# WIKI ENTITY SUMMARIZATION BENCHMARK

Anonymous authors

Paper under double-blind review

# Abstract

Entity summarization aims to compute concise summaries for entities in knowledge graphs. However, current datasets and benchmarks are often limited to only a few hundred entities and overlook knowledge graph structure. This is particularly evident in the scarcity of ground-truth summaries, with few labeled entities available for evaluation and training. We propose WIKES (Wiki Entity Summarization Benchmark), a large *benchmark* comprising of entities, their summaries, and their connections. Additionally, WIKES features a dataset *generator* to test entity summarization algorithms in different subgraphs of the knowledge graph. Importantly, our approach combines graph algorithms and NLP models, as well as different data sources such that WIKES does not require human annotation, rendering the approach cost-effective and generalizable to multiple domains. Finally, WIKES is scalable and capable of capturing the complexities of knowledge graphs in terms of topology and semantics. WIKES features existing *datasets* for comparison. Empirical studies of entity summarization methods confirm the usefulness of our benchmark. Data, code, and models are available at: https://anonymous.4open.science/ r/Wikes-2DDA/README.md.

024

000

001 002 003

004

005 006 007

800

010

011

012

013

014

015

016

017

018

019

020

021

# 1 INTRODUCTION

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

Yet, the size and complexity of these KGs progressively expand, rendering it increasingly challenging to convey the essential information about an entity in a concise and meaningful way Suchanek et al. (2007); 033 Vrandečić & Krötzsch (2014). This is where entity summarization becomes relevant. Entity summarization (ES) Liu et al. (2021) is the process of generating a concise and informative summary that captures the most 035 salient aspects of the entity, based on the information available in the KGs. In ES, the entity *description* refers 036 to all the triples involving such an entity. For instance, Figure 1 illustrates a set of relationships surrounding 037 the entity Ellen Johnson Sirleaf in a KG, along with a possible summary for this entity. Extensive 038 descriptions can overwhelm users and exceed the capacity of typical user interfaces, making it challenging to identify the most relevant triples. Entity summarization addresses this issue by computing an optimal 040 compact summary for an entity, selecting a size-constrained subset of triples Liu et al. (2021).

Despite advancements in entity summarization techniques Liu et al. (2021), their development and evaluation face significant limitations in current benchmarks and datasets Liu et al. (2020); Cheng et al. (2023). First, existing benchmarks comprise only a few hundred entities, limiting dataset size. Second, generating groundtruth summaries primarily relies on costly and time-consuming manual annotation, which can introduce bias based on the preferences and knowledge of a few annotators. Lastly, these benchmarks often overlook the rich information contained within the knowledge graph structure.

047	-
048	United States of 2011 Nobel 2011 Uberian
049	General general election
050	Shabara Provention doctorate Decomposition Contractio
051	Nobel
052	Peace Prize ward received Prize Pri
053	avaid received educated at encircled at Community of
054	Liberia contrit of clibereshin
055	place of birth Eller Johnson Sifeal Unrole U
056	Morroria Characteria Vige
057	President
058	significant event significant event sex or gender member of notificial out/
059	Paradise sward received
060	Grand Cross of Grand Scots of Grand
061	Honour
062	
063	Figure 1: KG subgraph of entity Ellen Johnson Sirleat: arrows depict the subgraph of relationships to other entities and labels indicate their roles. Selecting the hold edges as entity summaries of the most relevant triples may
064	reduce information overload while concisely describing the entity.
065	
066	To address the above limitations, we propose:
067	
068	• Novel wikes benchmark for ES based on summaries and graphs from wikidata and wikipedia.
069	• Subgraph extraction method preserving the complexity of real-world KOS, subsampling using random walks and proportionally preserving node degrees. WIKES captures the structure of the entities up to the
070	second-hop neighborhood thereby ensuring that the connections in WIKES accurately reflect those in the
071	source KG.
072	• Comprehensive summaries for any entity in the KG, ensuring that summaries are both relevant and
073	contextually rich by deriving them directly from corresponding Wikipedia abstracts, minimizing human
074	bias, as these abstracts are created and reviewed by several experts. In this manner, WIKES is scalable,
075	enabling it to generate large benchmark resources efficiently with high-quality annotation.
076	• Automatic entity summarization dataset generator allows for the creation of arbitrarily large datasets,
077	encompassing various domains of knowledge.
078	
079	2 EXISTING DATASETS
080	
081	Here, we review the existing datasets for entity summarization. Table 1 provides an overview and statistics of
082	the current datasets in this field. FACES and INFO datasets have a higher density than the entities in the
083	Entity Summarization Benchmark (ESBM). It is also clear that LMDB and FACES are not connected graphs,
084	that challenge graph-based learning methods where the information cannot easily propagate in disconnected
085	networks. Specifically, FACES consists of 12 connected components, which complicates the learning process
086	tor graph embedding methods by limiting the richness of information that can be leveraged from the graph.
087	

- 880
- We provide here a comprehensive description of each dataset or benchmark:
- 089 • ESBM Liu et al. (2020): The Entity Summarization Benchmark (ESBM) is the first benchmark to evaluate 090 the performance of entity summarization methods. ESBM has three versions; v1.2 is the latest and most 091 extensive version. This version comprises 175 entities, with 150 from DBpedia Lehmann et al. (2015) and 25 from LinkedMDB Hassanzadeh & Consens (2009). The summaries comprise triples selected by 30 092 "researchers and students" annotators. Each entity has exactly 6 summaries. Despite encompassing two 093

Table 1: Existing and WIKES datasets; number of entities  $|\mathcal{V}|$ , triples  $|\mathcal{E}|$ , ground-truth summaries, density , number of components, sampling method, min/max node degree, and time to generate the dataset excluding preprocessing. RW refers to the random walk sampling method.

(a) Existing Datasets

#### (c) WIKES Small Datasets

85346

136950 494

18e-6 RW

2172 91.9

(d) WIKES Large Datasets

WikiLitArt WikiCinema WikiPro WikiProFem

79825

19e-6 RW

2060 126.1

125912 493

70753

18e-6 RW

3005 118.0

126915 493

79926 123193 468

19e-6 RW

3142 177.6

Metric	DBpedia (ESBM)	LMDB (ESBM)	FACES	INFO	Metric
$ \mathcal{V} $	2721	1853	1379	1410	
$ \mathcal{E} $	4436	2148	2152	2019	E
Ground-truth	125	50	50	100	Ground-truth
Density	5e-4	6e-4	11e-4	10e-4	Density
Sampling	-	-	-	-	Sampling
Components	1	2	12	1	Components
Min Deg	1	1	1	1	Min Deg
Max Deg	125	208	88	100	Max Deg
Max Deg	125	208	88	100	Max Deg Generation

#### (b) WIKES Medium Datasets

107	(b) WIKES Medium Datasets					(d) WIKES Large Datasets				
108	Metric	WikiLitArt	WikiCinema	WikiPro	WikiProFem	Metric	WikiLitArt	WikiCinema	WikiPro	WikiProFem
109	$ \mathcal{V} $	128061	101529	119305	122728	$ \mathcal{V} $	239491	185098	230442	248012
440	E	220263	196061	198663	196838	E	466905	397546	412766	413895
110	Ground-truth	494	493	493	468	Ground	-truth 494	493	493	468
	Density	13e-6	19e-6	14e-6	13e-6	Density	8e-6	10e-6	8e-6	7e-6
111	Sampling	RW	RW	RW	RW	Samplin	ig RW	RW	RW	RW
	Components	1	1	1	1	Compor	ients 1	1	1	1
112	Min Deg	1	1	1	1	Min De	<u>o</u> 1	1	1	1
	Max Deg	3726	5124	3445	5282	Max De	g 8599	12189	7741	12939
113	Generation Time (s)	155.4	196.4	208.2	301.7	Generat	ion Time (s) 353.1	475.7	489.4	769.0

datasets, ESBM has several limitations. First, the entity sampling method is not explained. In particular, some triples in the neighborhood of the entity are missing in the datasets. Second, there are no connections among the entities in the neighborhood, nor any two-hop neighborhood. Third, the expertise and background of the annotators are not assessed nor disclosed. Due to the expensive annotation process, the dataset size is small.

• FACES Gunaratna et al. (2015) is a dataset from DBpedia (version 3.9)? and includes 50 randomly selected 119 entities, each with at least 17 different types of relations. Similar to ESBM, the FACES ground-truth is also 120 generated manually. 121

• INFO Cheng et al. (2023) features 100 randomly selected entities from 10 classes in DBpedia, including 122 two sets of ground-truth summaries: REF-E and REF-W. REF-E summaries are crafted from triples by five 123 experts with a 140-character limit, resembling Google search result snippets. In contrast, REF-W summaries 124 are derived from one expert who reads Wikipedia abstracts and selects closely related neighboring entities. 125 The number of ground-truth summaries per entity varies due to multiple evaluations by some experts, 126 complicating the evaluation process. Additionally, the expertise of the annotators is not specified.

127 In contrast, our benchmark uses Wikidata to automatically map entities from Wikipedia to Wikidata. This 128 automation allows us to efficiently generate summaries for any number of entities. Unlike previous work, 129 we use the Wikipedia abstract as a summary instead of manual annotators. Each abstract is a collaboration 130 of many users; as such, it should not introduce obvious biases. Additionally, with this process, we ensure 131 high-quality and cost-effective summaries. Furthermore, we present the characteristics of our dataset in 132 Table 1. The WIKES benchmark contains a significantly more entities and relations compared to existing datasets. Starting from approximately 500 target nodes, WIKES samples a connected graph, whereas existing 133 datasets include at most 125 target nodes. Besides, LMDB and FACES are not connected graphs. 134

135 136

137

094

095

096

098

100

101

102

104

105

106

114

115

116

117

118

#### THE WIKES BENCHMARK 3

138 A Knowledge Graph  $\mathcal{KG} = (\mathcal{V}, \mathcal{R}, \mathcal{T})$  is a directed multigraph consisting of entities  $\mathcal{V} = \{v_1, \ldots, v_n\}$ , 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 v_i, v_j \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{R} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{V} \land \exists r \in \mathcal{V} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{V} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{V} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{V} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{V} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \land \exists r \in \mathcal{V} \text{ s.t. } (v_i, r, v_j) \in \mathcal{V} \text{ s.t. } (v_i, r, v_$ 139  $\mathcal{T}$  contains pairs of nodes connected by a relationship. 140

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 ignoring edge directions.

144 A summary for an entity  $v_i$  is a subset  $S(v_i) \subseteq \Delta_t(v_i)$  of triples from the *t*-description of  $v_i$ , where the 145 *t*-description of an entity  $v_i \in \mathcal{V}$  in a knowledge graph  $\mathcal{KG}$  is the set  $\Delta_t(v_i) = \{(s, p, o) \in \mathcal{T} \mid s \in \mathcal{N}_t(v_i) \lor o \in \mathcal{N}_t(v_i)\}$  of triples in which one of the entities is in the *t*-hop neighborhood of  $v_i$ .

147 Entity summarization Liu et al. (2021) for an entity  $v_i \in \mathcal{V}$  in a knowledge graph  $\mathcal{KG}$  aims to find a summary 148  $\mathcal{S}(v_i)$  that maximizes a relevance score among all possible summaries for  $v_i$ , i.e.,

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

The scoring functions differ among entity summarization methods, with some focusing on centrality and diversity of neighbors Cheng et al. (2011) and others employ PageRank-like scores Thalhammer et al. (2016).

#### 155 156

149

150

151 152

153

154

# 3.1 EXTRACTING SUMMARIES FROM WIKIDATA USING WIKIPEDIA ABSTRACTS

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

Wikidata is a free and collaborative knowledge base that collects structured data to support Wikipedia and other Wikimedia projects. It includes descriptions and labels for entities. The descriptions offer in-depth details, while the labels serve as concise identifiers, facilitating efficient data retrieval and integration in subsequent steps. We load all Wikidata items XML dump files published on 2023/05/01<sup>1</sup> as entities  $\mathcal{V}$ alongside their properties as relationships  $\mathcal{R}$  into a graph database<sup>2</sup>. The result is a graph that connects all Wikidata items and statements. We include items if they (1) are not marked as redirects, (2) belong to the main Wikidata item and property, including labels and descriptions, into a relational database<sup>3</sup>.

Wikipedia pages contain infoboxes, abstracts, page content, categories, references, and more. Links to other Wikipedia pages are referred to as mentions. We detect these mentions in the abstracts and infoboxes of Wikipedia pages to use them later for labeling the summaries in Wikidata. We extract and load all the content from the XML dump files of Wikipedia pages, published on 2023/05/01<sup>4</sup>, into a relational database under the same conditions as Wikidata: the pages must be in English and not redirected.

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

Wikidata is a directed multigraph, which means that each entity (Wikidata item) can be connected to another
entity via multiple relations (Wikidata properties). Yet, links in Wikipedia are not labeled; as such, we need to
select one of the relations for the summary. To annotate the correct Wikidata property as part of the summary,
we employ the DistilBERT model Sanh et al. (2019). DistilBERT is a fast and lightweight model with a
reduced number of parameters compared to the original BERT model. This way, we can efficiently process
large amounts of data while maintaining high-quality embeddings for accurate relation selection.

<sup>184 &</sup>lt;sup>1</sup>https://dumps.wikimedia.org/wikidatawiki/

<sup>185 &</sup>lt;sup>2</sup>https://neo4j.com

<sup>186 &</sup>lt;sup>3</sup>https://www.postgresql.org/

<sup>187 &</sup>lt;sup>4</sup>https://dumps.wikimedia.org/enwiki/

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

# 194 3.2 CAPTURING THE GRAPH STRUCTURE

193

195

213

216

Here we introduce the WIKES generator algorithm. The main idea is to sample a connected graph that
preserves the original graph structure. To this end, we employ random walks Pearson (1905). The random
walk model is a straightforward yet effective method for preserving graph structure. While more recent
techniques may yield superior results, we choose to use this widely accepted and fundamentally sound
approach that exhibits good results even with 1% sampled nodes (Figure 3).

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

<sup>204</sup> The process is a Markov process, characterized by its memoryless property:

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

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

$$N_{i} = \frac{L_{i} - \min(\{L\})}{\max(\{L\}) - \min(\{L\})}$$
(3)

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

$$R_i = \text{round}\left(\min RW + N_i \times (\max RW - \min RW)\right) \tag{4}$$

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

Moreover, the random walk sampling process requires a set of seed nodes as a starting point. In our case, 223 the seed nodes represent the target entities we are interested in. The seed nodes can be any Wikidata Item 224 Identifier, Wikipedia title, or Wikipedia ID of the Wikipedia pages. We collect the seed nodes on the condition 225 that they have at least k (default k = 5) common entities with the abstract section and the infobox in the 226 Wikipedia pages. Therefore, this model is flexible, allowing you to choose any seed nodes from any domain 227 as an input. In the datasets that we generated, we collect seed nodes from Laouenan et al. (2022). This paper 228 has published information about individuals from various domains. The authors collected data from multiple 229 Wikipedia editions and Wikidata, using deduplication and cross-verification techniques to compile a database 230 of 1.6 million individuals with English Wikipedia pages. The seed nodes that we use include actor, athletic, 231 football, journalist, painter, player, politician, singer, sport, writer, lawyer, film, composer, novelist, poet, and 232 screenwriter. Using combinations of these seed nodes, we generate four sets of datasets, with each set having 233 small, medium, and large versions. In Table 8 in Section A in the supplementary material, we present the 234 seed nodes and their proportions for each dataset and their corresponding train-test-val splits.

#### 235 3.3 WIKES GENERATOR 236

237 We discuss how WIKES is created, and how further benchmarks can be generated without the need for manual 238 annotators. Algorithm 1 details the generator, which consists of the following steps.

239 Step1: Retrieve summaries of each seed node (explained in Section 3.1) 240

**Step2:** Expand the graph using the random walk method in Section 3.2. Set the random walk's length n241 (default n = 2), which means it explores up to the *n*-hop neighborhood of each seed node. We choose 242 n = 2 because extending beyond two hops risks introducing irrelevant entities, while our approach balances 243 efficiency and accuracy. This ensures scalability and relevance for large datasets like Wikidata, complementing 244 existing benchmarks Lissandrini et al. (2018). 245

Step3: Check if the graph is connected. If it is, done. If not, identify all disconnected components and sort 246 them by size, from largest to smallest. In each iteration, connect smaller components to the largest component 247 using h connections. Utilize the shortest path method, selecting paths that are equal to or less than a minimum 248 path length l. Continue connecting nodes from the smaller component to the larger one until h nodes are 249 connected. After each iteration, check graph connectivity again. If all components are connected to the largest 250 component, the algorithm ends. Otherwise, re-sort components and increase l by 1. Repeat until the graph is 251 a single connected component. 252

### Algorithm 1 WIKES Generator

253

270

271 272

254	Alg	
255	1:	<b>Input:</b> Graph $G$ , seed nodes $S$ , random walk length $n$ , minimum path length $l$
256	2:	Output: A connected graph
230	3. 4.	<b>procedure</b> GENERALEURAPH(G, S, $n, l)$
257	-+. 5.	summaries $\leftarrow$ REIREVEDUMMARIES(S)
258	5.	$G \leftarrow KANDOM WALKEAPANSION(G, S, n)$ included in section 5.2
200	7.	$v_{s}$ connected $\leftarrow$ CheckConnectiviti(G) while not is connected do
259	8:	while not $s_{control = control = 0}$
260	9. 9.	Sort components by size in descending order
	10:	$larest \leftarrow components[0]$
261	11:	for comp in components[1]: do
262	12:	Connect comp to largest using h connections via shortest paths $< l$
060	13:	$G \leftarrow \text{UPDATEGRAPH}(G, comp, largest)$
203	14:	$is\_connected \leftarrow CHECKCONNECTIVITY(G)$
264	15:	if is_connected then
265	16:	break
205	17:	end if
266	18:	end for
267	19:	$l \leftarrow l + 1$
201	20:	end while
268	21:	return G
269	22:	end procedure

#### 3.4 WIKES DATASETS

We generate three sizes for each of the four datasets, obtaining 12 datasets. We present their characteristics in 273 Table 1 in section A. The number of entities in the small datasets ranges from approximately 70k to 85k, and 274 the number of relations ranges from around 120k to 135k. In the medium datasets, the number of entities 275 ranges from 100k to 130k, and the number of relations ranges from 195k to 220k. The number of entities 276 in the large datasets ranges from approximately 185k to 250k, and the number of relations ranges from 277 around 397k to 470k. The average runtime for generating small graphs is approximately 128 seconds; for 278 medium-sized graphs, it is approximately 216 seconds; and for large graphs, it is approximately 512 seconds. 279 We construct the train-test-validation split for each dataset with 70% for training, 15% for testing, and 15%280 for validation. Detailed information about the run time, as well as the number of nodes and relations for these splits, is available on our GitHub repository. All graphs in each train-test-validation splits are connected. 281

### 4 EMPIRICAL EVALUATION

282

283 284

285

286

298 299 300

301

310

We study the quality of WIKES using the following metrics:

**F-Score.** Let  $S_m$  the summary obtained by a summarization method and  $S_h$  the ground-truth summary. We compare  $S_m$  with  $S_h$  using the F1-score based on precision P and recall R:

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

The F1 score lies within [0,1]. High F1 indicates that  $S_m$  is closer to the ground-truth  $S_h$ .

292 Mean Average Precision (MAP). This metric is particularly suitable for evaluating ranking tasks because 293 it takes into account the order of the predicted triples. MAP calculates precision at each position i in the 294 predicted summary and averages these values over all relevant summary triples. It reflects both the relevance 295 and the ranking quality of the predicted summaries. MAP, unlike F1-score, does not depend on a specific 296 value of k. This makes it a robust metric for assessing how well a summarization method ranks the relevant 297 triples.

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

where N is the total number of entities,  $S_h^{(n)}$  is the set of ground-truth summary triples for a particular entity 302 303  $v_n, S_m^{(n)}$  is the set of predicted summary triples for the entity  $v_n$ , Precision@i is the precision at the i-th 304 position in the predicted summary, and Rel(n, i) indicates whether the *i*-th predicted triple for entity  $v_n$  is 305 relevant (i.e., it belongs to  $S_h^{(n)}$ ). MAP scores are in the range [0,1], where a higher MAP indicates better 306 performance in terms of correctly predicting relevant summary triples. To account for the varying lengths 307 of the ground-truth summaries in real-world data, we also calculate MAP and F-score (which we refer to 308 as dynamic MAP and dynamic F-score) by setting the length of the generated summary ( $|S_m|$ ) equal to the 309 length of the corresponding ground-truth summary ( $|S_h|$ ).

We analyze our dataset and compare it with the ESBM benchmark using statistical measures such as frequency 311 and inverse frequency of entities and relations. We calculate the F-score and MAP score for the top-5 and 312 top-10 of both the ESBM dataset and our WikiProFem. We choose top-5 and top-10 because we only have 313 ground-truth summaries for top-5 and top-10 in the ESBM dataset. The F-score and MAP results for ESBM 314 are presented in Figure 2. The statistics show that for DBpedia, the F-score using inverse relation frequency 315 outperforms the random baseline by 0.15 for top-5 and by 0.34 for top-10. Furthermore, when using inverse 316 entity frequency, DBpedia achieves an even higher F-score, surpassing the random baseline by 0.07 for top-5 and by 0.15 for top-10. For LMDB, we observe a similar trend when using inverse frequency. The F-score 317 surpasses the random baseline by 0.10 for top-5 and by approximately 0.15 for top-10. Additionally, when 318 employing entity frequency, LMDB achieves an F-score that is around 0.17 higher than the baseline for top-5 319 and 0.07 higher for top-10. The results demonstrate that ESBM exhibits a strong bias towards entity, reverse 320 entity, and relation frequency. For Map score, we are exactly observing the same behavior for ESBM. We 321 believe that the bias comes from the fact that the datasets are small, their second-hop neighborhood is not 322 considered, and the relations between their first-hop neighbors are not considered. On the other hand, Figure 3 323 shows the F-score for top-5, top-10 and dynamic F-score on WIKES. Since the length of summaries varies 324 with the abstract, we calculate the F-score for each seed node based on its summary length. Results show that 325 WIKES F-score is close to random for different statistics, thus rejecting the hypothesis of obvious biases. We 326 observe a minor bias towards node frequency in small datasets. Yet, as we increase the size of the dataset, this bias disappears. We observe a similar behavior with MAP in Figure 11 (in appendix). Furthermore, we use *the entire* Wikidata to measure the F-score for our seed nodes. Thus, importantly, we observe that our 328

7

329

330

331

332

333

364

365

dataset's F-score trend is comparable to that of the entire data, especially our large dataset. We also extracted the first-hop neighborhood of all our seed nodes and observed a small bias in the F-score top-5 and dynamic F-score. We conclude that adding the two-hop neighborhood makes the sample follow the graph distribution. Thus, WIKES is an unbiased benchmark that retains the source KG distribution.



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

- **PageRank** Ma et al. (2008) is an unsupervised method that ranks nodes in a graph based on the structure of incoming links, with the idea that more important nodes are likely to receive more links from other nodes.
- RELIN Cheng et al. (2011) is another unsupervised approach, a weighted PageRank algorithm that evaluates the relevance of triples within a graph structure. We have re-implemented this model according to the specifications in the referenced paper. On our smaller dataset version, RELIN takes approximately 6 hours to compute all summaries.
- LinkSum Thalhammer et al. (2016), also an unsupervised approach, is a two-step, relevance-centric method that combines PageRank with an adaptation of the Backlink algorithm to identify relevant connected entities.
   We have re-implemented it according to the paper. The LinkSum method initially takes 10 hours to compute the backlinks for each node in the small version of our dataset. By parallelizing the implementation, we

reduced this to one hour. Additionally, the Backlink algorithm itself initially takes 100 minutes, but with
 parallelization, this was reduced to 10 minutes for the small version of our dataset.

GATES Firmansyah et al. (2021) is a recent *supervised* approach that integrates graph structure using
 Graph Attention Networks with knowledge graphs and text embeddings. We run GATES using the best
 performing hyperparameters of the original paper. GATES takes 20 minutes to run on the our small datasets.

We evaluate the methods on the smallest WIKES dataset due to their inefficiency. Table 2 shows that LinkSum, generally outperforms other models. Interestingly, GATES, despite being supervised, achieves lower accuracy compared to LinkSum. The deficiency in GATES may be due to its reliance on the frequency of nodes and relations, which are used as weights. As mentioned earlier, this frequency bias is present in ESBM but not in WIKES.These results highlight the significance of graph structure in summarizing entities within real-world knowledge graphs like WIKES, emphasizing the advantages of graph-based methods.

388 Efficiency concerns. The evaluation we conducted on various entity summarization models reveals significant 389 efficiency issues with many recent baselines. For example, BAFREC Kroll et al. (2018), a model highlighted 390 in a recent survey on unsupervised entity summarization, was unable to process a graph with 13 000 nodes which is  $5 \times$  smaller than our smallest dataset — even after running for two days. Similarly, MPSUM Wei 391 et al. (2020) did not finish after 15 days on the same graph. Additionally, models like INFO Cheng et al. 392 (2023), which depend on unavailable external resources, were excluded from our evaluation. These results 393 highlight the need for more scalable approaches that can efficiently handle large knowledge graphs without 394 sacrificing performance or accuracy. 395

		topK = 5		topK	= 10	Dynamic	
Model	Dataset	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
GATES	WikiLitArt	0.110	0.052	0.167	0.087	0.236	0.090
	WikiCinema	0.085	0.036	0.131	0.051	0.231	0.082
	WikiPro	0.149	0.074	0.225	0.118	0.313	0.149
	WikiProFem	0.128	0.062	0.227	0.097	0.243	0.114

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.

## 5 CONCLUSION

396

411

412 413

414 415 We introduce WIKES (Wiki Entity Summarization Benchmark), a benchmark for KG entity summarization 416 which provides a scalable dataset generator that eschews the need for costly human annotation. WIKES uses 417 Wikipedia abstracts for automatic summary generation, ensuring contextually rich and unbiased summaries. 418 It preserves the complexity and integrity of real-world KGs through a random walk sampling method that 419 captures the structure of entities down to their second-hop neighborhoods. Empirical evaluations demonstrate 420 that WIKES provides high-quality large-scale datasets for entity summarization tasks, and that it captures 421 the complexities of knowledge graphs in terms of topology, making it a valuable resource for evaluating and 422 improving entity summarization algorithms.

#### 423 424 REFERENCES

436

440

451

- Gong Cheng, Thanh Tran, and Yuzhong Qu. Relin: Relatedness and informativeness-based centrality for entity summarization. In *The Semantic Web ISWC 2011*, pp. 114–129, 2011.
- Gong Cheng, Qingxia Liu, and Yuzhong Qu. Generating characteristic summaries for entity descriptions. *TKDE*, 35(5):
   4825–4835, 2023.
- Asep Fajar Firmansyah, Diego Moussallem, and Axel-Cyrille Ngonga Ngomo. Gates: Using graph attention networks for entity summarization. In *Proceedings of the 11th Knowledge Capture Conference*, pp. 73–80, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450384575. doi: 10.1145/3460210.3493574. URL https://doi.org/10.1145/3460210.3493574.
- Kalpa Gunaratna, Krishnaprasad Thirunarayan, and Amit P Sheth. Faces: diversity-aware entity summarization using incremental hierarchical conceptual clustering. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pp. 116–122, 2015.
- 437 Oktie Hassanzadeh and Mariano P Consens. Linked movie data base. In *LDOW*, 2009.
- Sen Hu, Lei Zou, Jeffrey Xu Yu, Haixun Wang, and Dongyan Zhao. Answering natural language questions by subgraph
   matching over knowledge graphs. *IEEE Transactions on Knowledge and Data Engineering*, 30(5):824–837, 2018.
- Hermann Kroll, Denis Nagel, and Wolf-Tilo Balke. Bafrec: Balancing frequency and rarity for entity characterization in linked open data. In *Proceedings of the 1st International Workshop on Entity REtrieval (EntRE)*, 2018.
- Yunshi Lan, Shuohang Wang, and Jing Jiang. Multi-hop knowledge base question answering with an iterative sequence matching model. In *2019 IEEE International Conference on Data Mining (ICDM)*, pp. 359–368, 2019. doi: 10.1109/ICDM.2019.00046.
- Morgane Laouenan, Palaash Bhargava, Jean-Benoît Eyméoud, Olivier Gergaud, Guillaume Plique, and Etienne Wasmer.
   A cross-verified database of notable people, 3500bc-2018ad. *Scientific Data*, 9(1):290, 2022.
- Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, et al. Dbpedia–a large-scale, multilingual knowledge base extracted from wikipedia. *Semantic Web*, 6(2):167–195, 2015.
- Matteo Lissandrini, Davide Mottin, Themis Palpanas, and Yannis Velegrakis. Multi-example search in rich information graphs. In 2018 IEEE 34th International Conference on Data Engineering (ICDE), pp. 809–820. IEEE, 2018.
- 454 Qingxia Liu, Gong Cheng, Kalpa Gunaratna, and Yuzhong Qu. Esbm: an entity summarization benchmark. In
   455 *The Semantic Web: 17th International Conference, ESWC 2020, Heraklion, Crete, Greece, May 31–June 4, 2020,* 456 *Proceedings 17*, pp. 548–564. Springer, 2020.
- 457 Qingxia Liu, Gong Cheng, Kalpa Gunaratna, and Yuzhong Qu. Entity summarization: State of the art and future
   458 challenges. *Journal of Web Semantics*, 69:100647, 2021.
- 459
   460
   461
   461
   Nan Ma, Jiancheng Guan, and Yi Zhao. Bringing pagerank to the citation analysis. *Information Processing & Management*, 44(2):800–810, 2008.
- 462 Karl Pearson. The problem of the random walk. *Nature*, 72(1867):342–342, 1905.
- Hadas Raviv, Oren Kurland, and David Carmel. Document retrieval using entity-based language models. In *Proceedings* of the 39th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR
  '16, pp. 65–74, New York, NY, USA, 2016. Association for Computing Machinery. ISBN 9781450340694. doi: 10.1145/2911451.2911508.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *ArXiv*, abs/1910.01108, 2019. URL https://api.semanticscholar.org/CorpusID: 203626972.

470	
	Fabian M. Suchanek, Gjergji Kasneci, and Gerhard Weikum. Yago: a core of semantic knowledge. In <i>Proceedings of the</i>
471	16th International Conference on World Wide Web, WWW '07, pp. 697–706, New York, NY, USA, 2007. Association
472	for Computing Machinery. ISBN 9781595936547. doi: 10.1145/1242572.1242667.
473	

- Andreas Thalhammer, Nelia Lasierra, and Achim Rettinger. Linksum: Using link analysis to summarize entity data. In *ICWE*, pp. 244–261, 2016.
- Denny Vrandečić and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, 57 (10):78–85, 2014.
- Hongwei Wang, Fuzheng Zhang, Xing Xie, and Minyi Guo. Dkn: Deep knowledge-aware network for news recommendation. In *Proceedings of the 2018 World Wide Web Conference*, pp. 1835–1844. International World Wide Web Conferences Steering Committee, 2018. ISBN 9781450356398. doi: 10.1145/3178876.3186175.
  - Dongjun Wei, Shiyuan Gao, Yaxin Liu, Zhibing Liu, and Longtao Hang. Mpsum: entity summarization with predicatebased matching. arXiv preprint arXiv:2005.11992, 2020.