

Toward Fine-grained Causality Reasoning and CausalQA

Anonymous ACL submission

Abstract

Understanding causality is key to the success of NLP applications, especially in high-stakes domains. Causality comes in various perspectives such as *enable* and *prevent* that, despite their importance, have been largely ignored in the literature. This paper introduces a first-of-its-kind, fine-grained causal reasoning dataset that contains seven causal relations and defines a series of NLP tasks, from causality detection to event causality extraction and causal reasoning. Our dataset contains human annotations of 25K cause-effect event pairs and 24K question-answering pairs within multi-sentence samples, where each can contain multiple causal relationships. Through extensive experiments and analysis, we show that the complex relations in our dataset bring unique challenges to state-of-the-art methods across all three tasks and highlight potential research opportunities, especially in developing “causal-thinking” methods.

1 Introduction

Causality has received much research attention in recent years (Gao et al., 2019a; Schölkopf et al., 2021; Feder et al., 2021; Scherrer et al., 2021). It has been shown that causal reasoning entails a new goal of building more powerful AI systems beyond making predictions using statistical correlations (Kaushik et al., 2021; Srivastava et al., 2020; Li et al., 2021). In particular, understanding fine-grained causal relations between events in a document is an important step in language understanding and is beneficial to various NLP applications – information extraction, question answering, and machine reading comprehension, especially in high-stakes domains such as medicine and finance. Much work has been done on detecting a shallow “cause” relationship automatically for text (Khoo et al., 1998; Mirza et al., 2014; Chang and Chen, 2019; Mariko et al., 2020b). However, a single “cause” relationship cannot cover a plethora of causal concepts in the real-world scenarios reported

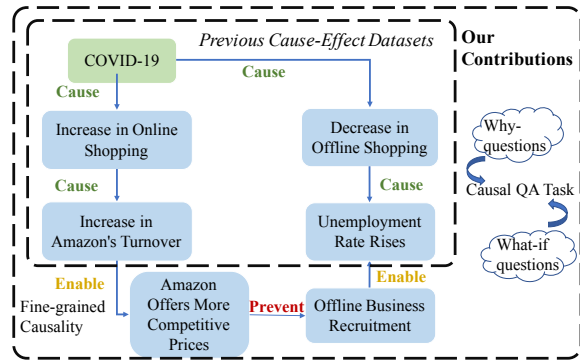


Figure 1: Comparisons between existing datasets and our dataset in the task of event causality analysis.

by the previous psychology research (Talmy, 1988; Wolff et al., 2005). For example, the spread of COVID-19 has led to the boom in online shopping – i.e., (cause) – but it also has deterred – i.e. (prevent) – people from going shopping-centres. According to classical psychology (Wolff and Song, 2003), it is important to understand possible fine-grained relationships between two events from three different causal perspectives, including *cause*, *enable*, and *prevent*.

In this paper, we construct a large-scale, hand-labeled, fine-grained causal reasoning (FineCR) dataset in the financial domain. A contrast between our dataset and the previous causality detection dataset is shown in Fig. 1. As can be seen, given the same passage “COVID-19 has accelerated change in online shopping, and given Amazon’s ... it will result in economic returns for years to come and offering more competitive prices compared to an offline business that brings pressures for the offline business recruitment.”, previous work can extract facts such as “COVID-19 causes an increase in online shopping”, yet cannot detect the subsequence for Amazon to “offer more competitive prices”, and further the negative influence on offline business recruitment, both of which can be valuable for predicting the future events.

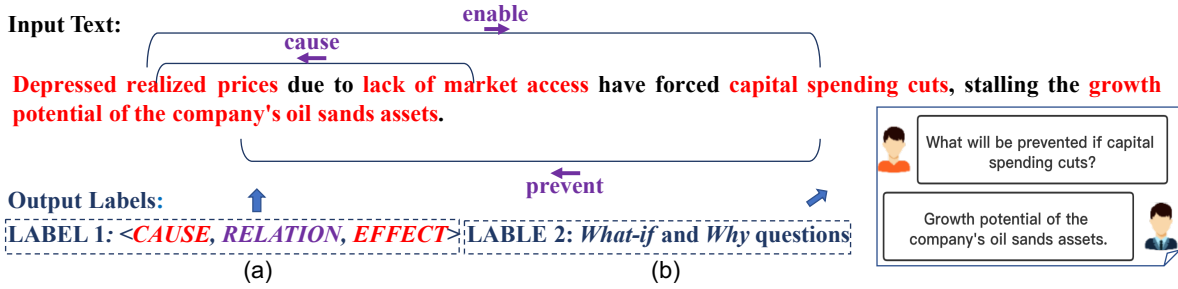


Figure 2: Illustration of our crowdsourcing tasks using an example that contains all three types of causal relationship.

069 Our dataset can also potentially benefit down-
070 stream applications such as financial analysis (Ding
071 et al., 2015) and BioNLP (Demner-Fushman et al.,
072 2021). As one practical application, we investigate
073 the utility of our dataset for the causal question
074 reasoning (CausalQA) task (Oh et al., 2016) in this
075 work. The resulting dataset, FineCR, consists of
076 25, 193 cause-effect pairs and 24, 486 question-
077 answering pairs, in terms of almost all questions
078 involving “why” and “what-if” scenarios belonging
079 to three fine-grained causalities. To establish the
080 benchmark performance on FineCR, which consists
081 of causality detection, fine-grained causality
082 extraction, and CausalQA tasks, we explore sev-
083 eral state-of-the-art neural models. Experimental
084 results show a significant gap between machine and
085 human ceiling performance (74.1% vs. 90.53% ac-
086 curacy in fine-grained classification). To the best of
087 our knowledge, FineCR is the first human-labeled
088 fine-grained event causality dataset, and we define
089 a novel CausalQA task based on that.

090 2 Dataset

091 We collected a financial analyst report dataset
092 from Yahoo Finance¹, which contains 6,786 well-
093 processed articles between December 2020 and
094 July 2021. Each instance corresponds to a specific
095 financial analyst report on a U.S. listed company,
096 which highlights the financial strengths and weak-
097 nesses of the company business.

098 2.1 Crowdsourcing

099 The original FineCR dataset consists of 6,786 arti-
100 cles in 54, 289 sentences. We employ editors from
101 a crowd-sourcing company to complete several hu-
102 man annotation tasks. Several pre-processing steps
103 required crowd-sourcing efforts were carried out
104 to prepare the raw dataset, including (1) A binary

¹We have received the written consent from the Yahoo Finance.

Metric	Counts
Causality Sentence Classification (Task1)	
#Positive Instances	21,046
#Negative Instances	29,979
#Multi-sentence Samples	846
#Average Token Length of POS Samples	42.8
#Average Token Length of NEG Samples	41.3
Cause-Effect Event Pairs (Task2)	
#Causal Text Chunks	45, 710
#Uni-causal Text Spans	18, 457
#Multi-causal Text Spans	3, 017
#Average Token Length of Cause Spans	16.0
#Average Token Length of Effect Spans	15.2
CausalQA Pairs (Task3)	
#Total Number of QA pairs	24, 486
#Average Token Length of Context	191.7
#Average Token Length of Questions	20.1
#Average Token Length of Answers	15.4
#Variance of Answer Length	69.15

Table 1: Statistics for causality detection, cause-effect pairs and QA pairs.

105 classification task for the causality detection; (2)
106 Mark the cause and effect formatted as text chunks
107 (a given instance may contain multiple causal rela-
108 tions), and give event pairs a fine-grained causal-
109 ity label, including *cause*, *cause_by*), *enable*, *en-
110 able_by*), *prevent*, *prevent_by*), and *irrelevant* re-
111 lation, where the suffix “_by” means the effect
112 comes before its cause; (3) Generating the follow-
113 ing question-answering dataset by using the labeled
114 event triples.

115 **Causality Detection.** We first focus on a bi-
116 nary classification task of the causality detection,
117 as such, removed sentences with outcome types of
118 non-causal relationships, leaving only those text se-
119 quences (one or two sentences) that are considered
120 containing at least a causal relation.

121 **Fine-grained Event Causality.** Given the sen-
122 tences each containing at least one event causal-
123 ity, human annotators are required to highlight all
124 the event causalities and give each instance a fine-
125 grained label. As shown in Fig. 2(a), a single
126 sentence can have more than one event causality,

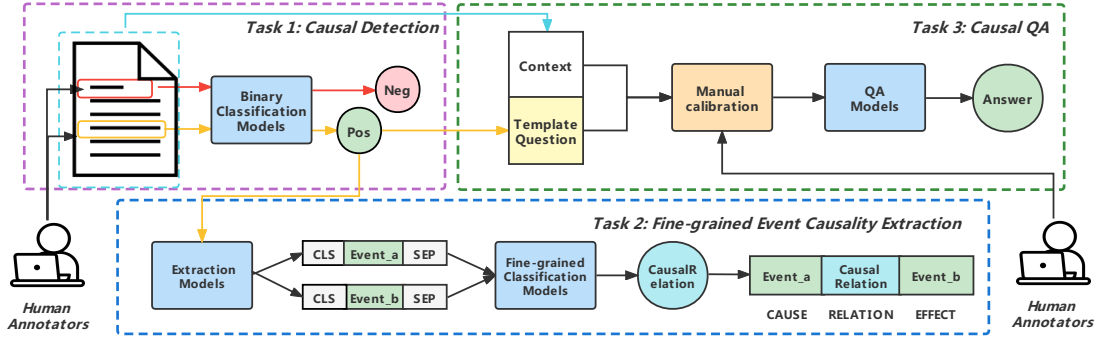


Figure 3: The pipeline of experiments based on the FineCR dataset.

127 which will be stored as triples containing \langle *cause*,
 128 *relation*, *effect* \rangle .

129 **CausalQA.** As shown in Fig. 2(b), we design
 130 a novel and challenging causal reasoning QA task
 131 based on the fine-grained causality labels. We ex-
 132 pand each \langle *cause*, *relation*, *effect* \rangle triple for
 133 generating a plausible question-answer pair. Differ-
 134 ent templates have been designed for different types
 135 of questions. For example, the active causal rela-
 136 tions – *CAUSE*, *ENABLE*, and *PREVENT* – could
 137 usually be used for generating why-questions while
 138 the corresponding passive causal relations could be
 139 used for generating what-if questions.

140 **Quality Control.** To ensure high quality, we
 141 restricted the participants to experienced human la-
 142 belers with relevant records. For each task, we con-
 143 ducted pilot tests before the crowd-sourcing work
 144 officially began, receiving feedback from quality
 145 inspectors and revising instructions accordingly.
 146 We filter out the sentences regarding the estima-
 147 tion of the stock price movement due to the nat-
 148 urally high-sensitive features and uncertainty of
 149 the complex financial market. After the first-round
 150 annotation (half of the data), we manually orga-
 151 nized spot checks for 10% samples in the dataset
 152 and revised the incorrect labels. After review, we
 153 revised roughly 3% of instances and refused the la-
 154 belers with above 10% error rate from participating
 155 in the second-round data annotation. Finally, the
 156 inter-annotators agreement ratio is 91% for fine-
 157 grained causality labels, and the F1 score of the
 158 inter-annotators agreement ratio is 0.94 for causal
 159 question-answer pairs.

160 Finally, we obtained a dataset of 51,025 in-
 161 stances (21,046 contain at least one causal relation)
 162 with fine-grained labels of cause-effect relations
 163 that were subsequently divided into training, valida-
 164 tion, and testing sets for the following experiments.
 165 It may be worth noting that we sort the dataset in

166 chronological order because the future data is not
 167 expected to be used for predictions.

168 2.2 Discussion

169 The primary data statistic of the FineCR dataset
 170 is shown in Table 1 for three different tasks. We
 171 observe that there is no significant difference in
 172 the average token numbers between positive and
 173 negative examples for Task 1 and 2, which shows
 174 that predictive models are difficult to learn from
 175 shortcut features (Sugawara et al., 2018, 2020; Lai
 176 et al., 2021) (e.g., the instance length) during the
 177 training process. Furthermore, our dataset con-
 178 tains 846 multi-sentence samples, and 3,017 text
 179 chunks contain more than one causal relation in
 180 one instance, which requires a complex reasoning
 181 process to get the correct answer, even for a hu-
 182 man. Most importantly, unlike other QA datasets
 183 (Sugawara et al., 2018, 2020) that can easily benefit
 184 from the test-train overlap as revealed by (Lewis
 185 et al., 2021a; Liu et al., 2021; Wang et al., 2021),
 186 our dataset is sorted in chronological order so that
 187 the future test data could be theoretically difficult
 188 to coincide with the training set. This allows us to
 189 obtain greater insight into what extent models can
 190 actually generalize.

191 2.3 Meta-information

192 Our dataset contains multi-sentence instances
 193 with fine-grained causality labels and the meta-
 194 information (company names and published dates).
 195 As shown in Fig. 4, we list the number of financial
 196 documents from different sectors, where the top
 197 three largest sectors belong to Consumer Cyclical,
 198 Industrial, and Technology. In contrast, compa-
 199 nies from the Utilities are the smallest group in
 200 our dataset. The use of the meta-information is
 201 two-fold. First, we choose the top three largest
 202 domains to perform out-of-domain evaluations (see

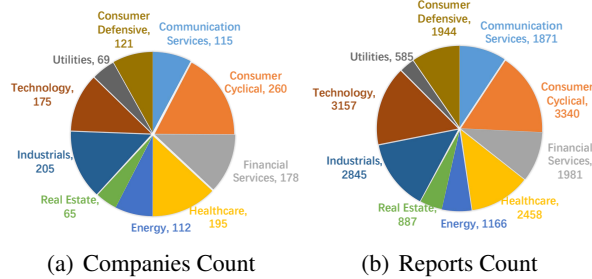


Figure 4: Sector distributions on companies and reports.

Appendix A). Second, company names would be used for generating question templates. Besides, the meta-information is crucial for benefiting the potential applications in NLP related to the domain of Finance (Tang et al., 2021; Chen et al., 2021).

3 Tasks and Methods

The pipeline of our experiments is shown in Fig. 3. We define three tasks on our FineCR dataset and build strong benchmark results for each task. First, as a prerequisite, models are evaluated on a binary classification task to predict whether a given text sequence contains a causal relation (**Task 1**). Second, we set up a joint event extraction and fine-grained causality task for identifying text chunks describing the cause and effect, respectively, and which fine-grained causality category it belongs to (**Task 2**). Finally, we design cause-effect question answering (**Task 3**).

3.1 Data Settings

For hold-out evaluation, we split our dataset into mutually exclusive training/validation/testing sets in the same ratio of 8:1:1 for all tasks. Predictive models and data splitting strategies have been kept the same among these tasks for building the benchmark results of each task. In line with the best practice, model hyper-parameters are tuned using the validation set. Both validation results and testing results will be reported in experiments.

3.2 Models

We consider using both classical deep learning models – CNN-Test (Kim, 2014) and HAN (Yang et al., 2016) – and Transformer-based models downloaded from Huggingface² – BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), and SpanBERT (Joshi et al., 2020) – as predictive models.

²<https://github.com/huggingface/models>

In addition, we perform a causal reasoning QA task by leveraging six Transformer-based pre-trained models provided by Huggingface (Wolf et al., 2020) on our dataset, including BERT-base, BERT-large, RoBERTa-base, RoBERTa-large, RoBERTa-base-with-squad, and RoBERTa-large-with-squad³. Furthermore, pre-trained seq2seq models such as T5 (Raffel et al., 2020), or BART (Lewis et al., 2020) are fine-tuned on QA-pairs as the benchmark methods of the generative QA tasks. In particular, we consider T5-small, T5-base, T5-large, BART-base, and BART-large models for building the benchmark results.

We use Adam as the optimizer and adopt the trick of decay learning-rate with the steps increase to train our model until converging for all models.

3.3 Methods

Our methods are built on the recently advanced Transformer architectures (Vaswani et al., 2017a) with the framework provided by Huggingface⁴. As follows, we introduce the detailed implementation of deep neural methods on three tasks completed on the FineCR dataset.

3.3.1 Causal Detection

The dataset consists of instances labeled with positive for the binary classification task if a given instance contains one causal relation and negative non-causal instances. The input data is extracted from the raw dataset directly, which contains 846 multi-sentence samples. We include multi-sentence samples besides a single sentence because causality could be found in multi-sentence contexts.

3.3.2 Fine-grained Event Causality Extraction

We use [CLS] and [SEP] to mark the event’s begin and end positions, respectively. For example, we have “[CLS] Better card analytics, increased capital markets and M&A offerings, and bolt-on acquisitions [SEP] should help drive [CLS] growth in fee income [SEP]”, in which “Better card analytics, increased capital markets and M&A offerings, and bolt-on acquisitions” is a “cause” event and “growth in fee come” is an “effect” event. In total, there are 33, 634 event triples with seven-class labels used for our experiments.

Then, we conduct the cause-effect extraction between samples. We consider cause-effect extrac-

³<https://huggingface.co/navteca/roberta-large-squad2>

⁴<https://github.com/huggingface/transformers>

Methods	Dev.		Test	
	F1	Acc	F1	Acc
CNN-Text	81.35	81.59	80.03	81.01
HAN	81.18	81.23	80.60	81.26
BERT-base	83.72	84.23	84.02	84.43
BERT-large	84.03	84.41	84.63	84.90
SpanBERT-base	84.09	84.38	84.51	84.72
SpanBERT-large	84.43	84.80	84.55	84.82
RoBERTa-base	84.59	85.16	84.31	84.76
RoBERTa-large	84.39	84.75	84.64	84.89
Human	-	-	94.32	95.94
Best Results @FinCausal	-	-	97.75	97.76

Table 2: The results of the causal sentence classification. 'F1' refers to the Macro F1. 'ACC.' is short for the accuracy.

tion a multi-span event extraction task as complex causal scenarios containing multiple causes or effects within a single instance is under consideration. We set the label for the first token of a cause or effect to "B", the rest of the tokens within the detected text chunks are given the label "I", and the other words in a given instance are set to "O". The results of event causality extraction are reserved for generating causal questions.

3.3.3 CausalQA

Both extractive methods and generative methods have been evaluated. For the extractive QA task, we adopt the same methods as the previous Transformer-based QA works (Kayesh et al., 2020). In particular, we first convert the context $C = (c_1, c_2, \dots, c_l)$ and question $Q = (q_1, q_2, \dots, q_{l'})$ into a single sequence $X = [CLS] c_1 c_2 \dots c_l [SEP] q_1 q_2 \dots q_{l'} [SEP]$, passing it to the pre-trained Transformer encoders for predicting the answer span boundary (start and end).

3.4 Metrics

The F1-score and accuracy are used for evaluating the event causality analysis task, and the exact match and F1-score are used for CausalQA.

The Macro F1-score is defined as the mean of label-wise F1-scores:

$$\text{Macro F1-score} = \frac{1}{N} \sum_{i=0}^N \text{F1-score}_i \quad (1)$$

where i is the label index and N is the number of classes.

4 Results and Discussion

We present the results of **Tasks 1-3** based on our dataset in this section.

4.1 Causality Detection

The causal detection result is shown in Table 2. We find that although Transformer-based methods achieve much better results than other methods – CNN and HAN using ELMO embeddings – on judging whether an instance contains at least a causal relationship (RoBERTa-Large can get the highest F1 Score – 84.64), it is still significantly below the human performance (84.64 vs. 94.32). The results of human performance are reported by quality inspectors from the crowdsourcing company. It is worth noting that the best results on the FinCausal (Mariko et al., 2020b) dataset can reach the human-level result (F1 = 97.75), providing indirect evidence that our dataset is more challenging caused by more complex causality instances.

4.2 Fine-grained Event Causality Extraction

The results of the fine-grained event causality extraction task are shown in Table 3. We find that SpanBERT and RoBERTa model can achieve the best performance for event causality extraction (F1 = 86.82 and EM = 60.26) and fine-grained classification (F1 = 68.99 and EM = 74.09), respectively. Nevertheless, all methods perform dramatically worse on the more challenging joint task, where the prediction is judged true only if event extraction and classification results exactly match the ground truth. Although the SpanBERT-large model can achieve the highest 21.78 EM on the test set, there is still much room for improvement.

We find that the large Transformer-based models (Vaswani et al., 2017b) with larger parameter sizes could not improve the performance on these tasks based on the FineCR dataset by comparing the test performance of BERT-base (63.72 in F1, 71.72 in ACC) with BERT-large (60.24 in F1, 69.85 in ACC) on the task of the fine-grained classification. It sheds new light that increasing the parameter size could not be helpful for causal reasoning tasks.

A more detailed error analysis by using the best-performed RoBERTa-Large model is given in Table 4. The model performs well in terms of the F1 score when predicting simple causal relations – Irrelevant (84.40), Cause (74.00), Cause_by (79.62), and Prevent (76.34). In contrast, complex relations – Enable (62.61) and Enable_by (41.49) – and the category with few examples – Prevent_by (64.46) – are not well predicted.

Model	Event Causality Extraction				Fine-grained Classification				Joint Evaluation	
	Dev.		Test		Dev.		Test		Dev.	Test
	F1	EM	F1	EM	F1	ACC	F1	ACC	EM	EM
BERT-base	84.37	51.48	85.30	53.53	71.74	70.43	63.72	71.72	21.21	20.15
BERT-large	85.13	50.34	86.93	52.88	70.90	64.16	60.24	69.85	17.54	21.73
RoBERTa-base	85.67	53.41	86.32	56.04	73.09	68.37	65.99	71.63	20.45	19.08
RoBERTa-large	85.12	54.70	85.95	56.77	74.54	71.99	68.99	74.09	20.46	19.77
SpanBERT-base	85.84	55.40	86.82	57.26	71.18	68.40	63.73	70.52	21.17	21.09
SpanBERT-large	85.50	57.40	86.33	60.26	73.65	68.15	64.43	72.93	23.01	21.78
Human Performance	-	-	94.32	81.34	-	-	88.61	90.53	-	-

Table 3: The results of the joint event causality detection (task2), 'F1' refers to the Macro F1. 'ACC.' is short for the accuracy, 'EM' refers to exact match and spe.

Category	Counts	Dev.		Test	
		F1	Acc	F1	Acc
Irrelevant	8,441	84.17	86.49	84.40	85.55
Cause	8,428	73.60	73.93	74.00	76.61
Cause_By	7,437	80.21	84.94	79.62	83.69
Enable	5,506	63.42	60.91	62.61	58.70
Enable_By	2,367	47.66	41.06	<i>41.49</i>	<i>35.68</i>
Prevent	1,086	79.79	71.43	76.34	67.87
Prevent_By	369	55.88	52.78	64.46	65.00

Table 4: Error analysis for fine-grained classifications.

CausalQA	Dev.		Test	
	F1	EM	F1	EM
BERT-base	79.90	55.52	79.33	55.70
BERT-large	82.48	59.24	82.37	58.71
RoBERTa-base	82.96	60.13	83.11	60.33
<i>SQuAD2.0-only</i>	<i>64.87</i>	<i>26.71</i>	<i>65.20</i>	<i>27.36</i>
SQuAD2.0-enhanced	84.39	61.22	84.34	61.17
RoBERTa-large	84.28	61.69	84.35	61.76
<i>SQuAD2.0-only</i>	<i>63.99</i>	<i>26.02</i>	<i>63.82</i>	<i>25.26</i>
SQuAD2.0-enhanced	84.65	61.63	84.65	61.58
Generative Methods				
BART-base	74.34	35.81	74.35	36.16
BART-large	65.52	27.24	65.70	26.48
T5-small	75.98	42.31	76.40	41.61
T5-Large	81.95	48.17	81.77	47.43

Table 5: The results of causal reasoning QA using both extractive methods and generative methods. 'SQuAD2.0' refers to the evaluation results using the model trained with the training set of SQuAD2.0⁵ only.

4.3 CausalQA

We provide both quantitative analysis and qualitative analysis for CausalQA. In addition, we compare the best performance on our dataset and other popular QA datasets.

4.3.1 Quantitative Analysis

The results of CausalQA are given in Table 5, where the bold values indicate the best performance while the italic values show the results of transfer learning methods trained by the SQuAD2.0 training data only. We find that the best-performing generative model – T5-Large – can achieve comparable results with the RoBERTa-large in terms of the F1

(81.77 vs. **84.35**). Meanwhile, the average EM of generative methods is largely below the extractive methods using the same training data. Second, the results of models trained with SQuAD2.0 data are much worse than those models trained with the original FineCR training set in terms of the F1 score (65.20 vs. **83.11** for RoBERTa-base and 63.82 vs. **84.35** for RoBERTa-large). On the other hand, we note a distinct improvement of using SQuAD2.0 data for initially training for both RoBERTa-base (from 83.11 to **84.34**) and RoBERTa-large (from 84.35 to **84.65**), which indicates that the training with additional well-labelled data could bring significant benefits for CausalQA. This may hint that the current QA data sources are still helpful for improving the performance of the causal reasoning QA task, although further research is required, as to what extent models can actually benefit from the additional data for the generalization is hard to be evaluated.

4.3.2 Qualitative Analysis for Answers

Table 6 presents a qualitative analysis for CausalQA, where we highlight the question and answer parts extracted from the raw context. Human labelers label the gold answers while the BERT-based model generates the output answers. The first three questions are answered correctly by the model, while the last two instances show two typical patterns prone to errors. In the first incorrect example, the model outputs “*targeted marketing*” using the keyword “*through*” but fails to give the gold answer “*analyzing the data and applying artificial intelligence*”. This could be because the model fails to identify the difference between the same word appearing in two different positions. The last example shows that the model tends to output the answer closer to the question in the context instead of observing the whole sentence. The real reasons – “*equity and credit markets*” and “*Brexit*”

Context	Question	Gold Answer	Output Answer
(Relation: Cause) Amazon’s 2017 purchase of Whole Foods remains a threat ... The COVID-19 outbreak has lifted near-term revenue as shoppers spend more time at home.	Why the COVID-19 outbreak has lifted near-term revenue for Amazon?	Shoppers spend more time at home	Shoppers spend more time at home
(Relation: Enable) As a first mover in the local-market daily deals space, Groupon has captured a leadership position, but not robust profitability.	What enable Groupon capture a leadership position?	A first mover in the local-market daily deals space	A first mover in the local-market daily deals space
(Relation: Prevent_By) In neurology, RNA therapies can reach their intended targets via intrathecal administration into spinal fluid, directly preventing the production of toxic proteins	What will be prevented if intrathecal administration into spinal fluid?	The production of toxic proteins	The production of toxic proteins
Examples of Incorrect Predictions			
(Relation: Enable_By) ... Through analyzing the data and applying artificial intelligence, the advertisers can improve the efficiency of advertisements through targeted marketing for Tencent ...	What can help advertisers to improve the efficiency of advertisements?	Analyzing the data and applying artificial intelligence	Targeted marketing
(Relation: Cause_By) Given expectations for more volatile equity and credit markets, as well as some disruption as Brexit moves forward, it remain doubtful that flows will improve too dramatically, a negative 3%-5% annual organic growth...	Why a negative 3%-5% annual organic growth happened?	Given expectations ... as well as some disruption as Brexit moves forward	It remains doubtful that flows will improve too dramatically.

Table 6: Qualitative analysis of “Why” and “What-if” questions answering tasks based on the best-performed RoBERTa-Large model. The **company name** can be found in the meta-information of our dataset. **Cause** and **Effect** are extracted from the original context. The inputs of models consist with the context and question.

Dataset	Method	F1	ACC	EM
SQuAD1.1 (Rajpurkar et al., 2016)	LUKE (Yamada et al., 2020)	95.7	-	90.6
SQuAD2.0 (Rajpurkar et al., 2018)	IE-Net (Gao et al., 2019b)	93.2	-	90.9
DROP (Dua et al., 2019)	QDGAT (Chen et al., 2020)	88.4	-	-
HotpotQA (Yang et al., 2018)	BigBird-etc (Zaheer et al., 2020)	95.7	-	90.6
Reasoning Based Datasets				
LogiQA(Liu et al., 2020)	DAGAN (Huang et al., 2021)	-	39.3	-
CausalQA (Ours)	RoBERTa-SQuAD	84.7	85.6	61.6

Table 7: The comparison of best performance between our dataset and other popular QA datasets.

– are ignored as it is relatively away for the question position.

4.3.3 Challenges by CausalQA

We are interested in better understanding the difficulty of the CausalQA task compared to other popular datasets regarding prediction performance. We list the best-performing model of several popular datasets in Table 7. In general, we find that reasoning-based tasks are more complex than other tasks in terms of the relatively low accuracy achieved by the state-of-the-art method. LogiQA is more challenging than our dataset (39.3 vs. 85.6 in accuracy) because it requires heavy logical reasoning rather than identifying causal relations from text. Moreover, we find that the state-of-the-art result on our dataset (RoBERTa-SQuAD) is dramatically worse than the best performance on other datasets (EM = 90.9 on SQuAD2.0 while EM = 61.6 on CausalQA). This may suggest that the model tends to output the partially right answer but fails to output the utterly correct answer, al-

though further research is required, as the model still could be easily perturbed by the length of an event. Meanwhile, the human performance is still ahead of the best-performing model’s result in the causal reasoning QA task. Thus, we argue that CausalQA is worth investigating by using more “causal-thinking” methods in the future.

5 Related Work

This paper brings together two interesting ideas – event causality and causal question answering – and in what follows, we briefly introduce the existing relevant works and datasets to the present work.

Event Causality. There is a deep literature on causal inference techniques using non-text datasets (Pearl, 2009; Morgan and Winship, 2015; Keith et al., 2020; Feder et al., 2021), and a line of work focusing on discovering the causal relationship between events from textual data (Gordon et al., 2012; Mirza and Tonelli, 2016; Du et al., 2021). Previous efforts lie on the graph-based event causality detec-

Datasets	Event Extraction	Causal Reasoning	Fine-grained Causality	Span-based QA	Sources
FinCausal (Mariko et al., 2020b)	✓	✓	✗	✗	Finance
COPA (Roemmele et al., 2011)	✗	✓	✗	✗	Open
SQuAD (Rajpurkar et al., 2018)	✗	✗	✗	✓	Wikipedia
LogiQA (Liu et al., 2020)	✗	✓	✗	✗	Examination
HotpotQA (Yang et al., 2018)	✗	✗	✗	✓	Wikipedia
DROP (Dua et al., 2019)	✗	✗	✗	✓	Wikipedia
DREAM (Sun et al., 2019)	✗	✓	✗	✗	Examination
RACE (Lai et al., 2017)	✗	✓	✗	✗	Examination
FineCR (Ours)	✓	✓	✓	✓	Finance

Table 8: Comparisons of our fin-grained causal reasoning dataset and related public datasets.

tion tasks (Tanon et al., 2017; Li et al., 2020; Du et al., 2021) and the event-level causality detection tasks (Mariko et al., 2020a; El-Haj et al., 2021; Gusev and Tikhonov, 2021). However, causal reasoning for text data with a special focus on fine-grained causality between events has been relatively little considered. For this reason, we build a fine-grained causality dataset in the financial domain and expect to see whether the state-of-the-art models can achieve human-like accuracy on several causal reasoning tasks, and if not, to what extent.

Datasets. Table 8 compares our dataset with datasets in the domain of both event causality and question answering (QA). FinCausal (Mariko et al., 2020a) dataset is the most relevant to ours, which developed a relatively small dataset from the Edgar Database⁶ focusing on the simple “cause” relation only and do not contain QA tasks. In addition, existing popular question answering datasets (Sun et al., 2019; Liu et al., 2020; Cui et al., 2020) mainly focus on *what*, *who*, *where* and *when* questions, making their usage scenarios somewhat limited. SQuAD (Rajpurkar et al., 2016, 2018) consists of factual questions concerning Wikipedia articles, and some unanswerable questions are involved in SQuAD2.0. Although there are some datasets contain the causal reasoning tasks (Lai et al., 2017; Sun et al., 2019; Cui et al., 2020), none of them consider answering questions by text span. Span-based question answering problems have gained wide interest in recent years (Yang et al., 2018; Huang et al., 2019; Lewis et al., 2021b). HotpotQA (Yang et al., 2018) focuses on multi-hop QA where the question can only be answered through analyzing multiple documents. The answers in the (Dua et al., 2019) may come from different spans of a passage and require some combination technolo-

⁶<https://www.sec.gov/edgar/searchedgar/>

gies to get the correct answer. Compared with these datasets, that none of them have features of causal reasoning and span-based QA simultaneously, our dataset is the first to leverage fine-grained human-labeled causality for designing the CausalQA task consisting with “Why” and “What-if” questions.

Our task is similar to the machine reading comprehension setting (Huang et al., 2019) where the algorithms make a multiple-choice selection given a passage and a question. Nevertheless, we focus on causal questions, which turn out to be more challenging. To the best of our knowledge, we are the first to evaluate models on the event causality analysis and causal question answering (CausalQA) tasks based on the fine-grained causality dataset. We will release our code and dataset on Github.

6 Conclusion

We explored the efficacy of current state-of-the-art methods for causal reasoning tasks by considering a novel fine-grained reasoning setting and developing a dataset with rich human labels. Experimental results using the state-of-the-art neural language models provide the evidence that there is still much room for improvement on causal reasoning tasks, where there is a need for designing better solutions to correlation discovery related to event causality analysis and Why/What-if QA tasks.

7 Ethical Statement

This paper honors the ACL Code of Ethics. Public available financial analysis reports are used to extract fine-grained event relationships. No private data or non-public information was used. We obtain permission from Yahoo Finance for non-profit research. All annotators have received labor fees corresponding to their amount of annotated corpus.

References

- 529
- 530 Ting-Yun Chang and Yun-Nung Chen. 2019. **What does**
531 **this word mean? explaining contextualized embed-**
532 **dings with natural language definition.** In *Proceed-*
533 *ings of the 2019 Conference on Empirical Methods*
534 *in Natural Language Processing and the 9th Inter-*
535 *national Joint Conference on Natural Language Pro-*
536 *cessing (EMNLP-IJCNLP)*, pages 6064–6070, Hong
537 Kong, China. Association for Computational Linguistics.
538
- 539 Chung-Chi Chen, Hen-Hsen Huang, and Hsin-Hsi Chen.
540 2021. *From Opinion Mining to Financial Argu-*
541 *ment Mining.* Springer Briefs in Computer Science.
542 Springer.
- 543 Kunlong Chen, Weidi Xu, Xingyi Cheng, Zou Xi-
544 aochuan, Yuyu Zhang, Le Song, Taifeng Wang, Yuan
545 Qi, and Wei Chu. 2020. Question directed graph
546 attention network for numerical reasoning over text.
547 *arXiv preprint arXiv:2009.07448.*
- 548 Leyang Cui, Yu Wu, Shujie Liu, Yue Zhang, and Ming
549 Zhou. 2020. Mutual: A dataset for multi-turn dia-
550 logue reasoning. In *Proceedings of the 58th Confer-*
551 *ence of the Association for Computational Linguis-*
552 *tics.* Association for Computational Linguistics.
- 553 Dina Demner-Fushman, Kevin Bretonnel Cohen,
554 Sophia Ananiadou, and Junichi Tsujii, editors. 2021.
555 *Proceedings of the 20th Workshop on Biomedical*
556 *Language Processing.* Association for Computational
557 Linguistics, Online.
- 558 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
559 Kristina Toutanova. 2019. **BERT: Pre-training of**
560 **deep bidirectional transformers for language under-**
561 **standing.** In *Proceedings of the 2019 Conference of*
562 *the North American Chapter of the Association for*
563 *Computational Linguistics: Human Language Tech-*
564 *nologies, Volume 1 (Long and Short Papers)*, pages
565 4171–4186, Minneapolis, Minnesota. Association for
566 Computational Linguistics.
- 567 Xiao Ding, Yue Zhang, Ting Liu, and Junwen Duan.
568 2015. Deep learning for event-driven stock predic-
569 tion. In *Proceedings of the 24th International Con-*
570 *ference on Artificial Intelligence*, page 2327–2333,
571 Buenos Aires, Argentina.
- 572 Li Du, Xiao Ding, Kai Xiong, Ting Liu, and Bing Qin.
573 2021. **ExCAR: Event graph knowledge enhanced**
574 **explainable causal reasoning.** In *Proceedings of the*
575 *59th Annual Meeting of the Association for Comput-*
576 *ational Linguistics and the 11th International Joint*
577 *Conference on Natural Language Processing (Vol-*
578 *ume 1: Long Papers)*, pages 2354–2363, Online. As-
579 sociation for Computational Linguistics.
- 580 Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel
581 Stanovsky, Sameer Singh, and Matt Gardner. 2019.
582 Drop: A reading comprehension benchmark re-
583 quiring discrete reasoning over paragraphs. *arXiv*
584 *preprint arXiv:1903.00161.*
- Mahmoud El-Haj, Paul Rayson, and Nadhem Zmandar,
editors. 2021. *Proceedings of the 3rd Financial Nar-*
rative Processing Workshop. Association for Compu-
tational Linguistics, Lancaster, United Kingdom.
- Amir Feder, Katherine A Keith, Emaad Manzoor, Reid
Pryzant, Dhanya Sridhar, Zach Wood-Doughty, Jacob
Eisenstein, Justin Grimmer, Roi Reichart, Margaret E
Roberts, et al. 2021. Causal inference in natural lan-
guage processing: Estimation, prediction, interpreta-
tion and beyond. *arXiv preprint arXiv:2109.00725.*
- Lei Gao, Prafulla Kumar Choubey, and Ruihong Huang.
2019a. **Modeling document-level causal structures**
for event causal relation identification. In *Proceed-*
ings of the 2019 Conference of the North American
Chapter of the Association for Computational Lin-
guistics: Human Language Technologies, Volume 1
(Long and Short Papers), pages 1808–1817, Min-
neapolis, Minnesota. Association for Computational
Linguistics.
- Yuan Gao, Zixiang Cai, and Lei Yu. 2019b. Intra-
ensemble in neural networks. *arXiv preprint*
arXiv:1904.04466.
- Andrew Gordon, Zornitsa Kozareva, and Melissa Roem-
mele. 2012. Semeval-2012 task 7: Choice of plau-
sible alternatives: An evaluation of commonsense
causal reasoning. In *SEM 2012: The First Joint Con-*
ference on Lexical and Computational Semantics-
Volume 1: Proceedings of the main conference and
the shared task, and Volume 2: Proceedings of the
Sixth International Workshop on Semantic Evalua-
tion (SemEval 2012), pages 394–398.
- Suchin Gururangan, Swabha Swayamdipta, Omer Levy,
Roy Schwartz, Samuel Bowman, and Noah A Smith.
2018. Annotation artifacts in natural language infer-
ence data. In *Proceedings of the 2018 Conference of*
the North American Chapter of the Association for
Computational Linguistics: Human Language Tech-
nologies, Volume 2 (Short Papers), pages 107–112.
- Ilya Gusev and Alexey Tikhonov. 2021. Headlinecause:
A dataset of news headlines for detecting casualties.
ArXiv, abs/2108.12626.
- Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and
Yejin Choi. 2019. Cosmos qa: Machine reading com-
prehension with contextual commonsense reasoning.
arXiv preprint arXiv:1909.00277.
- Yinya Huang, Meng Fang, Yu Cao, Liwei Wang, and
Xiaodan Liang. 2021. Dagn: Discourse-aware
graph network for logical reasoning. *arXiv preprint*
arXiv:2103.14349.
- Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld,
Luke Zettlemoyer, and Omer Levy. 2020. Spanbert:
Improving pre-training by representing and predict-
ing spans. *Transactions of the Association for Com-*
putational Linguistics, 8:64–77.

639	Divyansh Kaushik, Eduard Hovy, and Zachary C Lipton.	Zhongyang Li, Xiao Ding, Kuo Liao, Ting Liu, and Bing	696
640	2020. Learning the difference that makes a difference	Qin. 2021. Causalbert: Injecting causal knowledge	697
641	with counterfactually-augmented data. In <i>Proceed-</i>	into pre-trained models with minimal supervision.	698
642	<i>ings of the International Conference on Learning</i>	<i>ArXiv</i> , abs/2107.09852.	699
643	<i>Representations (ICLR).</i>		
644	Divyansh Kaushik, Amrith Setlur, Eduard Hovy, and	Zhongyang Li, Xiao Ding, Ting Liu, J. Edward Hu,	700
645	Zachary C Lipton. 2021. Explaining the efficacy of	and Benjamin Van Durme. 2020. Guided generation	701
646	counterfactually augmented data . In <i>International</i>	of cause and effect . In <i>Proceedings of the Twenty-</i>	702
647	<i>Conference on Learning Representations</i> .	<i>Ninth International Joint Conference on Artificial</i>	703
648		<i>Intelligence, IJCAI-20</i> , pages 3629–3636. Interna-	704
649	Humayun Kayesh, Md. Saiful Islam, Junhu Wang,	tional Joint Conferences on Artificial Intelligence	705
650	Shikha Anirban, A.S.M. Kayes, and Paul Watters.	Organization. Main track.	706
651	2020. Answering binary causal questions: A transfer		
652	learning based approach . In <i>2020 International Joint</i>	Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang,	707
653	<i>Conference on Neural Networks (IJCNN)</i> , pages 1–9.	Yile Wang, and Yue Zhang. 2020. Logiqa: A	708
654		challenge dataset for machine reading compre-	709
655	Katherine Keith, David Jensen, and Brendan O’Connor.	hension with logical reasoning. <i>arXiv preprint</i>	710
656	2020. Text and causal inference: A review of using	<i>arXiv:2007.08124</i> .	711
657	text to remove confounding from causal estimates.		
658	In <i>Proceedings of the 58th Annual Meeting of the As-</i>	Linqing Liu, Patrick Lewis, Sebastian Riedel, and Pon-	712
659	<i>sociation for Computational Linguistics</i> , pages 5332–	tus Stenetorp. 2021. Challenges in generalization	713
660	5344.	in open domain question answering. <i>arXiv preprint</i>	714
661		<i>arXiv:2109.01156</i> .	715
662	Christopher SG Khoo, Jaklin Kornfilt, Robert N Oddy,		
663	and Sung Hyon Myaeng. 1998. Automatic extrac-	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	716
664	tion of cause-effect information from newspaper text	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	717
665	without knowledge-based inferencing. <i>Literary and</i>	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	718
666	<i>Linguistic Computing</i> , 13(4):177–186.	RoBERTa: A Robustly Optimized BERT Pretraining	719
667		Approach . <i>arXiv e-prints</i> , page arXiv:1907.11692.	720
668	Yoon Kim. 2014. Convolutional neural networks for		
669	sentence classification. In <i>Proceedings of the 2014</i>	Dominique Mariko, Hanna Abi-Akl, Estelle Labidurie,	721
670	<i>Conference on Empirical Methods in Natural Lan-</i>	Stephane Durfort, Hugues De Mazancourt, and Mah-	722
671	<i>guage Processing (EMNLP)</i> , pages 1746–1751.	moud El-Haj. 2020a. The financial document causal-	723
672		ity detection shared task (FinCausal 2020) . In <i>Pro-</i>	724
673	Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang,	<i>ceedings of the 1st Joint Workshop on Financial</i>	725
674	and Eduard Hovy. 2017. Race: Large-scale reading	<i>Narrative Processing and MultiLing Financial Sum-</i>	726
675	comprehension dataset from examinations. <i>arXiv</i>	<i>marisation</i> , pages 23–32, Barcelona, Spain (Online).	727
676	<i>preprint arXiv:1704.04683</i> .	COLING.	728
677			
678	Yuxuan Lai, Chen Zhang, Yansong Feng, Quzhe Huang,	Dominique Mariko, Estelle Labidurie, Yagmur Ozturk,	729
679	and Dongyan Zhao. 2021. Why machine read-	Hanna Abi Akl, and Hugues de Mazancourt. 2020b.	730
680	ing comprehension models learn shortcuts? <i>arXiv</i>	Data processing and annotation schemes for fincausal	731
681	<i>preprint arXiv:2106.01024</i> .	shared task. <i>arXiv preprint arXiv:2012.02498</i> .	732
682			
683	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan	Paramita Mirza, Rachele Sprugnoli, Sara Tonelli, and	733
684	Ghazvininejad, Abdelrahman Mohamed, Omer Levy,	Manuela Speranza. 2014. Annotating causality in	734
685	Veselin Stoyanov, and Luke Zettlemoyer. 2020.	the tempeval-3 corpus. In <i>EACL 2014 Workshop on</i>	735
686	BART: Denoising sequence-to-sequence pre-training	<i>Computational Approaches to Causality in Language</i>	736
687	for natural language generation, translation, and com-	<i>(CatoCL)</i> , pages 10–19. Association for Computa-	737
688	prehension . In <i>Proceedings of the 58th Annual Meet-</i>	tional Linguistics.	738
689	<i>ing of the Association for Computational Linguistics</i> ,		
690	pages 7871–7880, Online. Association for Computa-	Paramita Mirza and Sara Tonelli. 2016. Catena: Causal	739
691	tional Linguistics.	and temporal relation extraction from natural lan-	740
692		guage texts. In <i>Proceedings of COLING 2016, the</i>	741
693	Patrick Lewis, Pontus Stenetorp, and Sebastian Riedel.	<i>26th International Conference on Computational Lin-</i>	742
694	2021a. Question and answer test-train overlap in	<i>guistics: Technical Papers</i> , pages 64–75.	743
695	open-domain question answering datasets. In <i>Pro-</i>		
696	<i>ceedings of the 16th Conference of the European</i>	Stephen L Morgan and Christopher Winship. 2015.	744
697	<i>Chapter of the Association for Computational Lin-</i>	<i>Counterfactuals and causal inference</i> . Cambridge	745
698	<i>guistics: Main Volume</i> , pages 1000–1008.	University Press.	746
699			
700	Patrick Lewis, Yuxiang Wu, Linqing Liu, Pasquale Min-	Jong-Hoon Oh, Kentaro Torisawa, Chikara Hashimoto,	747
701	ervini, Heinrich Küttler, Aleksandra Piktus, Pontus	Ryu Iida, Masahiro Tanaka, and Julien Kloetzer.	748
702	Stenetorp, and Sebastian Riedel. 2021b. Paq: 65 mil-	2016. A Semi-Supervised Learning Approach to	749
703	lion probably-asked questions and what you can do	Why-Question Answering . <i>Proceedings of the AAAI</i>	750
704	with them . <i>arXiv</i> .	<i>Conference on Artificial Intelligence</i> , 30(1).	751
705			

752	Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. 2019. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. <i>Advances in Neural Information Processing Systems</i> , 32:13991–14002.	Kai Sun, Dian Yu, Jianshu Chen, Dong Yu, Yejin Choi, and Claire Cardie. 2019. Dream: A challenge data set and models for dialogue-based reading comprehension. <i>Transactions of the Association for Computational Linguistics</i> , 7:217–231.	809
753			810
754			811
755			812
756			813
757			
758			
759	Judea Pearl. 2009. <i>Causality</i> . Cambridge university press.	Leonard Talmy. 1988. Force dynamics in language and cognition. <i>Cognitive science</i> , 12(1):49–100.	814
760			815
761	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of Machine Learning Research</i> , 21(140):1–67.	Tsun-Hsien Tang, Chung-Chi Chen, Hen-Hsen Huang, and Hsin-Hsi Chen. 2021. Retrieving implicit information for stock movement prediction. In <i>Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR ’21</i> , page 2010–2014, New York, NY, USA. Association for Computing Machinery.	816
762			817
763			818
764			819
765			820
766			821
767	Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don’t know: Unanswerable questions for SQuAD. In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 784–789, Melbourne, Australia. Association for Computational Linguistics.	Thomas Pellissier Tanon, Daria Stepanova, Simon Razniewski, Paramita Mirza, and Gerhard Weikum. 2017. Completeness-aware rule learning from knowledge graphs. In <i>International Semantic Web Conference</i> , pages 507–525. Springer.	822
768			823
769			824
770			825
771			826
772			827
773			
774	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017a. Attention is all you need. In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.	828
775			829
776			830
777			831
778			832
779			833
780	Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. 2011. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In <i>2011 AAAI Spring Symposium Series</i> .	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017b. Attention is all you need. In <i>NIPS</i> .	834
781			835
782			836
783			837
784	Nino Scherrer, Olexa Bilaniuk, Yashas Annadani, Anirudh Goyal, Patrick Schwab, Bernhard Schölkopf, Michael C Mozer, Yoshua Bengio, Stefan Bauer, and Nan Rosemary Ke. 2021. Learning neural causal models with active interventions. <i>arXiv preprint arXiv:2109.02429</i> .	Cunxiang Wang, Pai Liu, and Yue Zhang. 2021. Can generative pre-trained language models serve as knowledge bases for closed-book QA? In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3241–3251, Online. Association for Computational Linguistics.	838
785			839
786			840
787			841
788			842
789			843
790	Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. 2021. Toward causal representation learning. <i>Proceedings of the IEEE</i> , 109(5):612–634.	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	844
791			845
792			846
793			847
794			848
795	Megha Srivastava, Tatsunori Hashimoto, and Percy Liang. 2020. Robustness to spurious correlations via human annotations. In <i>International Conference on Machine Learning</i> , pages 9109–9119. PMLR.		849
796			850
797			851
798			852
799	Saku Sugawara, Kentaro Inui, Satoshi Sekine, and Akiko Aizawa. 2018. What makes reading comprehension questions easier? In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 4208–4219.		853
800			854
801			855
802			856
803			857
804	Saku Sugawara, Pontus Stenetorp, Kentaro Inui, and Akiko Aizawa. 2020. Assessing the benchmarking capacity of machine reading comprehension datasets. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 34, pages 8918–8927.	Phillip Wolff, Bianca Klettke, Tatyana Ventura, and Grace Song. 2005. Expressing causation in english and other languages. <i>Journal of experimental psychology: General</i> .	858
805			859
806			860
807			861
808		Phillip Wolff and Grace Song. 2003. Models of causation and the semantics of causal verbs. <i>Cognitive psychology</i> , 47(3):276–332.	862

865 Ikuya Yamada, Akari Asai, Hiroyuki Shindo, Hideaki
866 Takeda, and Yuji Matsumoto. 2020. Luke: deep con-
867 textualized entity representations with entity-aware
868 self-attention. *arXiv preprint arXiv:2010.01057*.

869 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Ben-
870 gio, William W Cohen, Ruslan Salakhutdinov, and
871 Christopher D Manning. 2018. Hotpotqa: A dataset
872 for diverse, explainable multi-hop question answer-
873 ing. *arXiv preprint arXiv:1809.09600*.

874 Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He,
875 Alex Smola, and Eduard Hovy. 2016. Hierarchical at-
876 tention networks for document classification. In *Pro-
877 ceedings of the 2016 conference of the North Ameri-
878 can chapter of the association for computational lin-
879 guistics: human language technologies*, pages 1480–
880 1489.

881 Manzil Zaheer, Guru Guruganesh, Kumar Avinava
882 Dubey, Joshua Ainslie, Chris Alberti, Santiago On-
883 tanon, Philip Pham, Anirudh Ravula, Qifan Wang,
884 Li Yang, et al. 2020. Big bird: Transformers for
885 longer sequences. In *NeurIPS*.

A Appendix: Out-of-domain Test

Sec	Consumer		Industrial		Technology	
	F1	Acc	F1	Acc	F1	Acc
Con	59.69	69.01	50.03	66.47	47.95	63.74
Ind	50.60	67.39	51.14	64.61	47.49	64.09
Tec	48.65	65.87	47.70	63.56	50.39	61.12

Table 9: Out-of-domain test results of the BERT-base model for the fine-grained causality classification task.

Sec	Consumer		Industrial		Technology	
	F1	EM	F1	EM	F1	EM
Con	86.26	49.54	85.75	48.37	85.20	47.26
Ind	84.23	49.83	86.00	50.68	84.90	47.45
Tec	86.24	47.99	86.03	47.50	87.09	61.12

Table 10: Out-of-domain test results of the Span-Large model for the cause-effect extraction task.

It has been shown that sector-relevant features from a given domain could become spurious patterns on the other domains, leading to performance decay under distribution shift (Ovadia et al., 2019). We use instances from three sectors with the largest amounts of samples in our dataset for conducting out-of-domain generalization text. These observe in line with recent works revealing that current deep neural models mostly memorize training instances yet struggle to predict on the out-of-distribution data (Gururangan et al., 2018; Kaushik et al., 2020; Srivastava et al., 2020). To evaluate whether methods can generalize on the out-of-distribution data, and to what extent, the results of the out-of-domain test are shown in Table 9 and Table 10.

In particular, the model achieves the best performance when the training and test sets are extracted from the articles of the same domain companies. In the out-of-domain test, the model shows varying degrees of performance decay for both tasks. For example, in the fine-grained causality classification task, the model trained with the data from the Consumer Cyclical domain achieves 59.69 F1 Score when testing on the Consumer Cyclical data while decreasing to 47.95 when testing on technology companies. Moreover, in the cause-effect extraction task, the model trained with the data from the Consumer Cyclical domain achieves 86.26 F1 Score when testing on itself while decreasing to 85.20 when testing on Technology. This shows that the domain-relevant patterns learned by the model cannot transfer well between domains.

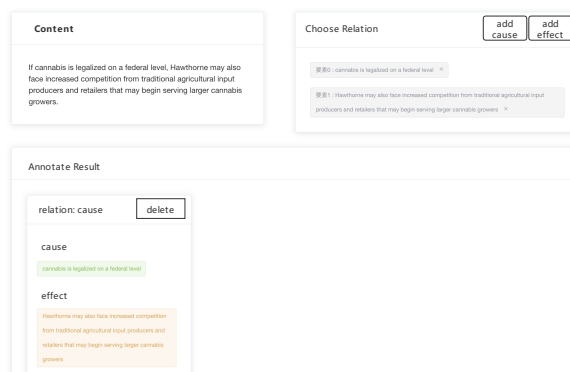


Figure 5: The annotation platform provided the crowdsourcing company for collecting annotations for fine-grained causality reasoning and CausalQA.

B Appendix: Annotation Instructions

The annotation platform used in this work is introduced in Fig. 5. As follows, we provide the detailed annotation instructions used for training the human labelers. Also, we show the annotation of some real examples stored in our dataset.

B.1 General Instruction

This is a annotation task related to event causality. In this task, you are asked to find all the cause-effect pairs and the fine-grained event relationship types from give passages.

B.2 Steps

1. Please read your assigned examples carefully.
2. A sentence is considered contains event causality if at least two events occur in it and the two events are causally related.
3. If a sentence contains causality, mark it as **Positive**; otherwise, mark it as **Negative**.
4. For the **positive** sentence, first find all the events that occur in the sentence, and then pair the events to see if they constitute a causal relationship. The relationship much be one of **Cause**, **Enable** and **Prevent**.
5. “A causes B” means B always happens if A happens. “A enables B” means A is a possible way for B to happen, but not necessarily. “A prevents B” means A and B cannot happen at the same time.
6. Remember to annotate all event causality pairs. If there is no more pairs, process to the next passage.

950 B.3 Examples

951 Here are some annotation examples, please read it
952 before starting your annotation.

953 **Example 1:** Moreover, we do not think that
954 DBK's investment banking operation has the neces-
955 sary scale and set-up to outcompete peers globally
956 or within Europe.

957 **Answer: # Negative**

958 **Explanation:** This is a sentence that contains
959 no causal relationship between events.

960 **Example 2:** In our view, **customers are likely**
961 **to stay with VMware** because of **knowledge of**
962 **its product ecosystem as well as the risks and**
963 **complexities associated with changing virtual**
964 **machine providers.**

965 **Answer: # Positive**

966 **Explanation:** This is a causal sentence. There
967 exist two events marked by **yellow** color. You
968 should first annotate the two events and then give
969 them the label according to their relationship, us-
970 ing one of **Cause**, **Enable** and **Prevent**. Here the
971 relationship is **Cause**.

972 **Example 3:** **Depressed realized prices** due to
973 **lack of market access** have forced **capital spend-**
974 **ing cuts**, stalling the **growth potential of the com-**
975 **pany's oil sands assets.**

976 **Answer: # Positive**

977 **Explanation:** This is a causal sentence and there
978 exist four events. You need to mark out all four of
979 these events and then pair them up to see if they're
980 related. If so, determine what kind of relationship
981 they belong to.