# Temporal Relational Reasoning of Large Language Models for Detecting Stock Portfolio Crashes

Kelvin J.L. Koa
National University of Singapore
kelvin.koa@u.nus.edu

Yunshan Ma*
Singapore Management University
ysma@smu.edu.sg

Yi Xu
National University of Singapore
yixu@u.nus.edu

Ritchie Ng
National University of Singapore
ritchieng@u.nus.edu

Huanhuan Zheng
National University of Singapore
sppzhen@nus.edu.sg

Tat-Seng Chua
National University of Singapore
dcscts@nus.edu.sg

## ABSTRACT

Stock portfolios are often exposed to rare consequential events (*e.g.,* 2007 global financial crisis, 2020 COVID-19 stock market crash), as they do not have enough historical information to learn from. Large Language Models (LLMs) now present a possible tool to tackle this problem, as they can generalize across their large corpus of training data and perform zero-shot reasoning on new events, allowing them to detect possible portfolio crash events without requiring specific training data. However, detecting portfolio crashes is a complex problem that requires more than reasoning abilities. Investors need to dynamically process the impact of each new piece of information found in news articles, analyze the relational network of impacts across different events and portfolio stocks, as well as understand the temporal context between impacts across time-steps, in order to obtain the aggregated impact on the target portfolio. In this work, we propose an algorithmic framework named Temporal Relational Reasoning (TRR). It seeks to emulate the spectrum of human cognitive capabilities used for complex problem-solving, which include *brainstorming*, *memory*, *attention* and *reasoning*. Through extensive experiments, we show that TRR is able to outperform state-of-the-art techniques on detecting stock portfolio crashes, and demonstrate how each of the proposed components help to contribute to its performance through an ablation study. Additionally, we further explore the possible applications of TRR by extending it to other related complex problems, such as the detection of possible global crisis events in Macroeconomics.

## CCS CONCEPTS

• **Computing methodologies** → *Information extraction*; • **Applied computing** → *Forecasting*; *Economics*.

## KEYWORDS

Large Language Models, Temporal Graphs, Stock Portfolio Crash

## 1 INTRODUCTION

In equity investing [31, 32], investors typically form stock portfolios [39] to diversify their risk across multiple stocks. This could be done by selecting the stocks from across different categories, based on geographical regions [46] and/or business sectors [12], *etc.* , in order to dampen the impacts of events that affect any specific category. However, there also exist rare, consequential events [38] that are unprecedented in history and cause the market to be increasingly interconnected [50], which can result in crashes (*e.g.,* 2007 global

financial crisis [20], 2020 COVID-19 stock market crash [40]). Even though stock portfolios have extensively considered various risks when they are curated, they are still often ill-prepared to handle these events [50], as they do not have past statistics or historical information to learn from. Due to this reason, there are currently limited works [51] on detecting portfolio crash events in literature.

Today, Large Language Models (LLMs) present a possible toolset for detecting these crash events, without requiring specific training data for the task. This stems from their known capabilities to perform zero-shot reasoning [33], which can be attributed to their ability to generalize [6, 63] across the large corpus of data they have previously been trained on. This allows them to identify repeating patterns on new emerging events, and potentially detect possible crashes before they happen. In this work, we explore the use of LLMs to predict possible stock portfolio crashes, by reasoning over publicly-available news information that can be found on the web.

Detecting portfolio crashes is a complex problem that requires more than basic reasoning abilities. Currently, there are some reasoning frameworks for LLMs to handle complex tasks: Thought-based frameworks (*e.g.,* Tree-of-Thoughts (ToT) [64], Graph-of-Thoughts (GoT) [7]) break down a task into generated sub-steps that can be merged to solve the task; Search-based frameworks (*e.g.,* Think-on-Graph (ToG) [49]) search through an existing Knowledge Graph of facts to find a reasoning path that can answer questions on a single entity. However, among these methods, we can still identify three challenges for our task. **(1)** The current methods focus mainly on tackling isolated problems, such as solving a task through sequential thoughts or answering questions from a static graph of information. However, these methods do not deal with the constantly evolving nature of news events, which would require constant and dynamic processing of new information. **(2)** Portfolio crashes are often caused by the unexpected interconnectivity of its constituent stocks in response to unprecedented events [20, 40]. While current frameworks combine thoughts or search for a single path on a Knowledge Graph (see Figure 1), they do not reason across *multiple* search paths, which could reveal these interconnectivity between news events and portfolio stocks within the graph. **(3)** It is also known from stock prediction works that there exists temporal context dependency [27, 62] between news events when considering their impacts on stock prices. While there are some LLM works on temporal graphs [60, 68], these works focus on performing question-answering on individual graphs with temporal information in the nodes or edges, and do not handle information spread across multiple graphs captured from different time-steps.

---

*Corresponding author.

(a) Thought-Based Works   (b) Search-Based Works   (c) TRR, Relational Only (Ours)   (d) TRR, Relational+Temporal (Ours)
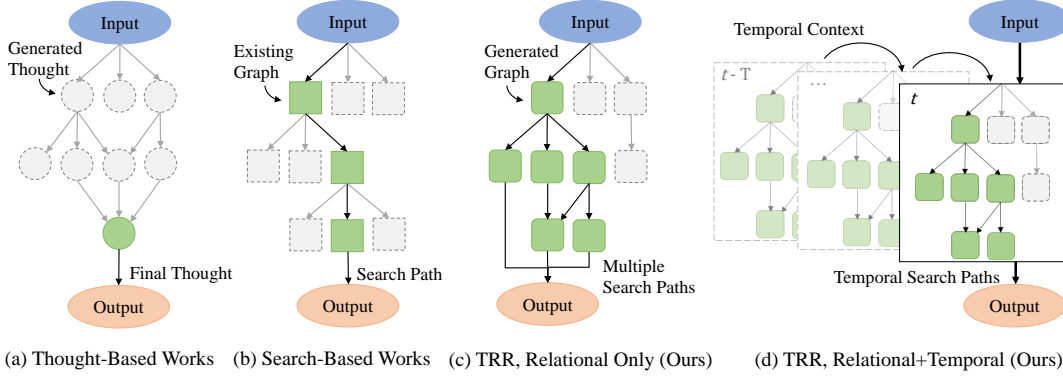
**Figure 1: Illustration of thought-based frameworks (*e.g.,* ToT [49], GoT [7]), search-based frameworks (*e.g.,* ToG [49]), and our proposed Temporal Relational Reasoning (TRR) framework, with Relational Only and Relational+Temporal variants.**

To tackle the above-mentioned problems, we propose an algorithmic framework named *Temporal Relational Reasoning* (TRR), which seeks to emulate the spectrum of human cognitive capabilities that are used for complex problem-solving (see Figure 2). Using the retrieved daily online news, we *brainstorm* for all possible impacts on the target portfolio of stocks by dynamically generating chains of sub-impacts through related entities. These are used to form a graph, where the nodes start from each news article, passes through the LLM-generated entities and end at the portfolio stocks, while the edges represent the direction of impacts. Next, to understand the temporal context of these impacts, we simulate the associative *memory* of humans by retrieving past related events which have affected the same entities, from a stored temporal "memory-bank" that decays over time [70]. Finally, to analyze the interconnectivity across the relational network of impacts, we first mimic the *attention* of investors by utilizing the PageRank algorithm [44], to filter the most important impact chains that will affect the target portfolio, and reduce the size of the final graph. We then emulate investors' *reasoning* process by reasoning across the filtered temporal-relational graph, to determine if a crash is likely to occur.

To demonstrate the effectiveness of TRR, we perform extensive experiments over multiple portfolios and time periods, and show that our method outperforms deep-learning models and other LLM reasoning frameworks in predicting portfolio crashes. Through an ablation study, we also demonstrate how each component of TRR helps to contribute to its performance. Furthermore, we explore the applications of our method to other related complex problems by extending it to a macroeconomic setting. By viewing the global economy as a *network of regional economies*, we determine if the set of news could result in possible global financial crises, by tracing their overall impacts using TRR. We find that TRR can also predict crisis events more effectively than the available economic methods.

The main contributions of this work can be summarized as below:

- We investigate the limitations of zero-shot LLMs on complex problems such as detecting portfolio crashes, which require dealing with information across a temporal-relational network.
- We propose an integrated framework that allows a LLM to reason over a self-generated temporal-relational graph. This is done through a fully algorithmic framework that emulates the set of human cognitive capabilities used for solving complex problems, which include *brainstorming*, *memory*, *attention* and *reasoning*.

- We conduct extensive experiments across multiple portfolio structures and time periods, and show that our TRR framework can detect portfolio crashes more effectively than state-of-the-art techniques. Given these results, we further explore other possible applications of TRR on complex problems, and utilize the framework to detect global crisis events in macroeconomics.

## 2 RELATED WORKS

In this section, we trace the progress in the use of relational graph techniques in stock-related works, and also explore various zero-shot reasoning frameworks for Large Language Models (LLMs).

**Relational Stock Prediction** Utilizing relational information to predict stock prices have been widely explored in multiple previous works. Early works [16, 17] have studied the use of relational tuples in the form of (*Actor*, *Action*, *Object*, *Timestamp*) to learn embeddings, such that similar events [16] or similar stock entities [17] would have similar representative vectors. The tuples are generated with rule-based techniques [67] from news headlines, as opposed to our LLM-based method, which can form multiple levels of relations.

Later works would improve on this using graph-based methods [18, 19, 47], by learning embeddings across a Knowledge Graph to represent stock entities. These works utilize stock relational information from external Knowledge Graphs such as Freebase [8] or Wikidata [53] to train their models. However, these models rely mainly on the static relation information retrieved from a central database, and do not consider possible changes in the connectivity between stocks, that could result from new company developments.

In a more recent work [13], LLMs were used to infer relations between stocks from news headlines, resulting in more dynamic relational data. This information is then used to generate stock embeddings using a Graph Neural Network, which are used to *train* a deep-learning model to do single stock prediction. In contrast to this, our work focuses on *zero-shot* reasoning frameworks in order to detect crashes across a portfolio of stocks, that often occur due to events that are *unprecedented* in the historical training data.

**Reasoning Frameworks for Large Language Models** LLMs are known for their zero-shot reasoning capabilities [33], which has been largely attributed to their ability to generalize knowledge across the large corpus of data it was trained on [11, 63]. To enhance this capability, researchers have proposed reasoning frameworks to tackle more advanced tasks, such as ToT [64], GoT [7]. In particular, these works were stated [64] to be inspired by general
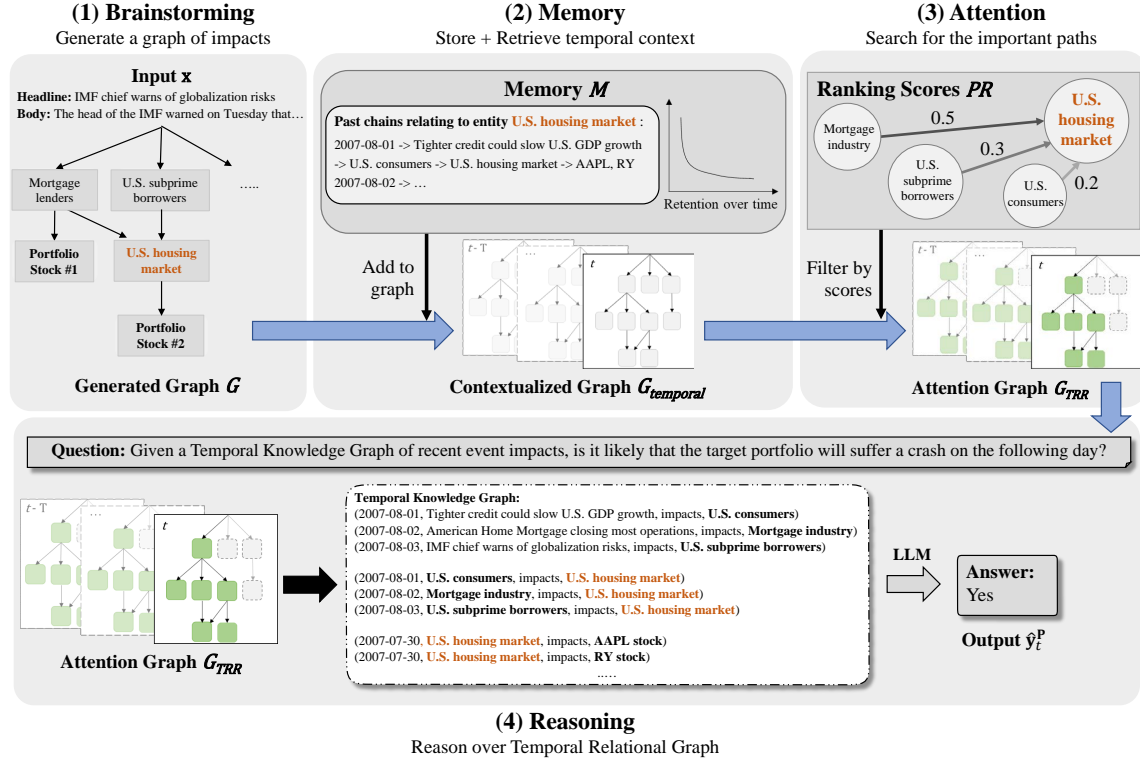
**Figure 2: The components of our Temporal Relational Reasoning (TRR) framework. TRR emulates the human cognitive capabilities used for solving complex problems, that include (1) *brainstorming*, (2) *memory*, (3) *attention* and (4) *reasoning* skills.**

problem-solving strategies from the 1950s [41], which can be seen as searching through a combinatorial problem space to find a task solution. Since then, the original source research has further led to works discussing more complex problems dealing with multiple entities [21, 48]. which serve as a source of inspiration for our work.

Another line of research deals with enhancing the reliability of LLMs' responses by using Knowledge Graphs containing external information, such as StructGPT [29] and ToG [49]. A key observation from these works is the ability of LLMs to reason over Knowledge Graph inputs, provided in the form of triplet tuples. However, while these works focus on extracting specific tuples from an *existing* Knowledge Graph of information to answer questions, we extract a sub-graph of tuples from a *generated* graph to discover any new correlated impacts and detect possible portfolio crashes.

## 3 TEMPORAL RELATIONAL REASONING

Our TRR framework seeks to emulate the spectrum of human cognitive capabilities that are used together for complex problem-solving (see Figure 2). It consists of four phases: 1. **Brainstorming**, which generates a graph of sub-impacts on affected entities; 2. **Memory**, which retrieves relevant past impact chains that contain the same entities; 3. **Attention**, which extracts the most important impact chains to form a new sub-graph; 4. **Reasoning**, which reasons over this sub-graph to to determine if a portfolio crash will occur.

For this task, we begin with a specified portfolio of $N$ stocks, $P = \{s_1, s_2, \cdots, s_N\}$, where $s_n$ is a single stock and $n \in N$. For each day, given a set of $J$ news articles $X = \{x_1, x_2, \cdots, x_J\}$ retrieved from the web, we aim to make a binary prediction on whether a crash will occur on portfolio $P$ on the following day, *i.e.,* $\hat{y} \in \{0, 1\}$.

### 3.1 Brainstorming

To obtain the overall impact of each news event on our selected portfolio, we first brainstorm for all possible chains of sub-impacts that lead to its constituent stocks. We model this as a directed graph $G = (Z, A)$. The set of vertices is $Z = \{x_1, \cdots, x_J, e_1, \cdots, e_H, s_1, \cdots, s_N\}$, which starts from each news article $x_j$, passes through all impacted sub-entities $e_h$ where $h \in H$, and ends at the portfolio stocks $s_n$. $H$ refers to the total number of entities generated by the LLM. The set of edges is $A$, which simply represent the direction of impact between the vertices. Some possible examples are (**American Home Mortgage closing most operations**, *impacts*, **Mortgage industry**) and (**Mortgage industry**, *impacts*, **U.S. housing market**).

To generate new vertices for the graph, we iteratively prompt an LLM to generate possible affected entities given each news article or previously affected sub-entity. Hence, for all vertices $z_i \in Z$ at iteration $i$, we have $[z^1, \cdots z^k] \sim p_\theta^{brainstorm}(z_{i+1}^{(1 \cdots k)} | z_i)$, where $k$ is the number of new vertices generated. This is done iteratively until the chain of impacts reaches a portfolio stock $s_n$, or the max number of iterations $I$ is reached. Repeated entities are merged as a single vertex on the graph, following works such as GoT [7]. Some examples of the prompts used can be found in Appendix D.

For notation, we also refer to each individual chain of impact as $C$, where $C \in G$. Each chain $C$ starts from a news article $x_j$, passes through any number of impacted entities, and ends at a stock $s_n$.

### 3.2 Memory

For understanding the temporal context of news events, we equip TRR with a memory module, which stores all previous mentions of impacted sub-entities. We denote the memory module as $\mathcal{M} =$

$\{M_{e_0}, M_{e_1}, \cdots\}$, where $M_{e_h}$ is the collection of all the previous impact chains $C_{e_h}$ that contains the entity $e_h$. With the memory module, we are then able to perform retrieve and store functions:

(a) On each day, for each impacted entity that is generated by the LLM, we search the memory module for its previous mentions and add them to our daily graph $G$, which give us a new temporal contextualized graph $G_{temporal}$. Given the impacted entities, we can obtain: $G_{temporal} = \bigcup_{\forall e_h \in Z}(G \cup M_{e_h})$.

(b) At the end of each day, we then update the memory module with the daily chains of impact, i.e., $M_{e_h} = \bigcup_{\forall e_h \in Z}(M_{e_h} \cup C_{e_h})$, which stores the temporal context for future time periods.

Here, the additional temporal context allows us to form a temporal relational graph $G \rightarrow G_{temporal}$, which represents the Relational+Temporal variant of our TRR framework (see Figure 1).

In addition, human collective memory on news events tend to fade over time [2, 23], which can lessen their impact on the market. The temporal decay of memory has previously been modelled with an exponential decay in both the social sciences [23, 36] and LLM works [70]. Following this, we also track the memory retention of each impact with a variable $R_{u,v} = \exp(-\frac{t_{u,v}}{\lambda})$, where $t_{u,v}$ is the time-step when entity $u$ impacted entity $v$, and $\lambda$ is a decay rate constant to be determined. The variable $R_{u,v}$ will be used in the next phase to decide if an impact is considered in the market context.

### 3.3 Attention

The overall temporal-contextualized graph $G_{temporal}$ is too large to be used in LLMs, which have fixed token limits. While other works deal with this by merging thoughts [7] or finding a single answering path [49], we want to maintain a relational graph of information (see Figure 1) in order to provide the LLM with a holistic overview of the market. In a similar fashion, the amount of news from the web that investors can process each day is also limited, and their attention is usually focused on more important information [25, 61].

$G_{temporal}$ contains a network of impact chains, with varying impact strengths on the target portfolio. To obtain the most important information on this network, we draw inspiration from the PageRank algorithm [44] by assigning ranking scores to each entity. This is done by iteratively transferring scores across the entities following the direction of impact, until the convergence of scores. Furthermore, the scores are weighted [59] based on their retention in memory from the previous Memory phase. For an entity $e_h \in Z$ and the set of its parent vertices $B_{e_h}$, the ranking score it receives in each iteration can be formulated as $PR(e_h) = \sum_{b \in B_{e_h}} \frac{PR(b)}{L_b} \cdot R_{b,e_h}$, where $L_b$ is the number of outgoing impacts from an entity $b$.

Using the ranking scores $PR(e_h)$, we then filter for all the impact chains $C$ containing the top-$q$ highest scoring entities, where $q$ is to be defined. These are used to form a new sub-graph $G_{TRR}$, which represents the most important information for each day that investors would pay attention to, which likely impact the market.

### 3.4 Reasoning

Finally, to determine if a portfolio crash will occur, we reason over the generated temporal-relational graph $G_{TRR}$. This emulates the reasoning process of investors, who will assess their portfolio risks by considering the most relevant news impacts and how the constituent stocks would be interconnected [4, 14]. Following previous graph-based LLM works [49, 66], we let a LLM reason on the graph

$G_{TRR}$ in the form of relational tuples [65]. Each tuple can be formalized as $(t, z_s, a, z_o)$, where $t$ is the time-step when the impact was generated, $z_s$ and $z_o$ are the subject and object entities, and $a$ is the direction of impact. Given the stock portfolio $P$ and the graph $G_{TRR}$ in the form of tuples, we prompt the LLM to generate a crash prediction for the next day. This step can be formalized as $\hat{y} \sim p_\theta^{reason}(\hat{y}|P, G_{TRR})$. The prompt can be found in Appendix D.

## 4 EXPERIMENTS

We extensively evaluate TRR across multiple portfolios and time periods to demonstrate its effectiveness. We form two *diversified* portfolios using common investor strategies: **(1)** Country-neutral portfolio, where each constituent stock company is based in a different country [46]; **(2)** Sector-neutral portfolio, where each stock company is from a different market sector [12] (see Appendix A).

To determine if we can accurately predict crashes, we select three notable time periods containing events which have caused a big impact on the stock market: **(a)** June-August 2007 (Global financial crisis); **(b)** March-May 2010 (Greek government debt crisis); **(c)** January-March 2020 (COVID-19 stock market crash). Each time period consists of three months, and the mentioned events can be found towards its middle, which allows us to capture news impacts from both the stable (*i.e.,* before portfolio crash) and crash periods.

To evaluate our model under more typical market conditions, we also run our model on **(d)** January-December 2012, which is relatively stable year. Given that crash events are rare, this experiment highlights the performance of our model over a longer, regular period, and sees if it would generate any erroneous False Positives.

To ensure that the results were not affected by data contamination, we also run experiments over **(e)** October 2021-March 2022, which is after the post-knowledge cut-off date of GPT-3.5-turbo [42], the selected model for TRR. This period also coincides with the start of the Russo-Ukrainian War, which resulted in a brief stock market decline. This experiment ensures that the reported performance is attributed to TRR, and not its internal knowledge.

**Table 1: Statistics of the dataset. As the portfolios are widely diversified and the events that have significant impacts are usually rare, the percentage of data labelled as Crash are low.**

| Dataset | Time Period | Description | % Crash (Country-Neu) | % Crash (Sector-Neu) |
|---------|-------------|-------------|-----------------------|----------------------|
| **2007** | Jun 2007-Aug 2007 | Global Financial Crisis | 0.108 | 0.046 |
| **2010** | Mar 2010-May 2010 | Greek government debt crisis | 0.078 | 0.063 |
| **2020** | Jan 2020-Mar 2020 | COVID-19 stock market crash | 0.161 | 0.210 |
| **2012** | Jan 2012-Dec 2012 | Supplementary (Stable Period) | 0.020 | 0.008 |
| **2022** | Oct 2021-Mar 2022 | Supplementary (Post-Knowledge) | 0.047 | 0.094 |

### 4.1 Dataset and Evaluation Metrics

For news data, we use the Reuters financial news dataset [16], which we also extend to the year 2020 to cover the selected time periods. The dataset contains general financial news from Reuters which are *not* filtered by any stock or country. This allows the LLM to decide by itself if each article is relevant to the target portfolio.

**Table 2: Overall performance comparison of TRR against deep-learning and LLM baselines. Each individual result represents the average AUROC and standard deviation (presented in subscript) across 5 runs. The best results are presented in boldface.**

| Models | | 2007 | | 2010 | | 2020 | | 2022 | |
|---|---|---|---|---|---|---|---|---|---|
| | | Country-Neu | Sector-Neu | Country-Neu | Sector-Neu | Country-Neu | Sector-Neu | Country-Neu | Sector-Neu |
| Deep-Learning Models | GRU+GAT | $0.458_{\pm 0.0061}$ | $0.472_{\pm 0.0116}$ | $0.509_{\pm 0.0038}$ | $0.501_{\pm 0.0010}$ | $0.503_{\pm 0.0207}$ | $0.510_{\pm 0.0085}$ | $0.495_{\pm 0.0170}$ | $0.493_{\pm 0.0158}$ |
| | GPT+GNN | $0.514_{\pm 0.0127}$ | $0.504_{\pm 0.0039}$ | $0.508_{\pm 0.0081}$ | $0.504_{\pm 0.0089}$ | $0.500_{\pm 0.0013}$ | $0.499_{\pm 0.0014}$ | $0.539_{\pm 0.0310}$ | $0.515_{\pm 0.0351}$ |
| | DAN+FCL | $0.607_{\pm 0.0210}$ | $0.468_{\pm 0.0096}$ | $0.597_{\pm 0.0092}$ | $0.505_{\pm 0.0320}$ | $0.441_{\pm 0.0220}$ | $0.452_{\pm 0.0084}$ | $0.523_{\pm 0.0563}$ | $0.503_{\pm 0.0375}$ |
| LLM Frameworks | IO | $0.465_{\pm 0.0349}$ | $0.440_{\pm 0.0104}$ | $0.500_{\pm 0.0000}$ | $0.483_{\pm 0.0009}$ | $0.503_{\pm 0.0289}$ | $0.542_{\pm 0.0102}$ | $0.454_{\pm 0.0356}$ | $0.511_{\pm 0.0400}$ |
| | CoT | $0.480_{\pm 0.0328}$ | $0.504_{\pm 0.0123}$ | $0.508_{\pm 0.0506}$ | $0.494_{\pm 0.0048}$ | $0.519_{\pm 0.0533}$ | $0.564_{\pm 0.0114}$ | $0.504_{\pm 0.0387}$ | $0.518_{\pm 0.0739}$ |
| | GoT | $0.500_{\pm 0.0489}$ | $0.577_{\pm 0.0854}$ | $0.537_{\pm 0.0129}$ | $0.522_{\pm 0.0709}$ | $0.573_{\pm 0.0352}$ | $0.600_{\pm 0.0514}$ | $0.512_{\pm 0.0108}$ | $0.531_{\pm 0.0230}$ |
| | ToG | $0.502_{\pm 0.0589}$ | $0.629_{\pm 0.0453}$ | $0.585_{\pm 0.0160}$ | $0.588_{\pm 0.0118}$ | $0.607_{\pm 0.0540}$ | $0.625_{\pm 0.0558}$ | $0.579_{\pm 0.0423}$ | $0.587_{\pm 0.0399}$ |
| | TRR (Ours) | $\mathbf{0.690_{\pm 0.0426}}$ | $\mathbf{0.684_{\pm 0.0361}}$ | $\mathbf{0.610_{\pm 0.0625}}$ | $\mathbf{0.598_{\pm 0.0787}}$ | $\mathbf{0.657_{\pm 0.0231}}$ | $\mathbf{0.644_{\pm 0.0350}}$ | $\mathbf{0.638_{\pm 0.0348}}$ | $\mathbf{0.601_{\pm 0.0425}}$ |

To generate the portfolio crash labels, we first retrieve the price data of each constituent stock from Yahoo Finance, and calculate their daily percentage returns. Next, we average them to obtain the portfolio returns for each day. To find the portfolio crashes, we set a threshold to capture sharp falls in value [34]. We label returns $\leq -2.0\%$ as a crash, which represents the bottom $95^{th}$ percentile of the overall returns series [10, 55]. In addition, as the occurrences of these crashes are rare (see Table 1), the dataset is largely imbalanced and prediction models that predict all False (*i.e.,* no crashes) would produce a high accuracy score. Following works that deal with highly imbalanced classification in Object Detection [43] and Medical Imaging [52] tasks, we use the Area Under Receiver Operating Characteristics curve (AUROC) as our metric, which captures the trade-off between True and False Positives across all thresholds.

## 4.2 Baselines

As the task of detecting portfolio crashes is not widely explored currently, we compare with multiple zero-shot LLM reasoning frameworks, such as standard IO prompting, Chain-of-Thought (CoT) [54], Graph-of-Thoughts (GoT) [7] and Think-on-Graph (ToG) [49]. Descriptions of their implementation can be found in Appendix B.

In addition, we also compare with various relevant deep-learning models, which are trained on past Reuters news data from the same dataset, that is four times of each respective test data size (resembling a 7:1:2 train-valid-test split). While this additional training seemingly gives them an "advantage" over zero-shot LLM reasoning, we argue that they are not able to handle unprecedented events that has not previously occurred, *e.g.,* COVID-19. These models are:

- **GRU+GAT** [47]: In this model, a Gated Recurrent Unit (GRU) network is enhanced with a Graph Attention Network (GAT) to predict stock movements using text and relational data. Here, the stock relational graph is *static*, retrieved from Wikidata [19].
- **GPT+GNN** [13]: In this work, ChatGPT was used to generate the relational graph *dynamically* from daily news headlines. A Graph Neural Network (GNN) is then used to generate embeddings for stock prediction. As it requires training, this model would still be susceptible to the imbalanced dataset and unprecedented events.
- **DAN+FCL** [35]: This method was designed to deal with imbalanced dataset in stock prediction. The model first generates sentiment embeddings with FinBERT [1]. To make predictions, it then uses a Deep Averaging Network (DAN) [22], coupled with a novel Focal Calibration Loss (FCL) to handle class imbalance.

## 4.3 Parameter Settings

For all LLM experiments, we use OpenAI GPT-3.5-turbo to generate the responses, with a temperature setting of 0.0 to maximize replicability. In the last reasoning phase, we repeat the prediction prompt 5 times for each model and report the average AUROC and standard deviation. For the main experiments, we set $\lambda$ to 1 and $q$ to 6. These parameters will be explored further in the model study.

## 5 RESULTS

Table 2 reports the main results for our task over each of the crisis periods. From the table, we can make the following observations:

- The first two deep-learning models (*i.e.,* GRU+GAT, GPT+GNN) show results that lie close to 0.5. This is because the models predict mostly False (Note that All-True or All-False predictions produce an AUROC of exactly 0.5). Portfolio crashes are often rare, which causes a large bias towards False predictions when training the model. In addition, it is likely that the models cannot handle events that are previously unseen in the train set, causing them to predict mostly False on the unprecedented crash events.
- The DAN+FCL model shows results that are further from 0.5, as it can now deal with the imbalanced dataset. However, it is likely still unable to handle crash events not seen in its training data.
- On the other hand, the zero-shot LLM frameworks do not have these limitations. Furthermore, among the thought-based frameworks (*i.e.,* IO, CoT and GoT), we can see a rising trend in the AUROC. This shows that it is beneficial to break down the task of predicting portfolio crashes into smaller thought processes.
- Going further, we observe that the search-based ToG was able to outperform these models. As the input dataset was not manually filtered, it is likely that there are numerous news articles that contain noisy information not relevant to the specified portfolio. By first searching for an impact path from the articles to the portfolio, ToG was able to find the most relevant information that can help it to decide if a possible portfolio crash will occur.
- Finally, our TRR framework was able to outperform all models, by an average of 8.93% over the strongest baseline (ToG). By considering *multiple* impact paths that are relevant to the portfolio and also the relationship between these paths, TRR was able to get a more holistic overview of the various market forces on the portfolio, which can help it to predict crashes more accurately.
- In addition, TRR was also able to perform well in the 2022 dataset, which is after the post-knowledge cut-off date. This shows that the performance can be attributed to the reasoning process of the TRR framework, and not the internal knowledge of the LLM.

**Table 3: Performance comparison of TRR on stable periods.**

| Models | 2012 (Country-Neu) | | | | |
|---|---|---|---|---|---|
| | AUROC | TP | TN | FP | FN |
| GRU+GAT | 0.500 | 0.000 | 0.980 | 0.000 | 0.020 |
| ToG | 0.433 | 0.004 | 0.652 | 0.328 | 0.016 |
| TRR | 0.600 | 0.004 | 0.980 | 0.000 | 0.016 |
| Models | 2012 (Sector-Neu) | | | | |
| | AUROC | TP | TN | FP | FN |
| GRU+GAT | 0.500 | 0.000 | 0.992 | 0.000 | 0.008 |
| ToG | 0.401 | 0.000 | 0.796 | 0.196 | 0.008 |
| TRR | 0.500 | 0.000 | 0.992 | 0.000 | 0.008 |

Table 3 reports the performance of TRR over a longer, stable period. Here, it is harder to compare the AUROC, given that our framework also predicts mostly No Crash in this period. This experiment serves to highlight the robustness of TRR on normal days, as it does not produce False Positives. Note that for an event to trigger a crash prediction, it has to impact an entity over a period of time (Temporal component) and impact many portfolio stocks (Relational component), which is unlikely to occur on a regular day.

## 5.1 Ablation Study

To investigate the effectiveness of the TRR framework design, we perform ablation studies over its components and LLMs used.

Firstly, we remove the individual components in TRR, which include its relational, temporal (memory) and the memory decay. For removing the relational component, we repeat the results from the ToG experiment, given that generating a relational graph forms the backbone of the TRR model. In our implementation, ToG was used to search for a single path across the graph to make its prediction, and hence does not consider the relations between multiple paths. For removing the temporal component, we remove the memory module (we set $G_{temporal} = G$), hence not providing any past temporal context to the LLM for reasoning. Finally, for removing the decay component, we set the memory retention to a constant value (we set $R_{u,v} = 1$ for all $u, v$). Secondly, we evaluate TRR across its implementations using various state-of-the-art LLMs, to gauge their efficiency in the framework. These experiments were conducted over the dataset for year 2007, on the Country-Neutral portfolio.

**Table 4: Ablation study over the components in TRR.**

| Variant | AUROC |
|---|---|
| No relational | $0.502_{\pm 0.0589}$ |
| No temporal | $0.524_{\pm 0.0496}$ |
| No decay | $0.630_{\pm 0.0459}$ |
| TRR (Ours) | $\mathbf{0.690_{\pm 0.0426}}$ |

**Table 5: Ablation study over the LLM models for TRR.**

| Models | AUROC |
|---|---|
| Mixtral v0.3 | $0.483_{\pm 0.0000}$ |
| Gemma 2 | $0.500_{\pm 0.0000}$ |
| Llama 3.1 | $0.505_{\pm 0.0695}$ |
| GPT-4o | $0.509_{\pm 0.0000}$ |
| Qwen2 | $0.612_{\pm 0.0144}$ |
| GPT-3.5 Turbo | $\mathbf{0.690_{\pm 0.0426}}$ |

Table 4 reports the AUROC over the different components in TRR. From the table, we can observe that each component helped to contribute to the overall performance of the model. Both the relational and temporal components provide additional graph paths to the LLM during the reasoning phase without degrading the performance, which shows that useful information was extracted from the graph as opposed to noisy data. These paths help to provide additional contextual information regarding each news impact,

which allows the LLM to better determine if a crash will occur. We can also observe that the memory decay component helps to improve the AUROC results. It is likely that events that are further back in history would have a lesser impact on the stock prices, as these events start to fade away from public memory [36], making it effective to decay their importance weightage to reduce the noise.

Table 5 reports the AUROC across different LLMs used in TRR. At a low temperature setting, some models have more consistent results than others, resulting in zero standard deviation. We find that the newer LLMs, such as Llama 3.1 and GPT-4o, tend to ignore the graph of impacts and generalize directly on the news articles (*i.e.,* any mention of troubles in the housing market would lead to a Crash prediction, regardless of its relational or temporal context), which results in poorer overall performance. From the ablation study, we select GPT-3.5 Turbo to implement our TRR framework.

## 5.2 Parameter Selection

For choosing the parameters, we conduct a model study over different values of $\lambda$ and $q$, which determine the memory decay rate and top-$q$ entities that investors would pay attention to respectively, which affects the final graph $G_{TRR}$. The experiments were conducted over the 2007 dataset, on the Country-Neutral portfolio.

| $\lambda$ | AUROC |
|---|---|
| 0.1 | $0.615_{\pm 0.0311}$ |
| 0.5 | $0.627_{\pm 0.0295}$ |
| 1 | $\mathbf{0.690_{\pm 0.0426}}$ |
| 2 | $0.665_{\pm 0.0326}$ |
| 10 | $0.642_{\pm 0.0452}$ |

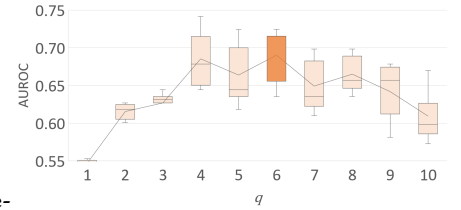**Table 6: Parameter selection of $\lambda$.**



**Figure 3: Parameter selection of $q$.**

Table 6 reports the AUROC over different values of the decay rate constant $\lambda$. Here, it is observed that the AUROC drops as $\lambda$ both increases and decreases. Note that in the memory module, the memory retention of each impact was tracked by the variable $R_{u,v} = \exp(-\frac{t_{u,v}}{\lambda})$, which is then used to weigh the ranking scores $PR(e_h)$. As $\lambda$ decreases, the memory retention tends towards very small values, *i.e.,* $\lim_{\lambda \to 0} R_{u,v} = 0$, which also causes the ranking scores for all entities to shift towards an equal value of zero. Because of this, the top entities for the attention graph $G_{att}$ would be chosen more randomly, resulting in poorer prediction performance. On the other hand, as $\lambda$ increases, the memory retention tends towards a value of one, *i.e.,* $\lim_{\lambda \to \infty} R_{u,v} = 1$. This causes the ranking scores to be weighted more equally, and the top entities selection for $G_{att}$ will become less affected by the temporal information. At very high values of $\lambda$, the weights will remain constant at 1, which is equivalent to having no decay in the memory component. Higher $\lambda$ causes the AUROC value to fall, highlighting the usefulness of the memory decay process. Through the ablation study, we set $\lambda = 1$.

Figure 3 shows the range of AUROC values over 5 runs, across different values of $q$. The value of $q$ determines the number of top $q$ entities to be considered in the attention graph $G_{att}$, which the LLM will perform reasoning on to make a crash prediction. We can observe that at small values of $q$, too limited information was provided to the LLM, resulting in poorer AUROC performance that is fairly consistent. However, as $q$ increases, more relevant information is provided for the LLM to perform reasoning over. This results
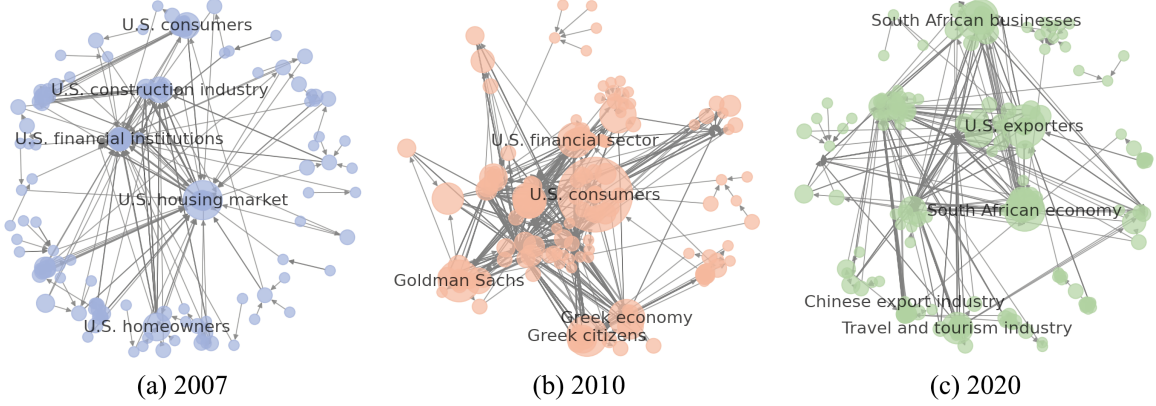
(a) 2007       (b) 2010       (c) 2020

**Figure 4: Examples of generated graphs during the crash periods in the 2007, 2010 and 2020 dataset.**

**Table 7: Portfolio Results for the three main datasets during the crash periods.**

| Portfolio Metrics | 2007 (Country-Neu) | | | | | 2010 (Country-Neu) | | | | | 2020 (Country-Neu) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | IO | ToG | TRR (Ours) | S&P500 | 1/N | IO | ToG | TRR (Ours) | S&P500 | 1/N | IO | ToG | TRR (Ours) | S&P500 | 1/N |
| Cumulative Ret | -0.0385 | -0.0015 | **0.1459** | -0.0406 | 0.0292 | -0.0871 | -0.0375 | **-0.0089** | -0.0113 | -0.0350 | -0.0650 | -0.0327 | **-0.0001** | -0.1937 | -0.1979 |
| Max Drawdown | 0.1329 | 0.0733 | **0.0356** | 0.0943 | 0.1548 | 0.1955 | 0.1188 | **0.1089** | 0.1227 | 0.1495 | 0.2233 | 0.1436 | **0.0496** | 0.3392 | 0.3337 |
| Sharpe Ratio | -0.0397 | 0.0000 | **0.2202** | -0.0486 | 0.0272 | -0.1040 | -0.0579 | -0.0086 | **-0.0078** | -0.0376 | -0.0266 | -0.0211 | **0.0026** | -0.0803 | -0.0963 |

in better AUROC but with a larger standard deviation, as the LLM could choose different parts of the information to focus on. Finally, at very high values of $q$, there is a dip in the AUROC performance, given that there might now be too much noisy information provided. This also highlights the importance of forming an attention sub-graph, instead of providing all information directly into the LLM, which could affect its performance. Here, we set $q = 6$.

## 5.3 Graph Analysis

In addition, we explore the generated graphs by visualizing an example from each dataset over the crash periods. For each graph, we project the vertex sizes based on the number of incoming edges. To prevent overcrowding of the labelled entities, we label only the top few vertices with the *highest* number of incoming edges.

Figure 4 showcases the graphs $G_{TRR}$ generated using the series of news articles, which are used by the LLM to detect portfolio crashes for the next day. From the graphs, we can qualitatively determine that TRR was able to highlight the most important information that caused the portfolio crashes in real life. Within the 2007 dataset, TRR was able to find that the U.S. housing market was impacted *more* than other entities from the given news, as shown from its higher number of incoming edges. This coincides with the global financial crisis in 2007, which was caused by the housing bubble. From the 2010 dataset, TRR was able to capture the impact on the Greece citizens and the Greek economy, which aligns with the Greek government debt crisis. In the 2020 dataset, the top impacted entities were less obvious as they were spread out over various entities. However, the impacted entities, such as the export and tourism industries, show the impacts that was caused by COVID-19.

## 5.4 Portfolio Analysis

To study the real-world performance of TRR, we also evaluate TRR on common portfolio metrics. To form our portfolio, we sell all the constituent stocks when the prediction switches to True (*i.e.*, Crash),

and buy all when the prediction switches to False. We use 0.02 as the transaction cost for each reallocation. We compare against some LLM methods, a market index (S&P500) and the $1/N$ portfolio [15]. For the deep-learning models, as they predict all False, they will hold all stocks without selling, which is similar to the $1/N$ portfolio.

Table 7 reports the portfolio metrics over the crash periods. We see that the TRR portfolio was able to avoid the most losses, achieving the highest cumulative returns and maximum drawdowns in a market downturn. The model did under-perform the market index on the Sharpe ratio over the 2010 dataset, but the values are close.

**Table 8: Portfolio Results for the stable period dataset.**

| Portfolio Metrics | 2012 (Country-Neu) | | | | |
|---|---|---|---|---|---|
| | IO | ToG | TRR (Ours) | S&P500 | 1/N |
| Cumulative Ret | 0.1335 | 0.1232 | **0.2697** | 0.2384 | 0.0982 |
| Max Drawdown | 0.1804 | **0.0849** | 0.1474 | 0.1474 | 0.0994 |
| Sharpe Ratio | 0.0541 | 0.0712 | **0.0975** | 0.0787 | 0.0515 |

Table 8 reports the portfolio metrics over a stable period. The TRR portfolio also beats the baselines on returns and Sharpe Ratio. For this dataset, we also note that there was only a single True Positive made by TRR (see Table 3). Since the model predicts large falls in value, a single correct prediction resulted in a big difference in the cumulative returns (13.1% increase over next best baseline).

## 5.5 Additional Experiments

Additionally, we explore the generalizability of TRR to a closely-related task. In Macroeconomics, it is a crucial task to develop warning indicators for economic crisis events [3, 45, 69], in order for policymakers to take preemptive measures to mitigate these events. By viewing the global economy as a *network of regional economies*, we can utilize our TRR framework to trace the impacts on each individual economy, then reason over these impacts to
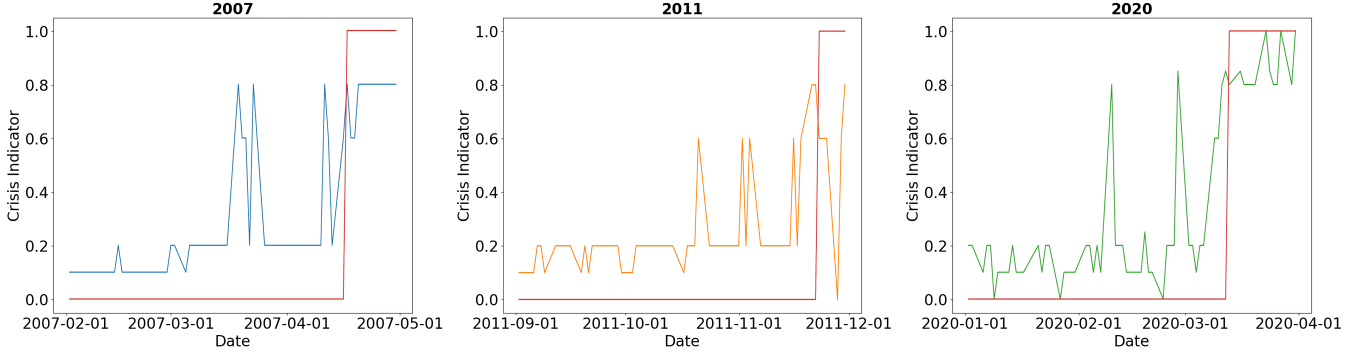
**Figure 5: Examples of the economic crisis indicator generated from our TRR framework. The line coloured red represents the ground truth crisis labels, *i.e.,* a value of 1 represent a crisis; 0 otherwise. The other coloured lines represent the crisis indicators.**

detect possible widespread crisis globally. For this task, instead of a binary prediction, we prompt the LLM to output the *probability* of a global crisis, for its application as a continuous warning indicator.

For this experiment, we now set our "portfolio" $P$ as a series of economies, *e.g.,* the American economy, European economy, Asian economy, *etc.* , and use TRR to trace the impacts of news events to these entities. For the crisis labels, we use the TED spread, which is the difference between the interest rates on interbank loans and short-term government debt. It has been shown in empirical works that a TED spread above 48 basis points is indicative of economic crisis [9], as lenders switch to safer government investments when they believe the risk of default on interbank loans is rising. We label the TED spread above 0.48 as a crisis. In addition, we also provide the past 5 days of TED spread data in the LLM prompt for context.

For the baselines, we use the Financial Stress Indicator (FSI) [45], the volatility index (VIX) [28] and the yield curve (Yield) values [24] from the previous day as indicators. In particular, the FSI is one of the earliest works in economics that use news articles as a financial indicator. This is done by searching the article headlines for financial keywords, such as "economy", "gold" or "railroads". The reporting frequency of these terms were found to correlate heavily with crisis events in the work [45]. More details on the implementation, dataset and baselines can be found in Appendix C.

**Table 9: Performance comparison on the global crisis detection task using AUROC. The baselines are calculated using a deterministic equation, resulting in no standard deviation.**

| Dataset | 2007 | 2011 | 2020 |
|---|---|---|---|
| FSI | $0.550_{\pm0.000}$ | $0.407_{\pm0.000}$ | $0.178_{\pm0.000}$ |
| VIX | $0.430_{\pm0.000}$ | $0.329_{\pm0.000}$ | $0.997_{\pm0.000}$ |
| Yield | $0.821_{\pm0.000}$ | $0.260_{\pm0.000}$ | $\mathbf{0.999_{\pm0.000}}$ |
| TRR (Ours) | $\mathbf{0.940_{\pm0.023}}$ | $\mathbf{0.747_{\pm0.013}}$ | $0.961_{\pm0.044}$ |

Table 9 reports the results for the global crisis detection task:

- Among most of the models (including TRR), the 2011 dataset typically produce the lowest AUROC results. The 2011 dataset corresponds with the Greek government debt crisis, whose impact was mostly contained within the European economy. While there were still some spillover effects on the other economies, they were not as prevalent as those from the global financial crisis or COVID-19. For our model, this makes it harder to capture the interconnectedness between the entities on the graphs.

- The keyword-based FSI seems to drop greatly in performance in the 2020 dataset. This is likely because the keywords used, such as "gold" or railroads", were not relevant in this period, which correspond to the COVID-19 event. Given that it is hard for humans to predict what event would cause the next crisis, it is also difficult to know what keywords to search for in advance. Hence, it becomes crucial to utilize tools that can do *zero-shot* reasoning over unseen events, such as LLMs, in this specific task.

- Our proposed TRR framework outperform most baselines, except for the 2020 dataset, where the economic indicators (*i.e.,* VIX and yield curve) showed exceptional predictive performance. However, for this case, our performance is still comparable to these methods, even without the use of statistical information. This highlights TRR as a possible useful tool for practitioners.

Figure 5 provides some qualitative examples of the crisis indicator generated from our TRR framework. In addition to the AUROC performance, we note that our indicator shows a peak at the *start* of the crisis date, which is also an important consideration. On the other hand, as the values for each day are generated individually, we note that the LLM does not have a sense of the continuity or smoothness in the data, which could limit its application as a continuous indicator. This can be further studied in a future work.

## 6 CONCLUSION AND FUTURE WORK

In this work, we explored the main task of portfolio crash detection, which was difficult to solve before the advent of LLMs, due to the unprecedented nature of crash-related events. We proposed our TRR framework, which is able to do zero-shot reasoning across relational and temporal information through a set of human cognitive capabilities. Through extensive experiments, we showed that TRR is able to outperform state-of-the-art frameworks on detecting portfolio crashes. Furthermore, we also explored the generalizability of TRR by using it to develop a global crisis warning indicator.

The results of this work open some possible future directions for research. Firstly, each component in TRR can be further improved. For example, the memory component in TRR can be augmented with a more advanced symbolic database [26]; the PageRank algorithm is also dated and can be replaced with newer information retrieval-based methods [58]. Secondly, for the crisis detection task, more baselines could be studied [3, 30], such as government debt, external trade, *etc.* These statistical indicators could also be used together with TRR in an ensemble system, which could help to improve the prediction capability of the overall warning indicator.

# REFERENCES

[1] D Araci. 2019. FinBERT: Financial Sentiment Analysis with Pre-trained Language Models. *arXiv preprint arXiv:1908.10063* (2019).

[2] Ching-man Au Yeung and Adam Jatowt. 2011. Studying how the past is remembered: towards computational history through large scale text mining. In *Proceedings of the 20th ACM international conference on Information and knowledge management*. 1231–1240.

[3] Jan Babecky, Tomas Havranek, Jakub Mateju, Marek Rusnak, Katerina Smidkova, and Borek Vasicek. 2013. Leading indicators of crisis incidence: Evidence from developed countries. *Journal of International Money and Finance* 35 (2013), 1–19.

[4] Eduard Baitinger and Jochen Papenbrock. 2017. Interconnectedness risk and active portfolio management. *Journal of Investment Strategies* (2017).

[5] Scott R Baker, Nicholas Bloom, Steven J Davis, Kyle Kost, Marco Sammon, and Tasaneeya Viratyosin. 2020. The unprecedented stock market reaction to COVID-19. *The review of asset pricing studies* 10, 4 (2020), 742–758.

[6] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big?. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*. 610–623.

[7] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. 2024. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 17682–17690.

[8] Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*. 1247–1250.

[9] Kris Boudt, Ellen CS Paulus, and Dale WR Rosenthal. 2017. Funding liquidity, market liquidity and TED spread: A two-regime model. *Journal of Empirical Finance* 43 (2017), 143–158.

[10] Chris Brooks and Gita Persand. 2000. Value-at-risk and market crashes. *Journal of Risk* 2 (2000), 5–26.

[11] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.

[12] Yuezhang Che, Shuyan Chen, and Xin Liu. 2022. Sparse index tracking portfolio with sector neutrality. *Mathematics* 10, 15 (2022), 2645.

[13] Zihan Chen, Lei Nico Zheng, Cheng Lu, Jialu Yuan, and Di Zhu. 2023. ChatGPT informed graph neural network for stock movement prediction. *arXiv preprint arXiv:2306.03763* (2023).

[14] Marcos Lopez De Prado. 2016. Building diversified portfolios that outperform out of sample. *The Journal of Portfolio Management* 42, 4 (2016), 59–69.

[15] Victor DeMiguel, Lorenzo Garlappi, and Raman Uppal. 2009. Optimal versus naive diversification: How inefficient is the 1/N portfolio strategy? *The review of Financial studies* 22, 5 (2009), 1915–1953.

[16] Xiao Ding, Yue Zhang, Ting Liu, and Junwen Duan. 2014. Using structured events to predict stock price movement: An empirical investigation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 1415–1425.

[17] Xiao Ding, Yue Zhang, Ting Liu, and Junwen Duan. 2015. Deep learning for event-driven stock prediction. In *Twenty-fourth international joint conference on artificial intelligence*.

[18] Xiao Ding, Yue Zhang, Ting Liu, and Junwen Duan. 2016. Knowledge-driven event embedding for stock prediction. In *Proceedings of coling 2016, the 26th international conference on computational linguistics: Technical papers*. 2133–2142.

[19] Fuli Feng, Xiangnan He, Xiang Wang, Cheng Luo, Yiqun Liu, and Tat-Seng Chua. 2019. Temporal relational ranking for stock prediction. *ACM Transactions on Information Systems (TOIS)* 37, 2 (2019), 1–30.

[20] Ian Goldin and Tiffany Vogel. 2010. Global governance and systemic risk in the 21st century: Lessons from the financial crisis. *Global Policy* 1, 1 (2010), 4–15.

[21] Klaus Greff, Sjoerd Van Steenkiste, and Jürgen Schmidhuber. 2020. On the binding problem in artificial neural networks. *arXiv preprint arXiv:2012.05208* (2020).

[22] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. 2017. On calibration of modern neural networks. In *International conference on machine learning*. PMLR, 1321–1330.

[23] Maurice Halbwachs. 2020. *On collective memory*. University of Chicago press.

[24] Campbell Russell Harvey. 1986. *Recovering expectations of consumption growth from an equilibrium model of the term structure of interest rates*. The University of Chicago.

[25] David Hirshleifer, Sonya S Lim, and Siew Hong Teoh. 2011. Limited investor attention and stock market misreactions to accounting information. *The Review of Asset Pricing Studies* 1, 1 (2011), 35–73.

[26] Chenxu Hu, Jie Fu, Chenzhuang Du, Simian Luo, Junbo Zhao, and Hang Zhao. 2023. Chatdb: Augmenting llms with databases as their symbolic memory. *arXiv preprint arXiv:2306.03901* (2023).

[27] Ziniu Hu, Weiqing Liu, Jiang Bian, Xuanzhe Liu, and Tie-Yan Liu. 2018. Listening to chaotic whispers: A deep learning framework for news-oriented stock trend prediction. In *Proceedings of the eleventh ACM international conference on web search and data mining*. 261–269.

[28] George J Jiang and Yisong S Tian. 2007. Extracting model-free volatility from option prices: An examination of the VIX index. *Journal of Derivatives* 14, 3 (2007).

[29] Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Xin Zhao, and Ji-Rong Wen. 2023. StructGPT: A General Framework for Large Language Model to Reason over Structured Data. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. 9237–9251.

[30] Òscar Jordà, Moritz Schularick, and Alan M Taylor. 2017. Macrofinancial history and the new business cycle facts. *NBER macroeconomics annual* 31, 1 (2017), 213–263.

[31] Kelvin J.L. Koa, Yunshan Ma, Ritchie Ng, and Tat-Seng Chua. 2023. Diffusion Variational Autoencoder for Tackling Stochasticity in Multi-Step Regression Stock Price Prediction. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*. 1087–1096.

[32] Kelvin J.L. Koa, Yunshan Ma, Ritchie Ng, and Tat-Seng Chua. 2024. Learning to Generate Explainable Stock Predictions using Self-Reflective Large Language Models. In *Proceedings of the ACM on Web Conference 2024*. 4304–4315.

[33] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems* 35 (2022), 22199–22213.

[34] Lancine Kourouma, Denis Dupre, Gilles Sanfilippo, and Ollivier Taramasco. 2010. Extreme value at risk and expected shortfall during financial crisis. *Available at SSRN 1744091* (2010).

[35] Wenhao Liang, Zhengyang Li, and Weitong Chen. 2024. Enhancing Financial Market Predictions: Causality-Driven Feature Selection. *arXiv preprint arXiv:2408.01005* (2024).

[36] Philipp Lorenz-Spreen, Bjarke Mørch Mønsted, Philipp Hövel, and Sune Lehmann. 2019. Accelerating dynamics of collective attention. *Nature communications* 10, 1 (2019), 1759.

[37] Tim Loughran and Bill McDonald. 2011. When is a liability not a liability? Textual analysis, dictionaries, and 10-Ks. *The Journal of finance* 66, 1 (2011), 35–65.

[38] Benoit B Mandelbrot. 2001. Scaling in financial prices: I. Tails and dependence. *Quantitative Finance* 1, 1 (2001), 113.

[39] Harry Markowitz. 1952. Portfolio Selection. *The Journal of Finance* 7, 1 (1952), 77–91. https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1540-6261.1952.tb01525.x

[40] Muhammad Abubakr Naeem, Saba Sehrish, and Mabel D Costa. 2021. COVID-19 pandemic and connectedness across financial markets. *Pacific Accounting Review* 33, 2 (2021), 165–178.

[41] Allen Newell, John C Shaw, and Herbert A Simon. 1959. Report on a general problem solving program. In *IFIP congress*, Vol. 256. Pittsburgh, PA, 64.

[42] OpenAI. 2023. GPT-3.5-turbo — OpenAI Platform Documentation. https://platform.openai.com/docs/models/gpt-3.5-turbo.

[43] Rafael Padilla, Sergio L Netto, and Eduardo AB Da Silva. 2020. A survey on performance metrics for object-detection algorithms. In *2020 international conference on systems, signals and image processing (IWSSIP)*. IEEE, 237–242.

[44] Lawrence Page, Sergey Brin, Rajeev Motwani, Terry Winograd, et al. 1999. The pagerank citation ranking: Bringing order to the web. (1999).

[45] Lukas Püttmann. 2018. Patterns of panic: Financial crisis language in historical newspapers. *Available at SSRN 3156287* (2018).

[46] K Geert Rouwenhorst. 1998. International momentum strategies. *The journal of finance* 53, 1 (1998), 267–284.

[47] Ramit Sawhney, Shivam Agarwal, Arnav Wadhwa, and Rajiv Shah. 2020. Deep attentive learning for stock movement prediction from social media text and company correlations. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 8415–8426.

[48] Evangelia Spiliopoulou. 2022. *Modeling Event Implications via Multi-faceted Entity Representations*. Ph. D. Dissertation. Carnegie Mellon University.

[49] Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Heung-Yeung Shum, and Jian Guo. 2024. Think-on-graph: Deep and responsible reasoning of large language model with knowledge graph. In *International Conference on Learning Representations*.

[50] N.N. Taleb. 2007. *The Black Swan: The Impact of the Highly Improbable*. Random House Publishing Group. https://books.google.com.sg/books?id=gWW4SkJjM08C

[51] Nassim Nicholas Taleb, Yaneer Bar-Yam, and Pasquale Cirillo. 2022. On single point forecasts for fat-tailed variables. *International Journal of Forecasting* 38, 2 (2022), 413–422.

[52] Akhil Vaid, Kipp W Johnson, Marcus A Badgeley, Sulaiman S Somani, Mesude Bicak, Isotta Landi, Adam Russak, Shan Zhao, Matthew A Levin, Robert S Freeman, et al. 2022. Using deep-learning algorithms to simultaneously identify right and left ventricular dysfunction from the electrocardiogram. *Cardiovascular Imaging* 15, 3 (2022), 395–410.

[53] Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Commun. ACM* 57, 10 (2014), 78–85.

[54] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.

[55] Florian Weigert. 2016. Crash aversion and the cross-section of expected stock returns worldwide. *The Review of Asset Pricing Studies* 6, 1 (2016), 135–178.

[56] Robert E Whaley. 2000. The investor fear gauge. *Journal of portfolio management* 26, 3 (2000), 12.

[57] Robert E Whaley. 2009. Understanding the VIX. *Journal of Portfolio Management* 35, 3 (2009), 98–105.

[58] Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie. 2021. Self-supervised graph learning for recommendation. In *Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval*. 726–735.

[59] Wenpu Xing and Ali Ghorbani. 2004. Weighted pagerank algorithm. In *Proceedings. Second Annual Conference on Communication Networks and Services Research, 2004.* IEEE, 305–314.

[60] Siheng Xiong, Ali Payani, Ramana Kompella, and Faramarz Fekri. 2024. Large Language Models Can Learn Temporal Reasoning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 10452–10470.

[61] Liao Xu, Xuan Zhang, and Jing Zhao. 2023. Limited investor attention and biased reactions to information: Evidence from the COVID-19 pandemic. *Journal of Financial Markets* 62 (2023), 100757.

[62] Yumo Xu and Shay B Cohen. 2018. Stock movement prediction from tweets and historical prices. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 1970–1979.

[63] Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Shaochen Zhong, Bing Yin, and Xia Hu. 2023. Harnessing the power of llms in practice: A survey on chatgpt and beyond. *ACM Transactions on Knowledge Discovery from Data* (2023).

[64] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2024. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems* 36 (2024).

[65] Chenhan Yuan, Qianqian Xie, Jimin Huang, and Sophia Ananiadou. 2024. Back to the future: Towards explainable temporal reasoning with large language models. In *Proceedings of the ACM on Web Conference 2024*. 1963–1974.

[66] Yichi Zhang, Zhuo Chen, Wen Zhang, and Huajun Chen. 2023. Making large language models perform better in knowledge graph completion. *arXiv preprint arXiv:2310.06671* (2023).

[67] Yue Zhang and Stephen Clark. 2011. Syntactic processing using the generalized perceptron and beam search. *Computational linguistics* 37, 1 (2011), 105–151.

[68] Zeyang Zhang, Xin Wang, Ziwei Zhang, Haoyang Li, Yijian Qin, and Wenwu Zhu. 2024. LLM4DyG: Can Large Language Models Solve Spatial-Temporal Problems on Dynamic Graphs?. In *Conference on Knowledge Discovery and Data Mining (ACM SIGKDD)*.

[69] Huanhuan Zheng. 2020. Coordinated bubbles and crashes. *Journal of Economic Dynamics and Control* 120 (2020), 103974.

[70] Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. 2024. Memorybank: Enhancing large language models with long-term memory. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 19724–19731.

## A  PORTFOLIO CONSTRUCTION

To construct our portfolios, we select the top market capitalization stocks from each country or sector that have historical data since 2007 (*i.e.,* newer big companies such as Alibaba or Spotify were not considered due to lack of data). We limit the portfolio sizes to 10 due to the LLM token limits, at the time of the experiments. The constituent stocks of the portfolios are found in Table 10 and 11.

### Table 10: Country-Neutral Portfolio.

| Company | Country |
| --- | --- |
| Apple Inc. | U.S. |
| Royal Bank of Canada | Canada |
| NetEase, Inc. | China |
| HDFC Bank Limited | India |
| AstraZeneca PLC | UK |
| TotalEnergies SE | France |
| SAP SE | Germany |
| Toyota Motor Corporation | Japan |
| BHP Group Limited | Australia |
| Accenture PLC | Ireland |

### Table 11: Sector-Neutral Portfolio.

| Company | Sector |
| --- | --- |
| HSBC Holdings PLC | Financials |
| Novo Nordisk A/S | Healthcare |
| ASML Holding N.V. | Technology |
| General Electric Company | Industrials |
| Amazon.com, Inc. | Consumer Discretionary |
| Linde PLC | Materials |
| Alphabet Inc. | Communication Services |
| Shell PLC | Energy |
| Unilever PLC | Consumer Staples |
| National Grid PLC | Utilities |

## B  BASELINE IMPLEMENTATIONS

In our work, we compare TRR to relevant deep-learning models and zero-shot LLM reasoning frameworks to justify its effectiveness.

For the deep-learning models, additional data was provided for training. The models are trained in the same manner as the task, where the labels are the imbalanced binary crashes. The training news data is taken from the same Reuters dataset, that is four times the size of the test data, on the dates just preceding it. For example, news from Feb 2009-Feb 2010 were used to train the models for Mar 2010-May 2010, while news from Jan 2008-Dec 2011 were used to train the model for Jan 2012-Dec 2012. The selected models are:

**GRU+GAT** [47]: In this model, historical price data for 5 days and company relational data from Wikidata [19] was provided, in additional to the historical news for 5 days. The model was trained over a batch size of 32 for 100 epochs, with a learning rate of 0.01.

**GPT+GNN** [13]: In this model, historical price data and news for 5 days were provided. A daily relational graph is dynamically generated from the news using ChatGPT, and is used to generate embeddings using GNN. The concatenated price and relational embeddings are then trained using Long Short-Term Memory (LSTM) networks to generate the final predictions. The model was trained over a batch size of 32 for 500 epochs, with a learning rate of 0.01.

**DAN+FCL** [35]: In this model, FinBERT [1] was used to generate aggregate sentiment scores from the news, which are used to train a 3-layer DAN [22] to generate predictions. The model was trained over a batch size of 32 for 100 epochs, with a learning rate of 0.01.

For Large Language Models, the baseline models have limited capabilities to tackle this task on their own (*e.g.,* limited token capacity for direct LLM prompting, no existing knowledge graph for ToG, *etc.* ), which also highlights the novelty of TRR. As such, some modifications have to be made to them for a fair comparison. We highlight our implementations of the LLM baselines below.

**Input-output (IO) prompting:** For this LLM baseline, we simply use the news articles as input and prompt the LLM to generate a crash prediction for the portfolio. However, as the full articles would contain a huge amount of text which goes beyond any LLM's token limits, we use only the headlines of the news articles as input.

**Chain-of-Thoughts (CoT)** [54]**:** This baseline largely follows the same methodology as IO prompting, but includes an additional line of prompt which instructs the LLM to "think step-by-step".

**Graph-of-Thoughts (GoT)** [7]**:** In this LLM baseline, we first provide the LLM with the retrieved news articles, and use GoT to divide and merge thoughts to arrive at a crash prediction for the target portfolio. This is done in a similar fashion as our TRR framework - we first split the portfolio into individual stocks, and prompt the LLM to discover the impact on each stock using sub-thoughts. These thoughts are then combined to find the overall impact on the portfolio. However, the key difference between this method and ours is that the thoughts are not used to form a structured temporal Knowledge Graph, and that there was no pruning of the combined information, which was done in the Attention phase of TRR.

**Think-on-Graph (ToG)** [49]**:** The ToG framework requires an existing factual Knowledge Graph, in order to search for a reasoning path that can tackle the task. To adapt it for our problem, we use the graph of impacts that was formed in our Brainstorming Phase. We then use ToG to find the best reasoning path on this graph to answer the prompt. Following the original work, this is done by identifying the most relevant paths at each depth via a beam search process, and checking if they are sufficient to make a crash prediction. This is repeated iteratively over each depth until the LLM respond that it has sufficient information, or the maximum depth of the graph is reached. When this happens, the most relevant path identified is then provided to the LLM to make its prediction.

## C  DETAILS OF ADDITIONAL EXPERIMENTS

We detail the experimental setup for the crisis detection task in this section. For this task, we redefine our portfolio of target entities as a set of regional economies, *i.e., $P$ = {American economy, European economy, Asian economy, African economy, Australian economy}*.

**Table 12: Statistics of the dataset for the crisis detection task. The time period was shifted to capture the TED cut-off.**

| Dataset | Time Period | Description | % Crisis (TED >0.48) |
|---|---|---|---|
| **2007** | 01 Feb 2007-30 Apr 2007 | Global Financial Crisis | 0.206 |
| **2011** | 01 Sep 2011-30 Nov 2011 | Greek government debt crisis | 0.138 |
| **2020** | 01 Jan 2020-31 Mar 2020 | COVID-19 stock market crash | 0.210 |

Using the same dataset of Reuters financial news articles, we then use our TRR framework to trace their impacts to these entities.

On this dataset, the selected time period was slightly shifted to capture the timestamps where the TED spread first goes above 48 points, resulting from the same crisis events as the portfolio crash detection task. The statistics of this dataset are found in Table 12.

For the baselines of this task, we choose a news article-based indicator and some financial indicators that are commonly used by practitioners to measure the overall economy health. We provide more information on each of these selected indicators below.

**Financial Stress Indicator (FSI)** [45]**:** The FSI is one of the earliest works in economics that uses newspaper articles as a financial indicator. The motivation was to observe what people are thinking *before* the crisis event (and hence reported in the news), and not what researchers find out in hindsight. The indicator is formed by first searching headlines for keywords from a list of 120 words and phrases [45] on economics and financial markets. Sentiment analysis [37] is then performed on the headlines to obtain a sentiment score on the economy for each day, which forms the indicator.

**Volatility Index (VIX)** [28]**:** The VIX is a measure of the stock market's volatility expectations by tracking the expected annualized change in the S&P 500 index for the next 30 days. It is also referred to by investors as the *fear gauge* [56], as it reflects investors' expectations of near-term market volatility. The VIX indicator has been consistently shown to spike over periods of crisis events [5, 57].

**Yield Curve** [24]**:** For the yield curve indicator, we use the difference between interest rates of the 10-year Treasury bond and the 2-year Treasury bond. An inverted yield curve (*i.e.,* when the 2-year bond yield exceeds the 10-year bond yield) typically signals an upcoming economic crisis [24], as more investors begin to believe that the risk of holding long-term government debt is higher.

The TED spread, VIX and yield curve data are all taken from the Federal Reserve Economic Data[1] (FRED), retrieved from the web.

## D  EXAMPLES OF LLM PROMPTS

We provide some examples of the LLM prompts and responses used in this work, in order for better replicability of our experiments.

Figure 6 and 7 show the initial and subsequent iterations of the Brainstorm prompt respectively. Given each news article from the dataset, the Brainstorm prompts will first generate the direct impacts, then trace their following knock-on effects repeatedly. This is done until the the chain of impacts reaches a portfolio stock, or the max number of iterations is reached. The responses are then used to form our graph of impacts. When forming our graph, we keep only the impacted entities and remove the explanations, which were only used to elicit a better reasoning process in the LLM [54].

Figure 8 shows the final generated temporal-relational graph $G_{TRR}$ in its tuple form, and the Reasoning prompt used for analyzing this graph. The tuples are organized first by their relational level in the graph, then temporally by their dates. Similar as above, the generated explanations are used for eliciting a better reasoning process in the LLM, and are not explicitly evaluated in our work. However in this case, the explanations are also used to ensure that the LLM does not refer to past events (*e.g.,* "*A portfolio crash will occur because the dates of the provided events correspond to the 2007 Global Financial Crisis...*"), but contains actual reasoning on the graph, to highlight the usefulness of the provided information.

---
[1]https://fred.stlouisfed.org/

Given a news article, name the entities that will be directly impacted negatively, and directly impacted positively. Give your response in the following format:
[Entity 1]: [Explanation]
...
[Entity n]: [Explanation]
For each entity, note any specific countries and industries, e.g. U.S. consumers or automobile industry.

Here is an example:
[…]
(END OF EXAMPLE)

Article:
American Home Mortgage tumbles on liquidity issues
NEW YORK  (Reuters) - American Home Mortgage Investment Corp AHM.N shares fell sharply on Monday after the company delayed its quarterly dividend, announced "major" write-downs, and said lenders were demanding it put up more cash.
[…]

Entities directly impacted negatively:
*American Home Mortgage Investment Corp: The company's shares fell sharply after it delayed its quarterly dividend, announced "major" write-downs, and said lenders were demanding it put up more cash. The company is struggling from "unprecedented" disruption in credit markets, and the announcement fed investor worries that bad loans are extending beyond "subprime" lenders to providers of home loans to borrowers considered to be good credit risks. The liquidity crisis gripping the mortgage funding markets has made it "very difficult for American Home to conduct an orderly flow of business."*

*Mortgage lenders: […]*

**Figure 6: Initial iteration of the Brainstorming prompt given each news article from the dataset. Italicized text refers to the LLM-generated response. [...] refers to text that are truncated in the given example.**

Given an impact on an entity, name the entities that will be impacted negatively, and impacted positively, by knock-on effects. Give your response in the following format:
[Entity 1]: [Explanation]
...
[Entity n]: [Explanation]
For each entity, note any specific countries and industries, e.g. U.S. consumers or automobile industry.
Do not repeat the same entities.

Here is an example:
[…]
(END OF EXAMPLE)

Entity: American Home Mortgage Investment Corp
Impact: The company's shares fell sharply after it delayed its quarterly dividend, announced "major" write-downs, and said lenders were demanding it put up more cash. The company is struggling from "unprecedented" disruption in credit markets, and the announcement fed investor worries that bad loans are extending beyond "subprime" lenders to providers of home loans to borrowers considered to be good credit risks. The liquidity crisis gripping the mortgage funding markets has made it "very difficult for American Home to conduct an orderly flow of business."

Entities impacted negatively by knock-on effects:
*U.S. housing market: The struggles of American Home Mortgage Investment Corp could further contribute to the instability in the housing market, potentially leading to decreased home sales and declining property values.*
*U.S. financial institutions: […]*

**Figure 7: Subsequent iterations of the Brainstorming prompt given each self-generated impact. Italicized text refers to the LLM-generated response. [...] refers to text that are truncated in the given example.**

Target Portfolio: ['AAPL', 'RY', 'NTES', 'HDB', 'AZN', 'TTE', 'SAP', 'TM', 'BHP', 'ACN']

Given a Temporal Knowledge Graph of recent event impacts, is it likely that the target portfolio will suffer a crash on the following day? Give your answer in this format:

Crash: [Yes/No]
Explanation: [Explanation]

Use your own reasoning on the Knowledge Graph and do not make reference to past events.

Temporal Knowledge Graph:
(2007-07-24, Cracks showing in keeping housing woes contained, impacts, U.S. home builders)
(2007-07-26, HSBC braces for $9 billion in mortgage price hikes, impacts, HSBC Holdings Plc)
(2007-07-26, Mortgage delinquencies seen peaking in 2008, impacts, Investors in mortgage-backed securities)
(2007-07-26, Wall St indexes sell off on housing, credit fears, impacts, Financial shares)
(2007-07-26, Wall St indexes sell off on housing, credit fears, impacts, U.S. home builders)
(2007-07-27, Stocks drop more than 1 percent on credit concerns, impacts, U.S. stocks)
(2007-07-30, American Home Mortgage tumbles on liquidity issues, impacts, American Home Mortgage Investment Corp)
(2007-07-30, American Home Mortgage tumbles on liquidity issues, impacts, Mortgage lenders)
[…]
(2007-08-08, GM lowers 2007 forecast for U.S. industry sales, impacts, U.S. automakers)
(2007-08-08, Home loan demand surges in week as rates fall, impacts, U.S. housing market)
(2007-08-08, Home sales to slip further this year, impacts, Investors in the mortgage finance market)
(2007-08-08, Toyota sees lower industrywide '07 U.S. auto sales, impacts, U.S. automakers)

(2007-08-08, American Home Mortgage, impacts, American Home Mortgage employees)
(2007-08-08, American Home Mortgage Investment Corp, impacts, U.S. housing market)
(2007-08-08, American Home Mortgage Investment Corp, impacts, U.S. real estate market)
(2007-08-08, American Home Mortgage Investment Corp., impacts, Subprime mortgage lenders)
(2007-08-08, American Home Mortgage employees, impacts, U.S. housing market)
(2007-08-08, Mortgage lenders, impacts, U.S. financial institutions)
[…]
(2007-08-08, U.S. automakers, impacts, U.S. automotive workers)
(2007-08-08, U.S. automakers, impacts, U.S. labor market)
(2007-08-08, U.S. automakers, impacts, U.S. oil and gas industry)
(2007-08-08, U.S. automakers, impacts, U.S. steel industry)

(2007-08-08, U.S. housing market, impacts, ['RY', 'AZN', 'SAP', 'HDB', 'ACN', 'TM', 'TTE', 'AAPL', 'NTES'])
(2007-08-08, U.S. financial institutions, impacts, ['RY', 'HDB', 'AAPL'])
[…]
(2007-08-08, U.S. oil and gas industry, impacts, ['TM', 'TTE'])
(2007-08-08, U.S. steel industry, impacts, ['TM', 'AAPL'])

Crash: *Yes*
*Explanation: The Knowledge Graph contains multiple references to the U.S. housing market, U.S. housing sector, U.S. homeowners, U.S. construction industry, U.S. real estate market, and U.S. mortgage lenders. These references indicate a significant impact on the housing and mortgage sectors in the U.S. Given that the target portfolio includes companies that are heavily tied to the U.S. economy, it is likely that the portfolio will suffer a crash.*

**Figure 8: The generated graph $G_{TRR}$ in tuple form and the Reasoning prompt. Italicized text refers to the LLM-generated response. [...] refers to text that are truncated in the given example.**