

# Text2NKG: Fine-Grained N-ary Relation Extraction for N-ary relational Knowledge Graph Construction

Haoran Luo<sup>1</sup>, Haihong E<sup>1\*</sup>, Yuhao Yang<sup>2</sup>, Tianyu Yao<sup>1</sup>,  
Yikai Guo<sup>3</sup>, Zichen Tang<sup>1</sup>, Wentai Zhang<sup>1</sup>, Shiyao Peng<sup>1</sup>, Kaiyang Wan<sup>1</sup>,  
Meina Song<sup>1</sup>, Wei Lin<sup>4</sup>, Yifan Zhu<sup>1</sup>, Luu Anh Tuan<sup>5</sup>

<sup>1</sup>School of Computer Science, Beijing University of Posts and Telecommunications, China

<sup>2</sup>School of Automation Science and Electrical Engineering, Beihang University, China

<sup>3</sup>Beijing Institute of Computer Technology and Application <sup>4</sup>Inspur Group Co., Ltd., China

<sup>5</sup>College of Computing and Data Science, Nanyang Technological University, Singapore

{luohaoran, ehaihong, yifan\_zhu}@bupt.edu.cn, anhtuan.luu@ntu.edu.sg

## Abstract

Beyond traditional binary relational facts, n-ary relational knowledge graphs (NKGs) are comprised of n-ary relational facts containing more than two entities, which are closer to real-world facts with broader applications. However, the construction of NKGs remains at a coarse-grained level, which is always in a single schema, ignoring the order and variable arity of entities. To address these restrictions, we propose Text2NKG, a novel fine-grained n-ary relation extraction framework for n-ary relational knowledge graph construction. We introduce a span-tuple classification approach with hetero-ordered merging and output merging to accomplish fine-grained n-ary relation extraction in different arity. Furthermore, Text2NKG supports four typical NKG schemas: *hyper-relational schema*, *event-based schema*, *role-based schema*, and *hypergraph-based schema*, with high flexibility and practicality. The experimental results demonstrate that Text2NKG achieves state-of-the-art performance in  $F_1$  scores on the fine-grained n-ary relation extraction benchmark. Our code and datasets are publicly available<sup>1</sup>.

## 1 Introduction

Modern knowledge graphs (KGs), such as Freebase [2], Google Knowledge Vault [7], and Wikidata [21], utilize a multi-relational graph structure to represent knowledge. Because of the advantage of intuitiveness and interpretability, KGs find various applications in question answering [28], query response [1], logical reasoning [4], and recommendation systems [29]. Traditional KGs are mostly composed of binary relational facts (*subject, relation, object*), which represent the relationship between two entities [3]. However, it has been observed [20] that over 30% of real-world facts involve n-ary relation facts with more than two entities ( $n \geq 2$ ). As shown in Figure 1, an n-ary relational knowledge graph (NKG) is composed of many n-ary relation facts, offering richer knowledge expression

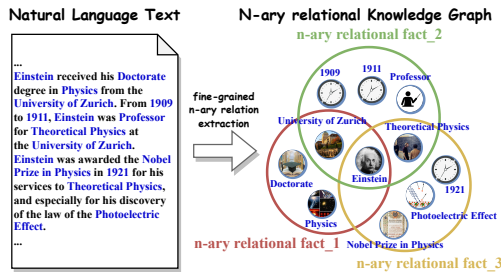


Figure 1: An example of NKG construction.

As shown in Figure 1, an n-ary relational knowledge graph (NKG) is composed of many n-ary relation facts, offering richer knowledge expression

\* Corresponding author.

<sup>1</sup> <https://github.com/LHRLAB/Text2NKG>

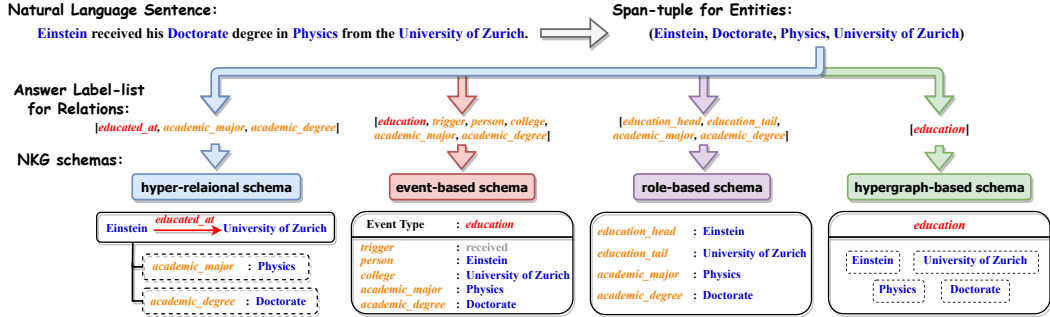


Figure 2: Taking a real-world textual fact as an example, we can extract a four-arity structured span-tuple for entities (Einstein, University of Zurich, Doctorate, Physics) with an answer label-list for relations accordingly as a 4-ary relational fact from the sentence through n-ary relation extraction.

and wider application capabilities. As a key step of constructing NKGs, n-ary relation extraction (n-ary RE) is a task of identifying n-ary relations among entities in natural language texts. Compared to binary relational facts, n-ary relational facts in NKGs have more diverse schemas for different scenarios. For example, Wikidata utilizes n-ary relational facts in a *hyper-relational schema* [20, 10, 23], i.e.,  $(s, r, o, \{(k_i, v_i)\}_{i=1}^{n-2})$  which adds  $(n - 2)$  key-value pairs to the main triple to represent auxiliary information. In addition to the *hyper-relational schema*, the existing NKG schemas also include *event-based schema*  $(r, \{(k_i, v_i)\}_{i=1}^n)$  [11, 16], *role-based schema*  $(\{(k_i, v_i)\}_{i=1}^n)$  [12, 15], and *hypergraph-based schema*  $(r, \{v_i\}_{i=1}^n)$  [26, 8], as shown in Figure 2.

Currently, most existing NKGs in four schemas, such as JF17K [26], Wikipedea [12], WD50K [10], and EventKG [11], are manually constructed. Previous n-ary RE methods [13, 31] focus on extraction with a fixed number of entities in *hypergraph-based schema* or *role-based schema*. Existing event extraction methods [16, 17, 9] can achieve n-ary RE in *event-based schema*. Recently, CubeRE [5] introduce a cube-filling method, which is the only n-ary RE method in *hyper-relational schema*.

However, there are still three main challenges in automated n-ary RE for NKG construction, which remains at a coarse-grained level: **(1) Diversity of NKG schemas.** Previous methods could only perform N-ary RE based on a specific schema, but currently, there is no flexible method that can perform n-ary RE for arbitrary schema with different number of relations. **(2) Determination of the order of entities.** N-ary RE involves more possible entity orders than binary RE, for example, as shown in Figure 2, in a *hyper-relational schema*, there is an order issue regarding which entity is the head entity, tail entity, or auxiliary entity. Previous methods often ignored the joint impact of different entity orders, leading to inaccurate extraction. **(3) Variability of the arity of n-ary RE.** Previous methods usually output a fixed number of entities and are not adept at determining the variable number of entities forming an n-ary relational fact.

To tackle these challenges, we introduce **Text2NKG**, a novel fine-grained n-ary RE framework designed to automate the generation of n-ary relational facts from natural language text for NKG construction. Text2NKG employs a **span-tuple multi-label classification** method, which transforms n-ary RE into a multi-label classification task for span-tuples, including all combinations of entities in the text. Because the number of predicted relation labels corresponds to the chosen NKG schema, Text2NKG is adaptable to all NKG schemas, offering examples with *hyper-relational schema*, *event-based schema*, *role-based schema*, and *hypergraph-based schema*, all of which have broad applications. Moreover, Text2NKG introduces a **hetero-ordered merging** method, considering the probabilities of predicted labels for different entity orders to determine the final entity order. Finally, Text2NKG proposes an **output merging** method, which is used to unsupervisedly derive n-ary relational facts of any number of entities for NKG construction.

In addition, we extend the only n-ary RE benchmark for NKG construction, HyperRED [5], which is in the *hyper-relational schema*, to four NKG schemas. We’ve done sufficient n-ary RE experiments on HyperRED, and the experimental results show that Text2NKG achieves state-of-the-art performance in  $F_1$  scores of hyper-relational extraction. We also compared the results of Text2NKG in the other three schemas to verify applications.

## 2 Related Work

### 2.1 N-ary relational Knowledge Graph

An n-ary relational knowledge graph (NKG) consists of n-ary relational facts, which contain  $n$  entities ( $n \geq 2$ ) and several relations. The n-ary relational facts are necessary and cannot be replaced by combinations of some binary relational facts because we cannot distinguish which binary relations are combined to represent the n-ary relational fact in the whole KG. Therefore, NKG utilizes a schema in every n-ary relational fact locally and a hypergraph representation globally [18].

Firstly, the simplest NKG schema is hypergraph-based. [26] found that over 30% of Freebase [2] entities participate facts with more than two entities, first defined n-ary relations mathematically and used star-to-clique conversion to convert triple-based facts representing n-ary relational facts into the first NKG dataset JF17K in *hypergraph-based schema*  $(r, \{v_i\}_{i=1}^n)$ . [8] proposed FB-AUTO and M-FB15K with the same *hypergraph-based schema*. Secondly, [12] introduced role information for n-ary relational facts and extracted Wikipeople, the first NKG dataset in *role-based schema*  $(\{(k_i, v_i)\}_{i=1}^n)$ , composed of role-value pairs. Thirdly, Wikidata [21], the largest knowledge base, utilizes an NKG schema based on hyper-relation  $(s, r, o, \{(k_i, v_i)\}_{i=1}^{n-2})$ , which adds auxiliary key-value pairs to the main triple. [10] first proposed an NKG dataset in *hyper-relational schema* WD50K. Fourthly, as [11] pointed out, events are also n-ary relational facts. One basic event representation has an event type, a trigger, and several key-value pairs [16]. Regarding the event type as the main relation, the (trigger: value) as one of the key-value pairs, and the arguments as the rest key-value pairs, we can obtain an *event-based NKG schema*  $(r, \{(k_i, v_i)\}_{i=1}^n)$ .

Based on four common NKG schemas, we propose Text2NKG, the first method for extraction of structured n-ary relational facts from natural language text, which improves NKG representation and application.

### 2.2 N-ary Relation Extraction

Relation extraction (RE) is an important step of KG construction, directly affecting the quality, scale, and application of KGs. While most of the current n-ary relation extraction (n-ary RE) for NKG construction depends on manual construction [26, 12, 10] but not automated methods. Most automated RE methods target the extraction of traditional binary relational facts. For example, [22] proposes a table-filling method for binary RE, and [30, 27] propose span-based RE methods with levetated marker and packed levetated marker, respectively.

For automated n-ary RE, some approaches [13, 31] treat n-ary RE in *hypergraph-based schema* or *role-based schema* as a binary classification problem and predict whether the composition of n-ary information in a document is valid or not. However, these methods extract n-ary information in fixed arity, which are not flexible. Moreover, some event extraction methods [16, 17, 9] propose different event trigger and argument extraction techniques, which can achieve n-ary RE in *event-based schema*. Recently, CubeRE [5] proposes an automated n-ary RE method in *hyper-relational schema*, which extends the table-filling extraction method to n-ary RE with cube-filling. However, these methods can only model one of the useful NKG schemas with limited extraction accuracy.

In this paper, we propose the first fine-grained n-ary RE framework Text2NKG for NKG construction in four example schemas, proposing a span-tuple multi-label classification method with hetero-ordered merging and output merging to improve the accuracy of fine-grained n-ary RE extraction in all NKG schemas substantially.

## 3 Preliminaries

**Formulation of NKG.** An NKG  $\mathcal{G} = \{\mathcal{E}, \mathcal{R}, \mathcal{F}\}$  consists of an entity set  $\mathcal{E}$ , a relation set  $\mathcal{R}$ , and an n-ary fact ( $n \geq 2$ ) set  $\mathcal{F}$ . Each n-ary fact  $f^n \in \mathcal{F}$  consists of entities  $\in \mathcal{E}$  and relations  $\in \mathcal{R}$ . For hyper-relational schema [20]:  $f_{hr}^n = (e_1, r_1, e_2, \{r_{i-1}, e_i\}_{i=3}^n)$  where  $\{e_i\}_{i=1}^n \in \mathcal{E}$ ,  $\{r_i\}_{i=1}^{n-1} \in \mathcal{R}$ . For event-based schema [16]:  $f_{ev}^n = (r_1, \{r_{i+1}, e_i\}_{i=1}^n)$ , where  $\{e_i\}_{i=1}^n \in \mathcal{E}$ ,  $\{r_i\}_{i=1}^{n+1} \in \mathcal{R}$ . For role-based schema [12]:  $f_{ro}^n = (\{r_i, e_i\}_{i=1}^n)$ , where  $\{e_i\}_{i=1}^n \in \mathcal{E}$ ,  $\{r_i\}_{i=1}^n \in \mathcal{R}$ . For hypergraph-based schema [26]:  $f_{hg}^n = (r_1, \{e_i\}_{i=1}^n)$ , where  $\{e_i\}_{i=1}^n \in \mathcal{E}$ ,  $r_1 \in \mathcal{R}$ .

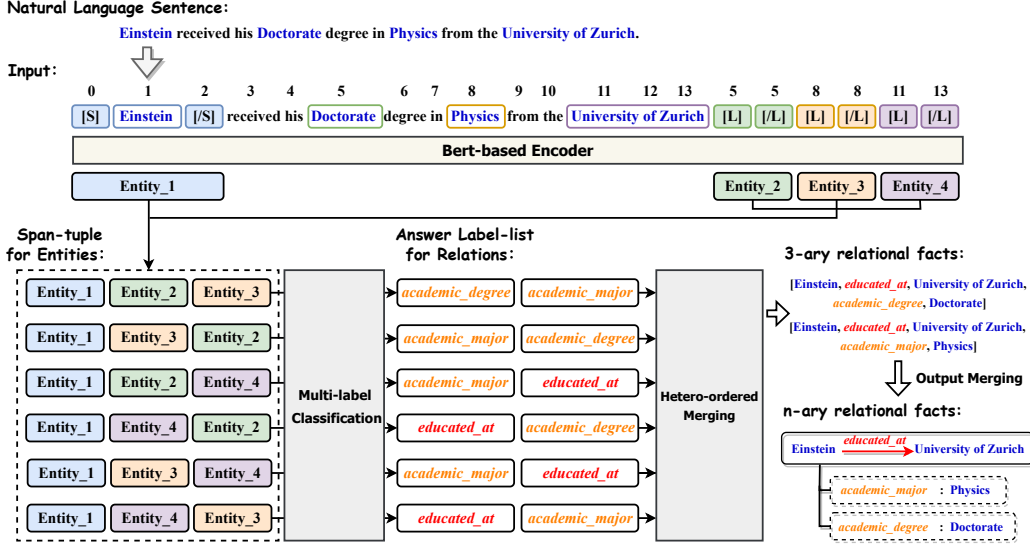


Figure 3: An overview of Text2NKG extracting n-ary relation facts from a natural language sentence in hyper-relational NKG schema for an example.

**Problem Definition.** Given an input sentence with  $l$  words  $s = \{w_1, w_2, \dots, w_l\}$ , an entity  $e$  is a consecutive span of words:  $e = \{w_p, w_{p+1}, \dots, w_q\} \in \mathcal{E}_s$ , where  $p, q \in \{1, \dots, l\}$ , and  $\mathcal{E}_s = \{e_j\}_{j=1}^m$  is the entity set of all  $m$  entities in the sentence. The output of n-ary relation extraction,  $R(\cdot)$ , is a set of n-ary relational facts  $\mathcal{F}_s$  in given NKG schema in  $\{f_{hr}^n, f_{ev}^n, f_{ro}^n, f_{hg}^n\}$ . Specifically, each n-ary relational fact  $f^n \in \mathcal{F}_s$  is extracted by multi-label classification of one of the ordered span-tuple for  $n$  entities  $[e_i]_{i=1}^n \in \mathcal{E}_s$ , forming an answer label-list for  $n_r$  relations  $[r_i]_{i=1}^{n_r} \in \mathcal{R}$ , where  $n$  is the arity of the extracted n-ary relational fact, and  $n_r$  is the number of answer relations in the fact, which is determined by the given NKG schema:  $R([e_i]_{i=1}^n) = [r_i]_{i=1}^{n-1}$ , when  $f^n = f_{hr}^n$ ,  $R([e_i]_{i=1}^n) = [r_i]_{i=1}^{n+1}$  when  $f^n = f_{ev}^n$ ,  $R([e_i]_{i=1}^n) = [r_i]_{i=1}^n$  when  $f^n = f_{ro}^n$ , and  $R([e_i]_{i=1}^n) = [r_1]$  when  $f^n = f_{hg}^n$ .

## 4 Methodology

In this section, we first introduce the overview of the Text2NKG framework, followed by the span-tuple multi-label classification, training strategy, hetero-ordered merging, and output merging.

### 4.1 Overview of Text2NKG

Text2NKG is a fine-grained n-ary relation extraction framework built for n-ary relational knowledge graph (NKG) construction. The input to Text2NKG is natural language text tokens labeled with entity span in sentence units. First, inspired by [27], Text2NKG encodes the entities using BERT-based Encoder [6] with a packaged leviatated marker for embedding. Then each arrangement of ordered span-tuple with three entity embeddings will be classified with multiple labels, and the framework will be learned by the weighted cross-entropy with a null-label bias. In the decoding stage, in order to filter the n-ary relational facts whose entity compositions have isomorphic hetero-ordered characteristics, Text2NKG proposes a hetero-ordered merging strategy to merge the label probabilities of  $3! = 6$  arrangement cases of span-tuples composed of the same entities and filter out the output 3-ary relational facts existing non-conforming relations. Finally, Text2NKG combines the output 3-ary relational facts to form the final n-ary relational facts with output merging.

### 4.2 Span-tuple Multi-label Classification

For the given sentence token  $s = \{w_1, w_2, \dots, w_l\}$  and the set of entities  $\mathcal{E}_s$ , in order to perform fine-grained n-ary RE, we need first to encode a span-tuple  $(e_1, e_2, e_3)$  consisting of every arrangement of three ordered entities, where  $e_1, e_2, e_3 \in \mathcal{E}_s$ . Due to the high time complexity of training every

span-tuple as one training item, inspired by [27], we achieve the reduction of training items by using packed levitated markers that pack one training item with each entity in  $\mathcal{E}_s$  separately. Specifically, in each packed training item, a pair of solid tokens, [S] and [/S], are added before and after the packed entity  $e_S = \{w_{p_S}, \dots, w_{q_S}\}$ , and  $(|\mathcal{E}_s| - 1)$  pairs of levitated markers, [L] and [/L], according to other entities in  $\mathcal{E}_s$ , are added with the same position embeddings as the beginning and end of their corresponding entities span  $e_{L_i} = \{w_{p_{L_i}}, \dots, w_{q_{L_i}}\}$  to form the input token  $\mathbf{X}$ :

$$\mathbf{X} = \{w_1, \dots, [S], w_{p_S}, \dots, w_{q_S}, [/S], \dots, w_{p_{L_i}} \cup [L], \dots, w_{q_{L_i}} \cup [/L], \dots, w_l\}. \quad (1)$$

We encode such token by the BERT-based pre-trained model encoder [6]:

$$\{h_1, h_2, \dots, h_t\} = \text{BERT}(\mathbf{X}), \quad (2)$$

where  $t = |\mathbf{X}|$  is the input token length,  $\{h_i\}_{i=1}^t \in \mathbb{R}^d$ , and  $d$  is embedding size.

There are several span-tuples  $(A, B, C)$  in a training item. The embedding of first entity  $h_A \in \mathbb{R}^{2d}$  in the span-tuple is obtained by concat embedding of the solid markers, [S] and [/S], and the embeddings of second and third entities  $h_B, h_C \in \mathbb{R}^{2d}$  are obtained by concat embeddings of levitated markers, [L] and [/L] with all  $A_{m-1}^2$  arrangement of any other two entities in  $\mathcal{E}_s$ . Thus, we obtain the embedding representation of the three entities to form  $A_{m-1}^2$  span-tuples in one training item. Therefore, every input sentence contains  $m$  training items with  $m A_{m-1}^2 = A_m^3$  span-tuples for any ordered arrangement of three entities.

We then define  $n_r$  linear classifiers, each of which consists of 3 feedforward neural networks  $\{\text{FNN}_i^k\}_{i=1}^{n_r}, k = 1, 2, 3$ , to classify the span-tuples for multiple-label classification. Each classifier targets the prediction of one relation  $r_i$ , thus obtaining a probability lists  $(\mathbf{P}_i)_{i=1}^{n_r}$  with all relations in given relation set  $\mathcal{R}$  plus a null-label:

$$\mathbf{P}_i = \text{FNN}_i^1(h_A) + \text{FNN}_i^2(h_B) + \text{FNN}_i^3(h_C), \quad (3)$$

where  $\text{FNN}_i^k \in \mathbb{R}^{2d \times (|\mathcal{R}|+1)}$ , and  $\mathbf{P}_i \in \mathbb{R}^{(|\mathcal{R}|+1)}$ .

### 4.3 Training Strategy

To train the  $n_r$  classifiers for each relation prediction more accurately, we propose a data augmentation strategy for span-tuples. Taking the *hyper-relational schema* as an example, given a hyper-relational fact  $(A, r_1, B, r_2, C)$ , we consider swapping the head and tail entities, and changing the main relation to its inverse  $(B, r_1^{-1}, A, r_2, C)$ , as well as swapping the tail entities with auxiliary values, and the main relation with the auxiliary key  $(A, r_2, C, r_1, B)$ , also as labeled training span-tuple cases. Thus  $R_{hr}(A, B, C) = (r_1, r_2)$  can be augmented with  $3! = 6$  orders of span-tuples:

$$\begin{cases} R_{hr}(A, B, C) = (r_1, r_2), \\ R_{hr}(B, A, C) = (r_1^{-1}, r_2), \\ R_{hr}(A, C, B) = (r_2, r_1), \\ R_{hr}(B, C, A) = (r_2, r_1^{-1}), \\ R_{hr}(C, A, B) = (r_2^{-1}, r_1), \\ R_{hr}(C, B, A) = (r_1, r_2^{-1}). \end{cases} \quad (4)$$

For other schemas, we can also obtain 6 fully-arranged cases of labeled span-tuples in a similar way, as described in Appendix A. If no n-ary relational fact exists between the three entities of span-tuples, then relation labels are set as null-label.

Since most cases of span-tuple are null-label, we set a weight hyperparameter  $\alpha \in (0, 1]$  between the null-label and other labels to balance the learning of the null-label. We jointly trained the  $n_r$  classifiers for each relations by cross-entropy loss  $\mathcal{L}$  with a null-label weight bias  $\mathbf{W}_\alpha$ :

$$\mathcal{L} = - \sum_{i=1}^{n_r} \mathbf{W}_\alpha \log \left( \frac{\exp(\mathbf{P}_i[r_i])}{\sum_{j=1}^{|\mathcal{R}|+1} \exp(\mathbf{P}_i[j])} \right), \quad (5)$$

where  $\mathbf{W}_\alpha = [\alpha, 1.0, 1.0, \dots, 1.0] \in \mathbb{R}^{(|\mathcal{R}|+1)}$ .

Dataset	#Ent	#R_hr	#R_ev	#R_ro	#R_hg	All		Train		Dev		Test	
						#Sentence	#Fact	#Sentence	#Fact	#Sentence	#Fact	#Sentence	#Fact
HyperRED	40,293	106	232	168	62	44,840	45,994	39,840	39,978	1,000	1,220	4,000	4,796

Table 1: Dataset statistics, where the columns indicate the number of entities, relations with four schema, sentences and n-ary relational facts in all sets, train set, dev set, and test set, respectively.

#### 4.4 Hetero-ordered Merging

In the decoding stage, since Text2NKG labels all 6 different arrangement of the same entity composition, we design a hetero-ordered merging strategy to merge the corresponding labels of these 6 hetero-ordered span-tuples into one to generate non-repetitive n-ary relational facts unsupervisedly. For *hyper-relational schema* ( $n_r = 2$ ), we combine the predicted probabilities of two labels  $\mathbf{P}_1, \mathbf{P}_2$  in 6 orders to  $(A, B, C)$  order as follows:

$$\begin{cases} \mathbf{P}_1 = \mathbf{P}_1^{(ABC)} + I(\mathbf{P}_1^{(BAC)}) + \mathbf{P}_2^{(ACB)} \\ \quad + I(\mathbf{P}_2^{(BCA)}) + \mathbf{P}_2^{(CAB)} + \mathbf{P}_1^{(CBA)}, \\ \mathbf{P}_2 = \mathbf{P}_2^{(ABC)} + \mathbf{P}_2^{(BAC)} + \mathbf{P}_1^{(ACB)} \\ \quad + \mathbf{P}_1^{(BCA)} + I(\mathbf{P}_1^{(CAB)}) + I(\mathbf{P}_2^{(CBA)}), \end{cases} \quad (6)$$

where  $I()$  is a function for swapping the predicted probability of relations and the corresponding inverse relations. Then, we take the maximum probability to obtain labels  $r_1, r_2$ , forming a 3-ary relational fact  $(A, r_1, B, r_2, C)$  and filter it out if there are null-label in  $(r_1, r_2)$ . If there are inverse relation labels in  $(r_1, r_2)$ , we can also transform the order of entities and relations as equation 4. For *event-based schema*, *role-based schema*, and *hypergraph-based schema*, all can be generated by hetero-ordered merging according to this idea, as shown in Appendix B.

#### 4.5 Output Merging

After hetero-ordered merging, we merge the output 3-ary relational facts to form higher-arity facts, with *hyper-relational schema* based on the same main triple, *event-based schema* based on the same main relation (event type), *role-based schema* based on the same key-value pairs, and *hypergraph-based schema* based on the same hyperedge relation. This way, we can unsupervisedly obtain n-ary relational facts with dynamic number of arity numbers for NKG construction. More details are discussed in Appendix G.2 and Appendix G.3.

## 5 Experiments

This section presents the experimental setup, results, and analysis. We answer the following research questions (RQs): **RQ1**: Does Text2NKG outperform other n-ary RE methods? **RQ2**: Whether Text2NKG can cover NKG construction for various schemas? **RQ3**: Does the main components of Text2NKG work? **RQ4**: How does the null-label bias hyperparameter in Text2NKG affect performance? **RQ5**: Can Text2NKG get complete n-ary relational facts in different arity? **RQ6**: How is Text2NKG’s computational efficiency? **RQ7**: How does Text2NKG perform in specific case study? **RQ8**: What is the future development of Text2NKG in the era of large language models?

### 5.1 Experimental Setup

**Datasets.** The existing fine-grained n-ary RE dataset is **HyperRED** [5] only in *hyper-relational schema* with annotated extracted entities. Therefore, we expand the HyperRED dataset to four schemas as standard fine-grained n-ary RE benchmarks and conduct experiments on them. The statistics of the HyperRED with four schemas are shown in Table 1 and the construction detail is in Appendix C.

**Baselines.** We compare Text2NKG against **Generative Baseline** [14], **Pipeline Baseline** [24], and **CubeRE** [5] in fine-grained n-ary RE task of *hyper-relational schema*. For n-ary RE in the other three schemas, we compared Text2NKG with event extraction models such as **Text2Event** [16], **UIE** [17],

Model	PLM	<i>hyper-relational schema / Dev</i>			<i>hyper-relational schema / Test</i>		
		Precision	Recall	$F_1$	Precision	Recall	$F_1$
<b>Unsupervised Method</b>							
ChatGPT	gpt-3.5-turbo	12.0583	11.2764	11.6542	11.4021	10.9134	11.1524
GPT-4	gpt-4	15.7324	15.2377	15.4811	15.8187	15.4824	15.6487
<b>Supervised Method</b>							
Generative Baseline		63.79 $\pm$ 0.27	59.94 $\pm$ 0.68	61.80 $\pm$ 0.37	64.60 $\pm$ 0.47	59.67 $\pm$ 0.35	62.03 $\pm$ 0.21
Pipeline Baseline		69.23 $\pm$ 0.30	58.21 $\pm$ 0.57	63.24 $\pm$ 0.44	69.00 $\pm$ 0.48	57.55 $\pm$ 0.19	62.75 $\pm$ 0.29
CubeRE		66.14 $\pm$ 0.88	64.39 $\pm$ 1.23	65.23 $\pm$ 0.82	65.82 $\pm$ 0.84	64.28 $\pm$ 0.25	65.04 $\pm$ 0.29
Text2NKG w/o DA	BERT-base (110M)	76.02 $\pm$ 0.50	72.28 $\pm$ 0.68	74.10 $\pm$ 0.55	73.55 $\pm$ 0.81	70.63 $\pm$ 1.40	72.06 $\pm$ 0.34
Text2NKG w/o $\alpha$		88.77 $\pm$ 0.85	78.39 $\pm$ 0.47	83.26 $\pm$ 0.70	88.09 $\pm$ 0.69	76.64 $\pm$ 0.45	81.97 $\pm$ 0.58
Text2NKG w/o HM		61.74 $\pm$ 0.34	76.97 $\pm$ 0.44	68.52 $\pm$ 0.69	61.07 $\pm$ 0.73	76.16 $\pm$ 0.59	67.72 $\pm$ 0.48
Text2NKG (ours)		<b>91.26 <math>\pm</math> 0.69</b>	<b>79.36 <math>\pm</math> 0.51</b>	<b>84.89 <math>\pm</math> 0.44</b>	<b>90.77 <math>\pm</math> 0.60</b>	<b>77.53 <math>\pm</math> 0.32</b>	<b>83.63 <math>\pm</math> 0.63</b>
Generative Baseline		67.08 $\pm$ 0.49	65.73 $\pm$ 0.78	66.40 $\pm$ 0.47	67.17 $\pm$ 0.40	64.56 $\pm$ 0.58	65.84 $\pm$ 0.25
Pipeline Baseline	BERT-large (340M)	70.58 $\pm$ 0.78	66.58 $\pm$ 0.66	68.52 $\pm$ 0.32	69.21 $\pm$ 0.55	64.27 $\pm$ 0.24	66.65 $\pm$ 0.28
CubeRE		68.75 $\pm$ 0.82	68.88 $\pm$ 1.03	68.81 $\pm$ 0.46	66.39 $\pm$ 0.96	67.12 $\pm$ 0.69	66.75 $\pm$ 0.28
Text2NKG (ours)		<b>91.90 <math>\pm</math> 0.79</b>	<b>79.43 <math>\pm</math> 0.42</b>	<b>85.21 <math>\pm</math> 0.69</b>	<b>91.06 <math>\pm</math> 0.81</b>	<b>77.64 <math>\pm</math> 0.46</b>	<b>83.81 <math>\pm</math> 0.54</b>

Table 2: Comparison of Text2NKG with other baselines in the hyper-relational extraction on HyperRED. Results of the supervised baseline models are mainly taken from the original paper [5]. The best results in each metric are in **bold**.

Model	PLM	<i>event-based schema</i>			<i>role-based schema</i>			<i>hypergraph-based schema</i>		
		Precision	Recall	$F_1$	Precision	Recall	$F_1$	Precision	Recall	$F_1$
<b>Unsupervised Method</b>										
ChatGPT	gpt-3.5-turbo	10.4678	11.1628	10.8041	11.4387	10.4203	10.9058	11.2998	11.7852	11.5373
GPT-4	gpt-4	13.3681	14.6701	13.9888	13.6397	12.5355	13.0643	13.0907	13.6701	13.3741
<b>Supervised Method</b>										
Text2Event		73.94 $\pm$ 0.76	70.56 $\pm$ 0.58	72.21 $\pm$ 1.25	72.73 $\pm$ 0.79	68.45 $\pm$ 1.34	70.52 $\pm$ 0.62	73.68 $\pm$ 0.88	70.37 $\pm$ 0.51	71.98 $\pm$ 0.92
UIE	T5-base (220M)	76.51 $\pm$ 0.28	73.02 $\pm$ 0.66	74.72 $\pm$ 0.18	72.17 $\pm$ 0.29	69.84 $\pm$ 0.11	70.98 $\pm$ 0.31	72.03 $\pm$ 0.41	68.74 $\pm$ 0.13	70.34 $\pm$ 1.07
LasUIE		79.62 $\pm$ 0.27	78.04 $\pm$ 0.75	78.82 $\pm$ 0.26	77.01 $\pm$ 0.20	74.26 $\pm$ 0.25	75.61 $\pm$ 0.24	76.21 $\pm$ 0.07	73.75 $\pm$ 0.17	74.96 $\pm$ 0.42
Text2NKG	BERT-base (110M)	<b>86.20 <math>\pm</math> 0.57</b>	<b>79.25 <math>\pm</math> 0.33</b>	<b>82.58 <math>\pm</math> 0.20</b>	<b>86.72 <math>\pm</math> 0.80</b>	<b>78.94 <math>\pm</math> 0.59</b>	<b>82.64 <math>\pm</math> 0.38</b>	<b>83.53 <math>\pm</math> 1.18</b>	<b>86.59 <math>\pm</math> 0.38</b>	<b>85.03 <math>\pm</math> 0.86</b>
Text2Event		75.58 $\pm$ 0.53	72.39 $\pm$ 0.82	73.97 $\pm$ 1.19	73.21 $\pm$ 0.45	70.85 $\pm$ 0.67	72.01 $\pm$ 0.31	75.28 $\pm$ 0.93	72.73 $\pm$ 1.07	73.98 $\pm$ 0.49
UIE	T5-large (770M)	79.38 $\pm$ 0.28	74.69 $\pm$ 0.61	76.96 $\pm$ 0.95	74.47 $\pm$ 1.42	71.84 $\pm$ 0.77	73.14 $\pm$ 0.38	74.57 $\pm$ 0.64	71.93 $\pm$ 0.86	73.22 $\pm$ 0.19
LasUIE		81.29 $\pm$ 0.83	79.54 $\pm$ 0.26	80.40 $\pm$ 0.65	79.37 $\pm$ 0.92	76.63 $\pm$ 0.44	77.97 $\pm$ 0.76	77.49 $\pm$ 0.35	74.96 $\pm$ 0.60	76.20 $\pm$ 0.87
Text2NKG	BERT-large (340M)	<b>88.47 <math>\pm</math> 0.95</b>	<b>80.30 <math>\pm</math> 0.75</b>	<b>84.19 <math>\pm</math> 1.29</b>	<b>86.87 <math>\pm</math> 0.87</b>	<b>80.86 <math>\pm</math> 0.29</b>	<b>83.76 <math>\pm</math> 1.17</b>	<b>85.06 <math>\pm</math> 0.33</b>	<b>86.72 <math>\pm</math> 0.36</b>	<b>85.89 <math>\pm</math> 0.69</b>

Table 3: Comparison of Text2NKG with other baselines in the n-ary RE in event-based, role-based, and *hypergraph-based schemas* on HyperRED. The best results in each metric are in **bold**.

and **LasUIE** [9]. Furthermore, we utilized different prompts to test the currently most advanced large-scale pre-trained language models **ChatGPT** [25] and **GPT-4** [19] in an unsupervised manner, specifically for the extraction performance across the four schemas. The detailed baseline settings can be found in Appendix D.

**Ablations.** To evaluate the significance of Text2NKG’s three main components, data augmentation (DA), null-label weight hyperparameter ( $\alpha$ ), and hetero-ordered merging (HM), we obtain three simplified model variants by removing any one component (**Text2NKG w/o DA**, **Text2NKG w/o  $\alpha$** , and **Text2NKG w/o HM**) for comparison.

**Evaluation Metrics.** We use the  $F_1$  score with precision and recall to evaluate the dev set and the test set. For a predicted n-ary relational fact to be considered correct, the entire fact must match the ground facts completely.

**Hyperparameters and Environment.** We train 10 epochs on HyperRED using the Adam optimizer. All experiments were done on a single NVIDIA A100 GPU, and all experimental results were derived by averaging 5 random seed experiments. Appendix E shows Text2NKG’s optimal hyperparameter settings. Appendix F shows training details.

## 5.2 Main Results (RQ1)

The experimental results of proposed Text2NKG and other baselines with both BERT-base and BERT-large encoders can be found in Table 2 for the fine-grained n-ary RE in *hyper-relational schema*. We can observe that Text2NKG shows a significant improvement over the existing optimal model CubeRE on both the dev and test datasets of HyperRED. The  $F_1$  score is improved by 19.66 percentage points in the dev set and 18.60 percentage points in the test set with the same BERT-base encoder, and 16.40 percentage points in the dev set and 17.06 percentage points in the test set with

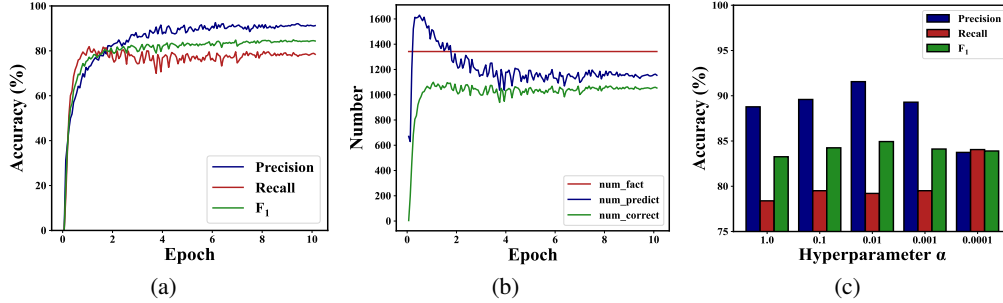


Figure 4: (a) Precision, Recall, and  $F_1$  changes in the dev set during the training of Text2NKG. (b) The changes of the number of true facts, the number of predicted facts, and the number of predicted accurate facts during the training of Text2NKG. (c) Precision, Recall, and  $F_1$  results on different null-label hyperparameter ( $\alpha$ ) settings.

the same BERT-large encoder, reflecting Text2NKG’s excellent performance. Figure 4(a) and 4(b) intuitively show the changes of evaluation metrics and answers of facts in the dev set during the training of Text2NKG. It is worth noting that Text2NKG exceeds 90% in precision accuracy, which proves that the model can obtain very accurate n-ary relational facts and provides a good guarantee for the quality of fine-grained NKG construction.

### 5.3 Results on Various NKG Schemas (RQ2)

As shown in Table 3, besides *hyper-relational schema*, Text2NKG also accomplishes the tasks of fine-grained n-ary RE in three other different NKG schemas on HyperRED, which demonstrates good utility. In the added tasks of n-ary RE for event-based, role-based, and *hypergraph-based schemas*, since no model has done similar experiments at present, we used event extraction or unified extraction methods such as Text2Event [16], UIE [17], and LasUIE [9] for comparison. Text2NKG still works best in these schemas, which demonstrates good versatility.

### 5.4 Ablation Study (RQ3)

Data augmentation (DA), null-label weight hyperparameter ( $\alpha$ ), and hetero-ordered merging (HM) are the three main components of Text2NKG. For the different Text2NKG variants as shown in Table 2, DA,  $\alpha$ , and HM all contribute to the accurate results of our complete model. By comparing the differences, we find that HM is most effective by combining the probabilities of labels of different orders, followed by DA and  $\alpha$ .

### 5.5 Analysis of Null-label Weight Hyperparameters (RQ4)

We compared the effect for different null-label weight hyperparameters ( $\alpha$ ). As shown in Figure 4(c), the larger the  $\alpha$ , the greater the learning weight of null-label compared with other labels, the more relations are predicted as null-label. After filtering out the facts having null-label, fewer facts are extracted, so the precision is generally higher, and the recall is generally lower. The smaller the  $\alpha$ , the more relations are predicted as non-null labels, thus extracting more n-ary relation facts, so the recall is generally higher, and the precision is generally lower. Comparing the results of  $F_1$  values for different  $\alpha$ , it is found that  $\alpha = 0.01$  works best, which can be adjusted in practice according to specific needs to obtain the best results.

### 5.6 Analysis of N-ary Relation Extraction in Different Arity (RQ5)

Figure 5(a) shows the number of n-ary relational facts extracted after output merging and the number of the answer facts in different arity during training of Text2NKG on the dev set. We find that, as the training proceeds, the final output of Text2NKG converges to the correct answer in terms of the number of complete n-ary relational facts in each arity, achieving implementation of n-ary RE in indefinite arity unsupervised, with good scalability.



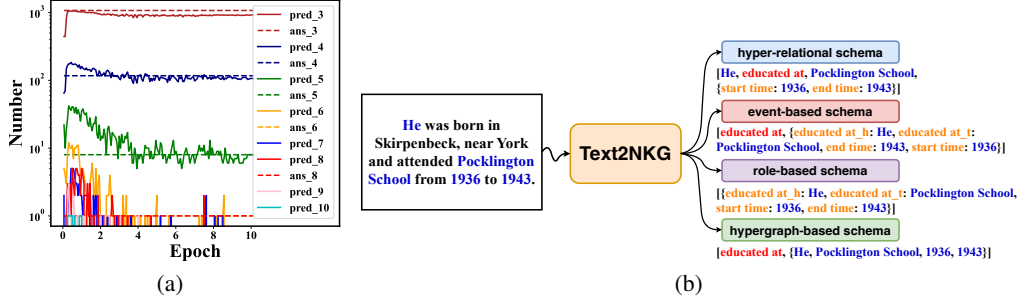


Figure 5: (a) The changes of the number of extracted n-ary RE in different arity, where "pred\_n" represents the number of extracted n-ary facts with different arities by Text2NKG, and "ans\_n" represents the ground truth. (b) Case study of Text2NKG's n-ary relation extraction in four schemas on HyperRED.

## 5.7 Computational Efficiency (RQ6)

As mentioned in Section 4.2, the main computational consumption of Text2NKG is selecting every span-tuple of three ordered entities to encode them and get the classified labels in multiple-label classification part. If we adopt an traversal approach with each span-tuple in one training items, the time complexity will be  $O(m^3)$ . To reduce the high time complexity of training every span-tuple as one training item, Text2NKG uses packed levitated markers that pack one training item with each entity in  $\mathcal{E}_s$  separately. We obtain the embedding representation of the three entities to form  $A_{m-1}^2$  span-tuples in one training item. Every input sentence contains  $m$  training items with  $mA_{m-1}^2 = A_m^3$  span-tuples for any ordered arrangement of three entities for multiple-label classification. Therefore, the time complexity decreased from  $O(m^3)$  to  $O(m)$ .

## 5.8 Case Study (RQ7)

Figure 5(b) shows a case study of n-ary RE by a trained Text2NKG. For a sentence, "He was born in Skirpenbeck, near York and attended Pocklin.", four structured n-ary RE can be obtained by Text2NKG according to the requirements. Taking the *hyper-relational schema* for an example, Text2NKG can successfully extract one n-ary relational fact consisting of a main triple [He, educated at, Pocklington], and two auxiliary key-value pairs {start time:1936}, {end time:1943}. This intuitively validates the practical performance of Text2NKG on fine-grained n-ary RE to better contribute to NKG construction.

## 5.9 Comparison with ChatGPT (RQ8)

As shown in Table 2 and Table 3, we compared the extraction effects under four NKG schemas of the supervised Text2NKG with the unsupervised ChatGPT and GPT-4. We found that these large language models cannot accurately distinguish the closely related relations in the fine-grained NKG relation repository, resulting in their F1 scores ranging around 10%-15%, which is much lower than the performance of Text2NKG. On the other hand, the limitation of Text2NKG is that its performance is confined within the realm of supervised training. Therefore, in future improvements and practical applications, we suggest combining small supervised models with large unsupervised models to balance solving the cold-start and fine-grained extraction, which is detailed in Appendix G.1.

## 6 Conclusion

In this paper, we introduce Text2NKG, a novel framework designed for fine-grained n-ary relation extraction (RE) aimed at constructing N-ary Knowledge Graphs (NKGs). Our extensive experiments demonstrate that Text2NKG outperforms all existing baseline models across a wide range of fine-grained n-ary RE tasks. Notably, it excels in four distinct schema types: hyper-relational, event-based, role-based, and hypergraph-based. Furthermore, we have extended the HyperRED dataset, transforming it into a comprehensive fine-grained n-ary RE benchmark that supports all four schemas.

## Acknowledgments

This work is supported by the National Science Foundation of China (Grant No. 62176026, Grant No. 62406036, Grant No. 62473271, and Grant No. 62076035). This work is also supported by the SMP-Zhipu.AI Large Model Cross-Disciplinary Fund, the BUPT Excellent Ph.D. Students Foundation (No. CX2023133), the BUPT Innovation and Entrepreneurship Support Program (No. 2024-YC-A091 and No. 2024-YC-T022), and the Engineering Research Center of Information Networks, Ministry of Education.

## References

- [1] Erik Arakelyan, Daniel Daza, Pasquale Minervini, and Michael Cochez. Complex query answering with neural link predictors. In *International Conference on Learning Representations*, 2021.
- [2] Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. Freebase: A collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*, SIGMOD '08, page 1247–1250, New York, NY, USA, 2008. Association for Computing Machinery.
- [3] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013.
- [4] Xuelu Chen, Ziniu Hu, and Yizhou Sun. Fuzzy logic based logical query answering on knowledge graphs. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(4):3939–3948, Jun. 2022.
- [5] Yew Ken Chia, Lidong Bing, Sharifah Mahani Aljunied, Luo Si, and Soujanya Poria. A dataset for hyper-relational extraction and a cube-filling approach. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10114–10133, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- [7] Xin Dong, Evgeniy Gabrilovich, Jeremy Heitz, Wilko Horn, Ni Lao, Kevin Murphy, Thomas Strohmman, Shaohua Sun, and Wei Zhang. Knowledge vault: A web-scale approach to probabilistic knowledge fusion. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '14, page 601–610, New York, NY, USA, 2014. Association for Computing Machinery.
- [8] Bahare Fatemi, Perouz Taslakian, David Vazquez, and David Poole. Knowledge hypergraphs: Prediction beyond binary relations. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*, IJCAI'20, 2021.
- [9] Hao Fei, Shengqiong Wu, Jingye Li, Bobo Li, Fei Li, Libo Qin, Meishan Zhang, Min Zhang, and Tat-Seng Chua. Lasuie: Unifying information extraction with latent adaptive structure-aware generative language model. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 15460–15475. Curran Associates, Inc., 2022.
- [10] Mikhail Galkin, Priyansh Trivedi, Gaurav Maheshwari, Ricardo Usbeck, and Jens Lehmann. Message passing for hyper-relational knowledge graphs. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7346–7359, Online, November 2020. Association for Computational Linguistics.

- [11] Saiping Guan, Xueqi Cheng, Long Bai, Fujun Zhang, Zixuan Li, Yutao Zeng, Xiaolong Jin, and Jiafeng Guo. What is event knowledge graph: A survey. *IEEE Transactions on Knowledge and Data Engineering*, pages 1–20, 2022.
- [12] Saiping Guan, Xiaolong Jin, Yuanzhuo Wang, and Xueqi Cheng. Link prediction on n-ary relational data. In *The World Wide Web Conference, WWW '19*, page 583–593, New York, NY, USA, 2019. Association for Computing Machinery.
- [13] Robin Jia, Cliff Wong, and Hoifung Poon. Document-level n-ary relation extraction with multiscale representation learning. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3693–3704, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- [14] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online, July 2020. Association for Computational Linguistics.
- [15] Yu Liu, Quanming Yao, and Yong Li. Role-aware modeling for n-ary relational knowledge bases. In *Proceedings of the Web Conference 2021, WWW '21*, page 2660–2671, New York, NY, USA, 2021. Association for Computing Machinery.
- [16] Yaojie Lu, Hongyu Lin, Jin Xu, Xianpei Han, Jialong Tang, Annan Li, Le Sun, Meng Liao, and Shaoyi Chen. Text2Event: Controllable sequence-to-structure generation for end-to-end event extraction. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 2795–2806, Online, August 2021. Association for Computational Linguistics.
- [17] Yaojie Lu, Qing Liu, Dai Dai, Xinyan Xiao, Hongyu Lin, Xianpei Han, Le Sun, and Hua Wu. Unified structure generation for universal information extraction. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5755–5772, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [18] Haoran Luo, Haihong E, Yuhao Yang, Yikai Guo, Mingzhi Sun, Tianyu Yao, Zichen Tang, Kaiyang Wan, Meina Song, and Wei Lin. HAHE: Hierarchical attention for hyper-relational knowledge graphs in global and local level. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8095–8107, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [19] OpenAI. Gpt-4 technical report, 2023.
- [20] Paolo Rosso, Dingqi Yang, and Philippe Cudré-Mauroux. Beyond triplets: Hyper-relational knowledge graph embedding for link prediction. In *Proceedings of The Web Conference 2020, WWW '20*, page 1885–1896, New York, NY, USA, 2020. Association for Computing Machinery.
- [21] Denny Vrandečić and Markus Krötzsch. Wikidata: A free collaborative knowledgebase. *Commun. ACM*, 57(10):78–85, sep 2014.
- [22] Jue Wang and Wei Lu. Two are better than one: Joint entity and relation extraction with table-sequence encoders. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1706–1721, Online, November 2020. Association for Computational Linguistics.
- [23] Quan Wang, Haifeng Wang, Yajuan Lyu, and Yong Zhu. Link prediction on n-ary relational facts: A graph-based approach. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 396–407, Online, August 2021. Association for Computational Linguistics.

- [24] Yijun Wang, Changzhi Sun, Yuanbin Wu, Hao Zhou, Lei Li, and Junchi Yan. UniRE: A unified label space for entity relation extraction. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 220–231, Online, August 2021. Association for Computational Linguistics.
- [25] Xiang Wei, Xingyu Cui, Ning Cheng, Xiaobin Wang, Xin Zhang, Shen Huang, Pengjun Xie, Jinan Xu, Yufeng Chen, Meishan Zhang, Yong Jiang, and Wenjuan Han. Zero-shot information extraction via chatting with chatgpt, 2023.
- [26] Jianfeng Wen, Jianxin Li, Yongyi Mao, Shini Chen, and Richong Zhang. On the representation and embedding of knowledge bases beyond binary relations. In Subbarao Kambhampati, editor, *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI 2016, New York, NY, USA, 9-15 July 2016*, pages 1300–1307. IJCAI/AAAI Press, 2016.
- [27] Deming Ye, Yankai Lin, Peng Li, and Maosong Sun. Packed levitated marker for entity and relation extraction. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4904–4917, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [28] Scott Wen-tau Yih, Ming-Wei Chang, Xiaodong He, and Jianfeng Gao. Semantic parsing via staged query graph generation: Question answering with knowledge base. In *Proceedings of the Joint Conference of the 53rd Annual Meeting of the ACL and the 7th International Joint Conference on Natural Language Processing of the AFNLP*. ACL - Association for Computational Linguistics, July 2015. Outstanding Paper Award.
- [29] Fuzheng Zhang, Nicholas Jing Yuan, Defu Lian, Xing Xie, and Wei-Ying Ma. Collaborative knowledge base embedding for recommender systems. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16*, page 353–362, New York, NY, USA, 2016. Association for Computing Machinery.
- [30] Zexuan Zhong and Danqi Chen. A frustratingly easy approach for entity and relation extraction. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 50–61, Online, June 2021. Association for Computational Linguistics.
- [31] Yuchen Zhuang, Yinghao Li, Junyang Zhang, Yue Yu, Yingjun Mou, Xiang Chen, Le Song, and Chao Zhang. ReSel: N-ary relation extraction from scientific text and tables by learning to retrieve and select. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 730–744, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.

## Appendix

### A Supplement to Data Augmentation

In addition to the *hyper-relational schema*, the data augmentation strategies for other schemas are as follows:

For *event-based schema*, given an event-based fact  $(r_1, r_2, A, r_3, B, r_4, C)$ , we consider keeping the main relation  $r_1$  unchanged, and swapping other key-value pairs,  $\{r_2, A\}$ ,  $\{r_3, B\}$ , and  $\{r_4, C\}$ , positionally, also as labeled training span-tuple cases. Thus  $R_{ev}(A, B, C) = (r_1, r_2, r_3, r_4)$  can be augmented with 6 orders of span-tuples:

$$\left\{ \begin{array}{l} R_{ev}(A, B, C) = (r_1, r_2, r_3, r_4), \\ R_{ev}(B, A, C) = (r_1, r_3, r_2, r_4), \\ R_{ev}(A, C, B) = (r_1, r_2, r_4, r_3), \\ R_{ev}(B, C, A) = (r_1, r_3, r_4, r_2), \\ R_{ev}(C, A, B) = (r_1, r_4, r_2, r_3), \\ R_{ev}(C, B, A) = (r_1, r_4, r_3, r_2). \end{array} \right. \quad (7)$$

For *role-based schema*, given a role-based fact  $(r_1, A, r_2, B, r_3, C)$ , we consider swapping key-value pairs,  $\{r_1, A\}$ ,  $\{r_2, B\}$ , and  $\{r_3, C\}$ , positionally, also as labeled training span-tuple cases. Thus  $R_{ro}(A, B, C) = (r_1, r_2, r_3)$  can be augmented with 6 orders of span-tuples:

$$\left\{ \begin{array}{l} R_{ro}(A, B, C) = (r_1, r_2, r_3), \\ R_{ro}(B, A, C) = (r_2, r_1, r_3), \\ R_{ro}(A, C, B) = (r_1, r_3, r_2), \\ R_{ro}(B, C, A) = (r_2, r_3, r_1), \\ R_{ro}(C, A, B) = (r_3, r_1, r_2), \\ R_{ro}(C, B, A) = (r_3, r_2, r_1). \end{array} \right. \quad (8)$$

For *hypergraph-based schema*, given a hypergraph-based fact  $(r_1, A, B, C)$ , we consider keeping the main relation  $r_1$  unchanged, and swapping entities,  $A$ ,  $B$ , and  $C$ , positionally, also as labeled training span-tuple cases. Thus  $R_{hg}(A, B, C) = (r_1)$  can be augmented with 6 orders of span-tuples:

$$\left\{ \begin{array}{l} R_{hg}(A, B, C) = (r_1), \\ R_{hg}(B, A, C) = (r_1), \\ R_{hg}(A, C, B) = (r_1), \\ R_{hg}(B, C, A) = (r_1), \\ R_{hg}(C, A, B) = (r_1), \\ R_{hg}(C, B, A) = (r_1). \end{array} \right. \quad (9)$$

### B Supplement to Hetero-ordered Merging

In addition to the *hyper-relational schema*, the hetero-ordered merging strategies for other schemas are as follows:

For *event-based schema* ( $n_r = 4$ ), we combine the predicted probabilities of four labels  $\mathbf{P}_1, \mathbf{P}_2, \mathbf{P}_3, \mathbf{P}_4$  in 6 orders to  $(A, B, C)$  order as follows:

$$\left\{ \begin{array}{l} \mathbf{P}_1 = \mathbf{P}_1^{(ABC)} + \mathbf{P}_1^{(BAC)} + \mathbf{P}_1^{(ACB)} \\ \quad + \mathbf{P}_1^{(BCA)} + \mathbf{P}_1^{(CAB)} + \mathbf{P}_1^{(CBA)}, \\ \mathbf{P}_2 = \mathbf{P}_2^{(ABC)} + \mathbf{P}_3^{(BAC)} + \mathbf{P}_2^{(ACB)} \\ \quad + \mathbf{P}_4^{(BCA)} + \mathbf{P}_3^{(CAB)} + \mathbf{P}_4^{(CBA)}, \\ \mathbf{P}_3 = \mathbf{P}_3^{(ABC)} + \mathbf{P}_2^{(BAC)} + \mathbf{P}_4^{(ACB)} \\ \quad + \mathbf{P}_2^{(BCA)} + \mathbf{P}_4^{(CAB)} + \mathbf{P}_3^{(CBA)}, \\ \mathbf{P}_4 = \mathbf{P}_4^{(ABC)} + \mathbf{P}_4^{(BAC)} + \mathbf{P}_3^{(ACB)} \\ \quad + \mathbf{P}_3^{(BCA)} + \mathbf{P}_2^{(CAB)} + \mathbf{P}_2^{(CBA)}. \end{array} \right. \quad (10)$$

Then, we take the maximum probability to obtain labels  $r_1, r_2, r_3, r_4$ , forming a 3-ary relational fact  $(r_1, r_2, A, r_3, B, r_4, C)$  and filter it out if there are null-label in  $(r_1, r_2, r_3, r_4)$ .

For *role-based schema* ( $n_r = 3$ ), we combine the predicted probabilities of three labels  $\mathbf{P}_1, \mathbf{P}_2, \mathbf{P}_3$  in 6 orders to  $(A, B, C)$  order as follows:

$$\left\{ \begin{array}{l} \mathbf{P}_1 = \mathbf{P}_1^{(ABC)} + \mathbf{P}_2^{(BAC)} + \mathbf{P}_1^{(ACB)} \\ \quad + \mathbf{P}_3^{(BCA)} + \mathbf{P}_2^{(CAB)} + \mathbf{P}_3^{(CBA)}, \\ \mathbf{P}_2 = \mathbf{P}_2^{(ABC)} + \mathbf{P}_1^{(BAC)} + \mathbf{P}_3^{(ACB)} \\ \quad + \mathbf{P}_1^{(BCA)} + \mathbf{P}_3^{(CAB)} + \mathbf{P}_2^{(CBA)}, \\ \mathbf{P}_3 = \mathbf{P}_3^{(ABC)} + \mathbf{P}_3^{(BAC)} + \mathbf{P}_2^{(ACB)} \\ \quad + \mathbf{P}_2^{(BCA)} + \mathbf{P}_1^{(CAB)} + \mathbf{P}_1^{(CBA)}. \end{array} \right. \quad (11)$$

Then, we take the maximum probability to obtain labels  $r_1, r_2, r_3$ , forming a 3-ary relational fact  $(r_1, A, r_2, B, r_3, C)$  and filter it out if there are null-label in  $(r_1, r_2, r_3)$ .

For *hypergraph-based schema* ( $n_r = 1$ ), we combine the predicted probabilities of one label  $\mathbf{P}_1$  in 6 orders to  $(A, B, C)$  order as follows:

$$\left\{ \begin{array}{l} \mathbf{P}_1 = \mathbf{P}_1^{(ABC)} + \mathbf{P}_1^{(BAC)} + \mathbf{P}_1^{(ACB)} \\ \quad + \mathbf{P}_1^{(BCA)} + \mathbf{P}_1^{(CAB)} + \mathbf{P}_1^{(CBA)}. \end{array} \right. \quad (12)$$

Then, we take the maximum probability to obtain labels  $r_1$ , forming a 3-ary relational fact  $(r_1, A, B, C)$  and filter it out if  $r_1$  is null-label.

## C Construction of Dataset

Based on the original *hyper-relational schema* on HyperRED dataset [5], we construct other three schemas (event-based, role-based, and hypergraph-based) for fine-grained n-ary RE. Firstly, we view the main relation in the *hyper-relational schema* as the event type in the *event-based schema*, combine the head entity and tail entity with two extra head key and tail key to convert them into two key-value pairs, and remain the auxiliary key-value pairs in the *hyper-relational schema*. Taking ‘Einstein received his Doctorate degree in Physics from the University of Zurich.’ as an example, it can be represented as  $(Einstein, educated, University\ of\ Zurich, \{academic\_major, Physics\}, \{academic\_degree, Doctorate\})$  in the *hyper-relational schema* and  $(education, \{trigger, received\}, \{person, Einstein\}, \{college, University\ of\ Zurich\}, \{academic\_major, Physics\}, \{academic\_degree, Doctorate\})$  in the *event-based schema*. Secondly, we remove the event type in the *event-based schema* to obtain the *role-based schema*. Thirdly, we remove all the keys in key-value pairs and remain the relation to build the *hypergraph-based schema*.

## D Baseline Settings

Firstly, for the original *hyper-relational schema* of HyperRED, we adopted the same baselines as in the CubeRE paper [5] to compare with Text2NKG:

**Generative Baseline:** Generative Baseline uses BART [14], a sequence-to-sequence model, to transform input sentences into a structured text sequence.

**Pipeline Baseline:** Pipeline Baseline uses UniRE [24] to extract relation triplets in the first stage and a span extraction model based on BERT-Tagger [6] to extract value entities and corresponding qualifier labels in the second stage.

**CubeRE:** CubeRE [5] is the only hyper-relational extraction model that uses a cube-filling model inspired by table-filling approaches and explicitly considers the interaction between relation triplets and qualifiers.

Secondly, for the *event-based schema*, *role-based schema*, and *hypergraph-based schema*, we added the following baselines to further validate the effect of Text2NKG on the fine-grained N-ary relation fact extraction task in the HyperRED dataset:

**Text2Event:** Text2Event [5] is a classic model in the Event extraction domain. However, it is not applicable to extractions of the *hyper-relational schema*. For the *role-based schema* extraction, we retained the key without referring to the main relation, while for the *hypergraph-based schema* extraction, we retained the main relation without referring to the key to get the final result for comparison.

**UIE / LasUIE:** UIE [17] and LasUIE [9] are unified information extraction models that can handle most tasks like NER, RE, EE, etc. However, they are still only suitable for event extraction in the multi-relational extraction domain and are not applicable to extractions of the *hyper-relational schema*. Therefore, we adopted the same approach as with Text2Event to compare with Text2NKG.

Thirdly, under the impact of the wave of large-scale language models brought about by ChatGPT on traditional natural language processing tasks, we added unsupervised large models as baselines to compare with Text2NKG in the n-ary RE tasks of the four schemas.

**ChatGPT / GPT4:** Using different prompts, we tested the latest state-of-the-art large-scale pre-trained language models ChatGPT [25] and GPT-4 [19] in an unsupervised manner, evaluating their performance on the extraction of the four schemas.

## E Hyperparameter Settings

We use the grid search method to select the optimal hyperparameter settings for both Text2NKG with Bert-base and Bert-large. We use the same hyperparameter settings in Text2NKG with different encoders. The hyperparameters that we can adjust and the possible values of the hyperparameters are first determined according to the structure of our model in Table 4. Afterward, the optimal hyperparameters are shown in **bold**.

Hyperparameter	HyperRED
$\alpha$	{1.0, 0.1, <b>0.01</b> , 0.001}
Train_batch_size	{2, 4, <b>8</b> , 16}
Eval_batch_size	{ <b>1</b> }
Learning rate	{ $1e-5$ , <b><math>2e-5</math></b> , $5e-5$ }
Max_sequence_length	{128, <b>256</b> , 512, 1024}
Weight decay	{ <b>0.0</b> , 0.1, 0.2, 0.3}

Table 4: Hyperparameter Selection.

## F Model Training Details

We train 10 epochs on HyperRED with the optimal combination of hyperparameters. Text2NKG and all its variants have been trained on a single NVIDIA A100 GPU. Using our optimal hyperparameter settings, the time required to complete the training on HyperRED is 4h with BERT-base encoder and 10h with BERT-large encoder.

## G Further Discussions

### G.1 How does ChatGPT perform in Fine-grained N-ary RE tasks?

We have tried to use LLM APIs such as ChatGPT and GPT to do similar n-ary RE tasks, i.e., prompting model input and output formats for extraction. The advantage of ChatGPT is that it can perform similar tasks in a few-shot situation, however, for building high-quality knowledge graphs, the performance and the fineness of the n-ary RE are much lower than Text2NKG. This is because ChatGPT is not good at multi-label classification tasks that contain less semantic interpretation. When the number of labels of relations in our relation collection is very large, we need to write a very long prompt to tell the LLM about our label candidate collection, which again leads to the problem of forgetting. Therefore, we have tried numerous prompt templates to enhance the extraction effect of ChatGPT, however, on fine-grained n-ary RE task, the best result of ChatGPT can only reach about 10% of  $F_1$  value on HyperRED, which is much lower than the result of 80%+  $F_1$  value of Text2NKG.

However, advanced LLMs such as ChatGPT are a good idea for training dataset generation for Text2NKG in such tasks to save some manual labor to only verify and correct the training items generated. For future work, we will continue our research in this direction and try to combine large language models with Text2NKG-like supervised models for automated fine-grained n-ary RE for n-ary relational knowledge graph construction.

### G.2 Why first Extracting 3-ary facts and then Merging them into N-ary Facts ?

We use output merging to address the dynamic changes in the number of elements in n-ary relational facts. The atomic unit of an n-ary fact includes a 3-ary fact with three entities. For instance, in the hyper-relational fact (*Einstein, educated\_at, University of Zurich, degree: Doctorate degree, major: Physics*), the Text2NKG algorithm allows us to extract two 3-ary atomic facts: (*Einstein, educated\_at, University of Zurich, degree: Doctorate degree*) and (*Einstein, educated\_at, University of Zurich, major: Physics*). These are then merged based on the same primary triple (*Einstein, educated\_at, University of Zurich*) to form a 4-ary fact. The same principle applies to facts of higher arities.

As another example demonstrating the problem with merging binary relations: consider the statement "*Einstein received his Bachelor's degree in Mathematics and his Doctorate degree in Physics.*" When represented as binary relations, the facts become (*Einstein, degree, Doctorate degree*), (*Einstein, major, Physics*), (*Einstein, degree, Bachelor*), and (*Einstein, major, Mathematics*). With this representation, we cannot merge these binary relation facts effectively because there's no way to determine whether *Einstein's doctoral major* was *Physics* or *Mathematics*. This necessitates the use of NKG's n-ary relationship facts to represent this information, as seen in (*Einstein, degree, Doctorate degree, major, Physics*).

Therefore, using binary facts, we can't merge them into n-ary facts based on shared elements within these facts. On the other hand, using facts with four entities or more makes it challenging to effectively extract 3-ary atomic facts.

In Section 5.6 and Figure 5(a), we also analyzed the effects and detailed insights of unsupervised extraction of arbitrary-arity facts.

### G.3 How Text2NKG can address Long Contexts with Relations spread across Various Sentences ?

As long as the text to be extracted is a lengthy piece with entities annotated, it can undergo long-form n-ary relation extraction. The maximum text segment size for our proposed method depends on the maximum text length that a transformer-based encoder can accept, such as Bert-base and Bert-large, which have a maximum limit of 512. To extract from larger documents, we simply need to switch to encoders with larger context length, which all serve as the encoder portion of Text2NKG and are entirely decoupled from the n-ary relation extraction technique we propose. This is one of the advantages of Text2NKG. Its primary focus is to address the order and combination issues of multi-ary relationships. We can seamlessly combine a transformer encoder that supports long texts with Span-tuple Multi-label Classification to process n-ary relation extraction in long chapters.



## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction accurately reflect the paper's contributions and scope. The claims made are supported by the detailed methodology and experimental results sections. Text2NKG demonstrates improvements in fine-grained n-ary relation extraction, supports multiple NKG schemas, and achieves state-of-the-art performance as claimed. The public availability of code and datasets further supports the paper's transparency and reproducibility.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The paper discusses several limitations, including the constraint of supervised training, the suggestion to combine supervised and unsupervised models, and the handling of long contexts with appropriate encoders. These points are addressed in the comparison with large language models and in discussions about future work and handling long texts.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best

judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: The paper provides a comprehensive set of assumptions and proofs for its theoretical results. These are detailed in both the main sections and the appendices. For instance, the hetero-ordered merging strategy and the output merging methodology are explained with equations and detailed descriptions of the processes involved. Additionally, the assumptions for the span-tuple multi-label classification method and the handling of null-label weights are clearly stated, with mathematical formulations provided to support the theoretical claims.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The paper includes detailed information necessary to reproduce the main experimental results. It provides a comprehensive description of the datasets used, the baselines for comparison, the experimental setup, and the evaluation metrics. Additionally, the paper mentions the use of a single NVIDIA A100 GPU and details the hyperparameters and training environment in the appendices. The inclusion of the link to the anonymous GitHub repository ensures that the code and datasets are accessible for reproducibility.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.

- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [\[Yes\]](#)

Justification: The paper provides open access to the data and code, including a link to an anonymous GitHub repository. It also includes detailed instructions on how to reproduce the main experimental results, covering data access, preparation, and the specific commands and environment needed to run the experiments. These details are described in the supplemental material and appendices.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [\[Yes\]](#)

Justification: The paper specifies all the necessary training and test details, including data splits, hyperparameters, and the type of optimizer used. This information is provided in the main text and further elaborated in the appendices, ensuring that the experimental setup is fully transparent and reproducible.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: The paper reports error bars and provides appropriate information about the statistical significance of the experiments. The factors of variability captured by the error bars are clearly stated, and the method for calculating the error bars is explained. The assumptions made and whether the error bar represents the standard deviation or the standard error of the mean are also clearly mentioned.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The paper provides detailed information on the compute resources used for the experiments. It specifies that all experiments were conducted on a single NVIDIA A100 GPU. The training details include the time required for training on HyperRED, which is 4 hours with the BERT-base encoder and 10 hours with the BERT-large encoder. The optimal hyperparameter settings and the use of the Adam optimizer are also described, along with the number of epochs and batch sizes.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.

- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: The research described in the paper conforms to the NeurIPS Code of Ethics. We have ensured transparency, reproducibility, and ethical use of the data and models. We have also made the code and datasets publicly available to support open science and reproducibility. There is no indication of any ethical violations or concerns regarding the research process or outcomes.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The paper discusses potential positive societal impacts, such as improving the quality of knowledge graphs and enabling more accurate information extraction from text. It also considers potential negative impacts, such as the possibility of misuse of the technology for disinformation or unfair decision-making. We suggest that combining small supervised models with large unsupervised models could help mitigate some of these risks by improving the accuracy and robustness of the extraction process.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [Yes]

Justification: The paper addresses the potential misuse of the proposed methods and models. It emphasizes the importance of responsible use and includes strategies to mitigate risks. We suggest combining small supervised models with large unsupervised models to balance solving the cold-start problem and fine-grained extraction. They highlight that unsupervised large language models like ChatGPT and GPT-4 cannot accurately distinguish closely related relations, which implies careful consideration for controlled use. They also recommend specific usage guidelines and restrictions to ensure safe and ethical deployment.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The paper properly credits the creators and original owners of the assets used. It cites the relevant datasets and models, including their sources and versions. We ensure that the licenses and terms of use are explicitly mentioned and respected. This information is detailed in the references and the supplemental material, where the datasets and models used are listed along with their corresponding licenses.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The new assets introduced in the paper are well documented. The documentation includes details about the training data, model architecture, and usage instructions. We provide structured templates for communicating the dataset/code/model details, including training procedures, licenses, limitations, and consent obtained from people whose data is used. This comprehensive documentation is available alongside the assets in the supplemental material and the anonymous GitHub repository.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

#### 14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing experiments or research with human subjects. Therefore, this question is not applicable to the current work.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

#### 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing experiments or research with human subjects. Therefore, IRB approvals or equivalent reviews are not applicable to the current work.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.