

Entity Alignment For Knowledge Graphs: Progress, Challenges, and Empirical Studies

Deepak Chaurasiya^{1†}, Anil Surisetty^{1†}, Nitish Kumar¹, Alok Singh¹, Vikrant Dey^{1,2}, Aakarsh Malhotra¹, Gaurav Dhama¹, Ankur Arora¹

¹AI Garage, Mastercard India

²IIT Roorkee

Abstract

Entity Alignment (EA) identifies entities across databases that refer to the same entity. Knowledge graph-based embedding methods have recently dominated EA techniques. Such methods map entities to a low-dimension space and align them based on their similarities. With the corpus of EA methodologies growing rapidly, this paper presents a comprehensive analysis of various existing EA methods, elaborating their applications and limitations. Further, we distinguish the methods based on their underlying algorithms and the information they incorporate to learn entity representations. Based on challenges in industrial datasets, we bring forward 4 research questions (RQs). These RQs empirically analyse the algorithms from the perspective of *Hubness*, *Degree distribution*, *Non-isomorphic neighbourhood*, and *Name bias*. For Hubness, where one entity turns up as the nearest neighbour of many other entities, we define an *h*-score to quantify its effect on the performance of various algorithms. Additionally, we try to level the playing field for algorithms that rely primarily on name-bias existing in the benchmarking open-source datasets by creating a low name bias dataset. We further create an open-source repository for 14 embedding-based EA methods and present the analysis for invoking further research motivations in the field of EA.

1 Introduction

The rise of the internet and multimedia has generated an exponential amount of data. To interpret and use this data efficiently, it needs to be stored in a structured manner. Knowledge Graphs ($\mathcal{KG}s$) [4; 17; 10; 12] are one such mechanism to store data, leading to its widespread usage. Due to asynchronous data sources, an entity often appears with different aliases in different $\mathcal{KG}s$. This has led to the rise of EA task, which aims to align entities from different $\mathcal{KG}s$ alluding to the same real-world entity. Lack of generalization of early methods [11; 19; 35; 52] and advance-

ment in \mathcal{KG} representation learning [20; 32; 13; 18; 22; 43; 29] has led to a rise in embedding based EA methods.

This paper details the working of various existing embedding-based EA methods and compares their applications and limitations. Traditionally these methods are classified into two classes: (i) Trans-based methods and (ii) GNN-based methods. Trans-based methods learn the representations based on the translation constraint of *head+relation = tail* [5]. In contrast, GNN-based methods use Graph Neural Networks to learn structural entity embeddings. In this paper, along with the traditional classification, we propose a categorisation based on the type of network information they exploit— (i) Structure aware (S), (ii) Structure and Relation aware (SR), and (iii) Structure and Attribute aware (SA).

This paper complements existing surveys [3; 49; 50] by providing a comprehensive categorisation of the latest embedding-based EA methods along with inferences based on four RQs, with a focus on EA in the industrial setting. An existing study by Sun et al. [41] provides a survey of EA methods but lacks discussion around the presence of biases in industrial datasets and its effect on models' performance. Our study proposes a novel graph sampling method to generate a low name bias dataset required for rigorous evaluation of methods and provides a metric to quantify the effect of performance adversary phenomena like hubness. Moreover, we provide a comparative evaluation of the latest EA methods specific to the listed RQs to show which methods are better suited for a given industrial setting.

- **RQ1. Does the method address the hubness problem present in $\mathcal{KG}s$?** Hubness [19] in EA is the phenomenon where some entities (hubs) get aligned to many other entities in the vector space. In this study, we propose an *h*-score that measures hubness's extent for an alignment result.
- **RQ2. Does the method align entities well irrespective of their degree?** The higher the degree of an entity, the higher the structural information present for the same and vice versa. In this paper, we demonstrate the effect of degree on the alignment accuracy of each algorithm.
- **RQ3. Does the method align entities having non-isomorphic subgraphs in the two $\mathcal{KG}s$?** Due to distinct sources, an identical entity in two $\mathcal{KG}s$ has non-identical neighbourhoods or, equivalently, non-isomorphic sub-

[†]Equal Contribution

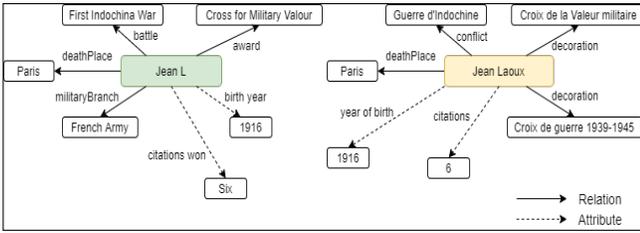


Figure 1: English (left) and French (right) sub-graph for *Jean L*

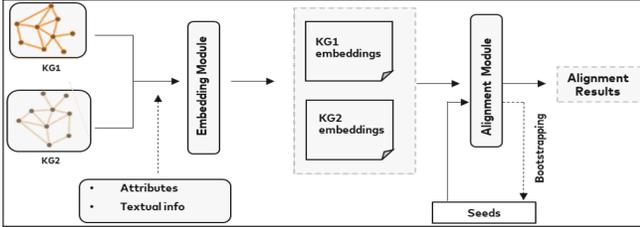


Figure 2: General framework of embedding based EA methods.

graphs. We quantify an algorithm’s effectiveness against non-isomorphism by analyzing accurately aligned entities having different neighbourhoods.

- **RQ4. Does name bias lead to an algorithm’s high accuracy?** An algorithm’s matching result shows name bias when the alignment relies primarily on linguistic information in entities’ names. We propose a cross-lingual benchmarking dataset with low name bias to highlight an algorithm’s effectiveness.

Our major contributions are listed as follows:

- A survey of existing methodologies in chronological order and their classification based on the information they capture for EA task.
- Extensive experiments to compare the efficacy of existing methods in an industrial setting by drawing empirical inferences based on the listed RQs.
- Propose a new graph sampling technique that generates a \mathcal{KG} by sampling entity pairs with low-name bias while keeping graph properties unchanged. Further, we generate and publish a cross-lingual low name-bias dataset for robust testing of EA methods.
- To ensure reproducibility, we publish the implementation of the listed methods in an open-source library*.

2 Problem Definition

Here, we introduce preliminary notation that will be consistent for the rest of the paper. A \mathcal{KG} is denoted by $\mathcal{KG}=(\mathcal{T}, \mathcal{A}, \mathcal{V}, \mathcal{E}, \mathcal{R})$, where \mathcal{T} represents the union of relation and attribute triples, \mathcal{A} represents the set of attributes, and \mathcal{E} represents that of entities. All attribute values and relations form the set \mathcal{V} and \mathcal{R} , respectively. General notations used in this study are defined in Table 1.

*Link for the open-source library will be provided post review

Table 1: Notations

Notation	Description
h, r, t	Relation triple ($h, t \in \mathcal{E}, r \in \mathcal{R}$)
$f_t(h, r, t)$	$\ h + r - t\ $
T	Set of positive triples i.e $\mathcal{T}_1 \cup \mathcal{T}_2$
T'	Set of negative triples
h', r', t'	A triple $\in T'$
$l_{ti}(h, r, t)$	$[\gamma + f_t(h, r, t) - f_t(h', r', t')]_+$
$l_{cl}(h, r, t)$	$[f_t(h, r, t) - \gamma_1]_+ + [\gamma_2 - f_t(h', r', t')]_+$
\mathcal{L}_{ti}	$\sum_{(h,r,t) \in T} \beta \sum_{h',r',t' \in T'} l_{ti}(h, r, t)$
\mathcal{L}_{cl}	$\sum_{(h,r,t) \in T} \beta \sum_{h',r',t' \in T'} l_{cl}(h, r, t)$

Table 2: Classification of EA methods

	Methods	Struc	Rel	Attr	Name	Boot [†]
S	GMNN [48]	✓	-	-	-	-
	AliNet [40]	✓	-	-	-	-
SR	MTransE [6]	✓	✓	-	-	-
	IPTransE [51]	✓	✓	-	-	✓
	BootEA [39]	✓	✓	-	-	✓
	RSN4EA [15]	✓	✓	-	-	-
	HGCN [47]	✓	✓	-	✓	-
	RDGCN [46]	✓	✓	-	✓	-
	HyperKA [37]	✓	-	-	-	-
	MRAEA [25]	✓	✓	-	-	✓
SA	DualAMN [24]	✓	✓	-	-	✓
	JAPE [38]	✓	✓	✓	-	-
	AttrE [42]	✓	✓	✓	✓	-
	GCNAlign [44]	✓	✓	✓	-	-

Given two $\mathcal{KG}s$, \mathcal{KG}_1 and \mathcal{KG}_2 , EA aims to identify pairs (e_1, e_2) where $e_1 \in \mathcal{E}_1$ and $e_2 \in \mathcal{E}_2$ such that e_1 and e_2 denote the same entity. Some methods make use of pre-aligned set of entity pairs as training labels which hereby is referred to as seeds. Fig. 1 depicts sub-graphs centered around *Jean L* in English and French \mathcal{KG} . EA aims to align *Jean L* and *Jean Laoux* by using knowledge about *Jean* in both the $\mathcal{KG}s$.

3 Categorisation of Techniques

For a real-world EA task, it is essential to select an algorithm that optimally captures the information present in the real-world \mathcal{KG} . For instance, when entity names are absent in a real-world \mathcal{KG} , one would not prefer to use methods like HGCN [47] or RDGCN [46] since they primarily rely on entity name information for aligning entities. Hence, to ease the process of selecting a suitable method for the problem at hand, we propose categorising all such methodologies based on the information they utilise to align the entities. The methods mentioned above are segregated into 3 categories: (i) Structure aware (S), (ii) Structure and Relation aware (SR), and (iii) Structure and Attribute aware (SA). We now describe and compare existing EA methods belonging to each of these categories. These methods use the general framework of an embedding based EA model, as depicted in Fig. 2.

[†]Bootstrapping

Structure aware methods:

Few popular studies include GMNN [48] and AliNet [40] which use only structure information of $\mathcal{KG}s$ to align entities. GMNN introduced the concept of topic entity matching to EA, which aligns their neighbourhood subgraphs instead of entities. Its embedding module uses a 2-layer stacked GCN to encode a node’s neighbourhood structural information, and similar to [45] employs a cross-graph attention mechanism to learn cross-lingual \mathcal{KG} representations. These nodes’ representations are aggregated to generate a topic’s representation, used for final matching. In contrast, AliNet focuses on aligning entities with non-isomorphic neighbourhoods and proposes a multi-hop aggregation embedding module to capture the local and distant neighbourhoods. Its embedding module employs a GCN to generate the local structure representation. To generate distant structure representation, it utilizes an attention mechanism to calculate the contribution of the n -hop neighbourhood of an entity. These local and distant representations are aggregated using a gated layer [34]. Both methods do not consider relational and attribute information.

Structure and Relation aware methods:

This category includes all Trans-based algorithms since they utilise the *head + relation = tail* constraint while generating entity and relation embeddings. The embedding module of such methods, including MTransE [6], IPTransE [51] and BootEA [39] propose variations on the use of the above constraint. MTransE [6] embeds the nodes and relations of the two $\mathcal{KG}s$ by optimizing a constraint-based objective i.e. $f_t()$ loss, defined in Table 1. IPTransE [51], inspired from PTransE [23], proposed an enhanced learning method by capturing indirect relations between entities using multiple relational paths. It used a triplet-based loss function \mathcal{L}_{tl} where the negative samples are sampled uniformly. The embedding module of BootEA is similar to that of MTransE [39]. However, it utilises a contrastive loss function \mathcal{L}_{cl} , where the negative samples are generated using ϵ -truncated negative sampling [39]. Unlike MTransE, IPTransE and BootEA utilise bootstrapping where entity pairs are aligned during training and appended to the training set iteratively.

Differing from Trans-based methods, Guo et al. [15] proposed RSN4EA, which captures long term relation dependencies of entities. Inspired from [14], RSN4EA generates relational paths by performing biased random walks. These paths are encoded using Recurrent Skipping Network (RSN), with skip connections existing only for entities to distinguish between entities and relations. Guo et al. proposed type based NCE [16], to maximise the likelihood of predicting the following entity in the relational path using the current entity’s representation and the learned hidden state for this step.

Superior performance achieved by edge-aware GNNs in graph representation learning tasks has led to the rise of relation aware GNN based EA methods which include HGCN [47], RDGCN [46], HyperKA [37], MRAEA [25], and DualAMN [24]. The embedding module of HGCN uses a 2-layer gated GCN [21] to learn the structural embedding of entities, where a relation is represented as the average of associated entities’ embeddings. HGCN represents a node’s relational embedding as the average embedding of its associ-

ated relations to capture relational information. These structural and relational embeddings are concatenated to generate a node’s output representation. Inspired from [28], RDGCN attempts to jointly learn node and relation embeddings by transforming a \mathcal{KG} (primal) into its dual where each node denotes a particular primal relation, and each edge weight depicts the likelihood of two relations sharing the same head or tail entity in the primal graph. Its embedding module uses a primal-dual attention mechanism, where the primal attention layer captures edge weights using dual graph’s node representation and vice versa. Sun et al. proposed HyperKA [37] to learn representation in a low-dimensional hyperbolic manifold, which enables to capture hierarchies between entities. HyperKA has a translation based layer followed by n -GCN layers, learning representations in a hyperbolic manifold.

To further improve the captured relation information, Mao et al. [25] proposed MRAEA. MRAEA learns entity and relational representation by capturing relation’s semantic properties (direction, type, and inverse). Further, a joint representation of an entity (meta relation aware embedding) is represented as the average of neighbouring entities’ representations concatenated with associated relations’ average embedding. The embedding module of MRAEA employs a relation aware attention aggregation mechanism using meta relation aware node embeddings to perform a relation weighted neighbourhood aggregation to generate output embeddings.

Mao et al. further proposed a time-efficient EA method named DualAMN [24]. Its embedding uses a simplified relation aggregation layer that utilises a relation reflection operator [26], significantly decreasing the number of parameters required to learn relation-aware representations. Further, DualAMN learns cross \mathcal{KG} interaction by introducing n proxy nodes representing cross-graph alignment relation, which act as reference points for all the nodes in the two $\mathcal{KG}s$. This transforms the learning of node-node cross- \mathcal{KG} interaction problem to learning node-proxy interaction for every node, reducing time complexity from $O(|E_1||E_2|)$ to $O(|E_1| + |E_2|)$. The alignment module of DualAMN utilises LogExpSum [33; 36], an approximation of the ϵ -truncated negative sampling technique. LogExpSum enables parallelisation in negative sample generation, improving the overall time complexity.

Structure and Attribute aware methods:

The presence of non-isomorphism in $\mathcal{KG}s$ has led to the use of supplementary information, i.e. attributes of entities for EA. Recent studies include JAPE [38], AttrE [42] and GCNAlign [44]. Considering that attributes of pre-aligned entities are correlated, JAPE [38] generates attribute representation using a skip-gram [27]. The objective function aims at maximising attribute co-occurrence. The embedding module of JAPE is similar to that of MTransE, additionally utilising the attribute similarities for refining the learned embeddings by penalising L_2 distance between highly similar entities.

On the other hand, Trisdeya et al. proposed AttrE [42], an unsupervised EA approach. AttrE treats an attribute triple similar to a relation triple while learning representations. It identifies pseudo aligned entities and relations by employing a string matching algorithm on their names and use these

aligned entities as true alignment for its model. The embedding module of AttrE initialises attribute nodes with the machine-translated vector of its value and employs MTransE based embedding module to generate entity, relation and attribute embeddings. GCNAlign [44] also refrains from differentiating between relation and attribute triples and generates representation for both types of nodes in a unified space. It generates a relation and attribute-based adjacency matrix and pre-processes them to represent each edge by its tf-idf importance. GCNAlign’s embedding module utilises two parallel embedding modules, one for learning structure embeddings and one for learning attribute embeddings. The independent embedding modules are optimised separately, whereas the alignment module uses both structure and attribute embeddings for final alignment.

4 Implementation

In our study, we use a total of five datasets to analyze and answer our study’s RQs. These include two sparse cross-lingual datasets: (i) EN-FR_{V1} and (ii) EN-DE_{V1}, a dense cross-lingual dataset: (i) EN-FR_{V2}, a large cross-lingual dataset: (i) EN-FR_{100K}, and a low name bias cross-lingual dataset: (i) EN-FR_{LNB}. EN, FR, and DE stand for English, French and German languages, respectively. The first four datasets are from the study Sun et al. [41] where authors generate the same using Iterative Degree Sampling on DBpedia knowledge base [1], whose statistics are shown in Table 3.

We create the fifth dataset EN-FR_{LNB} (Low Name Bias) to answer RQ4. The entities in EN-FR_{LNB} have relatively low BERT [9] embedding similarity with their counterparts. The sampled dataset’s degree distribution should be close to source \mathcal{KG} for a fair comparison. Hence, entities are binned based on their degree, and a set of them are sampled uniformly from every bin. The size of the sampled set is proportional to the bin size. To reduce name bias, nodes in the selected set are dropped with a probability proportional to the similarity of its BERT embedding with its counterpart. Degree distributions of all the datasets are depicted in Fig.3(a).

All the methods are executed with 20% seeds for training, 10% for validation, and 70% for testing. For a fair comparison, we use a batch size of 5000 relation triples, maximum training epochs as 1000, and early stopping on validation Hits@1 with the patience of 10 epochs for all the methods. For other method-specific hyperparameters such as the number of layers, margins for triplet, and contrastive loss, values proposed by the respective author are followed.

5 Results

In this study, we demonstrate our results on 2 key metrics used in EA- **Hits@k and Mean Reciprocal Rank (MRR)**[8]. Table 4 lists the performance of existing methods on the aforementioned benchmarking datasets.

Dense v/s Sparse Dataset: In general, relation-aware methods perform better on EN-FR_{V2} due to their ability to better incorporate multiple relations associated with every entity. An exception to this is MTransE, where, as the dataset

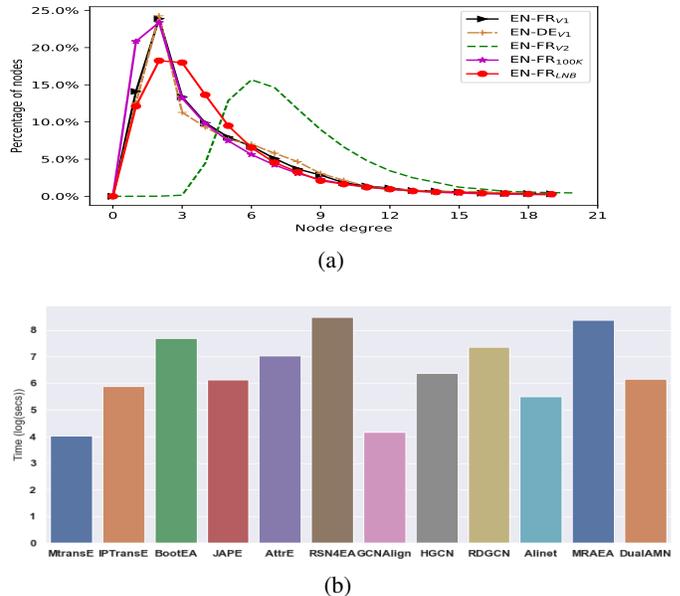


Figure 3: (a) Degree distributions of the five datasets. (b) Training time of existing methodologies on EN-FR_{V1}

becomes dense, the proportion of entities with multiple relations increases. This results in a drop in performance due to TransE’s inability to distinguish entities with similar multi-mapping relations. The performance of all GNN-based methods also improved due to their ability to capture complex structures around entities.

Training Time: Fig. 3(b) illustrates the training time of different methods. RSN4EA’s training time is relatively high (86× training time of MTransE) as it trains on a large number of triples created using biased random walks (5× relation triples in EN-FR_{V1}). BootEA’s negative sampling and bootstrapping are responsible for its high training time (39× training time of MTransE). In the case of RDGCN, alternating over the primal-dual graphs is the primary cause for high training time. Techniques like MTransE and GCNAlign are the fastest as their embedding modules do not incorporate attention or multi-hop relation triples. Amongst SOTA, DualAMN is the fastest due to fewer parameters, proxy cross graph attention and the use of an approximation of truncated uniform negative sampling. In general, supplementary information such as attributes, multi-hop paths, or techniques like ϵ -negative sampling boosts the performance but increases the training time simultaneously.

6 Analysis

In this section, we present our analysis towards the four RQs.

RQ1. Does the method address hubness problem in \mathcal{KG} s?

Most of the existing methods learn representations in high dimensional vector spaces. Consequently, the curse of dimensionality [2] leads to the formation of hub nodes [31] that become nearest neighbours of many smaller degree nodes. In

Table 3: Description of EN-FR_{V1}, EN-DE_{V1}, EN-FR_{V2}, EN-FR_{100K}, and EN-FR_{LNB}.

S. no.	Datasets	KGs	#Ent	#Rel	#Attr	#Rel triples	#Attr triples
1	EN-FR _{V1}	EN	15K	267	308	47,334	73,121
		FR	15K	210	404	40,864	67,167
2	EN-DE _{V1}	EN	15K	215	286	47,676	83,755
		DE	15K	131	194	50,419	156,150
3	EN-FR _{V2}	EN	15K	193	189	96,318	66,899
		FR	15K	166	221	80,112	68,779
4	EN-FR _{100K}	EN	100K	400	466	309,607	497,729
		FR	100K	300	519	258,285	426,672
5	EN-FR _{LNB}	EN	15K	284	273	56,088	75,805
		FR	15K	223	395	47,720	65,149

Table 4: Results of EA methods in terms of Hits@1, Hits@5 and MRR

	Dataset	EN-FR _{V1}			EN-DE _{V1}			EN-FR _{V2}			EN-FR _{100K}		
		Method	H@1	H@5	MRR	H@1	H@5	MRR	H@1	H@5	MRR	H@1	H@5
S	GMNN [48]	74.3	83.8	0.79	83.7	90.7	0.85	80.1	89.7	0.84	68.3	81.5	0.75
	AliNet [40]	36.9	59.3	0.47	60.1	75.8	0.67	54.5	80.5	0.65	33.6	57.1	0.44
SR	MTransE [6]	24.3	45.5	0.34	27.7	49.0	0.38	23.5	42.8	0.33	13.2	25.8	0.20
	IPTransE [51]	17.4	32.9	0.25	30.5	49.6	0.40	19.9	40.9	0.30	14.2	26.8	0.21
	BootEA [39]	50.5	71.6	0.60	65.5	80.8	0.72	66.1	85.4	0.75	36.1	55.4	0.47
	RSN4EA [15]	39.2	59.1	0.48	59.2	75.8	0.67	57.1	74.6	0.65	49.3	66.9	0.57
	HGCN [47]	77.9	85.5	0.81	82.7	87.5	0.84	87.5	93.1	0.90	74.3	77.1	0.76
	RDGCN [46]	75.1	86.4	0.80	84.1	91.5	0.87	81.9	90.9	0.86	71.0	78.2	0.74
	HyperKA [37]	43.5	66.6	0.54	52.2	74.4	0.62	62.2	85.4	0.76	41.8	63.4	0.51
	MRAEA [25]	53.9	77.2	0.64	70.2	86.6	0.77	74.0	92.4	0.82	38.0	58.8	0.47
	DualAMN [24]	58.2	78.2	0.67	77.4	92.0	0.83	78.8	94.6	0.85	49.2	68.9	0.58
SA	JAPE [38]	26.4	49.3	0.37	29.8	53.3	0.41	30.3	52.7	0.41	30.3	52.7	0.41
	AttrE [42]	45.2	63.6	0.54	48.0	65.9	0.56	47.0	69.0	0.57	45.4	62.8	0.53
	GCNAlign [44]	34.4	60.4	0.46	55.2	73.8	0.64	41.9	70.7	0.55	26.5	46.2	0.36

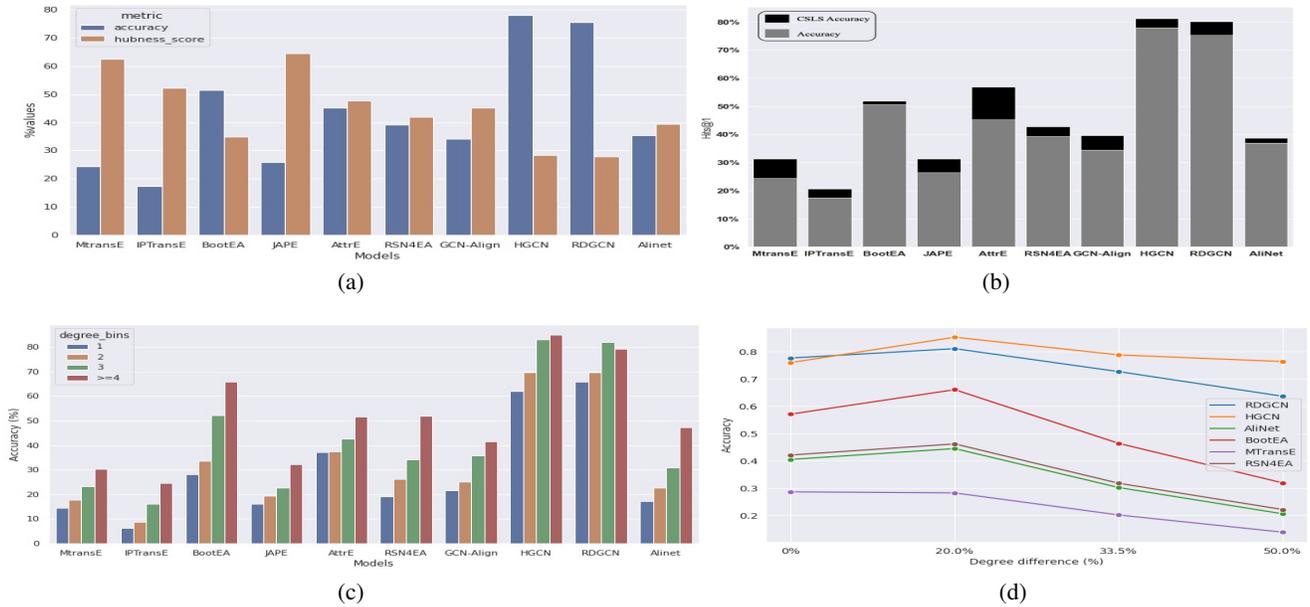


Figure 4: (a) Hits@1% and h -score on EN-FR_{V1} (b) Hits@1% with and without CSLS on EN-FR_{V1} (c) Hits@1% w.r.t. entity degree on EN-FR_{V1} (d) Hits@1% w.r.t. degree difference on EN-FR_{V1}

EA, we observe this phenomenon occurring where an entity gets aligned to many smaller degree entities. For instance,

in the alignment results of BootEA on the EN-FR_{V1} dataset, the French entity “Beillé” gets aligned to 31 English entities. To analyze the effect of hubness in the alignment results, we formulate a h -score (hubness score) as:

$$h\text{-score} = \sum_{z \in \mathbf{Z}} \text{hub}(z); \text{hub}(z) = \frac{\#\text{entities aligned with } z}{\#\text{entities in } \mathcal{KG}_i} \quad (1)$$

where, i is the \mathcal{KG} to which z belongs to and \mathbf{Z} is the set of largest 10% hubs (based on $\text{hub}()$). A high h -score indicates that a higher proportion of entities align with very few hub nodes, which results in low alignment accuracy in injective mapping. Fig. 4(a) lists the prediction accuracy and h -score of various existing EA methods. Trans-family EA algorithms, such as MTransE and JAPE, rely only on structural and relational information and have a high h -score, which shows a pronounced effect of hubness. Whereas algorithms that use external information (RDGCN, HGCN) decrease the effect of hubness due to the use of lingual information to differentiate an actual and a potential match (hubs). To overcome hubness, Conneau et al. [7] proposed Cross-domain Similarity Local Scaling (CSLS). CSLS aims to decrease the similarity associated with vectors lying in dense areas and is effective against hubness. The same can be observed in Fig. 4(b). To conclude, external information and modifications in the alignment strategy like CSLS can help suppress the effect of hubness in EA, consequently improving the overall accuracy.

RQ2. Does the method align the entities well irrespective of their degree?

Degree distribution of a \mathcal{KG} , shown in Fig. 3(a), reflects that 50% of the entities have degree < 4 . This negatively skewed distribution indicates that most entities have a low degree, implying less structural information to learn the embeddings. For a highly connected entity, many neighbouring entities influence their structural information, resulting in generalized representation, and for low degree entities, the representation is sub-optimal as only a few neighbours influence it.

Fig. 4(c) plots the alignment accuracy of algorithms for entities with a specific degree. The general trend reaffirms our hypothesis that higher degree entities have better alignment accuracy. Also, we observe that algorithms that use supplementary information of entities show a marginal difference in alignment accuracy of smaller and larger degree nodes. These include RDGCN and AttrE. Hence, we conclude that algorithms primarily relying on the structure can lead to sub-optimal performance for very sparse datasets. The use of node attributes and linguistic information to supplement structure can improve alignment results.

RQ3. Does the method align entities having non-isomorphic neighbourhood in the two \mathcal{KG} s?

\mathcal{KG} s are usually incomplete [30], leading to different degrees of identical entities in the two \mathcal{KG} s resulting in a non-isomorphic neighbourhood. We study the effectiveness of existing methods in aligning entities with different degrees. A negative degree difference would imply that a less connected entity in \mathcal{KG}_1 appears as a highly connected entity in \mathcal{KG}_2 and vice-versa.

Fig.4(d) shows how accuracy varies for different algorithms when aligning nodes with the non-isomorphic neighbourhood. RDGCN and HGCN perform degree-difference invariant alignment as they primarily utilize the name linguistic similarity of entities. Further, AliNet’s performance varies marginally with increasing non-isomorphism as the former uses attribute predicates and structural information while the latter looks at the multi-hop neighbourhood of an entity. Contrary to them, MTransE shows a drop in alignment accuracy, which reassures our hypothesis that algorithms relying solely on structural information struggle to align non-isomorphic entities. Hence, the use of complementary information about entities such as attribute triples, lingual information, or global network properties, along with structural information, is essential in handling non-isomorphic subgraphs.

Table 5: Hits@1% on EN-FR_{V1} & EN-FR_{LNB}

	AttrE	HGCN	RDGCN
EN-FR _{V1}	45.25	77.93	75.17
EN-FR _{LNB}	38.77	65.63	61.10
Difference	6.48	12.30	14.07

RQ4. Does name-bias in the data lead to an algorithm’s high accuracy?

Techniques that use textual information cannot be adjudged with other algorithms because of high similarities of linguistic information with counterparts. For instance, nearest neighbour search using BERT embeddings on EN-FR_{V1} gives Hits@1 of 0.69, indicating a very high name bias. Table 5 quantifies the sensitivity of EA techniques towards entity names. It compares the results on EN-FR_{V1} dataset against a low name biased EN-FR_{LNB} dataset having 0.57 Hits@1 with BERT embedding. The significant dip in the performance of RDGCN and HGCN reflects their strong dependency on entity names. We believe that these methods are unsuited for scenarios where the name similarity of entities in the two complementary \mathcal{KG} s does not bear any significance. For instance, consider the task of aligning debit cards, identified by their number, which belongs to the same cardholder. Such algorithms struggle in this scenario as the two debit card identifiers contains no linguistic similarity. Contrarily, AttrE sees a relatively less drop in performance as it uses multiple attribute information along with entity names. Therefore, we conclude that for industrial datasets with encryptions or low-name bias, methods such as AttrE should be preferred.

7 Conclusion

In this paper, we present a comprehensive overview of the various methods and techniques in the field of entity alignment between \mathcal{KG} s and conduct a benchmarking study of the representative approaches. We investigate the existing methods on four research questions inspired by limitations in real-world datasets. We propose a novel graph sampling technique to ensure robust testing of methods to sample a low name-bias dataset. Further, we develop an open-source library with multiple EA approaches and publish their implementations along with the aforementioned low-name bias dataset. This study

provides a more profound understanding of EA techniques and the effect of different experimental settings on their performance. The challenges and inferences presented in this work pave the path for future direction for the research community.

References

- [1] Auer, S., Bizer, C., Kobilarov, G., Lehmann, J., Cyganiak, R., Ives, Z.: Dbpedia: A nucleus for a web of open data. In: ISWC (2007)
- [2] Bellman, R.: Dynamic programming. In: AAAS (1966)
- [3] Berrendorf, M., Faerman, E., Melnychuk, V., Tresp, V., Seidl, T.: Knowledge graph entity alignment with gcns: Lessons learned. In: ECIR (2020)
- [4] Bollacker, K., Evans, C., Paritosh, P., Sturge, T., Taylor, J.: Freebase: A collaboratively created graph database for structuring human knowledge. In: SIGMOD (2008)
- [5] Bordes, A., Usunier, N., Garcia-Durán, A., Weston, J., Yakhnenko, O.: Translating embeddings for modeling multi-relational data. In: NIPS (2013)
- [6] Chen, M., Tian, Y., Yang, M., Zaniolo, C.: Multi-lingual knowledge graph embeddings for cross-lingual knowledge alignment. In: IJCAI (2016)
- [7] Conneau, A., Lample, G., Ranzato, M., Denoyer, L., Jégou, H.: Word translation without parallel data. In: ICLR (2018)
- [8] Craswell, N.: Mean reciprocal rank. In: EDS (2009)
- [9] Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. In: NAACL (2019)
- [10] Dong, X.L., Gabrilovich, E., Heitz, G., Horn, W., Lao, N., Murphy, K., Strohmann, T., Sun, S., Zhang, W.: Knowledge vault: A web-scale approach to probabilistic knowledge fusion. In: SIGKDD (2014)
- [11] El-Roby, A., Aboulnaga, A.: Automatic link exploration in linked data. In: ICDE (2016)
- [12] Färber, M.: The microsoft academic knowledge graph: A linked data source with 8 billion triples of scholarly data. In: SEMWEB (2019)
- [13] Gesese, G.A., Biswas, R., Sack, H.: A comprehensive survey of knowledge graph embeddings with literals: Techniques and applications. In: ESWC (2019)
- [14] Grover, A., Leskovec, J.: node2vec: Scalable feature learning for networks. In: KDD (2016)
- [15] Guo, L., Sun, Z., Hu, W.: Learning to exploit long-term relational dependencies in knowledge graphs. In: ICML (2019)
- [16] Gutmann, M., Hyvärinen, A.: Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In: AIS (2010)
- [17] Hoffart, J., Suchanek, F.M., Berberich, K., Weikum, G.: Yago2: A spatially and temporally enhanced knowledge base from wikipedia: Extended abstract. In: Artificial Intelligence (2013)
- [18] Hogan, A., Blomqvist, E., Cochez, M., d'Amato, C., de Melo, G., Gutiérrez, C., Gayo, J.E.L., Kirrane, S., Neumaier, S., Polleres, A., Navigli, R., Ngomo, A.N., Rashid, S.M., Rula, A., Schmelzeisen, L., Sequeda, J.F., Staab, S., Zimmermann, A.: Knowledge graphs (2020)
- [19] Hu, W., Chen, J., Qu, Y.: A self-training approach for resolving object coreference on the semantic web. In: WWW (2011)
- [20] Ji, S., Pan, S., Cambria, E., Marttinen, P., Yu, P.S.: A survey on knowledge graphs: Representation, acquisition, and applications. In: IEEE TNNLS (2022)
- [21] Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. In: ICLR (2017)
- [22] Lin, Y., Han, X., Xie, R., Liu, Z., Sun, M.: Knowledge representation learning: A quantitative review (2018)
- [23] Lin, Y., Liu, Z., Luan, H., Sun, M., Rao, S., Liu, S.: Modeling relation paths for representation learning of knowledge bases. In: EMNLP (2015)
- [24] Mao, X., Wang, W., Wu, Y., Lan, M.: Boosting the speed of entity alignment 10 ×: Dual attention matching network with normalized hard sample mining. In: WWW (2021)
- [25] Mao, X., Wang, W., Xu, H., Lan, M., Wu, Y.: Mraea: An efficient and robust entity alignment approach for cross-lingual knowledge graph. In: WSDM (2020)
- [26] Mao, X., Wang, W., Xu, H., Wu, Y., Lan, M.: Relational reflection entity alignment. In: CIKM (2020)
- [27] Mikolov, T., Chen, K., Corrado, G., Dean, J.: Efficient estimation of word representations in vector space. In: CoRR (2013)
- [28] Monti, F., Shchur, O., Bojchevski, A., Litany, O., Günnemann, S., Bronstein, M.M.: Dual-primal graph convolutional networks. In: CoRR (2018)
- [29] Paulheim, H.: Knowledge graph refinement: A survey of approaches and evaluation methods. In: SEMWEB (2017)
- [30] Pujara, J., Miao, H., Getoor, L., Cohen, W.W.: Knowledge graph identification. In: SEMWEB (2013)
- [31] Radovanovic, M., Nanopoulos, A., Ivanovic, M.: Hubs in space: Popular nearest neighbors in high-dimensional data. In: JMLR (2010)
- [32] Rossi, A., Firmani, D., Matinata, A., Merialdo, P., Barbosa, D.: Knowledge graph embedding for link prediction: A comparative analysis. In: TKDD (2021)
- [33] Song, H.O., Xiang, Y., Jegelka, S., Savarese, S.: Deep metric learning via lifted structured feature embedding. In: CVPR (2015)
- [34] Srivastava, R.K., Greff, K., Schmidhuber, J.: Highway networks. In: ICML (2015)
- [35] Suchanek, F.M., Abiteboul, S., Senellart, P.: Paris: Probabilistic alignment of relations, instances, and schema. In: PVLDB (2011)

- [36] Sun, Y., Cheng, C., Zhang, Y., Zhang, C., Zheng, L., Wang, Z., Wei, Y.: Circle loss: A unified perspective of pair similarity optimization. In: CVPR (2020)
- [37] Sun, Z., Chen, M., Hu, W., Wang, C., Dai, J., Zhang, W.: Knowledge association with hyperbolic knowledge graph embeddings. In: EMNLP (2020)
- [38] Sun, Z., Hu, W., Li, C.: Cross-lingual entity alignment via joint attribute-preserving embedding. In: CoRR (2017)
- [39] Sun, Z., Hu, W., Zhang, Q., Qu, Y.: Bootstrapping entity alignment with knowledge graph embedding. In: IJCAI (2018)
- [40] Sun, Z., Wang, C., Hu, W., Chen, M., Dai, J., Zhang, W., Qu, Y.: Knowledge graph alignment network with gated multi-hop neighborhood aggregation. In: AAAI (2019)
- [41] Sun, Z., Zhang, Q., Hu, W., Wang, C., Chen, M., Akrami, F., Li, C.: A benchmarking study of embedding-based entity alignment for knowledge graphs. In: VLDB (2020)
- [42] Trisedya, B.D., Qi, J., Zhang, R.: Entity alignment between knowledge graphs using attribute embeddings. In: AAAI (2019)
- [43] Wang, Q., Mao, Z., Wang, B., Guo, L.: Knowledge graph embedding: A survey of approaches and applications. In: TKDE (2017)
- [44] Wang, Z., Lv, Q., Lan, X., Zhang, Y.: Cross-lingual knowledge graph alignment via graph convolutional networks. In: ACL (2018)
- [45] Wang, Z., Hamza, W., Florian, R.: Bilateral multi-perspective matching for natural language sentences. In: IJCAI (2017)
- [46] Wu, Y., Liu, X., Feng, Y., Wang, Z., Yan, R., Zhao, D.: Relation-aware entity alignment for heterogeneous knowledge graphs. In: IJCAI (2019)
- [47] Wu, Y., Liu, X., Feng, Y., Wang, Z., Zhao, D.: Jointly learning entity and relation representations for entity alignment. In: EMNLP (2019)
- [48] Xu, K., Wang, L., Yu, M., Feng, Y., Song, Y., Wang, Z., Yu, D.: Cross-lingual knowledge graph alignment via graph matching neural network. In: ACL (2019)
- [49] Zhang, R., Trisedya, B.D., Li, M., Jiang, Y., Qi, J.: A comprehensive survey on knowledge graph entity alignment via representation learning (2021)
- [50] Zhao, X., Zeng, W., Tang, J., Wang, W., Suchanek, F.: An experimental study of state-of-the-art entity alignment approaches. In: IEEE TKDE (2020)
- [51] Zhu, H., Xie, R., Liu, Z., Sun, M.: Iterative entity alignment via joint knowledge embeddings. In: IJCAI (2017)
- [52] Zhuang, Y., Li, G., Zhong, Z., Feng, J.: Hike: A hybrid human-machine method for entity alignment in large-scale knowledge bases. In: CIKM (2017)