

Cross-lingual Inference with A Chinese Entailment Graph

Anonymous ACL submission

Abstract

Predicate entailment detection is a crucial task for question-answering from text, where previous work has explored unsupervised learning of entailment graphs from typed open relation triples. In this paper, we present the first pipeline for building Chinese entailment graphs, which involves a novel high-recall open relation extraction (ORE) method and the first Chinese fine-grained entity typing dataset under the FIGER type ontology. Through experiments on the Levy-Holt dataset and a boolean QA task, we verify the strength of our Chinese entailment graph, and reveal the cross-lingual complementarity: on the parallel Levy-Holt dataset, an ensemble of Chinese and English entailment graphs beats both monolinguals, and raises unsupervised SOTA by 4.7 AUC points¹.

1 Introduction

Predicate entailment detection is important for many tasks of natural language understanding (NLU), including reading comprehension and semantic parsing. Suppose we wish to answer a question by finding a relation V between entities A and B . Often, V cannot be found directly from the reference passage or database, but another relation U can be found between A and B , where U entails V (for instance, suppose U is *buy*, V is *own*). If we can identify this with predicate entailment detection, we can then answer the question.

To detect predicate entailments, previous work has explored unsupervised learning of typed entailment graphs (Szpektor and Dagan, 2008; Berant et al., 2011, 2015; Hosseini et al., 2018, 2019, 2021). Entailment graphs are directed graphs, where each node represents the predicate of a relation, and an edge from node U to node V denotes “ U entails V ”. Entailment graphs are built based on the Distributional Inclusion Hypothesis (DIH) (Dagan et al., 1999; Geffet and Dagan, 2005;

Herbelot and Ganesalingam, 2013; Kartsaklis and Sadrzadeh, 2016). Predicates are disambiguated according to their arguments’ types, predicates taking the same types of arguments go into one subgraph.

While previous work on entailment graphs has mostly been limited to English, building entailment graphs in other languages is interesting and challenging. The importance is two-fold: for that language, a native entailment graph would facilitate NLU in it; for cross-lingual inference, entailment graphs in different languages host exploitable complementary information. In particular, we argue that by jointly consulting strong entailment graphs in multiple languages, improvements can be gained for inference in **all** participant languages.

In this paper, we choose Chinese as our target language to build entailment graphs, as it is distant enough from English to exhibit rich complementarity, while relatively high-resource. The main challenge for building Chinese entailment graphs, is to extract reliable **typed relation triples** from raw corpora as strong input. This involves open relation extraction (ORE) and fine-grained entity typing (FET), which we discuss below.

ORE extracts predicate-argument triples from sentences, where previous work has used rule-based methods over syntactic parsers either directly (Fader et al., 2011; Etzioni et al., 2011; Angeli et al., 2015), or for distant supervision (Cui et al., 2018; Stanovsky et al., 2018; Kolluru et al., 2020). The challenge in ORE can be largely attributed to the poor definition of “open relations”. The situation worsens in Chinese, as the parts of speech are more ambiguous and many linguistic indicators of relations are poorly represented. Previous work on Chinese ORE (Qiu and Zhang, 2014; Jia et al., 2018) has defined narrow sets of open relations, failing to identify many relational constructions. Conversely, we propose a novel dependency-based ORE method, which we claim provides comprehensive coverage of relational constructions.

¹Our codes and data-sets will be available on Github.

FET assigns types to the arguments of extracted relations, so that word-senses of predicates can be disambiguated. The challenge in Chinese FET lies mainly in the lack of datasets over a suitable type ontology: too coarse a type set would be insufficient for disambiguation, too granular a type set would result in disastrous sparsity in the entailment graph. Following [Hosseini et al. \(2018\)](#), we use the popular FIGER type set ([Ling and Weld, 2012](#)), and build CFIGER, the first FIGER-labelled Chinese FET dataset via label mapping. Entity typing models built on this dataset show satisfactory accuracy and are helpful for predicate disambiguation.

We evaluate our Chinese entailment graph on the Levy-Holt entailment dataset ([Levy and Dagan, 2016](#); [Holt, 2019](#)) (via translation) and a natively-Chinese boolean QA task following [McKenna et al. \(2021\)](#). Results show that our Chinese entailment graph outperforms baselines by large margins, and is comparable to the English graph. We verify our cross-lingual complementarity argument on Levy-Holt dataset: by ensembling English and Chinese graphs, we show a clear advantage over both monolingual graphs, and set a new SOTA.

Our contributions are as follows: 1) we present a novel Chinese ORE method sensitive to a much wider range of relations than previous SOTA, and a Chinese FET dataset, the first under the FIGER type ontology; 2) we construct the first Chinese entailment graph, comparable to its English counterpart; 3) we reveal the cross-lingual complementarity of entailment graphs with an ensemble.

2 Background and Related Work

Predicate entailment detection has been an area of active research. [Lin \(1998\)](#); [Weeds and Weir \(2003\)](#); [Szpektor and Dagan \(2008\)](#) proposed various count-based entailment scores; [Berant et al. \(2011\)](#) proposed to “globalize” typed entailment graphs by closing them with transitivity constraints; [Hosseini et al. \(2018\)](#) proposed a more scalable global learning approach with soft transitivity constraints; [Hosseini et al. \(2019, 2021\)](#) further refined the entailment scores with link prediction.

Our work is closely related to [Hosseini et al. \(2018\)](#), with key adaptations for Chinese in ORE and FET. Their ORE method is based on a CCG parser ([Reddy et al., 2014](#)), while ours is based on a dependency parser ([Zhang et al., 2020](#)); their FET is done by linking entities to Wikipedia entries, while we use neural entity typing for the task.

Dependency parses are less informative than CCG parses, and require heavier adaptation. However, Chinese dependency parsers are currently more reliable than CCG parsers ([Tse and Curran, 2012](#)). Previous Chinese ORE methods ([Qiu and Zhang, 2014](#); [Jia et al., 2018](#)) are based on dependency parsers, but they omit many common constructions essential to ORE. In §3, we present the most comprehensive Chinese ORE method so far.

Linking-based entity-typing can be more accurate than neural methods, since the type labels are exact as long as linking is correct. However, current Chinese entity linking methods require either translation ([Pan et al., 2019](#)) or search logs ([Fu et al., 2020](#)). Both hurt linking accuracy, and the latter grows prohibitively expensive with scale. On the other hand, since the seminal work of [Ling and Weld \(2012\)](#), neural fine-grained entity typing has developed rapidly ([Yogatama et al., 2015](#); [Shimaoka et al., 2017](#); [Chen et al., 2020](#)), with a common interest in FIGER type set. For Chinese, [Lee et al. \(2020\)](#) built an ultra-fine-grained entity typing dataset, based on which we build our CFIGER dataset via label mapping.

[Weber and Steedman \(2019\)](#) aligned English and German entailment graphs, and showed that the English graph can help with German entailment detection. Yet it was uncertain whether this effect comes from genuine complementarity or the mere fact that the English graph is stronger. We take one step further, and show that complementarity can be exploited in both directions: for English, the higher resource language, entailment detection can also benefit from the ensemble to reach new heights.

As a related resource, [Ganitkevitch et al. \(2013\)](#) created a multi-lingual database for symmetric paraphrases; in contrast, entailment graphs are directional. More recently, [Schmitt and Schütze \(2021\)](#) proposed to fine-tune language models on predicate entailment datasets via prompt learning. In contrast to our entailment graphs, their approach is supervised, which carries the danger of overfitting to dataset artifacts ([Gururangan et al., 2018](#)).

Another related strand of research, e.g. SNLI ([Bowman et al., 2015](#)), is concerned with the more general NLI task, including hypernymy detection and logic reasoning like $A \wedge B \rightarrow B$, but rarely covers the cases requiring external knowledge of predicate entailment. Conversely, entailment graphs aim to serve as a robust resource for directional predicate entailments induced from textual corpora.

3 Chinese Open Relation Extraction

We build our ORE method based on DDParser (Zhang et al., 2020), a SOTA Chinese dependency parser. We mine relation triples from its output by identifying patterns in the dependency paths.

Depending on the semantics of the head verb, instances of a dependency pattern can range from being highly felicitous to marginally acceptable as a relation. Motivated by our downstream task of entailment graph construction, we go for higher recall and take them in based on the **Relation Frequency Assumption**: less felicitous relations occur less frequently, and are less likely to take part in entailments when they do occur, thus they are negligible.

Due to the lack of a commonly accepted benchmark or criterion for “relations”, we did not perform an intrinsic evaluation for our Chinese ORE method; its significant benefit to our EG_{Zh} graph (§6.3, §7) should suffice to demonstrate its strength.

3.1 Parsing for Chinese ORE

The task of open relation extraction on top of LM-driven dependency parsers, is really the task of binding the relations in surface forms to the underlying relation structures. Though trivial at first sight, the definition of these underlying and essentially semantic relations demands detailed analysis.

Jia et al. (2018) is the latest to propose an ORE method on dependency paths. They defined a set of rules to extract relations patterns, which they call dependency semantic normal forms (DSNFs)².

However, their set of DSNFs is inexhaustive and somewhat inaccurate. We show below that many linguistic features of Chinese demand a more principled account, more constructions need to be considered as relations, some to be ruled out. These observations are made from a multi-source news corpus, which we use to build entailment graphs (§5)³. Below, we highlight 5 additional constructions we identify, explained with examples⁴.

²We refer readers to Appendix A for a brief summary.

³A referee has commented that when refining our ORE method, we might have inadvertently or unconsciously fine-tuned the system for the evaluation tasks. However, as entailment graph construction is fully unsupervised, this source corpus is independent of the evaluations in §6 and §7. Particularly, the Levy-Holt dataset used in §6 consists of short sentences, which is a vastly different genre, involving much simpler structures, with a single relation per sentence and few subordinating constructions discussed above (see Appendix L for supporting statistics); the QA dataset used in §7 is built from a separate news corpus, strictly excluding overlaps with those used to develop the parser and the entailment graphs.

⁴We refer readers to Appendix J for diagram illustrations.

A. PP Modifiers as “De” Structures One key feature of Chinese is its prevalent use of “De” structures in the place of prepositional phrases, where “De” can be seen as roughly equivalent to the possessive clitic 的. For instance, in “咽炎(*pharyngitis*) 成为(*becomes*) 发热(*fever*) 的(*De*) 原因(*cause*); *Pharyngitis becomes the cause of fever*”, the root clause in Chinese is (Pharyngitis, becomes, cause), but we *additionally* extract the underlying relation (**pharyngitis, becomes-X·De-cause, fever**), where the true object “fever” is a **nominal** attribute of the direct object “cause”, and the true predicate subsumes the direct object⁵.

The same also applies to the subject, though somewhat more restricted. For sentences like “苹果(*Apple*) 的(*De*) 创始人(*founder*) 是(*is*) 乔布斯(*Jobs*); *The founder of Apple is Jobs*”, we additionally extract the relation (**Apple, founder-is, Jobs**), where the true subject “Apple” is a **nominal** attribute of the direct subject “founder”, and the true predicate subsumes the direct subject⁶.

B. Bounded Dependencies In Chinese, bounded dependencies, especially control structures, are expressed with a covert infinitival marker, equivalent to English “to”. We capture the following phenomena in addition to direct relations:

- Sequences of VPs: for sentences like “我(*I*) 去(*go-to*) 诊所(*clinic*) 打(*take*) 疫苗(*vaccine*); *I go to the clinic to take the vaccine*”, the two verb phrases “去(*go-to*) 诊所(*clinic*)” and “打(*take*) 疫苗(*vaccine*)” are directly concatenated, with no overt connection words. Here we additionally extract the relation (**I, take, vaccine**) by copying the subject of the head verb to subsequent verbs.
- Subject-control verbs: for the famous example “我(*I*) 想(*want*) 试图(*try*) 开始(*begin*) 写(*write*) 一个(*a*) 剧本(*play*); *I want to try to begin to write a play*”, again the verbs are directly concatenated; this time, all verbs but the first one behaves as infinitival complements to their direct antecedents. In such cases, we extract sequences of relations like (**I, want, try**), (**I, want·try, begin**), (**I, want·try·begin, write**), (**I, want·begin·try·write, a play**).

Notably, the above relations are different from Jia

⁵Here and below, examples are paired with English metaphrases, and when necessary, paraphrases; relation triples are presented as English metaphrases (inflections ignored).

⁶These relations are more felicitous with frequent predicate argument combinations, and less so for the infrequent ones. As in line with the Relation Frequency Assumption, less felicitous relations are also less statistically significant.

et al. (2018)’s conjunctions in Table 4: the event sequences here involve subordination (control) rather than coordination, thus need a separate account.

C. Relative Clauses Relative Clauses also take the form of modification structures in Chinese, for which additional relations should also be extracted. For example, in “他(*he*) 解决(*solve*) 了(*-ed*) 困扰(*puzzle*) 大家(*everyone*) 的(*De*) 问题(*problem*); *He solved the problem that puzzled everyone*”, we extract not only the direct relation (**he, solve, problem**), but also the relation embedded in the modification structure (**problem, puzzle, everyone**).

D. Nominal Compounds Relations can be extracted from nominal compounds, where an NP has two consecutive “ATT” modifiers: in “德国(*Germany*) 总理(*Chancellor*) 默克尔(*Merkel*); *German Chancellor Merkel*”, “Germany” modifies “Chancellor”, and “Chancellor” modifies “Merkel”. Jia et al. (2018) extracted relations like (**Germany, Chancellor, Merkel**) for these NPs.

However, they overlooked the fact that prepositional compounds in Chinese with omitted “De” take exactly the same form (see construction A). For example, in NPs with nested PP modifiers like “手续(*formalities*) 办理(*handle*) 时效(*timeliness*); *Timeliness of the handling of formalities*”, we have the same structure, but it certainly does not mean “*the handling of formalities is timeliness*”!

We take a step back and put restrictions on such constructions: only when all 3 words in the NP are nominals (but not pronouns), the third word is the head, the second is a ‘PERSON’ or ‘TITLE’, and the first is a ‘PERSON’, then it is a relation, like (**Merkel, is·X·De·Chancellor, Germany**). Otherwise, such NPs rarely host felicitous relations.

E. Copula with Covert Objects The copula is sometimes followed by modifiers ending with “De”. Examples are “玉米(*Corn*) 是(*is*) 从(*from*) 美国(*US*) 引进(*introduce*) 的(*De*); *Corn is introduced from US*”, “设备(*device*) 是(*is*) 木头(*wood*) 做(*make*) 的(*De*); *The device is made of wood*”.

In these cases, there exists an object following the indicator “的(*De*)”, but the object is an empty *pro* considered inferable from context. In the absence of the true object, the *VOB* label is given to “的(*De*)”, leading to direct relations like (**Corn, is, De**). However, the true predicates are rather “*is introduced from*” or “*is made of*”. To fix this, we replace the direct relations with ones like (**Corn, is·from·X·introduce·De·pro, America**), reminiscent of the constructions A.

3.2 Our ORE Method

With the above constructions taken into account, we build our ORE method on top of DDParse. For part-of-speech labels, we use the POS-tagger in Stanford CoreNLP (Manning et al., 2014). We detect negations by looking for negation keywords in the adjunct modifiers of predicates: for predicates with an odd number of negation matches, we insert a negation indicator to them, treating them as separate predicates from the non-negated ones.

4 Chinese Fine-Grained Entity Typing

As shown in previous work (Berant et al., 2011; Hosseini et al., 2018), the types of a predicate’s arguments are helpful for disambiguating a predicate in context. To this end, we need a fine-grained entity typing model to classify the arguments into sufficiently discriminative yet populous types.

Lee et al. (2020) presented CFET dataset, an ultra-fine-grained entity typing dataset in Chinese. They labelled entities in sentence-level context, into around 6,000 free-form types and 10 general types. Unfortunately, their free-form types are too fragmented for predicate disambiguation, and their general types are too ambiguous.

We turn to FIGER (Ling and Weld, 2012), a commonly used type set: we re-annotate the CFET dataset with FIGER types through label mapping. Given that there are around 6,000 ultra-fine-grained types and only 112 FIGER types (49 in the first layer), we can reasonably assume that each ultra-fine-grained type can be unambiguously mapped to a single FIGER type. For instance, the ultra-fine-grained type “湖 (lake)” is unambiguously mapped to the FIGER label “location / body_of_water”.

Based on this assumption, we manually create a mapping between the two, and re-annotate CFET dataset with the mapping. We call the re-annotated dataset **CFIGER**, as it is the first in Chinese with FIGER labels. As with CFET, this dataset consists of 4.8K crowd-annotated data (equally divided into crowd-train, crowd-dev and crowd-test) and 1.9M distantly supervised data from Wikipedia⁷.

For training set we combine the crowd-train and Wikipedia subsets; for dev and test sets we use crowd-dev and crowd-test respectively. We train two baseline models: *CFET*, the baseline model with CFET dataset; *HierType* (Chen et al., 2020), a SOTA English entity typing model.

⁷For detailed statistics, please refer to Appendix B.

Macro F1 (%)	dev	test
<i>CFET</i> with CFET dataset	-	24.9
<i>CFET</i> with CFIGER dataset	75.7	75.7
<i>HierType</i> with FIGER dataset	-	82.6
<i>HierType</i> with CFIGER dataset	74.8	74.5

Table 1: F1 scores of baseline models for CFIGER dataset, compared with the results on the datasets where they were proposed. Macro-F1 scores are reported because it is available in both baselines.

Results are shown in Table 1: the F1 score of *HierType* model is slightly lower on CFIGER dataset than on FIGER dataset in English; contrarily, thanks to fewer type labels, the F1 score of *CFET* baseline increases on CFIGER, bringing it on par with the more sophisticated *HierType* model. This means our CFIGER dataset is valid for Chinese fine-grained entity typing, and may contribute to a benchmark for cross-lingual entity typing.

For downstream applications, we nevertheless employ the *HierType* model, as empirically it generalizes better to our news corpora. As shown in later sections, the resulting FET model can substantially help with predicate disambiguation.

5 The Chinese Entailment Graph

We construct the Chinese entailment graph from the Webhose corpus⁸, a multi-source news corpus crawled from 133 news websites in October 2016. Similarly to the NewsSpike corpus used in Hosseini et al. (2018), the Webhose corpus contains multi-source non-fiction articles from a short period of time. This means it is also rich in reliable and diverse relation triples over a focused set of events, ideal for building entailment graphs.

For the 313K valid articles in Webhose, we get their CoreNLP POS tags and feed them into our ORE method in §3, to extract the open relation triples. Then, with *HierType* model (Chen et al., 2020) in §4, we type all arguments of the extracted relations; we type each predicate with its subject-object type pair, such as *person-event* or *food-law*.

We finally employ the entailment graph construction method in Hosseini et al. (2018), taking in only binary relations. The detailed statistics of our Chinese entailment graph are shown in Table 2: compared with EG_{En} , our graph is built on just over half the number of articles, yet we have extracted 70% the number of relation triples, and built a graph

⁸<https://webhose.io/free-datasets/chinese-news-articles/>

	EG_{Zh}	EG_{En}
# of articles taken	313,718	546,713
# of triples used	7,621,994	10,978,438
# of predicates	363,349	326,331
# of type pairs where:		
subgraph exists	942	355
$ \text{subgraph} > 100$	442	115
$ \text{subgraph} > 1,000$	149	27
$ \text{subgraph} > 10,000$	26	7

Table 2: Stats of our Chinese entailment graph (EG_{Zh}) compared with the English graph in Hosseini et al. (2018) (EG_{En}). $|\cdot|$ denotes the number of predicates.

involving even more predicates. In general, our EG_{Zh} is of comparable size to EG_{En} . We encourage interested readers to check Appendix D for details of our graph-building process and a quick introduction to Hosseini et al. (2018).

6 Evaluation by Entailment Detection

6.1 Benchmark and Baselines

To evaluate the quality of our Chinese entailment graph, we first perform an intrinsic evaluation on the predicate entailment detection task. Our experiments are based on the popular Levy-Holt dataset (Levy and Dagan, 2016; Holt, 2019), with the same dev/test configuration as Hosseini et al. (2018). We convert the Levy-Holt dataset to Chinese through machine translation, then do evaluation on the translated premise-hypothesis pairs.

In Levy-Holt dataset, the task is: to take as input a pair of relation triples about the same arguments, one premise and one hypothesis, and judge whether the premise entails the hypothesis. To convert Levy-Holt dataset into Chinese, we concatenate each relation triple into a pseudo-sentence, use Google Translate to translate the pseudo-sentences into Chinese, then parse them back to Chinese relation triples with our ORE method in §3. If multiple relations are returned, we retrieve the most representative ones, by considering only those relations whose predicate covers the HEAD word.⁹

To type the Chinese relation triples, we again use *HierType* model to collect their subject-object type-pairs. The premise and hypothesis need to take the same types of arguments, so we take the intersection of their possible type-pairs¹⁰. We search the entailment subgraphs of these type-pairs, for entailment edges from the premise to the hypothesis,

⁹See Appendix C for more details.

¹⁰Unless the intersection is empty, then we take the union.

and return the entailment scores associated with these edges. When edges are found from multiple subgraphs, we take their maximum score¹¹.

We compare our Chinese entailment graph with a few strong baselines:

BERT: We take the translated pseudo-sentence pairs, and compute the cosine similarity between their pretrained BERT representations at [CLS] token. This is a strong baseline but symmetric;

Jia: We build entailment graph in the same way as §5, but with the baseline ORE method by Jia et al. (2018); accordingly, Jia et al. (2018) method is also used in parsing the translated Levy-Holt;

DDPORE: Similar to Jia baseline, but with the baseline ORE method from DDParse (2020).

6.2 Cross-lingual Ensembles

In order to examine the complementarity between our Chinese entailment graph (EG_{Zh}) and the English graph (EG_{En}) (2018), we ensemble the predictions from the two graphs, $pred_{en}$ and $pred_{zh}$ ¹². We experiment with four ensemble strategies: lexicographic orders from English to Chinese and Chinese to English, max pooling and average pooling:

$$pred_{en_zh} = pred_{en} + \gamma * \Theta(pred_{en}) * pred_{zh}$$

$$pred_{zh_en} = \gamma * pred_{zh} + \Theta(pred_{zh}) * pred_{en}$$

$$pred_{max} = MAX(pred_{en}, \gamma * pred_{zh})$$

$$pred_{avg} = AVG(pred_{en}, \gamma * pred_{zh})$$

where $\Theta(\cdot)$ is the boolean function *IsZero*, γ is the relative weight of Chinese and English graphs. γ is a hyperparameter tuned on Levy-Holt dev set, searched between 0.0 and 1.0 with step size 0.1.

For instance, suppose our premise is “he, shopped in, the store”, and our hypothesis is “he, went to, the store”, then our Chinese relations, by translation, would be “他, 在·X·购物, 商店” and “他, 前往, 商店” respectively. Suppose we find in the English graph an edge from “shop in” to “go to”, scored $pred_{en} = 0.6$, and we find in the Chinese graph an edge from “在·X·购物” to “前往”, scored $pred_{zh} = 0.7$. Then we would have $pred_{en_zh} = 0.6$, $pred_{zh_en} = 0.7$, $pred_{max} = 0.7$, $pred_{avg} = 0.65$.

In addition to ensembling with EG_{En} , we also ensembled our entailment graph with the SOTA English graph EG_{En++} (2021). We call the latter ones **Ensemble++** here and below.

¹¹We provide a human evaluation on the quality of the resulting Chinese Levy-Holt dataset in Appendix I.

¹²“zh” is the abbreviation for Chinese by convention.

AUC (%)	dev	test
<i>BERT</i> *	5.5	3.2
<i>Jia</i> (2018) *	0.9	2.4
<i>DDPORE</i> (2020) *	9.8	5.9
EG_{Zh} *	15.7	9.4
EG_{En} (2018) \diamond	20.7	16.5
EG_{En++} (2021) \diamond	23.3	19.5
Ensemble En_Zh \diamond	28.3 ($\gamma : 0.8$)	21.2
Ensemble Zh_En \diamond	27.4 ($\gamma : 0.9$)	21.5
Ensemble MAX \diamond	29.9 ($\gamma : 0.8$)	22.1
Ensemble AVG \diamond	30.0 ($\gamma : 1.0$)	22.1 †
Ensemble++ AVG \diamond	31.2 ($\gamma : 0.3$)	24.2 †
EG_{Zh} -type *	11.1	7.0
DataConcat En \diamond	20.6	17.8
DataConcat Zh *	19.0	14.2
DataConcat Esb \diamond	31.8	25.0
BackTrans Esb \diamond	23.0	17.5

Table 3: Area Under Curve values on Levy-Holt dataset, for Chinese entailment graph (EG_{Zh}), its baselines, ensembles with English graphs, and ablation studies. EG_{En} is the English graph from (Hosseini et al., 2018); EG_{En++} is the English graph from (Hosseini et al., 2021). Entries with * uses Chinese lemma baseline; entries with \diamond uses English lemma baseline; entries with † are the best ensemble strategies by dev set results.

6.3 Results and Discussions

To measure the performance of our constructed Chinese entailment graphs, we follow previous work in reporting the Precision-Recall (P-R) Curves plotted for successively lower confidence thresholds, and their Area Under Curves (AUC), for the range with > 50% precision.

A language-specific lemma baseline sets the left boundary of recall, by exact match over the lemmatized premise / hypothesis. For our Chinese entailment graph (EG_{Zh}) and its baselines, the boundary is set by Chinese lemma baseline. For the ensembles, in order to get commensurable AUC values with previous work instead of being over-optimistic, we use the English lemma baseline.

As shown in Table 3, on the Chinese Levy-Holt dataset, our EG_{Zh} graph substantially outperforms the BERT pre-trained baseline. EG_{Zh} is also far ahead of entailment graphs with baseline ORE methods, proving the superiority of our Chinese ORE method against previous SOTA.

EG_{Zh} and EG_{En} are built with the same algorithm (Hosseini et al., 2018), and evaluated on parallel datasets. Learnt from 57% the data, EG_{Zh} achieves an AUC exactly 57% of its English counterpart. Note that the Chinese entailment graph is under-

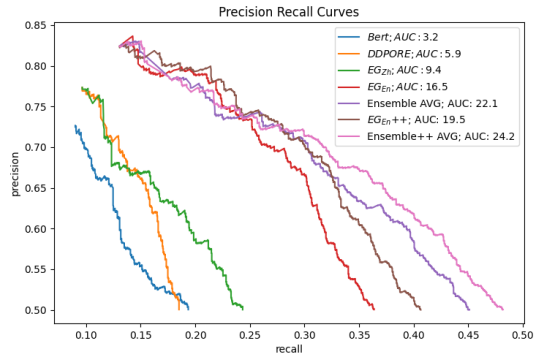


Figure 1: P-R Curves on Levy-Holt test set for EG_{Zh}, ensembles and baselines; Jia(2018) baseline is far behind others, thus omitted for the clarity of the figure.

510 estimated with the use of translated dataset: out
 511 of the 12,921 relation pairs in Levy-Holt test set,
 512 only 9,337 of them are parsed into valid Chinese bi-
 513 nary relations. This means, for Chinese entailment
 514 graphs, the upper bound for recall is not 100%, but
 515 rather 72.3%, as is the upper bound for AUC. Be-
 516 sides, the translationese language style in Chinese
 517 Levy-Holt also poses a gap in word-choice to the
 518 natively-built entailment graph, resulting in more
 519 mismatches. Considering the above extra noise,
 520 the performance of EG_{Zh} means our pipeline is
 521 utilizing the source corpus very well.

522 The ensemble between Chinese and English en-
 523 tailment graphs sets a new SOTA for unsupervised
 524 predicate entailment detection. With all 4 ensem-
 525 ble strategies, improvement is gained upon both
 526 monolingual graphs; with **Ensemble AVG**, the best
 527 on dev-set, the margin of test set improvement is
 528 more than 5 points. Moreover, when ensembling
 529 with EG_{En}++, we get a test-set AUC of 24.2 points
 530 (**Ensemble++ AVG**), raising SOTA by 4.7 points.

531 In Table 3, we additionally present three **abla-**
 532 **tion studies** to verify the solidity of our approach.

533 In the first ablation study, *EG_{Zh}-type*, we take
 534 away entity typing and train an untyped entailment
 535 graph. In this setting, we lose 2.4 AUC points. This
 536 means our entity typing method is indeed helpful
 537 for disambiguating predicates in entailment graphs.

538 In the second ablation study, the *DataConcat*
 539 settings, we disentangle cross-lingual complemen-
 540 tarity from the effect of extra data. We machine-
 541 translate NewsSpike corpus into Chinese, Webhose
 542 into English. We build an English graph “DataCon-
 543 cat En” using *NewsSpike + translated-Webhose*,
 544 and a Chinese graph “DataConcat Zh” using *Web-*
 545 *hose + translated-NewsSpike*. Results show that
 546 while both graphs improve with data from the other
 547 side, they are still far behind our **Ensemble** settings

548 above. Further, ensembling the two DataConcat
 549 graphs delivers an AUC of 25.0 points, 7.2 higher
 550 than DataConcat En, an even wider margin than our
 551 main setting. This suggests, the success of cross-
 552 lingual ensemble **cannot** be reproduced by sticking
 553 all the data together for a monolingual graph.

554 In the third case study, *BackTrans Esb*, we dis-
 555 entangle cross-lingual complementarity from the
 556 effect of machine-translation. Machine translation
 557 can be noisy, but it can also map synonyms in the
 558 source language to the same words in the target
 559 language. To single out this effect, we translate the
 560 Chinese Levy-Holt dataset back into English, and
 561 perform an ensemble between predictions on the
 562 original and the back-translated Levy-Holt. The
 563 gain in this case is only marginal, suggesting that
 564 cross-lingual complementarity is the reason for our
 565 success, while the synonym effect is not.

566 To further analyse the improvements with our
 567 ensembles, we conduct a case study over the differ-
 568 ence in predictions between our ensemble and the
 569 English monolingual, thresholded at 65% precision.
 570 We find that the majority of the improvements can
 571 be attributed to the additional evidence of entail-
 572 ment; we refer readers to Appendix E for details.

573 In conclusion, from the entailment detection ex-
 574 periment, we have learnt that: 1) our Chinese en-
 575 tailment graph is strong in the monolingual setting,
 576 with contributions from the ORE method and en-
 577 tity typing; 2) a cross-lingual complementarity is
 578 clearly shown between Chinese and English en-
 579 tailment graphs, where the effect of ensembles is
 580 most significant in the moderate precision range
 581 (see Figure 1). We expect that ensembling strong
 582 entailment graphs in more languages would lead to
 583 further improvements.

584 7 Evaluation by Question Answering

585 In addition to the entailment detection evaluation,
 586 we further demonstrate the strength of our Chinese
 587 entailment graph in application with an extrinsic
 588 question answering task, natively in Chinese. Since
 589 the existing Chinese QA datasets (Cui et al., 2019;
 590 Zheng et al., 2019; Sun et al., 2019) barely concern
 591 predicate entailments, we evaluate using a more ad-
 592 versarial boolean QA task following McKenna et al.
 593 (2021), inspired by Poon and Domingos (2009).

594 This task is designed as a boolean variant of *ma-*
 595 *chine reading at scale* (Chen et al., 2017): given a
 596 proposition and a pool of context articles, a model
 597 attempts to judge whether the queried proposi-

tion is true, based on the context pool. Our QA dataset is built off CLUE (Xu et al., 2020), a huge news corpus similar to Webhose¹³. The assumption is, frequently-mentioned predicates between frequently-mentioned arguments are high-quality events to be used as positives; absent predicates between frequently-mentioned arguments are probably not true, and can be used as negatives.

Articles in CLUE corpus are parsed into relation triples as in §3, and partitioned into 3-day time spans to get uniquely identifiable events. Those relation triples whose predicates appear over 30 times in the corpus, and whose argument-pairs appear in over 15 articles in their partitions, are selected as **high-quality positives**. For each positive, we generate negatives designed to be challenging for machines: following McKenna et al. (2021), we replace the positive predicates with their hyponyms / troponyms in Chinese WordNet (Wang and Bond, 2013). Since Chinese predicates are often multi-token and discontinuous, we look for substitutions of spans in predicates rather than entire predicates. If a substituted predicate is absent between this argument-pair in this partition, but exists elsewhere in the corpus with other argument-pairs, we consider it an **adversarial negative**: meaningful, but neither mentioned nor entailed in this context.

We use a sample of 40,000 positives along with their 70,583 adversarial negatives (a similar proportion to Levy-Holt) as our final QA dataset, split evenly into dev / test sets. We have omitted many details in dataset elicitation for the sake of brevity; we share this dataset as part of our release, and refer readers to McKenna et al. (2021) and Appendix H for more details and examples.

In our QA task, the positive / negative triples are concatenated into query propositions, the entire partition of articles in each query’s time-span¹⁴ is used as context pool, and a confidence score for each query is produced by each method. As in §6, Precision-Recall curves are drawn, AUC values with >50% precision are reported.

We compare our EG_{Zh} with the DDPORE baseline in §6 as well as BERT baselines in 3 different setups. Note that our QA dataset is monolingual in Chinese so ensembles are not involved. For all methods, the confidence score of a query is defined as its highest score w.r.t. any context sentences.

¹³Articles seen in Webhose are excluded for fairness of the experiment.

¹⁴Except the sentence hosting the query or its positive.

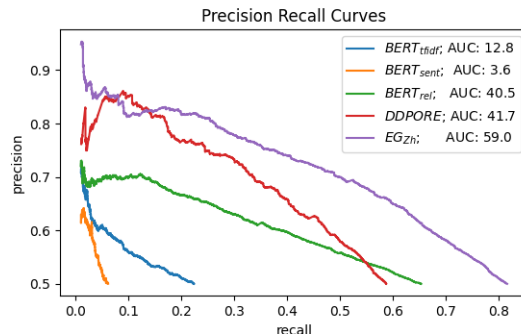


Figure 2: P-R Curves on QA evaluation test set for EG_{Zh} and baselines; AUC values are annotated in the legend.

BERT_{tfidf}: retrieves the top 5 most relevant articles by TF-IDF following (Chen et al., 2017), calculates cosine similarity between each retrieved sentence and the query at [CLS] token;

BERT_{sent}: retrieves the *host-sentences* of the relation triples involving the same arguments as the query, calculates cosine similarity between each host sentence and the query at [CLS] token;

BERT_{rel}: retrieves the *relation triples* involving the same arguments as the query, calculates cosine similarity between each retrieved triple (concatenated into a sentence) and the query at [CLS] token;

EG_{Zh} / DDPORE: retrieve the relation triples involving the same arguments as the query, return entailment scores from each retrieved context triple to the query triple (note that these are directional).

Results are shown in Figure 2. Our EG_{Zh} is again far above all baselines, further stressing the strength of our approach and the necessity of developing entailment graphs for directional inference. BERT_{rel} outperforms the other two BERT baselines, because of its focused context input of concatenated triples, in contrast to the more noisy news sentences for BERT_{tfidf} and BERT_{sent}.

8 Conclusion

We have presented a pipeline for building Chinese entailment graphs. Along the way, we proposed a novel high-recall open relation extraction method, and built a fine-grained entity-typing dataset via label mapping. As our main result, we have shown that: our Chinese entailment graph is comparable with English graphs, where unsupervised BERT baseline did poorly; an ensemble between Chinese and English entailment graphs substantially outperforms both monolinguals, and sets a new SOTA for unsupervised entailment detection. Directions for future work include multilingual EG alignment and alternative predicate disambiguation.

References

- 685 Gabor Angeli, Melvin Jose Johnson Premkumar, and
686 Christopher D. Manning. 2015. [Leveraging Lin-](#)
687 [guistic Structure For Open Domain Information Ex-](#)
688 [traction](#). In *Proceedings of the 53rd Annual Meet-*
689 *ing of the Association for Computational Linguistics*
690 *and the 7th International Joint Conference on Natu-*
691 *ral Language Processing (Volume 1: Long Papers)*,
692 pages 344–354, Beijing, China. Association for Com-
693 putational Linguistics.
- 694 Jonathan Berant, Noga Alon, Ido Dagan, and Jacob
695 Goldberger. 2015. [Efficient global learning of entail-](#)
696 [ment graphs](#). *Computational Linguistics*, 41(2):249–
697 291.
- 698 Jonathan Berant, Ido Dagan, and Jacob Goldberger.
699 2011. [Global Learning of Typed Entailment Rules](#).
700 In *Proceedings of the 49th Annual Meeting of the*
701 *Association for Computational Linguistics: Human*
702 *Language Technologies*, pages 610–619, Portland,
703 Oregon, USA. Association for Computational Lin-
704 guistics.
- 705 Samuel R. Bowman, Gabor Angeli, Christopher Potts,
706 and Christopher D. Manning. 2015. [A large anno-](#)
707 [tated corpus for learning natural language inference](#).
708 In *Proceedings of the 2015 Conference on Empiri-*
709 *cal Methods in Natural Language Processing*, pages
710 632–642, Lisbon, Portugal. Association for Compu-
711 tational Linguistics.
- 712 Danqi Chen, Adam Fisch, Jason Weston, and Antoine
713 Bordes. 2017. [Reading Wikipedia to Answer Open-](#)
714 [Domain Questions](#). *arXiv:1704.00051 [cs]*. ArXiv:
715 1704.00051.
- 716 Tongfei Chen, Yunmo Chen, and Benjamin Van Durme.
717 2020. [Hierarchical Entity Typing via Multi-level](#)
718 [Learning to Rank](#). In *Proceedings of the 58th Annual*
719 *Meeting of the Association for Computational Lin-*
720 *guistics*, pages 8465–8475, Online. Association for
721 Computational Linguistics.
- 722 Lei Cui, Furu Wei, and Ming Zhou. 2018. [Neural Open](#)
723 [Information Extraction](#). In *Proceedings of the 56th*
724 *Annual Meeting of the Association for Computational*
725 *Linguistics (Volume 2: Short Papers)*, pages 407–413,
726 Melbourne, Australia. Association for Computational
727 Linguistics.
- 728 Yiming Cui, Ting Liu, Wanxiang Che, Li Xiao, Zhipeng
729 Chen, Wentao Ma, Shijin Wang, and Guoping Hu.
730 2019. [A Span-Extraction Dataset for Chinese Ma-](#)
731 [chine Reading Comprehension](#). In *Proceedings of*
732 *the 2019 Conference on Empirical Methods in Natu-*
733 *ral Language Processing and the 9th International*
734 *Joint Conference on Natural Language Processing*
735 *(EMNLP-IJCNLP)*, pages 5883–5889, Hong Kong,
736 China. Association for Computational Linguistics.
- 737 Ido Dagan, Lillian Lee, and Fernando C. N. Pereira.
738 1999. [Similarity-Based Models of Word Cooccur-](#)
739 [rence Probabilities](#). *Machine Learning*, 34(1):43–69.
- Wayne Davis. 2019. [Implicature](#). In Edward N. Zalta,
editor, *The Stanford Encyclopedia of Philosophy*, fall
2019 edition. Metaphysics Research Lab, Stanford
University.
- Oren Etzioni, Anthony Fader, Janara Christensen,
Stephen Soderland, and Mausam Mausam. 2011.
Open information extraction: the second generation.
In *Proceedings of the Twenty-Second international*
joint conference on Artificial Intelligence - Volume
Volume One, IJCAI’11, pages 3–10, Barcelona, Cat-
alonia, Spain. AAAI Press.
- Anthony Fader, Stephen Soderland, and Oren Etzioni.
2011. [Identifying Relations for Open Information Ex-](#)
[traction](#). In *Proceedings of the 2011 Conference on*
Empirical Methods in Natural Language Processing,
pages 1535–1545, Edinburgh, Scotland, UK. Associ-
ation for Computational Linguistics.
- Xingyu Fu, Weijia Shi, Xiaodong Yu, Zian Zhao, and
Dan Roth. 2020. [Design Challenges in Low-resource](#)
[Cross-lingual Entity Linking](#). In *Proceedings of the*
2020 Conference on Empirical Methods in Natural
Language Processing (EMNLP), pages 6418–6432,
Online. Association for Computational Linguistics.
- Juri Ganitkevitch, Benjamin Van Durme, and Chris
Callison-Burch. 2013. [PPDB: The Paraphrase](#)
[Database](#). In *Proceedings of the 2013 Conference*
of the North American Chapter of the Association
for Computational Linguistics: Human Language
Technologies, pages 758–764, Atlanta, Georgia. As-
sociation for Computational Linguistics.
- Maayan Geffet and Ido Dagan. 2005. [The Distribu-](#)
[tional Inclusion Hypotheses and Lexical Entailment](#).
In *Proceedings of the 43rd Annual Meeting of the*
Association for Computational Linguistics (ACL’05),
pages 107–114, Ann Arbor, Michigan. Association
for Computational Linguistics.
- Suchin Gururangan, Swabha Swayamdipta, Omer Levy,
Roy Schwartz, Samuel Bowman, and Noah A. Smith.
2018. [Annotation Artifacts in Natural Language In-](#)
[ference 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,
New Orleans, Louisiana. Association for Computa-
tional Linguistics.
- Aurélie Herbelot and Mohan Ganesalingam. 2013. [Mea-](#)
[suring semantic content in distributional vectors](#). In
Proceedings of the 51st Annual Meeting of the Associ-
ation for Computational Linguistics (Volume 2: Short
Papers), pages 440–445, Sofia, Bulgaria. Association
for Computational Linguistics.
- Xavier Holt. 2019. [Probabilistic Models of Relational](#)
[Implication](#). *arXiv:1907.12048 [cs, stat]*. ArXiv:
1907.12048.
- Mohammad Javad Hosseini, Nathanael Chambers, Siva
Reddy, Xavier R. Holt, Shay B. Cohen, Mark John-

796	son, and Mark Steedman. 2018. Learning Typed Entailment Graphs with Global Soft Constraints . <i>Transactions of the Association for Computational Linguistics</i> , 6:703–717.	
797		
798		
799		
800	Mohammad Javad Hosseini, Shay B. Cohen, Mark Johnson, and Mark Steedman. 2019. Duality of Link Prediction and Entailment Graph Induction . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4736–4746, Florence, Italy. Association for Computational Linguistics.	
801		
802		
803		
804		
805		
806		
807	Mohammad Javad Hosseini, Shay B. Cohen, Mark Johnson, and Mark Steedman. 2021. Open-Domain Contextual Link Prediction and its Complementarity with Entailment Graphs . In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 2790–2802, Punta Cana, Dominican Republic. Association for Computational Linguistics.	
808		
809		
810		
811		
812		
813		
814	Shengbin Jia, Shijia E, Maozhen Li, and Yang Xiang. 2018. Chinese Open Relation Extraction and Knowledge Base Establishment . <i>ACM Transactions on Asian and Low-Resource Language Information Processing</i> , 17(3):1–22.	
815		
816		
817		
818		
819	Dimitri Kartsaklis and Mehrmoosh Sadrzadeh. 2016. Distributional Inclusion Hypothesis for Tensor-based Composition . In <i>Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers</i> , pages 2849–2860, Osaka, Japan. The COLING 2016 Organizing Committee.	
820		
821		
822		
823		
824		
825	Keshav Kolluru, Vaibhav Adlakha, Samarth Aggarwal, Mausam, and Soumen Chakrabarti. 2020. OpenIE6: Iterative Grid Labeling and Coordination Analysis for Open Information Extraction . <i>arXiv:2010.03147 [cs]</i> . ArXiv: 2010.03147.	
826		
827		
828		
829		
830	Chin Lee, Hongliang Dai, Yangqiu Song, and Xin Li. 2020. A Chinese Corpus for Fine-grained Entity Typing . In <i>Proceedings of the 12th Language Resources and Evaluation Conference</i> , pages 4451–4457, Marseille, France. European Language Resources Association.	
831		
832		
833		
834		
835		
836	Omer Levy and Ido Dagan. 2016. Annotating Relation Inference in Context via Question Answering . In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 249–255, Berlin, Germany. Association for Computational Linguistics.	
837		
838		
839		
840		
841		
842	Dekang Lin. 1998. Automatic Retrieval and Clustering of Similar Words . In <i>36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics, Volume 2</i> , pages 768–774, Montreal, Quebec, Canada. Association for Computational Linguistics.	
843		
844		
845		
846		
847		
848	Xiao Ling and Daniel S. Weld. 2012. Fine-grained entity recognition. In <i>Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence, AAAI’12</i> , pages 94–100, Toronto, Ontario, Canada. AAAI Press.	
849		
850		
851		
852		
	Christopher Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven Bethard, and David McClosky. 2014. The Stanford CoreNLP Natural Language Processing Toolkit . In <i>Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations</i> , pages 55–60, Baltimore, Maryland. Association for Computational Linguistics.	853 854 855 856 857 858 859 860
	Nick McKenna, Liane Guillou, Mohammad Javad Hosseini, Sander Bijl de Vroe, Mark Johnson, and Mark Steedman. 2021. Multivalent Entailment Graphs for Question Answering . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 10758–10768, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	861 862 863 864 865 866 867 868
	Xiaoman Pan, Thammie Gowda, Heng Ji, Jonathan May, and Scott Miller. 2019. Cross-lingual Joint Entity and Word Embedding to Improve Entity Linking and Parallel Sentence Mining . In <i>Proceedings of the 2nd Workshop on Deep Learning Approaches for Low-Resource NLP (DeepLo 2019)</i> , pages 56–66, Hong Kong, China. Association for Computational Linguistics.	869 870 871 872 873 874 875 876
	Hoifung Poon and Pedro Domingos. 2009. Unsupervised Semantic Parsing . In <i>Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing</i> , pages 1–10, Singapore. Association for Computational Linguistics.	877 878 879 880 881
	Likun Qiu and Yue Zhang. 2014. ZORE: A Syntax-based System for Chinese Open Relation Extraction . In <i>Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1870–1880, Doha, Qatar. Association for Computational Linguistics.	882 883 884 885 886 887
	Siva Reddy, Mirella Lapata, and Mark Steedman. 2014. Large-scale Semantic Parsing without Question-Answer Pairs . <i>Transactions of the Association for Computational Linguistics</i> , 2:377–392.	888 889 890 891
	Martin Schmitt and Hinrich Schütze. 2021. Language Models for Lexical Inference in Context . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 1267–1280, Online. Association for Computational Linguistics.	892 893 894 895 896 897
	Sonse Shimaoka, Pontus Stenetorp, Kentaro Inui, and Sebastian Riedel. 2017. Neural Architectures for Fine-grained Entity Type Classification . In <i>Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers</i> , pages 1271–1280, Valencia, Spain. Association for Computational Linguistics.	898 899 900 901 902 903 904
	Gabriel Stanovsky, Julian Michael, Luke Zettlemoyer, and Ido Dagan. 2018. Supervised Open Information Extraction . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 885–895,	905 906 907 908 909 910

911	New Orleans, Louisiana. Association for Computational Linguistics.	Chujie Zheng, Minlie Huang, and Aixin Sun. 2019. ChID: A Large-scale Chinese IDiom Dataset for Cloze Test . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 778–787, Florence, Italy. Association for Computational Linguistics.	967
912			968
913	Kai Sun, Dian Yu, Dong Yu, and Claire Cardie. 2019. Investigating Prior Knowledge for Challenging Chinese Machine Reading Comprehension . <i>arXiv:1904.09679 [cs]</i> . ArXiv: 1904.09679.		969
914			970
915			971
916			972
917	Idan Szpektor and Ido Dagan. 2008. Learning Entailment Rules for Unary Templates . In <i>Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008)</i> , pages 849–856, Manchester, UK. Coling 2008 Organizing Committee.		
918			
919			
920			
921			
922	Daniel Tse and James R. Curran. 2012. The Challenges of Parsing Chinese with Combinatory Categorical Grammar . In <i>Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 295–304, Montréal, Canada. Association for Computational Linguistics.		
923			
924			
925			
926			
927			
928			
929	Shan Wang and Francis Bond. 2013. Building the Chinese Open Wordnet (COW): Starting from Core Synsets . In <i>Proceedings of the 11th Workshop on Asian Language Resources</i> , pages 10–18, Nagoya, Japan. Asian Federation of Natural Language Processing.		
930			
931			
932			
933			
934			
935	Sabine Weber and Mark Steedman. 2019. Construction and Alignment of Multilingual Entailment Graphs for Semantic Inference . pages 77–79.		
936			
937			
938	Julie Weeds and David Weir. 2003. A General Framework for Distributional Similarity . In <i>Proceedings of the 2003 Conference on Empirical Methods in Natural Language Processing</i> , pages 81–88.		
939			
940			
941			
942	Liang Xu, Hai Hu, Xuanwei Zhang, Lu Li, Chenjie Cao, Yudong Li, Yechen Xu, Kai Sun, Dian Yu, Cong Yu, Yin Tian, Qianqian Dong, Weitang Liu, Bo Shi, Yiming Cui, Junyi Li, Jun Zeng, Rongzhao Wang, Weijian Xie, Yanting Li, Yina Patterson, Zuoyu Tian, Yiwen Zhang, He Zhou, Shaowei Hua Liu, Zhe Zhao, Qipeng Zhao, Cong Yue, Xinrui Zhang, Zhengliang Yang, Kyle Richardson, and Zhenzhong Lan. 2020. CLUE: A Chinese Language Understanding Evaluation Benchmark . In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 4762–4772, Barcelona, Spain (Online). International Committee on Computational Linguistics.		
943			
944			
945			
946			
947			
948			
949			
950			
951			
952			
953			
954			
955	Dani Yogatama, Daniel Gillick, and Nevena Lazic. 2015. Embedding Methods for Fine Grained Entity Type Classification . In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)</i> , pages 291–296, Beijing, China. Association for Computational Linguistics.		
956			
957			
958			
959			
960			
961			
962			
963	Shuai Zhang, Lijie Wang, Ke Sun, and Xinyan Xiao. 2020. A Practical Chinese Dependency Parser Based on A Large-scale Dataset . <i>arXiv:2009.00901 [cs]</i> . ArXiv: 2009.00901.		
964			
965			
966			

A A Brief Summary of Jia et al. (2018)

In Table 4 are the 7 rules from Jia et al. (2018) which they call Dependency Structure Normal Forms. The first rule corresponds to nominal compounds which we elaborated in constructions **D** in §3.1; the second rule corresponds to direct S-V-O relations; the third rule attends to the semantic objects hidden in adjuncts, which are always preverbs in Chinese; the fourth rule subsumes complements of head verbs into the predicate; the fifth rule handles the coordination of subjects, the sixth handles coordination of object, and the seventh handles coordination of predicates. These rules are reflected in our ORE method as well, but for the sake of brevity, only the constructions which have never been covered by previous work are listed in §3.1.

德国 总理 默克尔 。 German Chancellor Merkel . (German, Chancellor, Merkel)
我 看到 你 。 I see you . (I, see, you)
他 在 家 玩 游 戏 。 He at home play game . (He, play-game, home)
我 走 到 图 书 馆 。 I walk to library . (I, walk-to, library)
我 和 你 去 商 店 。 I and you go-to shop . (I, go-to, shop) (you, go-to, shop)
我 吃 汉 堡 和 薯 条 。 I eat burger and chips . (I, eat, burger) (I, eat, chips)
罪 犯 击 中 、 杀 死 了 他 。 Criminal shot, kill -ed him . (criminal, shot, him) (criminal, kill, him)

Table 4: Set of DSNFs from Jia et al. (2018) exemplified. In each box, at top is an example sentence, presented in Chinese and its English metaphrase (inflection ignored); below are the relations they extract.

B Detailed Statistics of the CFIGER dataset

To test our assumption that each ultra-fine-grained type can be unambiguously mapped to a single FIGER type, we inspect the number of FIGER type labels to which each ultra-fine-grained type

is mapped through manual labelling. Among the 6273 ultra-fine-grained types in total, 5622 of them are mapped to exactly one FIGER type, another 510 are not mapped to any FIGER types; only 134 ultra-fine-grained types are mapped to 2 FIGER types, and 7 mapped to 3 FIGER types. No ultra-fine-grained types are mapped to more than 3 FIGER types. Therefore, it is safe to say that our no-ambiguity assumption roughly holds.

We further inspected the number of FIGER types each mention is attached with. It turns out that among the 1,913,197 mentions in total, 59,517 of them are mapped to no FIGER types, 1,675,089 of them are mapped to 1 FIGER type, 160,097 are mapped to 2 FIGER types, 16,309 are mapped to 3 FIGER types, 1,952 are mapped to 4 FIGER types, 200 are mapped to 5 FIGER types, and 33 are mapped to 6 FIGER types. No mentions are mapped to more than 6 FIGER types. Note that each mention can be mapped to more than one ultra-fine-grained types from the start, so these numbers are not in contradiction with the above numbers.

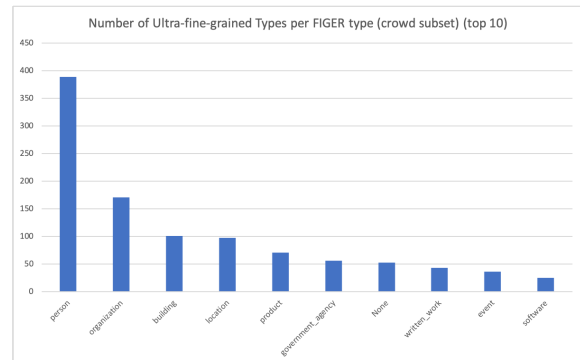


Figure 3: Number of ultra-fine-grained types in crowd-annotated subset mapped to each FIGER type; only the FIGER types with top 10 number of ultra-fine-grained types are displayed.

We also looked at the number of ultra-fine-grained types each FIGER type is mapped to, so as to understand the skewness of our mapping. Results are shown in Figure 3 and 4. Unsurprisingly, the most popular ultra-fine-grained labels are highly correlated with the ones that tend to appear in coarse-grained type sets, with “PERSON” label taking up a large portion. This distribution is largely consistent between crowd-annotated and Wikipedia subsets.

Another set of stats are the number of mentions that corresponds to each FIGER type, shown in Figure 5 and 6. The winners in terms of the number of mentions are consistent with that of the number of

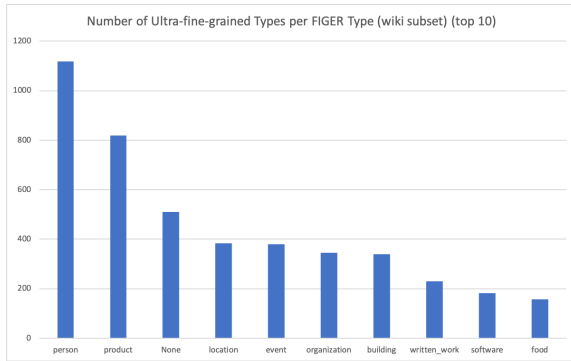


Figure 4: Number of ultra-fine-grained types in wikipedia distantly supervised subset mapped to each FIGER type; only the FIGER types with top 10 number of ultra-fine-grained types are displayed.

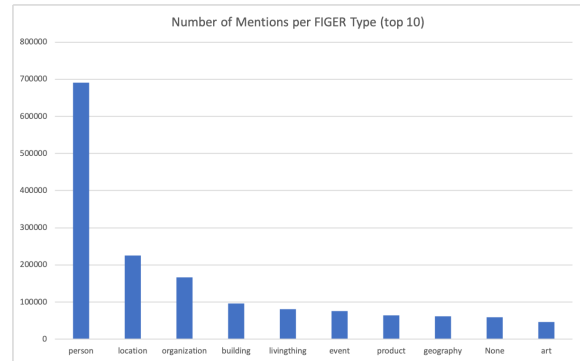


Figure 6: Number of mentions in wikipedia distantly supervised subset labelled as each FIGER type; only the FIGER types with top 10 number of mentions are displayed.

1032 ultra-fine-grained types, and also consistent among
1033 themselves (between the two subsets).

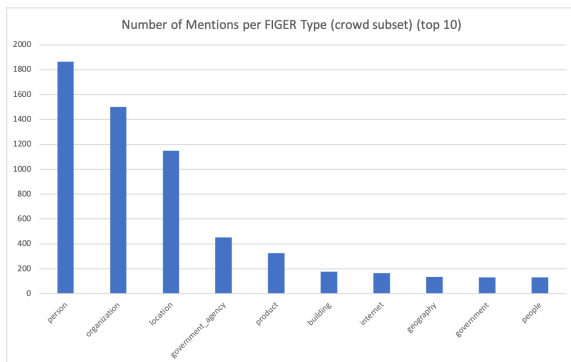


Figure 5: Number of mentions in crowd-annotated subset labelled as each FIGER type; only the FIGER types with top 10 number of mentions are displayed.

1034 C Selecting Relation Triples for 1035 Translated Levy-Holt

1036 To retrieve the relation triple most likely reflecting
1037 the meaning of the whole sentence, we follow this
1038 order when determining which relation triple to
1039 select:

- 1040 • For the amended relations, if the predicate of
1041 any of them cover the word with HEAD token
1042 in DDParse dependency parse, we randomly
1043 choose one of these;
- 1044 • If none is found, but the predicate of any non-
1045 amended relations cover the word with HEAD
1046 token in DDParse dependency parse, we ran-
1047 domly choose one of these;
- 1048 • If none is found, but there are any other rela-
1049 tions, we randomly choose one of these;

- Finally, if still none is found, we assign
PREMISE_PLACEHOLDER to the premise and
HYPOTHESIS_PLACEHOLDER to the hypothe-
sis, so that no entailment relation would ever
be detected between them.

D Implementation Details for Entailment Graph Construction

D.1 Corpus and Preprocessing

The original Webhose news corpus that we use con-
sists of 316K news articles. We cut the articles into
sentences by punctuations, limiting the maximum
sentence length to 500 characters (the maximum
sequence length for Chinese Bert). We discard the
sentences shorter than 5 characters, and the articles
whose sentences are all shorter than 5 characters.
After applying the filter, we are left with 313,718
articles, as shown in Table 2.

In the process of entity typing, following previ-
ous work, we consider only the first-layer FIGER
types; when multiple type labels are outputted, we
consider all combinations as valid types for that
predicate.

We have also considered using another, larger
corpus for building the Chinese entailment graphs,
but couldn't finish due to limits on computational
resources. We have referred to the larger corpus
as the CLUE corpus in §7: the larger corpus is
developed by Xu et al. (2020). It is eight times
the size of the Webhose corpus and originally in-
tended for training Chinese language models. We
provide the typed relation triples extracted from the
CLUE corpus as a part of our release, and encour-
age interested readers to build their own Chinese
entailment graph on this larger corpus, as we ex-
pect it to exhibit stronger performance, and present

an interesting comparison to the language-model driven models pre-trained with the same corpus.

D.2 Entailment Graph Construction

We have used the same entailment graph construction algorithm as Hosseini et al. (2018) to build our Chinese entailment graph from the pool of typed relation triples. When building our entailment graphs, we only feed in the relation triples whose predicate and arguments both appear at least 2 times¹⁵. Their approach of building entailment graphs comes in two steps, in the paragraphs below we will briefly summarize each step and discuss our implementation details.

The first step is local learning. In this step, instances of relation triples are grouped into clusters based on the arguments they take. Relations (predicates) that are seen with the same arguments of the same types are considered to have co-occurred. For each pair of predicates, based on the co-occurrence information, a few different entailment scores have been proposed, of which the BInc score (Szpektor and Dagan, 2008) was found to have the best empirical performance in (Hosseini et al., 2018). Following them, we also use the BInc score in the local learning step of our Chinese entailment graphs. Note that after the local learning step, the entailment scores between each pair of predicates are independently calculated, and there are no interactions between entailment subgraphs of different type pairs, thus the name **local** learning.

The second step is global learning. In this step, global transitivity constraint is “softly” applied to the local graphs as an optimization problem: paraphrase predicates are encouraged to have the same pattern of entailment; different typed subgraphs are encouraged to have the same entailment score for the same (ignoring type) pair of predicates; finally, the global scores are encouraged to stay similar to the local scores as a measure of regularization. In *Jia* baseline, the local graphs are too weak for global learning to be helpful; in *DDPORE* baseline, the best dev set AUC (as reported in Table 3) is achieved after 2 epochs; in EG_{Zh} , the best dev set AUC is achieved after 3 epochs.

E Case Study for Entailment Detection

In order to further verify the source of our improvements, we analyse our ensemble with a case

¹⁵We experimented with 2-2, 2-3, 3-2 and 3-3, among which this 2-2 setting is empirically favoured.

study: we compare the predictions of our Ensemble_AVG to that of the English monolingual EG_{En} , both thresholded over 65% precision. We categorize the prediction differences into 4 classes: *True Positives*, *False Positives*, *True Negatives*, *False Negatives*. *Positives* are cases where the ensemble switched the prediction label from negative to positive, vice versa for *negatives*; *True* means that the switch is correct, *False*, that the switch is incorrect.

In Table 5, we break down each class of differences according to the direct cause of EG_{Zh} making a different prediction than EG_{En} ¹⁶¹⁷:

- **same sentence after translation:** The premise and hypothesis become identical in relation structure; this can only happen with *positives*;
- **translation error:** The premise or hypothesis becomes unparsable into relations due to translation error; this can only happen with *negatives*;
- **lexicalization:** The difference in predictions is attributed to the cross-lingual difference in the lexicalization of complex relations;
- **ORE error:** After translation, the true relations in premise and hypothesis have the same arguments, but are mistaken due to ORE error;
- **evidence of entailment:** The difference is attributed to the different evidence of entailment in the two graphs; this is most relevant to our EG_{Zh} .

As shown, the majority of our performance gain comes from the additional evidence of entailment in EG_{Zh} ; surprisingly, translation played a positive role in the ensemble, though not a major contributor. We attribute this to the fact that MT systems tend to translate semantically similar sentences to the same target sentence, though this similarity is still symmetric, not directional. We have singled out this effect in the “BackTrans Esb” ablation study in §6.3, and have confirmed that this effect is marginal to our success.

In Table 5, for both the differences from evidence of entailment, and differences in TOTAL, the precision of *positives* is lower than that of *negatives*. Namely, $TP/(TP + FP)$ is lower than $TN/(TN + FN)$. This is no surprise, as *positives* and *negatives* have different baselines to start with: *Positives* attempt to correct the false negatives from EG_{En} , where 17% of all negatives are false; *Negatives* attempt to correct the false positives, where 35% of all positives are false (as dictated in the

¹⁶since the switch in the ensemble is driven by EG_{Zh} .

¹⁷examples of each class of cause are given in Appendix F.

Direct causes of EG _{Zh} 's different prediction	TP (+)	FP (-)	TN (+)	FN (-)	+/-
translation-related causes, among which:	+52	-28	+42	-47	+19
· <i>same sentence after translation</i>	+52	-28	0	0	+24
· <i>translation error</i>	0	0	+42	-47	-5
lexicalization	+29	-54	+16	-12	-21
ORE error	+8	-20	+8	-5	-9
evidence of entailment	+109	-95	+86	-40	+60
TOTAL	+198	-197	+152	-104	+49

Table 5: Breakdown of the different predictions between our ensembles and English monolingual graph. “TP”, “FP”, “TN”, “FN” represent *True Positive*, *False Positive*, *True Negative* and *False Negative* respectively; in the column “+/-” is the overall impact of each factor.

setting of our case study). In this context, it is expectable that our evidence of entailment gets $109/(109 + 95) = 53\%$ correct for *positives*, while a much better $86/(86 + 40) = 68\%$ correct for *negatives*. These results support the solidarity of our contributions.

F Examples of Different Predictions in Case Study by Category of Direct Cause

In this section, we provide one example for each class of direct cause, as described in the above Appendix E. Chinese sentences and relations in the examples are presented in the same format as §3.1.

Same sentence after translation

- Premise - English: (magnesium sulfate, relieves, headache)
- Hypothesis - English: (magnesium sulfate, alleviates, headaches)
- Premise - Chinese translation: “硫酸镁(magnesium) 缓解(relieves) 头痛(headache)”
- Hypothesis - Chinese translation: “硫酸镁(magnesium) 缓解(alleviates) 头痛(headache)”

The two sentences are translated to the same surface form in Chinese, as the predicates are in many cases synonyms. There are more true positives than false positives, because synonyms are simultaneously more likely true entailments and more likely translated to the same Chinese word.

Translation Error

- Premise - English: (Refuge, was attacked by, terrorists)
- Hypothesis - English: (Terrorists, take, refuge)
- Premise - Chinese translation: “避难所(refuge) 遭到(suffered) 恐怖分子(terrorists) 袭

击(attack); Refuge suffered attack from terrorists.”

- Hypothesis - Chinese translation: “恐怖分子(terrorists) 避难(take-shelter); Terrorists take shelter.”

The hypothesis is supposed to mean “The terrorists took over the refuge”. However, with translation, the hypothesis in Chinese is mistaken as an intransitive relation where take-refuge is considered a predicate.

Lexicalization

- Premise - English: (Granada, is located near, mountains)
- Hypothesis - English: (Granada, lies at the foot of, mountains)
- Premise - Chinese translation: “格拉纳达(Granada) 靠近(is-near) 山脉(mountains)”
- Hypothesis - Chinese translation: “格拉纳达(Granada) 位于(is-located-at) 山脚下(hillfoot)”

When the hypothesis is translated into Chinese, the lexicalization of the relation changed, the part of the predicate hosting the meaning of ‘the foot of’ is absorbed into the object. Therefore, while in English “is located near” does not entail “lies at the foot of”, in Chinese “is-near” is considered to entail “is-located-at”. In this way, an instance of *false positive* comes into being.

ORE Error

- Premise - English: (A crow, can eat, a fish)
- Hypothesis - English: (A crow, feeds on, fish)
- Premise - Chinese translation: “乌鸦(crow) 可以(can) 吃(eat) 鱼(fish)”
- Hypothesis - Chinese translation: “乌鸦(crow) 以(take) 鱼(fish) 为(as) 食(food)”

1250 • Premise - extracted Chinese relation: (crow, eat,
1251 fish)

1252 • Hypothesis - extracted Chinese relation: (crow,
1253 take-X-as-food, fish)

1254 While the translations for this pair of relations
1255 is correct, in the subsequent Chinese open relation
1256 extraction, our ORE method failed to recognize “可
1257 以(can)” as an important part of the predicate. To
1258 avoid sparsity, most adjuncts of the head verb are
1259 discarded, and modals are part of them. While the
1260 original premise “can eat” does not entail “feeds
1261 on”, the Chinese premise “eat” does in a way entail
1262 “feeds on”, where another instance of *false positive*
1263 arises.

1264 Evidence of Entailment

- 1265 • Premise - English: (quinine, cures, malaria)
- 1266 • Hypothesis - English: (quinine, is used for the
1267 treatment of, malaria)
- 1268 • Premise - Chinese translation: “奎宁(quinine) 治
1269 疗(cure) 疟疾(malaria)”
- 1270 • Hypothesis - Chinese translation: “奎宁(quinine)
1271 用于(is-used-to) 治疗(cure) 疟疾(malaria)”
- 1272 • Premise - extracted Chinese relation: (quinine,
1273 cure, malaria)
- 1274 • Hypothesis - extracted Chinese relation: (quinine,
1275 is-used-to-cure, malaria)

1276 In the above example, sufficiently strong evi-
1277 dence for “cure” entailing “is used for the treat-
1278 ment of” is not found in the English graph, whereas
1279 strong evidence for “治疗(cure)” entailing “用
1280 于·治疗(is-used-to-cure)” is found in the Chinese
1281 graph. In this way we get an instance of *true posi-*
1282 *tive*.

1283 G More Precision-Recall Curves

1284 In this section, we present more precision-recall
1285 curves from the baselines and ablation studies in
1286 Table 3. These curves contain more details explain-
1287 ing the AUC values in the table.

1288 Figure 7 contains the curves for the ablation
1289 study of DataConcat. Here all three models ulti-
1290 mately come from the same corpus, so the per-
1291 formance difference can be fully attributed to
1292 the cross-lingual complementarity of entailment
1293 graphs.

1294 Figure 8 contains the curves for two sets of ab-
1295 lation studies: EG_{Zh} with or without entity typing;
1296 EG_{En} ensembled with back-translation predictions

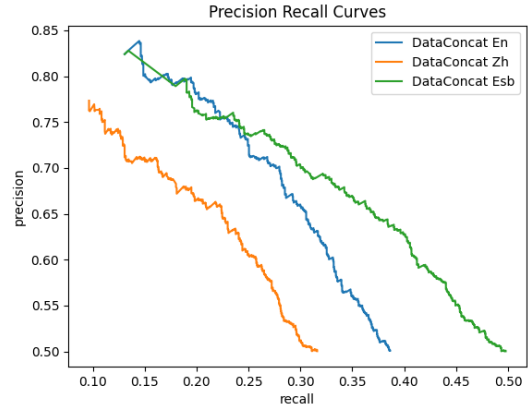


Figure 7: P-R Curves on Levy-Holt test set for Data-
Concat ablation study.

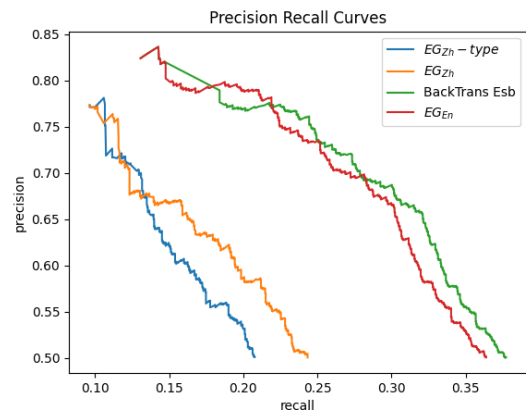


Figure 8: P-R Curves on Levy-Holt test set for EG_{Zh}
-type, BackTrans Esb, in comparison to EG_{Zh} and
 EG_{En} respectively.

1297 or not. The former study shows the clear benefit
1298 of our entity typing system, while the latter study
1299 shows that ensembling with back-translated pre-
1300 dictions only results in a marginal gain, therefore
1301 the synonym effect from translation is not a ma-
1302 jor contributor to the success of our ensembling
1303 method.

1304 H Implementation Details for Our 1305 Question Answering Evaluation

1306 Our Chinese boolean QA dataset is constructed
1307 from the CLUE Chinese news corpus (Xu et al.,
1308 2020), a huge Chinese news corpus of 2.4M news
1309 articles. The CLUE corpus has 8 times the number
1310 of articles as the Webhose corpus¹⁸.

1311 We partition the corpus into 122 disjoint 3-day
1312 time spans. We look for frequent predicates be-
1313 tween frequent typed-argument-pairs in each parti-
1314 tion. Since we want the typed-argument-pairs that

¹⁸the one from which we built our Chinese entailment
graphs.

are reported by multiple sources, we look for typed-argument-pairs that appear in at least 15 articles and with at least 15 predicates within a partition; for predicates, we just want to make sure they are felicitous, so we look for those predicates that appear at least 30 times in an arbitrary number of articles anywhere in the corpus. The motivation for these thresholds is the following trade-off: lower thresholds lead to noisier datasets; higher thresholds lead to more biased datasets.

Relation triples satisfying the above criteria are reformatted into textual propositions, and selected as “positives”: these predicates are frequently-mentioned, so they are felicitous relations; these argument-pairs are mentioned in many articles within the time-span, so the “positive” predicates should be inferrable from the other mentions of these argument-pairs in the time-span. In order to balance the dataset, at most one positive is chosen from each sentence.

One naive approach to generating negatives, is to substitute random predicates into the positive propositions. However, that is not adversarial enough as a test for directional inference: unrelated words can be easily detected by symmetric similarity measures like Bert similarities, without involving any wisdom in directionality.

Following (McKenna et al., 2021), we replace the positive predicates with their hyponyms / troponyms in Chinese WordNet (Wang and Bond, 2013). These replacements are semantically related to the positives, but do not logically follow the positives. We select those replacements that are absent with the argument pair of its corresponding positive in the corresponding partition. We also require that the replacements appear elsewhere in the corpus at least 5 times. Based on the Gricean cooperative principle of communication (Davis, 2019), we assume that the collection of news articles would report all and only the facts that are known. It is then implied, that these selected replacements are felicitous predicates, but untrue (or, not confirmed to be true) in the contexts of their corresponding positives. Thus, they can be used as adversarial negatives. As discussed in §7, we look for substitution of spans rather than entire predicates, to deal with the multi-token and discontinuous feature of Chinese predicates; to balance the dataset, at most two negatives are chosen from each sentence; in order for quality control, only those positives from which some negatives can be generated are kept.

For instance, for a positive proposition “约翰(John) 在(at) 乐购(Tesco) 购物(shop)”; John shopped at Tesco”, the predicate in this positive is “在-X-购物 (at-X-shop)”. By replacing random predicates, for example “起诉(sue)” into the positive proposition, we would get negatives like “约翰(John) 起诉(sue) 乐购(Tesco); John sued Tesco”, which is irrelevant to the positive and easy to guess for symmetric measures. On the other hand, using Chinese WordNet, for the subspan “购物(shop)” in the positive predicate, we can find a troponym “买日用品(go marketing)”, thus we can get negatives like “约翰(John) 在(at) 乐购(Tesco) 买日用品(go marketing); John went marketing at Tesco”. This replacement is much more semantically related than the random one, and, unless otherwise mentioned in the context, it can be assumed that **we don’t know John went marketing at Tesco**. Therefore, this latter replacement is still a negative, and a much more challenging one.

Notably, both for this evaluation and for the entailment detection evaluation in §6, we use the actual arguments in the sentences for BERT similarities, not the types of the arguments. This is because we empirically find that by replacing the actual arguments with their types, the language models get confused, and their performances drop.

As our QA task is a *machine reading at scale* task, each method uses the entire partition of news articles as context, and produces a score of whether the queried proposition is true according to the context. Method design choices can be categorized into two dimensions: how to retrieve the relevant context from the context pool, and how to calculate a truthfulness score for the query based on the retrieved relevant context.

Along the first dimension, our $Bert_{tfidf}$ baseline uses TF-IDF matching to retrieve the relevant articles, and use the sentences in these articles as relevant context; the other 4 methods in our experiment use exact match of argument-pairs to identify the “related relation triples”, among them, $BERT_{sent}$ baseline retrieves the host-sentences of these relation triples, while $Bert_{rel}$, EG_{Zh} and DDPORE use these relation triples themselves.

Along the second dimension, all methods take each context sentence or relation individually, and calculate the score as “whether the query proposition can be inferred from any context sentence/relation retrieved”. The three BERT baselines calculate the cosine similarity between the

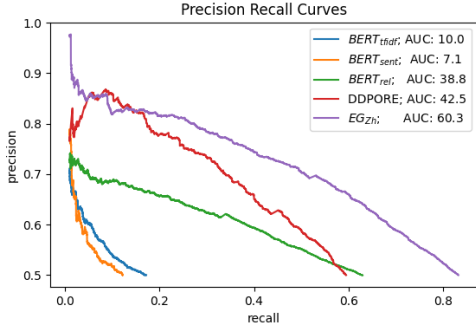


Figure 9: P-R Curves on QA evaluation **dev** set for EG_{Zh} and baselines; AUC values are annotated in the legend.

QA eval AUC (%)	dev	test
BERT _{tfidf}	10.0	12.8
BERT _{sent}	7.1	3.6
BERT _{rel}	38.8	40.5
DDPORE	42.5	41.7
EG _{Zh}	60.3	59.0

Table 6: Area Under Curve (AUC) on QA evaluation, for Chinese entailment graph (EG_{Zh}) and its baselines.

BERT representations of each context sentence and the query at the [CLS] token; the DDPORE baseline and our EG_{Zh} retrieve the entailment scores from each context triple to the query triple, from the corresponding typed entailment sub-graphs.

In addition to the test set Precision-Recall curves reported in §7, we present the dev set Precision-Recall curves as well, in Figure 9; we also summarize the AUC values in Table 6. The observations here are consistent with the conclusions in §7. It is to be noticed, that we did not do any hyperparameter tuning on the dev set, the best settings in Levy-Holt dev set are directly applied here. Nevertheless, we present this dev set along with the test set, to form a complete dataset on the task of Chinese boolean machine reading at scale, which we advocate as a solid benchmark for directional inference.

I Manual Examination of Chinese Levy-Holt

In order to provide a quantified evaluation for the quality of our Chinese Levy-Holt dataset from a human perspective, we manually labelled 100 proposition pairs in the Chinese Levy-Holt dev set (1-29, 1124-1136, 2031-2059, 3091-3122, 4061-4089, excluding the entries which are not parsed back into binary relation triples).

In this evaluation, we aim to answer the question

of “how accurate is the translate-then-parse procedure when it claims to have successfully converted an evaluation entry”. We label each Chinese entry along two dimensions: semantic consistency, whether it has preserved the meaning of the English entry; label consistency, whether the entailment label remains correct.

Along the first dimension of semantic consistency, we summarize our findings as follows:

- **Correct:** 74/100. These are the Chinese entries whose Chinese **predicates** precisely reflects the meaning of the English entry¹⁹;
- **Metaphors:** 3/100. These are the cases where the English entry involves metaphorical word-senses of predicates, but such metaphorical senses of these words are infelicitous in Chinese context;
- **Adjuncts:** 9/100. These are the cases where a part of an English predicate is translated into an adjunct to the Chinese head-verb, and is not included in the Chinese predicate (as in the example for ORE Errors in Appendix F); examples of missed-out adjuncts are ‘widely’, ‘should’ and ‘may’;
- **Lexical:** 5/100. These are the cases where the word-segmentation of the Chinese sentence is incorrect (as Chinese sentences come with no spaces between words);
- **Errors:** 7/100. These are the cases where, although the Chinese ORE method outputs some binary relation triples for the translation, that relation triple is not the true relation for the sentence;
- **Translation:** 2/100. These are the cases where, although the translation can be parsed into some binary relation triples by our Chinese ORE method, the translation is incorrect, thus everything downstream is wrong.

Along the second dimension of label consistency, we find that: in 89 / 100 entries, the actual labels in Chinese are consistent with the English labels; in 10 / 100 entries, the actual labels in Chinese are inconsistent with the English labels; in the remaining 1 / 100 entry, the actual label in Chinese is consistent with the actual label in English, but the provided English label is corrupted.

In summary, for the portion where the conversion is successful, the entries in Chinese Levy-Holt

¹⁹Arguments are allowed to be translated to different senses of the words, as long as the entailment label between the predicates is not affected.

preserves the meaning of the English entries reasonably well; more importantly, the labels of the Chinese Levy-Holt dataset remains robust.

J Diagram Illustrations of Our Syntactic Analysis

In this section, we present for interested readers a set of diagram illustrations of the set of constructions, as involved in our syntactic analysis in §3.1. For each construction, we draw a diagram to illustrate its dependency structure, an example to instantiate the dependency structure, and in the following lines, all the relations that we extract from this construction (one relation per line). Each relation comes in the form of triple-of-types (consistent with the diagram) and triple-of-words (as in the example), separated by semi-colons. The diagrams are presented in Table 8, Table 9 and Table 10.

K Ethics Considerations

Below we discuss the ethics considerations in our work.

The limitation to our work is two-fold. Firstly, our Chinese entailment graphs focus on the task of predicate entailment detection, and does not attempt to independently solve the more general problem of reasoning and inference: this more general task would also involve other resources including argument hypernymy detection, quantifier identification and co-reference resolution. These are out of the scope of this work. Secondly, while we have shown the effect of cross-lingual complementarity, adding in more languages to the ensemble is not directly straight-forward: this would require linguistic expertise and NLP infrastructure in the respective languages; including more languages, and eventually including arbitrary languages, is one of the directions for our future work.

The risk of our work mostly stems from our use of large-scale news corpora: if the media coverage itself is biased toward certain aspects of the world or certain groups of people, then these biases would be inherited by our entailment graphs. Our response to this is to include as many diverse news sources as possible to reduce such biases to the minimum: our source corpus for building Chinese entailment graphs includes 133 different news sources from a variety of countries and regions.

For the computational cost of building Chinese entailment graphs, the algorithm for open relation extraction takes roughly 140 CPU hours to process

Stats	Webhose	Levy-Holt
AVG sentence length (in # of Chinese characters)	24.9	10.1
AVG # of relations per sentence	15.6	2.72
Percentage of relations from our additional patterns in §3.1	48%	32%

Table 7: Some key statistics of Webhose corpus and Chinese Levy-Holt dataset.

the entirety of Webhose corpus; the entity typing model takes roughly 180 GPU hours on NVidia 1080Ti GPUs to do inference on the entirety of Webhose corpus; the local learning process takes less than one hour, and, the global learning process, our major computational bottleneck, takes roughly 800 CPU hours to finish.

The major datasets of use, namely, Webhose corpus, CLUE dataset and the CFET dataset, are open corpora with no specified licenses, thus our academic use is allowed; no license was specified for the Levy-Holt dataset as well; our own CFIGER dataset as well as the constructed entailment graphs can be distributed under the MIT license.

L Comparison Between Webhose Corpus and Levy-Holt Dataset

In this section, we report some key statistics of the Webhose corpus in comparison to the Levy-Holt dataset, which highlight their difference in genre.

As shown in Table 7, the Webhose corpus has much longer sentences than the Chinese Levy-Holt dataset, and on average, a much larger number of open relations can be extracted from the sentences in Webhose corpus. More importantly, the relation patterns which we additionally identified in §3.1 are much better represented (constituting 48% of all relations) than in Chinese Levy-Holt (32%). Thus, it is clear that: 1) our ORE method in §3 was not tuned on the test data, namely Chinese Levy-Holt; 2) tuning on Chinese Levy-Holt would not help with building better ORE methods for news corpora. On the other hand, as a large-scale multi-source news corpus of 5 million sentences, Webhose corpus can be believed to accurately reflect the distribution of linguistic patterns in the entirety of the news genre.

Construction ID	Diagrams and Examples
A.1	
	<p>Example: “咽炎(<i>pharyngitis</i>) 成为(<i>becomes</i>) 发热(<i>fever</i>) 的(<i>De</i>) 原因(<i>cause</i>); <i>Pharyngitis becomes the cause of fever</i>”</p>
	<p>Relation 1: (Subj, Pred, Direct_Object); (咽炎(<i>pharyngitis</i>), 成为(<i>becomes</i>), 原因(<i>cause</i>))</p>
	<p>Relation 2: (Subj, Pred·X·DE·Direct_Object, True_Object); (咽炎(<i>pharyngitis</i>), 成为·X·的·原因(<i>becomes·X·DE·cause</i>), 发烧(<i>fever</i>))</p>
A.2	
	<p>Example: “苹果(<i>Apple</i>) 的(<i>De</i>) 创始人(<i>founder</i>) 是(<i>is</i>) 乔布斯(<i>Jobs</i>); <i>The founder of Apple is Jobs</i>”</p>
	<p>Relation 1: (Direct_Subject, Pred, Object); (创始人(<i>founder</i>), 是(<i>is</i>), 乔布斯(<i>Jobs</i>))</p>
	<p>Relation 2: (True_Subject, Direct_Subject·Pred, Object); (苹果(<i>Apple</i>), 创始人·是(<i>founder·is</i>), 乔布斯(<i>Jobs</i>))</p>
B.1	
	<p>Example: “我(<i>I</i>) 去(<i>go-to</i>) 诊所(<i>clinic</i>) 打(<i>take</i>) 疫苗(<i>vaccine</i>); <i>I go to the clinic to take the vaccine</i>”</p>
	<p>Relation 1: (Subject, Pred_1, Object_1); (我(<i>I</i>), 去(<i>go-to</i>), 诊所(<i>clinic</i>))</p>
	<p>Relation 2: (Subject, Pred_2, Object_2); (我(<i>I</i>), 打(<i>take</i>), 疫苗(<i>vaccine</i>))</p>

Table 8: The syntactic analysis in §3.1 illustrated with diagrams, examples and their extracted relations.

Construction ID	Diagrams and Examples
B.2	
	Example: “我(I) 想(want) 试图(try) 开始(begin) 写(write) 一个(a) 剧本(play); <i>I want to try to begin to write a play</i> ”
	Relation 1: (Subject, Pred_1, Pred_2); (我(I), 想(want-to), 试图(try))
	Relation 2: (Subject, Pred_1·Pred_2, Pred_3); (我(I), 想·试图(want-to-try), 开始(begin))
 Relation K: (Subject, Pred_1...·Pred_K, Object); (我(I), 想·试图·开始·写(want-to-try·begin·write), 一个剧本(A play))
C	
	Example: “他(he) 解决(solve) 了(-ed) 困扰(puzzle) 大家(everyone) 的(De) 问题(problem); <i>He solved the problem that puzzled everyone</i> ”
	Relation 1: (Subject, Pred_1, Object_1); (他(He), 解决(solved), 问题(problem)) Relation 2: (Object_1, Pred_2, Object_2); (问题(Problem), 困扰(puzzled), 大家(everyone))
D	Analysis in construction D removes the infelicitous instances of the Nominal Compound construction; for the illustration of this construction, we refer readers to Jia et al. (2018) and do not repeat here.

Table 9: More syntactic analysis in §3.1 illustrated with diagrams, examples and their extracted relations.

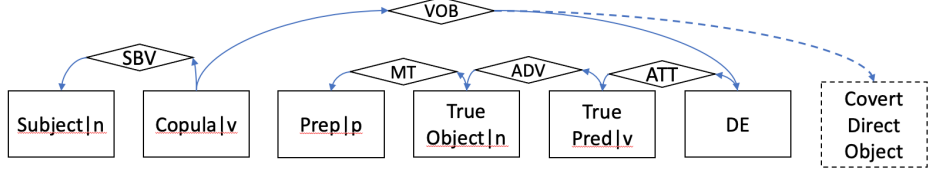
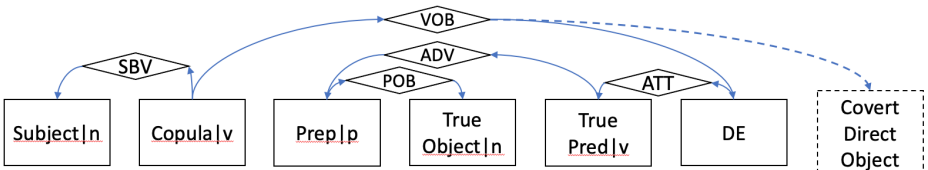
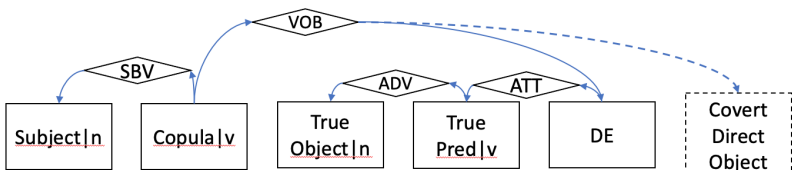
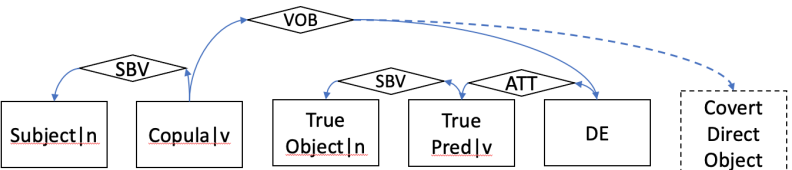
Construction ID	Diagrams and Examples
E.1	 <p>Example: “玉米(<i>Corn</i>) 是(<i>is</i>) 从(<i>from</i>) 美国(<i>US</i>) 引进(<i>introduce</i>) 的(<i>De</i>); <i>Corn is introduced from US</i>”</p> <p>Relation 1: (Subject, Copula-Prep-X-True_Pred-DE, True_Object); (玉米(<i>Corn</i>), 是·从·X·引进·的(<i>is-from-X-introduced-DE</i>), 美国(<i>US</i>))</p>
E.2	 <p>Example: “设备(<i>device</i>) 是(<i>is</i>) 用(<i>from</i>) 木头(<i>wood</i>) 做(<i>make</i>) 的(<i>De</i>); <i>The device is made of wood</i>”</p> <p>Relation 1: (Subject, Copula-Prep-X-True_Pred-DE, True_Object); (设备(<i>device</i>), 是·用·X·做·的(<i>is-from-X-made</i>), 木头(<i>wood</i>))</p>
E.3	 <p>Example: “设备(<i>device</i>) 是(<i>is</i>) 木头(<i>wood</i>) 做(<i>make</i>) 的(<i>De</i>); <i>The device is made of wood</i>”</p> <p>Relation 1: (Subject, Copula-X-True_Pred-DE, True_Object); (设备(<i>device</i>), 是·X·做·的(<i>is-X-made</i>), 木头(<i>wood</i>))</p>
E.4	 <p>Example: “设备(<i>device</i>) 是(<i>is</i>) 木匠(<i>carpenter</i>) 做(<i>make</i>) 的(<i>De</i>); <i>The device is made by a carpenter</i>”</p> <p>Relation 1: (Subject, Copula-X-True_Pred-DE, True_Object); (设备(<i>device</i>), 是·X·做·的(<i>is-X-made-DE</i>), 木匠(<i>carpenter</i>))</p>

Table 10: Yet more syntactic analysis in §3.1 illustrated with diagrams, examples and their extracted relations.