

---

# Wiki Entity Summarization Benchmark

---

**Saeedeh Javadi\***  
Polytechnic University of Turin  
saeedeh.javadi@studenti.polito.it

**Atefeh Moradan\***  
Aarhus University  
atefeh.moradan@cs.au.dk

**Mohammad Sorkhpar\***  
Indiana State University  
msorkhpar@sycamores.indstate.edu

Klim Zaporjets  
Aarhus University  
klim@cs.au.dk

Davide Mottin  
Aarhus University  
davide@cs.au.dk

Ira Assent  
Aarhus University  
ira@cs.au.dk

## Abstract

1 Entity summarization aims to compute concise summaries for entities in knowledge  
2 graphs. Existing datasets and benchmarks are often limited to a few hundred  
3 entities and discard graph structure in source knowledge graphs. This limitation  
4 is particularly pronounced when it comes to ground-truth summaries, where there  
5 exist only a few labeled summaries for evaluation and training. We propose WIKES  
6 (Wiki Entity Summarization Benchmark), a comprehensive *benchmark* comprising  
7 of entities, their summaries, and their connections. Additionally, WIKES features  
8 a dataset *generator* to test entity summarization algorithms in different areas of the  
9 knowledge graph. Importantly, our approach combines graph algorithms and NLP  
10 models, as well as different data sources such that WIKES does not require human  
11 annotation, rendering the approach cost-effective and generalizable to multiple  
12 domains. Finally, WIKES is scalable and capable of capturing the complexities of  
13 knowledge graphs in terms of topology and semantics. WIKES features existing  
14 *datasets* for comparison. Empirical studies of entity summarization methods  
15 confirm the usefulness of our benchmark. Data, code, and models are available at:  
16 <https://github.com/msorkhpar/wiki-entity-summarization>.

## 17 1 Introduction

18 *Knowledge Graphs* (KGs) are a valuable information representation: interconnected networks of  
19 entities and their relationships enable machine reasoning to empower question answering Hu et al.  
20 [2018], Lan et al. [2019], recommender systems Wang et al. [2018], information retrieval Raviv et al.  
21 [2016]. KGs may comprise millions of entities representing real-world objects, concepts, or events.

22 Yet, the size and complexity of these KGs progressively expand, rendering it increasingly challenging  
23 to convey the essential information about an entity in a concise and meaningful way Suchanek et al.  
24 [2007], Vrandečić and Krötzsch [2014]. This is where entity summarization becomes relevant. *Entity*  
25 *summarization* (ES) Liu et al. [2021] is the process of generating a concise and informative summary  
26 that captures the most salient aspects of the entity description, based on the information available in

---

\*Equal contribution

27 the KGs. In ES, the entity *description* refers to all the triples involving such an entity. For instance,  
 28 Figure 1 illustrates a set of relationships surrounding the entity Ellen Johnson Sirleaf in a KG,  
 29 along with a possible summary for this entity. Extensive descriptions can overwhelm users and  
 30 exceed the capacity of typical user interfaces, making it challenging to identify the most relevant  
 31 triples. Entity summarization addresses this issue by computing an optimal compact summary for an  
 entity, selecting a size-constrained subset of triples Liu et al. [2021].

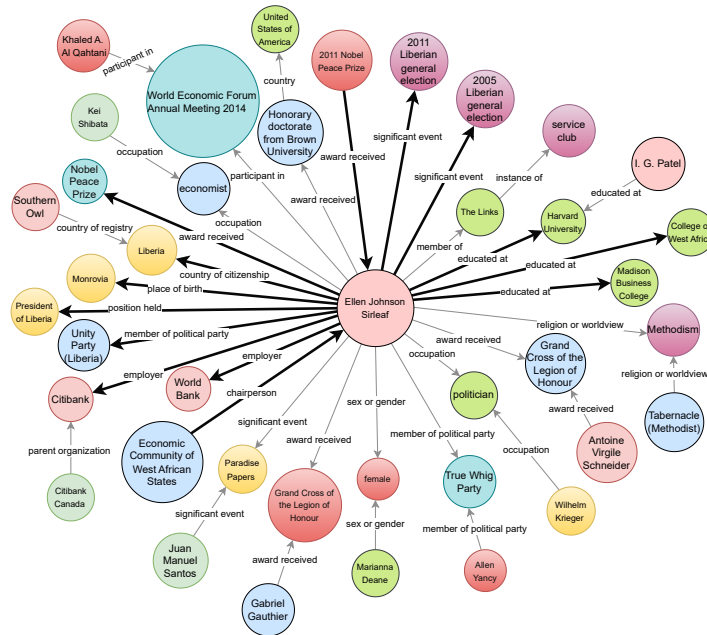


Figure 1: KG subgraph of entity Ellen Johnson Sirleaf: arrows depict the subgraph of relationships to other entities, and labels indicate their roles. Selecting the bold edges as entity summaries of the most relevant triples may reduce information overload while concisely describing the entity.

32  
 33 Despite advances in entity summarization techniques Liu et al. [2021], the development and evaluation  
 34 of these methods are hindered by a number of limitations in the benchmarks and datasets Liu et al.  
 35 [2020], Cheng et al. [2023]. The first limitation of the current benchmarks is the small dataset size,  
 36 encompassing only a few hundred entities. Second, the generation of ground-truth summaries for  
 37 testing mostly relies on expensive and lengthy manual annotation. Moreover, the dependence on a few  
 38 human annotators often biases the data towards the annotators' preferences and knowledge. Third,  
 39 existing benchmarks often disregard the wealth of information in the knowledge graph structure.

40 To address the above limitations, we propose:

- 41 • **Novel WIKES benchmark for ES** based on summaries and graphs from Wikidata and Wikipedia.
- 42 • **Subgraph extraction method** preserving the complexity of real-world KGs; subsampling using  
 43 random walks and proportionally preserving node degrees, WIKES captures the structure of the  
 44 entities up to the second-hop neighborhood, thereby ensuring that the connections in WIKES  
 45 accurately reflect those in the source KG.
- 46 • **Comprehensive summaries for any entity in the KG**, ensuring that summaries are both relevant  
 47 and contextually rich by deriving them directly from corresponding Wikipedia abstracts, minimizing  
 48 human bias, as these abstracts are created and reviewed by several experts. In this manner, WIKES is  
 49 scalable, enabling it to generate large benchmark resources efficiently with high-quality annotation.
- 50 • **Automatic entity summarization dataset generator** allows for the creation of arbitrarily large  
 51 datasets, encompassing various domains of knowledge.

## 52 2 Existing Datasets

53 Here, we review the existing datasets for entity summarization. Table 1 provides an overview and  
 54 statistics of the current datasets in this field. **FACES** and **INFO** datasets have a higher density  
 55 than the entities in the Entity Summarization Benchmark (ESBM). It is also clear that **LMDB** and  
 56 **FACES** are not connected graphs, that challenge graph-based learning methods where the information  
 57 cannot easily propagate in disconnected networks. Specifically, **FACES** consists of 12 connected  
 58 components, which complicates the learning process for graph embedding methods by limiting the  
 59 richness of information that can be leveraged from the graph.

Table 1: Entity summarization datasets in terms of number of entities  $|\mathcal{V}|$ , triples  $|\mathcal{E}|$ , number of ground-truth summaries (target entities), density as  $|\mathcal{E}|/\binom{|\mathcal{V}|}{2}$ , graph connectivity, number of components, sampling method to select entities and subgraph, and minimum / maximum node degree.

Metric	DBpedia (ESBM)	LMDB (ESBM)	FACES	INFO
Entities ( $ \mathcal{V} $ )	2 721	1 853	1 379	1 410
Relations ( $ \mathcal{E} $ )	4 436	2 148	2 152	2 019
Target Entities	125	50	50	100
Density	0.0005	0.0006	0.0011	0.0010
Sampling method	Not specified	Not specified	Not specified	Not specified
Connected-graph	Yes	No	No	Yes
Num-comp	1	2	12	1
Min Degree	1	1	1	1
Max Degree	125	208	88	100

60 We provide here a comprehensive description of each dataset or benchmark:

- 61 • **ESBM** Liu et al. [2020]: The Entity Summarization Benchmark (ESBM) is the first benchmark to  
 62 evaluate the performance of entity summarization methods. ESBM has three versions; v1.2 is the  
 63 latest and most extensive version. This version comprises 175 entities, with 150 from DBpedia  
 64 Lehmann et al. [2015] and 25 from LinkedMDB Hassanzadeh and Consens [2009]. The summaries  
 65 comprise triples selected by 30 “researchers and students“ annotators. Each entity has exactly 6  
 66 summaries. Despite encompassing two datasets, ESBM has several limitations. First, the entity  
 67 sampling method is not explained. In particular, some triples in the neighborhood of the entity are  
 68 missing in the datasets. Second, there are no connections among the entities in the neighborhood,  
 69 nor any two-hop neighborhood. Third, the expertise and background of the annotators are not  
 70 assessed nor disclosed. Due to the expensive annotation process, the dataset size is small.
- 71 • **FACES** Gunaratna et al. [2015] is a dataset from DBpedia (version 3.9) Auer et al. [2007] and  
 72 includes 50 randomly selected entities, each with at least 17 different types of relations. Similar to  
 73 ESBM, the FACES ground-truth is also generated manually.
- 74 • **INFO** Cheng et al. [2023] contains 100 randomly selected entities from 10 classes in DBpedia. It  
 75 comprises two sets of ground-truth summaries, REF-E and REF-W. REF-E summaries comprise a  
 76 selection of triples from five experts adhering to a 140-character limit, similar to typical Google  
 77 search result snippets. REF-W summaries are obtained by one expert who reads the abstract  
 78 sections of the respective entities on Wikipedia and selects neighboring entities that closely match  
 79 the Wikipedia abstracts. The number of ground-truth summaries per entity varies, as some experts  
 80 evaluate multiple entities. This inconsistency complicates the evaluation process. The expertise of  
 81 the annotators remains unspecified.

82 In contrast, our benchmark uses Wikidata to automatically map entities from Wikipedia to Wikidata.  
 83 This automation allows us to efficiently generate summaries for any number of entities. Unlike  
 84 previous work, we use the Wikipedia abstract as a summary instead of manual annotators. Each  
 85 abstract is a collaboration of many users; as such, it should not introduce obvious biases. Additionally,  
 86 with this process, we ensure high-quality and cost-effective summaries. Furthermore, we present the  
 87 characteristics of our dataset in Table 3. The WIKES benchmark includes a larger number of entities

88 and relations than existing datasets. It is a connected graph containing approximately 500 seed nodes.  
89 Further details regarding the specific characteristics of our dataset are provided in Section 3.4.

### 90 3 The WIKES Benchmark

91 A *Knowledge Graph*  $\mathcal{KG} = (\mathcal{V}, \mathcal{R}, \mathcal{T})$  is a directed multigraph consisting of entities  $\mathcal{V} =$   
92  $\{v_1, \dots, v_n\}$ , relationships  $\mathcal{R}$ , and triples  $\mathcal{T} \subseteq \mathcal{V} \times \mathcal{R} \times \mathcal{V}$ . The set of edges  $\mathcal{E} = \{(i, j) \mid$   
93  $v_i, v_j \in \mathcal{V} \wedge \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{T}\}$  contains pairs of nodes connected by a relationship.

94 The *t-hop neighborhood*  $\mathcal{N}_t(v_i)$  of node  $v_i$  is the set of nodes reachable from  $v_i$  within  $t$  edges when  
95 ignoring edge directions.

96 A *summary* for an entity  $v_i$  is a subset  $\mathcal{S}(v_i) \subseteq \Delta_t(v_i)$  of triples from the  $t$ -description of  $v_i$ , where  
97 the *t-description* of an entity  $v_i \in \mathcal{V}$  in a knowledge graph  $\mathcal{KG}$  is the set  $\Delta_t(v_i) = \{(s, p, o) \in \mathcal{T} \mid$   
98  $s \in \mathcal{N}_t(v_i) \vee o \in \mathcal{N}_t(v_i)\}$  of triples in which one of the entities is in the  $t$ -hop neighborhood of  $v_i$ .

99 **Entity summarization** for an entity  $v_i \in \mathcal{V}$  in a knowledge graph  $\mathcal{KG}$  aims to find a summary  $\mathcal{S}(v_i)$   
100 that maximizes some score among all possible summaries for  $v_i$ , i.e.,

$$\arg \max_{\substack{\mathcal{S}(v_i) \subseteq \Delta_t(v_i) \\ |\mathcal{S}(v_i)|=k}} \text{score}(\mathcal{S}(v_i)), \quad (1)$$

#### 101 3.1 Extracting Summaries from Wikidata using Wikipedia Abstracts

102 We extract summaries for each Wikidata item using Wikipedia abstracts and infoboxes. Each abstract  
103 is a joint effort of many users and experts, which ensures quality and accuracy. Leveraging Wikipedia,  
104 we avoid time-consuming manual annotation and enable the automatic generation of large-scale  
105 datasets.

106 **Wikidata** is a free and collaborative knowledge base that collects structured data to support Wikipedia  
107 and other Wikimedia projects. It includes descriptions and labels for entities. The descriptions offer  
108 in-depth details, while the labels serve as concise identifiers, facilitating efficient data retrieval  
109 and integration in subsequent steps. We load all Wikidata items XML dump files published on  
110 2023/05/01<sup>2</sup> as entities  $\mathcal{V}$  alongside their properties as relationships  $\mathcal{R}$  into a graph database<sup>3</sup>. The  
111 result is a graph that connects all Wikidata items and statements. We include items if they (1) are not  
112 marked as redirects, (2) belong to the main Wikidata namespace, and (3) have an English label or  
113 description. Additionally, we load metadata for each Wikidata item and property, including labels  
114 and descriptions, into a relational database<sup>4</sup>. **Wikipedia** pages contain infoboxes, abstracts, page  
115 content, categories, references, and more. Links to other Wikipedia pages are referred to as mentions.  
116 We detect these mentions in the abstracts and infoboxes of Wikipedia pages to use them later for  
117 labeling the summaries in Wikidata. We extract and load all the content from the XML dump files of  
118 Wikipedia pages, published on 2023/05/01<sup>5</sup>, into a relational database under the same conditions as  
119 Wikidata: the pages must be in English and not redirected.

120 **Summary annotation.** We annotate the summaries in Wikidata using the corresponding Wikipedia  
121 pages. For each Wikipedia page corresponding to a Wikidata entity, we iterate through all connected  
122 Wikidata items using Wikidata properties. If a connected Wikidata item is mentioned in the Wikipedia  
123 abstract and infobox, we annotate the Wikidata item with the corresponding Wikidata property as  
124 part of the summary.

125 Wikidata is a directed multigraph, which means that each entity (Wikidata item) can be connected to  
126 another entity via multiple relations (Wikidata properties). Yet, links in Wikipedia are not labeled;  
127 as such, we need to select one of the relations for the summary. To annotate the correct Wikidata

<sup>2</sup><https://dumps.wikimedia.org/wikidatawiki/>

<sup>3</sup><https://neo4j.com>

<sup>4</sup><https://www.postgresql.org/>

<sup>5</sup><https://dumps.wikimedia.org/enwiki/>

128 property as part of the summary, we employ the DistilBERT model Sanh et al. [2019]. DistilBERT is  
 129 a fast and lightweight model with a reduced number of parameters compared to the original BERT  
 130 model. This way, we can efficiently process large amounts of data while maintaining high-quality  
 131 embeddings for accurate relation selection.

132 Concretely, we first embed the abstract of the Wikidata item for which we are generating summaries  
 133 using DistilBERT. We then calculate the cosine similarity between the embedding of the abstract  
 134 and the embeddings of each candidate relation. Finally, we add the relation with the highest cosine  
 135 similarity to the abstract embedding to the summary. This approach ensures that the most relevant  
 136 Wikidata property is selected for the summary based on its semantic similarity to the Wikipedia  
 137 abstract.

### 138 3.2 Capturing the Graph Structure

139 Here we introduce the WIKES generator algorithm. The main idea is to sample a connected graph  
 140 that preserves the original graph structure. To this end, we employ random walks Pearson [1905].

141 A random walk is a stochastic process defined as a sequence of steps, where the direction and  
 142 magnitude of each step are determined by the random variable  $X_{t+1} = X_t + S_t$  where  $X_t$  represents  
 143 the position at time  $t$ , and  $S_t$  is the step taken from position  $X_t$ .

144 The process is a Markov process, characterized by its memoryless property:

$$P(X_{t+1} = x | X_t = x_t, X_{t-1} = x_{t-1}, \dots, X_0 = x_0) = P(X_{t+1} = x | X_t = x_t) \quad (2)$$

145 In adapting this concept to our work, we redefine the number of random walks assigned to nodes  
 146 based on their degrees, ensuring the distribution remains proportional to real data. This is achieved  
 147 through logarithmic transformation and normalization. The logarithmic transformation is applied to  
 148 reduce the impact of high-degree nodes and also low-degree nodes, making it more manageable for  
 149 the random walk. Given a graph with node degrees  $\{d_1, d_2, \dots, d_i\}$ , the log-transformed degree for  
 150 node  $i$  is  $L_i = \log(d_i)$ . These values are then normalized:

$$N_i = \frac{L_i - \min(\{L\})}{\max(\{L\}) - \min(\{L\})} \quad (3)$$

151 where  $N_i$  is the normalized logarithmic degree of node  $i$ . Finally, the number of random walks  $R_i$   
 152 assigned to each node is:

$$R_i = \text{round}(\text{minRW} + N_i \times (\text{maxRW} - \text{minRW})) \quad (4)$$

153 Here, minRW and maxRW are the user-defined minimum and maximum limits for random walks.  
 154 This adaptation ensures that the random walks are proportional to the normalized logarithmic degree  
 155 of each node, reflecting the true structure of the network. For a small dataset we set minRW = 100  
 156 and maxRW = 300; for a medium dataset minRW = 150 and maxRW = 600; for a large dataset,  
 157 minRW = 300 and maxRW = 1800. This ensures that the random walks are tailored to both the  
 158 scale and the complexity of the dataset. Importantly, our approach can be used to extract further  
 159 subgraphs at the scale needed for benchmarking in a given scenario.

160 Moreover, the random walk sampling process requires a set of seed nodes as a starting point. In our  
 161 case, the seed nodes represent the target entities we are interested in. The seed nodes can be any  
 162 Wikidata Item Identifier, Wikipedia title, or Wikipedia ID of the Wikipedia pages. We collect the seed  
 163 nodes on the condition that they have at least  $k$  (default  $k = 5$ ) common entities with the abstract  
 164 section and the infobox in the Wikipedia pages. Therefore, this model is flexible, allowing you to  
 165 choose any seed nodes from any domain as an input. In the datasets that we generated, we collect  
 166 seed nodes from Laouenan et al. [2022]. This paper has published information about individuals  
 167 from various domains. The authors collected data from multiple Wikipedia editions and Wikidata,  
 168 using deduplication and cross-verification techniques to compile a database of 1.6 million individuals  
 169 with English Wikipedia pages. The seed nodes that we use include actor, athletic, football, journalist,  
 170 painter, player, politician, singer, sport, writer, lawyer, film, composer, novelist, poet, and screenwriter.  
 171 Using combinations of these seed nodes, we generate four sets of datasets, with each set having small,

172 medium, and large versions. In Table 4 in Section 6 in the supplementary material, we present the  
173 seed nodes and their proportions for each dataset and their corresponding train-test-val splits.

### 174 3.3 WIKES Generator

175 We discuss how WIKES is created, and how further benchmarks can be generated without the need  
176 for manual annotators. Algorithm 1 details the generator, which consists of the following steps.

177 **Step1:** Retrieve summaries of each seed node (explained in Section 3.1)

178 **Step2:** Expand the graph using the random walk method in Section 3.2. Set the random walk’s length  
179  $n$  (default  $n = 2$ ), which means it explores up to the  $n$ -hop neighborhood of each seed node.

180 **Step3:** Check if the graph is connected. If it is, done. If not, identify all disconnected components  
181 and sort them by size, from largest to smallest. In each iteration, connect smaller components to the  
182 largest component using  $h$  connections. Utilize the shortest path method, selecting paths that are equal  
183 to or less than a minimum path length  $l$ . Continue connecting nodes from the smaller component  
184 to the larger one until  $h$  nodes are connected. After each iteration, check graph connectivity again.  
185 If all components are connected to the largest component, the algorithm ends. Otherwise, re-sort  
components and increase  $l$  by 1. Repeat until the graph is a single connected component.

---

**Algorithm 1** WIKES Generator

---

```
1: Input: Graph  $G$ , seed nodes  $S$ , random walk length  $n$ , minimum path length  $l$ 
2: Output: A connected graph
3: procedure GENERATEGRAPH( $G, S, n, l$ )
4:    $summaries \leftarrow$  RETRIEVESUMMARIES( $S$ )
5:    $G \leftarrow$  RANDOMWALKEXPANSION( $G, S, n$ ) mentioned in section 3.2
6:    $is\_connected \leftarrow$  CHECKCONNECTIVITY( $G$ )
7:   while not  $is\_connected$  do
8:      $components \leftarrow$  FINDCOMPONENTS( $G$ )
9:     Sort  $components$  by size in descending order
10:     $largest \leftarrow components[0]$ 
11:    for  $comp$  in  $components[1 : ]$  do
12:      Connect  $comp$  to  $largest$  using  $h$  connections via shortest paths  $\leq l$ 
13:       $G \leftarrow$  UPDATEGRAPH( $G, comp, largest$ )
14:       $is\_connected \leftarrow$  CHECKCONNECTIVITY( $G$ )
15:      if  $is\_connected$  then
16:        break
17:      end if
18:    end for
19:     $l \leftarrow l + 1$ 
20:  end while
21:  return  $G$ 
22: end procedure
```

---

186

### 187 3.4 WIKES Datasets

188 We generate three sizes for each of the four datasets, obtaining 12 datasets. We present their  
189 characteristics in Table 3 in section 6. The number of entities in the small datasets ranges from  
190 approximately  $70k$  to  $85k$ , and the number of relations ranges from around  $120k$  to  $135k$ . In the  
191 medium datasets, the number of entities ranges from  $100k$  to  $130k$ , and the number of relations  
192 ranges from  $195k$  to  $220k$ . The number of entities in the large datasets ranges from approximately  
193  $185k$  to  $250k$ , and the number of relations ranges from around  $397k$  to  $470k$ . The average runtime for  
194 generating small graphs is approximately 128 seconds; for medium-sized graphs, it is approximately  
195 216 seconds; and for large graphs, it is approximately 512 seconds. We construct the train-test-  
196 validation split for each dataset with 70% for training, 15% for testing, and 15% for validation.  
197 Detailed information about the run time, as well as the number of nodes and relations for these splits,  
198 is available on our GitHub repository. All graphs in each train-test-validation splits are connected.

199 **4 Empirical Evaluation**

200 We study the quality of WIKES using the following metrics:

201 **F-Score.** Let  $\mathcal{S}_m$  the summary obtained by a summarization method and  $\mathcal{S}_h$  the ground-truth  
 202 summary. We compare  $\mathcal{S}_m$  with  $\mathcal{S}_h$  using the F1-score based on precision  $P$  and recall  $R$ :

$$\text{F1} = \frac{2 \cdot P \cdot R}{P + R}, \text{ where } P = \frac{|\mathcal{S}_m \cap \mathcal{S}_h|}{|\mathcal{S}_m|} \text{ and } R = \frac{|\mathcal{S}_m \cap \mathcal{S}_h|}{|\mathcal{S}_h|} \quad (5)$$

203 The F1 score lies within  $[0,1]$ . High F1 indicates that  $\mathcal{S}_m$  is closer to the ground-truth  $\mathcal{S}_h$ .

204 **Mean Average Precision (MAP).** This metric is particularly suitable for evaluating ranking tasks  
 205 because it takes into account the order of the predicted triples. MAP calculates precision at each  
 206 position  $i$  in the predicted summary and averages these values over all relevant summary triples. It  
 207 reflects both the relevance and the ranking quality of the predicted summaries. MAP, unlike F1-score,  
 208 does not depend on a specific value of  $k$ . This makes it a robust metric for assessing how well a  
 209 summarization method ranks the relevant triples.

$$\text{MAP} = \frac{1}{N} \sum_{n=1}^N \frac{\sum_{i=1}^{|\mathcal{S}_m^{(n)}|} \begin{cases} \text{Precision}@i(\mathcal{S}_h^{(n)}) & \text{if } \text{Rel}(n, i) \\ 0 & \text{otherwise} \end{cases}}{|\mathcal{S}_h^{(n)}|} \quad (6)$$

210 where  $N$  is the total number of entities,  $\mathcal{S}_h^{(n)}$  is the set of ground-truth summary triples for a particular  
 211 entity  $v_n$ ,  $\mathcal{S}_m^{(n)}$  is the set of predicted summary triples for the entity  $v_n$ ,  $\text{Precision}@i$  is the precision  
 212 at the  $i$ -th position in the predicted summary, and  $\text{Rel}(n, i)$  indicates whether the  $i$ -th predicted triple  
 213 for entity  $v_n$  is relevant (i.e., it belongs to  $\mathcal{S}_h^{(n)}$ ). MAP scores are in the range  $[0,1]$ , where a higher  
 214 MAP indicates better performance in terms of correctly predicting relevant summary triples. To  
 215 account for the varying lengths of the ground-truth summaries in real-world data, we also calculate  
 216 MAP and F-score (which we refer to as dynamic MAP and dynamic F-score) by setting the length of  
 217 the generated summary ( $|\mathcal{S}_m|$ ) equal to the length of the corresponding ground-truth summary ( $|\mathcal{S}_h|$ ).

218 We analyze our dataset and compare it with the ESBM benchmark using statistical measures such as  
 219 frequency and inverse frequency of entities and relations. We calculate the F-score and MAP score  
 220 for the top-5 and top-10 of both the ESBM dataset and our WikiProfem. We choose top-5 and top-10  
 221 because we only have ground-truth summaries for top-5 and top-10 in the ESBM dataset. The F-score  
 222 and MAP results for ESBM are presented in Figure 2. The statistics show that for DBpedia, the  
 223 F-score using inverse relation frequency outperforms the random baseline by 0.15 for top-5 and by  
 224 0.34 for top-10. Furthermore, when using inverse entity frequency, DBpedia achieves an even higher  
 225 F-score, surpassing the random baseline by 0.07 for top-5 and by 0.15 for top-10. For Lmdb, we  
 226 observe a similar trend when using inverse frequency. The F-score surpasses the random baseline by  
 227 0.10 for top-5 and by approximately 0.15 for top-10. Additionally, when employing entity frequency,  
 228 Lmdb achieves an F-score that is around 0.17 higher than the baseline for top-5 and 0.07 higher  
 229 for top-10. The results demonstrate that ESBM exhibits a strong bias towards entity, reverse entity,  
 230 and relation frequency. For Map score, we are exactly observing the same behavior for ESBM. We  
 231 believe that the bias comes from the fact that the datasets are small, their second-hop neighborhood  
 232 is not considered, and the relations between their first-hop neighbors are not considered. On the  
 233 other hand, Figure 3 shows the F-score for top-5, top-10 and dynamic F-score on WIKES. Since  
 234 the length of summaries varies with the abstract, we calculate the F-score for each seed node based  
 235 on its summary length. Results show that WIKES F-score is close to random for different statistics,  
 236 thus rejecting the hypothesis of obvious biases. We observe a minor bias towards node frequency in  
 237 small datasets. Yet, as we increase the size of the dataset, this bias disappears. We observe a similar  
 238 behavior with MAP in Figure 4 Furthermore, we use *the entire* Wikidata to measure the F-score for  
 239 our seed nodes. Thus, importantly, we observe that our dataset’s F-score trend is comparable to that  
 240 of the entire data, especially our large dataset. We also extracted the first-hop neighborhood of all our  
 241 seed nodes and observed a small bias in the F-score top-5 and dynamic F-score. We conclude that

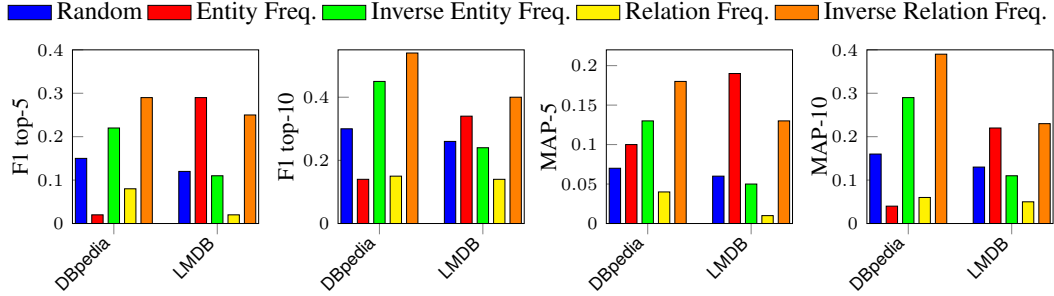


Figure 2: F1 score and MAP for frequency statistics on ESBM datasets.

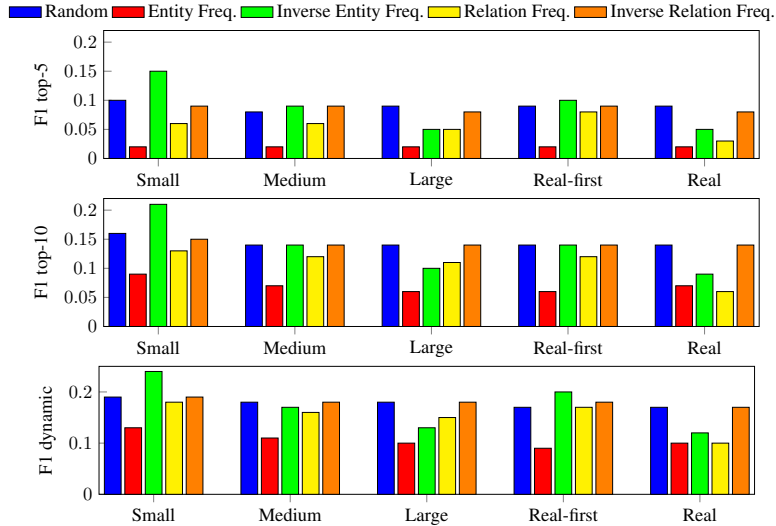


Figure 3: F1 for frequency statistics on WikiProFem.

242 adding the two-hop neighborhood makes the sample follow the graph distribution. Thus, WIKES is  
 243 an unbiased benchmark that retains the source KG distribution.

244 We evaluate the performance of different entity summarization methods on our benchmark, and  
 245 provide all implementations in the WIKES GitHub repository.

246 • **PageRank** Ma et al. [2008] ranks nodes in a graph based on the structure of incoming links, with  
 247 the idea that more important nodes are likely to receive more links from other nodes.

248 • **RELIN** Cheng et al. [2011] is a weighted PageRank algorithm that evaluates the relevance of  
 249 triples within a graph structure. We have re-implemented this model according to the specifications  
 250 in the referenced paper. On our smaller dataset version, RELIN takes approximately 6 hours to  
 251 compute all summaries.

252 • **LinkSum** Thalhammer et al. [2016] is a two-step, relevance-centric method that combines PageR-  
 253 ank with an adaptation of the Backlink algorithm to identify relevant connected entities. We have  
 254 re-implemented it according to the paper. The LinkSum method initially takes 10 hours to compute  
 255 the backlinks for each node in the small version of our dataset. By parallelizing the implementation,  
 256 we reduced this to one hour. Additionally, the Backlink algorithm itself initially takes 100 minutes,  
 257 but with parallelization, this was reduced to 10 minutes for the small version of our dataset.

258 Due to the inefficiency of the methods, we use a smaller version of WIKES for evaluation. The results  
 259 in Table 2 show that LinkSum outperforms both RELIN and PageRank. These findings suggest that



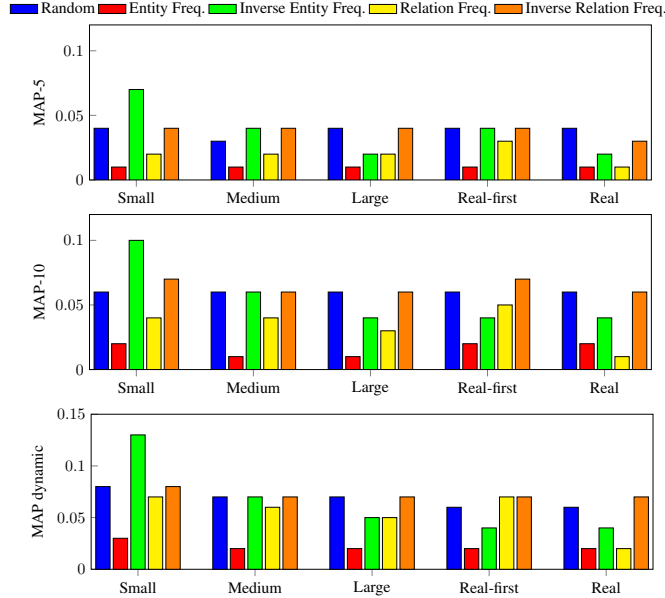


Figure 4: MAP for frequency statistics on WikiProFem.

260 models capable of exploiting the graph structure while handling large-scale datasets and maintaining  
 261 high accuracy in entity summarization are valuable for such real-world KGs, such as WIKES.

Model	Dataset	topK = 5		topK = 10		Dynamic	
		F-Score	MAP	F-Score	MAP	F-Score	MAP
PageRank	WikiLitArt	0.024	0.01	0.081	0.02	0.175	0.046
	WikiCinema	0.003	0.001	0.041	0.005	0.146	0.028
	WikiPro	0.060	0.02	0.169	0.049	0.288	0.109
	WikiProFem	0.032	0.01	0.093	0.024	0.145	0.036
RELIN	WikiLitArt	0.093	0.035	0.148	0.054	0.208	0.080
	WikiCinema	0.071	0.023	0.127	0.038	0.209	0.068
	WikiPro	0.125	0.053	0.200	0.086	0.273	0.127
	WikiProFem	0.111	0.050	0.179	0.081	0.219	0.095
LinkSum	WikiLitArt	0.184	0.080	0.239	0.109	0.225	0.127
	WikiCinema	0.119	0.048	0.152	0.060	0.135	0.068
	WikiPro	0.249	0.127	0.347	0.190	0.350	0.242
	WikiProFem	0.195	0.097	0.236	0.127	0.213	0.136

Table 2: Performance comparison of entity summarization models on the small version of WIKES. The models are evaluated with different topK values (5 and 10) and a dynamic setting.

## 262 5 Conclusion

263 We introduce WIKES (Wiki Entity Summarization Benchmark), a benchmark for KG entity summarization  
 264 which provides a scalable dataset generator that eschews the need for costly human annotation.  
 265 WIKES uses Wikipedia abstracts for automatic summary generation, ensuring contextually rich and  
 266 unbiased summaries. It preserves the complexity and integrity of real-world KGs through a random  
 267 walk sampling method that captures the structure of entities down to their second-hop neighborhoods.  
 268 Empirical evaluations demonstrate that WIKES provides high-quality large-scale datasets for entity  
 269 summarization tasks, and that it captures the complexities of knowledge graphs in terms of topology,  
 270 making it a valuable resource for evaluating and improving entity summarization algorithms.

271 **References**

- 272 Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. Dbpedia:  
273 A nucleus for a web of open data. In *The semantic web*, pages 722–735. Springer, 2007.
- 274 Gong Cheng, Thanh Tran, and Yuzhong Qu. Relin: Relatedness and informativeness-based centrality for entity  
275 summarization. In *The Semantic Web – ISWC 2011*, pages 114–129, 2011.
- 276 Gong Cheng, Qingxia Liu, and Yuzhong Qu. Generating characteristic summaries for entity descriptions. *TKDE*,  
277 35(5):4825–4835, 2023.
- 278 Kalpa Gunaratna, Krishnaprasad Thirunarayan, and Amit P Sheth. Faces: diversity-aware entity summarization  
279 using incremental hierarchical conceptual clustering. In *Proceedings of the AAAI Conference on Artificial*  
280 *Intelligence (AAAI)*, pages 116–122, 2015.
- 281 Oktie Hassanzadeh and Mariano P Consens. Linked movie data base. In *LDOW*, 2009.
- 282 Sen Hu, Lei Zou, Jeffrey Xu Yu, Haixun Wang, and Dongyan Zhao. Answering natural language questions by  
283 subgraph matching over knowledge graphs. *IEEE Transactions on Knowledge and Data Engineering*, 30(5):  
284 824–837, 2018.
- 285 Yunshi Lan, Shuohang Wang, and Jing Jiang. Multi-hop knowledge base question answering with an iterative  
286 sequence matching model. In *2019 IEEE International Conference on Data Mining (ICDM)*, pages 359–368,  
287 2019. doi: 10.1109/ICDM.2019.00046.
- 288 Morgane Laouenan, Palaash Bhargava, Jean-Benoît Eyméoud, Olivier Gergaud, Guillaume Plique, and Etienne  
289 Wasmer. A cross-verified database of notable people, 3500bc-2018ad. *Scientific Data*, 9(1):290, 2022.
- 290 Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian  
291 Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, et al. Dbpedia—a large-scale, multilingual  
292 knowledge base extracted from wikipedia. *Semantic Web*, 6(2):167–195, 2015.
- 293 Qingxia Liu, Gong Cheng, Kalpa Gunaratna, and Yuzhong Qu. ESBM: an entity summarization benchmark. In  
294 *The Semantic Web: 17th International Conference, ESWC 2020, Heraklion, Crete, Greece, May 31–June 4,*  
295 *2020, Proceedings 17*, pages 548–564. Springer, 2020.
- 296 Qingxia Liu, Gong Cheng, Kalpa Gunaratna, and Yuzhong Qu. Entity summarization: State of the art and future  
297 challenges. *Journal of Web Semantics*, 69:100647, 2021.
- 298 Nan Ma, Jiancheng Guan, and Yi Zhao. Bringing pagerank to the citation analysis. *Information Processing &*  
299 *Management*, 44(2):800–810, 2008.
- 300 Karl Pearson. The problem of the random walk. *Nature*, 72(1867):342–342, 1905.
- 301 Hadas Raviv, Oren Kurland, and David Carmel. Document retrieval using entity-based language models. In  
302 *Proceedings of the 39th International ACM SIGIR Conference on Research and Development in Information*  
303 *Retrieval, SIGIR ’16*, page 65–74, New York, NY, USA, 2016. Association for Computing Machinery. ISBN  
304 9781450340694. doi: 10.1145/2911451.2911508.
- 305 Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller,  
306 faster, cheaper and lighter. *ArXiv*, abs/1910.01108, 2019. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:203626972)  
307 [CorpusID:203626972](https://api.semanticscholar.org/CorpusID:203626972).
- 308 Fabian M. Suchanek, Gjergji Kasneci, and Gerhard Weikum. Yago: a core of semantic knowledge. In *Proceedings*  
309 *of the 16th International Conference on World Wide Web, WWW ’07*, page 697–706, New York, NY, USA,  
310 2007. Association for Computing Machinery. ISBN 9781595936547. doi: 10.1145/1242572.1242667.
- 311 Andreas Thalhammer, Nelia Lasierra, and Achim Rettinger. Linksum: Using link analysis to summarize entity  
312 data. In *ICWE*, pages 244–261, 2016.
- 313 Denny Vrandečić and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. *Communications of the*  
314 *ACM*, 57(10):78–85, 2014.
- 315 Hongwei Wang, Fuzheng Zhang, Xing Xie, and Minyi Guo. Dkn: Deep knowledge-aware network for news  
316 recommendation. In *Proceedings of the 2018 World Wide Web Conference*, page 1835–1844. International  
317 World Wide Web Conferences Steering Committee, 2018. ISBN 9781450356398. doi: 10.1145/3178876.  
318 3186175.

Table 3: Generated Datasets in terms of number of entities  $|\mathcal{V}|$ , triples  $|\mathcal{E}|$ , ground-truth summaries, density as  $|\mathcal{E}|/\binom{|\mathcal{V}|}{2}$ , graph connectivity, number of components, sampling method to select the entities and the subgraph, minimum and maximum node degree and, running time.

(a) Small Datasets

Metric	WikiLitArt	WikiCinema	WikiPro	WikiProFem
Entities ( $ \mathcal{V} $ )	85 346	70 753	79 825	79 926
Relations ( $ \mathcal{E} $ )	136 950	126 915	125 912	123 193
Target Entities	494	493	493	468
Density	0.000018	0.000018	0.000019	0.000019
Sampling method	Random Walk	Random Walk	Random Walk	Random Walk
Connected-graph	Yes	Yes	Yes	Yes
Num-comp	1	1	1	1
Min Degree	1	1	1	1
Max Degree	2172	3005	2060	3142
Run-time (seconds)	91.934	118.014	126.119	177.63

(b) Medium Datasets

Metric	WikiLitArt	WikiCinema	WikiPro	WikiProFem
Entities ( $ \mathcal{V} $ )	128 061	101 529	119 305	122 728
Relations ( $ \mathcal{E} $ )	220 263	196 061	198 663	196 838
Target Entities	494	493	493	468
Density	0.000013	0.000019	0.000014	0.000013
Sampling method	Random Walk	Random Walk	Random Walk	Random Walk
Connected-graph	Yes	Yes	Yes	Yes
Num-comp	1	1	1	1
Min Degree	1	1	1	1
Max Degree	3726	5124	3445	5282
Run-time (seconds)	155.36	196.413	208.157	301.718

(c) Large Datasets

Metric	WikiLitArt	WikiCinema	WikiPro	WikiProFem
Entities ( $ \mathcal{V} $ )	239 491	185 098	230 442	248 012
Relations ( $ \mathcal{E} $ )	466 905	397 546	412 766	413 895
Target Entities	494	493	493	468
Density	0.000008	0.00001	0.000008	0.000007
Sampling method	Random Walk	Random Walk	Random Walk	Random Walk
Connected-graph	Yes	Yes	Yes	Yes
Num-comp	1	1	1	1
Min Degree	1	1	1	1
Max Degree	8599	12189	7741	12939
Run-time (seconds)	353.113	475.679	489.409	768.99

Dataset	Seed Nodes Categories
WikiLitArt	<b>Entire graph:</b> actor=150, composer=35, film=41, novelist=24, painter=59, poet=39, screenwriter=17, singer=72, writer=57
	<b>Train:</b> actor=105, composer=24, film=29, novelist=17, painter=42, poet=27, screenwriter=12, singer=50, writer=40
	<b>Val:</b> actor=23, composer=5, film=6, novelist=4, painter=9, poet=6, screenwriter=2, singer=11, writer=8
	<b>Test:</b> actor=22, composer=6, film=6, novelist=3, painter=8, poet=6, screenwriter=3, singer=11, writer=9
WikiCinema	<b>Entire graph:</b> actor=405, film=88
	<b>Train:</b> actor=284, film=61
	<b>Val:</b> actor=59, film=14
	<b>Test:</b> actor=62, film=13
WikiPro	<b>Entire graph:</b> actor=58, football=156, journalist=14, lawyer=16, painter=23, player=25, politician=125, singer=27, sport=21, writer=28
	<b>Train:</b> actor=41, football=109, journalist=10, lawyer=11, painter=16, player=17, politician=87, singer=19, sport=15, writer=20
	<b>Val:</b> actor=9, football=23, journalist=2, lawyer=3, painter=3, player=4, politician=19, singer=4, sport=3, writer=4
	<b>Test:</b> actor=8, football=24, journalist=2, lawyer=2, painter=4, player=4, politician=19, singer=4, sport=3, writer=4
WikiProFem	<b>Entire graph:</b> actor=141, athletic=25, football=24, journalist=16, painter=16, player=32, politician=81, singer=69, sport=18, writer=46
	<b>Train:</b> actor=98, athletic=18, football=17, journalist=9, painter=13, player=22, politician=57, singer=48, sport=14, writer=34
	<b>Val:</b> actor=21, athletic=4, football=3, journalist=4, painter=1, player=5, politician=13, singer=11, sport=1, writer=5
	<b>Test:</b> actor=22, athletic=3, football=4, journalist=3, painter=2, player=5, politician=11, singer=10, sport=3, writer=7

Table 4: Seed nodes categories for each dataset. "Entire graph" refers to using the seed nodes and generating the data without train-test-val splits. In train-test-val, each of the datasets is a single weakly connected graph.

320 Table 5 presents the versions of the technologies and configurations that we use in this work.

Table 5: Technology and Configuration Details for Dataset Generations

(a) Technologies Used: Software Versions and Data Sources

Technology	Version/Details
Java	Version 21
Spring Boot	Version 3
Docker	Version 24.0.8
Python	Version 3.10
PostgreSQL	Version 16.3
Neo4j	Version 5.20.0-community
Wikipedia XML Article Dump Files	Published by Wikimedia on 2023/05/01
Wikidata XML Article Dump Files	Published by Wikimedia on 2023/05/01

(b) Pre-processing Setup: Specifications of the AWS EC2 Instance (r5a.4xlarge) Used for Dataset Pre-processing

Specification	Details
vCPU	16 (AMD EPYC 7571, 16 MiB cache, 2.5 GHz)
Memory	128 GB (DDR4, 2667 MT/s)
Storage	500 GB (EBS, 2880 Max Bandwidth)

## 321 **NeurIPS Paper Checklist**

### 322 **1. Claims**

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

325 Answer: [Yes]

326 Justification: To support our claims in the introduction and abstract, we provide experiments in  
327 section 4

328 Guidelines:

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

### 338 **2. Limitations**

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

340 Answer: [No]

341 Justification: We do not have this information because assessing the limitations of a dataset can be  
342 challenging. One clear limitation is that our data is generated from an encyclopedic knowledge graph,  
343 and we are uncertain about its suitability for specific domains. However, we have made a concerted  
344 effort to diversify the topics covered.

345 Guidelines:

- 346 • The answer NA means that the paper has no limitation while the answer No means that the paper  
347 has limitations, but those are not discussed in the paper.
- 348 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 349 • The paper should point out any strong assumptions and how robust the results are to violations of  
350 these assumptions (e.g., independence assumptions, noiseless settings, model well-specification,  
351 asymptotic approximations only holding locally). The authors should reflect on how these  
352 assumptions might be violated in practice and what the implications would be.
- 353 • The authors should reflect on the scope of the claims made, e.g., if the approach was only tested  
354 on a few datasets or with a few runs. In general, empirical results often depend on implicit  
355 assumptions, which should be articulated.
- 356 • The authors should reflect on the factors that influence the performance of the approach. For  
357 example, a facial recognition algorithm may perform poorly when image resolution is low or  
358 images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide  
359 closed captions for online lectures because it fails to handle technical jargon.
- 360 • The authors should discuss the computational efficiency of the proposed algorithms and how  
361 they scale with dataset size.
- 362 • If applicable, the authors should discuss possible limitations of their approach to address problems  
363 of privacy and fairness.
- 364 • While the authors might fear that complete honesty about limitations might be used by reviewers  
365 as grounds for rejection, a worse outcome might be that reviewers discover limitations that  
366 aren't acknowledged in the paper. The authors should use their best judgment and recognize  
367 that individual actions in favor of transparency play an important role in developing norms that  
368 preserve the integrity of the community. Reviewers will be specifically instructed to not penalize  
369 honesty concerning limitations.

### 370 **3. Theory Assumptions and Proofs**

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

373 Answer:[N/A]

374 Justification: We do not have proofs, as our focus is on empirical evaluations. We compare our dataset  
375 with real-world data in section 4.

376 Guidelines:

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

#### 386 4. Experimental Result Reproducibility

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

390 Answer: [Yes]

391 Justification: The model, dataset, and instructions for running the models are available in our GitHub  
 392 repository which is public.

393 Guidelines:

- 394 • The answer NA means that the paper does not include experiments.
- 395 • If the paper includes experiments, a No answer to this question will not be perceived well by the
- 396 reviewers: Making the paper reproducible is important, regardless of whether the code and data
- 397 are provided or not.
- 398 • If the contribution is a dataset and/or model, the authors should describe the steps taken to make
- 399 their results reproducible or verifiable.
- 400 • Depending on the contribution, reproducibility can be accomplished in various ways. For
- 401 example, if the contribution is a novel architecture, describing the architecture fully might suffice,
- 402 or if the contribution is a specific model and empirical evaluation, it may be necessary to either
- 403 make it possible for others to replicate the model with the same dataset, or provide access to
- 404 the model. In general, releasing code and data is often one good way to accomplish this, but
- 405 reproducibility can also be provided via detailed instructions for how to replicate the results,
- 406 access to a hosted model (e.g., in the case of a large language model), releasing of a model
- 407 checkpoint, or other means that are appropriate to the research performed.
- 408 • While NeurIPS does not require releasing code, the conference does require all submissions
- 409 to provide some reasonable avenue for reproducibility, which may depend on the nature of the
- 410 contribution. For example
  - 411 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to
  - 412 reproduce that algorithm.
  - 413 (b) If the contribution is primarily a new model architecture, the paper should describe the
  - 414 architecture clearly and fully.
  - 415 (c) If the contribution is a new model (e.g., a large language model), then there should either be
  - 416 a way to access this model for reproducing the results or a way to reproduce the model (e.g.,
  - 417 with an open-source dataset or instructions for how to construct the dataset).
  - 418 (d) We recognize that reproducibility may be tricky in some cases, in which case authors are
  - 419 welcome to describe the particular way they provide for reproducibility. In the case of
  - 420 closed-source models, it may be that access to the model is limited in some way (e.g.,
  - 421 to registered users), but it should be possible for other researchers to have some path to
  - 422 reproducing or verifying the results.

#### 423 5. Open access to data and code

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

426 Answer: [Yes]

427 Justification: The model, dataset, and instructions for running the models are available in our GitHub  
 428 repository which is public.

429 Guidelines:

- 430 • The answer NA means that paper does not include experiments requiring code.
- 431 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 432

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

## 448 6. Experimental Setting/Details

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

451 Answer: [Yes]

452 Justification: Hyperparameters and data splits are explained both in the paper and our github repository.  
453 We detail the characteristics of the methods and provide implementations.

454 Guidelines:

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

## 459 7. Experiment Statistical Significance

460 Question: Does the paper report error bars suitably and correctly defined or other appropriate informa-  
461 tion about the statistical significance of the experiments?

462 Answer: [No]

463 Justification: We did not repeat the experiments to report the confidence intervals.

464 Guidelines:

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

## 484 8. Experiments Compute Resources

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

487 Answer: [Yes]

488 Justification: In 3.4 and the appendix, we provide information about the running time for producing  
489 datasets. Additionally, our GitHub repository contains detailed information about the technologies we  
490 used.

491 Guidelines:

- 492 • The answer NA means that the paper does not include experiments.
- 493 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud  
494 provider, including relevant memory and storage.
- 495 • The paper should provide the amount of compute required for each of the individual experimental  
496 runs as well as estimate the total compute.
- 497 • The paper should disclose whether the full research project required more compute than the  
498 experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into  
499 the paper).

## 500 9. Code Of Ethics

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

503 Answer:[Yes]

504 Justification:We have reviewed the NeurIPS code of ethics.

505 Guidelines:

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

## 511 10. Broader Impacts

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

514 Answer: [N/A]

515 Justification: The paper does not involve societal impacts as it primarily focuses on foundational  
516 research in knowledge graph entity summarization, without direct application to scenarios that would  
517 cause societal impact.

518 Guidelines:

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

## 540 11. Safeguards

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



544 Answer: [N/A]  
545 Justification: The paper poses no risks associated with the release of data or models, as it focuses  
546 on foundational research in knowledge graph entity summarization without generating or releasing  
547 high-risk data or models.  
548 Guidelines:  
549 • The answer NA means that the paper poses no such risks.  
550 • Released models that have a high risk for misuse or dual-use should be released with necessary  
551 safeguards to allow for controlled use of the model, for example by requiring that users adhere to  
552 usage guidelines or restrictions to access the model or implementing safety filters.  
553 • Datasets that have been scraped from the Internet could pose safety risks. The authors should  
554 describe how they avoided releasing unsafe images.  
555 • We recognize that providing effective safeguards is challenging, and many papers do not require  
556 this, but we encourage authors to take this into account and make a best faith effort.

557 **12. Licenses for existing assets**  
558 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper,  
559 properly credited and are the license and terms of use explicitly mentioned and properly respected?  
560 Answer: [Yes]  
561 Justification: We cite the paper and re-implement the techniques.  
562 Guidelines:  
563 • The answer NA means that the paper does not use existing assets.  
564 • The authors should cite the original paper that produced the code package or dataset.  
565 • The authors should state which version of the asset is used and, if possible, include a URL.  
566 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.  
567 • For scraped data from a particular source (e.g., website), the copyright and terms of service of  
568 that source should be provided.  
569 • If assets are released, the license, copyright information, and terms of use in the package should  
570 be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for  
571 some datasets. Their licensing guide can help determine the license of a dataset.  
572 • For existing datasets that are re-packaged, both the original license and the license of the derived  
573 asset (if it has changed) should be provided.  
574 • If this information is not available online, the authors are encouraged to reach out to the asset's  
575 creators.

576 **13. New Assets**  
577 Question: Are new assets introduced in the paper well documented and is the documentation provided  
578 alongside the assets?  
579 Answer: [Yes]  
580 Justification: We provide our Github repository and datasets publicly.  
581 Guidelines:  
582 • The answer NA means that the paper does not release new assets.  
583 • Researchers should communicate the details of the dataset/code/model as part of their sub-  
584 missions via structured templates. This includes details about training, license, limitations,  
585 etc.  
586 • The paper should discuss whether and how consent was obtained from people whose asset is  
587 used.  
588 • At submission time, remember to anonymize your assets (if applicable). You can either create an  
589 anonymized URL or include an anonymized zip file.

590 **14. Crowdsourcing and Research with Human Subjects**  
591 Question: For crowdsourcing experiments and research with human subjects, does the paper include  
592 the full text of instructions given to participants and screenshots, if applicable, as well as details about  
593 compensation (if any)?  
594 Answer: [N/A]  
595 Justification: The paper does not involve crowdsourcing or research with human subjects.  
596 Guidelines:  
597 • The answer NA means that the paper does not involve crowdsourcing nor research with human  
598 subjects.

599 • Including this information in the supplemental material is fine, but if the main contribution of the  
600 paper involves human subjects, then as much detail as possible should be included in the main  
601 paper.

602 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other  
603 labor should be paid at least the minimum wage in the country of the data collector.

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

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

608 Answer: [N/A]

609 Justification: The paper does not involve crowdsourcing or research with human subjects.

610 Guidelines:

611 • The answer NA means that the paper does not involve crowdsourcing nor research with human  
612 subjects.

613 • Depending on the country in which research is conducted, IRB approval (or equivalent) may be  
614 required for any human subjects research. If you obtained IRB approval, you should clearly state  
615 this in the paper.

616 • We recognize that the procedures for this may vary significantly between institutions and  
617 locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for  
618 their institution.

619 • For initial submissions, do not include any information that would break anonymity (if applica-  
620 ble), such as the institution conducting the review.