

Coarse-to-fine Knowledge Graph Domain Adaptation based on Distantly-supervised Iterative Training

Hongmin Cai
hmcai@scut.edu.cn
South China University of Technology
China

Wenxiong Liao
cswxliao@mail.scut.edu.cn
South China University of Technology
China

Zhengliang Liu
zl18864@uga.edu
University of Georgia
USA

Yiyang Zhang
zyyinyourarea@163.com
South China University of Technology
China

Xiaoke Huang
csxkhuang@mail.scut.edu.cn
South China University of Technology
China

Siqi Ding
d15995291636@163.com
South China University of Technology
China

Hui Ren
hren2@mgh.harvard.edu
Massachusetts General Hospital and
Harvard Medical School
USA

Zihao Wu
zihao.wu1@uga.edu
University of Georgia
USA

Haixing Dai
hd54134@uga.edu
University of Georgia
USA

Sheng Li
vga8uf@virginia.edu
University of Virginia
USA

Lingfei Wu
Teddy.lfwu@gmail.com
Pinterest
USA

Ninghao Liu
ninghao.liu@uga.edu
University of Georgia
USA

Quanzheng Li
Li.Quanzheng@mgh.harvard.edu
Massachusetts General Hospital and
Harvard Medical School
USA

Tianming Liu
tliu@cs.uga.edu
University of Georgia
USA

Xiang Li
XLI60@mgh.harvard.edu
Massachusetts General Hospital and
Harvard Medical School
USA

ABSTRACT

The knowledge graph (KG) is a highly needed basis to support the high-fidelity, high-interpretability modeling of various tasks in healthcare artificial intelligence. In this work, we focus on constructing an oncology knowledge graph that will be used in downstream cancer research and solution development. Modern supervised learning for knowledge graph construction requires a large amount of manually labeled data, which makes the process time-consuming and labor-intensive. Although there exists multiple research on named entity recognition and relation extraction based on distantly supervised learning, constructing a domain-specific knowledge graph from large collections of textual data without manual annotations is still an urgent problem to be solved. In response, we propose an integrated framework for adapting and re-learning knowledge graphs from a general domain (biomedical in our case) to a fine-defined domain (oncology). In this framework,

we apply distant-supervision on cross-domain knowledge graph adaptation. Consequently, no manual data annotation is required to train the model. We introduce a novel iterative training strategy to facilitate the discovery of domain-specific named entities and triplets. Experimental results indicate that the proposed framework can perform domain adaptation and construction of knowledge graphs efficiently.

KEYWORDS

Knowledge Graph Domain Adaptation, Knowledge Graph Construction, Named Entity Recognition, Relationship Extraction

ACM Reference Format:

Hongmin Cai, Wenxiong Liao, Zhengliang Liu, Yiyang Zhang, Xiaoke Huang, Siqi Ding, Hui Ren, Zihao Wu, Haixing Dai, Sheng Li, Lingfei Wu, Ninghao Liu, Quanzheng Li, Tianming Liu, and Xiang Li. 2018. Coarse-to-fine Knowledge Graph Domain Adaptation based on Distantly-supervised Iterative Training. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

1 INTRODUCTION

In healthcare, the development of robust and interpretable clinical decision support systems and the corresponding research requires both a substantial amount of data and effective modeling of the medical domain knowledge [32]. Knowledge graphs (KG) have been

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 ACM 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](https://permissions.acm.org).
Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06... \$15.00
<https://doi.org/XXXXXXXX.XXXXXXX>

explored in healthcare to represent the underlying relationships representing the domain knowledge, including Unified Medical Language System (UMLS) [15] and Google Healthcare Knowledge Graph [21]. In this study, we focus on developing the KG for oncology, an important branch of medicine that studies cancer treatment and prevention. Delineating the relationships between cancer sub-types, symptoms, comorbidities, genetic factors, and treatment with the oncology KG would provide a powerful basis for the downstream task in clinical decision support, such as patient diagnosis, prognosis, phenotyping, and treatment optimization.

However, existing approaches are not effective for constructing domain-specific KGs, especially for oncology, where limited access to oncological expertise hinders the supply of labeled training data. Insufficient labeled data typically leads to suboptimal performance. In fact, the dependence on sizeable training data significantly diminishes the real-world potential of data-driven KG construction methods based on supervised learning. In addition, although rule-based methods (based on resources such as Stanford CORE NLP) do not have stringent data demands, they typically suffer from suboptimal hand-crafted feature designs and the absence of helpful fine-grained connections to the domain data. Consequently, automatically constructing knowledge graphs directly from natural texts has attracted close attention from researchers in recent years [10, 20, 25].

In order to address these challenges, we investigate the coarse-to-fine learning for constructing an oncology knowledge graph that leverages knowledge from general biomedical KGs, especially the distantly-supervised interactive training to achieve knowledge graph domain adaptation. In distantly-supervised learning, fine-domain KGs are derived from the general-domain KGs. For example, a biomedical KG that covers broad concepts and common sense knowledge in the biomedical domain can serve as the base KG for a specialized oncology KG. Therefore, the KG in the coarse domain can be used as a knowledge base for distant supervision, thus avoiding the need for extensive manual annotations. However, only using the KG of the coarse domain as the knowledge base might limit the model's ability to discover domain-specific named entities and triples in the fine domain, further limiting the construction of the fine domain KG. Thus in this paper, we propose a novel coarse-to-fine knowledge graph domain adaptation (KGDA) framework. Our KGDA framework utilizes an iterative training strategy to enhance the model's ability to discover fine-domain entities and triples, thereby facilitating fast and effective coarse-to-fine KG domain adaptation.

Overall, the contributions of our work are as follows:

- An integrated framework for adapting and re-learning KG from coarse-domain to fine-domain is proposed. As a case study, the biomedical domain and oncology domain are considered the coarse domain and fine domain, respectively.
- Our model does not require human annotated samples with distant-supervision for cross-domain KG adaptation, and the iterative training strategy is applied to discovering domain-specific named entities and new triples.
- The proposed method can be adapted to various pre-trained language models (PLMs) and can be easily applied to different

coarse-to-fine KGDA tasks. It is so far the simplest data-driven approach for learning a KG from free text data, with the help of the coarse domain KG.

- Experimental results demonstrate the effectiveness of the proposed KGDA framework. We will release the source code and the data used in this work to fuel further research. The constructed oncology KG will be hosted as a web service to be used by the general public.

2 BACKGROUND AND RELATED WORK

Automatic KG construction from text generally involves two primitive steps: named entity recognition (NER) and relation extraction (RE). Named entity recognition aims to identify the types of entities mentioned in text sequences, such as people, places, etc. in the open domain; or diseases, medicine, disease symptoms, etc. in the biomedical domain. Relation extraction is also known as triplet extraction, which aims to identify the relationship between two entities, such as the birthplace relationship between people and places, or the therapeutic relationship between drugs and diseases in the biomedical domain. NER and RE are the necessary steps for information extraction to construct KG from text.

Distant supervision [24] is an intuitive way to transfer general-domain KG to fine domains. Distant-supervision provides labels for data with the help of an external knowledge base, which saves the time of manual labeling. For distantly-supervised NER, we can build distant labels by matching unlabeled sentences with external semantic dictionaries or knowledge bases. The matching strategies usually include string matching [34], regular expressions [5], and some heuristic rules. The distantly-supervised RE holds an assumption [18]: if two entities participate in a relation, then any sentence that contains those two entities might express that relation. Following this assumption, any sentence mentioning a pair of entities that have a relation according to the knowledge base will be labeled with this relation [24].

2.1 Pipeline-based methods for KG construction

The pipeline-based methods apply carefully-crafted linguistic and statistical patterns to extract the co-occurred noun phrases as triples. There are many off-the-shelf toolkits available, for example, Stanford CoreNLP [16], NLTK [26], and spaCy, which can be used for the NER tasks; Reverb [4], OLLIE [23], and Stanford OpenIE [2] can be used for the information extraction task. There have been multiple pipelines [17, 20] developed as well, consisting of modules targeting different functionalities needed for the KG construction. However, the pre-defined rules of off-the-shelf toolkits are generally tailored to specific domains, such methods are not domain-agnostic, and a new set of rules will be needed for a new domain.

2.2 Data-driven methods for KG construction

With the development of representation learning in language models, researchers began to apply data-driven models to solve the KG construction tasks. Based on how the model is trained, such work can be divided into three categories: fully-supervised methods [12, 34], semi-supervised methods [29], and weakly-supervised methods [28]. We will introduce the methods of fully-supervised

and weakly-supervised in this section. Specifically, the NER, RE, and entity linking tasks in the KG construction pipeline can all be solved by fully-supervised learning methods such as long short-term memory neural network (LSTM) [7, 30]. Graph neural network methods have also been applied for domain-specific NER tasks [3] and document-level RE [33]. The bidirectional encoder representation from transformers (BERT) [9], a widely-used pretrained language model (PLM), can also tackle the NER [8], RE [22], and entity linking [11] tasks. While the advancement of deep learning-based methods has greatly improved the effectiveness of KG construction, fully-supervised learning requires a large amount of human-annotated data text. Furthermore, the annotation can only be domain-specific, making it difficult to transfer the KG construction work to a new domain, and ultimately limiting the scalability and efficiency of the research in KG.

On the other hand, distant supervision, a weakly supervised learning method, can replace manual annotation with an existing and remote knowledge base. Previous studies have applied remote supervised learning to deal with NER [35], and RE [27, 31] tasks. Thus in this work, we adopted the distant-supervision scheme in the proposed KGDA framework. It should be noted that KG of the coarse domain (e.g., biomedical) generally will not contain the complete knowledge of its finer sub-domains (e.g., oncology). So when we use the coarse-domain KG for distant supervision, labels of the target domain will be limited by the source domain, making it less effective to discover new knowledge. To address this issue, we introduced an iterative strategy to gradually update the model via distant supervision while at the same time using the partially-trained model to discover new entities and relations from the data of the target fine domain.

3 METHODOLOGY

3.1 Notation and task definition

An unstructured sentence $s = [w_1, w_2, w_3, \dots, w_n]$ indicates a sequence of tokens, where n is its length. A dataset \mathbb{D} is a collection of unstructured sentences (i.e. $\mathbb{D} = \{s_1, s_2, s_3, \dots, s_m\}$). The knowledge graph, denoted as \mathbb{K} , is a collection of triples $t = (e_i, r_j, e_k)$, where $e_i \in \mathbb{E}$ and $e_k \in \mathbb{E}$ are the head entity and the tail entity respectively, and $r_j \in \mathbb{V}$ is the relation between e_i and e_k . Here we denote coarse-domain KG as \mathbb{K}_c and fine-domain KG as \mathbb{K}_f .

In a typical scenario of KG domain adaptation, we will have an existing coarse-domain KG and a large amount of unlabeled text in the fine domain. For example, when constructing the oncology KG, we can utilize the existing biomedical KG and collect oncology-related literature as unlabeled text. KG constructed from the fine domain data would then include overlapping triples with the coarse-domain KG and new triples representing domain-specific knowledge. Specifically, the fine-domain KG contains the following three types of triples:

- **Overlapping triples** \mathbb{T}_O : Triples that also existed in the coarse-domain KG, indicating knowledge overlapping between the coarse and fine domains.
- **Triples of new relations but overlapping entities** \mathbb{T}_R : Triples with both entity pairs existing in the coarse-domain KG but no indicated relationships between these entity pairs.

- **Triples of new entities** \mathbb{T}_E : Triples with at least one entity not existing in the coarse-domain KG. Consequently, the relationship is also unknown in the coarse domain.

Both \mathbb{T}_R and \mathbb{T}_E belong to the specific knowledge of the fine domain. The goal of the coarse-to-fine KGDA task is to adapt the KG from the coarse domain to the fine domain and leverage the knowledge from the coarse domain to guide the mining of new knowledge specific to the fine domain. Finally, we will keep the definition of entity types and relation types from coarse-domain KG when constructing the fine-domain KG.

3.2 Iterative training framework

While it is trivial to identify the overlapping entities \mathbb{E}_O and triples \mathbb{T}_O by distant supervision, if the NER and RE models are trained on the entire corpus, they will not be able to recognize the fine domain-specific named entities and triples (\mathbb{T}_R and \mathbb{T}_E). Because the distant-supervision labels are generated by matching \mathbb{K}_c . Thus we introduce an iterative training strategy to construct \mathbb{T}_R and \mathbb{T}_E from the text and adapt the knowledge from \mathbb{K}_c to \mathbb{K}_f .

The overall framework of the iterative training scheme is shown in Fig. 1, and the detailed pseudo code can be found in Algorithm 1. Rather than performing distant-supervision training on the whole unlabeled text corpus, the core mechanism of the proposed iterative training is to split the whole unlabeled dataset into n sub-datasets without intersection. Before building distant-supervision corpus, the trained model is used to predict the text corpus for getting specific knowledge of fine-domain, which is conducive to mining \mathbb{T}_R and \mathbb{T}_E of the fine-domain.

As shown in Figure 1, firstly, it is necessary to preprocess the acquired text corpus in the fine domain. Preprocessing operations include: handling special characters, word segmentation, filtering sentences using human-defined rules (such as sentence length), etc. Then, our framework involves two neural network models: NER model and RE model. We replace the PLM's output layer with a classifier head as NER model $model_N$ and fine-tune it by minimizing the cross-entropy loss on distant-supervision NER corpus. Additionally, we apply the BIO scheme [13] to generate NER sequence labels. For the RE task, we use the template to generate distant-supervision samples. The template we adopted is "[CLS] *head entity (head entity type)* [SEP] *tail entity (tail entity type)* [SEP] *sentence*". The RE model $model_R$ is defined as a PLM with a fully connected layer as a relation classifier. The feature of special token [CLS] fed into this fully connected layer and fine-tune $model_R$ by minimizing the cross-entropy loss on distant-supervision RE corpus.

We summarize the steps to achieve KGDA in Algorithm 1. For the first parts of the text corpus \mathbb{D}_1 , the distant-supervision method is applied to construct the NER training corpus $corp_N$ and RE training corpus $corp_R$, and the NER model $model_N$ and RE model $model_R$ are trained based on corpus $corp_N$ and $corp_R$, respectively. For other part of the text corpus \mathbb{D}_i , we apply the previously trained $model_N$ and $model_R$ to extract the entities and triples in the fine-domain, and select the high confidence entities \mathbb{E}_{conf} and high confidence triples \mathbb{T}_{conf} as the specific knowledge of the fine-domain (line 7). Then, we take \mathbb{K}_c , \mathbb{E}_{conf} , and \mathbb{T}_{conf} as the external knowledge base for constructing distant-supervision $corp_N$ and $corp_R$ (line 8).

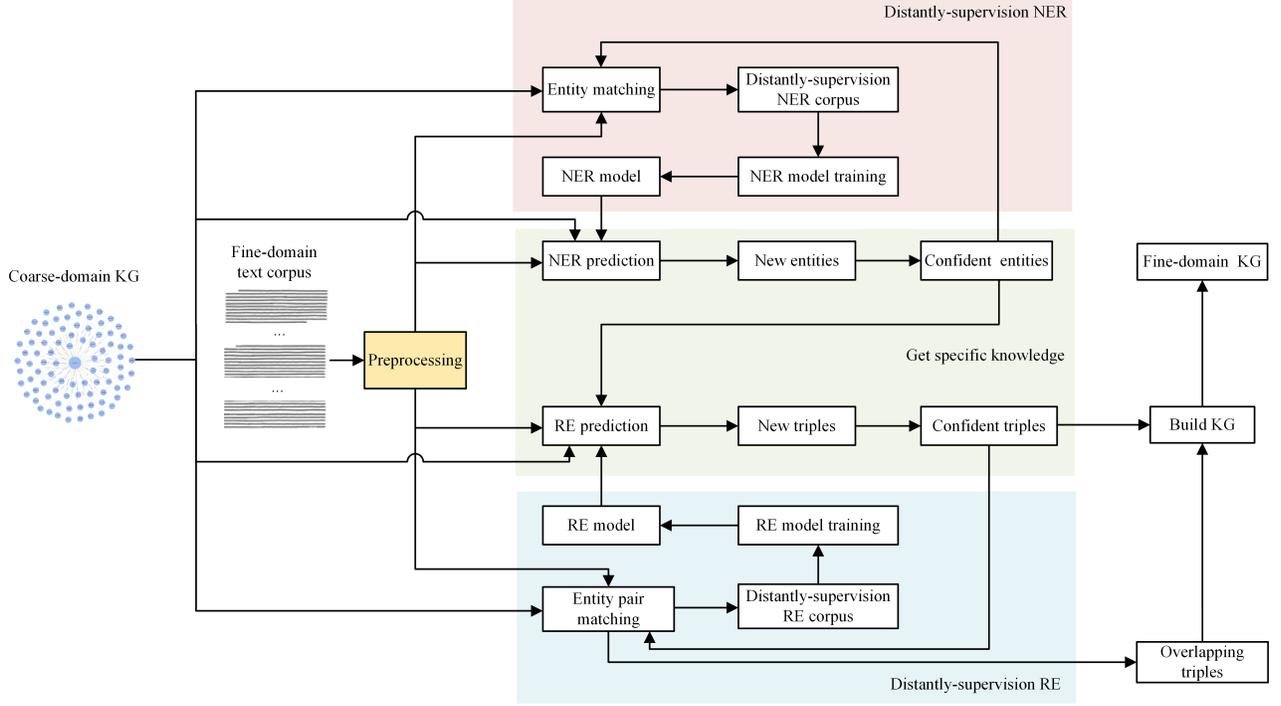


Figure 1: The overall framework of iterative training KGDA.

Finally, we use overlapping triples \mathbb{T}_O and high-confidence triples \mathbb{T}_{conf} to construct a knowledge graph of fine domains (line 17).

Next, we show the details of *get_distant_corpus* in Algorithm 2 and *get_specific_knowledge* in Algorithm 3.

3.3 Constructing distantly-supervised corpus

Through distant-supervision, we can only match entity pairs that have a relationship and use them as positive samples. We then construct negative samples with NULL relationship by the following two schemes: 1) randomly sampling two entities which have no relationship as defined in the coarse-domain; 2) randomly sampling a word from out-of-domain words (i.e., a word that is not an entity as defined in the coarse domain) \mathbb{W}_O as one of the entities. The parameter $ratio_n$ controls the ratio of negative samples (constructed by either schemes) to the total sample size. The parameter $ratio_o$ controls the ratio of entity pairs constructed by the second scheme (i.e., via sampling the words outside the domain) to the size of negative samples, respectively.

In addition to the \mathbb{K}_c in the source domain, we use both \mathbb{K}_c , \mathbb{E}_{conf} , and \mathbb{T}_{conf} as knowledge bases for constructing the remotely supervised corpus. This would ensure that the NER and RE models can identify the overlapping knowledge between \mathbb{K}_c and \mathbb{K}_f , while at the same time be guided to discover the new knowledge specific to the fine domain.

As shown in Algorithm 2, for building the distantly-supervised NER corpus $corp_N$, the sentence \mathbb{D}_i^j is firstly string-matched with the knowledge bases \mathbb{K}_c and \mathbb{E}_{conf} to extract the entities in the sentence (line 5). Afterward, the matched entities are merged into overlapping entities \mathbb{E}_O , and the NER label sequences are generated

through the BIO strategy to merge into $corp_N$ (line 6 and 7). For building the distantly-supervised RE corpus $corp_R$, we firstly take \mathbb{K}_c and \mathbb{T}_{conf} as knowledge bases and use entity pair matching to match the triples $triples_k$ based on \mathbb{K}_c and the triples $triples_c$ based on \mathbb{T}_{conf} appearing in the sentence \mathbb{D}_i^j (line 8). We then build negative triples with parameters $ratio_n$ and $ratio_o$ (line 10). Finally, we construct the RE corpus based on the triples $triples$, $triples_n$ and corresponding sentences through a pre-defined relationship sample template (line 11 and 12).

3.4 Discovering fine-domain specific knowledge

Recall that in the proposed iterative training framework, the whole unlabeled dataset is divided into n sub-dataset \mathbb{D}_i , $i = 1 \dots n$, the fine-domain specific knowledge discovery will be performed on each sub-dataset except the first one \mathbb{D}_i , $i = 2 \dots n$ (line 5 to 16 in Algorithm 1). For each new sub-dataset \mathbb{D}_i , $i = 2 \dots n$, we will use the previously-updated models $model_N$ and $model_R$ to predict the new entities and triples. Afterward, the sub-dataset will be used for updating $model_N$ and $model_R$ via distantly-supervised training. As noisy or incorrect entities and triples could be discovered during this procedure, we developed a filtering mechanism only to keep the entities and triples with higher confidence. Specifically, we design the rules for filtering the discovered entities and triples by: 1) probability of the new entities and triples predicted by the corresponding models should be greater than pre-defined thresholds th_{pe} and th_{pt} , respectively; 2) cumulative frequency of the new entities and triples discovered

Algorithm 1 Iterative training KGDA framework

Input: Text corpus $\mathbb{D} = \{\mathbb{D}_1, \mathbb{D}_2, \dots, \mathbb{D}_n\}$, coarse-domain KG \mathbb{K}_c , out-of-domain words \mathbb{W}_O

Parameter: Initialized NER model $model_N$, initialized RE model $model_R$

Output: fine-domain kg \mathbb{K}_f

- 1: Let new entities $\mathbb{E}_{new} = \{\}$, new entities with high confidence $\mathbb{E}_{conf} = \{\}$, new triples $\mathbb{T}_{new} = \{\}$, new triples with high confidence $\mathbb{T}_{conf} = \{\}$.
- 2: $corp_N, corp_R, \mathbb{E}_O, \mathbb{T}_O = \text{build_distant_corpus}(\mathbb{D}_1, \mathbb{K}_c, \mathbb{E}_{conf}, \mathbb{T}_{conf}, \mathbb{W}_O)$
- 3: $\text{train_NER}(model_N, corp_N)$
- 4: $\text{train_RE}(model_R, corp_R)$
- 5: $i = 2$
- 6: **while** $i \leq n$ **do**
- 7: $\mathbb{E}_{new}, \mathbb{E}_{conf}, \mathbb{T}_{new}, \mathbb{T}_{conf} = \text{get_specific_knowledge}(\mathbb{D}_i, \mathbb{K}_c, \mathbb{E}_{new}, \mathbb{E}_{conf}, \mathbb{T}_{new}, \mathbb{T}_{conf})$
- 8: $corp'_N, corp'_R, \mathbb{E}'_O, \mathbb{T}'_O = \text{get_distant_corpus}(\mathbb{D}_i, \mathbb{K}_c, \mathbb{E}_{conf}, \mathbb{T}_{conf}, \mathbb{W}_O)$
- 9: $corp_N = corp_N \cup corp'_N$
- 10: $corp_R = corp_R \cup corp'_R$
- 11: $\mathbb{E}_O = \mathbb{E}_O \cup \mathbb{E}'_O$
- 12: $\mathbb{T}_O = \mathbb{T}_O \cup \mathbb{T}'_O$
- 13: $\text{train_NER}(model_N, corp_N)$
- 14: $\text{train_RE}(model_R, corp_R)$
- 15: $i = i + 1$
- 16: **end while**
- 17: $\mathbb{K}_f = \text{build_kg}(\mathbb{T}_O, \mathbb{T}_{conf})$
- 18: **return** \mathbb{K}_f

from datasets \mathbb{D}_2 to \mathbb{D}_i should be greater than the pre-defined thresholds th_{fe} and th_{ft} , respectively.

As shown in Algorithm 3, for discovering new entities \mathbb{E}_{new} , we will apply the trained $model_N$ on dataset \mathbb{D}_i and obtain *entities* that are disjoint with \mathbb{K}_c (line 5 and 6). Then, we will merge *entities* with the previously-discovered entity set \mathbb{E}_{new} (line 7). Finally, we will select the "high-confident" entity as \mathbb{E}_{conf} based on the mechanism above by the prediction probability and cumulative frequency (line 10). For the discovery of new triples \mathbb{T}_{new} , we will enumerate entity pairs that are disjoint with the \mathbb{K}_c (line 13 - 15). We will then use the trained RE model and the predefined sample template to predict the relationship of the entity pairs and delete the triples whose predicted relationship is NULL (line 16). Other processing is similar to the discovery of new entities.

After Algorithm 3, discovered entities specific to the fine domain are stored in \mathbb{E}_{conf} . Discovered triples \mathbb{T}_R (new relation, overlapping entity) and \mathbb{T}_E (new relation, new entity) are stored in \mathbb{T}_{conf} . In the next iteration, Algorithm 2 will then use the updated \mathbb{E}_{conf} and \mathbb{T}_{conf} for building distant-supervision corpus. Such iterative design can facilitate the interoperability between the two competing tasks based on a fixed number of unannotated data samples in the fine target domain: distantly-supervised training of the NER and RE models versus the discovery of new knowledge using the trained

Algorithm 2 Constructing distantly-supervised corpus

Input: A part of text corpus text corpus \mathbb{D}_i , coarse-domain KG \mathbb{K}_c , new entities with high confidence \mathbb{E}_{conf} , new triples with high confidence \mathbb{T}_{conf} , out-of-domain words \mathbb{W}_O

Parameter: negative sample ratio $ratio_n$, out-of-domain sample ratio $ratio_o$

Output: Distant-supervision NER corpus $corp_N$, distant-supervision RE corpus $corp_R$, overlapping entities \mathbb{E}_O , overlapping triples \mathbb{T}_O

- 1: Let $corp_E = \{\}$, $corp_R = \{\}$, $\mathbb{E}_O = \{\}$, $\mathbb{T}_O = \{\}$.
- 2: $\text{sentence_num} = \text{len}(\mathbb{D}_i)$
- 3: $j = 1$
- 4: **while** $j \leq \text{sentence_num}$ **do**
- 5: $\text{entities} = \text{entity_matching}(\mathbb{D}_i^j, \mathbb{K}_c, \mathbb{E}_{conf})$
- 6: $\mathbb{E}_O = \mathbb{E}_O \cup \text{entities}$
- 7: $corp_N = corp_N \cup \text{build_NER_sample}(\mathbb{D}_i^j, \text{entities})$
- 8: $\text{triples}_k, \text{triples}_c = \text{entity_pair_matching}(\mathbb{D}_i^j, \mathbb{K}_c, \mathbb{T}_{conf})$
- 9: $\text{triples} = \text{triples}_k \cup \text{triples}_c$
- 10: $\text{triples}_n = \text{get_negative_triples}(\mathbb{D}_i^j, \mathbb{W}_O, \text{triples}, ratio_n, ratio_o)$
- 11: $corp_R = corp_R \cup \text{get_samples}(\text{triples})$
- 12: $corp_R = corp_R \cup \text{get_samples}(\text{triples}_n)$
- 13: $\mathbb{T}_O = \mathbb{T}_O \cup \text{triples}_k$
- 14: $j = j + 1$
- 15: **end while**
- 16: **return** $corp_N, corp_R, \mathbb{E}_O, \mathbb{T}_O$

NER and RE models, thus improve the efficiency of performing KG domain adaptation and construction without any annotation.

4 EXPERIMENTS

In this work, we used the adaptation of KG from the biomedical domain (coarse) to the oncology domain (fine) as an example to demonstrate the workflow of the KGDA framework, as well as to evaluate its effectiveness in practice. Implementation details of the experiment are also provided, along with the publicly-available data and the containerized environment in the released source code, for easy replication of the experiment and the development of other KG methods.

4.1 Dataset

The dataset we used in this paper is released by [19]. The source data is downloaded from 12 international journals in the oncology domain. PDF files of the papers were cleaned and converted to sentences. In total, we select 240,000 paragraphs as the unlabeled text corpus of the oncology domain \mathbb{D} . The coarse-domain KG \mathbb{K}_c used in this work is the biomedical KG¹, defines 18 entity types (anatomy, neoplastic process, microorganism, eukaryote, physiology, chemical or drug, diagnostic procedure, laboratory procedure, research activity or technique, therapeutic or preventive procedure, medical device, research device, pathology, disease or syndrome, anatomical abnormality, mental or behavioral dysfunction, injury or poisoning and sign, symptom or finding) and 19 relationship

¹<https://idea.edu.cn/bios.html>

Algorithm 3 Discovering fine-domain specific knowledge

Input: A part of text corpus \mathbb{D}_i , coarse-domain KG \mathbb{K}_c , new entities \mathbb{E}_{new} , new entities with high confidence \mathbb{E}_{conf} , new triples \mathbb{T}_{new} , new triples with high confidence \mathbb{T}_{conf}

Parameter: NER model $model_N$, RE model $model_R$, probability threshold of the entity th_{pe} , frequency threshold of the entity th_{fe} , probability threshold of the triple th_{pt} , frequency threshold of the triple th_{ft}

Output: \mathbb{E}_{new} , \mathbb{E}_{conf} , \mathbb{T}_{new} , \mathbb{T}_{conf}

```

1: Let  $corp_E = \{\}$ ,  $corp_R = \{\}$ ,  $\mathbb{E}_O = \{\}$ ,  $\mathbb{T}_O = \{\}$ .
2:  $sentence\_num = \text{len}(\mathbb{D}_i)$ 
3:  $j = 1$ 
4: while  $j \leq sentence\_num$ 
5:    $entities = \text{NER\_prediction}(\mathbb{D}_i^j, model_N)$ 
6:    $entities = \text{get\_new\_entities}(entities, \mathbb{K}_c)$ 
7:    $\mathbb{E}_{new} = \text{merge\_entity}(\mathbb{E}_{new}, entities)$ 
8:    $j = j + 1$ 
9: end while
10:  $\mathbb{E}_{conf} = \text{get\_confidence\_entity}(\mathbb{E}_{new}, th_{pe}, th_{fe})$ 
11:  $j = 1$ 
12: while  $j \leq sentence\_num$ 
13:    $entities = \text{entity\_matching}(\mathbb{D}_i^j, \mathbb{K}_c, \mathbb{E}_{conf})$ 
14:    $pairs = \text{enumerate\_pairs}(entities)$ 
15:    $pairs = \text{get\_new\_pairs}(pairs, \mathbb{K}_c)$ 
16:    $triples = \text{RE\_prediction}(\mathbb{D}_i^j, pairs, model_R)$ 
17:    $\mathbb{T}_{new} = \text{merge\_triple}(\mathbb{T}_{new}, triples)$ 
18:    $j = j + 1$ 
19: end while
20:  $\mathbb{T}_{conf} = \text{get\_confidence\_triple}(\mathbb{T}_{new}, th_{pt}, th_{ft})$ 
21: return  $\mathbb{E}_{new}$ ,  $\mathbb{E}_{conf}$ ,  $\mathbb{T}_{new}$ ,  $\mathbb{T}_{conf}$ 

```

types (is_a, reverse_is_a, is_part_of, reverse_is_part_of, may_treat, reverse_may_treat, found_in, reverse_found_in, may_cause, reverse_may_cause, expressed_in, is_expression_of, encodes, encoded_by, significant_drug_interaction, involved_in_biological_process, biological_process_involves, is_active_ingredient_in, has_active_ingredient), including 5.2 million English entities and 7.34 million triples.

4.2 Evaluation

Similar to the previous works [18], we evaluate our method in two schemes: held-out evaluation and manual evaluation. For the held-out evaluation, we reserved a part of the text corpus of \mathbb{D} as the test set. During the testing, we then compared the prediction results of the NER and RE models with the labels matched with \mathbb{K}_c , and calculated the precision, recall, and F1 of the held-out dataset. Specifically, we use seqvel² to evaluate the micro average precision, recall, F1 of NER. When evaluating the RE model, we perform relation classification prediction on the triples existing in \mathbb{K}_c and corresponding entity pairs appearing in the held-out corpus. Finally, weighted average precision, recall, and F1 from the held-out evaluation will be reported.

As the labels of testing samples in the held-out evaluation are all inferred by distant supervision from the coarse domain, such

scheme can only evaluate whether the trained model can capture the knowledge in the coarse domain, but cannot evaluate the ability of the models to discover new knowledge in the fine-domain. Therefore, we also adopted the manual evaluation scheme, consisting of the evaluations of: 1) the entities specific to fine domain \mathbb{E}_{conf} , which are not presented in \mathbb{K}_c ; 2) the triples of new relations \mathbb{T}_R ; 3) the triples of new entities \mathbb{T}_E . We randomly sampled 50 cases of \mathbb{E}_{conf} , \mathbb{T}_R , and \mathbb{T}_E respectively, then asked one physician to manually label them for whether the entities and triples are correct. As the number of name entities and triples instances that are expressed in the corpus is unknown, we cannot estimate the recall of fine-domain KG. Therefore, we only show the precision of \mathbb{E}_{conf} , \mathbb{T}_R , and \mathbb{T}_E . We fully recognize that the discovery of new knowledge in the fine-domain is an indispensable task for this work and we are recruiting more medical experts to conduct human reader study and performance evaluation for the proposed model.

4.3 Implementation settings

We divide the corpus \mathbb{D} into six equal subsets, and each subset contains around 40,000 sentences. We used \mathbb{D}_1 to \mathbb{D}_5 for model training and KG construction. We reserved \mathbb{D}_6 for held-out evaluation. We tested BERT [9], Bio_ClinicalBERT [1], biomed_RoBERTa [6] for initializing NER and RE models. Our experiments were run on an Ubuntu system computer with 4 NVIDIA A100 graphics cards. The learning rate, batch size, and epochs are set as 2E-05, 20, and 4, respectively. Hyperparameters $th_{fe}, th_{pe}, th_{ft}, th_{pt}$ are set as 2, 0.95, 3, 0.97. The parameters $ratio_n$ and $ratio_o$ that control negative sampling are set to 0.2 and 0.3.

4.4 Held-out evaluation

models	precision	recall	F1
BERT	0.908	0.900	0.904
Bio_ClinicalBERT	0.909	0.895	0.902
biomed_RoBERTa	0.908	0.901	0.905

Table 1: Held-out evaluation of NER model.

models	precision	recall	F1
BERT	0.987	0.949	0.967
Bio_ClinicalBERT	0.988	0.957	0.972
biomed_RoBERTa	0.987	0.959	0.972

Table 2: Held-out evaluation of RE model.

The results of the NER and RE models evaluated by the held-out dataset are shown in Table 1 and Table 2, respectively. The KGDA frameworks initialized by the three pre-trained language models (BERT, Bio_ClinicalBERT, and biomed_RoBERTa) all show good performance in held-out evaluations, demonstrating the robustness of our framework. Because Bio_ClinicalBERT and biomed_RoBERTa are pre-trained in biomedical data sets, their performance is better than BERT.

²<https://github.com/chakki-works/seqeval>

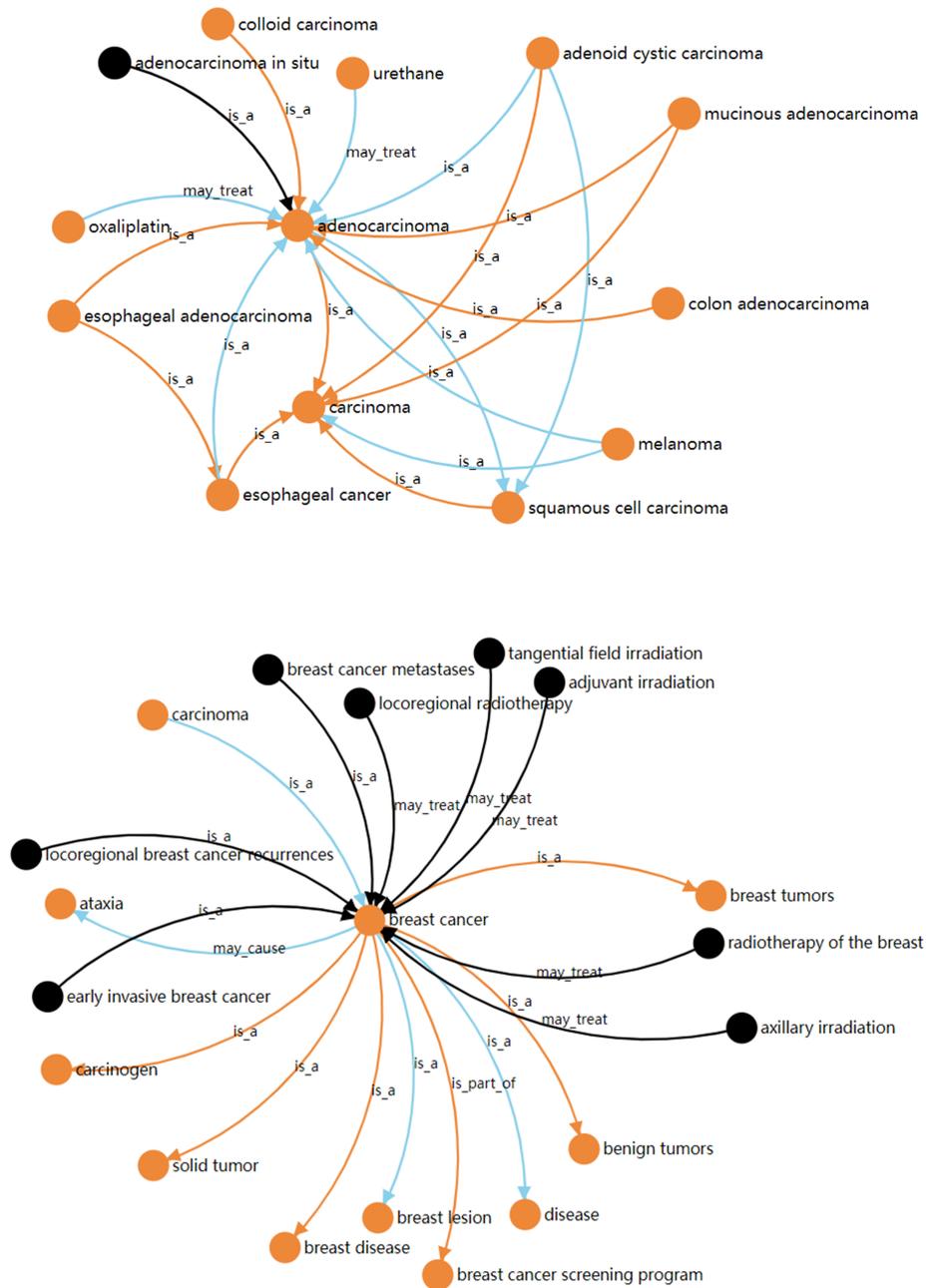


Figure 2: knowledge graph of the *adenocarcinoma* and *breast cancer*. The orange edges represent the overlapping triples \mathbb{T}_O , the blue edges denote the triples of new relations \mathbb{T}_R , and the black edges denote the triples of new entities \mathbb{T}_E , while the black node represents the new entity. It should be noted that in order to facilitate the display, we only show some associated triples, not all.

4.5 Manual evaluation

The number of all discovered entities (\mathbb{E}_O), triples (\mathbb{T}_O), new entities with high confidence (\mathbb{E}_{conf}), triples representing new relations

with overlapping entities (\mathbb{T}_R), and triples representing new relations with new entities (\mathbb{T}_E) are shown in Table 3, with each row belonging to one pre-trained language models used. Numbers of

models	# \mathbb{E}_O	# \mathbb{T}_O	# \mathbb{E}_{conf}	# \mathbb{T}_R	# \mathbb{T}_E
BERT	86741	26467	1378	36010	1178
Bio_ClinicalBERT	86801	26353	1541	39588	1413
biomed_RoBERTa	86821	26448	1444	37637	1110

Table 3: The number of entities and triples.

models	\mathbb{E}_{conf}	\mathbb{T}_R	\mathbb{T}_E
BERT	0.90	0.58	0.70
Bio_ClinicalBERT	0.90	0.66	0.62
biomed_RoBERTa	0.94	0.76	0.74

Table 4: results of manual evaluations.

\mathbb{E}_O and \mathbb{T}_O have minor differences among different pre-trained language models, possibly due to the conflicts in strings matching of knowledge bases. \mathbb{E}_{conf} , \mathbb{T}_R , and \mathbb{T}_E represent specific knowledge of the fine domain. We sampled 50 cases from \mathbb{E}_{conf} , \mathbb{T}_R , and \mathbb{T}_E for manual evaluation, and the results are shown in Table 4.

4.6 Knowledge graph construction in the fine domain

As our ultimate goal, we can construct the KG in the fine domain by combining \mathbb{T}_O , \mathbb{T}_R , and \mathbb{T}_E . We selected biomed_RoBERTa as the backbone language model for KGDA and constructed the knowledge graph correspondingly. An example of the KG we built is shown in the supplementary material. Knowledge related to *adenocarcinoma* and *breast cancer* are visualized in Figure 2.

4.7 Ablation study

models	NER precision	RE precision
BERT	0.908	0.987
w/o (cumulative)	0.888	0.985
w/o (iter)	0.894	0.984
Bio_ClinicalBERT	0.909	0.988
w/o (cumulative)	0.887	0.985
w/o (iter)	0.898	0.985
biomed_RoBERTa	0.908	0.987
w/o (cumulative)	0.888	0.983
w/o (iter)	0.894	0.986

Table 5: Results of ablation study.

We investigated the impact of techniques employed by KGDA on its held-out experiment performance by removing the corresponding component from the framework:

w/o (cumulative): When using corpus \mathbb{D}_i to train NER and RE models, the cumulative corpus is not used. i.e. delete lines 9 and 10 in Algorithm 1 and mark $corp'_N$ and $corp'_R$ in line 8 as $corp_N$ and $corp_R$ respectively.

w/o (iter): Remove the iterative training strategy and only use \mathbb{K}_C as an external knowledge base.

The results of the ablation analysis are shown in Table 5. Compared to the complete framework with w/o (cumulative), it can

be seen that the using of accumulated data through iterations is beneficial for improving the generalization ability of NER and RE models. The held-out performances of the model without iteration indicates that the iterative training strategy can not only discover the specific knowledge in the fine domain but also maintain the ability to discover overlapping knowledge between the coarse and fine domain.

4.8 Case study

In this case study, we illustrate a sample case of the model input (text from an abstract) and the model output (triplets extracted from the text) during the testing phase 3. It can be observed from the results that the proposed KGDA framework can recover name entities that are unique in the oncology domain (e.g., immune evasion and cd8+t cells). We have also asked our physicians to review the extracted relations, which are all clinically meaningful and accurate.

5 CONCLUSION AND DISCUSSION

In this paper, we propose an integrated, end-to-end framework for knowledge graph domain adaptation using distant supervision, which can be used to construct KG from fully unlabeled raw text data with the guidance of an existing KG. To deal with the potential challenges in distant supervision, which might limit the knowledge discovered from the new domain, we propose an iterative training strategy, which divides an unlabeled corpus into multiple corpuses. For each new corpus to the model, we then combine the knowledge in the coarse domain with the knowledge identified from the previous corpuses for distantly-supervised training. By adopting the iterative training strategy, our proposed KGDA framework can discover not only knowledge that overlaps with the coarse domain, but also knowledge specific to the fine domain and unknown to the coarse domain, thus enabling coarse-to-fine domain adaptation. We implemented the adaptation from biomedical KG to the oncology domain in our experiments and verified the effectiveness of the KGDA framework through held-out and manual evaluation.

Several limitations and challenges remain beyond the current work for more effective and accurate KG construction: Firstly, more thorough evaluation with human reader study is needed to validate that new knowledge relevant (not only correct) to the target domain can be discovered by KGDA. Secondly, it has been recognized by the field that distant supervision will inevitably introduce noisy labels [14, 31], thus the denoising step is usually needed but not implemented in the current version of KGDA. Thirdly, there has been existing KG constructed in the related domains of oncology and cancer research. We will investigate the scheme to allow adaption from multiple sources (not only the coarse domain) to leverage this existing knowledge better. Another type of crucial prior information for this work is clinical ontology, where we will integrate the relationships defined in ontology and entity description to enhance the model. Fourthly, an essential premise of the KGDA is that we assume the source and target domains share the same set of entity types and relation types, which can limit the knowledge discovered from the fine domain. We will investigate data mining techniques to adaptively add/remove entity and relation types in the fine domain. Finally, there have been many new large-scale pre-trained language models developed such as GPT-3 in recent years. While our model

Text:

The CD155/TIGIT axis can be co-opted during immune evasion in chronic viral infections and cancer. Pancreatic adenocarcinoma (PDAC) is a highly lethal malignancy, and immune-based strategies to combat this disease have been largely unsuccessful to date. We corroborate prior reports that a substantial portion of PDAC harbors predicted high-affinity MHC class I-restricted neopeptides and extend these findings to advanced/metastatic disease. Using multiple preclinical models of neoantigen-expressing PDAC, we demonstrate that intratumoral neoantigen-specific CD8+ T cells adopt multiple states of dysfunction, resembling those in tumor-infiltrating lymphocytes of PDAC patients. Mechanistically, genetic and/or pharmacologic modulation of the CD155/TIGIT axis was sufficient to promote immune evasion in autochthonous neoantigen-expressing PDAC. Finally, we demonstrate that the CD155/TIGIT axis is critical in maintaining immune evasion in PDAC and uncover a combination immunotherapy (TIGIT/PD-1 co-blockade plus CD40 agonism) that elicits profound anti-tumor responses in preclinical models, now poised for clinical evaluation.

Triplets:

Head entities	Tail entities	relations
axis	genetic	is_a
immune evasion	disease	is_a
immune evasion	dysfunction	is_a
immune evasion	genetic	is_a
malignancy	disease	is_a
cd8+ t cells	immune evasion	involved_in_biological_process
cd8+ t cells	lymphocytes	is_a
dysfunction	genetic	is_a

Figure 3: Text from a sample abstract used as input to the model testing and the corresponding triplets extracted from the text.

uses variations of BERT (BiomedRoBERTa and BioClinicalBERT) as backbone networks, we can easily adapt KGDA to other language models.

REFERENCES

- [1] Emily Alsentzer, John R Murphy, Willie Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, WA Redmond, and Matthew BA McDermott. 2019. Publicly Available Clinical BERT Embeddings. *NAACL HLT 2019* (2019), 72.
- [2] Gabor Angeli, Melvin Jose Johnson Premkumar, and Christopher D Manning. 2015. Leveraging linguistic structure for open domain information extraction. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 344–354.
- [3] Pei Chen, Haibo Ding, Jun Araki, and Ruihong Huang. 2021. Explicitly Capturing Relations between Entity Mentions via Graph Neural Networks for Domain-specific Named Entity Recognition. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*.
- [4] Anthony Fader, Stephen Soderland, and Oren Etzioni. 2011. Identifying relations for open information extraction. In *Proceedings of the 2011 conference on empirical methods in natural language processing*. 1535–1545.
- [5] Jason Fries, Sen Wu, Alex Ratner, and Christopher Ré. 2017. Swellshark: A generative model for biomedical named entity recognition without labeled data. *arXiv preprint arXiv:1704.06360* (2017).
- [6] Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. 2020. Don't Stop Pretraining: Adapt Language Models to Domains and Tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 8342–8360.
- [7] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [8] Chen Jia, Yuefeng Shi, Qinrong Yang, and Yue Zhang. 2020. Entity enhanced BERT pre-training for Chinese NER. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 6384–6396.
- [9] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of NAACL-HLT*. 4171–4186.
- [10] Natthawut Kertkeidkachorn and Ryutaro Ichise. 2017. T2kg: An end-to-end system for creating knowledge graph from unstructured text. In *Workshops at the Thirty-First AAAI Conference on Artificial Intelligence*.
- [11] Da Li, Ming Yi, and Yukai He. 2022. LP-BERT: Multi-task Pre-training Knowledge Graph BERT for Link Prediction. *arXiv preprint arXiv:2201.04843* (2022).
- [12] Nan Li, Qiang Shen, Rui Song, Yang Chi, and Hao Xu. 2022. MEDuKG: A Deep-Learning-Based Approach for Multi-Modal Educational Knowledge Graph Construction. *Information* 13, 2 (2022), 91.
- [13] Qi Li, Haibo Li, Heng Ji, Wen Wang, Jing Zheng, and Fei Huang. 2012. Joint bilingual name tagging for parallel corpora. In *Proceedings of the 21st ACM international conference on Information and knowledge management*. 1727–1731.
- [14] Chen Liang, Yue Yu, Haoming Jiang, Siawpeng Er, Ruijia Wang, Tuo Zhao, and Chao Zhang. 2020. Bond: Bert-assisted open-domain named entity recognition with distant supervision. In *Proceedings of the 26th ACM SIGKDD International*

- Conference on Knowledge Discovery & Data Mining*. 1054–1064.
- [15] Donald AB Lindberg, Betsy L Humphreys, and Alexa T McCray. 1993. The unified medical language system. *Yearbook of medical informatics* 2, 01 (1993), 41–51.
- [16] 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.
- [17] Aman Mehta, Aashay Singhal, and Kamalakar Karlapalem. 2019. Scalable knowledge graph construction over text using deep learning based predicate mapping. In *Companion Proceedings of The 2019 World Wide Web Conference*. 705–713.
- [18] Mike Mintz, Steven Bills, Rion Snow, and Dan 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*. 1003–1011.
- [19] Saeed Rezayi, Haixing Dai, Zhengliang Liu, Zihao Wu, Akarsh Hebbar, Andrew H Burns, Lin Zhao, Dajiang Zhu, Quanzheng Li, Wei Liu, et al. 2022. ClinicalRadioBERT: Knowledge-Infused Few Shot Learning for Clinical Notes Named Entity Recognition. In *Machine Learning in Medical Imaging: 13th International Workshop, MLMI 2022, Held in Conjunction with MICCAI 2022, Singapore, September 18, 2022, Proceedings*. Springer, 269–278.
- [20] Anderson Rossanez, Julio Cesar Dos Reis, Ricardo da Silva Torres, and Hélène de Ribapierre. 2020. KGen: a knowledge graph generator from biomedical scientific literature. *BMC medical informatics and decision making* 20, 4 (2020), 1–24.
- [21] Maya Rotmensch, Yoni Halpern, Abdulkhakim Tlimat, Steven Horng, and David Sontag. 2017. Learning a health knowledge graph from electronic medical records. *Scientific reports* 7, 1 (2017), 1–11.
- [22] Arpita Roy and Shimei Pan. 2021. Incorporating medical knowledge in BERT for clinical relation extraction. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. 5357–5366.
- [23] Michael Schmitz, Stephen Soderland, Robert Bart, Oren Etzioni, et al. 2012. Open language learning for information extraction. In *Proceedings of the 2012 joint conference on empirical methods in natural language processing and computational natural language learning*. 523–534.
- [24] Alisa Smirnova and Philippe Cudré-Mauroux. 2018. Relation extraction using distant supervision: A survey. *ACM Computing Surveys (CSUR)* 51, 5 (2018), 1–35.
- [25] Michael Stewart and Wei Liu. 2020. Seq2kg: an end-to-end neural model for domain agnostic knowledge graph (not text graph) construction from text. In *Proceedings of the International Conference on Principles of Knowledge Representation and Reasoning*, Vol. 17. 748–757.
- [26] Jalaj Thanaki. 2017. *Python natural language processing*. Packt Publishing Ltd.
- [27] Siheng Wei. 2021. Distantly Supervision for Relation Extraction via LayerNorm Gated Recurrent Neural Networks. In *2021 2nd International Conference on Computing and Data Science (CDS)*. IEEE, 94–99.
- [28] Haoze Yu, Haisheng Li, Dianhui Mao, and Qiang Cai. 2021. A domain knowledge graph construction method based on Wikipedia. *Journal of Information Science* 47, 6 (2021), 783–793.
- [29] Hamada M Zahera, Stefan Heindorf, and Axel-Cyrille Ngonga Ngomo. 2021. ASSET: A Semi-supervised Approach for Entity Typing in Knowledge Graphs. In *Proceedings of the 11th on Knowledge Capture Conference*. 261–264.
- [30] Donghuo Zeng, Chengjie Sun, Lei Lin, and Bingquan Liu. 2017. LSTM-CRF for drug-named entity recognition. *Entropy* 19, 6 (2017), 283.
- [31] Yue Zhang, Hongliang Fei, and Ping Li. 2021. ReadsRE: Retrieval-Augmented Distantly Supervised Relation Extraction. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2257–2262.
- [32] Yong Zhang, Ming Sheng, Rui Zhou, Ye Wang, Guangjie Han, Han Zhang, Chunxiao Xing, and Jing Dong. 2020. HKGB: An Inclusive, Extensible, Intelligent, Semi-auto-constructed Knowledge Graph Framework for Healthcare with Clinicians' Expertise Incorporated. *Information Processing and Management* 57, 6 (2020), 102324. <https://doi.org/10.1016/j.ipm.2020.102324>
- [33] Zhenyu Zhang, Bowen Yu, Xiaobo Shu, Tingwen Liu, Hengzhu Tang, Wang Yubin, and Li Guo. 2020. Document-level relation extraction with dual-tier heterogeneous graph. In *Proceedings of the 28th International Conference on Computational Linguistics*. 1630–1641.
- [34] Mingxiong Zhao, Han Wang, Jin Guo, Di Liu, Cheng Xie, Qing Liu, and Zhibo Cheng. 2019. Construction of an industrial knowledge graph for unstructured chinese text learning. *Applied Sciences* 9, 13 (2019), 2720.
- [35] Honghao Zheng, Hongtao Yu, Yinuo Hao, Yiteng Wu, and Shaomei Li. 2021. Distantly Supervised Named Entity Recognition with Spy-PU Algorithm. In *2021 IEEE 2nd International Conference on Pattern Recognition and Machine Learning (PRML)*. IEEE, 56–63.