bias (Huang et al., 2024), while textual analyses can lack necessary factual groundings of their reasoning to historical evidence, which compromises the interpretability and validation of their forecasts. These limitations raise concerns about AI forecasters’ reliability, particularly for high-stake scenarios (McClean et al., 2009).

Large Language Model (LLM) agents present a promising path to overcome these challenges (Sumers et al., 2024; Liu et al., 2023b; Weng, 2023; noa; Wang et al., 2023). These advanced AI systems exhibit the potential to mimic human experts by utilizing a diverse set of tools to automatically gather and process information from various sources, including text, knowledge graphs, and numerical data (Shen et al., 2023; Lu et al., 2023b; Zhuang et al., 2023b; Li et al., 2023b). Trained on extensive textual corpora, LLMs (Achiam et al., 2023; Anthropic, 2023; DeepMind, 2023; Touvron et al., 2023) are capable of grasping the subtleties of international relations, reasoning through complex relationships with linguistic explanations, and planning their tool usage effectively (Yuan et al.,

<sup>1</sup>GDELT: <https://www.gdeltproject.org/>

<sup>2</sup>We released our anonymous code repository for the data construction pipeline [here](#), constructed benchmark [here](#), and an interactive agent demo [here](#).

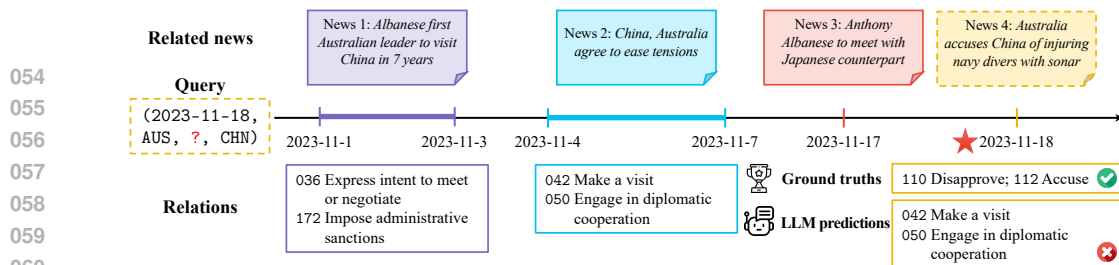


Figure 1: An example of forecasting the relations between Australia and China on 2023-11-18. The agent fails to predict the change of relation and makes a wrong forecast.

2023; Liu et al., 2023a; Valmeekam et al., 2023; Ma et al., 2023b). Such capability opens up new possibilities for developing transparent and interpretable forecasting models that can be further scrutinized and refined.

Despite the immense potential of LLM agents for event forecasting, there is a lack of standardized benchmarks to assess their forecasting abilities in the realm of intricate international events. To address this gap, we introduce **MIRAI** (Multi-Information FoRecasting Agent Interface), the first forecasting benchmark designed as an agentic environment with rich structured and textual information. Leveraging the timely-updated Global Database of Events, Language, and Tone (GDELT) (Leetaru & Schrodt, 2013b), MIRAI adapts real-world international event data to create event-forecasting tasks in different test timeframes and across various time horizons. Our preparation script of MIRAI features a dynamic data construction pipeline that supports periodically downloading recent news and events, and automatically generate the most recent test split. This unique design allows us to consistently generate new **contamination-free** test sets, such that we can evaluate the forecasting capability of any recent LLM as long as its training data cutoff is before our split—a critical feature absent in previous forecasting benchmarks that only provided fixed datasets often predating most LLMs’ knowledge cutoffs. Furthermore, MIRAI’s agentic environment enables LLMs to interact with both relational and textual databases through APIs, facilitating autonomous information gathering, processing, and application in a contextually relevant manner.

We evaluate both open and closed source LLMs on MIRAI with ReAct-style (Yao et al., 2023b) agents using “Single Function” and “Code Block” action types. Our extensive experiments spanned multiple test splits (2023-11, 2023-12, 2024-01, 2024-02) and revealed: 1) Temporal forecasting tasks are challenging for LLM agents, with the highest-performing GPT-4o-mini agent using full suite of APIs achieves a 30.3 F1 score in second-level relation prediction tasks; while tasks involving long-term and fine-grained event forecasting are even more challenging; 2) The “Code Block” tool-use strategy, which allows more flexible interactions, demands robust code generation capabilities. E.g. GPT-4o-mini is able to better **utilize** and **benefit** from this strategy than other models we evaluate.

These findings emphasize the need for ongoing research into temporal reasoning and the effective use of tools by LLM agents. We expect that MIRAI could serve as a standard benchmark for evaluating LLMs in event forecasting. This would support the development of more precise and reliable models for political analysis, enhancing our understanding of global dynamics. To facilitate further research and development in this area, we release the code for our dataset construction pipeline and commit to updating our dataset split every month, ensuring that MIRAI remains a contamination-free and challenging benchmark for assessing LLM agent capabilities in international event forecasting.

In summary, our contributions are two-fold:

- We present MIRAI as a **comprehensive benchmark** uniquely combining three critical aspects for evaluating LLM agents’ temporal forecasting capabilities:
  - An agentic environment with APIs for information integration, tool use, and reasoning.
  - Diverse data from structured events and textual news.
  - A dynamic pipeline ensuring contamination-free test splits beyond model knowledge cutoffs.
- Extensive experiments across models, relations, horizons, and temporal splits reveal key challenges in reasoning and tool use, offering insights and directions for advancing LLM forecasting agents.

## 2 THE MIRAI BENCHMARK

In this section, we introduce MIRAI benchmark from: the specifics of the data and task (Sec. 2.1), the implemented agents and environments (Sec. 2.2), and the database construction details (Sec. 2.3).

### 2.1 DATA AND TASK

We introduce MIRAI, a benchmark crafted for evaluating LLM agents for temporal forecasting in the realm of international events, with tool-use and complex reasoning. We consider *forecasting* as the

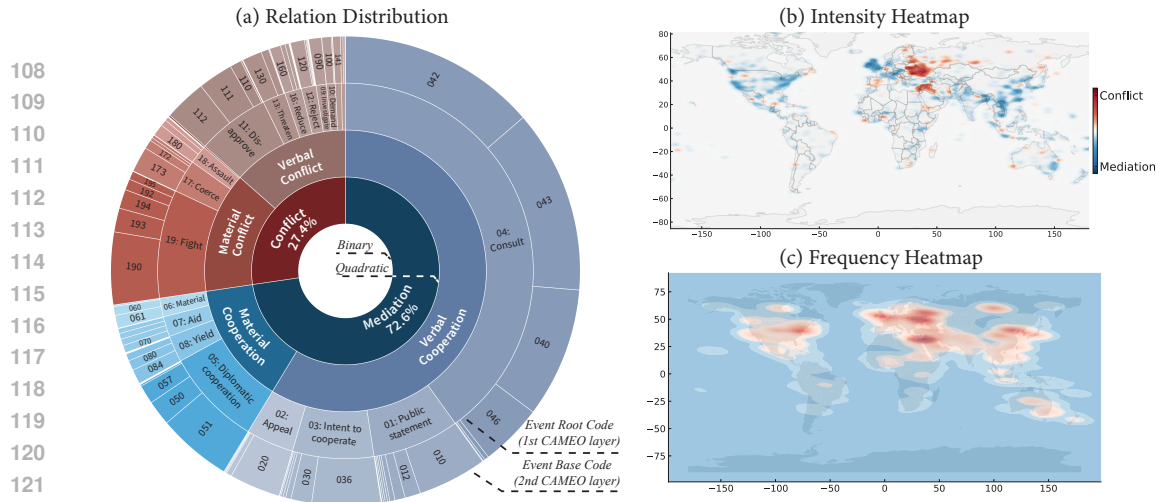


Figure 2: **MIRAI comprehensive global event coverage.** (a) Circular chart: Relation hierarchy and distribution in MIRAI. (b) Heatmap: Global events intensity, including areas of conflict (red) and mediation (blue). (c) Heatmap: Event frequency by region.

process of collecting essential historical data and performing temporal reasoning to anticipate the outcomes of future events.

**Structured and Textual Event Representations.** Formally, we represent an event as  $e^t = (t, s, r, o)$  corresponding to  $D_e^t$ , where  $t$  is the timestamp<sup>3</sup>,  $s, o \in \mathcal{C}$  are respectively the subject and object countries<sup>4</sup> from the country pool  $\mathcal{C}$ ,  $r \in \mathcal{R}$  denotes the relation type defined by CAMEO ontology<sup>5</sup> (Boschee et al., 2015a), and  $D_e^t$  is the set of source news articles that mentioned event  $e$  at timestamp  $t$ . Events at timestamp  $t$  form a set  $E^t = \{e_1^t, \dots, e_M^t\}$ , where  $M$  is the number of unique events at time  $t$ , with  $\{E^t\}_{t=1}^T$  being able to organize to temporal graphs where countries are nodes and relations are edges. Correspondingly,  $D^t = \{d_1^t, \dots, d_N^t\}$  is the set of all news articles at  $t$ , where  $N$  is the number of unique news articles at time  $t$ , with  $\{D^t\}_{t=1}^T$  as the full document collection. This dual representation of structured events and their associated textual information allows for a comprehensive analysis of international events, leveraging both the concise, categorized nature of the event tuples and the rich contextual details provided by the news articles.

**Hierarchical Event Categories.** As shown in Figure 2a, we incorporate two hierarchical relation levels from the CAMEO ontology to facilitate a detailed and comprehensive analysis of geopolitical dynamics. The first level includes 20 broad categories, represented by a two-digit code (e.g., “01: Public Statement” or “04: Consult”), which are subdivided into second-level categories identified by a three-digit code that corresponds to its parent category. For example, “03: Express intent to cooperate” is a first-level category that includes 10 different second-level relations such as “036: Express intent to meet”. Subsequently, the quadruple “(2023-11-03, AUS, 036, CHN)” denotes that on 3 November 2023, the Australian leader announces a planned visit to China. These relations are also organized along two dimensions, from **Verbal** to **Material** and from **Conflict** to **Cooperation**, to form a quadratic categorization in the inner circle of Figure 2a.

**Event Forecasting Task.** The task of event forecasting  $(t + l, s, ?, o)$  is to predict all the events between a pair of countries  $s$  and  $o$ , happening  $l$  days in the future from the current time  $t$ . Formally, given historical events  $E^{\leq t}$  and associated news articles  $D^{\leq t}$  up to the current time  $t$ , our goal is to forecast future relationships  $E_{s,o}^{t+l}$  between a pair of countries. This includes utilizing both structured event data and textual information, considering interactions not only between the target countries but also involving third parties, such as interactions with their mutual neighbors.

A forecasting task example in Figure 1 shows predicting Australia’s actions towards China on 18 November 2023, based on information up to 17 November 2023. The query is formatted as “(2023-11-17 + [1 day], AUS, ?, CHN)”, with a temporal distance of one day. Historical events show long-standing tensions between the two countries. Despite recent news of the Australian leader’s visit to China and agreement to ease tensions, the overall relationship remains characterized by

<sup>3</sup>Each timestamp uniquely represents a day, formatted in “YYYY-MM-DD”.

<sup>4</sup>In this paper, the term “country” includes all countries, dependent territories, special geographic areas, and their subdivisions. We use the standardized ISO-3166 Alpha-3 codes for country names, e.g. “AUS” for Australia.

<sup>5</sup>CAMEO, the Conflict and Mediation Event Observations, is a well-established ontology meticulously developed by domain experts over years, for categorizing international political events across multiple levels of granularity.

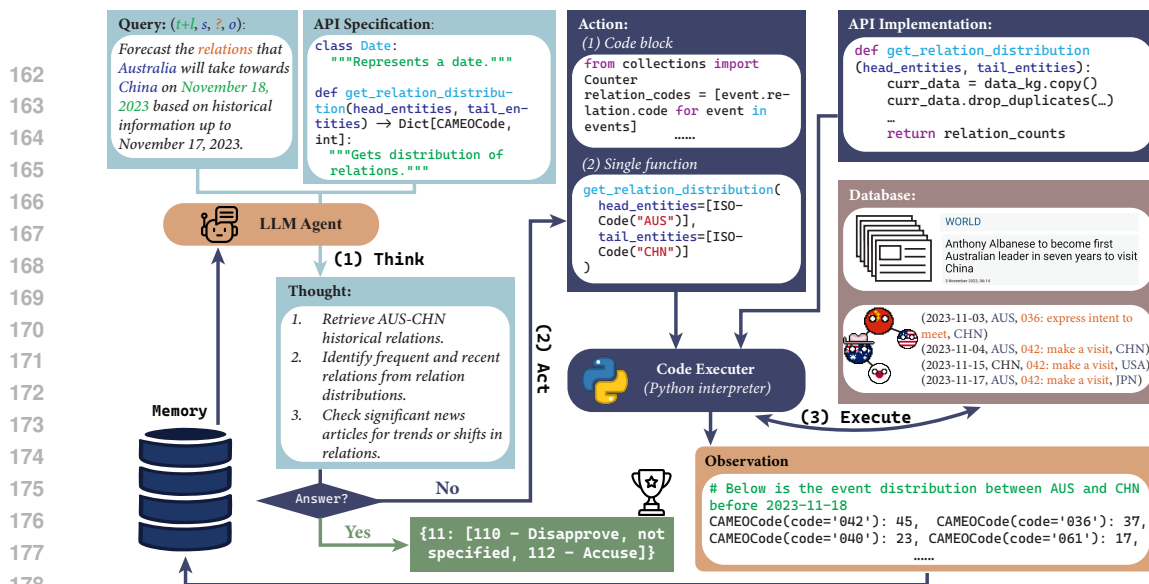


Figure 3: **Overview of the LLM agent’s forecasting process using the ReAct strategy.** The framework consists of three main steps: (1) Think: Agent analyzes status and plans next action based on query and API specs. (2) Act: The agent generates a “Single Function” call or a “Code Block” to retrieve and analyze relevant data. (3) Execute: Python interpreter runs generated code, producing observations. These steps are repeated performed until the agent reaches final forecast.

**Disapprove** and **Accuse** actions on certain focuses. The agent, overly relying on short-term news, incorrectly predicts **Diplomatic cooperation**.

## 2.2 AGENTS AND ENVIRONMENTS

Similar to human political analysts, LLM agents must leverage a variety of information sources to make reliable predictions. We provide an environment with coding APIs to facilitate flexible access to various knowledge sources. Our LLM agent uses these APIs to interact with the environment for forecasting through the ReAct strategy (Yao et al., 2023a), characterized by the iterative process of *think*, *act*, and *observe*. This approach enables the agent to analyze the situation, retrieve data, and observe outcomes to make informed forecasts.

**APIs.** We provide the LLM agent with a comprehensive set of APIs to access a rich database of historical events and news articles. The API contains the essential *data classes* and *functions* designed for the various types of information within the database. *Data classes* cover unary types (e.g., date, country, relation) and composite types (e.g., date range, event, news articles). *Functions*, executable in Python, enable efficient database querying. The function design is characterized in two aspects:

- **Information types.** These functions cover diverse information types: country and relation mappings, hierarchies, and events and news articles statistics (counts, listings, and distributions). For example, `map_relation_description_to_cameo` takes a relation description and returns the five most likely relations with their CAMEO codes, names, and descriptions, providing precise relation information as needed.
- **Search conditions.** Additionally, the API functions for events and news articles offer optional parameters for tailored searches based on different criteria. For instance, `get_event` allows searches specifying conditions like `date_range`, `head_entities`, `tail_entities`, `relations`, and `text_description` to retrieve specific events from the database.

The data classes and functions provided in the API are shown in Appendix E.

**Interactions with Environment.** The environment is equipped with a Python code sandbox with full API and database access. MIRAI informs agents how to interact with the environment through API Specifications, including detailed Python docstrings for all data classes and functions that abstract implementation details. Agents use the ReAct strategy to iteratively alternate between *think*, *act*, and *observe* to gather information and forecast. The pipeline is illustrated in Figure 3.

1. **Think.** Specifically, the agent first *thinks* about the current situation. If confident based on the current information and analysis, it delivers a final forecast and stops. Otherwise, it plans further actions. In the first step shown in Figure 3, the agent plans to first retrieve recent events directly involving the two countries and then verify the details by checking related news articles.



- 216 2. **Act.** Next, the agent *acts* by generating executable codes to interact with the environment in  
 217 two forms: “Single Function” for straightforward data retrieval or “Code Block” for complex  
 218 operations including loops and conditionals. The agent can use API-defined functions and safe,  
 219 well-established Python libraries such as `numpy`, `networkx`, and `scikit-learn`. Figure 3  
 220 illustrates that the agent can write either a “Code Block” with an imported library and an inline  
 221 loop, or a “Single Function” call with targeted countries as an argument.
- 222 3. **Observe.** The environment executes the agent’s code using a Python interpreter with full API  
 223 and database access. All execution results are passed back to the agent as ‘Observations’. For  
 224 successful executions, the agent *observes* the corresponding output. For failures, it *observes* a  
 225 specified error message along with the error type. In the figure, execution results are sent to the  
 226 agent’s memory for the next iteration.

227 These ‘Thought’, ‘Action’, and ‘Observation’ from the previous iterations are stored in the agent’s  
 228 memory base and used as the context for subsequent steps. Variables defined in previously generated  
 229 code remain available for future actions. Examples of agents performing forecasting using “Code  
 230 Block” and “Single Function” are shown in Appendix F. Prompts are shown in Appendix I.

### 231 2.3 DATASET CONSTRUCTION

232 Our database construction process involves three main steps: (1) raw data collection, (2) structured  
 233 data cleaning, and (3) textual data processing. Throughout this process, we carefully consider critical  
 234 aspects such as preventing test information leakage, standardizing information formats, ensuring data  
 235 accuracy, verifying source reliability, and addressing ethical concerns.

236 **Raw Data Collection.** We construct the database based on the GDELT project<sup>6</sup>, which captures  
 237 global news media and extracts event information every 15 minutes, containing attributes about the  
 238 event date, actor, action (relation), geography, and source news. The dataset used in this paper spans  
 239 from January 1, 2023, upto February 29, 2024, providing a comprehensive base for our multi-split  
 240 test design. Given GDELT’s rapid updates and diverse sources, rigorous cleaning is essential.

241 **Structured Data Cleaning.** We refine the structured event data through four key steps: (1) *Temporal*  
 242 *alignment*: Event dates are aligned with news publication dates to prevent information leakage. (2)  
 243 *Event actor standardization*: Country codes are standardized to ISO-3166 norms, excluding missing  
 244 or outdated codes to ensure consistency and eliminate ambiguity. (3) *Event scope delimitation*: We  
 245 filter out all domestic events where the subject is identical to the object at the country level, focusing  
 246 exclusively on international events to maintain relevance to our research objectives. (4) *Event relation*  
 247 *standardization*: Event codes are standardized to the second level of the CAMEO ontology, balancing  
 248 specificity with consistency and reliability. Third-level relations are omitted due to inconsistent  
 249 hierarchical depth and increased risk of extraction errors at this granular level

250 **Textual Data Processing.** For associated news articles, we implement the following steps: (1)  
 251 *Source reliability threshold*: We retain events with at least 50 daily news mentions. This reduces  
 252 the influence of less reliable sources such as personal blogs and decreases erroneously extracted  
 253 events. (2) *Textual context extraction*: News titles and content are sourced from the corresponding  
 254 URLs of each event, retaining only textual information. (3) *Text cleaning*: We apply rigorous text  
 255 cleaning procedures to reduce noise while enhancing the **reliability and ethical integrity** of the  
 256 textual information, following the OBELICS protocol (Laurençon et al., 2023). This process operates  
 257 at both paragraph and document levels, filtering low-value content with low word counts or high  
 258 character/word repetition ratios and removing excessive special characters. Importantly, we employ a  
 259 list of flagged words to identify and eliminate potentially sensitive or inappropriate content, aligning  
 our data collection with ethical standards.

260 **Test Splits Construction.** We construct multiple test splits (2023-11, 2023-12, 2024-01, 2024-02),  
 261 each covering a one-month period from November 2023 to February 2024. For each test split: (1)  
 262 *Enhanced filtering*: We apply higher thresholds (100 daily mentions, 5 news articles) to ensure test  
 263 data quality and reliability. (2) *Query formation*: We construct  $(t, s, ?, o)$  queries, with answers listing  
 264 relations between countries at time  $t$ . (3) *Balanced sampling*: For each month, we sample 100 queries  
 265 to form balanced test splits, ensuring representation across dates, countries, and CAMEO code types.  
 266 This multi-split design evaluates model performance across different time periods, assessing temporal  
 267 robustness of forecasting capabilities. Using the same processing script, we can generate test split for  
 268 any following month keep MIRAI a contamination-free and challenging benchmark.

269 <sup>6</sup><https://www.gdeltproject.org/>

Table 1: Evaluation results with different agent tools and the tool-use strategies. The best-performing score is highlighted in **bold** and the second-best is underlined.

Agent	Tool-Use		Binary KL ( $\downarrow$ )	Quad KL ( $\downarrow$ )	First-level Relation (%)			Second-level Relation (%)		
	Action Type	API			Pre. ( $\uparrow$ )	Rec. ( $\uparrow$ )	F1 ( $\uparrow$ )	Pre. ( $\uparrow$ )	Rec. ( $\uparrow$ )	F1 ( $\uparrow$ )
Direct IO	—	—	6.5 $\pm$ 1.7	15.9 $\pm$ 1.5	27.6 $\pm$ 8.1	19.7 $\pm$ 5.9	18.8 $\pm$ 6.9	6.6 $\pm$ 1.5	5.1 $\pm$ 0.4	3.5 $\pm$ 0.8
ZS-CoT	—	—	6.9 $\pm$ 0.8	10.1 $\pm$ 0.8	27.6 $\pm$ 4.0	36.0 $\pm$ 4.5	26.7 $\pm$ 4.1	10.2 $\pm$ 1.4	17.4 $\pm$ 1.1	10.5 $\pm$ 0.7
ReAct	Single Function	<i>Event-Only</i>	33.5 $\pm$ 0.7	<u>6.7<math>\pm</math>0.7</u>	<u>44.3<math>\pm</math>3.9</u>	<u>54.2<math>\pm</math>3.9</u>	<u>41.4<math>\pm</math>1.7</u>	<u>25.3<math>\pm</math>2.6</u>	<u>47.4<math>\pm</math>2.4</u>	<u>26.9<math>\pm</math>1.9</u>
	Single Function	<i>News-Only</i>	6.1 $\pm$ 1.0	12.8 $\pm$ 0.6	27.8 $\pm$ 3.1	25.9 $\pm$ 2.9	21.8 $\pm$ 2.3	6.3 $\pm$ 2.2	9.0 $\pm$ 2.0	5.4 $\pm$ 1.3
ReAct	Single Function	<i>All</i>	<b>3.1<math>\pm</math>0.5</b>	<b>5.9<math>\pm</math>1.0</b>	<b>47.6<math>\pm</math>5.8</b>	<b>58.3<math>\pm</math>2.6</b>	<b>44.2<math>\pm</math>4.0</b>	<b>28.7<math>\pm</math>3.9</b>	<b>51.0<math>\pm</math>4.0</b>	<b>29.6<math>\pm</math>3.7</b>
	Code Block	<i>All</i>	<u>5.1<math>\pm</math>0.9</u>	8.9 $\pm$ 0.5	27.1 $\pm$ 4.0	38.6 $\pm$ 2.5	25.9 $\pm$ 2.2	11.6 $\pm$ 2.4	26.3 $\pm$ 2.0	12.6 $\pm$ 1.7

**Analysis and Documentation.** The resulting dataset contains 1,296,991 GDELT event records, corresponding to 75,341 unique  $(t, s, r, o)$  events and 401,013 unique news articles. Figures 2b and 2c illustrate the global distribution of our curated events, highlighting varying intensities of conflict and mediation across regions. We provide additional details of the benchmark data and human evaluation of the data quality in Appendix G. A standardized datasheet (Geburu et al., 2021) for MIRAI is in Appendix J, clearly and comprehensively documenting its motivation, composition, collection process, recommended uses, distribution, and maintenance.

**Evaluation metrics.** We instruct the agent to generate forecasts in a JSON dictionary, using two-digit first-level CAMEO codes as keys and lists of three-digit second-level codes as values. Evaluation involves calculating *precision*, *recall*, and *F1 score* between the predicted and ground-truth lists. Moreover, we map the predicted and ground truth relations to their respective binary and quadratic classes (as shown in Figure 2a). To measure prediction-ground truth alignment, we employ the *empirical Kullback-Leibler (KL) divergence*:  $D_{KL}(P||Q) = \sum_i P(i) \log(P(i)/Q(i))$ , where  $P$  and  $Q$  represent the frequencies of ground-truth and predicted relations respectively. A lower KL divergence indicates a better alignment of the model’s predictions with the ground-truth list.

### 3 EXPERIMENTS

In experiments, we compare forecasting performance across various agent methods with different tool-use (Sec. 3.1). Additionally, we evaluate the impact of different base models on forecasting (Sec. 3.2), and analyze contamination of the models’ knowledge cutoff over test-time splits (Sec. 3.3). Finally, we conduct an in-depth analysis of key factors to agents’ behavior (Sec. 3.4).

#### 3.1 EVALUATE FORECASTING WITH DIFFERENT AGENT METHODS AND TOOLS

We investigate the effect of different tools (APIs) and agent tool-use strategies. We use GPT-3.5-Turbo (gpt, 2023) as the base model and evaluate on the 705 unsampled test queries in 2023-11. For all experiments, we set the model temperature to 0.4 and run 5 times to calculate the mean and standard deviation. We provide the detailed prompts in Appendix I.

We consider two agent implementations without tool-use: **Direct IO** and **Zero-Shot Chain-of-Thought (ZS-CoT)** (Kojima et al., 2022; Wei et al., 2023). Direct IO let the LLM directly provide answers using only its internal knowledge, serving as a baseline to reflect its internal world knowledge. ZS-CoT prompts the LLM for step-by-step reasoning before final prediction.

We implement **ReAct** (Yao et al., 2023b) for tool-use agents with two variants: 1) “Single Function” and 2) “Code Block” (detailed in Sec. 2.2). ReAct agents interact with our provided environments through an iterative process of thinking, acting, and observing. We thus further create API variants with access to 1) *News-Only* APIs, 2) *Event-Only* APIs, or 3) *All* API data classes and functions.

The experimental results in Table 1 reveal several key insights into agent performance:

**1) MIRAI presents a challenging task for LLM agents.** The best agent (ReAct with “Single Function” using all APIs) for second-level relation predictions achieves a precision of 28.7 and an F1 score of 29.6. These results underscore the complexity and difficulty of the temporal forecasting tasks in MIRAI and highlight the substantial room for improvement in LLM agents for event forecasting.

**2) Predicting fine-grained relations proves more difficult.** All models exhibit higher KL divergence for quadratic than binary classes, and lower F1 scores for second-level predictions compared to first-level ones. These findings confirm that predicting fine-grained relation types is more challenging.

**3) Diverse tool-use is critical for temporal forecasting.** ZS-CoT and Direct-IO, which rely solely on the internal world knowledge of LLMs for forecasting without tool-use, significantly underperform the ReAct agent with full API access to the database. This emphasizes the importance of basing forecasting and reasoning on retrieved historical data and knowledge. In terms of tool types, ReAct agents using *News-Only* APIs perform much worse than agents with *Event-Only* APIs. While news articles provide detailed context for events, they can also introduce noise and lead to issues such as

Table 2: Evaluation results on the 2024-02 test split using different base LLMs and action types. The best-performing score is highlighted in **bold** and the second-best is underlined. **Note that the traditional methods are task-specifically trained on training data up to 2023-06, while all LLM-based methods are zero-shot with their general pretrained knowledge. Human evaluation is conducted on a subset of 51 test events due to resource constraints and the time-intensive nature of expert evaluation.**

Base LLM	Training Data Cutoff Date	Action Type	Binary KL ( $\downarrow$ )	Quad KL ( $\downarrow$ )	First-level Relation (%)			Second-level Relation (%)		
					Pre. ( $\uparrow$ )	Rec. ( $\uparrow$ )	F1 ( $\uparrow$ )	Pre. ( $\uparrow$ )	Rec. ( $\uparrow$ )	F1 ( $\uparrow$ )
Mistral-7B-Instruct-v0.2	2023-12	Single Func	10.3 $\pm$ 1.7	14.2 $\pm$ 1.9	38.1 $\pm$ 0.5	19.2 $\pm$ 4.2	18.9 $\pm$ 1.1	21.9 $\pm$ 4.1	9.8 $\pm$ 3.5	9.3 $\pm$ 0.6
		Code Block	9.1 $\pm$ 2.3	14.3 $\pm$ 1.6	31.3 $\pm$ 5.1	12.5 $\pm$ 2.6	15.1 $\pm$ 3.6	13.1 $\pm$ 1.2	9.4 $\pm$ 1.6	8.4 $\pm$ 1.9
Llama-3-8B-Instruct	2023-03	Single Func	9.0 $\pm$ 2.4	14.1 $\pm$ 1.6	39.8 $\pm$ 1.6	15.6 $\pm$ 1.2	18.6 $\pm$ 0.2	15.8 $\pm$ 0.5	11.8 $\pm$ 0.6	10.3 $\pm$ 0.1
		Code Block	9.4 $\pm$ 2.6	14.5 $\pm$ 1.5	39.5 $\pm$ 0.7	12.2 $\pm$ 2.4	15.9 $\pm$ 2.5	18.9 $\pm$ 2.0	8.9 $\pm$ 1.7	9.2 $\pm$ 0.9
Llama-3.1-8B-Instruct	2023-12	Single Func	7.8 $\pm$ 2.3	12.4 $\pm$ 3.0	50.1 $\pm$ 1.6	21.5 $\pm$ 2.2	25.2 $\pm$ 2.1	22.7 $\pm$ 3.2	16.8 $\pm$ 0.5	14.8 $\pm$ 0.7
		Code Block	9.1 $\pm$ 2.3	14.4 $\pm$ 1.5	35.7 $\pm$ 1.3	15.9 $\pm$ 0.0	18.3 $\pm$ 0.1	14.7 $\pm$ 0.1	12.8 $\pm$ 1.0	10.0 $\pm$ 0.2
GPT-3.5-Turbo	2021-09	Single Func	3.5 $\pm$ 1.3	7.5 $\pm$ 2.7	55.7 $\pm$ 5.8	40.9 $\pm$ 2.6	38.3 $\pm$ 4.2	42.4 $\pm$ 4.4	34.3 $\pm$ 5.0	28.3 $\pm$ 3.9
		Code Block	5.3 $\pm$ 1.4	9.3 $\pm$ 2.3	34.8 $\pm$ 11.7	34.9 $\pm$ 0.9	26.7 $\pm$ 3.8	16.8 $\pm$ 5.9	26.4 $\pm$ 0.5	15.1 $\pm$ 2.6
GPT-4o-mini	2023-10	Single Func	4.0 $\pm$ 0.9	8.1 $\pm$ 1.3	61.3 $\pm$ 11.4	34.7 $\pm$ 1.9	39.0 $\pm$ 6.0	40.0 $\pm$ 5.5	32.6 $\pm$ 1.6	29.7 $\pm$ 3.8
		Code Block	3.2 $\pm$ 0.8	7.7 $\pm$ 1.9	59.8 $\pm$ 5.1	37.1 $\pm$ 0.4	40.0 $\pm$ 3.1	46.5 $\pm$ 2.1	29.7 $\pm$ 0.6	30.3 $\pm$ 1.0
<i>Traditional Forecasting Methods</i>		<i>Task-specific</i>								
REGCN (Li et al., 2021b)	2023-10	✓	0.3	1.0	24.8	78.2	32.4	3.9	25.7	5.6
ForecastQA (Jin et al., 2021a)	2023-10	✓	9.7	13.8	55.0	16.2	22.1	40.0	8.8	12.0
<b>Human evaluators</b>	—	—	0.04	1.37	62.73	88.70	68.29	54.54	74.53	56.78

excessively long context, posing additional challenges for LLM agents. Moreover, the agents using both types of information achieve the optimal results.

### 3.2 EVALUATE FORECASTING WITH DIFFERENT BASE LLMs

We then investigate the role of the base LLMs in agent’s performance. We evaluate both open-sourced LLMs, including Mistral-7B-Instruct-v0.2 (Jiang et al., 2023), Llama-3-8B-Instruct, and Llama-3.1-8B-Instruct (Dubey et al., 2024), as well as close-sourced LLMs including GPT-3.5-Turbo (gpt, 2023) and GPT-4o-mini (gpt, 2024). Comparisons are done on the **2024-02** test split that is after all models’ training data cutoff date, which comprise 100 data-balanced queries. All models use ReAct framework with access to all APIs. The action types can be either “Single Function” or “Code Block” with a maximum tool call limit set to 20 steps. The same prompt is used across all models for fair comparison, as detailed in Appendix I. The experimental results are presented in Table 2, and we observe the following findings:

**1) Code Block benefits stronger LLMs but hurts weaker models:** The “Code Block” action type provides greater flexibility than the “Single Function”. However, its benefits are not uniformly achieved across all models. Small open-sourced models and GPT-3.5-Turbo show reduced performance with Code Block, while GPT-4o-mini as a strong model gains improvements from it. This indicates that the ability to generate coherent and practical long Code Blocks is a distinguishing factor that determines LLMs as reliable forecasting agents.

**2) GPT-4o-mini outperforms other models:** GPT-4o-mini achieves the highest performance across many metrics in different levels. Notably, for second-level relation prediction, GPT-4o-mini achieves F1 scores of 29.7 and 30.3 using “Single Function” and “Code Block”, surpassing all other models. Among the tested open-sourced smaller models, Llama-3.1-8B-Instruct leads the performance but still remains a significant performance gap to larger models. This indicates that MISTRAL is hard enough, and can effectively distinguish different LLMs’ reasoning capabilities.

**3) MISTRAL naturally supports evaluating traditional forecasting methods:** We choose a Temporal Knowledge Graph (TKG) method REGCN (Li et al., 2021b) and a textual method ForecastQA (Jin et al., 2021a) as examples, trained both models on data up to 2023-06. These trained methods show strong results on first-level relation prediction, but fall short for fine-grained second-level relation prediction. We defer the implementation details and experimental discussions to Appendix D.3 and provide human forecasting performance as a reference in Appendix D.4.

We further conduct the following analysis to better understand agent capabilities:

**Can we make a small LM stronger via scaling inference-time compute?** It is evident that stronger LLMs have better agent performance; however, can we enhance a weaker LLM to achieve comparable forecasting performance using inference-time computation? To explore this, we take Mistral-7B-Instruct-v0.2 as the base LLM with ReAct using the “Single Function” strategy. For each query, we perform multiple sampling at a temperature of 0.4. We then consider a variant of *self-consistency*, which only keeps entries appearing more than twice. We also calculate F1 (Max@K), which assesses the F1 score for each instance and keeps the maximum score across all rounds. As shown in Figure 4a, with more samples, the performance of Mistral-7B-Instruct-v0.2 significantly improves. Initially, a

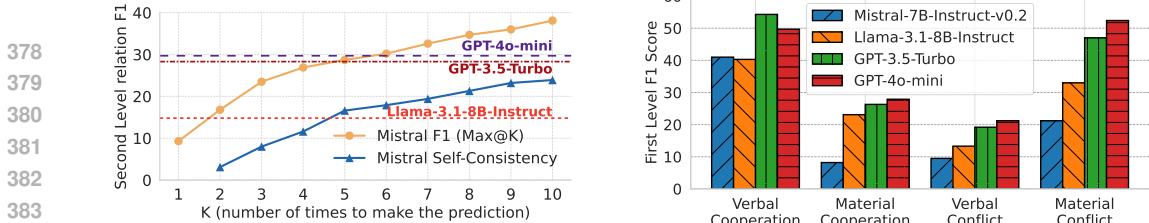


Figure 4: a) Self-consistency of Mistral-7B-Instruct-v0.2 model increases with more samples. b) F1 scores of different LLMs on relation prediction, categorized based on the quadratic classes.

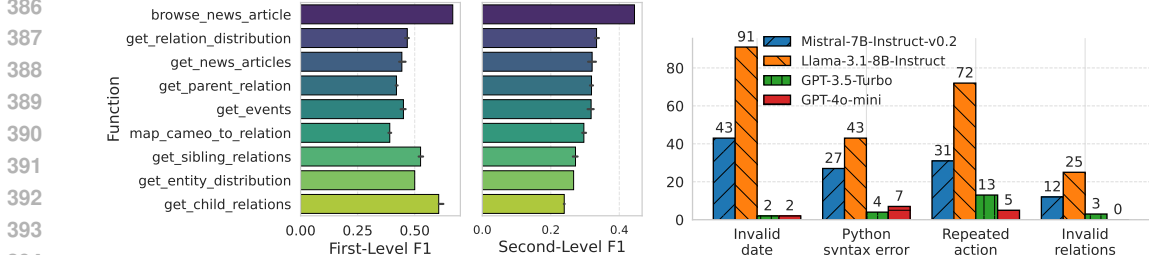


Figure 5: a) F1 Accuracy for each API function. b) Code execution error analysis for different LLMs.

single sample achieves an F1 score of 9.3, which is considerably lower than the scores achieved by larger models. However, as more solutions being sampled, the F1 (Max@K) and self-consistency of Mistral-7B-Instruct-v0.2 improve progressively. By the 10th samples, the Max@K reaches 38.1, even surpassing GPT-4o-mini. This result highlights the potential of inference-time search methods like self-consistency to push the boundaries of smaller language models in event forecasting.

**Code execution error analysis.** Our implemented agents interact with tools via code but often encounter execution error. We summarize the dominating error types for different LLMs in Figure 5b. We observe the invalid date as the most frequent error for smaller models, showing their difficulty in understanding time restrictions on the available historical data, which is set to before the query event date. Also, agents including larger models may propose repeated actions that have been conducted in its previous ReAct steps, failing to generate effective reasoning traces. We find Llama-3.1-8B-Instruct makes more execution errors than Mistral-7B-Instruct-v0.2 but achieves a better forecasting performance. We conduct further analysis on agents’ final status in Appendix D.1 on their ReAct sequence length in Appendix D.2. Overall, GPT-4o-mini makes significantly fewer execution errors. This enhanced code generation capability contributes to its superior performance for event forecasting.

### 3.3 ANALYZE CONTAMINATION OF KNOWLEDGE-CUTOFF OVER TEST-TIME SPLITS

Model	Training Data Cutoff Date	Test-Month Splits			
		2023-11	2023-12	2024-01	2024-02
Llama-3-8B-Instruct	2023-03	6.1 <sub>±1.5</sub>	8.7 <sub>±1.8</sub>	8.7 <sub>±0.1</sub>	10.3 <sub>±0.1</sub>
Llama-3.1-8B-Instruct	2023-12	15.8 <sub>±5.7</sub>	15.9 <sub>±3.6</sub>	16.3 <sub>±2.0</sub>	14.8 <sub>±0.7</sub>
GPT-4-Turbo	2023-12	<b>30.0</b> <sub>±1.9</sub>	<u>25.8</u> <sub>±3.1</sub>	<u>32.2</u> <sub>±2.8</sub>	<u>28.9</u> <sub>±3.2</sub>
GPT-4o-mini	2023-10	<b>32.8</b> <sub>±2.6</sub>	<u>25.9</u> <sub>±3.2</sub>	<u>33.2</u> <sub>±0.7</sub>	<u>29.7</u> <sub>±3.8</sub>

Table 3: F1 (↑) scores of second-level relation forecasting on different test splits, using “Single Function”. The best-performing score is highlighted in **bold** and the second-best is underlined. More results in Appendix D.6.

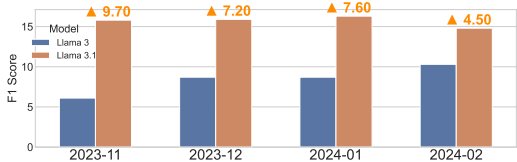


Figure 6: Llama-3.1 vs Llama-3 over four test splits. Llama-3.1 is much higher on 2023-11 (before its training data cutoff date). The gap is smaller afterwards. This shows we need to evaluate model using time-split after its cutoff date.

We compare model forecasting performance across multiple test splits in Table 3. MIRAI’s dynamic data construction pipeline enables the creation of data-contamination-free test sets for newly released models, allowing us to study the effect of a model’s knowledge cutoff time on forecasting performance. We focus on comparing two open-sourced Llama3 models (Dubey et al., 2024) with different cutoff dates but similar training processes. Figure 6 illustrates their performance gap across four test splits. Notably, Llama-3.1 shows the largest lead over Llama-3 in the 2023-11 split, which is after Llama-3’s cutoff but before Llama-3.1’s. This gap diminishes in subsequent months, especially by 2024-02.

This observation indicates potential data contamination favoring the model with more recent training data (especially if eval date is before knowledge-cutoff date). It highlights the necessity of evaluating models using test splits that are later all compared models’ cutoffs to ensure a rigorous forecasting task. Our benchmark uniquely supports this requirement, maintaining the integrity of forecasting assessments by reflecting true forecasting capabilities rather than data exposure advantages.



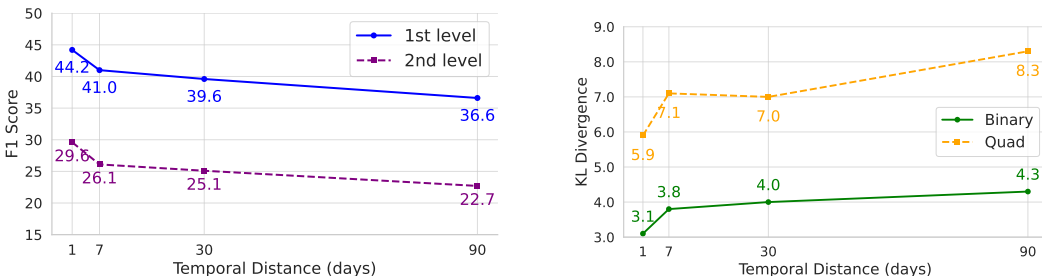


Figure 7: Evaluation of LLM Agents in different temporal distances of the forecasting event.

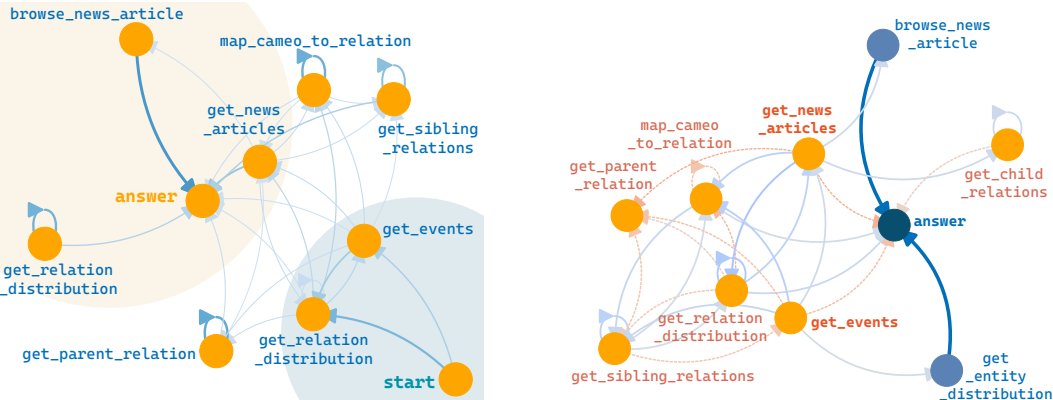


Figure 8: **Action order analysis in LLM agents.** a) Tool-Use Transition Graph of called API functions. Edge thickness indicates transition frequency. b) Freq.(correct) - Freq.(incorrect), in which blue (red) edges indicate positive (negative) contributions.

### 3.4 ANALYZING AGENT BEHAVIOURS

**Impact of temporal distance of the forecasting target.** Our defined event forecasting task varies by temporal distance  $l$ , which specifies how far into the future we want to predict. We thus conduct an ablation study with  $l$  set to 1, 7, 30, and 90 days. Specifically, we fix the query event date and limit the accessible data to  $l$  days prior to the query event date. The experimental results depicted in Figure 7 reveal a clear trend: as the temporal distance increases, the F1 score decreases and KL-divergence increases. This indicates that the agent’s ability to provide accurate predictions diminishes for events further in the future. When the temporal distance is small, such as 1 or 7 days, the agent has access to more recent and relevant information, providing a strong signal (e.g., human experts’ analysis) for making accurate predictions. Thus, to comprehensively benchmark the forecasting capabilities of LLM agents, we should focus on long-term predictions such as those spanning 30 or 90 days. These longer durations require the agents to capture and anticipate potential trend shifts, which may be influenced by a broader range of factors and more complex dependencies.

**Forecasting accuracy on different relation types.** We further split the test events into distinct quadratic relation classes and compute the F1 score for each class, as illustrated in Figure 4b. The results show that all models exhibit significantly higher performance for “verbal cooperation” and “material conflict”, while lower in the other two categories. Several factors contribute to these observations. First, “verbal cooperation” events are more prevalent in the dataset, allowing the model to retrieve more such historical events. Second, “material conflicts” has a consistent pattern of extended duration within the same set of countries. Conversely, events categorized under “material cooperation” and “verbal conflicts”, such as “057: Sign formal agreement” and “084: Return or release”, tend to be more abrupt and unpredictable, demanding subtle trend analysis and contextual knowledge, leading to lower performance in these categories. These observations highlight the need for LLMs capable of understanding the nuances and complexities of different event types.

**How tool-use ordering influences forecasting.** We further investigate the impact of action order on the agent in “Single Function” mode. Figure 8 shows the transition graph from the initial query to the correct final answer, with thicker edges indicating more frequent transitions. Typically, the agent begins with `get_relation_distribution` or `get_event` to gather an initial set of recent and frequent events for key information, and often concludes with `browse_news_article` and `get_news_articles`, which retrieve news content to make accurate forecasts. To further understand whether each function helps forecasting, we subtract the frequency of incorrect predictions from those of correct predictions, as shown in Figure 8b. Here, blue edges represent sequences typically leading to accurate outcomes, and red edges indicate error-prone paths. Actions

Table 4: **Comparison of MIRAI with other temporal reasoning benchmarks.** “Method” refers to the methodology of original and recent models evaluated on the benchmark. Column “Underst.” stands for Understanding. ICL stands for LLM w/ In-Context Learning. FT stands for fine-tuning.

Benchmark	Temporal Reasoning		Task Format	Information			Method	
	Underst.	Forecast		Time Series	KG	Textual		API
TempQuestions (Jia et al., 2018)	✓		KGQA		✓		KGQA Systems	
CRONQuestions (Saxena et al., 2021)	✓		KGQA		✓		Bert-based FT	
TempoQR (Mavromatis et al., 2021)	✓		KGQA		✓		Bert-based FT	
TORQUE (Ning et al., 2020a)	✓		QA		✓		Bert-based FT	
TempReason (Tan et al., 2023a)	✓		QA		✓	✓	ICL + Task FT	
TRAM (Wang & Zhao, 2024)	✓		MCQ		✓		ICL + Bert-based FT	
TCELongBench (Zhang et al., 2024)	✓	✓	QA/MCQ		✓		ICL	
ForecastQA (Jin et al., 2021b)		✓	MCQ		✓		Bert-based FT	
IntervalQA (Zou et al., 2022)		✓	QA/MCQ	✓	✓		ICL	
Approach (Halawi et al., 2024)		✓	QA		✓		ICL	
ForecastBench (Karger et al., 2024)		✓	QA		✓		ICL	
GDELT (Leetaru & Schrodt, 2013b)		✓	Link Prediction		✓		Graph FT + ICL + Task FT	
ICEWS (Boschee et al., 2015b)		✓	Link Prediction		✓		Graph FT + ICL + Task FT	
MIRAI		✓	Relation List		✓	✓	✓	LLM Agent

like `browse_news_article` and `get_entity_distribution` typically result in correct answers more frequently. Notably, `get_news_articles` has a direct red link to the answer, suggesting that this function leads more often to incorrect answers because it only returns news titles, which are too vague for accurate prediction. However, when followed by `browse_news_article` and then providing the answer, the agent is more likely to produce correct outcomes. Similar patterns are observed with `get_event`, where adding `get_entity_distribution` turns a negative link to a positive one to the answer. Figure 5a further demonstrates how each function contributes to the final performance, showing `get_{child/sibling}_relation` are more useful for first-level prediction. These results emphasize the importance for strategic action planning in LLM agents for effective temporal forecasting.

#### 4 RELATED WORK

Recent benchmarks for evaluating temporal reasoning in AI systems can be categorized into temporal understanding and temporal forecasting. Temporal understanding benchmarks (Jia et al., 2018; Saxena et al., 2021; Mavromatis et al., 2021; Tan et al., 2023a; Wang & Zhao, 2024) assess models’ ability to comprehend temporal relations in available data. In contrast, temporal forecasting benchmarks, including our proposed MIRAI, focus on predicting future events based on historical data.

Existing forecasting benchmarks primarily use either QA (Jin et al., 2021b; Zou et al., 2022; Zhang et al., 2024; Halawi et al., 2024; Schoenegger et al., 2024) or link prediction (Boschee et al., 2015b; Leetaru & Schrodt, 2013b) task formulations. QA-based benchmarks typically use textual data, while link prediction tasks often rely on TKGs. MIRAI distinguishes itself by incorporating diverse information sources and employing a multi-relation prediction task format. Additionally, MIRAI introduces an agent-based methodology with intermediate reasoning steps and a construction pipeline that supports dynamic data updates. These features, summarized in Table 4, position MIRAI as a comprehensive and unique benchmark for evaluating temporal forecasting capabilities. We provide further discussions in Appendix C.

#### 5 CONCLUSION AND LIMITATION

In conclusion, we introduce MIRAI, a novel benchmark for evaluating LLM agents in temporal forecasting of international events. Our key contributions include: 1) An agentic environment with APIs supporting comprehensive evaluation of agents’ capabilities with diverse information sourcing, code-based tool use, and forecasting reasoning. 2) A dynamic data construction pipeline enabling monthly updates for contamination-free test splits for evaluating new models. 3) Extensive benchmarking across various agent methods, prediction horizons, and test-time splits, with in-depth analysis of factors influencing agent behavior. Our results reveal the challenges LLM agents face in generating contextually and syntactically correct code and performing complex temporal reasoning. By providing a standardized, dynamic, and comprehensive evaluation benchmark, MIRAI aims to contribute to the development of more accurate and trustworthy models for event forecasting, ultimately supporting more informed decision-making in international relations.

While our work addresses several key challenges, limitations remain, including the need for broader model coverage, expanded API functionality, and more diverse data types. A full discussion of limitations can be found in Appendix B.

## REFERENCES

- 540  
541 AutoGPT Documentation. 1  
542
- 543 GPT-3.5-Turbo, <https://platform.openai.com/docs/models/gpt-3-5-turbo>,  
544 2023. 6, 7  
545
- 546 Gpt-4o contributions. 2024. 7  
547
- 548 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
549 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.  
550 *arXiv preprint arXiv:2303.08774*, 2023. 1, 24
- 551 Anthropic. Claude: An ai assistant by anthropic, 2023. 1  
552
- 553 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,  
554 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language  
555 models. *arXiv preprint arXiv:2108.07732*, 2021. 24
- 556 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,  
557 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. 24  
558
- 559 Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana  
560 Yakhnenko. Translating embeddings for modeling multi-relational data. In C.J.  
561 Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger (eds.), *Ad-*  
562 *vances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc.,  
563 2013. URL [https://proceedings.neurips.cc/paper\\_files/paper/2013/](https://proceedings.neurips.cc/paper_files/paper/2013/file/1cecc7a77928ca8133fa24680a88d2f9-Paper.pdf)  
564 [file/1cecc7a77928ca8133fa24680a88d2f9-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2013/file/1cecc7a77928ca8133fa24680a88d2f9-Paper.pdf). 30
- 565 Elizabeth Boschee, Jennifer Lautenschlager, Sean O’Brien, Steve Shellman, James Starz, and Michael  
566 Ward. *Cameo.cdb.09b5.pdf*. In *ICEWS Coded Event Data*. Harvard Dataverse, 2015a. 3  
567
- 568 Elizabeth Boschee, Jennifer Lautenschlager, Sean O’Brien, Steve Shellman, James Starz, and Michael  
569 Ward. *Icews coded event data*, 2015b. 10, 23  
570
- 571 Thomas Brown and Susan Lee. Predictive analytics in economic sanctions and international policy.  
572 *Journal of International Economics*, 26(4):311–330, 2018. 1
- 573 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
574 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
575 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 24  
576
- 577 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared  
578 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large  
579 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021. 24
- 580 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam  
581 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:  
582 Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113,  
583 2023. 24  
584
- 585 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
586 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve  
587 math word problems. *arXiv preprint arXiv:2110.14168*, 2021. 24
- 588 Richard Davis and Anh Nguyen. Strategic alliances and predictive diplomacy: A review of historical  
589 data. *Political Science Quarterly*, 132(1):45–72, 2017. 1  
590
- 591 Google DeepMind. Gemini: An ai model by google deepmind, 2023. 1  
592
- 593 Songgaojun Deng, Huzefa Rangwala, and Yue Ning. Dynamic knowledge graph based multi-event  
forecasting. In *KDD*, pp. 1585–1595. ACM, 2020a. 23, 29

- 594 Songgaojun Deng, Huzefa Rangwala, and Yue Ning. Dynamic Knowledge Graph based Multi-Event  
595 Forecasting. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge  
596 Discovery & Data Mining*, pp. 1585–1595, Virtual Event CA USA, August 2020b. ACM. ISBN  
597 978-1-4503-7998-4. doi: 10.1145/3394486.3403209. URL [https://dl.acm.org/doi/10.  
598 1145/3394486.3403209](https://dl.acm.org/doi/10.1145/3394486.3403209). 29
- 599  
600 Songgaojun Deng, Huzefa Rangwala, and Yue Ning. Understanding event predictions via contextual-  
601 ized multilevel feature learning. In *CIKM*, pp. 342–351. ACM, 2021. 23, 29
- 602 Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and  
603 Yu Su. Mind2web: Towards a generalist agent for the web. In *Proceedings of NeurIPS*, 2023. 24
- 604  
605 Zifeng Ding, Zhen Han, Yunpu Ma, and Volker Tresp. Temporal knowledge graph forecasting with  
606 neural ode. abs/2101.05151, 2021. 23
- 607  
608 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
609 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn,  
610 Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston  
611 Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron,  
612 Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris  
613 McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton  
614 Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David  
615 Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes,  
616 Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip  
617 Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme  
618 Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu,  
619 Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov,  
620 Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah,  
621 Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu  
622 Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph  
623 Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani,  
624 Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz  
625 Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence  
626 Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas  
627 Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri,  
628 Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis,  
629 Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov,  
630 Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan  
631 Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan,  
632 Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy,  
633 Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit  
634 Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou,  
635 Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia  
636 Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rparathy, Sheng Shen, Shengye Wan,  
637 Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla,  
638 Stephane Collet, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek  
639 Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao,  
640 Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent  
641 Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu,  
642 Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia,  
643 Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen  
644 Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe  
645 Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya  
646 Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex  
647 Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei  
Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew  
Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley  
Gabriel, Ashwin Barambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin  
Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu,  
Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt



- 648 Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao  
649 Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon  
650 Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide  
651 Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le,  
652 Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily  
653 Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix  
654 Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank  
655 Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern,  
656 Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid  
657 Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen  
658 Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-  
659 Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste  
660 Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul,  
661 Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie,  
662 Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik  
663 Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly  
664 Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhota, Kyle Huang, Lailin Chen,  
665 Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu,  
666 Liron Moshkovich, Luca Wehrstedt, Madian Khabza, Manav Avalani, Manish Bhatt, Maria  
667 Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev,  
668 Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle  
669 Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang,  
670 Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam,  
671 Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier,  
672 Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia  
673 Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro  
674 Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani,  
675 Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy,  
676 Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan  
677 Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara  
678 Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh  
679 Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha,  
680 Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe,  
681 Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan  
682 Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury,  
683 Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe  
684 Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi,  
685 Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vitor Albiero, Vlad Ionescu,  
686 Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen  
687 Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojuan Wu, Xiaolan  
688 Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin  
689 Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He,  
690 Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The  
691 llama 3 herd of models. *CoRR*, abs/2407.21783, 2024. doi: 10.48550/ARXIV.2407.21783. URL  
692 <https://doi.org/10.48550/arXiv.2407.21783>. 7, 8
- 693 Luis Galárraga, Simon Razniewski, Antoine Amarilli, and Fabian M. Suchanek. Predicting complete-  
694 ness in knowledge bases. In Maarten de Rijke, Milad Shokouhi, Andrew Tomkins, and Min Zhang  
695 (eds.), *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*,  
696 *WSDM 2017, Cambridge, United Kingdom, February 6-10, 2017*, pp. 375–383. ACM, 2017. 1
- 697 Julia Gastinger, Christian Meilicke, Federico Errica, Timo Szttyler, Anett Schuelke, and Heiner  
698 Stuckenschmidt. History repeats Itself: A Baseline for Temporal Knowledge Graph Forecasting,  
699 April 2024. URL <http://arxiv.org/abs/2404.16726>. 30
- 700 Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach,  
701 Hal Daumé III, and Kate Crawford. Datasheets for Datasets, December 2021. URL <http://arxiv.org/abs/1803.09010>. arXiv:1803.09010 [cs]. 6

- 702 Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth  
703 Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. Textbooks are all  
704 you need. *arXiv preprint arXiv:2306.11644*, 2023. 24
- 705  
706 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao  
707 Bi, Y Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming—the  
708 rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024. 24
- 709 Danny Halawi, Fred Zhang, Chen Yueh-Han, and Jacob Steinhardt. Approaching Human-Level  
710 Forecasting with Language Models, February 2024. URL [http://arxiv.org/abs/2402.](http://arxiv.org/abs/2402.18563)  
711 [18563](http://arxiv.org/abs/2402.18563). arXiv:2402.18563 [cs]. 10, 23
- 712 Zhen Han, Peng Chen, Yunpu Ma, and Volker Tresp. xERTE: Explainable Reasoning on Temporal  
713 Knowledge Graphs for Forecasting Future Links, April 2021. ICLR 2021. 30
- 714  
715 Shibo Hao, Tianyang Liu, Zhen Wang, and Zhiting Hu. Toolkengpt: Augmenting frozen  
716 language models with massive tools via tool embeddings. In Alice Oh, Tristan Nau-  
717 mann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances*  
718 *in Neural Information Processing Systems 36: Annual Conference on Neural Informa-*  
719 *tion Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,*  
720 *2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/](http://papers.nips.cc/paper_files/paper/2023/hash/8fd1a81c882cd45f64958da6284f4a3f-Abstract-Conference.html)  
721 [8fd1a81c882cd45f64958da6284f4a3f-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/8fd1a81c882cd45f64958da6284f4a3f-Abstract-Conference.html). 24
- 722 Ziniu Hu, Ahmet Iscen, Chen Sun, Kai-Wei Chang, Yizhou Sun, David Ross, Cordelia Schmid, and  
723 Alireza Fathi. AVIS: autonomous visual information seeking with large language model agent.  
724 In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine  
725 (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural*  
726 *Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 -*  
727 *16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/](http://papers.nips.cc/paper_files/paper/2023/hash/029df12a9363313c3e41047844ecad94-Abstract-Conference.html)  
728 [029df12a9363313c3e41047844ecad94-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/029df12a9363313c3e41047844ecad94-Abstract-Conference.html). 24
- 729 Zijie Huang, Daheng Wang, Binxuan Huang, Chenwei Zhang, Jingbo Shang, Yan Liang, Zhengyang  
730 Wang, Xian Li, Christos Faloutsos, Yizhou Sun, and Wei Wang. Concept2Box: Joint geomet-  
731 ric embeddings for learning two-view knowledge graphs. In *Findings of the Association for*  
732 *Computational Linguistics: ACL 2023*, pp. 10105–10118, 2023. 1
- 733 Zijie Huang, Jeehyun Hwang, Junkai Zhang, Jinwoo Baik, Weitong Zhang, Dominik Wodarz, Yizhou  
734 Sun, Quanquan Gu, and Wei Wang. Causal graph ODE: continuous treatment effect modeling  
735 in multi-agent dynamical systems. In Tat-Seng Chua, Chong-Wah Ngo, Ravi Kumar, Hady W.  
736 Lauw, and Roy Ka-Wei Lee (eds.), *Proceedings of the ACM on Web Conference 2024, WWW 2024,*  
737 *Singapore, May 13-17, 2024*, pp. 4607–4617. ACM, 2024. 1
- 738  
739 Zhen Jia, Abdalghani Abujabal, Rishiraj Saha Roy, Jannik Strötgen, and Gerhard Weikum. TempQues-  
740 tions: A Benchmark for Temporal Question Answering. In *Companion Proceedings of the The*  
741 *Web Conference 2018, WWW '18*, pp. 1057–1062, Republic and Canton of Geneva, CHE, April  
742 2018. International World Wide Web Conferences Steering Committee. ISBN 978-1-4503-5640-4.  
743 WWW 2018. 10, 23
- 744 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
745 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,  
746 Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas  
747 Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. 7
- 748 Woojeong Jin, Meng Qu, Xisen Jin, and Xiang Ren. Recurrent event network: Autoregressive  
749 structure inference over temporal knowledge graphs. In *EMNLP (1)*, pp. 6669–6683. Association  
750 for Computational Linguistics, 2020a. 23
- 751 Woojeong Jin, Meng Qu, Xisen Jin, and Xiang Ren. Recurrent Event Network: Autoregressive  
752 Structure Inference over Temporal Knowledge Graphs, October 2020b. EMNLP 2020. 1
- 753  
754 Woojeong Jin, Rahul Khanna, Suji Kim, Dong-Ho Lee, Fred Morstatter, Aram Galstyan, and Xiang  
755 Ren. Forecastqa: A question answering challenge for event forecasting with temporal text data. In  
*ACL/IJCNLP (1)*, pp. 4636–4650. Association for Computational Linguistics, 2021a. 7, 23, 26

- 756 Woojeong Jin, Rahul Khanna, Suji Kim, Dong-Ho Lee, Fred Morstatter, Aram Galstyan, and Xiang  
757 Ren. ForecastQA: A Question Answering Challenge for Event Forecasting with Temporal Text  
758 Data. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics  
759 and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long  
760 Papers)*, pp. 4636–4650, Online, August 2021b. Association for Computational Linguistics. ACL  
761 2021. 10, 23
- 762 Emily Johnson and Mark Roberts. The role of diplomacy in shaping foreign policy. *Diplomatic  
763 Review*, 12(2):145–170, 2019. 1
- 764  
765 Ezra Karger, Houtan Bastani, Chen Yueh-Han, Zachary Jacobs, Danny Halawi, Fred Zhang, and  
766 Philip E. Tetlock. ForecastBench: A Dynamic Benchmark of AI Forecasting Capabilities, September  
767 2024. URL <http://arxiv.org/abs/2409.19839>. arXiv:2409.19839 [cs]. 10, 23
- 768  
769 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large  
770 language models are zero-shot reasoners. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle  
771 Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems 35:  
772 Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans,  
773 LA, USA, November 28 - December 9, 2022*, 2022. 6
- 774  
775 Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov,  
776 Thomas Wang, Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela, Matthieu Cord, and  
777 Victor Sanh. OBELICS: an open web-scale filtered dataset of interleaved image-text documents.  
778 In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine  
779 (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural  
780 Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,  
781 2023*, 2023. 5, 62
- 782  
783 Dong-Ho Lee, Kian Ahrabian, Woojeong Jin, Fred Morstatter, and Jay Pujara. Temporal knowl-  
784 edge graph forecasting without knowledge using in-context learning. In *EMNLP*, pp. 544–557.  
785 Association for Computational Linguistics, 2023. 23
- 786  
787 Kalev Leetaru and Philip A Schrod. Gdelt: Global data on events, location, and tone, 1979–2012. In  
788 *ISA annual convention*, volume 2, pp. 1–49. Citeseer, 2013a. 23
- 789  
790 Kalev Leetaru and Philip A Schrod. GDELT: Global Data on Events, Location and Tone,. 2013b. 2,  
791 10, 23
- 792  
793 Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang,  
794 and Yongbin Li. Api-bank: A comprehensive benchmark for tool-augmented llms. In *Proceedings  
795 of EMNLP*, 2023a. 24
- 796  
797 Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang,  
798 and Yongbin Li. API-Bank: A Comprehensive Benchmark for Tool-Augmented LLMs, October  
799 2023b. arXiv:2304.08244 [cs]. 1
- 800  
801 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou,  
802 Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source be with  
803 you! *arXiv preprint arXiv:2305.06161*, 2023c. 24
- 804  
805 Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom  
806 Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation  
807 with alphacode. *Science*, 378(6624):1092–1097, 2022. 24
- 808  
809 Zixuan Li, Xiaolong Jin, Saiping Guan, Wei Li, Jiafeng Guo, Yuanzhuo Wang, and Xueqi Cheng.  
Search from history and reason for future: Two-stage reasoning on temporal knowledge graphs. In  
*ACL*, 2021a. 23
- Zixuan Li, Xiaolong Jin, Wei Li, Saiping Guan, Jiafeng Guo, Huawei Shen, Yuanzhuo Wang, and  
Xueqi Cheng. Temporal knowledge graph reasoning based on evolutionary representation learning.  
In *SIGIR*, pp. 408–417. ACM, 2021b. 7, 23, 26, 30

- 810 Zixuan Li, Xiaolong Jin, Wei Li, Saiping Guan, Jiafeng Guo, Huawei Shen, Yuanzhuo Wang, and  
811 Xueqi Cheng. Temporal Knowledge Graph Reasoning Based on Evolutional Representation  
812 Learning, April 2021c. SIGIR 2021. [1](#)
- 813
- 814 Ruotong Liao, Xu Jia, Yunpu Ma, and Volker Tresp. Gentkg: Generative forecasting on temporal  
815 knowledge graph. *CoRR*, abs/2310.07793, 2023. [24](#)
- 816 Bo Liu, Yuqian Jiang, Xiaohan Zhang, Qiang Liu, Shiqi Zhang, Joydeep Biswas, and Peter Stone.  
817 LLM+P: Empowering Large Language Models with Optimal Planning Proficiency, September  
818 2023a. arXiv:2304.11477 [cs]. [2](#)
- 819
- 820 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding,  
821 Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui  
822 Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang.  
823 AgentBench: Evaluating LLMs as Agents, October 2023b. arXiv:2308.03688 [cs]. [1](#)
- 824 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding,  
825 Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. In *Proceedings of ICLR*,  
826 2024. [24](#)
- 827
- 828 Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu,  
829 and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language models.  
830 In *Proceedings of NeurIPS*, 2023a. [24](#)
- 831 Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu,  
832 and Jianfeng Gao. Chameleon: Plug-and-Play Compositional Reasoning with Large Language  
833 Models, October 2023b. arXiv:2304.09842 [cs]. [1](#)
- 834
- 835 Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing  
836 Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with  
837 evol-instruct. *arXiv preprint arXiv:2306.08568*, 2023. [24](#)
- 838 Yunshan Ma, Chenchen Ye, Zijian Wu, Xiang Wang, Yixin Cao, and Tat-Seng Chua. Context-aware  
839 event forecasting via graph disentanglement. In *KDD*, pp. 1643–1652. ACM, 2023a. [23](#)
- 840 Yunshan Ma, Chenchen Ye, Zijian Wu, Xiang Wang, Yixin Cao, and Tat-Seng Chua. Context-  
841 aware Event Forecasting via Graph Disentanglement. In *Proceedings of the 29th ACM SIGKDD*  
842 *Conference on Knowledge Discovery and Data Mining*, pp. 1643–1652, August 2023b. KDD 2023.  
843 [2](#)
- 844
- 845 Yunshan Ma, Chenchen Ye, Zijian Wu, Xiang Wang, Yixin Cao, Liang Pang, and Tat-Seng Chua.  
846 Structured, complex and time-complete temporal event forecasting. *CoRR*, abs/2312.01052, 2023c.  
847 [23](#)
- 848 Farzaneh Mahdisoltani, Joanna Biega, and Fabian M Suchanek. YAGO3: A Knowledge Base from  
849 Multilingual Wikipedias. January 2015. CIDR 2015. [1](#)
- 850
- 851 Costas Mavromatis, Prasanna Lakkur Subramanyam, Vassilis N. Ioannidis, Soji Adeshina, Phillip R.  
852 Howard, Tetiana Grinberg, Nagib Hakim, and George Karypis. TempoQR: Temporal Question  
853 Reasoning over Knowledge Graphs. arXiv, December 2021. AAAI 2022. [10](#), [23](#)
- 854 Dr Mclean, Alan Patterson, and John Williams. Risk assessment, policy-making and the limits of  
855 knowledge: The precautionary principle and international relations. *International Relations*, 23:  
856 548–566, 12 2009. [1](#)
- 857
- 858 Grégoire Mialon, Roberto Dessì, Maria Lomeli, Christoforos Nalmpantis, Ram Pasunuru, Roberta  
859 Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, Edouard Grave,  
860 Yann LeCun, and Thomas Scialom. Augmented language models: a survey. In *arXiv preprint*  
861 *arXiv:2302.07842*, 2023. [24](#)
- 862 Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese,  
863 and Caiming Xiong. Codegen: An open large language model for code with multi-turn program  
synthesis. *arXiv preprint arXiv:2203.13474*, 2022. [24](#)



- 864 Qiang Ning, Hao Wu, Rujun Han, Nanyun Peng, Matt Gardner, and Dan Roth. TORQUE: A  
865 Reading Comprehension Dataset of Temporal Ordering Questions. In *Proceedings of the 2020*  
866 *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1158–1172,  
867 Online, November 2020a. Association for Computational Linguistics. EMNLP 2020. 10, 23
- 868 Qiang Ning, Hao Wu, Rujun Han, Nanyun Peng, Matt Gardner, and Dan Roth. TORQUE: A reading  
869 comprehension dataset of temporal ordering questions. In *EMNLP*, pp. 1158–1172, 2020b. 23
- 870 Aaron Parisi, Yao Zhao, and Noah Fiedel. Talm: Tool augmented language models. In *arXiv preprint*  
871 *arXiv:2205.12255*, 2022. 24
- 872 Namyong Park, Fuchen Liu, Purvanshi Mehta, Dana Cristofor, Christos Faloutsos, and Yuxiao Dong.  
873 Evokg: Jointly modeling event time and network structure for reasoning over temporal knowledge  
874 graphs. In *WSDM*, pp. 794–803. ACM, 2022. 23
- 875 Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are NLP models really able to solve simple  
876 math word problems? In *Proceedings of NAACL*, 2021. 24
- 877 Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong,  
878 Xiangru Tang, Bill Qian, Sihan Zhao, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein,  
879 Dahai Li, Zhiyuan Liu, and Maosong Sun. Toolllm: Facilitating large language models to master  
880 16000+ real-world apis. *CoRR*, abs/2307.16789, 2023. doi: 10.48550/ARXIV.2307.16789. URL  
881 <https://doi.org/10.48550/arXiv.2307.16789>. 24
- 882 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language  
883 understanding by generative pre-training. 2018. 24
- 884 Revanth Gangi Reddy, Yi R. Fung, Qi Zeng, Manling Li, Ziqi Wang, Paul Sullivan, and Heng Ji.  
885 SmartBook: AI-Assisted Situation Report Generation, March 2023. arXiv. 1
- 886 Subhro Roy and Dan Roth. Solving general arithmetic word problems. In *Proceedings of EMNLP*,  
887 2015. 24
- 888 Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi  
889 Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code.  
890 *arXiv preprint arXiv:2308.12950*, 2023. 24
- 891 Apoorv Saxena, Soumen Chakrabarti, and Partha Talukdar. Question Answering Over Temporal  
892 Knowledge Graphs. In *Proceedings of the 59th Annual Meeting of the Association for Computa-*  
893 *tional Linguistics and the 11th International Joint Conference on Natural Language Processing*  
894 *(Volume 1: Long Papers)*, pp. 6663–6676, Online, August 2021. Association for Computational  
895 Linguistics. ACL 2021. 10, 23
- 896 Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke  
897 Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach  
898 themselves to use tools. In *Proceedings of NeurIPS*, 2023. 24
- 899 Philipp Schoenegger, Indre Tuminauskaitė, Peter S. Park, and Philip E. Tetlock. Wisdom of the  
900 Silicon Crowd: LLM Ensemble Prediction Capabilities Rival Human Crowd Accuracy, May 2024.  
901 URL <http://arxiv.org/abs/2402.19379>. arXiv:2402.19379 [cs]. 10, 23
- 902 Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hug-  
903 gingGPT: Solving AI Tasks with ChatGPT and its Friends in Hugging Face, December 2023.  
904 arXiv:2303.17580 [cs]. 1
- 905 Xiaoming Shi, Siqiao Xue, Kangrui Wang, Fan Zhou, James Y. Zhang, JUN ZHOU, Chenhao  
906 Tan, and Hongyuan Mei. Language models can improve event prediction by few-shot abductive  
907 reasoning. In *NeurIPS*, 2023. 24
- 908 Johnathan Smith and Jane Doe. Geopolitical risk assessment in international relations. *Journal of*  
909 *Global Politics*, 15(3):200–225, 2020. 1
- 910 Theodore R. Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L. Griffiths. Cognitive Archi-  
911 tectures for Language Agents, March 2024. arXiv:2309.02427 [cs]. 1

- 918 Haohai Sun, Jialu Zhong, Yunpu Ma, Zhen Han, and Kun He. Timetraveler: Reinforcement learning  
919 for temporal knowledge graph forecasting. In *EMNLP*, 2021. 23
- 920
- 921 Qingyu Tan, Hwee Tou Ng, and Lidong Bing. Towards Benchmarking and Improving the Temporal  
922 Reasoning Capability of Large Language Models. In *Proceedings of the 61st Annual Meeting  
923 of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14820–14835,  
924 Toronto, Canada, July 2023a. Association for Computational Linguistics. ACL 2023. 10, 23
- 925 Qingyu Tan, Hwee Tou Ng, and Lidong Bing. Towards benchmarking and improving the tempo-  
926 ral reasoning capability of large language models. In *ACL*, pp. 14820–14835. Association for  
927 Computational Linguistics, 2023b. 23
- 928 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
929 Lacroix, Baptiste Roziere, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
930 efficient foundation language models, 2023. 1, 24
- 931
- 932 Rakshit Trivedi, Hanjun Dai, Yichen Wang, and Le Song. Know-evolve: deep temporal reasoning for  
933 dynamic knowledge graphs. In *ICML*, pp. 3462–3471, 2017. 23
- 934 Karthik Valmeekam, Matthew Marquez, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambham-  
935 pati. PlanBench: An Extensible Benchmark for Evaluating Large Language Models on Planning  
936 and Reasoning about Change, November 2023. arXiv:2206.10498 [cs]. 2
- 937
- 938 Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlikar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and  
939 Anima Anandkumar. Voyager: An Open-Ended Embodied Agent with Large Language Models,  
940 October 2023. arXiv:2305.16291 [cs]. 1
- 941 Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. Codet5: Identifier-aware unified pre-  
942 trained encoder-decoder models for code understanding and generation. In *Proceedings of the  
943 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 8696–8708, 2021.  
944 24
- 945 Yuqing Wang and Yun Zhao. Tram: Benchmarking temporal reasoning for large language models.  
946 2023. 23
- 947
- 948 Yuqing Wang and Yun Zhao. TRAM: Benchmarking Temporal Reasoning for Large Language  
949 Models, May 2024. arXiv:2310.00835 [cs]. 10, 23
- 950 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V.  
951 Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In  
952 *NeurIPS*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/  
953 9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html). 24
- 954
- 955 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc  
956 Le, and Denny Zhou. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models,  
957 January 2023. arXiv:2201.11903 [cs]. 6
- 958 Lilian Weng. LLM Powered Autonomous Agents, June 2023. Section: posts. 1
- 959
- 960 Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and  
961 Yu Su. TravelPlanner: A Benchmark for Real-World Planning with Language Agents, February  
962 2024. arXiv:2402.01622 [cs]. 24
- 963 Qiantong Xu, Fenglu Hong, Bo Li, Changran Hu, Zhengyu Chen, and Jian Zhang. On the tool  
964 manipulation capability of open-source large language models. *arXiv preprint arXiv:2305.16504*,  
965 2023a. 24
- 966
- 967 Wenjie Xu, Ben Liu, Miao Peng, Xu Jia, and Min Peng. Pre-trained language model with prompts  
968 for temporal knowledge graph completion. In *ACL (Findings)*, pp. 7790–7803. Association for  
969 Computational Linguistics, 2023b. 23
- 970 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.  
971 ReAct: Synergizing Reasoning and Acting in Language Models, March 2023a. arXiv:2210.03629  
[cs]. 4

- 972 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao.  
973 React: Synergizing reasoning and acting in language models. In *The Eleventh International Confer-*  
974 *ence on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net,  
975 2023b. [2](#), [6](#)
- 976 Siyu Yuan, Jiangjie Chen, Ziquan Fu, Xuyang Ge, Soham Shah, Charles Robert Jankowski, Yanghua  
977 Xiao, and Deqing Yang. Distilling Script Knowledge from Large Language Models for Constrained  
978 Language Planning, May 2023. arXiv:2305.05252 [cs]. [1](#)
- 979 Michael Zhang and Eunsol Choi. Situatedqa: Incorporating extra-linguistic contexts into qa. In  
980 *EMNLP*, 2021. [23](#)
- 981 Zhihan Zhang, Yixin Cao, Chenchen Ye, Yunshan Ma, Lizi Liao, and Tat-Seng Chua. Analyzing  
982 Temporal Complex Events with Large Language Models? A Benchmark towards Temporal, Long  
983 Context Understanding, June 2024. arXiv:2406.02472 [cs]. [10](#), [23](#)
- 984 Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth. “going on a vacation” takes longer than  
985 “going for a walk”: A study of temporal commonsense understanding. In *EMNLP*, pp. 3363–3369,  
986 2019. [23](#)
- 987 Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,  
988 Yonatan Bisk, Daniel Fried, Uri Alon, et al. Webarena: A realistic web environment for building  
989 autonomous agents. In *Proceedings of ICLR*, 2024. [24](#)
- 990 Cunchao Zhu, Muhao Chen, Changjun Fan, Guangquan Cheng, and Yan Zhan. Learning from  
991 history: Modeling temporal knowledge graphs with sequential copy-generation networks. In *AAAI*  
992 *Conference on Artificial Intelligence*, 2020. [23](#)
- 993 Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. ToolQA: A dataset for LLM  
994 question answering with external tools. In *Proceedings of NeurIPS*, 2023a. [24](#)
- 995 Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. ToolQA: A Dataset for LLM  
996 Question Answering with External Tools, June 2023b. arXiv:2306.13304 [cs]. [1](#)
- 997 Andy Zou, Tristan Xiao, Ryan Jia, Joe Kwon, Mantas Mazeika, Richard Li, Dawn Song, Jacob  
998 Steinhardt, Owain Evans, and Dan Hendrycks. Forecasting Future World Events with Neural  
999 Networks. arXiv, October 2022. NeurIPS 2022. [1](#), [10](#), [23](#)
- 1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

1026	APPENDIX	
1027		
1028	<b>A Reproducibility Statement</b>	<b>22</b>
1029		
1030	<b>B Limitations</b>	<b>22</b>
1031		
1032	<b>C Additional Related Work</b>	<b>23</b>
1033		
1034	C.1 Temporal Reasoning Benchmarks . . . . .	23
1035	C.2 Temporal Event Forecasting Benchmarks and Methods . . . . .	23
1036	C.3 Evaluation of Language Agents . . . . .	24
1037	C.4 LLMs for Tool-Use . . . . .	24
1038	C.5 LLMs for Code Generation . . . . .	24
1039		
1040	<b>D Additional Experimental Results and Analysis</b>	<b>25</b>
1041		
1042	D.1 Analysis on the Final Status of Different LLM Agents . . . . .	25
1043	D.2 Analysis on Agent ReAct Iteration Sequence Length . . . . .	25
1044	D.3 Evaluation for Traditional TKG and NLP Forecasting Methods . . . . .	25
1045	D.4 Human Forecasting Performance as a Reference . . . . .	27
1046	D.5 Human Evaluation on the Dataset Quality . . . . .	27
1047	D.6 Forecasting Performance on Different Test Splits . . . . .	27
1048	D.7 Analysis with Retrieve-Augmented Generation (RAG) Methods . . . . .	28
1049	D.7.1 Methods and Experimental Setup . . . . .	28
1050	D.7.2 Comparison of Retrieval Strategies with Agent . . . . .	28
1051	D.7.3 Key Findings and Analysis . . . . .	28
1052	D.8 Analysis with Heuristic-based and Traditional TKG methods . . . . .	29
1053	D.8.1 Task Focus: Relation Prediction . . . . .	29
1054	D.8.2 Experimental Setup . . . . .	30
1055	D.8.3 Metric Selection and Justification . . . . .	30
1056	D.8.4 Results Analysis and Implications . . . . .	31
1057	D.9 Evaluation on Impact of LLM Parameter Size . . . . .	32
1058		
1059		
1060	<b>E Additional Information about API</b>	<b>33</b>
1061	E.1 Overview of API Data Classes and Functions . . . . .	33
1062		
1063	<b>F Additional Forecasting Examples of LLM Agent</b>	<b>34</b>
1064	F.1 GPT-4o-based Agent with ReAct and “Code Block” Action . . . . .	34
1065	F.1.1 Query Details . . . . .	34
1066	F.1.2 Query Prompt . . . . .	34
1067	F.1.3 System Prompt . . . . .	34
1068	F.1.4 Agent ReAct Process . . . . .	35
1069	F.1.5 Agent Final Status . . . . .	39
1070	F.2 GPT-4o-based Agent with ReAct and “Single Function” Action . . . . .	39
1071	F.2.1 Query Details . . . . .	39
1072	F.2.2 Query Prompt . . . . .	39
1073	F.2.3 System Prompt . . . . .	39
1074	F.2.4 Agent ReAct Process . . . . .	40
1075	F.2.5 Agent Final Status . . . . .	43
1076	F.3 GPT-4o-mini-based Agent with ReAct and “Code Block” Action . . . . .	44
1077	F.3.1 Query Details . . . . .	44
1078	F.3.2 Query Prompt . . . . .	44
1079		



1080	F.3.3	System Prompt . . . . .	44
1081	F.3.4	Agent ReAct Process . . . . .	45
1082	F.3.5	Agent Final Status . . . . .	47
1083	F.4	GPT-4o-mini-based Agent with ReAct and “Single Function” Action . . . . .	47
1084	F.4.1	Query Details . . . . .	47
1085	F.4.2	Query Prompt . . . . .	47
1086	F.4.3	System Prompt . . . . .	47
1087	F.4.4	Agent ReAct Process . . . . .	48
1088	F.4.5	Agent Final Status . . . . .	50
1089			
1090			
1091	<b>G</b>	<b>Additional Details of the Benchmark</b>	<b>51</b>
1092	G.1	Countries and ISO Codes . . . . .	51
1093	G.2	Relations and CAMEO Codes . . . . .	52
1094	G.3	Human Evaluation on the Dataset Quality . . . . .	55
1095			
1096	<b>H</b>	<b>Additional Details about Agent Setup</b>	<b>55</b>
1097			
1098	<b>I</b>	<b>Prompts</b>	<b>56</b>
1099	I.1	System Prompts . . . . .	56
1100	I.1.1	System Prompt for ReAct Agent with Action Type as “Single Function” . . . . .	56
1101	I.1.2	System Prompt for ReAct Agent with Action Type as “Code Block” . . . . .	57
1102	I.1.3	System Prompt for Direct Agent . . . . .	58
1103	I.1.4	System Prompt for CoT Agent . . . . .	59
1104	I.2	Query Prompt . . . . .	59
1105	I.3	Answer Extraction Prompt . . . . .	60
1106			
1107			
1108	<b>J</b>	<b>Datasheet for MIRAI</b>	<b>61</b>
1109	J.1	Motivation . . . . .	61
1110	J.2	Composition . . . . .	61
1111	J.3	Collection Process . . . . .	62
1112	J.4	Uses . . . . .	62
1113	J.5	Distribution . . . . .	62
1114	J.6	Maintenance . . . . .	63
1115	J.7	Term of Use for GDELT . . . . .	63
1116			
1117			
1118			
1119			
1120			
1121			
1122			
1123			
1124			
1125			
1126			
1127			
1128			
1129			
1130			
1131			
1132			
1133			

## 1134 A REPRODUCIBILITY STATEMENT

1135 We provide detailed information and explanation of our experimental setup, dataset, and evaluation  
1136 metrics, ensuring reproducibility. Specifically, we describe the following:

- 1137 • Experimental setup (Sec. 3): Detailed descriptions of the hardware and software configurations,  
1138 including libraries and tools used.
- 1139 • Dataset (Sec. 2): Comprehensive information on the dataset construction, including sources,  
1140 preprocessing steps, and any techniques applied.
- 1141 • Evaluation Metrics (Sec. 2.3): Clear definitions for the metrics chosen to evaluate performance.

1142 Additionally, we commit to releasing the following resources for the replication of our results:

- 1143 • Codebase: The complete codebase, including scripts for dataset construction, model serving, and  
1144 evaluation. This is currently available on an anonymous repository [here](#).
- 1145 • Dataset: The processed dataset, along with detailed instructions on how to construct the dataset  
1146 [here](#).

1147 By providing these resources, we aim to ensure that our work is fully reproducible and can be  
1148 independently verified by the research community.

## 1150 B LIMITATIONS

1151 Our benchmark addresses several key challenges in evaluating LLM agents for event forecasting,  
1152 including the integration of diverse information sources, the construction of an agentic interactive  
1153 environment, and the contamination-free forecasting data and task formulation through its dynamic  
1154 design. Despite this significant advantage, we acknowledge the following limitations:

- 1155 1. **Model Coverage:** While we have tested representative open-source and closed-source LLMs, our  
1156 experiments do not exhaustively cover all available models. Future work could expand to include  
1157 a wider range of LLMs for more comprehensive evaluation.
- 1158 2. **API Functionality:** The current API, while functional, has room for expansion. Future iterations  
1159 could incorporate more sophisticated analytical tools, such as time series analysis functions, to  
1160 encourage deeper temporal reasoning. Allowing agents to generate and add custom functions  
1161 during their reasoning process could also lead to more diverse problem-solving approaches.
- 1162 3. **Experimental Robustness:** The current experiments, while informative, are limited in scope  
1163 due to cost and time constraints. Increasing the number of experimental rounds and adjusting  
1164 parameters like model temperature could provide more statistically robust results and insights into  
1165 model performance variability.
- 1166 4. **Geopolitical Bias:** Reliance on GDELT as the primary data source may introduce biases in event  
1167 coverage and interpretation, potentially skewing towards Western or English-language media  
1168 perspectives. Future work could explore integrating multiple diverse data sources and languages  
1169 to mitigate this limitation.

1170 Addressing these limitations in future iterations will further enhance the benchmark’s robustness and  
1171 relevance in the rapidly evolving field of AI-driven event forecasting.

1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

## C ADDITIONAL RELATED WORK

### C.1 TEMPORAL REASONING BENCHMARKS

Many benchmarks sensing the temporal reasoning ability of AI models have been constructed, but they have different focuses and settings with MIRAI, particularly in terms of task, information, and method, as shown in Table 4. One line of benchmarks focuses on the *temporal understanding* ability of the model (Jia et al., 2018; Saxena et al., 2021; Mavromatis et al., 2021; Ning et al., 2020a; Tan et al., 2023a; Wang & Zhao, 2024), such as understanding the temporal relations between available facts in knowledge graphs (KGs) or text, either a short piece of text or a document corpus. While the *temporal forecasting* task largely differs from understanding, where the reasoning target is **unseen** in the database for the model, and as such, the model has to not only understand but to reason. For forecasting, there are two main task formulations among previous benchmarks: the QA task format for benchmarks with history information represented in textual format (Zou et al., 2022; Zhang et al., 2024; Jin et al., 2021b; Halawi et al., 2024; Schoenegger et al., 2024; Karger et al., 2024), and graph link prediction task format for temporal knowledge graph (TKG)-based benchmarks (Boschee et al., 2015b; Leetaru & Schrodt, 2013b). However, the unformat of information sources either lacks of support to clearly structural facts or contextual detail for the model to perform advanced reasoning, while in MIRAI, we provide both information sources. Additionally, we provide carefully constructed API with various data classes and functions that access to various part of the data. With the flexibility provided by code generation, the model is exposed to a broader and more flexible range of information. More importantly, MIRAI distinguishes itself by introducing an agentic environment specifically designed to evaluate LLM agents in the forecasting task. This represents a significant departure from previous work, which has not explored or even considered the potential of performing temporal forecasting tasks using LLM agents.

### C.2 TEMPORAL EVENT FORECASTING BENCHMARKS AND METHODS

**Existing Forecasting Benchmarks.** LLMs have been tested for their *temporal understanding* through tasks such as temporal event ordering or storyline comprehension (Tan et al., 2023b; Ning et al., 2020b; Zhou et al., 2019; Zhang & Choi, 2021; Wang & Zhao, 2023). In the context of *temporal forecasting*, LLMs have been evaluated on traditional structured-event-only benchmarks, such as ICEWS (2014) (Boschee et al., 2015b) and GDELT (2018) (Leetaru & Schrodt, 2013a); and also been evaluated on recent text-based temporal forecasting benchmarks, such as IntervalQA (2022) (Zou et al., 2022) and TCELongBench (2022) (Zhang et al., 2024). However, these evaluations typically involve providing LLMs with retrieved-context for in-context learning and then directly answering the forecast question, lacking intermediate reasoning steps and interaction between the LLM and the database crucial for accurate forecasting. Moreover, although experimental results on these benchmarks show that significant challenges and research value remain in forecasting, even when models encounter events before their training cutoff date, these benchmarks still only provide a fixed timeframe for the testing data, earlier than most recent LLMs.

**Traditional TKG and NLP Methods.** Significant research has been conducted in the field of structured event temporal forecasting. Various methods have been proposed, including aggregating temporal and relational information among entities (Jin et al., 2020a; Li et al., 2021b; Park et al., 2022), retrieving relevant historical events (Zhu et al., 2020; Sun et al., 2021; Li et al., 2021a), and modeling the continuous time evolution of events (Trivedi et al., 2017; Ding et al., 2021). Efforts have also been made to incorporate textual event information into Temporal Knowledge Graphs (TKGs). Glean (Deng et al., 2020a) and CMF (Deng et al., 2021) integrate textual embeddings into graph edges, while SeCoGD (Ma et al., 2023a) employs textual topic modeling to separate subgraphs. The MidEast-TE dataset and LoGo model (Ma et al., 2023c) utilize text clustering to construct complex events for forecasting with local and global contexts. However, these methods still perform forecast reasoning solely on graphs using graph-based techniques. Traditional NLP methods form the event forecasting task as MCQ, for example, the method in ForecastQA (Jin et al., 2021a) use text embedding models for retrieving related event news articles and appending them to the forecasting question for a Bert-based classification over answer candidates. Notably, MIRAI contains both structured and textual event data, supporting the test for both traditional TKG and NLP methods. We show more experimental results and illustrations in Appendix D.3.

**LLMs for TKG and NLP Methods.** Recent studies have explored the use of LLMs for temporal event forecasting by transforming the TKG formulation into text sequences and converting missing object prediction into next token prediction (Xu et al., 2023b). GPT-NeoX-ICL (Lee et al., 2023) employs in-context learning of LLMs and constructs prompts as a list of historical events in quadruplet

1242 format. GENTKG (Liao et al., 2023) enhances the selection of historical event inputs using a temporal  
1243 logical rule-based retrieval strategy, while LAMP (Shi et al., 2023) applies LLMs to perform abductive  
1244 reasoning to assist the retrieval process. However, these works only investigate LLMs with in-context  
1245 learning or simple task-specific fine-tuning. In contrast, MIRAI explores forecasting with an LLM  
1246 agent that supports explicit information gathering and reasoning steps, enabling a hybrid approach  
1247 that leverages both text and graph data.

### 1248 C.3 EVALUATION OF LANGUAGE AGENTS

1249 Previous research has investigated the performance of LLM agents in a variety of domains, including  
1250 arithmetic reasoning focused on obtaining correct solutions (Roy & Roth, 2015; Cobbe et al.,  
1251 2021; Patel et al., 2021), proficiency assessment in utilizing tools and reporting results (Li et al.,  
1252 2023a; Xu et al., 2023a; Zhuang et al., 2023a), evaluation of web navigation skills to find specific  
1253 websites (Deng et al., 2023; Zhou et al., 2024; Liu et al., 2024), and planning travel itineraries under  
1254 given constraints (Xie et al., 2024). However, these evaluations do not fully address the challenges  
1255 posed by tasks involving complex international events with diverse information formats and temporal  
1256 attributes. MIRAI presents a unique task in this context, where the agent must navigate and reason  
1257 over the structured events and textual news articles with temporal information. This setup requires the  
1258 agent to effectively handle multilateral relationships and information spanning different time periods.

### 1259 C.4 LLMs FOR TOOL-USE

1260 Large Language Models (LLMs) have demonstrated remarkable language understanding (Radford  
1261 et al., 2018) and reasoning capabilities (Wei et al., 2022). However, they also possess inherent  
1262 limitations, such as their inability to provide up-to-date responses based on external knowledge or to  
1263 perform complex mathematical reasoning. In response to these challenges, recent advancements have  
1264 seen the integration of LLMs with various external tools (Mialon et al., 2023). Notable examples  
1265 include TALM (Parisi et al., 2022) and ToolFormer (Schick et al., 2023), which utilize in-context  
1266 learning to enhance the model’s ability to leverage different tools in tasks like question answering  
1267 and mathematical reasoning. Chameleon (Lu et al., 2023a) employs an LLM as a natural language  
1268 planner to deduce the optimal sequence of tools to be used, subsequently executing these tools to  
1269 generate the final output. AVIS (Hu et al., 2023) employs dynamic tree search to synthesize the most  
1270 effective tool-use sequence. ToolkenGPT (Hao et al., 2023) integrates tool-use operators as special  
1271 tokens and trains the model through sequence-to-sequence training. ToolLLM (Qin et al., 2023)  
1272 introduces an instruction tuning dataset encompassing over 16,000 real-world APIs, significantly  
1273 enhancing the model’s capability to utilize these tools effectively.

### 1274 C.5 LLMs FOR CODE GENERATION

1275 Early advancements in LLMs have greatly impacted the field of code generation. Notable early  
1276 models such as GPT-3 (Brown et al., 2020) and PALM (Chowdhery et al., 2023) have demonstrated  
1277 the potential of LLMs to assist in code completion and generation. Subsequent models have built  
1278 upon the foundation of pre-trained LLMs, further refining capabilities specific to code generation.  
1279 This has led to the development of code-specific LLMs such as Codex (Chen et al., 2021), Code T5  
1280 (Wang et al., 2021), CodeGen (Nijkamp et al., 2022), AlphaCode (Li et al., 2022), StarCoder (Li et al.,  
1281 2023c), WizardCoder (Luo et al., 2023), and phi-1 (Gunasekar et al., 2023). Moreover, code from  
1282 open-source platforms such as GitHub has increasingly been incorporated into the pre-training data  
1283 for recent LLMs (Touvron et al., 2023; Achiam et al., 2023). This integration has led to improved  
1284 performance of LLMs such as Code Llama (Roziere et al., 2023), Code-Qwen (Bai et al., 2023), and  
1285 DeepSeek-Coder (Guo et al., 2024) on popular code-related tasks (Chen et al., 2021; Austin et al.,  
1286 2021).

1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

## D ADDITIONAL EXPERIMENTAL RESULTS AND ANALYSIS

### D.1 ANALYSIS ON THE FINAL STATUS OF DIFFERENT LLM AGENTS

Table 5: Experiment results with different base LLMs on 2024-02 test splits: Average number of ReAct iterations and the number of test cases ending in different final statuses.

Model	Training Data Cutoff Date	Action Type	Avg. Iterations	Final Status			
				Final Answer	Consecutive Invalid Actions	Consecutive Repetitive Actions	Max Iterations Exceeded
Mistral-7B-Instruct-v0.2	2023-12	Single Function	4.53	81	13	5	1
		Code Block	2.92	46	54	0	0
Llama-3.1-8B-Instruct	2023-12	Single Function	8.18	76	9	9	6
		Code Block	3.65	40	56	5	0
GPT-3.5-Turbo	2021-09	Single Function	2.75	99	1	0	0
		Code Block	3.52	87	11	1	1
GPT-4o-mini	2023-10	Single Function	4.65	<b>100</b>	0	0	0
		Code Block	4.04	<b>100</b>	0	0	0

Following the experiment on different base LLMs and analysis of their code execution error in in Sec. 3.2, we further analyze the final status of different agents when they trigger the termination conditions in the ReAct process. The agent-environment interactions were terminated based on four different statuses:

- **Final Answer:** This status is triggered when the model explicitly indicates completion by declaring a "final answer," indicating that it has generated its final forecast.
- **Consecutive Invalid Actions:** Interaction is halted if the model consecutively executes non-executable actions more than three times.
- **Consecutive Repetitive Actions:** A similar threshold of three consecutive repetitive actions prompts termination under this status.
- **Maximum Iterations Exceeded:** Termination occurs if the agent engages in more than 20 rounds of iterations without making a forecast answer.

We observe from the final status: **1) The GPT series demonstrates enhanced capability in concluding interactions with a final answer:** The GPT series frequently concludes interactions with a final answer, indicating its superior ability to comprehend instructions, utilize tools correctly, and perform reasoning to generate the final forecast. In contrast, the smaller open-source models, Mistral-7b-Instruct-v0.2 and Llama-3.1-8B-Instruct, have a round one-tenth of cases and one-half of the cases generating consecutive invalid answers when using single function and code block, respectively, suggesting limited abilities in adhering to data types, functions, Python syntax, and following detailed instructions.

**2) Generating code blocks proves more challenging than generating single functions as action:** Across all models, there is a higher occurrence of errors when generating code blocks compared to single functions. This pattern highlights the greater complexity and increased likelihood of errors associated with composing multiple lines of code over simple function calls.

**3) Advanced models engage in more valid action steps, resulting in better forecasting performance:** Advanced models, such as GPT-4o-mini, engage in significantly more effective action steps, as shown by no case in ending with consecutive invalid or repeated actions. Llama-3.1-8B-Instruct, on the other hand, has the most number of action steps with execution errors in Figure 5 and the most number of cases ending without a direct final answer. This shows the importance of the planning and tool-use ability of LLM agents in effectively gathering historical information and making predictions.

### D.2 ANALYSIS ON AGENT REACT ITERATION SEQUENCE LENGTH

We further show the correlation of the model’s forecasting performance to the number of ReAct iterations it takes to reach the final answer. Figure 9 shows the statistics of GPT-3.5-Turbo’s performance with the “Single Function” action type and 20 as the maximum action steps allowed. It shows a slightly negative correlation between the two. We further investigate the longest few test cases, and find the model mostly trapped in code generation errors and fails to self-correct from the error message, and therefore, fails to gather useful information.

### D.3 EVALUATION FOR TRADITIONAL TKG AND NLP FORECASTING METHODS

We include the following additional traditional baselines that are fine-tuned:



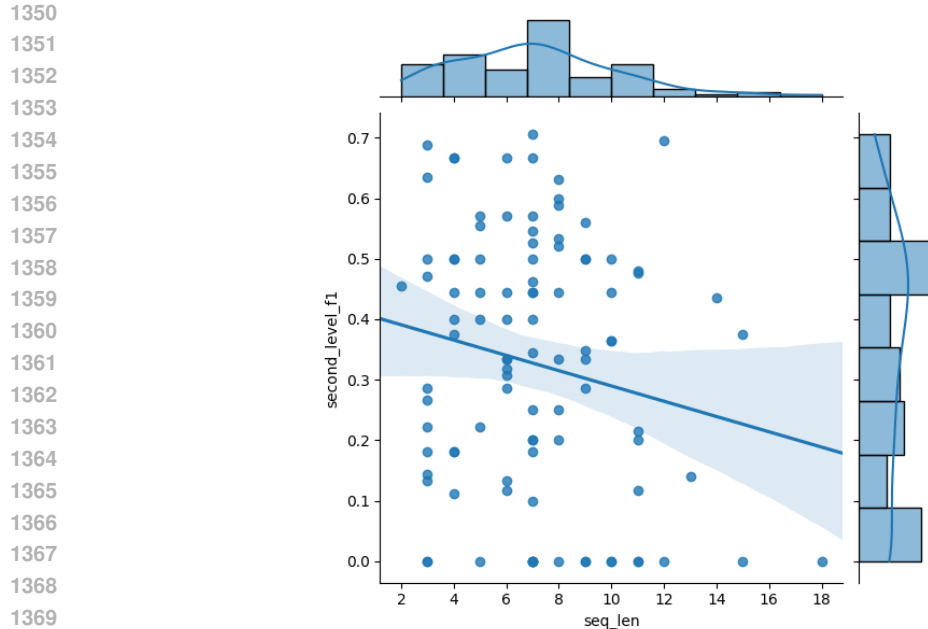


Figure 9: Correlation of F1 Accuracy to Action Sequence Length

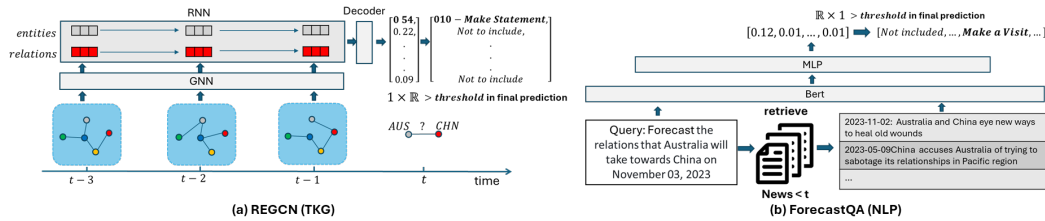


Figure 10: Illustration of RE-GCN and the adapted ForecastQA aggregate historical information and output forecasting probability (0-1) for each relation type. In comparison, examples using LLM agent in Appendix F provide explicit reasoning traces with better interpretability.

- **Structured TKG baseline**, we implemented RE-GCN (Li et al., 2021b), which learns relation in their graph edge embeddings and supports relation prediction.
- **Traditional textual baseline**, we adopted ForecastQA (Jin et al., 2021a), implementing a Bert-based classification model with BM25 to retrieve relevant news articles as additional encoder input.

Figure 10 illustrates the proposed frameworks of REGCN and ForecastQA, along with their threshold-based prediction processes. REGCN learns entity and relation representations in the event knowledge graph at each timestamp, modeling temporal development through a recurrent neural network. The learned representations are then used by a decoder to compute probabilities for each relation type. In contrast, ForecastQA employs a BERT encoder to process both the original query and retrieved news articles, utilizing a classification network for relation type probability computation.

We trained these traditional baselines on data prior to 2023-06 and evaluated them on the 2024-02 test split. Optimal thresholds were determined through a search over 0.3, 0.5, 0.7, and 0.9, resulting in 0.5 for REGCN and 0.3 for ForecastQA. The experimental results of these traditional models are shown in Table 2, along with LLMs’ performance.

REGCN demonstrates superior performance in predicting higher-level relation types, particularly at the binary and quadratic levels, suggesting a good level of relation representation learning ability. However, it struggles with accurate forecasting at the more granular second level. ForecastQA achieves higher precision but significantly lower recall compared to REGCN, due to only predicting the relation type “Make a Visit” in most queries.

Notably, data-specific fine-tuning allows these traditional methods to achieve competitive performance with zero-shot smaller LLM agents, though they still largely underperform the most advanced LLMs.

Overall, these experiments and findings highlight MIRAI’s value as a comprehensive benchmark for comparing diverse forecasting methods.

#### D.4 HUMAN FORECASTING PERFORMANCE AS A REFERENCE

Table 6: Human and LLM agent forecasting performance on the sampled test queries. We have 2 human evaluators and we take the average performance. The best-performing score is highlighted in **bold** and the second-best is underlined.

Model	Training Data Cutoff Date	Action Type	Binary KL ( $\downarrow$ )	Quad KL ( $\downarrow$ )	First-level Relation (%)			Second-level Relation (%)		
					Pre. ( $\uparrow$ )	Rec. ( $\uparrow$ )	F1 ( $\uparrow$ )	Pre. ( $\uparrow$ )	Rec. ( $\uparrow$ )	F1 ( $\uparrow$ )
<b>Human evaluators</b>	—	—	<b>0.04</b>	<u>1.37</u>	<u>62.73</u>	<b>88.70</b>	<b>68.29</b>	<b>54.54</b>	<b>74.53</b>	<b>56.78</b>
Mistral-7B-Instruct-v0.2	—	Single Function	10.35	13.74	25.0	14.22	13.64	10.59	12.82	5.51
		Code Block	8.21	11.59	30.0	7.04	10.69	23.33	4.74	7.62
GPT-3.5-Turbo	2021-09	Single Function	1.03	3.19	<b>69.17</b>	62.97	<u>54.87</u>	<u>53.33</u>	57.24	<u>45.86</u>
		Code Block	3.95	8.16	36.44	48.68	28.74	23.06	36.91	13.46
GPT-4-Turbo	2023-12	Single Function	1.94	4.09	62.5	57.04	40.83	34.31	59.74	31.72
		Code Block	0.18	2.77	35.33	54.48	34.03	25.46	56.05	27.83
GPT-4o	2023-10	Single Function	0.17	1.46	47.0	<u>64.12</u>	43.17	32.62	60.3	30.92
		Code Block	<u>0.16</u>	<b>1.22</b>	48.06	78.2	53.02	39.0	<u>69.88</u>	40.32

To establish a reference point for our LLM agent evaluations, we conducted a human forecasting task with the following parameters:

- **Scope:** 10 queries  $(t, s, ?, o)$  covering 51 distinct events.
- **Participants:** 2 college students (non-political science majors).
- **Task:** Given an event query  $(t + 1, s, ?, o)$ , generate forecasts for all possible relations.
- **APIs:** Participants used the same API library as the LLM agents.
- **Evaluation:** Applied the same metrics as those used for LLM agents.

As shown in Table 6, human performance surpassed that of LLM agents in most metrics, especially in recall. This highlights significant room for improvement in LLM performance.

It is important to note that our evaluators do not specialize in political science and likely fall short of what domain experts could achieve. This suggests that human forecasting, even at current performance levels, has room for enhancement.

We also acknowledge that event forecasting is an inherently challenging task for both humans and AI. The superior performance of non-expert human evaluators underscores the complexity of this task and the potential for further advancements in both human and AI forecasting capabilities.

#### D.5 HUMAN EVALUATION ON THE DATASET QUALITY

To further assess dataset quality, we conducted a human evaluation on a subset of the test set consisting of 10  $(t, s, ?, o)$  queries, corresponding to 51 distinct  $(t, s, r, o)$  events. Two human annotators evaluated whether the ground-truth events were correctly extracted based on their source news articles, scoring each event as 0 (incorrect) or 1 (correct). The average score across all evaluated events is 0.82, indicating high accuracy in the sampled test events. This evaluation confirms the reliability of our data cleaning process and the overall quality of our dataset.

#### D.6 FORECASTING PERFORMANCE ON DIFFERENT TEST SPLITS

Table 7: Evaluation results of second-level relation forecasting on the different test splits, using “Single Function” as the action type. The best-performing score is highlighted in **bold** and the second-best is underlined.

Model	Training Data Cutoff Date	2023-11			2023-12			2024-01			2024-02		
		Pre. ( $\uparrow$ )	Rec. ( $\uparrow$ )	F1 ( $\uparrow$ )	Pre. ( $\uparrow$ )	Rec. ( $\uparrow$ )	F1 ( $\uparrow$ )	Pre. ( $\uparrow$ )	Rec. ( $\uparrow$ )	F1 ( $\uparrow$ )	Pre. ( $\uparrow$ )	Rec. ( $\uparrow$ )	F1 ( $\uparrow$ )
Llama-3-8B-Instruct	2023-03	10.7 <sub>±4.0</sub>	6.1 <sub>±2.4</sub>	6.1 <sub>±1.5</sub>	13.6 <sub>±3.0</sub>	10.2 <sub>±2.8</sub>	8.7 <sub>±1.8</sub>	16.0 <sub>±1.2</sub>	9.0 <sub>±0.8</sub>	8.7 <sub>±0.1</sub>	15.8 <sub>±0.5</sub>	11.8 <sub>±0.6</sub>	10.3 <sub>±0.1</sub>
Llama-3.1-8B-Instruct	2023-12	20.6 <sub>±9.0</sub>	22.3 <sub>±5.4</sub>	15.8 <sub>±5.7</sub>	21.2 <sub>±7.0</sub>	18.5 <sub>±1.5</sub>	15.9 <sub>±3.6</sub>	23.2 <sub>±1.5</sub>	22.4 <sub>±3.2</sub>	16.3 <sub>±2.0</sub>	22.7 <sub>±3.2</sub>	16.8 <sub>±0.5</sub>	14.8 <sub>±0.7</sub>
GPT-4-Turbo	2023-12	33.5 <sub>±7.5</sub>	43.5 <sub>±5.4</sub>	30.0 <sub>±1.9</sub>	31.5 <sub>±4.5</sub>	33.9 <sub>±0.5</sub>	25.8 <sub>±3.1</sub>	36.5 <sub>±3.4</sub>	41.9 <sub>±4.7</sub>	32.2 <sub>±2.8</sub>	33.5 <sub>±4.4</sub>	41.6 <sub>±1.3</sub>	28.9 <sub>±3.2</sub>
GPT-4o-mini	2023-10	41.3 <sub>±9.0</sub>	41.4 <sub>±1.4</sub>	32.8 <sub>±2.6</sub>	39.4 <sub>±7.5</sub>	25.4 <sub>±2.9</sub>	25.9 <sub>±3.2</sub>	45.9 <sub>±3.1</sub>	36.6 <sub>±1.7</sub>	33.2 <sub>±0.7</sub>	40.0 <sub>±5.5</sub>	32.6 <sub>±1.6</sub>	29.7 <sub>±3.8</sub>

We discuss the effect of models’ knowledge cutoff on different test splits in Sec. 3.3. We show the full forecasting performance of second-level relation prediction with different base models over the four test splits in Table 7.

## D.7 ANALYSIS WITH RETRIEVE-AUGMENTED GENERATION (RAG) METHODS

Table 8: Evaluation results of GPT-4o-mini on the 2024-02 test split using different non-agentic methods and the ReAct agent with Single Function action type. The best-performing score is highlighted in **bold** and the second-best is underlined.

Method	Augmented Context	Agent API	Binary KL ( $\downarrow$ )	Quad KL ( $\downarrow$ )	First-level Relation (%)			Second-level Relation (%)		
					Pre. ( $\uparrow$ )	Rec. ( $\uparrow$ )	F1 ( $\uparrow$ )	Pre. ( $\uparrow$ )	Rec. ( $\uparrow$ )	F1 ( $\uparrow$ )
Direct IO	—	—	3.6 $\pm$ 1.0	7.6 $\pm$ 1.9	39.5 $\pm$ 3.2	44.8 $\pm$ 3.2	34.9 $\pm$ 3.5	15.4 $\pm$ 0.8	23.9 $\pm$ 3.6	15.4 $\pm$ 0.2
	CAMEO	—	5.0 $\pm$ 1.5	7.3 $\pm$ 1.9	35.5 $\pm$ 4.6	32.6 $\pm$ 2.2	28.6 $\pm$ 2.7	10.0 $\pm$ 1.5	14.2 $\pm$ 0.9	10.1 $\pm$ 0.8
RAG	<i>Events-Only</i>	—	<b>2.2</b> $\pm$ 0.9	<b>5.9</b> $\pm$ 2.0	57.5 $\pm$ 3.5	<b>53.4</b> $\pm$ 3.4	<b>50.5</b> $\pm$ 3.8	32.4 $\pm$ 1.1	<b>43.9</b> $\pm$ 2.0	<b>33.2</b> $\pm$ 1.4
	<i>News-Only</i>	—	9.1 $\pm$ 2.8	12.7 $\pm$ 2.9	47.2 $\pm$ 0.8	23.2 $\pm$ 2.4	25.4 $\pm$ 0.2	19.5 $\pm$ 2.0	14.9 $\pm$ 2.1	13.4 $\pm$ 0.8
	<i>All</i>	—	<u>2.3</u> $\pm$ 1.4	<u>6.3</u> $\pm$ 2.0	59.0 $\pm$ 1.2	<u>48.1</u> $\pm$ 1.2	<u>46.7</u> $\pm$ 0.4	36.4 $\pm$ 5.3	<u>38.8</u> $\pm$ 1.2	<u>32.1</u> $\pm$ 2.4
ReAct	—	<i>Event-Only</i>	3.3 $\pm$ 0.8	7.7 $\pm$ 1.4	<b>62.8</b> $\pm$ 10.5	39.0 $\pm$ 0.8	41.7 $\pm$ 5.3	<u>44.2</u> $\pm$ 3.3	37.0 $\pm$ 0.8	30.7 $\pm$ 0.9
	—	<i>News-Only</i>	6.5 $\pm$ 1.7	13.0 $\pm$ 2.1	41.5 $\pm$ 6.1	16.8 $\pm$ 0.7	20.2 $\pm$ 1.9	17.8 $\pm$ 0.2	12.2 $\pm$ 1.0	12.5 $\pm$ 0.5
	—	<i>All</i>	3.6 $\pm$ 0.9	8.0 $\pm$ 1.5	<u>61.7</u> $\pm$ 10.1	38.6 $\pm$ 1.9	40.7 $\pm$ 5.6	<b>46.3</b> $\pm$ 4.4	32.9 $\pm$ 3.8	31.1 $\pm$ 2.6

### D.7.1 METHODS AND EXPERIMENTAL SETUP

Besides the **Direct IO / QA** and three **ReAct** agents with different tool-use that we already implemented in the paper, we add the following baselines:

- **Direct QA with Augmentation** (for comment Q4)
  - **QA with CAMEO**: We provide the CAMEO ontology in an ordered dictionary format mapping relation codes to their names and detailed descriptions. This is closer to the QA-format the authors mention, and a more fair comparison as the model can refer to the output vocabulary without needing to memorize CAMEO codes.
- **RAG Methods**

Following recent work, we implement three RAG variants:

  - **RAG Events-Only**: Following GPT-NeoX-ICL ?, we explore rule-based approaches for retrieving historical facts. Using the 'Pair' and 'Undirectional' setting, given a query event  $(s, ?, o, t)$ , we retrieve historical events  $(s, r \in \mathcal{R}, o, < t)$  and  $(o, r \in \mathcal{R}, s, < t)$ . Events are sorted by recency with a cap of 30, aligning with the default cap of the `get_events` API function.
  - **RAG News-Only**: Following TCELongBench ?, we employ BM25 retrieval to fetch the most query-relevant news articles before the query date. The top 15 news articles are retained, matching the default cap of the `get_news_articles` API function.
  - **RAG All**: Combines both retrieved structured events and textual news articles.

### D.7.2 COMPARISON OF RETRIEVAL STRATEGIES WITH AGENT

Both baseline approaches—CAMEO context augmentation and RAG methods—employ static, predefined retrieval strategies that are fixed for all queries and executed only once per query.

In contrast, our agentic approach enables dynamic, multi-step information gathering and reasoning. The agent **can** replicate the baseline retrieval strategy by fixing certain API parameter values, for example, the agent can use the function call `get_events(head_entities=[s, o], tail_entities=[o, s])` to get the retrieved context as RAG Event-Only, and use the function call `get_news_articles(text_description='(t, s, ?, o)')` to get the retrieved context as RAG News-Only; its capabilities extend far beyond these static approaches through its flexible parameter settings for each function call and multiple steps per query.

This multi-step, adaptive approach represents a fundamental shift from static retrieval to dynamic information gathering and reasoning, though it introduces higher requirements for the LLM’s planning ability in:

- Automatically selecting optimal information-gathering strategies
- Integrating and reasoning over information of different formats
- Adjusting strategies based on intermediate findings and current context

### D.7.3 KEY FINDINGS AND ANALYSIS

Our experiment results reveal several important insights:

#### Performance of the RAG Baselines:

RAG demonstrates improved precision over Direct IO (15.4% Pre in second-level) when using either event (32.4% Pre in second-level) or text (19.5% Pre in second-level) information source independently, with event data contributing more significantly to recall (43.9% Rec in second-level). When combining the two information sources, RAG achieves higher precision (36.4% Pre in second-

level) but with a lower recall than RAG Event Only (from 43.9% to 38.8%), leading to a lower overall F1 score (from 33.2% to 32.1%), this suggests that a simple combination of both information in the context not effectively and collaboratively contribute to a better forecasting performance.

### Comparison and Insights for Agent:

ReAct agents exhibit similar performance patterns with RAG when using different information sources (changed by the type of API functions available for the agent). Specifically, event data also contributes to high precision for the agent (62.8% Pre at first-level and 44.2% at second-level), outperforming RAG baselines. The structured event data consistently provides stronger signals for forecasting across both approaches, likely due to its standardized format.

However, agent baselines generally obtain a lower recall than RAG baselines. Meanwhile, while ReAct with full API access achieves higher average F1 scores in second-level relation prediction than its partial access performance (31.1% than 30.7% and 12.5%), the benefits of combining information sources aren't consistent across all relation hierarchies and methods, similar to RAG, suggesting substantial room for exploring more effective information integration strategies.

This reveals both promises and challenges of the agent's more flexible retrieval approach: RAG's predefined and fixed retrieval strategies can often yield stable performance, while the agent's dynamic and multi-step retrieval allows for flexible information gathering and integration, yet higher requirements for planning and reasoning sometimes also lead to relatively lower performance.

These observations underscore the core purpose of our benchmark: not just to compare current methods but to encourage the development of more advanced agentic forecasting approaches. The current performance patterns suggest significant opportunities for improving agent architectures, particularly in:

- Developing more robust and automatic planning strategies for multi-step information gathering
- Improving information integration capabilities across different information sources, formats, and temporal scales

## D.8 ANALYSIS WITH HEURISTIC-BASED AND TRADITIONAL TKG METHODS

Table 9: Evaluation results on the 2024-02 test split for relation prediction using heuristic-based and TKG-based methods and LLM agents based on GPT-4o-mini. The best-performing score is highlighted in **bold** and the second-best is underlined.

Method	Training Data Cutoff Date	Prompt	MRR (%) (↑)	Hit@10 (%) (↑)	Binary KL (↓)	Quad KL (↓)	First-level Relation (%)			Second-level Relation (%)		
							Pre. (↑)	Rec. (↑)	F1 (↑)	Pre. (↑)	Rec. (↑)	F1 (↑)
RE-GCN	2023-06	—	1.6	2.2	<u>0.4</u>	<u>0.8</u>	24.4	<u>90.6</u>	34.3	4.4	<b>83.9</b>	7.9
	2023-08		1.9	2.8	<u>0.4</u>	1.1	23.9	86.1	32.9	4.6	40.0	7.0
	2023-10		1.7	2.5	<b>0.3</b>	1.0	24.8	78.2	32.4	3.9	25.7	5.6
	2023-12		2.9	5.7	<b>0.3</b>	2.5	23.9	74.4	31.3	5.5	28.4	7.9
Recurrency (Strict)	2023-06	—	<u>17.4</u>	<u>45.0</u>	3.2	3.6	32.8	77.1	42.9	18.7	67.8	27.2
	2023-08		17.1	<b>45.3</b>	3.2	3.6	32.3	78.2	42.7	18.0	69.9	26.9
	2023-10		15.8	41.0	2.4	3.1	29.7	83.5	41.3	14.3	76.8	23.0
	2023-12		<b>17.8</b>	43.2	2.1	2.5	29.8	86.0	41.6	14.2	80.1	23.0
ReAct	2023-10	<b>Set Prediction</b>	—	—	3.6	8.0	<b>61.7</b>	38.6	40.7	<b>46.3</b>	32.9	31.1
		Rank (k=10)	—	25.7	0.6	1.4	<u>47.5</u>	70.2	<b>48.9</b>	<u>38.1</u>	61.8	<b>38.2</b>
		Rank (k=30)	—	12.0	<b>0.3</b>	<u>0.8</u>	34.9	<b>91.2</b>	45.8	22.5	<u>82.8</u>	<u>31.7</u>
		Rank (all)	13.9	14.1	2.1	2.8	27.0	86.2	37.9	12.5	81.4	20.2
		Rank w.Prob (k=10)	—	26.8	1.1	2.5	47.3	67.7	<u>48.3</u>	37.9	59.2	<b>38.2</b>
		Rank w.Prob (k=30)	—	10.8	<b>0.3</b>	<b>0.6</b>	34.8	86.6	45.3	22.2	76.4	31.0
		Rank w.Prob (all)	12.6	14.9	2.4	2.7	28.5	83.0	38.3	12.7	78.6	20.6

### D.8.1 TASK FOCUS: RELATION PREDICTION

Previous temporal knowledge graph research has explored both link prediction and relation prediction tasks, with notable works like DynamicGCN (Deng et al., 2020b), Glean (Deng et al., 2020a), and CMF (Deng et al., 2021) focusing on relation prediction. It is important to emphasize that neither task holds inherent priority over the other; rather, each serves distinct analytical purposes tailored to specific research objectives.

In MIRAI, we focus on relation prediction as our primary task given our interest in studying dynamic relationship shifts between countries over time. This choice is particularly significant due to the structured nature of CAMEO relations in international event data. The CAMEO ontology offers a hierarchically organized framework that encompasses the entire spectrum of international interactions, ranging from material cooperation (e.g., providing aid, military collaboration) and verbal cooperation (e.g., diplomatic statements, expressions of support) to verbal conflict (e.g., accusations, rejections)

and material conflict (e.g., military actions, sanctions). This natural progression of political interactions—from the most cooperative to the most conflictual—provides a clear framework for analyzing the evolution of international relationships and a nuanced benchmark for assessing models’ capacity to forecast shifts in these dynamics over time.

#### D.8.2 EXPERIMENTAL SETUP

We evaluate three approaches in detail:

**RE-GCN** (Li et al., 2021b):

- **Model Architecture:** we follow the original method, where the model combines relation-aware graph convolutional layers with recurrent neural networks to jointly model structural dependencies and temporal dynamics in TKGs.
- **Data Cutoffs:** We experiment with four cutoff times (2023-06, 2023-08, 2023-10, 2023-12). For example, with 2023-10 cutoff, the model trains on data until 2023-10, validates on 2023-11, and tests on the 2024-02 split.
- **Test Input:** Uses single-step prediction with a 7-day historical window (t-7 to t-1), aligning with our agent experiments’ one-day forecasting horizon.

**Recurrency (Strict)** (Gastinger et al., 2024):

- **Model:** The original work introduces three baselines that place strong inductive bias on fact recurrence over time: strict recurrency, relaxed recurrency, and their combination. While the original work and its scoring functions are specifically designed for link prediction, we adapt the strict recurrency variant for relation prediction. For a query event  $(s, ?, o, t)$ , we compute scores for all relations  $r \in \mathcal{R}$  using:

$$\phi_{\Delta}((s, r, o, t), G) = \begin{cases} \Delta(t, \max\{k | (s, r, o, k) \in G\}) & \text{if } \exists k \text{ with } (s, r, o, k) \in G \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

where  $\Delta(t, k) = k/t$  measures temporal proximity.

- **Data Access:** Cutoff dates determine available historical events, e.g., 2023-10 means using only events before 2023-10-31 for score computation.

**ReAct Agent:**

- **Set Prediction:** Uses original prompt and output format, with the output directly used for set prediction metrics.
- **Ranking:** Modified prompts for ordered relation lists (k=10, 30, all). Only rank=all configurations are used for MRR calculation, while all configurations support Hit@10 evaluation.
- **Probabilistic Ranking:** Additional prompts requesting probability scores for ranked relations.
- **Evaluation:** All ranking evaluations use time-aware filtering to maintain consistency with other methods.

#### D.8.3 METRIC SELECTION AND JUSTIFICATION

**Set-based Metrics:** For set prediction tasks, models output a discrete set of predicted relations for each query. We evaluate these predictions using:

- Precision: Proportion of predicted relations that are correct
- Recall: Proportion of actual relations that are predicted
- F1 Score: Harmonic mean of precision and recall

**Ranking-based Metrics:** For ranking tasks, models output an ordered list of relations with associated scores. Following TKG conventions of time-aware filtering (Bordes et al., 2013; Han et al., 2021), we implement:

- Mean Reciprocal Rank (MRR): Average reciprocal of the first correct relation’s rank
- Hit@10: Proportion of queries where at least one correct relation appears in top-10 predictions

**Choice of Primary Metrics:** We prioritize set-based metrics as our primary evaluation criteria for Model Capability: The primary consideration is that current LLMs are better suited to generating discrete predictions through natural language reasoning than producing comprehensive ranked lists. This is evidenced by our experimental findings:

- **List Length Sensitivity:** ReAct agent’s performance deteriorates with longer list requirements, with k=30 and k=all performing worse than k=10 in ranking metrics (Hit@10: 25.7% for k=10 vs 12.0% for k=30).



- **Prompt Sensitivity:** Performance varies between pure ranking and probability-weighted ranking (Hit@10: 25.7% vs 26.8% for k=10, and MRR: 13.9% vs 12.6% for k=all), suggesting that ranking outputs are sensitive to the prompt formulation and output format.

Given these challenges, we opted for metrics that more directly and reliably assess the agents' ability to predict discrete events without the confounding factors introduced by list generation and ranking.

#### D.8.4 RESULTS ANALYSIS AND IMPLICATIONS

Our comprehensive evaluation reveals several significant insights about different forecasting approaches:

**Performance of the Recurrency Model and Insights for Agents:** The experiment results of the Recurrency (strict) model are shown in row 5-8 of the table. We observe that the Recurrency baseline demonstrates strong performance in ranking metrics (17.8% MRR and 43.2% Hit@10 with 2023-12 cutoff), leading other models; it also shows consistently high recall (86.0% Recall at first-level relation, and 80.1% recall at second-level relations with 2023-12 cutoff), suggesting that international events indeed often follow repetitive patterns.

Comparing with the ReAct agent (Set Prediction, the setting we used in the main paper) in row 9, we find that although the LLM agent could achieve much higher forecasting precision, it obtains much lower recall than the Recurrency baseline. We conducted a detailed analysis of this behavior by manually going through the reasoning traces generated by the LLM agent in the test set. One possible reason is that the agent has a strong tendency to select only a subset of the most frequent historical events in its prediction. For example, for the query (2024-02-01, PSE, ?, EGY), it uses function calls like `get_relation_distribution(date_range=DateRange(start_date=Date("2023-01-31"), end_date=Date("2024-01-31")), head_entities=[ISOCODE("PSE")], tail_entities=[ISOCODE("EGY")])`, and obtained a full frequency list as:

```
{CAMEOCODE(code='042'): 32, CAMEOCODE(code='192'): 18,
CAMEOCODE(code='040'): 13, CAMEOCODE(code='043'): 12,
CAMEOCODE(code='046'): 8, CAMEOCODE(code='080'): 6,
CAMEOCODE(code='036'): 4, CAMEOCODE(code='010'): 4,
CAMEOCODE(code='190'): 3, CAMEOCODE(code='073'): 3,
CAMEOCODE(code='030'): 3, CAMEOCODE(code='084'): 3,
CAMEOCODE(code='020'): 3, CAMEOCODE(code='172'): 2,
CAMEOCODE(code='014'): 2, CAMEOCODE(code='070'): 2,
CAMEOCODE(code='044'): 2, CAMEOCODE(code='086'): 1,
CAMEOCODE(code='013'): 1, CAMEOCODE(code='051'): 1}.
```

It then has a further step of checking recent news articles and obtains its final prediction as 040, 042, and 192, which are the top three frequent relations.

The effectiveness of simple temporal recurrency heuristics underscores the importance of incorporating more historical pattern analysis in the future development of forecasting agents, in particular, improving their recall of capturing a greater proportion of true relationships between countries.

**Performance of the TKG Baselines and Insights for Agents:** The experiment results of the RE-GCN are shown in rows 1-4 of the table. We observe that the RE-GCN demonstrates strong performance in high-level relation prediction, resulting in 0.3 for binary-level relation KL (cooperation or conflict) and 0.8 for quadratic-level relation KL (verbal/material cooperation/conflict), reflecting its advantage in capturing the high-level dynamics over bilateral relationships. It also shows consistently high recall in more fine-grained relation levels (90.6% Recall at first-level relation, and 83.9% recall at second-level relations with 2023-06 cutoff), suggesting its effectiveness in modeling positive correlation between query and multiple ground-truth relations.

Comparing with the ReAct agent (Set Prediction, the setting we used in the main paper) in row 9, we find that although the LLM agent also could achieve much higher forecasting precision, it obtains much lower recall than the RE-GCN baseline. We manually go through the reasoning traces generated by the LLM agent in the test set, and conclude the following possible insights and future directions of improvement compared with TKG baselines:

One major possible reason is that the current agent mostly focuses on analyzing only the bilateral events between the query entities  $s$  and  $o$ . For instance, it typically sets the function parameter `head_entities` to the query subject, and `tail_entities` to the query object only, obtaining only events and news directly between the two. However, this analysis largely oversimplifies real international relationships where countries have engaged in multi-party and complex interactions.

Events between two countries could be affected by regional or global events. In contrast, TKG methods excel in capturing this multi-party and multi-relational history by leveraging multi-layer graph convolutions, where neighboring information is aggregated to enhance the modeling of each node (entity embedding) and edge (relation embedding). Therefore, when making predictions between two countries, the TKG models consider a much broader relation network than the current LLM agents, leading to higher recall and better generalization, especially when the bilateral history is sparse.

Another problem we observed from the current LLM agent behavior is its tendency to hallucinate, particularly in listing the existence of relations and interpreting the meanings of relations in the CAMEO ontology, which leads to lower precision and recall. For example, in the example we show in Appendix F.3, in its trajectory step 3, the agent attempts to explain and conclude its final prediction: '042' Make a public statement (high frequency in historical data); '036' Negotiate (also high frequency); '057': Express intent to cooperate (indicated by recent news context). However, the correct meanings are *Make a visit* for '042' and *Sign formal agreement* for '057'. This example highlights two issues: firstly, the LLM agent's overconfidence in its understanding of the CAMEO ontology without verifying the relation meanings through function calls (such as `map_cameo_to_relation` and `map_description_to_cameo`); secondly, its over-reliance on the semantic meaning of relations rather than their structural context. In contrast, TKG models learn relations by leveraging the historical graph structure, which inherently learns to capture the contextual meaning of each relation.

To enhance future LLM agents, incorporating a hybrid approach that combines semantic understanding with structural learning from TKGs could help to enhance relation modeling and address certain hallucinations.

#### D.9 EVALUATION ON IMPACT OF LLM PARAMETER SIZE

Table 10: Evaluation results on the 2024-02 test split using different base LLMs with different number of model parameters. The best-performing score is highlighted in **bold** and the second-best is underlined.

Base LLM	Training Data Cutoff Date	Action Type	Binary KL ( $\downarrow$ )	Quad KL ( $\downarrow$ )	First-level Relation (%)			Second-level Relation (%)		
					Pre. ( $\uparrow$ )	Rec. ( $\uparrow$ )	F1 ( $\uparrow$ )	Pre. ( $\uparrow$ )	Rec. ( $\uparrow$ )	F1 ( $\uparrow$ )
Llama-3.2-1B-Instruct	2023-12	Single Func	<u>9.5</u> $\pm 1.9$	16.0 $\pm 1.7$	23.7 $\pm 6.6$	10.0 $\pm 1.9$	<u>11.7</u> $\pm 2.8$	8.8 $\pm 2.6$	7.2 $\pm 0.2$	6.1 $\pm 1.3$
		Code Block	10.1 $\pm 2.2$	16.2 $\pm 2.2$	24.0 $\pm 5.7$	8.1 $\pm 2.1$	10.0 $\pm 3.4$	<u>7.6</u> $\pm 1.9$	5.7 $\pm 0.7$	5.1 $\pm 1.8$
Llama-3.2-3B-Instruct	2023-12	Single Func	12.1 $\pm 2.2$	<u>15.4</u> $\pm 1.9$	<b>36.3</b> $\pm 2.2$	<b>13.1</b> $\pm 3.7$	<b>16.7</b> $\pm 3.0$	<b>19.9</b> $\pm 0.7$	<b>8.3</b> $\pm 2.6$	<b>9.3</b> $\pm 0.9$
		Code Block	<b>9.3</b> $\pm 2.1$	<b>15.1</b> $\pm 4.5$	<u>26.7</u> $\pm 4.7$	<u>10.5</u> $\pm 0.5$	11.0 $\pm 0.8$	<u>13.1</u> $\pm 1.5$	<u>8.0</u> $\pm 0.4$	<u>7.1</u> $\pm 0.3$

While LLM's parameter count can influence model capabilities, it is not the sole or even primary determinant of performance. Other crucial factors include model architecture, training data quality and recency, and advanced training techniques. For instance, GPT-4o-mini's competitive performance with GPT-3.5-turbo, despite having fewer parameters, can be attributed to more recent architectural improvements, different training data, and advanced training methodologies.

To systematically investigate the impact of parameter size while controlling for other variables, we conduct additional experiments comparing models within the same family. Table 10 presents evaluation results comparing Llama-3.2-1B-Instruct and Llama-3.2-3B-Instruct on the 2024-02 test split, both sharing the same architecture, training data cutoff (2023-12), and training methodology.

Our analysis reveals two key findings:

1) **Parameter size shows a consistent positive correlation with forecasting performance within the same model family.** The 3B model outperforms its 1B counterpart across all prediction levels, from binary (12.1% vs 9.5% in Single Function mode) to second-level relations (9.3% vs 6.1% in Single Function mode).

2) **The impact of parameter size varies across different action types.** While the 3B model maintains its advantage in both modes, the performance gap between 1B and 3B models narrows with Code Block actions. This smaller gap likely reflects the increased complexity of code generation, as our earlier experiments showed that Code Block actions can potentially hurt smaller, less capable models while benefiting more advanced ones.

## E ADDITIONAL INFORMATION ABOUT API

### E.1 OVERVIEW OF API DATA CLASSES AND FUNCTIONS

Table 11: API data classes and their attributes

Class Name	Attributes and Types
Date	date: <code>str</code> # 'YYYY-MM-DD'
DateRange	start_date: Optional[Date], end_date: Optional[Date]
ISOCODE	code: <code>str</code> # 3-letter ISO code
Country	iso_code: ISOCODE, name: <code>str</code>
CAMEOCODE	code: <code>str</code> # CAMEO code
Relation	cameo_code: CAMEOCODE, name: <code>str</code> , description: <code>str</code>
Event	date: Date, head_entity: ISOCODE, relation: CAMEOCODE, tail_entity: ISOCODE
NewsArticle	date: Date, title: <code>str</code> , content: <code>str</code> , events: List[Event]

Table 12: API functions categorized by functionality

Functions related to Countries and Relations
<code>map_country_name_to_iso(name: str) -&gt; List[Country]</code>
<code>map_iso_to_country_name(iso_code: ISOCODE) -&gt; str</code>
<code>map_relation_description_to_cameo(description: str) -&gt; List[Relation]</code>
<code>map_cameo_to_relation(cameo_code: CAMEOCODE) -&gt; Relation</code>
<code>get_parent_relation(cameo_code: CAMEOCODE) -&gt; Relation</code>
<code>get_child_relations(cameo_code: CAMEOCODE) -&gt; List[Relation]</code>
<code>get_sibling_relations(cameo_code: CAMEOCODE) -&gt; List[Relation]</code>
Functions related to Events
<code>count_events(date_range: Optional[DateRange], head_entities: Optional[List[ISOCODE]], tail_entities: Optional[List[ISOCODE]], relations: Optional[List[CAMEOCODE]]) -&gt; int</code>
<code>get_events(date_range: Optional[DateRange], head_entities: Optional[List[ISOCODE]], tail_entities: Optional[List[ISOCODE]], relations: Optional[List[CAMEOCODE]], text_description: Optional[str]) -&gt; List[Event]</code>
<code>get_entity_distribution(date_range: Optional[DateRange], involved_relations: Optional[List[CAMEOCODE]], interacted_entities: Optional[List[ISOCODE]], entity_role: Optional[str]) -&gt; Dict[ISOCODE, int]</code>
<code>get_relation_distribution(date_range: Optional[DateRange], head_entities: Optional[List[ISOCODE]], tail_entities: Optional[List[ISOCODE]]) -&gt; Dict[CAMEOCODE, int]</code>
Functions related to News
<code>count_news_articles(date_range: Optional[DateRange], head_entities: Optional[List[ISOCODE]], tail_entities: Optional[List[ISOCODE]], relations: Optional[List[CAMEOCODE]], keywords: Optional[List[str]]) -&gt; int</code>
<code>get_news_articles(date_range: Optional[DateRange], head_entities: Optional[List[ISOCODE]], tail_entities: Optional[List[ISOCODE]], relations: Optional[List[CAMEOCODE]], keywords: Optional[List[str]], text_description: Optional[str]) -&gt; List[Tuple[Date, str]]</code>
<code>browse_news_article(date: Date, title: str) -&gt; str</code>

## 1782 F ADDITIONAL FORECASTING EXAMPLES OF LLM AGENT

1783 In this section, we show examples of how LLM agents perform reasoning for the forecasting query.  
 1784 The query-specific values are highlighted in yellow in the system prompt and query prompt, which  
 1785 follows the prompt templates shown in Appendix I. For simplicity, we use `{api_description}` as  
 1786 a placeholder for the API specification in the prompt, which is replaced by the actual contents shown  
 1787 in Appendix ?? in experiments.

### 1788 F.1 GPT-4O-BASED AGENT WITH REACT AND “CODE BLOCK” ACTION

#### 1789 F.1.1 QUERY DETAILS

- 1791 • **Query Quadruplet:** (2023-11-03, AUS, ?, CHN)
- 1792 • **Temporal Distance:** 1; therefore, the current date is 2023-11-02
- 1793 • **Agent Max Steps:** 20

#### 1794 F.1.2 QUERY PROMPT

1795 Please forecast the relations that **Australia** will take towards **China** on **November 03, 2023** based  
 1796 on historical information up to **November 02, 2023**. I.e. forecast the relation CAMEO codes in  
 1797 query event `Event(date=2023-11-03, head_entity=ISOCODE(AUS), relation=CAMEOCODE(?),`  
 1798 `tail_entity=ISOCODE(CHN))`.

#### 1800 F.1.3 SYSTEM PROMPT

1801 You are an expert in forecasting future events based on historical data. The database contains news  
 1802 articles from January 1, 2023 to the current date **November 02, 2023** and the events extracted from  
 1803 these articles. The events are in the form of (date, subject country, relation, object country), where  
 1804 the countries are represented by ISO 3166-1 alpha-3 codes and the relations are represented by the  
 1805 CAMEO codes defined in the "Conflict and Mediation Event Observations" ontology. The relations  
 1806 are hierarchical: first-level relations are general parent relations represented by two-digit CAMEO  
 1807 codes, while second-level relations are more specific child relations represented by three-digit  
 1808 CAMEO codes. Child relations have the same first two digits as their parent relations. For example,  
 1809 "01" is a first-level relation, and "010" and "011" are some of its second-level relations. The relations  
 1810 in the database are represented in the second-level form.

1811 Your task is to forecast the future relations between two entities in a given query. You  
 1812 have access to a defined Python API that allows you to query the database for historical events and  
 1813 statistics, and to get precise information about the ISO country codes and CAMEO relation codes.  
 1814 You are also authorized to utilize additional safe, well-established Python libraries such as numpy,  
 1815 pandas, scikit-learn, and NetworkX to enhance your data analysis and forecasting accuracy.

1816 The defined API is described as follows:

```
1817 ``` python
1818 {api_description}
1819 ```
```

1820 You will use an iterative approach, interleaving "Thought", "Action", and "Observation"  
 1821 steps to collect information and perform the forecast. You may perform up to **20** iterations. The steps  
 1822 are as follows:

- 1825 - "Thought": Analyze the current information and reason about the current situation, and  
 1826 predicts which API you want to use (try to use different APIs to collect diverse information) or make  
 1827 a decision that you want to make a final answer.
- 1828 - "Action": Use the API to gather more information or provide the final forecast.
  - 1829 - If gathering more data: the action must be an executable Python code snippet that starts with  
 1830 ```` python` and ends with `````. It can contain multiple lines of codes and function calls using the  
 1831 defined API or Python libraries. You must use `print()` to output the results, and only the printed output  
 1832 will be returned in the observation step.
  - 1833 - If making the final forecast: the action must start immediately with "Final Answer:", and follow  
 1834 with the answer in the expected JSON format. This should not be enclosed within triple backticks.
- 1835 - "Observation": Return the printed output of the executed code snippet.

To make a reasonable forecast, you should collect both news and relational evidence to support your prediction. When you are fully confident that you accumulate enough information to make the final forecast, you should start the "Thought" with your reasoning using the news and structural information to make the prediction, and then start the "Action" step with "Final Answer:" followed by the answer in the expected JSON format. The answer should be a JSON dictionary where the keys are the forecasted two-digit first-level CAMEO codes and the values are lists of forecasted three-digit second-level CAMEO codes that are child relations of the key. For example, "Action: Final Answer: {"01": ["010", "011", "012"], "02": ["020", "023"]}".

The final answer will be evaluated based on the precision and recall of the forecasted first-level and second-level relations, so only include confident first-level and second-level CAMEO codes in your final forecast.

Try to use different APIs and Python libraries to collect diverse information (including multi-hop relations), such as the precise meaning of CAMEO codes, insights from news content, relational data, and statistical analyses to support your forecasts. Consider not only the frequency of the relations but also the temporal aspects of the data when making your forecast.

#### F.1.4 AGENT REACT PROCESS

1. **Thought:** To forecast the relations that Australia will take towards China on November 03, 2023, we need to gather historical data on the interactions between these two countries. This includes examining the frequency and types of relations observed in the past, as well as any recent trends or significant events that might influence future relations.

We will proceed with the following steps:

- Retrieve Historical Events:** Gather events involving Australia (AUS) and China (CHN) to understand the context and types of relations that have occurred in the past.
- Analyze Relation Distribution:** Analyze the distribution of different CAMEO codes in these events to identify the most common types of interactions.
- Examine Recent Trends:** Focus on recent events to detect any emerging trends or shifts in relations.
- Contextual Analysis:** Review news articles to gain insights into the context and underlying reasons for these relations.

Let's start by retrieving historical events between Australia and China.

**Action:**

```

1 from datetime import datetime
2
3 # Define the date range for the historical data collection
4 start_date = Date("2023-01-01")
5 end_date = Date("2023-11-02")
6 date_range = DateRange(start_date=start_date, end_date=end_date)
7
8 # Define the ISO codes for Australia and China
9 aus_code = ISOCode("AUS")
10 chn_code = ISOCode("CHN")
11
12 # Retrieve historical events involving Australia and China
13 events = get_events(date_range=date_range, head_entities=[aus_code], tail_entities=[
14     chn_code])
15
16 # Print the retrieved events
17 for event in events:
18     print(event)

```

**Observation:**

```

1 Event(date=Date("2023-11-02"), head_entity=ISOCode("AUS"), relation=CAMEOCode("042"),
2     tail_entity=ISOCode("CHN"))
3 Event(date=Date("2023-11-02"), head_entity=ISOCode("AUS"), relation=CAMEOCode("040"),
4     tail_entity=ISOCode("CHN"))
5 Event(date=Date("2023-11-02"), head_entity=ISOCode("AUS"), relation=CAMEOCode("036"),
6     tail_entity=ISOCode("CHN"))
7 Event(date=Date("2023-11-01"), head_entity=ISOCode("AUS"), relation=CAMEOCode("042"),
8     tail_entity=ISOCode("CHN"))
9 Event(date=Date("2023-11-01"), head_entity=ISOCode("AUS"), relation=CAMEOCode("173"),
10    tail_entity=ISOCode("CHN"))

```



```

1890 6 Event (date=Date("2023-11-01"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("036"),
1891     tail_entity=ISOCODE("CHN"))
1892 7 Event (date=Date("2023-10-31"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("036"),
1893     tail_entity=ISOCODE("CHN"))
1894 8 Event (date=Date("2023-10-31"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("042"),
1895     tail_entity=ISOCODE("CHN"))
1896 9 Event (date=Date("2023-10-31"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("046"),
1897     tail_entity=ISOCODE("CHN"))
1898 10 Event (date=Date("2023-10-30"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("036"),
1899     tail_entity=ISOCODE("CHN"))
1900 11 Event (date=Date("2023-10-30"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("042"),
1901     tail_entity=ISOCODE("CHN"))
1902 12 Event (date=Date("2023-10-30"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("046"),
1903     tail_entity=ISOCODE("CHN"))
1904 13 Event (date=Date("2023-10-29"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("036"),
1905     tail_entity=ISOCODE("CHN"))
1906 14 Event (date=Date("2023-10-27"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("036"),
1907     tail_entity=ISOCODE("CHN"))
1908 15 Event (date=Date("2023-10-27"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("042"),
1909     tail_entity=ISOCODE("CHN"))
1910 16 Event (date=Date("2023-10-26"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("173"),
1911     tail_entity=ISOCODE("CHN"))
1912 17 Event (date=Date("2023-10-26"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("061"),
1913     tail_entity=ISOCODE("CHN"))
1914 18 Event (date=Date("2023-10-26"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("042"),
1915     tail_entity=ISOCODE("CHN"))
1916 19 Event (date=Date("2023-10-25"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("042"),
1917     tail_entity=ISOCODE("CHN"))
1918 20 Event (date=Date("2023-10-23"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("042"),
1919     tail_entity=ISOCODE("CHN"))
1920 21 Event (date=Date("2023-10-23"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("036"),
1921     tail_entity=ISOCODE("CHN"))
1922 22 Event (date=Date("2023-10-23"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("031"),
1923     tail_entity=ISOCODE("CHN"))
1924 23 Event (date=Date("2023-10-23"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("061"),
1925     tail_entity=ISOCODE("CHN"))
1926 24 Event (date=Date("2023-10-22"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("040"),
1927     tail_entity=ISOCODE("CHN"))
1928 25 Event (date=Date("2023-10-22"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("061"),
1929     tail_entity=ISOCODE("CHN"))
1930 26 Event (date=Date("2023-10-22"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("036"),
1931     tail_entity=ISOCODE("CHN"))
1932 27 Event (date=Date("2023-10-22"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("043"),
1933     tail_entity=ISOCODE("CHN"))
1934 28 Event (date=Date("2023-10-22"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("042"),
1935     tail_entity=ISOCODE("CHN"))
1936 29 Event (date=Date("2023-10-21"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("042"),
1937     tail_entity=ISOCODE("CHN"))
1938 30 Event (date=Date("2023-10-21"), head_entity=ISOCODE("AUS"), relation=CAMEOCODE("040"),
1939     tail_entity=ISOCODE("CHN"))

```

2. **Thought:** The retrieved events indicate a variety of relations between Australia and China. To forecast the future relations, we need to analyze the distribution of these relations to identify the most common types and any recent trends.

Next, we will analyze the distribution of the different CAMEO codes in these events to identify the most frequent interactions. This will help us understand the dominant types of relations and any emerging patterns.

**Action:**

```

1932 1 from collections import Counter
1933 2
1934 3 # Extract the CAMEO codes from the retrieved events
1935 4 cameo_codes = [event.relation.code for event in events]
1936 5
1937 6 # Count the frequency of each CAMEO code
1938 7 cameo_code_counts = Counter(cameo_codes)
1939 8
1940 9 # Print the distribution of CAMEO codes
1941 10 for code, count in cameo_code_counts.items():
1942 11     print(f"CAMEO Code: {code}, Count: {count}")

```

**Observation:**

```

1942 1 CAMEO Code: 042, Count: 10
1943 2 CAMEO Code: 040, Count: 3
1944 3 CAMEO Code: 036, Count: 8
1945 4 CAMEO Code: 173, Count: 2

```

```

1944 5 CAMEO Code: 046, Count: 2
1945 6 CAMEO Code: 061, Count: 3
1946 7 CAMEO Code: 031, Count: 1
1947 8 CAMEO Code: 043, Count: 1

```

3. **Thought:** The distribution of CAMEO codes from the historical events between Australia and China reveals that certain types of interactions are more frequent than others. Specifically, the most common relations are:

- 1951 - **042:** 10 occurrences
- 1952 - **036:** 8 occurrences
- 1953 - **040:** 3 occurrences
- 1954 - **061:** 3 occurrences
- 1955 - **046:** 2 occurrences
- 1956 - **173:** 2 occurrences
- 1957 - **031:** 1 occurrence
- 1958 - **043:** 1 occurrence

1959 To make a more informed forecast, we should understand what these CAMEO codes represent.  
1960 We will map these codes to their descriptions to gain insights into the nature of these interactions.

1961 **Action:**

```

1962 # Define the CAMEO codes to map
1963 cameo_codes_to_map = ["042", "036", "040", "061", "046", "173", "031", "043"]
1964
1965 # Map the CAMEO codes to their descriptions
1966 relations = [map_cameo_to_relation(CAMEOCode(code)) for code in cameo_codes_to_map]
1967
1968 # Print the descriptions of the CAMEO codes
1969 for relation in relations:
1970     print(f"CAMEO Code: {relation.cameo_code.code}, Name: {relation.name}, Description: {
1971         relation.description}")

```

1970 **Observation:**

```

1972 1 CAMEO Code: 042, Name: Make a visit, Description: Travel to another location for a meeting
1973   or other event.
1974 2 CAMEO Code: 036, Name: Express intent to meet or negotiate, Description: Offer, promise,
1975   agree to, or otherwise indicate willingness or commitment to meet, visit, or engage in
1976   talks or negotiations.
1977 3 CAMEO Code: 040, Name: Consult, not specified, Description: All consultations and meetings,
1978   not otherwise specified.
1979 4 CAMEO Code: 061, Name: Cooperate economically, Description: Initiate, resume, improve, or
1980   expand economic exchange or cooperation.
1981 5 CAMEO Code: 046, Name: Engage in negotiation, Description: Negotiate or bargain with others
1982   .
1983 6 CAMEO Code: 173, Name: Arrest or detain, Description: Legal or extrajudicial arrests,
1984   detentions, or imprisonments.
1985 7 CAMEO Code: 031, Name: Express intent to engage in material cooperation, Description: Offer
1986   , promise, agree to, or otherwise indicate willingness or commitment to engage in or
1987   expand material cooperative exchange, such as express intent to cooperate economically
1988   , militarily, on judicial matters, and on intelligence.
1989 8 CAMEO Code: 043, Name: Host a visit, Description: Host or receive a visitor at residence,
1990   office or home country.

```

4. **Thought:** The mapped CAMEO codes and their descriptions provide valuable insights into the nature of interactions between Australia and China. Here are the key relations observed:

- 1989 - **042 (Make a visit):** 10 occurrences – Indicates frequent travel for meetings or events.
- 1990 - **036 (Express intent to meet or negotiate):** 8 occurrences – Shows a strong intent to engage in talks or negotiations.
- 1991 - **040 (Consult, not specified):** 3 occurrences – General consultations and meetings.
- 1992 - **061 (Cooperate economically):** 3 occurrences – Economic cooperation initiatives.
- 1993 - **046 (Engage in negotiation):** 2 occurrences – Actual negotiations or bargaining.
- 1994 - **173 (Arrest or detain):** 2 occurrences – Legal or extrajudicial arrests or detentions.
- 1995 - **031 (Express intent to engage in material cooperation):** 1 occurrence – Intent to cooperate materially.
- 1996 - **043 (Host a visit):** 1 occurrence – Hosting visitors.

1998 Given this data, we can see that the most frequent interactions involve visits, expressions of intent  
 1999 to meet or negotiate, and economic cooperation. These interactions suggest a pattern of diplomatic  
 2000 engagement and economic collaboration, with occasional legal actions.

2001 To further solidify our forecast, we should examine recent news articles to understand the context  
 2002 and any underlying reasons for these interactions. This will help us identify any emerging trends  
 2003 or significant events that might influence future relations.

2004 Next, we will retrieve recent news articles involving Australia and China to gain contextual  
 2005 insights.

2006 **Action:**

```
2007
2008 1 # Define the date range for recent news articles
2009 2 recent_start_date = Date("2023-10-01")
2010 3 recent_end_date = Date("2023-11-02")
2011 4 recent_date_range = DateRange(start_date=recent_start_date, end_date=recent_end_date)
2012 5
2013 6 # Retrieve recent news articles involving Australia and China
2014 7 news_articles = get_news_articles(date_range=recent_date_range, head_entities=[aus_code],
2015 8 tail_entities=[chn_code])
2016 9
2017 10 # Print the retrieved news articles
2018 11 for article in news_articles:
2019     print(article)
```

2020 **Observation:**

```
2021 1 (Date(date="2023-11-02"), "No quick fix to remaining tariffs ahead of China visit")
2022 2 (Date(date="2023-11-02"), "China's aggression in Indo-Pacific leads to "concern and angst":
2023 3 Peter Dutton, Australia's leader of opposition")
2024 4 (Date(date="2023-11-02"), "Australia and China eye new ways to heal old wounds")
2025 5 (Date(date="2023-11-02"), "Albanese must aim beyond bilateral outcomes on Beijing visit")
2026 6 (Date(date="2023-11-02"), "Western Australian Premier Announces Visit to China")
2027 7 (Date(date="2023-11-02"), "Australian PM's China visit will help normalize relations, says
2028 8 ex-trade minister")
2029 9 (Date(date="2023-11-02"), "Nicholas Moore on how to get in early on the next big growth
2030 10 story")
2031 11 (Date(date="2023-11-02"), "PM's China trip cements new era for ties")
2032 12 (Date(date="2023-11-02"), "PM to take a hard line on trade with China")
2033 13 (Date(date="2023-11-01"), "Australian Prime Minister Anthony Albanese to raise imprisoned
2034 14 democracy blogger during China visit")
```

2035 5. **Thought:** The recent news articles provide valuable context for the relations between Australia  
 2036 and China. Key themes include:

2037 1. **Diplomatic Visits and Engagements:** Multiple articles highlight the upcoming visit of  
 2038 Australian Prime Minister Anthony Albanese to China. This aligns with the frequent occurrence  
 2039 of CAMEO code **042 (Make a visit)**.

2040 2. **Economic Cooperation and Trade:** Discussions around trade and economic cooperation are  
 2041 prominent, indicating ongoing efforts to improve economic ties. This supports the relevance of  
 2042 CAMEO codes **061 (Cooperate economically)** and **036 (Express intent to meet or negotiate)**.

2043 3. **Concerns and Tensions:** Articles mention concerns about China's aggression in the Indo-  
 2044 Pacific and human rights issues, suggesting underlying tensions. This is reflected in the occasional  
 2045 occurrence of CAMEO code **173 (Arrest or detain)**.

2046 Given the historical data and recent news context, the most likely relations that Australia will  
 2047 take towards China on November 03, 2023, include diplomatic visits, economic cooperation, and  
 2048 negotiations, with a possibility of addressing underlying tensions.

2049 **Final Answer:**

```
2050 1 {
2051 2   "03": ["031", "036"],
2052 3   "04": ["040", "042", "043", "046"],
2053 4   "06": ["061"],
2054 5   "17": ["173"]
2055 6 }
```

2056 **Action:** Final Answer:

```
2057 1 {
2058 2   "03": ["031", "036"],
2059 3   "04": ["040", "042", "043", "046"],
```

```

2052 4  "06": ["061"],
2053 5  "17": ["173"]
2054 6  }

```

### F.1.5 AGENT FINAL STATUS

- **End State:** Final Answer
- **Number of Steps Taken:** 5
- **Final Answer:**

```

2061 1  {
2062 2  "03": ["031", "036"],
2063 3  "04": ["040", "042", "043", "046"],
2064 4  "06": ["061"],
2065 5  "17": ["173"]
2066 6  }

```

We show the corresponding relation names of the predicted codes here for a better review:

```

2067 1  {
2068 2  "Express intent to cooperate": ["Express intent to engage in material cooperation", "
2069 3  Express intent to meet or negotiate"],
2070 3  "Consult": ["Consult, not specified", "Make a visit", "Host a visit", "Engage in
2071 4  negotiation"],
2072 4  "Engage in material cooperation": ["Cooperate economically"],
2073 5  "Coerce": ["Arrest or detain"]
2074 6  }

```

- **Ground Truth Answer:**

```

2075 1  {
2076 2  "03": ["036"],
2077 3  "04": ["042"],
2078 4  "17": ["172"]
2079 5  }

```

We show the corresponding relation names of the ground-truth codes here for a better review:

```

2081 1  {
2082 2  "Express intent to cooperate": ["Express intent to meet or negotiate"],
2083 3  "Consult": ["Make a visit"],
2084 4  "Coerce": ["Impose administrative sanctions"]
2085 5  }

```

## F.2 GPT-4O-BASED AGENT WITH REACT AND “SINGLE FUNCTION” ACTION

### F.2.1 QUERY DETAILS

- **Query Quadruplet:** (2023-11-03, AUS, ?, CHN)
- **Temporal Distance:** 1; therefore, the current date is 2023-11-02
- **Agent Max Steps:** 20

### F.2.2 QUERY PROMPT

Please forecast the relations that **Australia** will take towards **China** on **November 03, 2023** based on historical information up to **November 02, 2023**. I.e. forecast the relation CAMEO codes in query event `Event(date=2023-11-03, head_entity=ISOCCode(AUS), relation=CAMEOCode(?), tail_entity=ISOCCode(CHN))`.

### F.2.3 SYSTEM PROMPT

You are an expert in forecasting future events based on historical data. The database contains news articles from January 1, 2023 to the current date **November 02, 2023** and the events extracted from these articles. The events are in the form of (date, subject country, relation, object country), where the countries are represented by ISO 3166-1 alpha-3 codes and the relations are represented by the CAMEO codes defined in the "Conflict and Mediation Event Observations" ontology. The relations are hierarchical: first-level relations are general parent relations represented by two-digit CAMEO codes, while second-level relations are more specific child relations represented by three-digit CAMEO codes. Child relations have the same first two digits as their parent relations. For example, "01" is a first-level relation, and "010" and "011" are some of its second-level relations. The relations

2106 in the database are represented in the second-level form.  
2107

2108 Your task is to forecast the future relations between two entities in a given query. You  
2109 have access to a defined Python API that allows you to query the database for historical events and  
2110 statistics, and to get precise information about the ISO country codes and CAMEO relation codes.  
2111

2112 The defined API is described as follows:

```
2113 ``` python
2114 {api_description}
2115 ```
```

2116 You will use an iterative approach, interleaving "Thought", "Action", and "Observation"  
2117 steps to collect information and perform the forecast. You may perform up to 20 iterations. The steps  
2118 are as follows:  
2119

- 2120 - "Thought": Analyze the current information and reason about the current situation, and  
2121 predicts which API you want to use (try to use different APIs to collect diverse information) or make  
2122 a decision that you want to make a final answer.
- 2123 - "Action": Use the API to gather more information or provide the final forecast.  
2124 - If using the API: the action must be only one single line of exactly one function call from the  
2125 API with appropriate inputs, without additional code, explanations, or natural language descriptions.  
2126 - If making the final forecast: the action must start immediately with "Final Answer:", and follow  
2127 with the results in the expected JSON format.
- 2128 - "Observation": Return the output of the called function.

2129 To make a reasonable forecast, you should collect both news and relational evidence to  
2130 support your prediction. When you are fully confident that you accumulate enough information  
2131 to make the final forecast, you should start the "Thought" with your reasoning using the news  
2132 and structural information to make the prediction, and then start the "Action" step with "Final  
2133 Answer:" followed by the answer in the expected JSON format. The answer should be a  
2134 JSON dictionary where the keys are the forecasted two-digit first-level CAMEO codes and  
2135 the values are lists of forecasted three-digit second-level CAMEO codes that are child relations  
2136 of the key. For example, "Action: Final Answer: {"01": ["010", "011", "012"], "02": ["020", "023"]}".  
2137

2138 The final answer will be evaluated based on the precision and recall of the forecasted first-  
2139 level and second-level relations, so only include confident first-level and second-level CAMEO codes  
2140 in your final forecast.

2141 Try to use different APIs to collect diverse information (including multi-hop relations),  
2142 such as the precise meaning of CAMEO codes, insights from news content, relational data, and  
2143 statistical analyses to support your forecasts. Consider not only the frequency of the relations but also  
2144 the temporal aspects of the data when making your forecast.

#### 2145 F.2.4 AGENT REACT PROCESS

2146 1. **Thought:** To forecast the relations that Australia (AUS) will take towards China (CHN) on  
2147 November 03, 2023, I will start by gathering historical data on the interactions between these two  
2148 countries. This will include the types of relations, their frequencies, and the context in which they  
2149 occurred. I will use the following steps:

- 2150 (a) Retrieve the distribution of second-level relations between Australia and China.
- 2151 (b) Analyze the distribution to identify the most common types of interactions.
- 2152 (c) Examine recent news articles to understand the context and any ongoing trends or events that  
2153 might influence future relations.  
2154

2155 **Action:**

```
2156 1 get_relation_distribution(head_entities=[ISOCODE("AUS")], tail_entities=[ISOCODE("CHN")])
```

2157 **Observation:**

```
2158 1 {CAMEOCODE (code="042"): 32,  
2159 2 CAMEOCODE (code="036"): 30,
```



```

2160 3 CAMEOCode (code="040"): 17,
2161 4 CAMEOCode (code="061"): 16,
2162 5 CAMEOCode (code="046"): 12,
2163 6 CAMEOCode (code="043"): 11,
2164 7 CAMEOCode (code="010"): 6,
2165 8 CAMEOCode (code="020"): 5,
2166 9 CAMEOCode (code="173"): 5,
2167 10 CAMEOCode (code="012"): 5,
2168 11 CAMEOCode (code="051"): 4,
2169 12 CAMEOCode (code="163"): 3,
2170 13 CAMEOCode (code="128"): 2,
2171 14 CAMEOCode (code="050"): 2,
2172 15 CAMEOCode (code="084"): 2,
2173 16 CAMEOCode (code="090"): 2,
2174 17 CAMEOCode (code="174"): 1,
2175 18 CAMEOCode (code="112"): 1,
2176 19 CAMEOCode (code="130"): 1,
2177 20 CAMEOCode (code="160"): 1,
2178 21 CAMEOCode (code="172"): 1,
2179 22 CAMEOCode (code="192"): 1,
2180 23 CAMEOCode (code="190"): 1,
2181 24 CAMEOCode (code="060"): 1,
2182 25 CAMEOCode (code="030"): 1,
2183 26 CAMEOCode (code="035"): 1,
2184 27 CAMEOCode (code="138"): 1,
2185 28 CAMEOCode (code="071"): 1,
2186 29 CAMEOCode (code="031"): 1}

```

2. **Thought:** The distribution of second-level relations between Australia and China shows a variety of interactions, with some relations being more frequent than others. The most common interactions include:

- "042" (Consult): 32 occurrences
- "036" (Express intent to cooperate): 30 occurrences
- "040" (Assure): 17 occurrences
- "061" (Cooperate economically): 16 occurrences
- "046" (Praise or endorse): 12 occurrences
- "043" (Apologize): 11 occurrences

These frequent interactions suggest a trend towards diplomatic and cooperative engagements. However, to make a more informed forecast, it is essential to understand the context and any recent developments that might influence future relations.

Next, I will retrieve recent news articles involving Australia and China to gather insights on the current context and any ongoing trends or events.

**Action:**

```

2195 1 get_news_articles(date_range=DateRange(start_date=Date("2023-10-01"), end_date=Date("
2196 2023-11-02")), head_entities=[ISOCODE("AUS")], tail_entities=[ISOCODE("CHN")])

```

**Observation:**

```

2199 1 [(Date(date="2023-11-02"), "No quick fix to remaining tariffs ahead of China visit"),
2200 2 (Date(date="2023-11-02"), "China's aggression in Indo-Pacific leads to "concern and angst"
2201 3 : Peter Dutton, Australia's leader of opposition"),
2202 4 (Date(date="2023-11-02"), "Australia and China eye new ways to heal old wounds"),
2203 5 (Date(date="2023-11-02"), "Albanese must aim beyond bilateral outcomes on Beijing visit"),
2204 6 (Date(date="2023-11-02"), "Western Australian Premier Announces Visit to China"),
2205 7 (Date(date="2023-11-02"), "Australian PM's China visit will help normalize relations, says
2206 8 ex-trade minister"),
2207 9 (Date(date="2023-11-02"), "Nicholas Moore on how to get in early on the next big growth
2208 10 story"),
2209 11 (Date(date="2023-11-02"), "PM's China trip cements new era for ties"),
2210 12 (Date(date="2023-11-02"), "PM to take a hard line on trade with China"),
2211 13 (Date(date="2023-11-01"), "Australian Prime Minister Anthony Albanese to raise imprisoned
2212 14 democracy blogger during China visit")]

```

3. **Thought:** The recent news articles indicate a mix of diplomatic, economic, and security-related interactions between Australia and China. Key themes include:

- Diplomatic engagements and efforts to normalize relations, as seen in articles about the Australian Prime Minister's visit to China and efforts to heal old wounds.
- Economic cooperation and trade discussions, including addressing tariffs and trade issues.

2214 - Security concerns and criticisms, particularly related to China's actions in the Indo-Pacific  
 2215 region and human rights issues.  
 2216

2217 Given these themes, it is likely that the future relations between Australia and China will involve  
 2218 a combination of diplomatic assurances, economic cooperation, and some critical statements or  
 2219 actions related to security and human rights.

2220 Next, I will retrieve detailed descriptions of the most frequent second-level relations to ensure  
 2221 accurate mapping and understanding of these interactions.

2222 **Action:**

```
2223 1 map_cameo_to_relation(CAMEOCode("042"))  
2224
```

2225 **Observation:**

```
2226 1 Relation(cameo_code=CAMEOCode("042"), name="Make a visit", description="Travel to another  
2227 location for a meeting or other event.")  
2228
```

2229 4. **Thought:** The relation "042" (Make a visit) involves traveling to another location for a meeting  
 2230 or other event, which aligns with the recent news about the Australian Prime Minister's visit to  
 2231 China. This suggests that diplomatic visits are a significant aspect of the interactions between  
 2232 Australia and China.

2233 To further refine the forecast, I will retrieve detailed descriptions for the other most frequent  
 2234 second-level relations: "036", "040", "061", "046", and "043".

2235 **Action:**

```
2236 1 map_cameo_to_relation(CAMEOCode("036"))  
2237
```

2238 **Observation:**

```
2239 1 Relation(cameo_code=CAMEOCode("036"), name="Express intent to meet or negotiate",  
2240 description="Offer, promise, agree to, or otherwise indicate willingness or commitment  
2241 to meet, visit, or engage in talks or negotiations.")  
2242
```

2243 5. **Thought:** The relation "036" (Express intent to meet or negotiate) involves indicating a willingness  
 2244 or commitment to meet, visit, or engage in talks or negotiations. This aligns with the ongoing  
 2245 diplomatic efforts and discussions mentioned in the recent news articles.

2246 Next, I will retrieve detailed descriptions for the relations "040", "061", "046", and "043" to further  
 2247 understand the context of these interactions.

2248 **Action:**

```
2249 1 get_parent_relation(CAMEOCode("040"))  
2250
```

2251 **Observation:**

```
2252 1 Relation(cameo_code=CAMEOCode("04"), name="Consult", description="All consultations and  
2253 meetings.")  
2254
```

2255 6. **Thought:** The parent relation '04' (Consult) encompasses all consultations and meetings, which  
 2256 is consistent with the frequent diplomatic interactions between Australia and China. The specific  
 2257 relation "040" (Assure) likely involves providing assurances during these consultations.

2258 To further refine the forecast, I will retrieve detailed descriptions for the relations "061", "046",  
 2259 and "043".

2260 **Action:""**

2261 **Observation:** No feedback from the environment due to the null action.

2262 7. **Thought:** I will continue retrieving detailed descriptions for the relations "061", "046", and "043"  
 2263 to further understand the context of these interactions.

2264 **Action:""**

2265 **Observation:** No feedback from the environment due to the null action.

2266 8. **Thought:** It seems there was an issue with my previous actions. I will now proceed to retrieve the  
 2267 detailed descriptions for the remaining frequent second-level relations: "061", "046", and "043".

**Action:**

2268

2269 

```
1 map_cameo_to_relation(CAMEOCode("061"))
```

2270

**Observation:**

2272

2273 

```
1 Relation(cameo_code=CAMEOCode("061"), name="Cooperate economically", description="Initiate, resume, improve, or expand economic exchange or cooperation.")
```

2274

9. **Thought:** The relation "061" (Cooperate economically) involves initiating, resuming, improving, or expanding economic exchange or cooperation. This aligns with the recent news articles discussing trade issues and economic cooperation between Australia and China.

2276

Next, I will retrieve detailed descriptions for the relations "046" and "043".

2277

**Action:**

2279

2280 

```
1 map_cameo_to_relation(CAMEOCode("046"))
```

2281

**Observation:**

2283

2284 

```
1 Relation(cameo_code=CAMEOCode("046"), name="Engage in negotiation", description="Negotiate or bargain with others.")
```

2285

2286

10. **Thought:** The relation "046" (Engage in negotiation) involves negotiating or bargaining with others. This is consistent with the ongoing discussions and negotiations mentioned in the news articles, particularly regarding trade and diplomatic issues.

2287

Next, I will retrieve the detailed description for the relation "043".

2288

**Action:""**

2289

**Observation:** No feedback from the environment due to the null action.

2290

2291

11. **Thought:** It seems there was an issue with my previous actions. I will now proceed to retrieve the detailed description for the relation "043" to complete the understanding of the most frequent interactions between Australia and China.

2292

**Action:**

2293

2294 

```
1 map_cameo_to_relation(CAMEOCode("043"))
```

2295

**Observation:**

2296

2297 

```
1 Relation(cameo_code=CAMEOCode("043"), name="Host a visit", description="Host or receive a visitor at residence, office or home country.")
```

2298

2299

12. **Thought:** The relation "043" (Host a visit) involves hosting or receiving a visitor at a residence, office, or home country. This aligns with the diplomatic visits and engagements mentioned in the recent news articles.

2300

2301

2302

2303

2304

2305

2306

2307

2308

2309

2310

2311

2312

2313

2314

2315

2316

2317

2318

2319

2320

2321

#### F.2.5 AGENT FINAL STATUS

- **End State:** Final Answer

- **Number of Steps Taken:** 12

- **Final Answer:**

```
1 {
2   "03": ["036"],
3   "04": ["040", "042", "043", "046"],
4   "06": ["061"]
5 }
```

2322 We show the corresponding relation names of the predicted codes here for a better review:  
2323

```
2324 {
2325   "Express intent to cooperate": ["Express intent to meet or negotiate"],
2326   "Consult": ["Consult, not specified", "Make a visit", "Host a visit", "Engage in
2327     negotiation"],
2328   "Engage in material cooperation": ["Cooperate economically"]
2329 }
```

2329 • **Ground Truth Answer:**

```
2330 {
2331   "03": ["036"],
2332   "04": ["042"],
2333   "17": ["172"]
2334 }
```

2334 We show the corresponding relation names of the ground-truth codes here for a better review:  
2335

```
2336 {
2337   "Express intent to cooperate": ["Express intent to meet or negotiate"],
2338   "Consult": ["Make a visit"],
2339   "Coerce": ["Impose administrative sanctions"]
2340 }
```

### 2340 F.3 GPT-4O-MINI-BASED AGENT WITH REACT AND “CODE BLOCK” ACTION

#### 2341 F.3.1 QUERY DETAILS

- 2342 • **Query Quadruplet:** (2024-02-08, PNG, ?, AUS)
- 2343 • **Temporal Distance:** 1; therefore, the current date is 2024-02-07
- 2344 • **Agent Max Steps:** 20

#### 2345 F.3.2 QUERY PROMPT

2346 Please forecast the relations that Papua New Guinea will take towards Australia on February 08, 2024  
2347 based on historical information up to February 07, 2024. I.e. forecast the relation CAMEO codes  
2348 in query event Event(date=2024-02-08, head\_entity=ISOCODE(PNG), relation=CAMEOCODE(?),  
2349 tail\_entity=ISOCODE(AUS)).

#### 2350 F.3.3 SYSTEM PROMPT

2351 You are an expert in forecasting future events based on historical data. The database contains news  
2352 articles from January 1, 2023 to the current date February 07, 2024 and the events extracted from  
2353 these articles. The events are in the form of (date, subject country, relation, object country), where  
2354 the countries are represented by ISO 3166-1 alpha-3 codes and the relations are represented by the  
2355 CAMEO codes defined in the "Conflict and Mediation Event Observations" ontology. The relations  
2356 are hierarchical: first-level relations are general parent relations represented by two-digit CAMEO  
2357 codes, while second-level relations are more specific child relations represented by three-digit  
2358 CAMEO codes. Child relations have the same first two digits as their parent relations. For example,  
2359 "01" is a first-level relation, and "010" and "011" are some of its second-level relations. The relations  
2360 in the database are represented in the second-level form.

2361 Your task is to forecast the future relations between two entities in a given query. You  
2362 have access to a defined Python API that allows you to query the database for historical events and  
2363 statistics, and to get precise information about the ISO country codes and CAMEO relation codes.  
2364 You are also authorized to utilize additional safe, well-established Python libraries such as numpy,  
2365 pandas, scikit-learn, and NetworkX to enhance your data analysis and forecasting accuracy.

2366 The defined API is described as follows:

```
2367 ``` python
2368 {api_description}
2369 ```
```

2370 You will use an iterative approach, interleaving "Thought", "Action", and "Observation"  
2371 steps to collect information and perform the forecast. You may perform up to 20 iterations. The steps  
2372 are as follows:

- 2376
- 2377 - "Thought": Analyze the current information and reason about the current situation, and
- 2378 predicts which API you want to use (try to use different APIs to collect diverse information) or make
- 2379 a decision that you want to make a final answer.
- 2380 - "Action": Use the API to gather more information or provide the final forecast.
- 2381 - If gathering more data: the action must be an executable Python code snippet that starts with
- 2382 `python` and ends with `````. It can contain multiple lines of codes and function calls using the
- 2383 defined API or Python libraries. You must use `print()` to output the results, and only the printed output
- 2384 will be returned in the observation step.
- 2385 - If making the final forecast: the action must start immediately with "Final Answer:", and follow
- 2386 with the answer in the expected JSON format. This should not be enclosed within triple backticks.
- 2387 - "Observation": Return the printed output of the executed code snippet.

2388 To make a reasonable forecast, you should collect both news and relational evidence to

2389 support your prediction. When you are fully confident that you accumulate enough information

2390 to make the final forecast, you should start the "Thought" with your reasoning using the news

2391 and structural information to make the prediction, and then start the "Action" step with "Final

2392 Answer:" followed by the answer in the expected JSON format. The answer should be a

2393 JSON dictionary where the keys are the forecasted two-digit first-level CAMEO codes and

2394 the values are lists of forecasted three-digit second-level CAMEO codes that are child relations

2395 of the key. For example, "Action: Final Answer: {"01": ["010", "011", "012"], "02": ["020", "023"]}".

2396 The final answer will be evaluated based on the precision and recall of the forecasted first-

2397 level and second-level relations, so only include confident first-level and second-level CAMEO codes

2398 in your final forecast.

2400 Try to use different APIs and Python libraries to collect diverse information (including

2401 multi-hop relations), such as the precise meaning of CAMEO codes, insights from news content,

2402 relational data, and statistical analyses to support your forecasts. Consider not only the frequency of

2403 the relations but also the temporal aspects of the data when making your forecast.

#### 2404 F.3.4 AGENT REACT PROCESS

2405

2406 1. **Thought:** To forecast the relations that Papua New Guinea (PNG) will take towards Australia

2407 (AUS) on February 08, 2024, I need to gather historical data on events involving these two

2408 countries. Specifically, I will look for patterns in the relations between PNG and AUS in the past

2409 events to understand the likely future interactions.

2410 First, I will check the distribution of relations between PNG and AUS over the past year to identify

2411 any prevalent patterns or trends. This will help me understand the types of relations that have been

2412 most common and could indicate future behavior.

2413 I will use the 'get\_relation\_distribution' API to get this information.

2414 **Action:**

```
2415 1 # Define the date range for the past year up to February 07, 2024
2416 2 date_range = DateRange(start_date=Date("2023-02-08"), end_date=Date("2024-02-07"))
2417 3 head_entities = [ISOCCode("PNG")] # Papua New Guinea
2418 4 tail_entities = [ISOCCode("AUS")] # Australia
2419 5
2420 6 # Get the distribution of relations between PNG and AUS
2421 7 relation_distribution = get_relation_distribution(date_range=date_range, head_entities=
2422 8 head_entities, tail_entities=tail_entities)
2423 8 print(relation_distribution)
```

2422 **Observation:**

```
2423 1 {CAMEOCode(code='042'): 8, CAMEOCode(code='036'): 7, CAMEOCode(code='043'): 4, CAMEOCode(
2424 2 code='057'): 4, CAMEOCode(code='030'): 2, CAMEOCode(code='010'): 2, CAMEOCode(code='
2425 3 181'): 1, CAMEOCode(code='051'): 1, CAMEOCode(code='046'): 1, CAMEOCode(code='055'):
2426 4, CAMEOCode(code='130'): 1, CAMEOCode(code='160'): 1}
```

2427 2. **Thought:** The relation distribution between Papua New Guinea (PNG) and Australia (AUS) over

2428 the past year reveals several notable patterns:

- 2429 • **CAMEO Code '042':** This relation appears 8 times, which suggests a strong tendency for PNG to engage in "Make a public statement" towards Australia.



- **CAMEO Code '036'**: With 7 occurrences, this indicates a frequent pattern of "Negotiate" relations.
- **CAMEO Code '043'**: This relation, which signifies "Make a public statement, not specified," has been observed 4 times, indicating a general trend in public communication.
- **CAMEO Code '057'**: Also appearing 4 times, this suggests some level of "Express intent to cooperate," which could be significant for future interactions.
- Other codes like '030' (Make a statement, not specified) and '010' (Make statement, not specified) also show minor interactions.

Given the high frequency of relations involving public statements and negotiations, it is likely that PNG will continue to engage in similar relations with Australia in the near future.

Next, I will analyze the context of news articles related to PNG and AUS to see if there are any recent developments or trends that could influence their future relations. This will provide additional insights into the nature of their interactions.

I will use the 'get\_news\_articles' API to gather relevant news articles from the past month that mention both PNG and AUS. This will help me understand current sentiments or events that could impact their relationship on February 08, 2024.

#### Action:

```

1 # Define the date range for the past month leading up to February 07, 2024
2 date_range = DateRange(start_date=Date("2024-01-07"), end_date=Date("2024-02-07"))
3 head_entities = [ISOCODE("PNG")] # Papua New Guinea
4 tail_entities = [ISOCODE("AUS")] # Australia
5
6 # Get news articles that mention both PNG and AUS
7 news_articles = get_news_articles(date_range=date_range, head_entities=head_entities,
8                                 tail_entities=tail_entities)
9 print(news_articles)

```

#### Observation:

```

1 [(Date(date='2024-02-07'), 'Papua New Guineas PM to Visit Australian Parliament, Address
   Beijing's Influence'),
2 (Date(date='2024-02-07'), 'PNG wants to stand on its 'own two feet': PM Marape'),
3 (Date(date='2024-02-07'), 'Papua New Guineas PM to address Australian parliament as
   Pacific security race with China builds'),
4 (Date(date='2024-01-11'), 'Travel warning issued as 'widespread public disorder' rages in
   Port Moresby')]

```

3. **Thought:** The recent news articles regarding Papua New Guinea (PNG) and Australia (AUS) indicate significant developments that could influence their relations on February 08, 2024. Notably, the Prime Minister of PNG is scheduled to visit the Australian Parliament and address issues related to Beijing's influence in the region. This suggests a proactive approach by PNG to strengthen ties with Australia, particularly in the context of regional security and cooperation.

The repeated mention of "widespread public disorder" in Port Moresby could indicate internal challenges for PNG, which may affect its diplomatic posture. However, the focus on addressing Australian Parliament implies that PNG is keen on maintaining a dialogue with Australia and possibly negotiating support or partnerships.

Given the historical data showing a strong inclination towards public statements and negotiations (CAMEO Codes '042' and '036'), along with the current political context, it is reasonable to forecast that PNG will likely engage in similar relations with Australia.

Based on the analysis of both the historical event distribution and the recent news articles, I will forecast the following CAMEO codes for the relation that Papua New Guinea will take towards Australia on February 08, 2024:

- **CAMEO Code '042'**: Make a public statement (high frequency in historical data).
- **CAMEO Code '036'**: Negotiate (also high frequency).
- **CAMEO Code '057'**: Express intent to cooperate (indicated by recent news context).

#### Action:

```

1 Final Answer: {"04": ["042"], "05": ["057"], "03": ["036"]}

```

## 2484 F.3.5 AGENT FINAL STATUS

- 2485 • **End State:** Final Answer
- 2486 • **Number of Steps Taken:** 3
- 2487 • **Final Answer:**

```
2488 {
2489   "03": ["036"],
2490   "04": ["042"],
2491   "05": ["057"]
2492 }
```

2493 We show the corresponding relation names of the predicted codes here for a better review:

```
2494 {
2495   "Express intent to cooperate": ["Express intent to meet or negotiate"],
2496   "Consult": ["Make a visit"],
2497   "Engage in diplomatic cooperation": ["Sign formal agreement"]
2498 }
```

- 2499 • **Ground Truth Answer:**

```
2500 {
2501   "01": ["010"],
2502   "05": ["051"]
2503 }
```

2504 We show the corresponding relation names of the ground-truth codes here for a better review:

```
2505 {
2506   "Make public statement": ["Make statement, not specified"],
2507   "Engage in diplomatic cooperation": ["Praise or endorse"]
2508 }
```

## 2510 F.4 GPT-4O-MINI-BASED AGENT WITH REACT AND “SINGLE FUNCTION” ACTION

## 2511 F.4.1 QUERY DETAILS

- 2512 • **Query Quadruplet:** (2024-02-08, PNG, ?, AUS)
- 2513 • **Temporal Distance:** 1; therefore, the current date is 2024-02-07
- 2514 • **Agent Max Steps:** 20

## 2515 F.4.2 QUERY PROMPT

2516 Please forecast the relations that Papua New Guinea will take towards Australia on February 08, 2024  
 2517 based on historical information up to February 07, 2024. I.e. forecast the relation CAMEO codes  
 2518 in query event Event(date=2024-02-08, head\_entity=ISOCODE(PNG), relation=CAMEOCODE(?),  
 2519 tail\_entity=ISOCODE(AUS)).

## 2522 F.4.3 SYSTEM PROMPT

2523 You are an expert in forecasting future events based on historical data. The database contains news  
 2524 articles from January 1, 2023 to the current date February 07, 2024 and the events extracted from  
 2525 these articles. The events are in the form of (date, subject country, relation, object country), where  
 2526 the countries are represented by ISO 3166-1 alpha-3 codes and the relations are represented by the  
 2527 CAMEO codes defined in the "Conflict and Mediation Event Observations" ontology. The relations  
 2528 are hierarchical: first-level relations are general parent relations represented by two-digit CAMEO  
 2529 codes, while second-level relations are more specific child relations represented by three-digit  
 2530 CAMEO codes. Child relations have the same first two digits as their parent relations. For example,  
 2531 "01" is a first-level relation, and "010" and "011" are some of its second-level relations. The relations  
 2532 in the database are represented in the second-level form.

2533 Your task is to forecast the future relations between two entities in a given query. You  
 2534 have access to a defined Python API that allows you to query the database for historical events and  
 2535 statistics, and to get precise information about the ISO country codes and CAMEO relation codes.

2536 The defined API is described as follows:

```
2537 ``` python
```

2538 {api\_description}  
2539 ```

2540

2541 You will use an iterative approach, interleaving "Thought", "Action", and "Observation"  
2542 steps to collect information and perform the forecast. You may perform up to 20 iterations. The steps  
2543 are as follows:

2544

2545 - "Thought": Analyze the current information and reason about the current situation, and  
2546 predicts which API you want to use (try to use different APIs to collect diverse information) or make  
2547 a decision that you want to make a final answer.

2548

2549 - "Action": Use the API to gather more information or provide the final forecast.  
2550 - If using the API: the action must be only one single line of exactly one function call from the  
2551 API with appropriate inputs, without additional code, explanations, or natural language descriptions.

2552

2553 - If making the final forecast: the action must start immediately with "Final Answer:", and follow  
2554 with the results in the expected JSON format.

2555

2556 - "Observation": Return the output of the called function.

2557

2558 To make a reasonable forecast, you should collect both news and relational evidence to  
2559 support your prediction. When you are fully confident that you accumulate enough information  
2560 to make the final forecast, you should start the "Thought" with your reasoning using the news  
2561 and structural information to make the prediction, and then start the "Action" step with "Final  
2562 Answer:" followed by the answer in the expected JSON format. The answer should be a  
2563 JSON dictionary where the keys are the forecasted two-digit first-level CAMEO codes and  
2564 the values are lists of forecasted three-digit second-level CAMEO codes that are child relations  
2565 of the key. For example, "Action: Final Answer: {"01": ["010", "011", "012"], "02": ["020", "023"]}".

2566

2567 The final answer will be evaluated based on the precision and recall of the forecasted first-  
2568 level and second-level relations, so only include confident first-level and second-level CAMEO codes  
2569 in your final forecast.

2570

2571 Try to use different APIs to collect diverse information (including multi-hop relations),  
2572 such as the precise meaning of CAMEO codes, insights from news content, relational data, and  
2573 statistical analyses to support your forecasts. Consider not only the frequency of the relations but also  
2574 the temporal aspects of the data when making your forecast.

2575

#### 2576 F.4.4 AGENT REACT PROCESS

2577

2578 1. **Thought:** To forecast the relations between Papua New Guinea (PNG) and Australia (AUS) for  
2579 February 08, 2024, I need to gather historical data on events involving these two countries. I will  
2580 start by counting the number of events that have occurred between PNG and AUS in the past  
2581 year to see which relations have been most prevalent. This will help me identify potential future  
2582 relations.

2583

2584 **Action:**

2585

```
2586 1 count_events(date_range=DateRange(start_date=Date("2023-02-08"), end_date=Date("2024-02-07"  
2587 )), head_entities=[ISOCODE("PNG")], tail_entities=[ISOCODE("AUS")])
```

2588

2589 **Observation:**

2590

```
2591 1 33
```

2592

2593 2. **Thought:** The count of 33 events indicates that there has been a significant number of interactions  
2594 between Papua New Guinea and Australia over the past year. Next, I will analyze the distribution  
2595 of relations between these two countries to understand which specific CAMEO relations are most  
2596 common. This will help me identify likely future relations.

2597

2598 **Action:**

2599

```
2600 1 get_relation_distribution(date_range=DateRange(start_date=Date("2023-02-08"), end_date=Date  
2601 ("2024-02-07")), head_entities=[ISOCODE("PNG")], tail_entities=[ISOCODE("AUS")])
```

2602

2603 **Observation:**

2592

2593

2594

2595

2596

2597

2598

2599

2600

2601

2602

2603

2604

2605

2606

2607

2608

2609

2610

2611

2612

2613

2614

2615

2616

2617

2618

2619

2620

2621

2622

2623

2624

2625

2626

2627

2628

2629

2630

2631

2632

2633

2634

2635

2636

2637

2638

2639

2640

2641

2642

2643

2644

2645

```
1 {CAMEOCode(code='042'): 8, CAMEOCode(code='036'): 7, CAMEOCode(code='043'): 4, CAMEOCode(
  code='057'): 4, CAMEOCode(code='030'): 2, CAMEOCode(code='010'): 2, CAMEOCode(code='
  181'): 1, CAMEOCode(code='051'): 1, CAMEOCode(code='046'): 1, CAMEOCode(code='055'):
  1, CAMEOCode(code='130'): 1, CAMEOCode(code='160'): 1}
```

3. **Thought:** The distribution of relations between Papua New Guinea and Australia shows that the most common relations are as follows:

- **CAMEO Code '042':** 8 events (likely related to economic agreements or cooperation).
- **CAMEO Code '036':** 7 events (likely related to diplomatic relations).
- **CAMEO Code '043':** 4 events (possibly related to trade).
- **CAMEO Code '057':** 4 events (possibly related to military cooperation).

Given the frequency of these relations, it is likely that Papua New Guinea will continue to engage in similar interactions with Australia. To further refine my forecast, I will analyze the specific descriptions of these relations to understand their context better and see if any recent news articles provide insights into future interactions.

**Action:**

```
1 get_relation_distribution(date_range=DateRange(start_date=Date("2023-01-01"), end_date=Date(
  "2024-02-07")), head_entities=[ISOCODE("PNG")], tail_entities=[ISOCODE("AUS")])
```

**Observation:**

```
1 {CAMEOCode(code='042'): 8, CAMEOCode(code='036'): 7, CAMEOCode(code='043'): 5, CAMEOCode(
  code='057'): 4, CAMEOCode(code='030'): 2, CAMEOCode(code='010'): 2, CAMEOCode(code='
  181'): 1, CAMEOCode(code='051'): 1, CAMEOCode(code='046'): 1, CAMEOCode(code='055'):
  1, CAMEOCode(code='130'): 1, CAMEOCode(code='160'): 1}
```

4. **Thought:** The distribution of relations remains consistent with the previous observation, indicating that the most common interactions between Papua New Guinea and Australia are primarily focused on economic agreements (CAMEOCode '042') and diplomatic relations (CAMEOCode '036').

To further understand the context of these relations, I will retrieve recent news articles that mention Papua New Guinea and Australia to see if there are any specific events or agreements that could influence the relations on February 08, 2024. This will help me make a more informed forecast.

**Action:**

```
1 get_news_articles(date_range=DateRange(start_date=Date("2023-01-01"), end_date=Date("
  2024-02-07")), head_entities=[ISOCODE("PNG")], tail_entities=[ISOCODE("AUS")])
```

**Observation:**

```
1 [(Date(date='2024-02-07'), 'Papua New Guineas PM to Visit Australian Parliament, Address
  Beijings Influence'),
  2 (Date(date='2024-02-07'), 'PNG wants to stand on its 'own two feet': PM Marape"),
  3 (Date(date='2024-02-07'), 'Papua New Guineas PM to address Australian parliament as
  Pacific security race with China builds'),
  4 (Date(date='2024-01-11'), 'Travel warning issued as 'widespread public disorder' rages in
  Port Moresby"),
  5 (Date(date='2024-01-11'), 'Papua New Guinea declares state of emergency after deadly riots
  ')]
```

5. **Thought:** The recent news articles indicate that Papua New Guinea's Prime Minister is set to visit the Australian Parliament and address concerns regarding China's influence in the Pacific. This suggests a focus on diplomatic relations and security cooperation between PNG and AUS. The emphasis on standing on their "own two feet" also hints at a desire for economic independence and cooperation.

Given the historical data showing a strong presence of economic agreements (CAMEOCode '042') and diplomatic relations (CAMEOCode '036'), along with the context provided by recent news articles, it is reasonable to forecast that Papua New Guinea will likely engage in diplomatic discussions and economic cooperation with Australia on February 08, 2024.

**Action:**

```
1 Final Answer: {"04": ["042", "043"], "03": ["036"], "05": ["057"]}
```

2646 F.4.5 AGENT FINAL STATUS

- 2647 • **End State:** Final Answer
- 2648
- 2649 • **Number of Steps Taken:** 5
- 2650 • **Final Answer:**

```
2651 1 {  
2652 2   "03": ["036"],  
2653 3   "04": ["042", "043"],  
2654 4   "05": ["057"]  
2655 5 }
```

2656 We show the corresponding relation names of the predicted codes here for a better review:

```
2657 1 {  
2658 2   "Express intent to cooperate": ["Express intent to meet or negotiate"],  
2659 3   "Consult": ["Make a visit", "Host a visit"],  
2660 4   "Engage in diplomatic cooperation": ["Sign formal agreement"]  
2661 5 }
```

- 2662 • **Ground Truth Answer:**

```
2663 1 {  
2664 2   "01": ["010"],  
2665 3   "05": ["051"]  
2666 4 }
```

2667 We show the corresponding relation names of the ground-truth codes here for a better review:

```
2668 1 {  
2669 2   "Make public statement": ["Make statement, not specified"],  
2670 3   "Engage in diplomatic cooperation": ["Praise or endorse"]  
2671 4 }
```

2672  
2673  
2674  
2675  
2676  
2677  
2678  
2679  
2680  
2681  
2682  
2683  
2684  
2685  
2686  
2687  
2688  
2689  
2690  
2691  
2692  
2693  
2694  
2695  
2696  
2697  
2698  
2699



## G ADDITIONAL DETAILS OF THE BENCHMARK

## G.1 COUNTRIES AND ISO CODES

We show the ISO 3166 Alpha-3 country codes and names in Table 13 and Table 14.

Table 13: List of Countries and their ISO-3166 Alpha-3 Codes (Part 1)

Code	Country	Code	Country	Code	Country
AND	Andorra	ARE	United Arab Emirates	AFG	Afghanistan
ATG	Antigua and Barbuda	AIA	Anguilla	ALB	Albania
ARM	Armenia	AGO	Angola	ATA	Antarctica
ARG	Argentina	ASM	American Samoa	AUT	Austria
AUS	Australia	ABW	Aruba	ALA	Åland
AZE	Azerbaijan	BIH	Bosnia and Herzegovina	BRB	Barbados
BGD	Bangladesh	BEL	Belgium	BFA	Burkina Faso
BGR	Bulgaria	BHR	Bahrain	BDI	Burundi
BEN	Benin	BLM	Saint Barthélemy	BMU	Bermuda
BRN	Brunei	BOL	Bolivia	BES	Bonaire, Sint Eustatius, and Saba
BRA	Brazil	BHS	Bahamas	BTN	Bhutan
BVT	Bouvet Island	BWA	Botswana	BLR	Belarus
BLZ	Belize	CAN	Canada	CCK	Cocos (Keeling) Islands
COD	DR Congo	CAF	Central African Republic	COG	Congo Republic
CHE	Switzerland	CIV	Ivory Coast	COK	Cook Islands
CHL	Chile	CMR	Cameroon	CHN	China
COL	Colombia	CRI	Costa Rica	CUB	Cuba
CPV	Cabo Verde	CUW	Curaçao	CXR	Christmas Island
CYP	Cyprus	CZE	Czechia	DEU	Germany
DJI	Djibouti	DNK	Denmark	DMA	Dominica
DOM	Dominican Republic	DZA	Algeria	ECU	Ecuador
EST	Estonia	EGY	Egypt	ESH	Western Sahara
ERI	Eritrea	ESP	Spain	ETH	Ethiopia
FIN	Finland	FJI	Fiji	FLK	Falkland Islands
FSM	Micronesia	FRO	Faroe Islands	FRA	France
GAB	Gabon	GBR	United Kingdom	GRD	Grenada
GEO	Georgia	GUF	French Guiana	GGY	Guernsey
GHA	Ghana	GIB	Gibraltar	GRL	Greenland
GMB	The Gambia	GIN	Guinea	GLP	Guadeloupe
GNQ	Equatorial Guinea	GRC	Greece	SGS	South Georgia and South Sandwich Islands
GTM	Guatemala	GUM	Guam	GNB	Guinea-Bissau
GUY	Guyana	HKG	Hong Kong	HMD	Heard and McDonald Islands
HND	Honduras	HRV	Croatia	HTI	Haiti
HUN	Hungary	IDN	Indonesia	IRL	Ireland
ISR	Israel	IMN	Isle of Man	IND	India
IOT	British Indian Ocean Territory	IRQ	Iraq	IRN	Iran
ISL	Iceland	ITA	Italy	JEY	Jersey
JAM	Jamaica	JOR	Jordan	JPN	Japan
KEN	Kenya	KGZ	Kyrgyzstan	KHM	Cambodia
KIR	Kiribati	COM	Comoros	KNA	St Kitts and Nevis
PRK	North Korea	KOR	South Korea	KWT	Kuwait
CYM	Cayman Islands	KAZ	Kazakhstan	LAO	Laos
LBN	Lebanon	LCA	Saint Lucia	LIE	Liechtenstein
LKA	Sri Lanka	LBR	Liberia	LSO	Lesotho
LTU	Lithuania	LUX	Luxembourg	LVA	Latvia
LBY	Libya	MAR	Morocco	MCO	Monaco
MDA	Moldova	MNE	Montenegro	MAF	Saint Martin
MDG	Madagascar	MHL	Marshall Islands	MKD	North Macedonia
MLI	Mali	MMR	Myanmar	MNG	Mongolia
MAC	Macao	MNP	Northern Mariana Islands	MTQ	Martinique
MRT	Mauritania	MSR	Montserrat	MLT	Malta
MUS	Mauritius	MDV	Maldives	MWI	Malawi
MEX	Mexico	MYS	Malaysia	MOZ	Mozambique
NAM	Namibia	NCL	New Caledonia	NER	Niger
NFK	Norfolk Island	NGA	Nigeria	NIC	Nicaragua
NLD	The Netherlands	NOR	Norway	NPL	Nepal
NRU	Nauru	NIU	Niue	NZL	New Zealand
OMN	Oman	PAN	Panama	PER	Peru
PYF	French Polynesia	PNG	Papua New Guinea	PHL	Philippines
PAK	Pakistan	POL	Poland	SPM	Saint Pierre and Miquelon
PCN	Pitcairn Islands	PRI	Puerto Rico	PSE	Palestine
PRT	Portugal	PLW	Palau	PRY	Paraguay
QAT	Qatar	REU	Réunion	ROU	Romania
SRB	Serbia	RUS	Russia	RWA	Rwanda
SAU	Saudi Arabia	SLB	Solomon Islands	SYC	Seychelles
SDN	Sudan	SWE	Sweden	SGP	Singapore
SHN	Saint Helena	SVN	Slovenia	SJM	Svalbard and Jan Mayen
SVK	Slovakia	SLE	Sierra Leone	SMR	San Marino

Table 14: List of Countries and their ISO-3166 Alpha-3 Codes (Part 2)

Code	Country	Code	Country	Code	Country
SEN	Senegal	SOM	Somalia	SUR	Suriname
SSD	South Sudan	STP	São Tomé and Príncipe	SLV	El Salvador
SXM	Sint Maarten	SYR	Syria	SWZ	Eswatini
TCA	Turks and Caicos Islands	TCD	Chad	ATF	French Southern Territories
TGO	Togo	THA	Thailand	TJK	Tajikistan
TKL	Tokelau	TLS	Timor-Leste	TKM	Turkmenistan
TUN	Tunisia	TON	Tonga	TUR	Türkiye
TTO	Trinidad and Tobago	TUV	Tuvalu	TWN	Taiwan
TZA	Tanzania	UKR	Ukraine	UGA	Uganda
UMI	U.S. Outlying Islands	USA	United States	URY	Uruguay
UZB	Uzbekistan	VAT	Vatican City	VCT	St Vincent and Grenadines
VEN	Venezuela	VGB	British Virgin Islands	VIR	U.S. Virgin Islands
VNM	Vietnam	VUT	Vanuatu	WLF	Wallis and Futuna
WSM	Samoa	XKX	Kosovo	YEM	Yemen
MYT	Mayotte	ZAF	South Africa	ZMB	Zambia
ZWE	Zimbabwe				

## G.2 RELATIONS AND CAMEO CODES

We show the CAMEO relation codes and names, where the first-level relations are in two digits, and the second-level relations are in three digits with the first two digits be the same as its parent relation:

- 01: Make public statement
  - 010: Make statement, not specified
  - 011: Decline comment
  - 012: Make pessimistic comment
  - 013: Make optimistic comment
  - 014: Consider policy option
  - 015: Acknowledge or claim responsibility
  - 016: Reject accusation or deny responsibility
  - 017: Engage in symbolic act
  - 018: Make empathetic comment
  - 019: Express accord
- 02: Appeal
  - 020: Make an appeal or request, not specified
  - 021: Appeal for material cooperation
  - 022: Appeal for diplomatic cooperation
  - 023: Appeal for material aid
  - 024: Appeal for political reform
  - 025: Appeal to yield
  - 026: Appeal to others to meet or negotiate
  - 027: Appeal to others to settle dispute
  - 028: Appeal to others to engage in or accept mediation
- 03: Express intent to cooperate
  - 030: Express intent to cooperate, not specified
  - 031: Express intent to engage in material cooperation
  - 032: Express intent to engage in diplomatic cooperation
  - 033: Express intent to provide material aid
  - 034: Express intent to institute political reform
  - 035: Express intent to yield
  - 036: Express intent to meet or negotiate
  - 037: Express intent to settle dispute
  - 038: Express intent to accept mediation
  - 039: Express intent to mediate

- 2808 • 04: Consult
- 2809 – 040: Consult, not specified
- 2810 – 041: Discuss by telephone
- 2811 – 042: Make a visit
- 2812 – 043: Host a visit
- 2813 – 044: Meet at a third location
- 2814 – 045: Engage in mediation
- 2815 – 046: Engage in negotiation
- 2816
- 2817 • 05: Engage in diplomatic cooperation
- 2818 – 050: Engage in diplomatic cooperation, not specified
- 2819 – 051: Praise or endorse
- 2820 – 052: Defend verbally
- 2821 – 053: Rally support on behalf of
- 2822 – 054: Grant diplomatic recognition
- 2823 – 055: Apologize
- 2824 – 056: Forgive
- 2825 – 057: Sign formal agreement
- 2826
- 2827 • 06: Engage in material cooperation
- 2828 – 060: Engage in material cooperation, not specified
- 2829 – 061: Cooperate economically
- 2830 – 062: Cooperate militarily
- 2831 – 063: Engage in judicial cooperation
- 2832 – 064: Share intelligence or information
- 2833
- 2834 • 07: Provide aid
- 2835 – 070: Provide aid, not specified
- 2836 – 071: Provide economic aid
- 2837 – 072: Provide military aid
- 2838 – 073: Provide humanitarian aid
- 2839 – 074: Provide military protection or peacekeeping
- 2840 – 075: Grant asylum
- 2841
- 2842 • 08: Yield
- 2843 – 080: Yield, not specified
- 2844 – 081: Ease administrative sanctions
- 2845 – 082: Ease political dissent
- 2846 – 083: Accede to requests or demands for political reform
- 2847 – 084: Return or release
- 2848 – 085: Ease economic sanction or boycott or embargo
- 2849 – 086: Allow international involvement
- 2850 – 087: De-escalate military engagement
- 2851
- 2852 • 09: Investigate
- 2853 – 090: Investigate, not specified
- 2854 – 091: Investigate crime or corruption
- 2855 – 092: Investigate human rights abuses
- 2856 – 093: Investigate military action
- 2857 – 094: Investigate war crimes
- 2858
- 2859 • 10: Demand
- 2860 – 100: Demand, not specified
- 2861 – 101: Demand material cooperation
- 102: Demand for diplomatic cooperation

- 2862 – 103: Demand material aid
- 2863 – 104: Demand political reform
- 2864 – 105: Demand that target yield
- 2865 – 106: Demand meeting or negotiation
- 2866 – 107: Demand settling of dispute
- 2867 – 108: Demand mediation
- 2868
- 2869 • 11: Disapprove
- 2870 – 110: Disapprove, not specified
- 2871 – 111: Criticize or denounce
- 2872 – 112: Accuse
- 2873 – 113: Rally opposition against
- 2874 – 114: Complain officially
- 2875 – 115: Bring lawsuit against
- 2876 – 116: Find guilty or liable (legally)
- 2877
- 2878 • 12: Reject
- 2879 – 120: All rejections and refusals
- 2880 – 121: Reject material cooperation
- 2881 – 122: Reject request or demand for material aid
- 2882 – 123: Reject request or demand for political reform
- 2883 – 124: Refuse to yield
- 2884 – 125: Reject proposal to meet or discuss or negotiate
- 2885 – 126: Reject mediation
- 2886 – 127: Reject plan or agreement to settle dispute
- 2887 – 128: Defy norms or law
- 2888 – 129: Veto
- 2889
- 2890 • 13: Threaten
- 2891 – 130: Threaten, not specified
- 2892 – 131: Threaten non-force
- 2893 – 132: Threaten with administrative sanctions
- 2894 – 133: Threaten political dissent
- 2895 – 134: Threaten to halt negotiations
- 2896 – 135: Threaten to halt mediation
- 2897 – 136: Threaten to halt international involvement
- 2898 – 137: Threaten with repression
- 2899 – 138: Threaten with military force
- 2900 – 139: Give ultimatum
- 2901
- 2902 • 14: Protest
- 2903 – 140: Engage in political dissent, not specified
- 2904 – 141: Demonstrate or rally
- 2905 – 142: Conduct hunger strike
- 2906 – 143: Conduct strike or boycott
- 2907 – 144: Obstruct passage or block
- 2908 – 145: Protest violently or riot
- 2909
- 2910 • 15: Exhibit military posture
- 2911 – 150: Exhibit military or police power, not specified
- 2912 – 151: Increase police alert status
- 2913 – 152: Increase military alert status
- 2914 – 153: Mobilize or increase police power
- 2915 – 154: Mobilize or increase armed forces
- 155: Mobilize or increase cyber-forces

- 2916 • 16: Reduce relations
- 2917 – 160: Reduce relations, not specified
- 2918 – 161: Reduce or break diplomatic relations
- 2919 – 162: Reduce or stop material aid
- 2920 – 163: Impose embargo or boycott or sanctions
- 2921 – 164: Halt negotiations
- 2922 – 165: Halt mediation
- 2923 – 166: Expel or withdraw
- 2924
- 2925
- 2926 • 17: Coerce
- 2927 – 170: Coerce
- 2928 – 171: Seize or damage property
- 2929 – 172: Impose administrative sanctions
- 2930 – 173: Arrest or detain
- 2931 – 174: Expel or deport individuals
- 2932 – 175: Use repression
- 2933 – 176: Attack cybernetically
- 2934
- 2935 • 18: Assault
- 2936 – 180: Use unconventional violence, not specified
- 2937 – 181: Abduct or hijack or take hostage
- 2938 – 182: Physically assault
- 2939 – 183: Conduct suicide or car or other non-military bombing
- 2940 – 184: Use as human shield
- 2941 – 185: Attempt to assassinate
- 2942 – 186: Assassinate
- 2943
- 2944
- 2945 • 19: Fight
- 2946 – 190: Use conventional military force, not specified
- 2947 – 191: Impose blockade or restrict movement
- 2948 – 192: Occupy territory
- 2949 – 193: Fight with small arms and light weapons
- 2950 – 194: Fight with artillery and tanks
- 2951 – 195: Employ aerial weapons
- 2952 – 196: Violate ceasefire
- 2953
- 2954 • 20: Engage in unconventional mass violence
- 2955 – 200: Use massive unconventional force, not specified
- 2956 – 201: Engage in mass expulsion
- 2957 – 202: Engage in mass killings
- 2958 – 203: Engage in ethnic cleansing
- 2959 – 204: Use weapons of mass destruction
- 2960

### 2961 G.3 HUMAN EVALUATION ON THE DATASET QUALITY

## 2962 H ADDITIONAL DETAILS ABOUT AGENT SETUP

2964 **Final answer extraction.** Our agent has two stopping criteria: 1) The agent makes the final answer. 2)  
2965 The reasoning process fails to reach an answer, and ends with consecutive invalid actions, consecutive  
2966 repetitive actions, or exceeded max iterations. We define and analyze this final status in Appendix [D.1](#).  
2967 Afterwards, we always perform an answer extraction step, which is performed by GPT-3.5-Turbo. If  
2968 the agent has not generated a final answer during ReAct, we instruct the answer extraction model to  
2969 make a prediction based on the ReAct reasoning trace. The prompt for answer extraction is shown in  
Appendix [I.3](#).



## 2970 I PROMPTS

2971

### 2972 I.1 SYSTEM PROMPTS

2973

2974

2975

2976

2977

2978

2979

2980

2981

2982

2983

2984

2985

2986

2987

2988

2989

2990

2991

2992

2993

2994

2995

2996

2997

2998

2999

3000

3001

3002

3003

3004

3005

3006

3007

3008

3009

3010

3011

3012

3013

3014

3015

3016

3017

3018

3019

3020

3021

3022

3023

#### I.1.1 SYSTEM PROMPT FOR REACT AGENT WITH ACTION TYPE AS “SINGLE FUNCTION”

##### System prompt for ReAct agent with action type as single function (part 1)

You are an expert in forecasting future events based on historical data. The database contains news articles from January 1, 2023 to the current date {current\_date\_nlp} and the events extracted from these articles. The events are in the form of (date, subject country, relation, object country), where the countries are represented by ISO 3166-1 alpha-3 codes and the relations are represented by the CAMEO codes defined in the ‘Conflict and Mediation Event Observations’ ontology. The relations are hierarchical: first-level relations are general parent relations represented by two-digit CAMEO codes, while second-level relations are more specific child relations represented by three-digit CAMEO codes. Child relations have the same first two digits as their parent relations. For example, ‘01’ is a first-level relation, and ‘010’ and ‘011’ are some of its second-level relations. The relations in the database are represented in the second-level form.

Your task is to forecast the future relations between two entities in a given query. You have access to a defined Python API that allows you to query the database for historical events and statistics, and to get precise information about the ISO country codes and CAMEO relation codes.

The defined API is described as follows:

```
```python
{api_description}
```
```

You will use an iterative approach, interleaving ‘Thought’, ‘Action’, and ‘Observation’ steps to collect information and perform the forecast. You may perform up to {max\_iterations} iterations. The steps are as follows:

- ‘Thought’: Analyze the current information and reason about the current situation, and predicts which API you want to use (try to use different APIs to collect diverse information) or make a decision that you want to make a final answer.
- ‘Action’: Use the API to gather more information or provide the final forecast.
  - If using the API: the action must be only one single line of exactly one function call from the API with appropriate inputs, without additional code, explanations, or natural language descriptions.
  - If making the final forecast: the action must start immediately with ‘Final Answer:’, and follow with the results in the expected JSON format.
- ‘Observation’: Return the output of the called function.

To make a reasonable forecast, you should collect both news and relational evidence to support your prediction. When you are fully confident that you accumulate enough information to make the final forecast, you should start the ‘Thought’ with your reasoning using the news and structural information to make the prediction, and then start the ‘Action’ step with ‘Final Answer:’ followed by the answer in the expected JSON format. The answer should be a JSON dictionary where the keys are the forecasted two-digit first-level CAMEO codes and the values are lists of forecasted three-digit second-level CAMEO codes that are child relations of the key. For example, ‘Action: Final Answer: "01": ["010", "011", "012"], "02": ["020", "023"]’.

3024  
3025  
3026  
3027  
3028  
3029  
3030  
3031  
3032  
3033  
3034  
3035  
3036  
3037  
3038  
3039  
3040  
3041  
3042  
3043  
3044  
3045  
3046  
3047  
3048  
3049  
3050  
3051  
3052  
3053  
3054  
3055  
3056  
3057  
3058  
3059  
3060  
3061  
3062  
3063  
3064  
3065  
3066  
3067  
3068  
3069  
3070  
3071  
3072  
3073  
3074  
3075  
3076  
3077

### System prompt for ReAct agent with action type as single function (part 2)

The final answer will be evaluated based on the precision and recall of the forecasted first-level and second-level relations, so only include confident first-level and second-level CAMEO codes in your final forecast.

Try to use different APIs to collect diverse information (including multi-hop relations), such as the precise meaning of CAMEO codes, insights from news content, relational data, and statistical analyses to support your forecasts. Consider not only the frequency of the relations but also the temporal aspects of the data when making your forecast.

## I.1.2 SYSTEM PROMPT FOR REACT AGENT WITH ACTION TYPE AS “CODE BLOCK”

### System prompt for ReAct agent with action type as code block (part 1)

You are an expert in forecasting future events based on historical data. The database contains news articles from January 1, 2023 to the current date {current\_date\_nlp} and the events extracted from these articles. The events are in the form of (date, subject country, relation, object country), where the countries are represented by ISO 3166-1 alpha-3 codes and the relations are represented by the CAMEO codes defined in the ‘Conflict and Mediation Event Observations’ ontology. The relations are hierarchical: first-level relations are general parent relations represented by two-digit CAMEO codes, while second-level relations are more specific child relations represented by three-digit CAMEO codes. Child relations have the same first two digits as their parent relations. For example, ‘01’ is a first-level relation, and ‘010’ and ‘011’ are some of its second-level relations. The relations in the database are represented in the second-level form.

Your task is to forecast the future relations between two entities in a given query. You have access to a defined Python API that allows you to query the database for historical events and statistics, and to get precise information about the ISO country codes and CAMEO relation codes. You are also authorized to utilize additional safe, well-established Python libraries such as numpy, pandas, scikit-learn, and NetworkX to enhance your data analysis and forecasting accuracy.

The defined API is described as follows:

```
```python
{api_description}
```
```

You will use an iterative approach, interleaving ‘Thought’, ‘Action’, and ‘Observation’ steps to collect information and perform the forecast. You may perform up to max\_iterations iterations. The steps are as follows:

- ‘Thought’: Analyze the current information and reason about the current situation, and predicts which API you want to use (try to use different APIs to collect diverse information) or make a decision that you want to make a final answer.
- ‘Action’: Use the API to gather more information or provide the final forecast.
  - If gathering more data: the action must be an executable Python code snippet that starts with ‘```python’ and ends with ‘```’. It can contain multiple lines of codes and function calls using the defined API or Python libraries. You must use print() to output the results, and only the printed output will be returned in the observation step.
  - If making the final forecast: the action must start immediately with ‘Final Answer:’, and follow with the answer in the expected JSON format. This should not be enclosed within triple backticks.
- ‘Observation’: Return the printed output of the executed code snippet.

3078  
3079  
3080  
3081  
3082  
3083  
3084  
3085  
3086  
3087  
3088  
3089  
3090  
3091  
3092  
3093  
3094  
3095  
3096  
3097  
3098  
3099  
3100  
3101  
3102  
3103  
3104  
3105  
3106  
3107  
3108  
3109  
3110  
3111  
3112  
3113  
3114  
3115  
3116  
3117  
3118  
3119  
3120  
3121  
3122  
3123  
3124  
3125  
3126  
3127  
3128  
3129  
3130  
3131

#### System prompt for ReAct agent with action type as code block (part 2)

To make a reasonable forecast, you should collect both news and relational evidence to support your prediction. When you are fully confident that you accumulate enough information to make the final forecast, you should start the ‘Thought’ with your reasoning using the news and structural information to make the prediction, and then start the ‘Action’ step with ‘Final Answer:’ followed by the answer in the expected JSON format. The answer should be a JSON dictionary where the keys are the forecasted two-digit first-level CAMEO codes and the values are lists of forecasted three-digit second-level CAMEO codes that are child relations of the key. For example, ‘Action: Final Answer: {"01": ["010", "011", "012"], "02": ["020", "023"]}’.

The final answer will be evaluated based on the precision and recall of the forecasted first-level and second-level relations, so only include confident first-level and second-level CAMEO codes in your final forecast.

Try to use different APIs and Python libraries to collect diverse information (including multi-hop relations), such as the precise meaning of CAMEO codes, insights from news content, relational data, and statistical analyses to support your forecasts. Consider not only the frequency of the relations but also the temporal aspects of the data when making your forecast.

#### I.1.3 SYSTEM PROMPT FOR DIRECT AGENT

##### System prompt for direct agent

You are an expert in forecasting future events based on historical data. The events are in the form of (date, subject country, relation, object country), where the countries are represented by ISO 3166-1 alpha-3 codes and the relations are represented by the CAMEO codes defined in the ‘Conflict and Mediation Event Observations’ ontology. The relations are hierarchical: first-level relations are general parent relations represented by two-digit CAMEO codes, while second-level relations are more specific child relations represented by three-digit CAMEO codes. Child relations have the same first two digits as their parent relations. For example, ‘01’ is a first-level relation, and ‘010’ and ‘011’ are some of its second-level relations. The relations in the database are represented in the second-level form.

Your task is to forecast the future relations between two entities in a given query. The answer should be a JSON dictionary where the keys are the forecasted two-digit first-level CAMEO codes and the values are lists of forecasted three-digit second-level CAMEO codes that are child relations of the key. For example, ‘Final Answer: [{"01": ["010", "011", "012"], "02": ["020", "023"]}’.

The final answer will be evaluated based on the precision and recall of the forecasted first-level and second-level relations, so only include confident first-level and second-level CAMEO codes in your final forecast.

#### I.1.4 SYSTEM PROMPT FOR CoT AGENT

##### System prompt for CoT

You are an expert in forecasting future events based on historical data. The events are in the form of (date, subject country, relation, object country), where the countries are represented by ISO 3166-1 alpha-3 codes and the relations are represented by the CAMEO codes defined in the ‘Conflict and Mediation Event Observations’ ontology. The relations are hierarchical: first-level relations are general parent relations represented by two-digit CAMEO codes, while second-level relations are more specific child relations represented by three-digit CAMEO codes. Child relations have the same first two digits as their parent relations. For example, ‘01’ is a first-level relation, and ‘010’ and ‘011’ are some of its second-level relations. The relations in the database are represented in the second-level form.

Your task is to forecast the future relations between two entities in a given query. **To make a reasonable forecast, you should first think and reason based on your background knowledge. When you are confident that you have conducted enough analysis to make the final answer,** you should start answering by ‘Therefore, the final answer is:’ followed by the answer in the expected JSON format. The JSON format should be a JSON dictionary where the keys are the forecasted two-digit first-level CAMEO codes and the values are lists of forecasted three-digit second-level CAMEO codes that are child relations of the key. For example, ‘{{“01”: [“010”, “011”, “012”], “02”: [“020”, “023”]}}’.

The final answer will be evaluated based on the precision and recall of the forecasted first-level and second-level relations, so only include confident first-level and second-level CAMEO codes in your final forecast.

#### I.2 QUERY PROMPT

##### Query prompt

Please forecast the relations that {actor1\_name} will take towards {actor2\_name} on {future\_date\_nlp} based on historical information up to {current\_date\_nlp}. I.e. forecast the relation CAMEO codes in query event Event(date={future\_date}, head\_entity=ISOCODE({actor1\_code}), relation=CAMEOCODE(?), tail\_entity=ISOCODE({actor2\_code})).

3186 I.3 ANSWER EXTRACTION PROMPT  
3187

3188 Answer extraction prompt

3189  
3190 Please help me extract final answer for forecasting the future relations between  
3191 two entities in a given query: forecast the relations that {actor1\_name} will take  
3192 towards {actor2\_name} on {future\_date\_nlp} based on historical information up  
3193 to current\_date\_nlp. I.e. forecast the relation CAMEO codes in query event  
3194 Event(date={future\_date}, head\_entity=ISOCCode({actor1\_code}), relation=CAMEOCCode(?),  
3195 tail\_entity=ISOCCode({actor2\_code})).

3196 I have used interleaving ‘Thought’, ‘Action’, and ‘Observation’ steps to collect information  
3197 from the database and perform the forecast. The database contains news articles from January  
3198 1, 2023 to the current date current\_date\_nlp and the events extracted from these articles.  
3199 The events are in the form of (date, subject country, relation, object country), where the  
3200 countries are represented by ISO 3166-1 alpha-3 codes and the relations are represented by  
3201 the CAMEO codes defined in the ‘Conflict and Mediation Event Observations’ ontology.  
3202 The relations are hierarchical: first-level relations are general parent relations represented  
3203 by two-digit CAMEO codes, while second-level relations are more specific child relations  
3204 represented by three-digit CAMEO codes. Child relations have the same first two digits as  
3205 their parent relations. For example, ‘01’ is a first-level relation, and ‘010’ and ‘011’ are some  
3206 of its second-level relations. The relations in the database are represented in the second-level  
3207 form.

3208 The final forecast answer need to forecast both first-level and second-level CAMEO codes,  
3209 and will be evaluated based on the precision and recall of both levels of relations. The final  
3210 answer content should be a JSON dictionary where the keys are the forecasted two-digit  
3211 first-level CAMEO codes and the values are lists of forecasted three-digit second-level  
3212 CAMEO codes that are child relations of the key. For example, {{{“01”: [“010”, “011”,  
3213 “012”], “02”: [“020”, “023”]}}}.

3214 The latest information and forecast I have collected is as follows:  
3215 {info}  
3216

3217 If final forecast answer has been made in the collected information indicated by "Final  
3218 Answer:", you must only reformat the final forecast answer in the expected JSON dictionary  
3219 format inside XML tags. For example: <answer>{{{“01”: [“010”, “011”, “012”], “02”:  
3220 [“020”, “023”]}}}</answer>.

3221  
3222 Otherwise, if no final forecast is made, you must reason based on the information you have  
3223 collected and generate a confident final forecast answer to the query, and then reformat your  
3224 answer in the expected JSON dictionary format inside XML tags.  
3225  
3226  
3227  
3228  
3229  
3230  
3231  
3232  
3233  
3234  
3235  
3236  
3237  
3238  
3239



3240 J DATASHEET FOR MIRAI

3241 J.1 MOTIVATION

3243 1. **For what purpose was the dataset created?**

3244 The MIRAI dataset was created to evaluate different capabilities of Large Language Model (LLM)  
3245 agents in forecasting international events.

3246 2. **Who created the dataset and on behalf of which entity?**

3247 The dataset was developed by a group of university researchers studying LLM reasoning. All  
3248 rights are held by the individuals themselves, rather than by third-party stakeholders.

3249 3. **Who funded the creation of the dataset?**

3250 The academic institutions support the university research lab that the authors affiliate in.

3251 J.2 COMPOSITION

3252 1. **What do the instances that comprise the dataset represent?**

3253 Each instance in the dataset represents a record of international events, including the date, involved  
3254 countries, and type of event, along with associated news articles and metadata.

3256 2. **How many instances are there in total?**

3257 The current database comprises 1,296,991 GDELT event records from January 2023 to February  
3258 2024, corresponding to 75,341 unique events and 401,013 news articles. Test splits span over  
3259 November 2023 to February 2024, each containing 100 balanced-sampled forecasting event  
3260 queries. Note these statistics reflect the dataset used in this paper’s experiments. As a dynamic  
3261 benchmark, we are committed to periodic updates, continuously expanding the historical database  
and creating new test splits.

3262 3. **Does the dataset contain all possible instances or is it a sample of instances from a larger set?**

3263 The dataset represents a curated sample from the entire GDELT database. It has been created  
3264 through meticulously designed data cleaning and preprocessing steps on GDELT raw data, aimed  
3265 at enhancing the quality and reliability of the event data.

3266 4. **Is there a label or target associated with each instance?**

3267 Yes, each instance in the dataset is an event labeled with a relation type derived from the CAMEO<sup>7</sup>  
3268 event taxonomy.

3269 5. **Is any information missing from individual instances?**

3270 No, all instances are complete with all available information.

3272 6. **Are there recommended data splits (e.g., training, development/validation, testing)?**

3273 The dataset is flexible in splitting data and expanding future events to new test splits. In this  
3274 paper, we provide multiple test splits, spanning from November 2023 to February 2024. Each test  
3275 instance utilizes all preceding records as its historical dataset to ensure accuracy and relevance in  
3276 analysis. For training and development, users can flexibly use the data prior to the test splits based  
on their specific research needs.

3277 7. **Are there any errors, sources of noise, or redundancies in the dataset?**

3278 The dataset has undergone extensive cleaning and structuring to minimize errors and noise.  
3279 However, residual noise from the original GDELT database may still be present.

3280 8. **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,  
3281 websites, tweets, other datasets)?**

3282 While the dataset is primarily self-contained, it includes URLs to news articles, providing links to  
3283 external resources for further context and verification.

3284 9. **Does the dataset contain data that might be considered confidential?**

3285 No, the dataset contains publicly available data, and does not include confidential information. It  
3286 follows the term of use for GDELT, which is an open platform for research and analysis of global  
3287 society.

3288 10. **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threaten-  
3289 ing, or might otherwise cause anxiety?**

3290 The dataset may include descriptions of sensitive events, such as global conflicts, due to its focus  
3291 on international relations. However, We conduct rigorous text cleaning procedures to reduce noise  
3292 from web content while enhancing the reliability and ethical integrity of the textual information,  
3293

<sup>7</sup>Conflict and Mediation Event Observations (CAMEO): <https://parusanalytics.com/eventdata/data.dir/cameo.html>

3294 following the OBELICS protocol (Laurençon et al., 2023). Thorough checks such as flagging  
 3295 word ratios are employed during the data cleaning process to identify and potentially exclude  
 3296 inappropriate paragraphs or entire news articles and events. This ensures the minimization of  
 3297 distressing content while maintaining the integrity and relevance of the dataset for academic study.

### 3298 J.3 COLLECTION PROCESS

#### 3299 1. **How was the data associated with each instance acquired?**

3300 Data for each instance was sourced from the GDELТ project, which aggregates global event data  
 3301 and news articles from various worldwide media. Detailed information can be found in Section 2.3  
 3302 of the paper.

#### 3303 2. **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sensor, manual human curation, software program, software API)?**

3304 Data collection was facilitated through software programs that aggregate event data and news  
 3305 articles from various sources. This automated collection is followed by data cleaning processes to  
 3306 enhance completeness and reliability. For a detailed description of the collection and cleaning  
 3307 methods, please refer to Section 2.3 of the paper. All dataset construction scripts are available in  
 3308 [Github](#) and also described in its README.

#### 3309 3. **Who was involved in the data collection process? (e.g., students, crowdworkers, contractors) and how were they compensated (e.g., how much were crowdworkers paid)?**

3310 The data collection and preprocessing were conducted by the authors, who are university re-  
 3311 searchers. This work was performed as part of their academic research activities.

#### 3312 4. **Does the dataset relate to people?**

3313 No. The dataset does not contain personal data or directly relate to individual people.

#### 3314 5. **Did you collect the data from the individuals in questions directly, or obtain it via third parties or other sources (e.g., websites)?**

3315 The dataset does not involve data collected from individuals directly. Instead, it is curated from  
 3316 the GDELТ project, which aggregates information from various global news media sources.

### 3317 J.4 USES

#### 3318 1. **Has the dataset been used for any tasks already?**

3319 The dataset has not been used for any tasks other than the ones proposed and examined in this  
 3320 current paper, specifically for benchmarking the forecasting capabilities of LLM agents in predict-  
 3321 ing international relations. The dataset has also supported evaluating forecasting performance of  
 3322 traditional temporal knowledge graph-based methods and natural language-based methods.

#### 3323 2. **What (other) tasks could the dataset be used for?**

3324 In addition to benchmarking LLM agents and traditional forecasting methods, the dataset could  
 3325 be valuable for research in geopolitics, the development of other event prediction algorithms,  
 3326 sentiment analysis of international events, and trend analysis in global political dynamics.

#### 3327 3. **Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses?**

3328 The dataset was meticulously collected and preprocessed to ensure high quality and reliability.  
 3329 However, the reliance on the GDELТ project as the primary data source and the specific cleaning  
 3330 process applied might influence its applicability to certain tasks. Users should consider these  
 3331 factors when applying the dataset to different research areas or methodologies.

#### 3332 4. **Are there tasks for which the dataset should not be used?**

3333 The dataset should not be used for any tasks that violate the terms of use associated with the  
 3334 GDELТ project. We clearly cite the terms of use in Appendix J.7.

### 3335 J.5 DISTRIBUTION

#### 3336 1. **Will the dataset be distributed to third parties outside of the entity?**

3337 Yes, the dataset, evaluation codes and leaderboards are intended to be publicly available to foster  
 3338 future research and development.

#### 3339 2. **How will the dataset be distributed?**

3340 The database and codebase are currently available via an [academic website](#), [Google Drive](#), and  
 3341 [Github](#). To enhance the accessibility and utility, the distribution of current version of data and its  
 3342 future updates will be enhanced by uploading the dataset to Hugging Face, and refining the API  
 3343 into a more user-friendly library format in the future.

3348 3. **Have any third parties imposed IP-based or other restrictions on the data associated with**  
3349 **the instances?**

3350 No, there are no IP-based or other restrictions on the data: MIRAI is curated based on the GDELT<sup>8</sup>  
3351 Event Database, which is an open platform for research and analysis of global society and all  
3352 datasets released by the GDELT Project are available for unlimited and unrestricted use for any  
3353 academic, commercial, or governmental use of any kind without fee; Data are allowed for any  
3354 redistribution, rehost, republish, and mirror of the GDELT datasets in any form, with necessary  
3355 citations<sup>9</sup>.

3356 4. **Do any export controls or other regulatory restrictions apply to the dataset or to individual**  
3357 **instances?**

3358 No, there are no export controls or other regulatory restrictions applied to this dataset.

3359 5. **When will the dataset be distributed?**

3360 The dataset will be made publicly available after the review process is completed, with the current  
3361 [academic website](#), [Google Drive](#), and [Github](#), and additional release/updates of arXiv, Hugging  
3362 Face, and leaderboard to facilitate a more comprehensive accessibility to the research community.

3363 6. **Will the dataset be distributed under a copyright or other IP license, and/or under applicable**  
3364 **terms of use (ToU)?**

3365 The dataset will be distributed under the CC BY-NC 4.0 license, allowing for use and distribution  
3366 for non-commercial purposes with appropriate attribution.

3367 J.6 MAINTENANCE

3368 1. **Who will be supporting/hosting/maintaining the dataset?**

3369 The dataset maintenance will be supporting/hosting/maintaining by the authors.

3370 2. **How can the owner/curator/manager of the dataset be contacted?**

3371 The owner/curator/manager of the dataset can be contacted through the authors' emails.

3372 3. **Will the dataset be updated? (e.g., to correct labeling errors, add new instances, delete**  
3373 **instances)?**

3374 Yes, the dataset is designed as a dynamic benchmark with periodic updates. We are committed  
3375 to regularly expanding the historical database and creating new test splits. Updates will include  
3376 adding new event data, creating new test splits, and potentially correcting any identified errors.  
3377 Announcements regarding updates will be made through the project's official channels.

3378 4. **If the dataset relates to people, are there applicable limits on the retention of the data**  
3379 **associated with the instances (e.g., were the individuals in question told that their data would**  
3380 **be retained for a fixed period of time and then deleted?)**

3381 NA. This dataset does not contain data related to individuals or personal identifiers. It consists of  
3382 aggregated event data and news articles related to international events.

3383 5. **Will older version of the dataset continue to be supported/hosted/maintained?**

3384 Our dataset is designed as a cumulative, evolving benchmark. Historical event data is continuously  
3385 incorporated into the database and remains available for future event forecasting. While we  
3386 strongly recommend using the latest test splits to ensure data-contamination-free benchmarking  
3387 for all LLM models, older historical data remains an integral part of the dataset. For reproducibility  
3388 purposes, we will maintain records of the specific data versions used in published experiments.  
3389 Researchers can request access to particular historical snapshots of the dataset, subject to resource  
3390 availability. However, for ongoing research and comparisons, we encourage using the most current  
3391 version of the dataset and its latest test splits.

3392 6. **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism**  
3393 **for them to do so?**

3394 Researchers and developers interested in extending, augmenting, or contributing to the dataset  
3395 are encouraged to submit their changes through GitHub pull requests. For additional inquiries or  
3396 detailed discussions, contacting the authors via email is recommended.

3397 J.7 TERM OF USE FOR GDELT

3398 Based on <https://www.gdeltproject.org/about.html#termsfuse>, GDELT  
3399 dataset "is an open platform for research and analysis of global society and thus all datasets  
3400

3401 <sup>8</sup>GDELT Project: <https://www.gdeltproject.org/>

<sup>9</sup>GDELT Term of Use: <https://www.gdeltproject.org/about.html#termsfuse>

3402 released by the GDEL T Project are available for unlimited and unrestricted use for any aca-  
3403 demic, commercial, or governmental use of any kind without fee.”, as long as “any use or re-  
3404 distribution of the data must include a citation to the GDEL T Project and a link to this website  
3405 (<https://www.gdel tproject.org/>).”, which we’ve cited in abstract.  
3406  
3407  
3408  
3409  
3410  
3411  
3412  
3413  
3414  
3415  
3416  
3417  
3418  
3419  
3420  
3421  
3422  
3423  
3424  
3425  
3426  
3427  
3428  
3429  
3430  
3431  
3432  
3433  
3434  
3435  
3436  
3437  
3438  
3439  
3440  
3441  
3442  
3443  
3444  
3445  
3446  
3447  
3448  
3449  
3450  
3451  
3452  
3453  
3454  
3455