

REFinD: Relation Extraction Financial Dataset

Simerjot Kaur*

simerjot.kaur@jpmchase.com
JPMorgan Chase and Co
Palo Alto, CA, USA

Charese Smiley*

charese.h.smiley@jpmchase.com
JPMorgan Chase and Co
Chicago, IL, USA

Akshat Gupta

akshat.x.gupta@jpmchase.com
JPMorgan Chase and Co
New York, NY, USA

Joy Sain

sain.9@wright.edu
Wright State University
Dayton, OH, USA

Dongsheng Wang

dongsheng.wang@jpmchase.com
JPMorgan Chase and Co
London, UK

Suchetha Siddagangappa

suchetha.siddagangappa@jpmchase.com
JPMorgan Chase and Co
New York, NY, USA

Toyin Aguda

toyin.d.aguda@jpmchase.com
JPMorgan Chase and Co
Chicago, IL, USA

Sameena Shah

sameena.shah@jpmchase.com
JPMorgan Chase and Co
New York, NY, USA

ABSTRACT

A number of datasets for Relation Extraction (RE) have been created to aid downstream tasks such as information retrieval, semantic search, question answering and textual entailment. However, these datasets fail to capture financial-domain specific challenges since most of these datasets are compiled using general knowledge sources, hindering real-life progress and adoption within the financial world. To address this limitation, we propose REFinD, the first large-scale annotated dataset of relations, with ~29K instances and 22 relations amongst 8 types of entity pairs, generated entirely over financial documents. We also provide an empirical evaluation with various state-of-the-art models as benchmarks for the RE task and highlight the challenges posed by our dataset. We observed that various state-of-the-art deep learning models struggle with numeric inference, relational and directional ambiguity.

CCS CONCEPTS

• **Information systems** → **Information extraction; Test collections**; • **Applied computing** → *Document searching*; Annotation.

KEYWORDS

relation extraction, finance, natural language processing, benchmarking, information retrieval, annotation datasets

ACM Reference Format:

Simerjot Kaur*, Charese Smiley*, Akshat Gupta, Joy Sain, Dongsheng Wang, Suchetha Siddagangappa, Toyin Aguda, and Sameena Shah. 2023. REFinD: Relation Extraction Financial Dataset. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '23)*, July 23–27, 2023, Taipei, Taiwan. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3539618.3591911>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
SIGIR '23, July 23–27, 2023, Taipei, Taiwan

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-9408-6/23/07...\$15.00
<https://doi.org/10.1145/3539618.3591911>

1 INTRODUCTION

The exponential progress of AI across multiple domains can largely be attributed to the availability of large datasets coupled with an increase in available compute power. Relation extraction (RE) from text is a fundamental problem in NLP and information retrieval, which facilitates various tasks like knowledge graph construction, question answering and semantic search. It has seen significant progress in recent years, thanks to advanced machine learning techniques and the availability of large-scale relation extraction datasets. However, most existing large-scale RE datasets are derived from general knowledge sources such as Wikipedia, web texts and news articles [7, 9, 16, 22, 23]. These datasets often fall short in addressing domain-specific challenges. Hence, various state-of-the-art models that perform competitively on such datasets fail to perform well in the financial domain (shown in Section 5). In particular, financial text documents, such as financial reports and various Securities and Exchange Commission (SEC) filings, differ significantly from standard English language documents. They necessitate the extraction of entities and relations that involve numbers, currencies, dates, legal facts, and claims, often embedded in longer and more complex sentences with substantial distances between entities. Figure 1 illustrates a prototypical sentence from a financial report, emphasizing unique relations like *acquired by*, *revenue of*.

Moreover, financial documents often necessitate more advanced numerical inference to identify relationships amongst entities. For instance, a company is deemed to have been acquired by another entity if the latter owns more than 50% of its shares. Additionally, there can be ambiguity amongst relations in financial text, such as when a person serves only on a company's board is just a member of company and not considered as an employee. These financial domain-specific challenges make relation extraction from such documents more difficult. However, the current absence of large-scale finance-specific relation extraction dataset impedes progress, benchmarking, and real-life adoption of various relation extraction algorithms within the financial industry.

*These authors contributed equally to this work

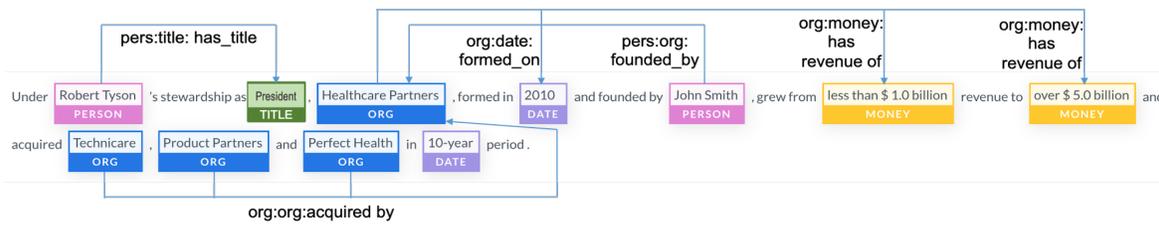


Figure 1: Example sentence from a financial report with highlighted entities and relations. In a single sentence, there are 8 relations present: *has title*, *formed in*, *founded by*, 2 instances of *has revenue of*, and 3 instances of *acquired by*.

To address this limitation, we have developed the largest relation extraction dataset for financial documents to date, REFinD¹, which contains ~29K instances and 22 relations among 8 types of entity pairs. REFinD is a domain-specific financial relation extraction dataset created using raw text from various 10-X reports (including 10-K, 10-Q, etc. broadly known as 10-X) of publicly traded companies obtained from US Securities and Exchange Commission (SEC)² website (detailed in Section 3). Although primarily built on financial reports and focused on financial domain-specific challenges, this dataset can also be leveraged by other domains such as legal, risk modeling, and econometrics.

In this work, we also highlight various challenges associated with creating a large-scale relation extraction dataset specifically over financial-domain. Since financial documents contain much longer and complex sentences and inferring relations from them involves a tremendous amount of financial-domain expertise, hence collecting seed labels, removing noisy text, and finally annotating such dataset becomes an extremely challenging task. Finally, we also provide benchmarks on REFinD dataset to identify and highlight challenges it poses, as well as to spur further research and improvements in the field of financial relation extraction. We observed that despite fine-tuning various state-of-the-art deep learning models on the REFinD dataset, their performance remains sub-optimal on finance-specific relations. Even specialized models like FinBERT and FLANG, which have been further trained on financial news articles and incorporate masking for financial-terms, do not demonstrate significant improvement. This is likely due to their lack of exposure to semantically complex financial documents, which hinders their ability to effectively handle intricate finance-specific scenarios.

This resource paper presents the following key contributions:

- Introducing REFinD, the first large-scale Relation Extraction Dataset over financial documents.
- Establishing benchmarks for various state-of-the-art models using the REFinD dataset.
- Identifying and highlighting the unique financial-domain specific challenges posed by the REFinD dataset.

2 RELATED WORK

Several datasets have been developed for RE using general knowledge sources such as Wikipedia³ and web articles including ACE

2003-2004 [18], SemEval2010 Task 8 [7], KBP37 [22], TACRED [23], FewRel [6], CrossRE [3], and MAVEN-ERE [19]. However, financial texts pose their own unique set of challenges and there has been limited attention paid towards creating RE datasets within the financial domain. Recently, a few datasets have been developed using financial news and earnings calls, including FinRED [16], CorpusFR [9] and Financial News Corpus [20]. Table 1 compares the number of relations and instances for recent general-purpose datasets (top) and financial datasets (bottom).

Dataset	Rels	Instances
SemEval-2010 Task 8	19	10,717
TACRED	42	119,474
KBP37	37	21,046
FinRED	29	6,767
Financial News Corpus	6	22,812
CorpusFR	20	1,754
REFinD	22	28,676

Table 1: Comparison between REFinD and prior RE datasets. REFinD is a large dataset, second only to TACRED, although it targets fewer relations.

Table 1 shows that TACRED, which annotates 42 relations, is the largest of these datasets. However, it should be noted that 79.5% of the TACRED dataset comprises NO_RELATION instances, whereas only 45.5% of REFinD is NO_RELATION. Consequently, while REFinD has a smaller total number of relations than TACRED, it has a higher

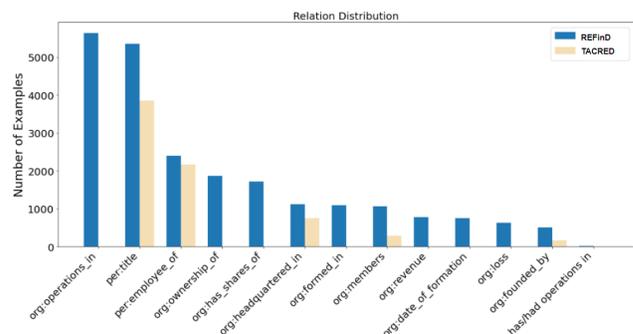


Figure 2: Relationship Distribution in REFinD and TACRED datasets

¹<https://www.jpmorgan.com/technology/artificial-intelligence/initiatives/refind-dataset/problem-motivation-outcome>

²<https://www.sec.gov/edgar.shtml>

³<https://www.wikipedia.org>

number of instances for each relation it covers. Moreover, REFinD has greater coverage than TACRED for the relations of interest in finance, as shown in Figure 2. Among financial datasets, both FinRED and CorpusFR are small compared to REFinD, and Financial News Corpus covers fewer relation types. Finally, REFinD is the first RE dataset to utilize SEC filings, which are a rich and complex data source.

3 DATASET OVERVIEW AND CONSTRUCTION

REFinD is specifically designed for use within the financial domain, using 10-X filings obtained from the SEC. The dataset targets 8 finance-specific entity pairs: PERSON-TITLE, PERSON-ORG, PERSON-UNIV, PERSON-GOV_AGE-NCY, ORG-GPE, ORG-ORG, and ORG-MONEY. Each entity pair group includes several possible finance-oriented relation types and Figure 3 illustrates these entity pair groups along with the 22 proposed relation types in REFinD dataset.

Entity (subj)	Relation	Entity (obj)
PERSON	title	TITLE
PERSON	employee_of member_of founder_of	ORGANIZATION
PERSON	employee_of member_of attended	UNIVERSITY
PERSON	member_of	GOV. AGENCY
ORGANIZATION	formed_on acquired_on	DATE
ORGANIZATION	headquartered_in operations_in formed_in	GPE
ORGANIZATION	shares_of subsidiary_of acquired_by agreement_with	ORGANIZATION
ORGANIZATION	revenue_of profit_of loss_of cost_of	MONEY
	no/other_relation	

Figure 3: REFinD Dataset. The dataset has 8 Entity pair groups and 22 Relation types.

Our annotations are at the instance-level, and each instance corresponds to the directed entity pair in a sentence that is annotated with one of the 22 relation labels. The dataset statistics and its distribution are provided in Table 5 (Appendix A.5). Additionally, we have included a snippet of the REFinD dataset in Appendix A.4. In total, the REFinD dataset contains ~29K instances with an average sentence length of 53 words and average contextual complexity of 11 words between the entity pairs.

The following sections provide detailed information on the collection process of the REFinD dataset, as well as the preprocessing steps that were taken to obtain high-quality data. Additionally, the annotation process used to label the instances with the 22 proposed relation types is explained in detail.

3.1 Document Collection

While most RE datasets have been constructed using sources such as Wikipedia, web articles, and financial news articles [7, 9, 16, 22, 23], the REFinD dataset is designed specifically for use within the financial domain and has been constructed from 10-X filings downloaded from the SEC website for the years 2016-2017. These publicly available financial regulatory reports are submitted to the SEC at regular intervals (monthly/quarterly/annually) by publicly traded companies and issuers of securities. They often provide detailed information about ownership, executive compensation, corporate structure, shares, trades, and other financial details.

3.2 Preprocessing

To prepare the 10-X reports for annotation, extensive data cleaning was required to address noise such as HTML tags and redundant spaces. Additionally, specific financial terms used to refer to legal entities presented a challenge, requiring entity tagging at the sentence-level to build the knowledge base.

Data Cleaning: To prepare the corpus for annotation, we first performed extensive data cleaning to remove header and footer information, tables, HTML tags, and redundant spaces. Additionally, financial text has unique characteristics, including the use of specific financial terms to refer to legal entities. To resolve these terms, we replaced pronouns and referring terms such as ‘we’, ‘our’, and ‘the company’ with the corresponding organization name. For example, the sentence such as *The Company had net revenues of \$14,720,545*, with organization name as *Technicare, LLC*, is preprocessed as *Technicare, LLC had net revenues of \$14,720,545*.

Named Entity Tagging: In order to build the knowledge base, we performed sentence tokenization and part-of-speech (POS) tagging using the spaCy library [8]. Named entity recognition (NER) was also carried out using spaCy for five entity types: PER, ORG, DATE, GPE, and MONEY. To capture the extensive use of job titles within financial documents, we employed Stanford CoreNLP [12] to tag entities as TITLE. Additionally, we introduced two new entity types: UNIV for educational establishments such as schools, colleges, universities, and institutes, and GOV for U.S. government agencies, using Gazetteer⁴ lists and regular expressions, respectively.

3.3 Dataset Construction

In order to construct this dataset, ~26K filings per year (2016-2017) were downloaded and after extracting and preprocessing we are left with millions of instances. Manual annotation of each instance is prohibitively costly and due to the sparsity of relations in the dataset (e.g. the location of a company’s headquarters may only be mentioned once per document), using a random sample of instances for annotation would lead to an extremely imbalanced dataset. Moreover, distantly supervised techniques such as Mintz et al. [13] work well for generating datasets, but rely on the use of an external knowledge base (KB) whereas financial text contains relatively unknown person/legal entities and finance-oriented relations which cannot be tracked by a general purpose KB.

Hence, to achieve purposeful annotation, we utilized a context-sensitive approach based on the construction of a set of phrases

⁴<https://census.gov>

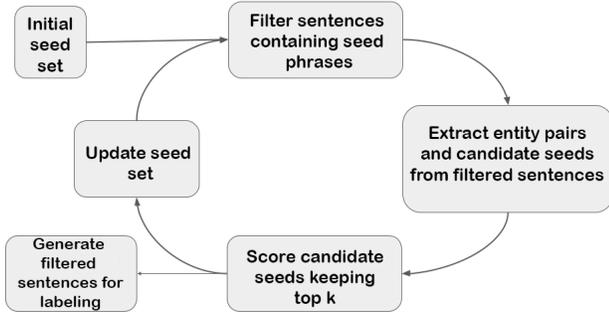


Figure 4: Corpus filtering using initial seed set.

[1, 24]. Specifically, we focused on pairs of entities e_1 and e_2 and a set of seed phrases p that were relevant to each relation type. For instance, we used phrases such as ‘net revenue’ for ORG-MONEY and ‘served as’ for PER-ORG. We retained instances that contained these phrases and collected all patterns (e_1, e_2, p) within the target relation group. We then used the HIT score, as in Zuo et al. [24], to rank these patterns and retained the top k patterns. We then calculate the convergence score $Conv$ to filter out candidates that co-occur with few entity pairs and retain those with a threshold greater than τ (see Appendix Eq. (3)).

Finally, to introduce diversity in our dataset, we expanded upon the approach used in [24] by collecting phrases from the shortest dependency path (SDP) between e_1 and e_2 . This captures longer contexts seen in financial texts and includes all the context surrounding the entities. We then created a new set of candidate seeds and retained the top k based on their HIT and $Conv$ scores. The selected seeds were then used in the next round of filtering, as shown in Figure 4. The process was repeated until no new seeds were added to the list or a maximum number of cycles was reached. The expanded seed list was then used to filter the instances and obtain the final set for annotation. This final set contained a more concentrated number of target relations, introduced more diverse expressions, and included near-misses that should be categorized as *no_relation*.

3.4 Dataset Annotation

To annotate our constructed dataset, we leverage Amazon Mechanical Turk (MTurk)⁵, a platform that allows us to crowdsource human intelligence tasks. Our dataset annotation task involves presenting a Human Intelligence Task (HIT) to a crowdworker. In each HIT, two entities, e_1 and e_2 , are highlighted, and the worker is asked to select the relation r that best describes the ordered entity pair (e_1, e_2) from a list of relations presented to them (see Figure 3). The list includes all possible entity pair relations, as well as an option for *no_relation* / *different_relation*. To ensure the quality of the annotations, we provide clear instructions and examples of how to complete the task, an example of how the annotation task is presented to the crowdworker can be found in Appendix A.1.

⁵<https://www.mturk.com>

3.5 Annotation Guidelines

To ensure high-quality annotations, we developed guidelines for the annotation task through an iterative process. This involved two preliminary rounds on a subset of data, followed by two official rounds for the 2016 data, one official round for the 2017 data, and one round for both years, all of which were reviewed by financial experts. During the preliminary rounds, we observed that workers experienced ambiguity when it came to the temporal aspect of relations in some instances. To address this, we provided a list of relations for each instance, along with the present and past tense verbs. For example, to show the relation between a PERSON and an ORG, we listed the relation as e_1 *is/was an employee of* e_2 , where entity markers e_1 and e_2 are replaced with the relevant entities.

Entity Group	Original Relation	Modified Relation
PER:ORG	is/was an employee of	is/was an employee of (e.g. CEO, President, Vice President, Manager)
PER:ORG	is/was a member of	is/was a (board, committee) member of
PER:TITLE	has/had job title	has/had job title (e.g. CEO, manager, director)
ORG:GPE	has/had operations in	has/had operations in (headquartered elsewhere)
ORG:ORG	has/had shares of	has/had % or number of shares of

Table 2: Relation labels modified during annotation rounds to enhance accuracy of annotations.

We also found that some workers would select the *founder_of* relation for any high-ranking position, such as the CEO or President, in cases where the roles of *founder_of* and *employee_of* may not necessarily overlap. To address this, we specified in our guidelines that workers should choose the relation that is most clearly stated in the displayed instance. Additionally, to increase clarity, we included examples in the relation text shown to workers. We modified the list of relations, and a complete list is provided in Table 2. Overall, our guidelines aim to reduce ambiguity and ensure consistent and accurate annotations.

For the 2016 data, we conducted two rounds of annotations during the official rounds. In the first round, we collected two judgements for all instances. For the second round, we collected one more judgement, but only for the instances where there was no consensus in the first round. However, we realized that this approach was time-consuming and required multiple rounds of labeling. To improve the efficiency of the annotated data collection process, for the 2017 data, we increased the number of judgements collected in the first round. We collected the number of judgements equal to the number of options within an entity pair group. This approach was more expensive, but it helped us save time by reducing the number of rounds needed for labeling. After collecting all the annotations from MTurk, we conducted one final official round. This round involved adjudication by experts in the finance domain to ensure that all the judgements were accurate and reliable.

4 DATASET ANALYSIS AND QUALITY

4.1 Contextual Complexity

To evaluate the contextual complexity of financial-domain instances, we compared the sentence length distributions for REFinD and the largest known Wikipedia relations dataset, TACRED. As illustrated in Figure 5, our analysis reveals that REFinD contains much longer sentences than TACRED wherein the average sentence length in TACRED dataset is 36.2 while REFinD is 53.7.

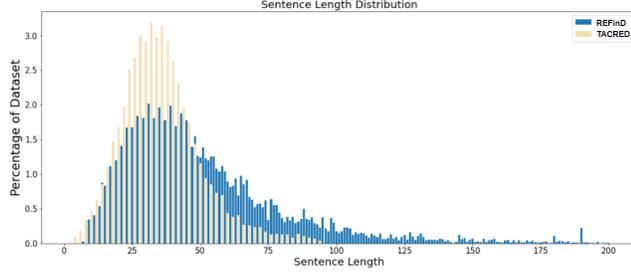


Figure 5: A detailed comparison of the sentence length distributions between REFinD and TACRED datasets.

We also compared the distance between entity pairs in REFinD and TACRED. As illustrated in Figure 6, our analysis indicates that REFinD includes more complex sentences than TACRED, with an average entity-pair distance of 11, compared to 8 in TACRED. This finding suggests that REFinD presents a greater level of difficulty in terms of identifying and linking entities within sentences, further emphasizing the dataset’s contextual complexity in the financial domain.

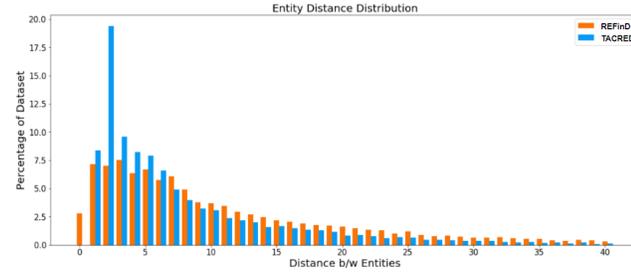


Figure 6: A detailed comparison of the contextual complexity between REFinD and TACRED datasets.

4.2 Label Aggregation

Our data collection efforts resulted in a total of 28,676 assessments from 1,209 assessors. To ensure the accuracy and reliability of the annotations, we aggregated the assessments using an internally calculated trust score that takes into account each assessor’s performance.

To calculate the trust score score for each assessor as follows: for each assessor w , we build a $|L| \times |L|$ confusion matrix F such that the (i, j) th entry is the count of times the assessor assigned label i to a label j . Thus, the accuracy is simply the sum of the

diagonal entries divided by the total number of assessments made by the assessor. However, this approach leads to a potential bias against the assessor, because an assessor could have used a label close to, but not identical to the ground truth label, e.g., *founder_of* versus *employee_of*. Therefore, to address the issue, we multiply the performance confusion matrix by a label similarity matrix Sim of size $|L| \times |L|$. This yields an adjusted confusion matrix $F^* = F \times Sim$, which takes into account the similarity between labels. While the label similarity matrix can theoretically take any form, we typically collaborate with domain experts who have reviewed the instances in the dataset to generate it. Then, we define the reliability of assessor $trust_w$ as the sum of all values of the diagonal entries of F^* .

Figure 7 shows the distribution of trust levels among assessors who annotated more than three instances. Out of a total of 1,209 assessors, 1,116 met this criterion. The range of trust levels is from 0 to 1, indicating variability in assessor performance. Some assessors may be spammers or lack sufficient domain knowledge to contribute effectively to the task. It’s also worth noting that assessors with lower trust levels may have assessed fewer instances.

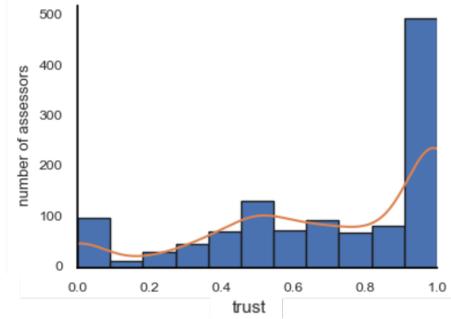


Figure 7: The trust distribution of the assessors, x-axis is the trust score $trust_w$ between 0 and 1.

Thus, the vote for an instance with label l is formulated with the consideration of reliability as,

$$vote(l) = \sum_{w \in W} trust_w(l) \quad (1)$$

Finally, for a given instance and pair of two entities, we assign the label with the maximum vote score,

$$(s, e_1, e_2) \mapsto \bar{l} \text{ where } \bar{l} = \operatorname{argmax}_{l \in L} vote(l) \quad (2)$$

To assess the impact of financial experts on the quality of the annotations, we plotted $vote(l)$ (Eq. 1) before and after their involvement. Figure 8 (i) and Figure 8 (ii) show the results before and after expert checking, respectively. We observed that, initially, 11.01% of the instances had a confidence value of less than 0.5. However, after expert checking, this percentage decreased to 3.49%. As a result, the average confidence score across all instances improved from 0.33 to 0.46. These findings suggest that financial experts significantly improved the quality of the annotations.

4.3 Annotation Agreement

To determine the inter-annotator agreement among the MTurk assessors, we utilized Fleiss Kappa Score [5]. In order to investigate

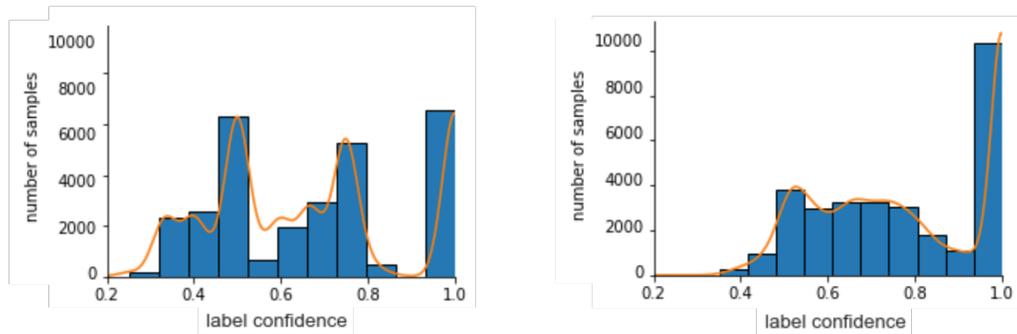


Figure 8: The normalized aggregated confidence for all samples (i) before expert checking and (ii) after expert checking, where x-axis is the confidence score between 0 and 1.

the impact of worker reliability on our results, we conducted additional analyses by varying the trust score threshold t from 0 to 1. Specifically, we excluded workers with trust scores below t and recalculated the kappa scores. As depicted in Figure 9, the Kappa score for each entity pair group varies with respect to the worker trust score threshold. Notably, we observe that removing assessments from workers with low trust scores leads to an increase in the Kappa score, suggesting that less reliable workers had a negative impact on the overall agreement among assessors. Hence, excluding less reliable workers improved the overall quality of our annotations and increased inter-annotator agreement.

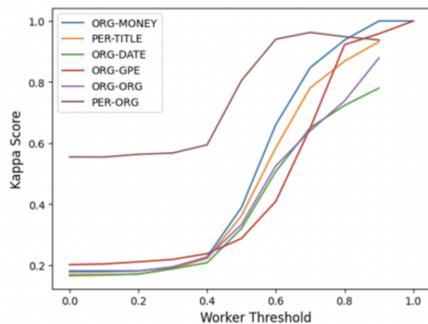


Figure 9: Fleiss Kappa Score. Inter-annotator agreement for each entity pair groups.

It is worth noting that, for this analysis, we have merged the entity groups PER-ORG, PER-UNIV, and PER-GOV into a single entity group PER-ORG, given that UNIV and GOV can be regarded as specific types of organizations.

4.4 Noise Rate

We conducted a noise rate analysis on our dataset by randomly selecting a 5% sample and verifying the annotations with the help of a linguist and a financial expert. The analysis revealed an overall noise rate of 6%. We found that the primary source of noise was the ORG-ORG and ORG-MONEY entity groups, which contain finance-specific relations such as *acquired_by*, *shares_of*, *revenue_of*, and *cost_of*. Annotating these groups is challenging and requires domain

expertise to understand and correctly annotate these relations. We leveraged this information to enhance the quality of our dataset by correcting the errors with the help of financial experts. We have included the noise rates for each entity group in Table 3.

Entity Pair	Noise Rate (%)
ORG:ORG	14
ORG:GPE	4
ORG:MONEY	13
ORG:DATE	1
PER:ORG	5
PER:UNIV	8
PER:GOV_AGY	8
PER:TITLE	3

Table 3: Noise Rate

5 BENCHMARK EXPERIMENTS AND RESULTS

In this section, we fine-tune various state-of-the-art deep learning models on the REFinD dataset for the relation extraction task, in order to assess and highlight the challenges posed by REFinD. To ensure a comprehensive evaluation of the benchmarks, we report both micro- and macro-F1 score metrics (evaluation details in Appendix A.3). We also evaluate the performance of each model on each entity pair group and the entire REFinD dataset to obtain a more detailed understanding of the models' strengths and weaknesses.

Parameter	Value
Classifier	1-layer FFNN
Loss	Cross Entropy
Optimizer	Adam optimizer
Learning rate	2e-5
Batch Size	32
Epochs	5

Table 4: Hyperparameters Setting. Model details for reproducibility of the baselines.

MODELS	Micro F1 / Macro F1 (%)								Micro F1 (%)	Macro F1 (%)
	ORG/ ORG	ORG/ GPE	ORG/ MONEY	ORG/ DATE	PER/ ORG	PER/ UNIV	PER/ GOV	PER/ TITLE	TOTAL (AVG/SD)	TOTAL (AVG/SD)
BB	35.6/25.1	84.0/50.4	75.5/46.9	79.7/45.9	65.6/29.1	50.9/59.9	21.5/14.7	89.6/45.0	73.7 (0.7)	53.9 (1.3)
BL	36.6/23.9	84.4/51.0	77.7/49.1	79.7/45.5	65.9/30.4	54.9/57.5	17.3/11.1	90.7/45.5	74.2 (0.9)	50.3 (2.6)
FB	37.7/24.9	82.9/49.7	74.5/48.4	80.1/46.0	65.3/31.5	51.0/60.7	20.7/12.9	89.7/45.0	73.7 (0.8)	49.4 (1.8)
SB	40.0/25.1	82.8/46.4	73.0/36.6	79.2/45.1	65.2/31.1	56.7/56.0	11.9/8.9	89.5/44.9	73.8 (0.6)	54.5 (3.4)
SL	39.5/26.3	83.1/50.6	76.5/41.3	79.5/45.3	67.1/33.6	56.0/63.3	8.8/6.7	90.1/45.2	74.9 (0.4)	50.7 (1.5)
FS	38.6/24.3	84.1/50.1	72.9/36.5	78.4/44.6	65.3/32.0	49.4/55.2	0.0/0.0	90.1/45.2	73.5 (0.7)	45.6 (3.4)
RB	39.8/26.5	85.2/49.2	78.2/49.7	80.8/46.2	66.4/32.7	56.1/64.7	19.5/12.9	90.1/45.1	74.4 (1.1)	55.6 (0.8)
RL	41.0/26.6	84.4/52.5	78.2/57.3	81.9/45.5	65.8/33.5	61.2/65.0	26.7/16.1	90.4/45.3	74.9 (0.6)	53.2 (2.8)
FR	38.4/25.2	83.3/51.3	74.6/46.1	79.8/46.3	64.8/29.9	54.3/61.3	12.3/8.9	90.3/45.3	74.8 (0.7)	47.4 (5.8)
FinB	38.1/25.0	82.5/49.8	76.5/51.6	80.7/46.2	65.9/30.6	49.6/57.8	15.9/11.1	89.4/44.9	73.3 (0.7)	53.8 (0.9)
MTB	39.9/25.5	82.2/50.0	75.0/41.9	80.3/45.3	64.6/33.5	62.4/66.5	19.0/11.7	90.7/45.5	74.2 (0.6)	54.4 (0.2)
CP	37.6/23.6	83.0/51.7	76.1/45.6	78.4/44.7	64.6/29.1	51.8/58.7	24.4/16.1	89.3/44.7	73.5 (0.8)	55.6 (2.3)
LB	38.9/25.8	84.2/52.5	77.9/52.8	79.5/45.6	66.1/33.1	53.5/61.1	22.1/13.9	90.7/45.4	74.9 (0.3)	56.3 (2.1)
LL	38.9/25.2	84.4/52.3	78.3/54.7	78.0/44.5	66.0/35.7	54.9/59.5	24.9/16.1	90.7/45.5	74.2 (0.5)	56.5 (1.1)

Figure 10: REFinD Baselines. Results achieved by the benchmark models. Reported are the % averages (AVG) and standard deviation (SD) over five random seeds. Here the models correspond to (i) BERT-base (BB) [4], (ii) BERT-large (BL) [4], (iii) FLANG-BERT (FB) [15], (iv) SpanBERT-base (SB) [10], (v) SpanBERT-large (SL) [10], (vi) FLANG-SpanBERT (FS) [15], (vii) Roberta-base (RB) [11], (viii) Roberta-large (RL) [11], (ix) FLANG-Roberta (FR) [15], (x) FinBERT (FinB) [2], (xi) Matching the Blanks (MTB) [17], (xii) Contrastive Pre-training (CP) [14], (xiii) Luke-base (LB) [21], and (xiv) Luke-large (LL) [21].

5.1 Experimental Setup

To perform relation extraction on our REFinD dataset, we adopt the architecture of Matching the Blanks [17]. First, we augment four entity markers, namely $[E1]$, $[/E1]$, $[E2]$, and $[/E2]$, to mark the beginning and end of each entity mention in a given relation instance s and ordered pair of entity mentions (e_1, e_2) . We then use a linear classifier based on the concatenated representation of the final hidden states corresponding to the start tokens of the two entities $[E1]$ and $[E2]$ to solve the task. For fine-tuning, we use hyperparameters as outlined in Table 4.

5.2 Benchmark Models

We evaluate our REFinD dataset using various pre-trained encoders from HuggingFace⁶, including BERT-base and -large [4], Roberta-base and -large [11], and Spanbert-base and -large [10], which have been pre-trained primarily on web-based articles. We also evaluate models that have been further trained on financial text, including FinBERT [2] and Flang encoders, namely Flang-BERT, -SpanBERT, and -Roberta [15]. Additionally, we assess the performance of several state-of-the-art models that have been specifically pre-trained for relation extraction, such as Matching the Blanks (MTB) [17], Contrastive Pre-training (CP) [14], and Luke-base and -large [21], as benchmarks for our REFinD dataset.

5.3 Results

The evaluation results of the benchmark models mentioned in Section 5.2 are presented in Figure 10. Among the benchmarks, Luke models exhibit the best overall performance, in terms of both micro- and macro-F1 scores. We also observe that entity groups like ORG-ORG, PER-GOV, and PER-ORG, which are more prevalent in the finance domain, present greater challenges than other entity groups. Specifically, finance-domain-specific relations such as

⁶<https://huggingface.co/>

shares_of, *cost_of*, and *member_of* (see Figure 10) exhibit F1 scores lower than 30%, whereas general relations such as *title*, *headquartered_in*, and *employee_of* achieve F1 scores greater than 70%. This could be attributed to the fact that models such as BERT, MTB, and Luke have primarily been pre-trained on web-based articles, which exhibit a different data distribution compared to financial datasets like REFinD. Furthermore, models like FinBERT and FLANG, which have undergone additional training on financial news articles and perform masking on financial terms, still fail to achieve better results, as they have never been exposed to semantically complex finance-specific documents. Thus, we must consider these experiments as starting points, and further improvements in the financial relation extraction task are likely to be achieved through increasing model capacity and architectural innovations.

5.4 Error Analysis

To provide insights for further improvements, we analyze the errors in the predictions of the relation extraction task, as shown in Figure 11. The most common types of errors made by the models can be categorized into three groups:

Numerical Inference: Most current LLM models treat numbers in text in the same way as other tokens, resulting in the most common error being the failure to capture numeracy. For example, the models predict *acquired_by* for both of these instances: $[ENT2]Company A[/ENT2]$ owns 31% equity interest in $[ENT1]Company B[/ENT1]$ and $[ENT1]Company A[/ENT1]$ has 100% equity interest in $[ENT2]Company B[/ENT2]$. However, the acquisition has only occurred in the second instance.

Semantic Ambiguity: Another challenge was that the models also struggled with distinguishing between similar relation types, such as *member_of*, *employee_of*, and *founder_of*, for entity pairs belonging to the PER-ORG, PER-UNIV, and PER-GOV groups. For example, the models often predicted *employee_of* for both of these

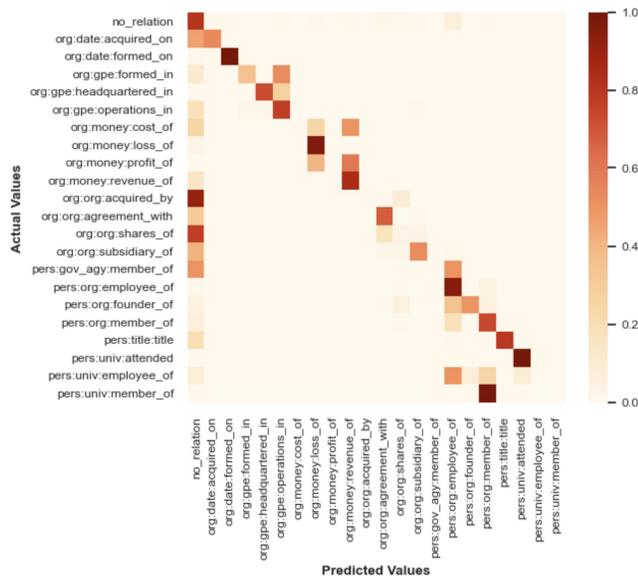


Figure 11: Confusion Matrix for each relation in REFinD based on predictions (normalized) over Luke-large encoder.

instances: $[ENT1] \text{ Jane Doe} [ENT1]$ is the CEO of $[ENT2]$ Company B $[ENT2]$ and $[ENT1] \text{ John Doe} [ENT1]$ is on the Board of Directors of $[ENT2]$ Company B $[ENT2]$. However, in the second instance, John Doe is only a member of Company B, not an employee. This suggests that the models may need to learn more nuanced distinctions between similar relation types, which could be addressed through improvements in training data, model architecture, or both.

Directional Ambiguity: The models exhibit uncertainty in determining the directional dependency between two entities for the relations *acquired_by* and *subsidiary_of*, despite being provided with the order of the entity pairs. This leads to errors where the models predict *subsidiary_of* for both instances in the phrases "Company A is the subsidiary of Company B" and "After purchasing Company A, through its wholly owned subsidiary, Company C purchased Company B". In reality, the second instance implies that Company B was acquired by Company A.

6 ETHICS STATEMENT

This work relies on the use of documents obtained from the SEC website. SEC filings are freely and publicly available, but do contain names and other identifying details such as current and former job titles, schools attended, and board or professional associations for typically high-ranking officers in publicly traded companies. This information is similar to what could be found about public figures mentioned in websites such as Wikipedia and thus we do not anticipate any harm to persons mentioned in our data beyond what could be learned from reading the public financial statements themselves. Additionally, we have made no attempt to aggregate or include non-public information about any individual or entity in this dataset. All data used for dataset construction was intentionally selected to be several years old as of this publication and thus we

do not anticipate any impact on financial markets with this release. REFinD is released under a license for non-commercial use.

7 CONCLUSION AND FUTURE WORK

In this paper, we introduced REFinD, which is the largest-scale annotated dataset of relations generated entirely over financial documents to-date, aimed to support the development of downstream applications within the finance domain. We highlighted the challenges involved in collecting such a large-scale relation extraction dataset specifically over financial text and emphasized the importance of financial-domain expertise in annotating such datasets. Furthermore, we fine-tuned various state-of-the-art deep learning models on the REFinD dataset to identify the challenges posed by our dataset. Our results showed that these models do not perform well on finance-specific relations, mainly because they have not been exposed to complex financial text and documents. We also identified three main challenges involved in relation extraction over financial text, namely numeracy, ambiguity amongst relations, and the direction of relations.

As for future directions, we plan to enhance the dataset by adding more types of finance-specific entity groups and hence more relations. Moreover, we aim to improve the dataset to capture the comparison between financial events over time, which is crucial in financial text analysis. For instance, we plan to include examples that capture the comparison between revenues earned by a company over different years, such as *Software LLC earned a net revenue of \$1.5 billion in 2018 as compared to \$1 billion in 2017*. These improvements will help us address the challenges involved in relation extraction over financial text and enable the development of more accurate and effective financial text analysis tools.

8 ACKNOWLEDGEMENTS

We would like to thank Natraj Raman, Daniel Borrajo, Armineh Nourbakhsh, Elena Kochkina, Zhiqiang Ma, and our anonymous reviewers for their thoughtful comments and feedback which greatly contributed to the quality of this work.

Disclaimer. This paper was prepared for informational purposes by the Artificial Intelligence Research group of JPMorgan Chase & Co. and its affiliates ("JP Morgan"), and is not a product of the Research Department of JP Morgan. JP Morgan makes no representation and warranty whatsoever and disclaims all liability, for the completeness, accuracy or reliability of the information contained herein. This document is not intended as investment research or investment advice, or a recommendation, offer or solicitation for the purchase or sale of any security, financial instrument, financial product or service, or to be used in any way for evaluating the merits of participating in any transaction, and shall not constitute a solicitation under any jurisdiction or to any person, if such solicitation under such jurisdiction or to such person would be unlawful.

A APPENDIX

A.1 MTurk Annotation Example

The entity pair (e_1, e_2) represents one instance. One sentence can have multiple entity pairs. Hence, multiple instances could be labeled from the same sentence. For example, in Figure 12, we see

REFERENCES

- [1] Eugene Agichtein and Luis Gravano. 2000. Snowball: Extracting relations from large plain-text collections. In *Proceedings of the fifth ACM conference on Digital libraries*. 85–94.
- [2] Dogu Araci. 2019. FinBERT: Financial Sentiment Analysis with Pre-trained Language Models. <https://doi.org/10.48550/ARXIV.1908.10063>
- [3] Elisa Bassignana and Barbara Plank. 2022. CrossRE: A Cross-Domain Dataset for Relation Extraction. In *Findings of the Association for Computational Linguistics: EMNLP 2022*. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 3592–3604. <https://aclanthology.org/2022.findings-emnlp.263>
- [4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. <https://doi.org/10.18653/v1/N19-1423>
- [5] Joseph L. Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological Bulletin* 76 (1971), 378–382.
- [6] Xu Han, Hao Zhu, Pengfei Yu, Ziyun Wang, Yuan Yao, Zhiyuan Liu, and Maosong Sun. 2018. FewRel: A Large-Scale Supervised Few-Shot Relation Classification Dataset with State-of-the-Art Evaluation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. 4803–4809.
- [7] Iris Hendrickx, Su Nam Kim, Zornitsa Kozareva, Preslav Nakov, Diarmuid Ó Séaghdha, Sebastian Padó, Marco Pennacchiotti, Lorenza Romano, and Stan Szpakowicz. 2010. SemEval-2010 Task 8: Multi-Way Classification of Semantic Relations between Pairs of Nominals. In *Proceedings of the 5th International Workshop on Semantic Evaluation*. Association for Computational Linguistics, Uppsala, Sweden, 33–38. <https://aclanthology.org/S10-1006>
- [8] Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. spaCy: Industrial-strength Natural Language Processing in Python. (2020). <https://doi.org/10.5281/zenodo.1212303>
- [9] Ali Jabbari, Olivier Sauvage, Hamada Zeine, and Hamza Chergui. 2020. A French corpus and annotation schema for named entity recognition and relation extraction of financial news. In *Proceedings of the 12th Language Resources and Evaluation Conference*. 2293–2299.
- [10] Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S Weld, Luke Zettlemoyer, and Omer Levy. 2020. Spanbert: Improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics* 8 (2020), 64–77.
- [11] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [12] Christopher D Manning, Mihai Surdeanu, John Bauer, Jenny Rose Finkel, Steven Bethard, and David McClosky. 2014. The Stanford CoreNLP natural language processing toolkit. In *Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations*. 55–60.
- [13] Mike Mintz, Steven Bills, Rion Snow, and Daniel Jurafsky. 2009. Distant supervision for relation extraction without labeled data. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*. Association for Computational Linguistics, Suntec, Singapore, 1003–1011. <https://aclanthology.org/P09-1113>
- [14] Hao Peng, Tianyu Gao, Xu Han, Yankai Lin, Peng Li, Zhiyuan Liu, Maosong Sun, and Jie Zhou. 2020. Learning from Context or Names? An Empirical Study on Neural Relation Extraction. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 3661–3672. <https://doi.org/10.18653/v1/2020.emnlp-main.298>
- [15] Raj Shah, Kunal Chawla, Dheeraj Eidnani, Agam Shah, Wendi Du, Sudheer Chava, Natraj Raman, Charese Smiley, Jiaao Chen, and Diyi Yang. 2022. When FLUE Meets FLANG: Benchmarks and Large Pretrained Language Model for Financial Domain. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 2322–2335. <https://aclanthology.org/2022.emnlp-main.148>
- [16] Soumya Sharma, Tapas Nayak, Arusarka Bose, Ajay Kumar Meena, Koustuv Dasgupta, Niloy Ganguly, and Pawan Goyal. 2022. FinRED: A Dataset for Relation Extraction in Financial Domain. In *Companion Proceedings of the Web Conference 2022*. 595–597.
- [17] Livio Baldini Soares, Nicholas FitzGerald, Jeffrey Ling, and Tom Kwiatkowski. 2019. Matching the Blanks: Distributional Similarity for Relation Learning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 2895–2905.
- [18] Stephanie M Strassel, Mark A Przybocki, Kay Peterson, Zhiyi Song, and Kazuaki Maeda. 2008. Linguistic Resources and Evaluation Techniques for Evaluation of Cross-Document Automatic Content Extraction.. In *LREC*. Citeseer.
- [19] Xiaozhi Wang, Yulin Chen, Ning Ding, Hao Peng, Zimu Wang, Yankai Lin, Xu Han, Lei Hou, Juanzi Li, Zhiyuan Liu, Peng Li, and Jie Zhou. 2022. MAVEN-ERE: A Unified Large-scale Dataset for Event Coreference, Temporal, Causal, and Subevent Relation Extraction. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 926–941. <https://aclanthology.org/2022.emnlp-main.60>
- [20] Haoyu Wu, Qing Lei, Xinyue Zhang, and Zhengqian Luo. 2020. Creating A Large-Scale Financial News Corpus for Relation Extraction. *2020 3rd International Conference on Artificial Intelligence and Big Data (ICAIBD)* (2020), 259–263.
- [21] Ikuya Yamada, Akari Asai, Hiroyuki Shindo, Hideaki Takeda, and Yuji Matsumoto. 2020. LUKE: Deep Contextualized Entity Representations with Entity-aware Self-attention. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 6442–6454. <https://doi.org/10.18653/v1/2020.emnlp-main.523>
- [22] Dongxu Zhang and Dong Wang. 2015. Relation classification via recurrent neural network. *arXiv preprint arXiv:1508.01006* (2015).
- [23] Yuhao Zhang, Victor Zhong, Danqi Chen, Gabor Angeli, and Christopher D Manning. 2017. Position-aware attention and supervised data improve slot filling. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. 35–45.
- [24] Zhe Zuo, Michael Loster, Ralf Krestel, and Felix Naumann. 2017. Uncovering Business Relationships: Context-sensitive Relationship Extraction for Difficult Relationship Types.. In *LWDA*. 271.