



**HAL**  
open science

## Knowledge graph embedding methods for entity alignment: experimental review

Nikolaos Fanourakis, Vasilis Efthymiou, Dimitris Kotzinos, Vassilis Christophides

### ► To cite this version:

Nikolaos Fanourakis, Vasilis Efthymiou, Dimitris Kotzinos, Vassilis Christophides. Knowledge graph embedding methods for entity alignment: experimental review. *Data Mining and Knowledge Discovery*, 2023, 37 (5), pp.2070-2137. 10.1007/S10618-023-00941-9 . hal-04319674

**HAL Id: hal-04319674**

**<https://hal.science/hal-04319674v1>**

Submitted on 2 Feb 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License



# Knowledge graph embedding methods for entity alignment: experimental review

Nikolaos Fanourakis<sup>1,2</sup> · Vasilis Efthymiou<sup>1</sup> · Dimitris Kotzinos<sup>3</sup> · Vassilis Christophides<sup>3</sup>

Received: 28 January 2022 / Accepted: 12 May 2023 / Published online: 29 June 2023  
© The Author(s) 2023

## Abstract

In recent years, we have witnessed the proliferation of knowledge graphs (KG) in various domains, aiming to support applications like question answering, recommendations, etc. A frequent task when integrating knowledge from different KGs is to find which subgraphs refer to the same real-world entity, a task largely known as the Entity Alignment. Recently, embedding methods have been used for entity alignment tasks, that learn a vector-space representation of entities which preserves their similarity in the original KGs. A wide variety of supervised, unsupervised, and semi-supervised methods have been proposed that exploit both factual (attribute based) and structural information (relation based) of entities in the KGs. Still, a quantitative assessment of their strengths and weaknesses in real-world KGs according to different performance metrics and KG characteristics is missing from the literature. In this work, we conduct the first meta-level analysis of popular embedding methods for entity alignment, based on a statistically sound methodology. Our analysis reveals statistically significant correlations of different embedding methods with various meta-features extracted by KGs and rank them in a statistically significant way according to their effectiveness across all real-world KGs of our test-bed. Finally, we study interesting trade-offs in terms of methods' effectiveness and efficiency.

**Keywords** Knowledge graph embeddings · Entity alignment · Knowledge graphs

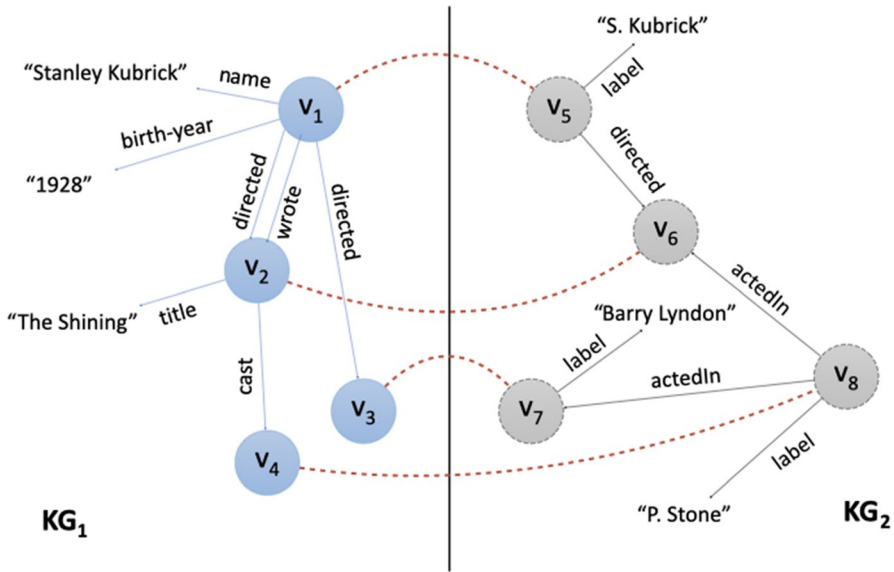
## 1 Introduction

In recent years, we have witnessed the proliferation of *knowledge graphs* (KGs) in various domains, aiming to support applications like entity search (Dong et al. 2014), question answering (Ahmetaj et al. 2021), and recommendations (Tarus et al.

---

Responsible editor: Aditya Prakash.

Extended author information available on the last page of the article



**Fig. 1** An example of two KGs whose entities can be aligned. The red dashed lines connect the aligned entities  $(v_1, v_5)$ ,  $(v_2, v_6)$ ,  $(v_3, v_7)$ , and  $(v_4, v_8)$

2018). Typically, KGs store machine-readable *descriptions* of real-world *entities* (e.g., people, movies, books) that capture both *relational* and *factual* information. In this work, we refer to an *entity description* as an identifiable set of property-value pairs that abstracts several data formats, such as relational, RDF, or property graphs.

As different KGs may independently describe the same real-world entity, a crucial task when integrating knowledge from several KGs is to align their entity descriptions. *Entity alignment* (EA), also known as *entity resolution* (Christophides et al. 2015, 2021), aims to identify pairs of descriptions from different KGs that refer to the same real-world entity, which we call matching pairs or simply, *matches*.

Figure 1 shows an example of two KGs, each containing four entity descriptions, represented as nodes, with their properties, represented as edges. The entities described in those KGs can be aligned. For instance, node  $v_1$  in  $KG_1$  is a description of the director Stanley Kubrick, providing his name and birth year as attributes, and the facts that he directed and wrote the movie The Shining, as well as that he also directed an entity represented by  $v_3$  as relations. Node  $v_1$  should be aligned with node  $v_5$  in  $KG_2$ , even if the “name” attribute is now called “label”, and its value, “S. Kubrick”, is slightly different than the name value, “Stanley Kubrick”, used in  $v_1$ . The birth year attribute is also missing in  $v_5$ , compared to  $v_1$ , and its only relation to  $v_6$  is “directed”, missing the relation “wrote” that was also provided in  $KG_1$  between  $v_1$  and  $v_2$ . Similarly, the other entity alignments in the example of Fig. 1 should be  $v_2$  with  $v_6$  (describing the movie “The Shining”),  $v_3$  with  $v_7$  (describing the movie “Barry Lyndon”), and  $v_4$  with  $v_8$  (describing the actor Philip Stone).

One way to implement entity alignment as a machine learning (ML) task, is to learn a vector-space representation of symbolic KGs, known as *embeddings*.

Numerical representations (embeddings) of KGs are preferred over symbolic ones in various ML tasks (link/node/subgraph prediction, matching, etc.), as they potentially mitigate the symbolic, linguistic and schematic heterogeneity of independently created KGs and thus aim to simplify knowledge reasoning. The idea is to embed the nodes (entities) and edges (relations or attributes) of a KG into a low-dimensional vector space. Particularly, we would like similar entities in the original KG to be close to each other in the embedding space and dissimilar entities to be far from each other. In this respect, both *positive* (i.e., actual KG edges) and *negative* (i.e., synthetic edges, non-existing in the actual KG) samples of the KGs are used. This way, by measuring their embeddings distance in the vector space, we can decide whether two entities are matching or not. Any available ground truth regarding the alignment of entities, called *seed alignment*, can be used for training and/or evaluating an embedding method.

Learning low-dimensional representations of KGs in a way such that the semantic relatedness of entities is captured by the geometrical structures of an embedding space is a challenging task that gave birth to numerous methods. *Embedding-based entity alignment* methods essentially exploit the *relational* (i.e., entity structural neighborhood) and the *factual* part (i.e., entity names/identities, attributes that represent literals) of descriptions. We refer to the former as *relation-based* methods, e.g., MTransE (Chen et al. 2017), MTransE+RotatE (Sun et al. 2020), RDGCN (Wu et al. 2019), RREA (Mao et al. 2020b), and to the latter as *attribute-based* methods, e.g., MultiKE (Zhang et al. 2019), AttrE (Trisedya et al. 2019), KDCoE (Chen et al. 2018), BERT\_INT (Tang et al. 2020). Although this research direction is rapidly growing, there are still several open questions regarding the underlying assumptions of methods, as well as, the efficiency and effectiveness of entity alignment in realistic settings. In particular, in this work we address the following missing insights from the literature:

- Q1. **Characteristics of methods.** What are the critical factors that affect the effectiveness of relation-based (e.g., negative sampling, range of neighborhood) and attribute-based methods (e.g., usage of literals) and how sensitive are the methods to hyperparameters tuning?
- Q2. **Families of methods.** What is the improvement in the effectiveness of embedding-based entity alignment methods if we consider not only the structural relations of entities, but also their attribute values?
- Q3. **Effectiveness vs Efficiency Tradeoff.** Is the runtime overhead of each method worth paying, with respect to the achieved effectiveness?
- Q4. **Characteristics of datasets.** To which characteristics of the datasets (e.g., sparsity, number of entity pairs in seed alignment, heterogeneity in terms of literals, predicate names and entity names) are supervised, semi-supervised and unsupervised methods sensitive?

Although several recent works (Zeng et al. 2021; Choudhary et al. 2021; Wang et al. 2017; Sun et al. 2020; Zhang et al. 2020; Zhao et al. 2022; Jiang et al. 2021; Wang et al. 2021; Leone et al. 2022; Zhang et al. 2022; Chaurasiya et al. 2022)

survey embedding-based entity alignment methods, only few of them (Sun et al. 2020; Zhang et al. 2020; Zhao et al. 2022; Leone et al. 2022; Zhang et al. 2022; Chaurasiya et al. 2022) conduct an experimental evaluation to obtain useful insights. The conclusions drawn from Zhang et al. (2020) are limited, as it leaves out some representative methods in embedding-based alignment such as MTransE (no negative sampling), KDCoE (semi-supervised exploiting long textual descriptions<sup>1</sup>), RREA (semi-supervised exploiting structural information), AttrE (unsupervised), BERT\_INT (supervised exploiting both structural and factual information), while neither RREA nor BERT\_INT were part of OpenEA (Sun et al. 2020), as both were published later. In addition, Leone et al. (2022) does not include MTransE, MTransE+RotatE, AttrE, KDCoE, while RREA is not included in Leone et al. (2022), Zhang et al. (2022), Chaurasiya et al. (2022) either. Moreover, benchmarking efforts such as Zhang et al. (2020), OpenEA (Sun et al. 2020), EAE (Zhao et al. 2022) and Leone et al. (2022), Zhang et al. (2022), Chaurasiya et al. (2022), do not shed light on questions Q1, Q2 and Q3, addressed in our work. Furthermore, only a subset of the dataset characteristics we study in our work, such as the density of the KGs, and the similarity of entity names, have been considered by previous works to answer question Q4.

To compare the effectiveness and efficiency of the methods in realistic settings, we have extended the testbed of datasets with pairs of KGs usually considered in related empirical studies. Specifically, OpenEA, EAE and Leone et al. (2022), Zhang et al. (2022), Chaurasiya et al. (2022) employ only datasets with a low number of entities featuring descriptions and literal values. In our testbed, we have included five additional datasets whose unique characteristics allow us to draw new insights regarding the evaluated EA methods, which were not previously reported in Sun et al. (2020), Zhang et al. (2020), Zhao et al. (2022), Leone et al. (2022), Zhang et al. (2022), Chaurasiya et al. (2022). More precisely, supervised methods like RDGCN exploiting KG relations, are outperformed by unsupervised (AttrE) and semi-supervised (KDCoE) methods that exploit the similarity of literals in datasets of decreasing density, but with rich factual information (i.e., attributes).

Rather than simply reporting the raw experimental results, we conduct a meta-level analysis, aiming to find statistically significant correlations between the methods and the dataset characteristics (meta-features). Furthermore, we consider the non-parametric Friedman test (Demsar 2006) and the post-hoc Nemenyi test (Nemenyi 1963), in order to perform a pairwise comparison of the methods and rank them based on their effectiveness across all datasets of our testbed. Finally, we are interested in the methods' training time curve and potential effectiveness vs efficiency trade-offs instead of simply reporting overall runtimes (as in OpenEA and EAE).

In a nutshell, the contributions of this work are the following:

---

<sup>1</sup> Although the term "literals" is inclusive of entity names and textual descriptions, in the EA literature we differentiate between them as follows: *entity names* are the suffixes of entity identifiers (which are sometimes meaningful), *textual descriptions* are the values of some pre-determined attributes that typically span several sentences, as opposed to the simple *literals*, used to describe all other factual data, typically of much shorter length (1-2 words).

- In Sect. 2, we present a qualitative comparison of state-of-the-art embedding-based entity alignment methods that span from supervised, i.e., MTransE (Chen et al. 2017), MTransE+RotatE (Sun et al. 2020), MultiKE (Zhang et al. 2019), RDGCN (Wu et al. 2019), RREA(basic) (Mao et al. 2020b), BERT\_INT (Tang et al. 2020), to unsupervised, i.e., AttrE (Trisedya et al. 2019), and semi-supervised, i.e., KDCoE (Chen et al. 2018), RREA(semi) (Mao et al. 2020b), paradigms. They are representative methods of different embedding families covering both relation- and attribute-based, but also considering one-hop and multi-hop neighborhoods in KGs, as well as different negative sampling strategies.
- In Sect. 3, we describe our framework for a fair empirical comparison of the different methods. We detail the extended testbed of datasets that exhibit diverse characteristics (w.r.t. KG density, entity naming, textual descriptions, etc.) usually encountered in reality along with the corresponding pre-processing pipelines. We additionally introduce the evaluation protocol and metrics capturing different aspects of the methods' effectiveness.
- In Sect. 4, we report and analyze the results of a series of experiments, including a comparison to a state-of-the-art non-embedding-based method, PARIS (Suchanek et al. 2011), conducted to answer the four open questions introduced previously, using a reliable, statistically sound methodology. First, we discover a statistical significant ranking of the methods according to their effectiveness across all real-world KGs of our testbed. Then, we study interesting trade-offs in terms of their effectiveness and efficiency. Last but not least, we extract statistically significant correlations between the methods' performance with various characteristics of our datasets (i.e., meta features).

Finally, the main conclusions drawn from our experiments, as well as the plans for future work, are discussed in Sect. 5.

## 2 Entity alignment with KG embeddings

In this section, we first formally define the entity alignment problem on Knowledge Graphs (KGs), along with some related constraints. Then, we provide a qualitative comparison of entity alignment methods based on KG embedding.

### 2.1 The entity alignment problem

Following the typical notation used in the literature (Zhang et al. 2019; Wang et al. 2020), we assume that entities (with the corresponding *entity names*),<sup>2</sup> are described in KGs by a collection of edges  $\langle h, r, t \rangle$ , whose head  $h$  is always an entity, and tail  $t$  may be either another entity, in which case we call this edge a *relation edge* and  $r$  a *relation*, or a literal (e.g., number, date, string), in which case we call this edge an *attribute edge* and  $r$  an *attribute* with its corresponding

---

<sup>2</sup> The string suffix after the last slash of a URI of an entity or an attribute (Zhang et al. 2019).

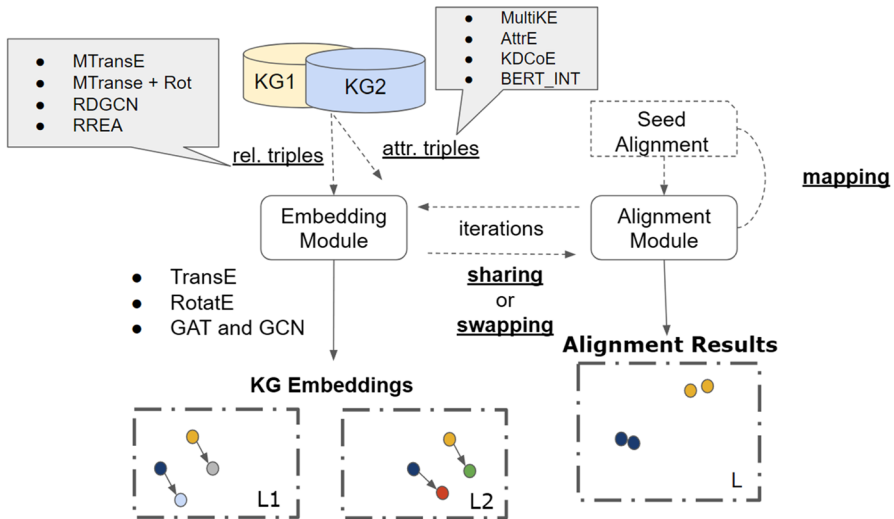
attribute name<sup>2</sup>. We represent a knowledge graph as  $KG = (E, R, A, L, X, Y)$ , where  $E$  is a set of entities,  $R$  is a set of relations,  $A$  is a set of attributes,  $L$  is a set of literals,  $X \subseteq (E \times R \times E)$  and  $Y \subseteq (E \times A \times L)$  are the sets of relation and attribute edges of the KG, respectively. Given a source  $KG_1 = (E_1, R_1, A_1, L_1, X_1, Y_1)$  and a target  $KG_2 = (E_2, R_2, A_2, L_2, X_2, Y_2)$ , the task of entity alignment is to find pairs of matching entities  $M = \{(e_i, e_j) \in E_1 \times E_2 \mid e_i \equiv e_j\}$ , where “ $\equiv$ ” denotes the equivalence relationship (Zhang et al. 2019; Wang et al. 2020). A subset  $\delta \subseteq M$  of the matching pairs may be used as a seed alignment for training. For instance, in Fig. 1, the entities of the two KGs are  $E_1 = \{v_1, v_2, v_3, v_4\}$  and  $E_2 = \{v_5, v_6, v_7, v_8\}$ . The relations are  $R_1 = \{cast, directed, wrote\}$  and  $R_2 = \{directed, actedIn\}$ , while the attributes are  $A_1 = \{name, birth-year, title\}$ ,  $A_2 = \{label\}$  and the literals are  $L_1 = \{\}\{StanleyKubrick'', \}\{1928'', \}\{TheShining''\}$ .  $L_2 = \{\}\{S.Kubrick'', \}\{BarryLyndon'', \}\{P.Stone''\}$ . The relation edges are  $X_1 = \{(v_1, directed, v_2), (v_2, directed, v_1), (v_1, directed, v_3), (v_2, directed, v_4)\}$ ,  $X_2 = \{(v_5, directed, v_6), (v_8, actedIn, v_6), (v_8, actedIn, v_7)\}$ , while the attribute edges are  $Y_1 = \{(v_1, name, \}\{StanleyKubrick''\}, (v_1, birth-year, \}\{1928''\}, (v_2, title, \}\{TheShining''\})$  and  $Y_2 = \{(v_5, label, \}\{S.Kubrick''\}, (v_8, label, \}\{P.Stone''\}, (v_7, label, \}\{BarryLyndon''\})$ . The task of entity alignment is to find the matches (denoted by dashed edges in Fig. 1)  $M = \{(v_1, v_5), (v_2, v_6), (v_3, v_7), (v_4, v_8)\}$ .

In practice, all the evaluated entity alignment methods rely on a number of assumptions/constraints, as listed below:

- Every entity is assumed to be the head of at least one relation edge (so we do not consider entities that are not part of a connected component of the KG):  
 $\forall e \in E, \exists r \in R, t \in E : (e, r, t) \in X$ .
- 1-to-1 constraint: Every entity in  $E_1$  should be matched to exactly one entity in  $E_2$ :  $\forall e_i \in E_1 (\exists e_j \in E_2 : (e_i, e_j) \in M) \wedge (\nexists e'_j \in E_2 : (e_i, e'_j) \in M)$  and vice versa  $\forall e_j \in E_2 (\exists e_i \in E_1 : (e_i, e_j) \in M) \wedge (\nexists e'_i \in E_1 : (e'_i, e_j) \in M)$ . This also implies that  $|M| = |E_1| = |E_2|$ .

## 2.2 Knowledge graph embeddings for entity alignment

KG embedding methods aim to learn a low-dimensional vector-space representation of symbolic KGs, known as *embeddings*. The idea is to embed the nodes (entities) and edges (relations or attributes) of a KG in an embedding space in a way that preserves their similarity in the original KG. Embedding methods have been proven to be effective in many machine learning tasks, such as node classification (Kipf and Welling 2017) that aims to assign entity types to KG nodes, or link prediction (Bordes et al. 2013; Sun et al. 2019) that aims to find missing relations between entities in a single KG. Lately, several embedding-based methods have been also proposed for entity alignment, exploiting either relation edges (*relation-based methods*), such as MTransE (Chen et al. 2017), MTransE + RotatE (Sun et al. 2020), RDGCN (Wu et al.



**Fig. 2** General framework of KG embeddings for entity alignment

2019), RREA(basic) (Mao et al. 2020b), RREA(semi) (Mao et al. 2020b) or attribute edges (*attribute-based methods*), such as MultiKE (Zhang et al. 2019), AttrE (Trisedya et al. 2019), KDCoE (Zhang et al. 2017), and BERT\_INT (Tang et al. 2020).

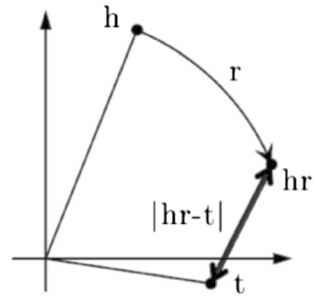
Figure 2 depicts the building blocks of embedding-based entity alignment methods: (i) The *embedding module*  $S_K$  that encodes the entities of each KG in an embedding space ( $L_1$  for  $KG_1$  and  $L_2$  for  $KG_2$ ) according to the *relational* (i.e., entity structural neighborhood) and/or the *factual* part (i.e., entity names/identities, literals/text) of descriptions. (ii) The *alignment module*  $S_A$  that aligns the produced entity embeddings using the *seed alignment* (supervised) or *attribute-values* similarity (unsupervised), or *both* (semi-supervised). It produces a common embedding space for the entities of two KGs, in order to generate the alignment result according to a distance metric (e.g., Euclidean), using three different techniques known as *sharing*, *swapping* and *mapping*. Sharing and Swapping, update directly the entity embeddings produced by the embedding module according to the available similarity evidence of entities, while Mapping essentially learns a linear transformation between the two embedding spaces of aligned KGs. In the rest of this section, we will detail popular methods that implement those modules.

### 2.2.1 Embedding module

KG embedding methods proposed for the task of link prediction are used to implement the embedding module of entity alignment methods. There are several families of KG embedding methods for link prediction have been proposed in the literature, e.g., Yang et al. (2015), Trouillon et al. (2016); Nickel and Kiela (2017). In this paper, we are interested in contrasting translational methods such



Fig. 3 Example of rotatE



as TransE (Bordes et al. 2013) and RotatE (Sun et al. 2019), with Graph Neural Networks such as Graph Convolutional Networks (GCNs) (Kipf and Welling 2017) and Graph Attention Networks (GATs) (Velickovic et al. 2018).

**2.2.1.1 Translational methods** Translational methods use distance-based scoring functions in order to optimize a margin-based loss function and learn the embeddings of entities in a KG. A distance-based scoring function is a function that measures the plausibility of a relation edge  $\langle h, r, t \rangle$  i.e., it measures the distance of the embedding of the head to the embedding of the tail entities, given the embedding of the relation. A margin-based loss function is a function that these methods aim to minimize, in order to minimize the distance of entity embeddings by a certain margin, computed by a distance-based scoring function. The key difference among all those translational methods is based on the degree that they are able to capture more complex graph structures such as cycles, by adopting the appropriate operator in the scoring function.

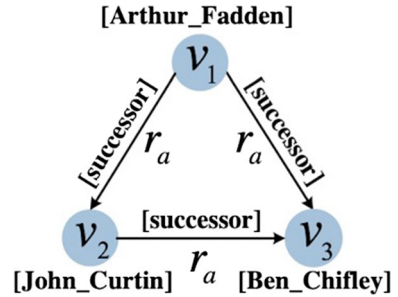
**TransE** (Bordes et al. 2013) is one of the most widely used translational KG embedding methods. In this method, both entities and relations are represented in the same vector space. The relation  $r$  is equivalent to the translation of vectors from head entity  $h$  to the tail entity  $t$ . If  $\langle h, r, t \rangle \in X$ , then the embedding  $\mathbf{t}$  of  $t$  should be close to the embedding  $\mathbf{h}$  of  $h$ , plus the vector  $\mathbf{r}$  of  $r$ , i.e.,  $\mathbf{h} + \mathbf{r} \approx \mathbf{t}$ . Formally, TransE minimizes the margin-based loss function:

$$J_{SE} = \sum_{x \in X} \sum_{x' \in X'} \max(0, \gamma + f(\mathbf{x}) - f(\mathbf{x}')), \tag{1}$$

where  $f(\mathbf{h}, \mathbf{r}, \mathbf{t}) = |\mathbf{h} + \mathbf{r} - \mathbf{t}|$  is the scoring function,  $X$  is the set of *positive relation edges* (relation edges that exists in the KG),  $X'$  is the set of *negative relation edges* (relation edges that do not exist in the KG), and  $\gamma$  is the margin hyperparameter. Each negative edge  $x' \in X'$  is created by replacing the head or the tail of a positive edge in  $X$  with a random entity, ensuring that  $x' \notin X$ .

**RotatE** (Sun et al. 2019) is a translation-based embedding model that, unlike TransE, infers various relation patterns, such as symmetries. Specifically, RotatE maps the entities and relations to the complex vector space and defines each relation as a rotation from the head entity to the tail entity (Fig. 3). Given a relation edge  $\langle h, r, t \rangle$ , we expect that  $\mathbf{t} = \mathbf{h} \circ \mathbf{r}$ , where  $\circ$  denotes the Hadamard (Million 2007) (element-wise) product.

**Fig. 4** Triangular structure (Wu et al. 2019)



This model aims to minimize the margin-based loss function

$$L = -\log \sigma(\gamma - d_r(\mathbf{h}, \mathbf{r}, \mathbf{t})) - \sum_{i=1}^n \frac{1}{k} \log \sigma(d_r(\mathbf{h}'_i, \mathbf{r}, \mathbf{t}'_i) - \gamma), \quad (2)$$

by maximizing the scores of positive relation edges and minimizing the scores of negative relation edges, where  $d_r(\langle h, r, t \rangle) = |\mathbf{h} \circ \mathbf{r} - \mathbf{t}|$  is a scoring function,  $\gamma$  is a fixed margin,  $\sigma$  is the sigmoid function, and  $(h'_i, r, t'_i)$  is the  $i$ -th negative edge. RotatE, like TransE, creates the negative relation edges by replacing the head and the tail of positive relation edges randomly.

**2.2.1.2 Graph Neural Network Methods** Translational methods cannot deal with various complex graph structures. For example, TransE (Bordes et al. 2013) cannot deal with triangular structures like the one in Fig. 4, because it requires the three equations  $v_1 + r_a \approx v_2$ ,  $v_2 + r_a \approx v_3$  and  $v_1 + r_a \approx v_3$  to hold at the same time. This is impossible, because for satisfying the former two equations we would have  $v_1 + 2r_a \approx v_3$  which is contradictory to the equation  $v_1 + r_a \approx v_3$ .

In order to cope with that, Graph Neural Network (GNN) methods have been proposed. GNNs learn entity embeddings, by recursively aggregating the representations of neighboring nodes. They essentially rely on *message passing*, according to which, each graph node recursively receives and aggregates features (node representations) from its neighbors in order to represent the local graph structure.

There is a range of GNN variants, that implement different aggregation strategies. In this section, we focus on standard graph convolutional networks (GCNs) (Kipf and Welling 2017) and graph attention networks (GATs) (Velickovic et al. 2018), since they are the core of both RDGCN (Wu et al. 2019) and RREA (Mao et al. 2020b); two of the proposed methods that we evaluate in this study and describe in Sect. 2.3.

**GCN** (Kipf and Welling 2017) takes as input the randomly initialized entity embeddings of the KG, which is treated as an undirected graph. Then, it learns a set of layer-specific weights, known as *filters* or *kernels*, that are multiplied with the input embeddings. In essence, it acts as a sliding window across the KG that learns entity features while preserving useful structural information from the neighborhoods. GCN uses the following function

$$H_i^{(l+1)} = \sigma \left( \sum_{j \in \mathcal{N}_i} \frac{1}{c_{ij}} W^{(l)} H_j^{(l)} \right) \tag{3}$$

to aggregate the entity embeddings of  $l$  layers, where  $\sigma$  is an activation function,  $\mathcal{N}_i$  is the set of the one-hop neighbors of the central entity  $i$  (including itself by adding a self-loop),  $c_{ij}$  is a normalization constant that defines isotropic average computation (each neighbor contributes equally to update the embedding of the central entity),  $W^{(l)}$  is a trainable layer-specific weighted matrix for feature transformation and  $H^{(l)}$  are the entity embeddings for layer  $l$ . More precisely, in order to learn the final embedding of a central entity, GCN sums its embedding with the neighbors embeddings.

*GAT* (Velickovic et al. 2018) expands the aggregation function of GCN, by an *attention mechanism* that assigns different weights to each neighbor of a central entity. GAT uses the following aggregation function

$$H_i^{(l+1)} = \sigma \left( \sum_{j \in \mathcal{N}_i} \alpha_{ij}^{(l)} z_j^{(l)} \right) \tag{4}$$

that aggregates the entity embeddings of  $l$  layers, where  $\sigma$  is an activation function,  $\mathcal{N}_i$  is the set of the one-hop neighbors of the central entity  $i$ ,  $H^{(l)}$  are the entity embeddings for layer  $l$ ,  $z_i^{(l)}$  is a transformation operation and  $\alpha_{ij}^{(l)}$  is the normalized coefficient score. Particularly,  $z_i^{(l)}$  and  $\alpha_{ij}^{(l)}$  are calculated as:

$$z_i^{(l)} = W^{(l)} H_i^{(l)} \tag{5}$$

and

$$\alpha_{ij}^{(l)} = \frac{\exp \left( e_{ij}^{(l)} \right)}{\sum_{k \in \mathcal{N}(i)} \exp \left( e_{ik}^{(l)} \right)}, \tag{6}$$

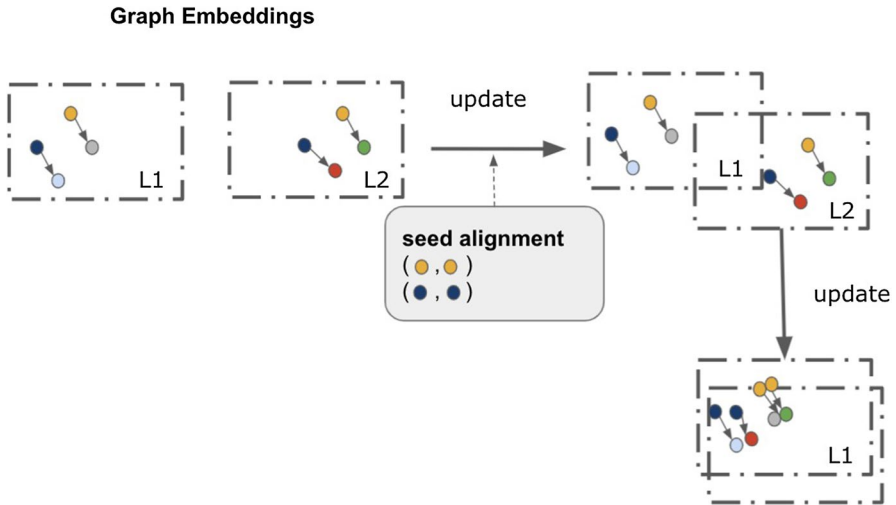
where

$$e_{ij}^{(l)} = \text{LeakyReLU} \left( \vec{\mathbf{a}}^{(l)T} \left( z_i^{(l)} \parallel z_j^{(l)} \right) \right), \tag{7}$$

$\vec{\mathbf{a}}^{(l)T}$  is a learnable weight vector, LeakyReLU a variant of the activation function ReLU (Parisi et al. 2022) and  $\parallel$  is the concatenation operation.

### 2.2.2 Alignment module

For the alignment module  $S_A$ , there exist three techniques: sharing, swapping and mapping (Sun et al. 2020). We describe them below, while we extensively compare them in Sect. 2.5.2.



**Fig. 5** Sharing alignment technique

**2.2.2.1 Sharing** Sharing aims to iteratively update the already produced entity embeddings, in order to minimize the embedding distance of each entity  $e$  and its aligned entity  $e'$  from the seed alignment  $\delta$ .

$$S_A = \sum_{(e,e') \in \delta(KG_i, KG_j)} \|e - e'\|. \quad (8)$$

In Fig. 5, we demonstrate the entity embeddings of  $KG_1$  and  $KG_2$  in  $L1$  and  $L2$  embedding spaces, respectively, while we also show the updates of the embeddings of the entities of seed alignment, in order to minimize their embedding distance. For simplification, we use a part of seed alignment, thus only the blue entities and the orange entities are considered as aligned. Therefore, by this technique, assuming the spatial similarity of aligned entities in two different KGs, we aim to adjust the axis of the two embedding spaces, so that entity vectors of the same entity in two KGs to overlap. It is worth mentioning that we started from two KGs encoded in two different embedding spaces (embeddings from the embedding module) and we ended up with two KGs encoded in an unified embedding space.

**2.2.2.2 Swapping** Swapping is a variation of sharing that produces extra positive edges, preserving the same objective as sharing. For instance, given two aligned entity pairs  $(h, h') \in \delta(KG_1, KG_2)$  and  $(t, t') \in \delta(KG_1, KG_2)$  and a relation edge  $(h, r, t)$  of  $KG_1$ , swapping produces two new positive edges  $(h', r, t)$  and  $(h, r, t')$  and feeds them in KG embedding models (embedding module) as positive relation edges, in order to increase the training data, benefiting the quality of the embeddings as we describe in Sects. 2.5.1.2 and 2.5.6. Swapping does not introduce a new loss function.

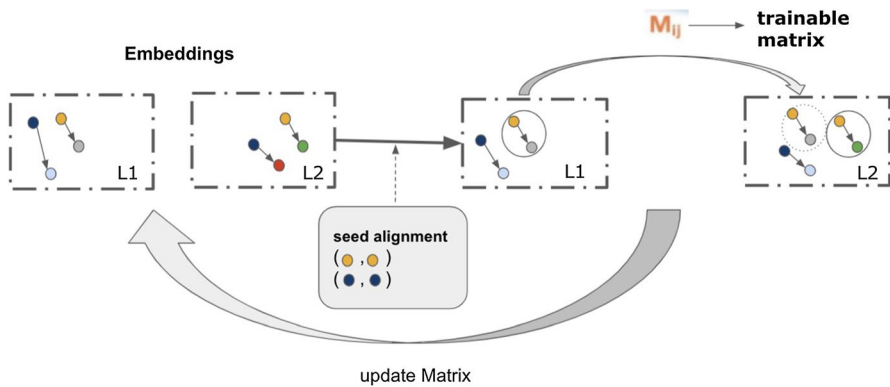


Fig. 6 Mapping alignment technique

**2.2.2.3 Mapping** Mapping aims to learn a matrix  $M$  as a linear transformation on entity vectors from  $L_i$  to  $L_j$ , in order to minimize the embedding distance of each linearly transformed entity  $e$  and its aligned entity  $e'$  from the seed alignment  $\delta$ :

$$S_A = \sum_{(e,e') \in \delta(KG_i, KG_j)} \|M_{ij}e - e'\|. \tag{9}$$

In Fig. 6, we demonstrate the entity embeddings of  $KG_1$  and  $KG_2$  in  $L_1$  and  $L_2$  embedding spaces respectively, while we also show the process in which we learn the matrix  $M_{i,j}$  that linearly transforms entities from  $L_1$  to  $L_2$ . During this process, the linearly transformed entities of  $KG_1$  should be close to their aligned entity of  $KG_2$  according to the seed alignment. For simplification, we use a part of the seed alignment, thus only the blue entities and the orange entities are considered as aligned. Mapping, in contrary to sharing and swapping, aims to learn the mappings between the two embedding spaces (deducing the linear transformation from  $L_1$  to  $L_2$ ), without assuming the similarity of spatial emergence. More precisely, it does not force the entity vectors of aligned entities to overlap, instead, it treats the learned mappings as topological transformations (one-to-one correspondence) from  $L_1$  to  $L_2$ , preserving the two KGs encoded in two different embedding spaces.

### 2.3 Knowledge graph embeddings using relations

In this section, we discuss relation-based KG embedding methods, all of which are supervised. These methods use only the structural information (relation edges) for learning the entity embeddings.

**MTransE** (Chen et al. 2017) is a translation-based model for multilingual KG embeddings, but it is also applicable to general-purpose KGs, capturing their structure. The objective is to minimize the loss function

$$J = S_K + \alpha S_A, \quad (10)$$

where  $S_K$  is the loss function of the embedding module,  $S_A$  is the loss function of the alignment module, and  $\alpha$  is a factor that weights  $S_K$  and  $S_A$ . As the loss function  $S_K$  of the embedding module, MTransE utilizes a simplified version of TransE (Eq. 1), in which no negative relation edges are considered, while as the loss function  $S_A$  of the alignment module, it uses mapping (Eq. 9).

**MTransE+RotatE** (Sun et al. 2020) is a variation of MTransE that uses RotatE (Eq. 2) as  $S_K$  in Eq. 10, and sharing (Eq. 8) as  $S_A$ , instead of TransE and mapping, respectively.

**RDGCN** (Wu et al. 2019) leverages GCNs (described in 2.2.1.2) to incorporate structural information in the entity embeddings. Particularly, given  $KG_1$  and  $KG_2$ , RDGCN constructs a primal (entity) graph  $G^e$  by merging  $KG_1$  and  $KG_2$ , and its dual (relation) graph  $G^r$ , by creating a node in  $G^r$  for every relation type of  $G^e$ , and connecting two nodes in  $G^r$  if the corresponding relations in  $G^e$  share the same head or tail entities.

Then, it uses a graph attention mechanism (a dual attention layer that assigns different importance to each neighbor's contribution) to make interactions between  $G^e$  and  $G^r$ , in order the resulting entity representations in  $G^e$  to capture the relation information, and then, to be fed to a GCN, capturing the structure of the neighborhood (Eq. 3). The resulting entity embeddings are refined using the mapping alignment technique (Sect. 2.2.2.3). The loss function that RDGCN aims to minimize is

$$L = \sum_{(e_i, e_j) \in \delta, (e'_i, e'_j) \notin \delta} \max\left(0, d(\mathbf{e}_i, \mathbf{e}_j) - d(\mathbf{e}'_i, \mathbf{e}'_j) + \gamma\right), \quad (11)$$

where  $(e_i, e_j)$  are entity pairs from the seed alignment  $\delta$ ,  $(e'_i, e'_j)$  are negative samples generated by replacing  $e_i$  or  $e_j$  with a random entity, and  $d$  is the embedding distance function used in mapping (Sect. 2.2.2.3).

**RREA** (Mao et al. 2020b) integrates GCNs and GATs (described in Sect. 2.2.1) with a *Relational Reflection Transformation*, in order to obtain relation-specific embeddings for KG entities. This transformation utilizes a matrix that, in contrary to standard GCN and GAT, is constrained to be orthogonal, in order to reflect entity embeddings across different relational hyperplanes. The orthogonal property of the aforementioned matrix keeps the norms and the relative distances of entities in the relational space unchanged.

More precisely, RREA stacks multiple GNN layers, in order to capture and aggregate multi-hop neighborhood information for each entity embedding. The output embedding of entity  $e_i$  from the  $l$ -th layer is obtained as follows:

$$\mathbf{H}_{e_i}^{l+1} = \text{ReLU} \left( \sum_{e_j \in \mathcal{N}_{e_i}^e} \sum_{r_k \in R_{ij}} \alpha_{ijk}^l \mathbf{M}_{r_k} \mathbf{h}_{e_j}^l \right), \quad (12)$$

where ReLU (Parisi et al. 2022) is an activation function,  $\mathcal{N}_{e_i}^e$  are the neighboring entities of  $e_i$ ,  $R_{ij}$  denotes the relations between  $e_i$  and  $e_j$ ,  $\mathbf{M}_{r_k}$  the relational reflection

matrix of  $r_k$ , and  $\alpha_{ijk}^l$  is a weight coefficient of  $M_{r_k}$  (similar to GAT). The final entity embedding comes from the concatenation of the embeddings of each layer. In addition, in order to include relational information around entities, RREA concatenates the summation of the relation embeddings with entity embeddings to get dual-aspect embeddings. The resulting entity embeddings are refined using the sharing alignment technique (Sect. 2.2.2.1). The loss function that RREA aims to minimize is the following:

$$L = \sum_{(e_i, e_j) \in P} \max \left( \text{dist}(e_i, e_j) - \text{dist}(e'_i, e'_j) + \lambda, 0 \right), \tag{13}$$

where  $e'_i$  and  $e'_j$  represent the negative pair of  $e_i$  and  $e_j$ , generated using *truncated uniform negative sampling* (Sun et al. 2018; Zhu et al. 2019; Cao et al. 2019) and *dist* is the embedding distance function used in sharing.

The methodology described above refers to the basic version of RREA, *RREA(basic)*. RREA also comes with a semi-supervised version, *RREA(semi)*, that proposes possibly aligned entity pairs in different iterations, in order to enrich the training set. According to Mao et al. (2020a), the entity pair  $(e_i, e_j)$  is proposed as aligned, if  $e_i$  and  $e_j$  are mutually nearest aligned.

## 2.4 Knowledge graph embeddings using attributes

In this section, we focus on attribute-based KG embedding methods. These methods utilize not only the structural information of the KGs (relation edges) to learn the entity embeddings, but also the attribute values (literals). In addition, in many methods the attribute embeddings help to enrich the seed alignment or even in refining the entity embeddings. We categorize the attribute-based methods depending on the usage of seed alignment as *supervised*, *semi-supervised* and *unsupervised*.

### 2.4.1 Supervised

*MultiKE* (Zhang et al. 2019) first constructs the embeddings of each literal  $l$

$$\phi(l) = \text{encode}([\text{LP}(o_1); \text{LP}(o_2); \dots; \text{LP}(o_n)]), \tag{14}$$

where  $\text{LP}(o_n)$  is the pre-trained word embedding of word  $o_n$ ,  $\text{encode}(\cdot)$  is the encoder that does the compression of the embeddings, and  $[\ ]$  is the concatenation operation. If  $o_n$  is an out-of-vocabulary word (i.e., there is no pre-trained embedding for this word), then MultiKE builds it by using pre-trained character embeddings. Then, it learns entity embeddings by exploring three different views: the *name view*  $\Theta^{(1)}$ , the *relation view*  $\Theta^{(2)}$  and the *attribute view*  $\Theta^{(3)}$ .

Given an entity  $e$ , its name view  $(\Theta^{(1)})$  is defined as

$$\mathbf{e}^{(1)} = \phi(\text{name}(e)), \tag{15}$$

where  $\text{name}(\cdot)$  is the name of the entity.

For the relation view  $\Theta^{(2)}$ , it adopts TransE to learn the entity embeddings of the two KGs, minimizing the following loss function

$$\mathcal{L}(\Theta^{(2)}) = \sum_{(h,r,t) \in X^+ \cup X^-} \log(1 + \exp(-\zeta_{(h,r,t)} f_{\text{rel}}(\mathbf{h}, \mathbf{r}, \mathbf{t}))), \tag{16}$$

where  $X^+ = X_1 \cup X_2$  are the relation edges that exist in the two KGs,  $X^-$  are relation edges that do not exist in the two KGs (negative relation edges, created as in TransE),  $f_{\text{rel}}$  is the scoring function of TransE (Eq. 1), and  $\zeta_{(h,r,t)} \in \{-1, 1\}$  denotes whether  $(h, r, t)$  is a positive or a negative edge.

For the attribute view  $\Theta^{(3)}$ , again it uses TransE to learn the embeddings exploiting the attributes and their values, aiming to minimize the loss function

$$\mathcal{L}(\Theta^{(3)}) = \sum_{(h,a,v) \in Y^+} \log(1 + \exp(-f_{\text{attr}}(\mathbf{h}^{(3)}, \mathbf{a}, \mathbf{l}))), \tag{17}$$

where  $Y^+ = Y_1 \cup Y_2$  are the attribute edges of the two KGs,  $f_{\text{attr}}(\mathbf{h}^{(3)}, \mathbf{a}, \mathbf{l}) = -|\mathbf{h}^{(3)} - \text{CNN}(\langle \mathbf{a}; \mathbf{l} \rangle)|$ , which is using a Convolution Neural Network (CNN) representation of an attribute  $a$  and its literal value  $l$ , as follows:

$$\text{CNN}(\langle \mathbf{a}; \mathbf{l} \rangle) = \sigma(\text{vec}(\sigma(\langle \mathbf{a}; \mathbf{l} \rangle * \Omega))\mathbf{W}), \tag{18}$$

where  $[\ ]$  denotes the concatenation operation,  $\mathbf{a}$  the embedding of an attribute  $a$ ,  $\mathbf{l}$  the embedding of a literal  $l$ ,  $\Omega$  the kernel of CNN,  $\sigma$  the activation function, and  $\mathbf{W}$  a trainable weighted matrix.

For refining the entity embeddings of the relation and the attribute views, MultiKE minimizes the following two loss functions, respectively:

$$\begin{aligned} \mathcal{L}_{\text{CE}}(\Theta^{(2)}) &= \sum_{(h,r,t) \in \mathcal{X}} \log(1 + \exp(-f_{\text{rel}}(\hat{\mathbf{h}}^{(2)}, \mathbf{r}, \mathbf{t}^{(2)}))) \\ &+ \sum_{(h,r,t) \in \mathcal{X}'} \log(1 + \exp(-f_{\text{rel}}(\mathbf{h}^{(2)}, \mathbf{r}, \hat{\mathbf{t}}^{(2)}))) \end{aligned} \tag{19}$$

and

$$\mathcal{L}_{\text{CE}}(\Theta^{(3)}) = \sum_{(h,a,l) \in \mathcal{Y}} \log(1 + \exp(-f_{\text{attr}}(\hat{\mathbf{h}}^{(3)}, \mathbf{a}, \mathbf{l}))), \tag{20}$$

where  $(h, \hat{h})$  and  $(t, \hat{t})$  are entity pairs in the seed alignment,  $\mathcal{X}'$  and  $\mathcal{X}''$  are the relation edges whose head and tail entities are in seed alignment, respectively, and  $\mathcal{Y}'$  are the attribute edges whose head entities are in seed alignment.

For the entity alignment (alignment module), MultiKE produces a set of aligned relations  $S_{\text{rel}}$  and a set of aligned attributes  $S_{\text{attr}}$  that are used to minimize the following loss function (cross-KG relation inference):

$$\mathcal{L}_{\text{CRA}}(\Theta^{(2)}) = \sum_{(h,r,t) \in \mathcal{X}''} \text{sim}(r, \hat{r}) \log(1 + \exp(-f_{\text{rel}}(\mathbf{h}^{(2)}, \hat{\mathbf{r}}, \mathbf{t}^{(2)}))), \tag{21}$$

where  $\mathcal{X}'''$  are the relation edges whose relations are in  $S_{\text{rel}}$  and



$$\text{sim}(r, \hat{r}) = \alpha_1 \text{cosine}(\phi(\text{name}(r)), \phi(\text{name}(\hat{r}))) + (1 - \alpha_1) \text{cosine}(\mathbf{r}, \hat{\mathbf{r}}) \quad (22)$$

is the similarity measure that is used to align or not two relations, based on their name similarity (from literal embeddings) and their semantic similarity (from relation embeddings).

Finally, MultiKE jointly learns the final entity embeddings from the different views in a unified embedding space, by minimizing the loss function

$$\mathcal{L}_{\text{ITC}}(\tilde{\mathbf{H}}, \mathbf{H}) = \sum_{i=1}^D \|\tilde{\mathbf{H}} - \mathbf{H}^{(i)}\|, \quad (23)$$

where  $\tilde{\mathbf{H}} = \bigcup_{i=1}^D \mathbf{H}^{(i)}$ ,  $\mathbf{H}$  is the view-specific entity embedding, and  $D$  is the number of views. For entity alignment, it uses swapping (Sect. 2.2.2.2).

**BERT\_INT** (Tang et al. 2020) utilizes a well-known language model, BERT (Devlin et al. 2019), in order to embed entities based on their factual information (e.g., descriptions, names) and an interaction model, in order to compute their interactions, instead of aggregating neighbors, which in many cases causes noisy matches.

More precisely, for each entity  $e$ , it applies a pre-trained basic BERT unit that accepts the factual information as input, aiming to minimize the following loss function

$$\mathcal{L} = \sum_{(e, e'^+) \in \mathcal{D}} \max \{0, g(e, e'^+) - g(e, e'^-) + m\} \quad (24)$$

to fine-tune BERT. Here,  $D$  is the seed alignment,  $e'^+$  is the correctly aligned entity known from seed alignment,  $e'^-$  is a randomly selected negative entity from the other KG - *truncated uniform negative sampling* (Sun et al. 2018),  $m$  the margin and  $g$  is the  $l1$  distance, used for measuring the similarity between the embeddings  $C(e)$  and  $C(e')$ .

Regarding the interaction model, it is divided into the name/description view, the neighbor-views and the attribute-view interactions. Firstly, as name/description interaction, it leverages the embeddings generated by BERT, calculating their cosine similarity. Then, neighbor-view interaction compares names/descriptions of each neighbor pair (considering also their neighboring relations and multi-hop neighbors), producing a similarity matrix. The similarity matrix is then processed by a dual aggregation function to extract the similarity vectors, i.e., the entity embeddings. Afterwards, rather than learning embeddings of entities by aggregating their attributes, it compares each attribute pair, learning similarly to neighbor-view the attribute similarity vectors. Finally, a unified dual aggregation function is applied to extract the features from the neighbor-view and attribute-view interactions and generate the final entity embeddings.

## 2.4.2 Semi-supervised

**KDCoE** (Chen et al. 2018) leverages a weakly aligned KG for semi-supervised entity alignment using long (typically from a couple of sentences) textual descriptions of entities. It co-trains iteratively two embedding models, one on the structure of the KG (KGEM) and another on the textual descriptions of the entities (DEM), respectively, given a small seed alignment. During each iteration, each embedding model proposes a new set of aligned entity pairs alternately, in order to enrich the seed alignment. The process runs until one of the models has no entity pair to propose.

KGEM is practically the same as MTransE (Eq. 10), using the TransE (Eq. 1) embedding module (including negative samples) and mapping (Eq. 9). At the end of this model, if the embedding distance between an entity  $e$  and its closest (based on the distance function of mapping) entity  $\hat{e}'$  is lower than a threshold, then the pair  $(e, \hat{e}')$  is proposed as aligned. DEM utilizes an encoder to process textual description sequences of vectors  $d_e$ , that are produced by pre-trained word embeddings, and learn the description embeddings. The learning objective of DEM is to maximize the log likelihood of an entity  $e$  and its counterpart  $e'$  in terms of description embeddings, by minimizing the following loss function:

$$S_D = \sum_{(e,e') \in \delta} -LL_1 - LL_2 = \sum_{(e,e') \in \delta} -\log(P(e | e')) - \log(P(e' | e)) \quad (25)$$

where  $\delta$  is the seed alignment and

$$LL_1 = \log \sigma(\mathbf{d}_e^\top \mathbf{d}_{e'}) + \sum_{k=1}^{|B_d|} \mathbb{E}_{e_k \sim U(e_k \in E_{L_i})} \left[ \log \sigma(-\mathbf{d}_e^\top \mathbf{d}_{e'}) \right] \quad (26)$$

$$LL_2 = \log \sigma(\mathbf{d}_e^\top \mathbf{d}_{e'}) + \sum_{k=1}^{|B_d|} \mathbb{E}_{e_k \sim U(e_k \in E_{L_j})} \left[ \log \sigma(-\mathbf{d}_e^\top \mathbf{d}_{e_k}) \right], \quad (27)$$

where  $d_e$  and  $d_{e'}$  are the embeddings of the textual descriptions of the two aligned entities  $e$  and  $e'$ , unrelated entities  $e_k$  are chosen randomly from a uniform distribution  $U$ , and  $|B_d|$  is the batched sampling size. Intuitively, the encoder aims to maximize the dot product of descriptions of aligned entities and decrease the dot product of descriptions of unrelated entities. At the end of this model, if the embedding distance between  $d_e$  and its closest entity  $d_{e'}$  is lower than a threshold (different than the one used for KGEM), then the pair  $(e, e')$  is proposed as aligned.

For generating entity embeddings from textual descriptions, KDCoE utilizes a self-attention Gated Recurrent Unit (GRU), in order to preserve the sequence of words in a textual description and remove information irrelevant for the prediction, while sharing information across the different descriptions. For more details on this, we refer the reader to Chen et al. (2018).

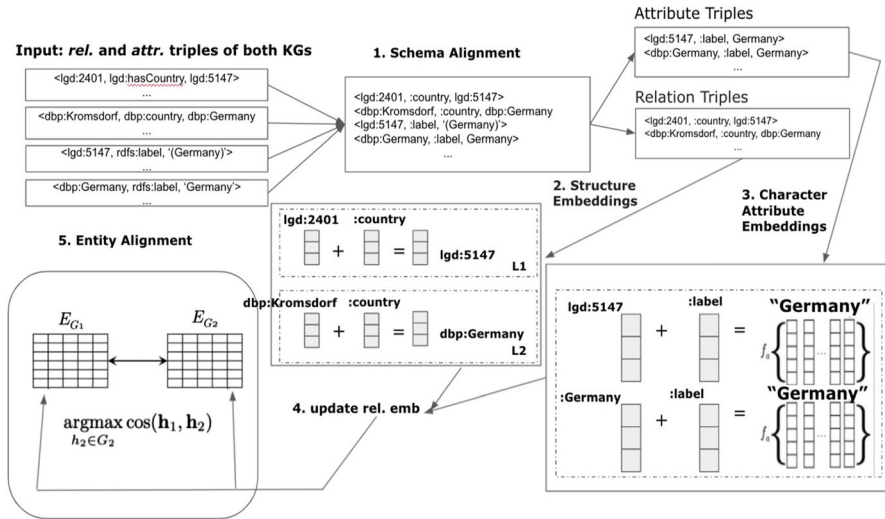


Fig. 7 Example of AttrE

### 2.4.3 Unsupervised

*AttrE* (Trisedya et al. 2019) is an unsupervised method that leverages both structural embeddings and attribute character embeddings for entity alignment. Instead of relying on a seed alignment to refine the structural embeddings, it uses the factual information to minimize the embedding distance between entities that have similar attribute character embeddings.

As shown in Fig. 7, AttrE consists of four modules: the schema alignment module (preprocessing step), the structure embedding module  $J_{SE}$ , the attribute character embedding module  $J_{CE}$ , and the alignment module  $J_{SIM}$

$$J = J_{SE} + J_{CE} + J_{SIM}. \tag{28}$$

The predicate alignment module merges the two KGs and renames predicates with a similar predicate name from  $KG_1$  and  $KG_2$ , with a unified naming schema. For example, *lgd:hasCountry* and *dbp:country* (with *lgd* and *dbp* being the prefixes of the Linked Geo Data and DBpedia namespaces, respectively) are converted to *:country*. To find the predicates with similar names, it computes the Levenshtein distance of the postfixes of the predicates' URIs. If this score is greater than a predefined threshold, then the two predicates are similar.

AttrE adopts TransE (Eq. 1) to learn the structure embeddings by minimizing the  $J_{SE}$  loss function, with  $\alpha = \frac{\text{count}(r)}{|X|}$ , where  $\text{count}(r)$  is the number of occurrences of relation  $r$ , and  $|X|$  is the total number of relation edges in the merged KG.

To learn the attribute character embedding, AttrE minimizes the following objective function:

$$J_{CE} = \sum_{y \in Y} \sum_{y' \in Y'} \max(0, [\gamma + w(f(y) - f(y'))]), \quad (29)$$

where

$$f(\langle h, a, l \rangle) = |\mathbf{h} + \mathbf{r} - f_a(l)| \quad (30)$$

and  $Y$  is the set of positive attribute edges,  $Y'$  is the set of negative attribute edges (generated by replacing the head entity with a random entity),  $w$  weights relation edges with aligned predicates, and  $f_a(l) = \sum_{n=1}^N \left( \frac{\sum_{i=1}^l \sum_{j=i}^n c_i}{l-i-1} \right)$  is an  $n$ -gram based compositional function that encodes literals  $l$  of attributes  $a$  in the embedding space. At the end of this process, entities that have similar attribute embeddings should also have similar entity embeddings. In the example of Fig. 7, `lgd:5147` has a similar embedding with `:Germany`, since their attributes `label:Germany` have similar embeddings.

Having learned the structure embeddings and the attribute embeddings, AttrE combines them to produce the final entity embeddings, by minimizing the loss function

$$J_{SIM} = \sum_{e \in E_1 \cup E_2} [1 - \cos(\mathbf{e}_{se}, \mathbf{e}_{ce})], \quad (31)$$

where  $\mathbf{e}_{se}$  the structure embedding of  $e$ , and  $\mathbf{e}_{ce}$  is the attribute character embedding of  $e$ . For example, the embedding of the entities `lgd:5147` and `:Germany` will be updated in order their structured and attribute embeddings to be close. However, the entities `lgd:2401` and `dbp:Kromsdorf` have tail entities with similar embeddings. At the end, the entities `lgd:2401` and `dbp:Kromsdorf` will end up having similar embeddings too. For generating the alignment results, AttrE reports two entities as aligned, if the cosine similarity of their entity embeddings is greater than a pre-defined threshold.

## 2.5 Qualitative comparison of embedding methods

In this section, we describe the general assumptions of the evaluated methods, while we also compare the embedding-based entity alignment methods from different perspectives. For this purpose, we summarize in Tables 1 and 2 the basic characteristics and techniques of the methods, with respect to eight main categories: embedding module, literal size, alignment module, learning, schema alignment, embedding initialization, and negative sampling on relations and attributes, as described next.

### 2.5.1 Embedding module

The type of information that each method exploits in order to learn the embeddings reveals one of the most important assumptions that differentiates the methods. All embedding-based entity alignment methods utilize the relational structure of entities (relations), assuming that structural similarity is the key to entity alignment.

**Table 1** Method categories

Methods	Embedding module				Lit. Values	Lit. Size	Align. Module	Learning	Schema Align.
	Ent. Names	Relations	Attr. Names	Ent. Names					
MTransE	✓	One-hop	✓	✓	-	Mapping	Supervised	✓	
MTransE+RotatE	✓	One-hop	✓	✓	-	Sharing	Supervised	✓	
RDGCN	✓	Multi-hop	✓	✓	-	Mapping	Supervised	✓	
RREA(basic)	✓	Multi-hop	✓	✓	-	Sharing	Supervised	✓	
RREA(semi)	✓	Multi-hop	✓	✓	-	Sharing	Semi-supervised	✓	
KDCoE	✓	One-hop	✓	Descriptions	4 words	Mapping	Semi-supervised	✓	
MultiKE	✓	One-hop	✓	Word and character	5 words	Swapping	Supervised	Name and Semantic	
AttrE	✓	One-hop	✓	Character	10 characters	Sharing	Unsupervised	Name	
BERT_INT	✓	Multi-hop	✓	Word and description	128 Words	Sharing	Supervised	✓	

**Table 2** Initialization, hidden layers and negative sampling

Method	Embedding initialization	Layers	Negative sampling on relations negative per positive	Negative sampling on attributes negative per positive
MTransE	Random	2	-	-
MTransE + RotatE	Random	2	10 (uniform)	-
DCoE	Random	2 + 2 GRU	1 (uniform)	1 (uniform)
AttrE	Random	5	1 (uniform)	1 (uniform)
MultiKE	Random	2 lit. + 4 + 3 conv.	10 (uniform)	-
RDGCN	Name embeddings	4 GAT + 2 GCN	125 (uniform)	-
RREA(semi)	Random	2 GAT	1 (truncated)	-
RREA(basic)	Random	2 GAT	1 (truncated)	-
BERT_INT	Random	3 MLP + 4 RBF	2 (truncated)	-

However, there are methods that also exploit the literal values of the attributes, as well as the entity and attribute names.

Among the evaluated methods, MultiKE, RDGCN and BERT\_INT are the only methods that use entity names in order to learn the final entity embeddings. Thus, for these methods, we expect an improved performance on KGs where there is some homogeneity in the naming of the entities.

**2.5.1.1 Entity names** Some embedding-based entity alignment methods learn the entity embeddings by encoding the entity names in the embedding space, so as entities with similar names to have similar embeddings. These methods work under the assumption that the same real-world entities have to follow the same or semantically similar naming. However, in many real-world KGs, this assumption is not holding as entity names are replaced with source-specific ids. For this reason, several EA methods either avoid to utilise the entity names for the entity embeddings, either use them jointly with the structural or factual information of KGs.

**2.5.1.2 Relations** The structural information of the KGs constitutes the main type of information that all embedding-based entity alignment methods use in order to learn the entity embeddings. Entity embedding methods exploit the structural information either in *one hop* or in *multiple hops*. The former use the local structural information of the entities, ignoring the impact of more distant neighbors. For these methods, the more relations we have per entity, the better results we expect to get, because entities have the opportunity to minimize their embedding distance with multiple similar entities (neighbors), exploiting multiple features to learn the final embeddings (more information for learning). In addition, the more relation edges we have, the more negative samples we have per entity, and as a result, entities move away from dissimilar entities in the embedding space. The latter (multi-hop methods) use the subgraph structure, exploiting a larger amount of relations between entities. This way, they focus on the extended (multi-hop) neighborhood and aggregate the embeddings of the multi-hop neighbors in order to learn their own.

Among the methods evaluated in this study, all of them use the relations to learn the entity embeddings, with most methods focusing on the one-hop neighborhood, except RDGCN, RREA(basic), RREA(semi), and BERT\_INT, which follow a multi-hop approach. This way, RDGCN and both versions of RREA increase the information that an entity can exploit and the final entity embeddings are more expressive. This benefit comes at the risk of incorporating more noisy information into the embeddings, due to considering some distant neighbors, irrelevant or less important for the entity alignment task. Thus, BERT\_INT, rather than aggregating the information of multi-hop neighbors, compares entity pairs from multi-hop neighbors, called interactions, based on factual information (e.g., textual descriptions, names).

**2.5.1.3 Attribute names** Some embedding-based entity alignment methods use the attribute names for learning the entity embeddings. Usually, attribute names are exploited by the methods in order to enhance the similarity measure of entities and

not as main source of information for learning the embeddings. That is because it requires KGs to have similar naming schemes, which is not usual in real-world KGs.

MultiKE is the only method that utilizes this type of information, requiring a homogeneous naming in terms of attributes. Thus, we expect this feature to perform better in datasets originating from the same or similar data sources (e.g., Wikipedia-based, or synthetic KGs).

**2.5.1.4 Literal values and literal size** Many embedding-based entity alignment methods use literal values as auxiliary information to learn the entity embeddings. Specifically, they learn the embeddings of the literal values and they use those embeddings for learning the entity embeddings, or for enriching the seed alignment, or even for alignment. In most methods, all literal values are used, but in some methods, only the literal values of the attributes with name “description” are used. In addition, there are two methods for learning the literal-value embeddings: character-based and word-based. The first one uses the characters of the literal to learn the final literal embedding, while in the second one, the literal is tokenized (split into words) and the final embedding is resulting by employing pre-trained word embeddings.

The character-based methods usually exploit a small part of the literals (e.g., the first few characters) and they are mostly used for short literals, such as dates. In addition, there are many methods to learn the character-based embeddings: the *n-gram*-based that sums the *n-gram* combinations of the literals and the method that computes the average of the pre-trained character embeddings. The word-based methods typically exploit more information from the literal value (e.g., the first few words) and for this reason, they are mostly used for longer literals, such as names and descriptions. However, the word-based methods assume the existence of the word embeddings in the pre-trained set for all the words that appear in the literal values, which is very frequently not the case, resulting in out-of-vocabulary errors. For this reason, there are also hybrid methods, that use the word-based method first, but if a word does not exist in the pre-trained set (vocabulary), they use one of the aforementioned character-based method, in order to learn the specific word’s embedding. Regarding the literal size, it is a very important hyperparameter that defines the size of the substring of the literal which will be used, including or excluding important words or characters for the final literal embeddings.

Among the evaluated methods, KDCoE, AttrE, MultiKE and BERT\_INT are the ones that use literal values for learning embeddings. KDCoE exploits the first 4 words of textual descriptions and pre-trained word embeddings, in order to learn the description embeddings and measure the entities similarity based on these embeddings. The objective is to enrich the training set (Sect. 3.4) with extra aligned entity pairs based on this similarity. Thus, if the similarity of an entity pair (based on description embeddings) is above a pre-defined threshold (95%), the entities are proposed as aligned, enriching the training set with an extra aligned entity pair that was not included in the training set before. AttrE, in contrary, exploits the first 10 characters of the literal values and follows an *n-gram* based method to learn the embeddings, for a combination of *n* values ranging from 1 to 10. The literal embeddings are used in the alignment module in order entities of the KGs with similar literal embeddings to be close in the embedding space. Specifically, AttrE uses literal



embeddings for the alignment, since, as an unsupervised method, it does not use the seed alignment during training process. Furthermore, MultiKE uses the first 5 words and pre-trained word embeddings for learning the literal embeddings. In case a word is not contained in the available pre-trained word embeddings, it averages the character embeddings of this word, i.e., it follows the hybrid approach described above. MultiKE uses the literal embeddings in many ways, as described in Sect. 2.4.1, e.g., for name view, attribute view and schema alignment. BERT\_INT utilizes long textual descriptions (limiting the words to 128), names, or literals for both BERT unit and interaction model, as described in Sect. 2.4.1.

All these methods require homogeneity in terms of literal values, which is not usual in real-world KGs. Finally, the character-based method that MultiKE uses as alternative to word-based methods suffers in the literals that contain the same characters in different order. For example, the number “5013” and the number “1350” will have the same embeddings. AttrE does not have this issue, since it is following an n-gram based method for character embeddings. BERT\_INT, leveraging a variety of factual information and giving priorities depending on their availability and quality, deals with heterogeneity issues (e.g., it prioritizes descriptions over entity names, since descriptions contain richer information).

### 2.5.2 Alignment module

The alignment module aims to produce a common embedding space for the entities of the two KGs, using three typical techniques described in Sect. 2.2.2. MTransE+RotatE, AttrE, both versions of RREA, and BERT\_INT utilize the sharing technique to calibrate the axis of the embedding spaces of the two KGs, in order the aligned entities to overlap. Consequently, we expect methods that rely on this technique to work well on KGs with dense and similar neighborhoods. This setting helps them to achieve the desired overlapping and preserve the initial structure of the KGs in the embedding space. On the other hand, MTransE, KDCoE and RDGCN, that rely on the mapping technique, learn a transformation matrix that aligns entities in two separate embedding spaces. In contrary to sharing, this technique does not assume the spatial similarity of the KGs. Finally, methods such as MultiKE, that unify multiple embedding spaces (multi-view) and consume more relation edges for training, rely on the swapping alignment technique.

### 2.5.3 Learning

We divide the evaluated methods into supervised, that require the seed alignment for the alignment module, semi-supervised, that use the seed for the alignment process, but also try to enrich it, and unsupervised, that do not use the seed alignment for the alignment process. All methods need seed alignment, either for training (alignment module) and testing (supervised and semi-supervised) or exclusively for testing (unsupervised). Supervised methods (MTransE, MTransE+RotatE, RDGCN, MultiKE RREA(basic) and BERT\_INT) assume the existence of seed alignment, which in many cases is hard to find (since it includes the matches of all the entities

of the two KGs), even in widely used Wikipedia-based KGs, hindering the EA task especially when the KGs scale up in terms of content and density.

In contrary to the aforementioned category, semi-supervised methods (e.g., KDCoE and RREA(semi)) work sufficiently well when a small percentage of aligned entities is available, which is more realistic in real-world KGs. However, these methods usually require auxiliary information to support the alignment task that involves the enrichment of the seed alignment. For example, KDCoE does this by assuming the existence of multiple, homogeneous textual descriptions. On the other hand, RREA(semi) does not use auxiliary information; it utilizes specific rules (e.g., the two entities have to be mutually nearest aligned, in order to be proposed as aligned).

Finally, unsupervised methods carry out a more difficult task, since they use no seed alignment for the alignment process. Particularly, the lack of seed alignment makes AttrE leverage the literals similarity for alignment purposes, assuming their homogeneity. It is worth noting that unsupervised methods, in contrary to supervised and semi-supervised, do not assume that every entity of  $KG_1$  should be matched to exactly one entity of  $KG_2$ .

#### 2.5.4 Schema alignment

We refer to schema alignment as the process of finding similar relations and/or attributes. The similarity of the relations or the attributes can be measured by two different ways: measuring their name similarity (AttrE) or measuring both their name similarity and their structural (semantic) similarity (MultiKE) in the embedding space. All these methods can benefit if the initial KGs follow similar naming schemes, making the process of finding similar relations and attributes easier.

AttrE and MultiKE are two methods that perform schema alignment. MultiKE captures the schema similarity not only in terms of relation or attribute name similarity, but also exploits the semantic similarity of relations and attributes. In the case that we have KGs with heterogeneous naming at the schema level, MultiKE would have better performance than AttrE, that only exploits the name similarity for the schema alignment.

#### 2.5.5 Embedding initialization

For initializing the entity embeddings, there are two practices: the random one, that initializes them by picking random numbers from a normal distribution, and the one that initializes them with the entity name embeddings. The latter has better results (Kocmi and Bojar 2017), because embeddings contain some information for the entities, but it works under the assumption that matching entities have similar entity names, which is not holding in many cases. RDGCN is the only method that initializes the embeddings with entity name embeddings.

### 2.5.6 Negative sampling

Negative sampling (Kamigaito and Hayashi 2022) is the process of generating  $n$  (a hyperparameter) negative examples of edges that do not exist in the KG, by replacing either the head or the tail entities of each positive edge with another random entity - uniform negative sampling (Bordes et al. 2013) - or with a highly similar neighbor - truncated negative sampling (Sun et al. 2018; Zhu et al. 2019; Cao et al. 2019). Truncated negative sampling, in contrary to uniform negative sampling, ensures difficult negative samples that contribute more on learning process than the easy ones. For example, if we sample a negative triple (Titanic, capitalOf, Greece) from (Athens, capitalOf, Greece), using uniform sampling, we can find that Athens and Titanic have very low similarity and the generated negative sample contributes little in the learning process. If, instead, we generate the negative triple (Thessaloniki, capitalOf, Greece) from (Athens, capitalOf, Greece), using truncated sampling, this negative triple would contribute more in the learning process, since Thessaloniki (the second biggest city in Greece) and Athens, are two very similar entities.

In general, negative sampling is widely used in KG embedding models on relations, as well as on attributes, in order to maximize the embedding distance of dissimilar entities. The more negative examples we have per positive example, the bigger the distance among dissimilar entities. In addition, an entity that participates in a high number of positive examples (high average relations per entity), will also participate in an increased number of negative examples. For instance, if an entity appears in the head or the tail of  $x$  relation edges, then the number of the negative examples for this entity can be as high as  $x * n$ , where  $n$  is the corresponding hyperparameter.

All the evaluated methods use negative sampling, with the exception of MTransE. The higher the ratio of negative per positive triples, the better the performance of the methods, since the embedding distance of dissimilar entities increases. Although requiring 10 or even 125 times bigger KGs (for 10 or 125 negative samples per positive, respectively) improves the performance, this comes at the cost of increased training time and less scalable methods. MultiKE further enriches the KGs by using the swapping alignment technique that generates additional positive examples, as described in Sect. 2.2.2.2. This technique not only increases the available information for training, but it also increases the negative examples for the entities indirectly, as entities will appear in more relation edges.

### 2.5.7 Neural network architectures

All evaluated methods are implemented on top of a neural network architecture with multiple hidden layers, each consisting of  $d$  neurons, where  $d$  are the dimensions of the embeddings (dim). More precisely, relation-based methods utilize either shallow neural networks (MTransE, MTransE+RotatE), or Graph Neural Network (RDGCN, RREA(basic), RREA(semi)) for learning the embeddings, while attribute-based methods use shallow neural networks for encoding relations (AttrE) or Graph Neural Networks for encoding literals (MultiKE) and textual descriptions (KDCoE) or Radial Basis Function kernel (Xiong et al. 2017) for the interactions (BERT\_INT).

Regarding relation-based methods, MTransE and MTransE+RotatE, utilize one layer for embedding entities and one layer for embedding entity relations. As shown in Table 7, the dimensions of both embeddings are the same (100). On the other hand, RDGCN consists of four Graph Attention Networks (GATs) (Velickovic et al. 2018), two primal layers and two dual layers. The final embeddings are output by two stacked GCN layers (Kipf and Welling 2017). Both versions of RREA utilize two stacked GATs (Relational Reflection Aggregate Layer). It is worth mentioning that both RDGCN and both versions of RREA use variations of graph attention mechanisms, not standard GATs.

Attribute-based methods, as shown in Table 2, use two standard layers (for entity and relation embeddings) and some extra layers to encode additional information of entity descriptions. For example, KDCoE utilizes two layers for embedding entities and their relations and two GRU layers for encoding their textual descriptions. It is worth mentioning that while all other examined methods learn the structure and attribute embeddings jointly, KDCoE is the only method that trains two different models alternatively (co-training). AttrE consists of five layers: one layer for embedding entities of relation edges, one layer for embedding entities of attribute edges, one layer for relations, one layer for attributes, and one layer for character embeddings. All of these layers contain the same number of neurons (100). MultiKE utilizes two layers for the word embeddings, three layers for entity, relation and attribute embeddings, one convolution layer for encoding literals, and one layer that combines them. Finally, it uses two additional convolution layers for the cross-KG entity inference based on relations and attributes. BERT\_INT utilizes four Radial Basis Function kernels (one layer each one) for the interaction model and three Multilayer Perceptron for obtaining the entity embeddings.

It is worth mentioning that the simpler the architecture, the more scalable the method is (e.g., as shown in Table 14, the translation-based method MTransE is more scalable than the GNN-based method RDGCN).

### 3 Experimental setting

In this section, we give an overview and statistics about the datasets, we describe the evaluation protocol and the metrics that we used. We also provide implementation details about this study. The source code of this work is publicly available.<sup>3</sup>

#### 3.1 Datasets

For the experimental evaluation we utilized nine datasets, each consisting of a pair of KGs to be aligned, four from OpenEA (Sun et al. 2020) and five datasets from Efthymiou et al. (2015), Obraczka et al. (2021) in order to enrich the available data characteristics and ensure a trustworthy meta-level analysis.

<sup>3</sup> <https://github.com/fanourakis/experimental-review-EA>.

For the datasets from OpenEA, three well-known KGs were used as sources: DBpedia (version 2016-10) (Lehmann et al. 2015), Wikidata (version 20160801) (Vrandečić and Krötzsch 2014) and YAGO3 (Rebele et al. 2016). The seed alignment is constructed by the *owl:sameAs* links, which were provided with those KGs. In addition, it is hard for embedding-based approaches to run on full KGs due to the large candidate space, hence the authors of OpenEA sampled real-world KGs, by their iterative degree-based sampling algorithm described in Sect. 3.4. The resulting KGs have a size of 15K entities each, constituting the datasets **D\_W\_15K\_V1** and **D\_Y\_15K\_V1**. In order to examine the behaviour of the methods with respect to the density of the KGs, the authors of OpenEA generated an additional dense version for each dataset (**D\_W\_15K\_V2** and **D\_Y\_15K\_V2**). For generating the dense version, they deleted entities (nodes) with low ( $\leq 5$ ) node degrees and they performed again their iterative degree-based sampling algorithm, described in Sect. 3.4. The above mentioned datasets are described as follows:

- **D\_W\_15K\_V1** and **D\_W\_15K\_V2** are the sparse and the dense datasets, respectively, that were constructed from DBpedia and Wikidata KGs, describing actors, musicians, writers, films, songs, cities, football players and football teams. We chose them in order to examine the influence of KG density variations in the relation- and attribute-based methods. In addition, these datasets have both low entity name similarity and a low number of entities that have long textual descriptions, undermining the methods that use them. In these datasets, relation and attribute names have been replaced with special ids, which influences in a negative way methods that perform schema alignment using predicate names.
- **D\_Y\_15K\_V1** and **D\_Y\_15K\_V2** are the sparse and the dense datasets, respectively, that were constructed from DBpedia and YAGO3. They contain the same entity types with **D\_W\_15K** but unlike **D\_W\_15K**, they have both high number of entities that have long textual descriptions and high entity name similarity.

Five additional datasets, from Efthymiou et al. (2015), Obraczka et al. (2021), were employed, that originate from the following eight KGs: the BTC2012 version of DBpedia,<sup>4</sup> BBCmusic (originating from Kasabi<sup>5</sup>), IMDb,<sup>6</sup> TVDB,<sup>7</sup> TMDb,<sup>8</sup> and two restaurant KGs.<sup>9</sup> None of these datasets originally cover the assumptions of embedding-based entity alignment methods, thus we pre-process them, as described in Sect. 3.4. The five extra datasets, built from those KGs, are the following:

- **BBC-DB** (Efthymiou et al. 2015) is a sparse dataset constructed by BBCmusic and the BTC2012 version of DBpedia. It contains various entity types such as musicians, their birth places and bands, while also it has the highest number of

<sup>4</sup> <http://km.aifb.kit.edu/projects/btc-2012/>.

<sup>5</sup> <https://archive.org/details/kasabi>.

<sup>6</sup> <https://www.imdb.com/>.

<sup>7</sup> <https://www.thetvdb.com/>.

<sup>8</sup> <https://www.themoviedb.org/>.

<sup>9</sup> <http://oaei.ontologymatching.org/2010/im/>.

**Table 3** Statistics per dataset

Datasets	#Rel_Triples	#Attr_Triples	#Rel_Types	#Attr_Types	Sym- metric_ Rels
D_W_V1 KG1	38,265	68,258	248	342	7
D_W_V1 KG2	42,746	138,246	169	649	2
D_W_V2 KG1	73,983	66,813	167	175	1
D_W_V2 KG2	83,365	175,686	121	457	0
D_Y_V1 KG1	30,291	71,716	165	257	4
D_Y_V1 KG2	26,638	132,114	28	35	2
D_Y_V2 KG1	68,063	65,100	72	90	0
D_Y_V2 KG2	60,970	131,151	21	20	0
BBC-DB KG1	15,478	18,165	9	4	3
BBC-DB KG2	45,561	149,720	98	723	1
imdb-tmdb KG1	4,747	10,279	2	13	0
imdb-tmdb KG2	4,857	8,947	2	29	0
imdb-tvdb KG1	3,389	5,648	2	13	0
imdb-tvdb KG2	1,051	3,535	2	8	0
tmdb-tvdb KG1	3,396	5,326	2	24	0
tmdb-tvdb KG2	938	4,096	2	8	0
Restaurants KG1	226	1,504	2	4	0
Restaurants KG2	1,504	3,760	2	4	0

average attributes per entity and the highest number of entities that have textual descriptions among other datasets. In addition, this dataset is better suited for the methods that perform schema alignment, since these KGs follow similar predicate naming conventions.

- **imdb-tmdb**, **imdb-tvdb**, and **tmdb-tvdb** (Obraczka et al. 2021) were constructed from IMDb, TMDb, and TVDB that contain descriptions about people, movies, tv series, episodes and companies. Those datasets have a low number of entities, which makes them interesting for examining the performance of supervised vs unsupervised methods with a small training dataset. In addition, they are much sparser than the previous datasets, with low average attributes per entity, while they exhibit high homogeneity in terms of literals and predicate names. tmdb-tvdb features higher homogeneity in terms of literals and lower similarity of long-text descriptions, compared to imdb-tvdb and tmdb-tvdb.
- **Restaurants**<sup>9</sup> contains descriptions of real restaurants and their addresses from two different KGs. It is the smallest dataset in terms of number of entities, relation types and attribute types, but its main challenge is that it does not cover the assumptions of all the embedding-based entity alignment methods. Specifically, not all entities of KG1 are matched with all the entities KG2 and vice versa. We chose this dataset to examine whether the evaluated methods can be generalized beyond datasets that satisfy the 1-to-1 assumption.

**Table 4** Statistics per dataset for the Meta-features

Datasets	Seed alignment size	Density			Heterogeneity						
		#Entity_Pairs (KG1/KG2)	#Avg_Rels_per_Entity (KG1/KG2)	#Avg_Attrs_per_Entity (KG1/KG2)	Sole_Rels (KG1/KG2)	Hyper_Rels (KG1/KG2)	#Ents_Descr	Descr_Sim	Ent_Name_Sim	Lit_Sim	Pred_Name_Sim
D_W_V1	15K / 15K	2.86 / 3.15	4.69 / 9.61	91 / 74	157 / 95	1,181 / 9,032	0.69	0.72	0.81	0	
D_W_V2	15K / 15K	5.39 / 5.77	4.52 / 11.97	53 / 44	114 / 77	278 / 3,579	0.73	0.78	0.83	0	
D_Y_V1	15K / 15K	2.8 / 3.58	4.89 / 8.8	53 / 5	112 / 23	3,455 / 12,000	0.88	0.99	0.78	47.37	
D_Y_V2	15K / 15K	6.12 / 10.98	4.38 / 8.74	17 / 6	55 / 15	2,945 / 12,000	0.87	0.91	0.79	47.99	
BBC-DB	9,396 / 9,396	1.72 / 6.03	1.94 / 20.15	9 / 20	0 / 78	2,430 / 2,133	0.62	0.34	0.87	80.93	
imdb-tmdb	1,933 / 1,933	2.65 / 2.71	5.31 / 4.62	2 / 2	0 / 0	753 / 753	0.35	0.48	0.92	88.49	
imdb-tvdb	1,076 / 1,076	3.22 / 1	5.25 / 3.28	2 / 2	0 / 0	472 / 473	0.38	0.43	0.91	79.96	
tmdb-tvdb	982 / 982	3.62 / 1	5.42 / 4.17	2 / 2	0 / 0	180 / 180	0.22	0.84	0.94	80.33	
Restaurants	112 / 112	1 / 2	1.66 / 1.66	2 / 2	0 / 0	90 / 73	0	0.95	0.96	80.11	

**Table 5** Meta-features of the datasets

Datasets	Seed alignment size		Density		Heterogeneity						
	#Entity_Pairs of seed alignment	Avg_Rels_per_Entity	Avg_Attrs_per_Entity	Sole_Rels	Hyper_Rels	#Ents_Descr	Descr sim	Ent_Name_Sim	Lit_Sim	Pred_Name_Sim	
D_W_V1	15,000	3.00	7.15	39.56	60.43	1,181	0.69	0.72	0.81	0	
D_W_V2	15,000	5.58	8.24	33.68	66.31	278	0.73	0.78	0.83	0	
D_Y_V1	15,000	3.19	6.84	30.05	69.94	3,455	0.88	0.99	0.79	47.37	
D_Y_V2	15,000	8.55	6.56	24.73	75.26	2,945	0.87	0.91	0.78	47.99	
BBC-DB	9,396	3.87	11.04	27.10	72.89	2,430	0.62	0.34	0.87	80.93	
imdb-tmdb	1,933	2.68	4.96	100	0	753	0.35	0.48	0.92	88.49	
imdb-tvdb	1,076	2.11	4.26	100	0	472	0.38	0.43	0.91	79.96	
tmdb-tvdb	982	2.31	4.79	100	0	180	0.22	0.84	0.94	80.33	
Restaurants	112	1.50	1.66	100	0	90	0	0.95	0.96	80.11	



**Table 6** Aggregation functions

Statistics	Aggregation functions
Avg_Rels_per_Entity	$(\text{avg\_rel\_per\_entity\_KG1} + \text{avg\_rel\_per\_entity\_KG2}) / 2$
Avg_Attrs_per_Entity	$(\text{avg\_attr\_per\_entity\_KG1} + \text{avg\_attr\_per\_entity\_KG2}) / 2$
Sole_Rels	$(\text{sole\_KG1} + \text{sole\_KG2}) / \text{rels.} * 100$
Hyper_Rels	$(\text{hyper\_KG1} + \text{hyper\_KG2}) / \text{rels.} * 100$
#Ents_Descr	$\text{ents\_with\_descr\_KG1}$

### 3.2 Statistics and meta-features

In order to quantify useful information about the knowledge graphs, we calculated some basic dataset statistics regarding the type and the number of relations and attributes that are presented in Table 3 and some additional dataset statistics that are shown in Table 4, categorized in the three following categories: seed alignment size, density and heterogeneity that were used for constructing the meta-features of Table 5.

*Seed alignment size* refers to the number of entity pairs in the seed alignment. *Density* consists of the average relations per entity and the average attributes per entity, which indicate the density of the knowledge graphs on structural and attribute level, respectively. In addition, it consists of the proportion of sole and hyper relation types of the two knowledge graphs, where a sole relation type is a relation type that does not co-occur with another relation type in any entity pair (Zhang et al. 2017), otherwise, a relation is called a hyper relation (Zhang et al. 2017). Finally, in order to measure the *heterogeneity* of the two knowledge graphs, we calculate the textual description similarity (Descr\_Sim) in the embedding space, of entities that are included in the test and validation sets, the number of entities of the two KGs that have textual descriptions (#Ents\_Descr), the similarity of entity names (Ent\_Name\_Sim), the similarity of literals (Lit\_Sim) and the similarity of predicate names (Pred\_Name\_Sim). Specifically, for Descr\_Sim, Ent\_Name\_Sim and Lit\_Sim, we report the average of MAX similarities of each description, entity name and literal, respectively. As for the Pred\_Name\_Sim, we measure the name similarity of relation and attribute types (predicates) using Levenshtein Distance, as suggested in AttrE (Trisedya et al. 2019). Specifically, we report the average of MAX similarities of each predicate (relation or attribute).

The meta-features of Table 5 were constructed by applying the aggregation functions of Table 6 on the statistics of Table 4, in order to construct features for each individual dataset instead of each individual KG of the dataset and use them in Sect. 4.3. More precisely, for Avg\_Rels\_per\_Entity, Avg\_Attrs\_per\_Entity, Sole\_Rels, Hyper\_Rels and #Ents\_Descr (based on KDCoE that proposes possibly aligned entities of the entities of the first KG) were used aggregation functions, while #Entity\_Pairs, Descr\_Sim, Ents\_Name\_Sim, Lit\_Sim and Pred\_Name\_Sim were used as meta-features without the application of any aggregation function.

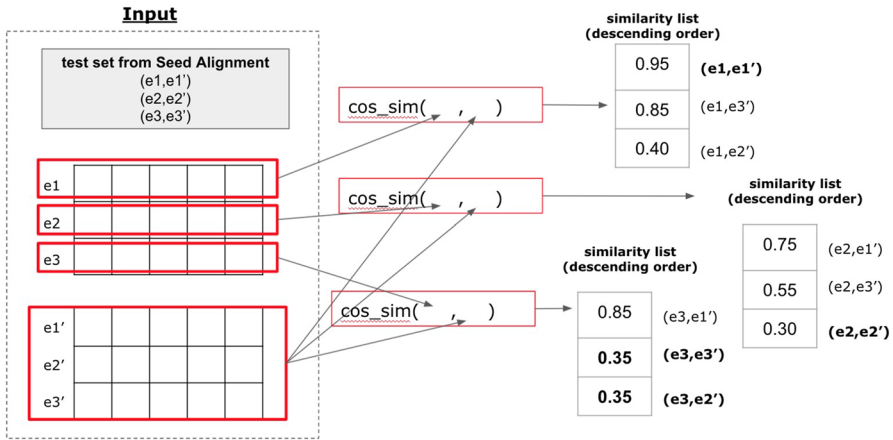


Fig. 8 Evaluation process

### 3.3 Evaluation protocol and metrics

In this section, we describe the evaluation protocol that embedding-based entity alignment methods use, the metrics used for their evaluation, while we also compare those metrics.

#### 3.3.1 Evaluation protocol

EA methods have been traditionally evaluated using classification-based metrics (e.g., precision and recall), comparing the matches proposed by those methods to the correct matches of a given ground truth (Leone et al. 2022). On the contrary, KG embedding-based EA methods, adopted rank-based evaluation metrics (see Sect. 3.3.2), possibly influenced by the recent literature on the embedding-based methods for link prediction (Bordes et al. 2013; Sun et al. 2019). In order to compare embedding-based EA methods on an equal basis and re-use the available open-source code as much as possible, we adopt the latter evaluation metrics and acknowledge that the criticism of Leone et al. (2022) is valid.

During the evaluation, the embedding-based EA methods use two embedding matrices, one for each KG, with the entity embeddings in their rows, and the test set<sup>10</sup> of the seed alignment (Fig. 8). The alignment methods calculate the similarity of every entity of the first matrix with every entity of the second matrix, in the embedding space using similarity measures such as *Euclidean distance*, *cosine similarity*, etc. The result of this process is a similarity list, that we sort in descending order and aim to find the index (ranking) of the aligned entity in the similarity list.

In the example of Fig. 8,  $e'_1$  is the aligned entity of entity  $e_1$ , known as *true entity*, according to seed alignment, and it is in the first position in the similarity list. In the

<sup>10</sup> The seed alignment is split into train (20%), validation (10%) and test (70%), as in (Sun et al. 2020).

last similarity list, there is a tie in the scores of the two candidate entities  $e'_3$  and  $e'_2$ , that is decisive for some evaluation metrics. According to Berrendorf et al. (2020), there are different behaviors for this phenomenon that are categorized as *optimistic*, *pessimistic*, *non-deterministic* and *realistic*. In optimistic behavior, it is assumed that the true entity is ranked first among other entities with the same score, in pessimistic, it is assumed that it is ranked last, and realistic is the mean of optimistic and pessimistic. In our case, all of the evaluated methods are non-deterministic, which means that handling ties depends on the inner working of the sorting algorithm that the methods use to sort each similarity list.

In the context of ensuring a reliable and unbiased experimental process, we conducted the experiments with 5-fold cross-validation (Sun et al. 2020), where each fold is split into train, validation and test<sup>10</sup>. Finally, the performance of each method, per dataset, for a specific metric is calculated by the average scores that the method achieves in the 5 folds for this metric.

### 3.3.2 Evaluation metrics

For the evaluation of the alignment methods, we use the following metrics, which depend on the individual ranks  $\mathcal{I}$  generated by each method, as described in the previous section.

**Hits@k** describes the fraction of hits (true entities) that appear in the first  $k$  ranks of the sorted similarity lists:

$$\text{Hits@}k = \frac{|\{r \in \mathcal{I} \mid r \leq k\}|}{|\mathcal{I}|}, \quad (32)$$

where  $\text{Hits@}k \in [0, 1]$ . In the experimental results, we report Hits@k values as percentages (i.e.,  $\text{Hits@}k \cdot 100\%$ ). Specifically, this metric measures the accuracy of the methods with adjustable error rate. For instance, Hits@10 allows a small error rate in contrary to Hits@1 that does not allow any errors. The weakness of this metric is that it considers only the first  $k$  position of the similarity list, and as a result, all the other positions have no effect to the final score. An advantage of this metric is its easy interpretation. In the example of Fig. 8,  $\text{Hits@}1 = 1/3 = 0.33$ , since among the three similarity lists ( $|\mathcal{I}| = 3$ ), only in one similarity list, a true entity pair is observed in the first rank.

**Mean Rank (MR)** computes the average ranks of the true entity pairs:

$$\text{MR} = \frac{1}{|\mathcal{I}|} \sum_{r \in \mathcal{I}} r, \quad (33)$$

where  $\text{MR} \in [1, |E|]$  (the lower the better). An advantage of MR is its sensitivity to any model performance changes, since it reflects the average performance and not only the first  $k$  ranks. However, it is also sensitive to outliers, because a very high or very low rank of a true entity pair affects the mean of the ranks. Finally, MR is an easily interpretable metric, but we should keep in mind the size of the candidate entities, because MR = 10 is better for 1,000,000 entities than for 20 entities. In the

example of Fig. 8, assuming that in the third similarity list the true pair is in the second rank,  $MR = (1 + 3 + 2)/3 = 2$ .

**Mean Reciprocal Rank (MRR)** is the inverse of harmonic mean rank:

$$MRR = \frac{1}{|\mathcal{I}|} \sum_{r \in \mathcal{I}} \frac{1}{r}, \quad (34)$$

where  $MRR \in (0, 1]$ . This metric is affected more by the top-ranked values rather than the bottom ones. For instance, the change is larger when moving from 1 to 2 compared to moving from 10 to 1,000. Thus, MRR is less sensitive to outliers. In the example of Fig. 8, assuming that in the third similarity list the true pair is in the second rank,  $MRR = (1/1 + 1/3 + 1/2)/3 = 0.61$ .

### 3.4 Pre-processing pipelines

The embedding-based entity alignment methods are working under some strict assumptions that make them inapplicable to the five new datasets. However, for a reliable experimental evaluation of the methods under the same conditions and data characteristics, we pre-process those datasets in order to conform to the assumptions described in Sect. 2.

**1-to-1 and structural-information assumptions** Firstly, we processed the seed alignment, in order to satisfy the 1-to-1 mapping assumption, as follows. If there are two entity pairs  $(e_i, e_j)$  and  $(e_i, e'_j)$  in the seed alignment, we need to keep one of them. In the dilemma of choosing which entity pair to remove, we examined two factors. The first factor is to keep the entity pair in which the true entity,  $e_j$  or  $e'_j$  in this example, is involved in at least one relation edge (one of the other assumptions of these methods). The second factor is to keep the entity pairs in which the true entity has the highest node degree (number of incoming and outgoing edges), in the effort to reduce as little as possible the initial KGs. In addition, we processed the initial KGs, by removing relation edges  $(h, r, t)$ , whose head ( $h$ ) or tail ( $t$ ) entities did not appear in the seed alignment. Then, we repeated the same process for the attribute edges, removing the edges whose head entity did not appear in the seed alignment, and also removing the attribute edges whose head entity was not involved in at least one relation edge.

**Pre-aligned predicates assumption** In addition to these assumptions, AttrE requires predicate alignment, as described in Sect. 2.4 and merging the two KGs into one KG. The predicate alignment step is performed by manual inspection of the data, as suggested by the authors of AttrE.<sup>11</sup> For example, given the attributes *:name* and *:hasName*, the second one is replaced by *:name*.

**Sampling.** For sampling real-world KGs, the authors of OpenEA proposed an iterative degree-based sampling (IDS) algorithm which simultaneously removes entities of the two KGs using seed alignment. The objective of this algorithm is to generate a dataset of a specific size from real-world KGs, such that the difference of

<sup>11</sup> <https://bitbucket.org/bayudt/kba/src/35c67d56a8f0e2bcc05e8d56fb98c4374e3ff542/>.

**Table 7** Hyperparameters of methods

Methods	Dim	Batch size	Max epochs
MTransE	100	5,000	2,000
MTransE+RotatE	100	5,000	2,000
RDGCN	300	5,000	2,000
RREA(basic)	100	5,000	1,200 (fix)
RREA(semi)	100	5,000	6,000 (fix)
KDCoE	100	5,000	2,000
MultiKE	100	5,000	2,000
AttrE	100	100	50 (fix)
BERT_INT	300	128	205 (fix)

the entity degree distribution, between the original KG and the sampled KG, does not exceed a certain value. In order to determine which entities to remove without greatly affecting the entity degree distribution, they delete entities with low PageRank scores, while assessing the divergence of the two degree distributions of the KGs using the Jensen-Shannon (JS) (Lin 1991) measure. The algorithm stops when the size of the datasets is the desired one (15K) and the divergence is  $\leq 5\%$ .

### 3.5 Implementation details

All experiments were performed in a server with 16 cores (AMD EPYC 7232P @ 3.1 GHz), 64 GB RAM, one RTX-4090 GPU (24 GB) and Ubuntu 18.04.5 LTS. We used the code of OpenEA<sup>12</sup> (2020-07-28) for the experiments of methods MTransE, MTransE + RotatE, KDCoE, MultiKE and RDGCN. For AttrE, we noticed that the implementation of OpenEA utilizes the seed alignment for training, which contradicts the unsupervised learning nature of AttrE (Trisedya et al. 2019). For this reason, we used the code provided by the authors of AttrE<sup>11</sup> (2019-01-11). Finally, for RREA (basic and semi), BERT\_INT and PARIS, we used the code provided by the authors.<sup>13,14,15</sup>

For conducting a fair experimental process, we used the same hyperparameter values as described in the works of OpenEA, AttrE, RREA and BERT\_INT, and are shown in Table 7. *Dim* is the dimensions of the embeddings (length of the vector representations). The higher dimensions an embedding has, the more features (expressiveness) it has. *Batch size* is the number of relation edges that we used in each iteration during training. All the examined methods adopt a batch size of 5,000, except AttrE which adopts 100 and BERT\_INT 128. *Max epochs* refers to the times the learning algorithm will work in the entire dataset. In our experiments,

<sup>12</sup> <https://github.com/nju-webssoft/OpenEA/tree/2a6e0b03ec8cdcad4920704d1c38547a3ad72abe>.

<sup>13</sup> <https://github.com/MaoXinn/RREA>.

<sup>14</sup> <https://github.com/kosugil1037/bert-int>.

<sup>15</sup> <https://github.com/epfl-dlab/entity-matchers>.

there are two ways to define the number of epochs. The first one is to give a fixed number (50, 205, 1,200 or 6,000 epochs), as AttrE, BERT\_INT, RREA(basic) and RREA(semi) do. The second one, *early stopping*, is to terminate the training process when Hits@1 starts dropping, based on the validation set, checked every 10 epochs. The latter ensures that the model does not overfit. It is worth noting RREA(semi) requires a fixed number of iterations (5), while KDCoE decides the number of iterations dynamically. RREA(basic) is the first iteration of RREA(semi).

## 4 Analysis of experimental results

In this section, we report and analyze the results of a series of experiments we conducted to answer the open questions Q1-Q4 posed in Sect. 1, regarding embedding-based entity alignment methods. More precisely, to shed light regarding Q1 (characteristics of methods) and Q2 (families of methods), we analyze the performance of the methods described in Sect. 2 across all datasets (Sect. 4.1). To answer Q3 (effectiveness-efficiency trade-off), we rank the methods in a statistically significant way, while we analyze their corresponding execution time over all datasets of our testbed, and compare the time they require to reach a relative effectiveness threshold (Sect. 4.2). Finally, in Sect. 4.3, we address Q4 (characteristics of datasets), by conducting a meta-level analysis to identify correlations between the methods and the various meta-features extracted from the KGs of our testbed.

### 4.1 Effectiveness of EA methods using different metrics

To answer Q1 and Q2, we experimentally compare five relation-based (MTransE, MTransE+RotatE, RDGCN, RREA(basic) and RREA(semi)) and three attribute-based methods (KDCoE, MultiKE and AttrE), using the protocol and metrics described in Sect. 3.3. Table 8 reports the different evaluation metrics per method over all datasets of our testbed.

#### 4.1.1 Relation-based EA methods

**MTransE** exploits one-hop entity neighborhoods and does not use negative sampling. As we can see in Table 8, the performance of MTransE remains relatively the same between the sparse version D\_W\_15K\_V1 with 3 average relations per entity<sup>16</sup> (Table 5), and the dense one D\_W\_15K\_V2 with 5.58 average relations per entity, as well as between the sparse D\_Y\_15K\_V1 with 3.19 average relations per entity, and the dense D\_Y\_15K\_V2 with 8.55 average relations per entity. In general, entities in dense KGs are more likely to minimize their embedding distance with multiple similar entities (neighbors), exploiting richer semantic information to learn the embeddings. For this reason, dense KGs are the best use cases for relation-based methods. However, the density of KGs is also reflected in the embedding

<sup>16</sup> A KG is considered “dense” in Sun et al. (2020) if it has more than an average of 5 relations per entity.

**Table 8** Performance of methods in different datasets and metrics

Datasets	Dataset	Relation-based methods				Attribute-based methods				
		MTransE	MTransE+RotatE	RDGCN	RREA(basic)	RREA(semi)	KDCoE	MultiKE	AttrE	BERT_INT
OpenEA	D_W_15K_V1	H@1 26.07	27.27	53.02	65.57	<b>71.82</b>	23.22	40.49	N/A	44.08
		H@10 54.05	56.98	72.96	88.49	<b>90.03</b>	46.29	56.71	N/A	48.96
		MR 352.92	501.58	494.69	46.32	57.71	868.85	320.29	N/A	<b>13.27*</b>
		MRR 0.35	0.36	0.59	0.74	<b>0.79</b>	0.30	0.45	N/A	0.45
	D_W_15K_V2	H@1 26.27	49.22	63.71	87.86	<b>93.72</b>	33.78	49.54	N/A	42.63
		H@10 57.45	81.69	82.33	98.66	<b>99.18</b>	64.38	72.76	N/A	48.50
OpenEA	D_Y_15K_V1	MR 145.53	58.53	210.94	<b>2.64</b>	2.74	169.46	37.44	N/A	10.38*
		MRR 0.36	0.59	0.69	0.92	<b>0.96</b>	0.44	0.56	N/A	0.44
		H@1 45.54	46.25	94.10	76.19	81.79	69.57	90.18	0.03	<b>99.23</b>
		H@10 72.78	72.94	97.68	91.31	91.95	88.02	95.03	0.2	<b>99.28</b>
		MR 245.02	574.13	14.87	29.79	73.34	143.78	19.69	5077.59	<b>1.0004*</b>
		MRR 0.55	0.55	0.95	0.82	0.86	0.75	0.91	0.001	<b>0.99</b>
OpenEA	D_Y_15K_V2	H@1 44.72	93.89	94.15	96.56	97.53	86.72	85.64	0.04	<b>99.37</b>
		H@10 71.21	98.93	97.73	99.46	<b>99.55</b>	98.06	92.70	0.3	99.37
		MR 79.47	3.46	11.50	1.32	1.39	2.28	9.86	4900.23	<b>1*</b>
		MRR 0.53	0.95	0.97	0.98	0.98	0.91	0.88	0.002	<b>0.99</b>

Table 8 (continued)

Dataset	Relation-based methods				Attribute-based methods								
	H@1	H@10	MR	MRR	MTransE	MTransE+RotatE	RDGCN	RREA(basic)	RREA(semi)	KDCoE	MultiKE	AttrE	BERT_INT
New datasets	BBC-DB	24.98	18.41	6.38	43.62	46.87	32.69	15.22	17.31	<b>92.50</b>			
		50.28	41.66	11.01	65.28	65.19	52.88	33.47	46.65	<b>93.78</b>			
		672.17	1343.30	2656.5	318.03	631.69	643.35	1299.81	760.98	<b>1.38*</b>			
		0.33	0.26	0.08	0.52	0.54	0.38	0.21	0.27	<b>0.93</b>			
imdb-tmdb	H@1	15.00	9.30	0.07	27.06	29.56	20.02	15.28	87.01	<b>99.70</b>			
	H@10	54.04	37.82	0.73	74.86	73.65	29.43	54.60	98.06	<b>100</b>			
	MR	52.99	131.41	677.50	16.25	61.83	199.59	59.72	5.96	<b>1.002</b>			
	MRR	0.27	0.18	0.005	0.42	0.44	0.23	0.27	0.90	<b>0.99</b>			
imdb-tvdb	H@1	13.63	7.18	0.13	20.53	19.41	15.80	8.61	84.54	<b>99.86</b>			
	H@10	43.82	24.16	1.32	52.76	46.82	24.11	33.21	98.99	<b>100</b>			
	MR	40.41	149.27	377.50	25.86	78.26	149.64	94.45	1.89	<b>1.001</b>			
	MRR	0.23	0.13	0.009	0.32	0.29	0.18	0.16	0.89	<b>0.99</b>			
tmdb-tvdb	H@1	10.63	5.60	0.14	14.13	14.51	20.00	12.26	34.41	<b>100</b>			
	H@10	49.68	29.73	1.45	58.31	54.27	27.09	52.23	65.42	<b>100</b>			
	MR	22.89	74.29	344.50	14.59	31.79	133.18	23.94	28.11	<b>1</b>			



Table 8 (continued)

Dataset	Relation-based methods			Attribute-based methods						
	MTransE	MTransE+RotatE	RDGCN	RREA(basic)	RREA(semi)	KDCoE	MultiKE	AtrE	BERT_INT	
Restaurants	MRR	0.22	0.13	0.01	0.27	0.26	0.22	0.24	0.44	<b>1</b>
	H@1	N/A	N/A	N/A	N/A	N/A	N/A	N/A	<b>94.6</b>	N/A
	H@10	N/A	N/A	N/A	N/A	N/A	N/A	N/A	<b>98.19</b>	N/A
	MR	N/A	N/A	N/A	N/A	N/A	N/A	N/A	<b>1.68</b>	N/A
	MRR	N/A	N/A	N/A	N/A	N/A	N/A	N/A	<b>0.95</b>	N/A

Best scores are given in bold

\* For memory reasons, the source code of BERT\_INT restricts the evaluation computations to the top-k ranks, which may yield better results in some cases (default k=50). We set k=1000, which is the largest k value for which we did not run into memory issues.

space. In case that the embedding space has limited dimensions and the KGs are dense, dissimilar entities end up close in the embedding space. The lack of negative sampling (that distances dissimilar entities in the embedding space) is the reason why MTransE is not improved in dense datasets.

In the new datasets (BBC-DB, imdb-tmdb, imdb-tvdb and tmdb-tvdb) that are all sparse, fewer entity pairs in seed alignment and fewer relation edges are used for training, resulting in worse performance, too, for this method (when comparing BBC-DB to tmdb-tvdb). For example, Hits@1 for BBC-DB with 61,039 relation edges is 24.98, while Hits@1 for imdb-tmdb with 9,604 relation edges is 15. On the contrary, MR is improved as can be seen in Table 8. This is attributed to the dependency of MR on the size of the similarity lists (see the evaluation protocol in Sect. 3.3.1). Specifically, as the entity pairs of seed alignment are reduced, the size of similarity lists is reduced too. The smaller the similarity lists, the lower the ranks for true entity pairs in the test set, resulting to an improved MR metric.

**MTransE+RotatE** is a variation of MTransE that exploits RotatE as embedding module, as well as negative sampling, and was originally proposed for better handling symmetrical relations in KGs (Sun et al. 2020). Since such relations are not encountered in datasets imdb\_tmdb, imdb\_tvdb, and tmdb\_tvdb (see Table 3), we investigate only the effect of negative sampling on the performance of the method. MTransE+RotatE exhibits a similar performance to MTransE in the sparse versions of D\_W\_15K\_V1 and D\_Y\_15K\_V1 (e.g., in dataset D\_W\_15K\_V1, MTransE Hits@10 is 72.78 and MTransE+RotatE Hits@10 is 72.94). However, the negative sampling used by MTransE+RotatE significantly improves its performance in dense datasets compared to MTransE. For example, for D\_Y\_15K\_V2, the MR of MTransE is 79.47 and that of MTransE+RotatE is 3.46, and for D\_W\_15K\_V2, MTransE Hits@1 is 26.27 and MTransE+RotatE Hits@1 is 49.22. This means that MTransE+RotatE represents similar entities closer in the embedding space and dissimilar entities farther, exploiting the rich semantic information that dense KGs offer.

In the new datasets, MTransE+RotatE is outperformed by MTransE (e.g., for tmdb-tvdb Hits@1 of MTransE is 10.63, while Hits@1 of MTransE+RotatE drops to 5.60). This is due to the small size of the new datasets in terms of seed entity pairs (e.g., 982 in tmdb-tvdb) and relation edges (e.g., 3,396 for  $KG_1$  vs only 938 for  $KG_2$  in tmdb-tvdb) that, in conjunction with the high dimensionality of the embedding space, place similar entities far away. Negative sampling seems to amplify this behavior.

**RDGCN** exploits multi-hop entity neighborhoods, as well as a higher number of negative samples and dimensions in the embedding space than the other methods (see Table 7). As shown in Table 8, RDGCN has the highest Hits@1 and Hits@10, compared to MTransE and MTransE+RotatE in datasets D\_W\_15K and D\_Y\_15K, both for sparse (V1) and dense (V2) versions (e.g., for D\_W\_15K\_V2 and D\_Y\_15K\_V1, RDGCN Hits@1 is respectively 63.71 and 94.10, while MTransE Hits@1 is 26.27 and 45.54 and MTransE+RotatE Hits@1 is 49.22 and 46.25). This behavior can be attributed to the fact that RDGCN exploits an extended neighborhood of entities (multi-hop) for constructing the embeddings. Of course, this richer semantic context incurs the risk of incorporating noisy information into

the embeddings, due to potentially irrelevant or less important distant neighbors for aligning entities. For these reasons, RDGCN requires both high-dimensional embeddings to capture more features, as well as increased negative samples per positive edge that place irrelevant neighbors far away in the embedding space. It is worth mentioning that RDGCN outperforms the other methods in all OpenEA datasets, demonstrating the added value of multi-hop neighbors in embeddings. On the other hand, RDGCN is outperformed by MTransE+RotatE in datasets BBC-DB, imdb-tmdb, imdb-tvdb and tmdb-tvdb (e.g., in BBC-DB Hits@1 drops from 18.41 to 6.38). This behavior stems from the same problem that negative sampling causes to MTransE+RotatE in these datasets. As RDGCN uses more negative samples and higher dimensions in the embedding space than MTransE+RotatE, it aggravates the problem.

**RREA(basic)**, like RDGCN, exploits multi-hop entity neighborhoods, having the same issue about irrelevant or less important neighbors. Contrary to RDGCN, that needs higher number of negative samples and dimensions to solve the aforementioned issue, RREA(basic) utilizes an attention mechanism that assigns high weights to relevant and important neighbors, outperforming RDGCN, but also MTransE and MTransE+RotatE in all datasets, regardless of the metric (e.g., in D\_W\_15K\_V2, RREA(basic)'s Hits@1 is 87.86 and for MTransE it is 26.27, in D\_W\_15K\_V1, RREA(basic)'s Hits@10 is 88.49, while for MTransE+RotatE it is 56.98, and in D\_Y\_15K\_V2, RREA(basic)'s MR is 1.32, while the corresponding MR for RDGCN is 11.50). As we can see from D\_Y\_15K\_V1, RREA(basic) fails to differentiate the importance of relevant and irrelevant neighbors, and RDGCN achieves higher performance in all metrics (e.g., Hits@1 for RDGCN is 94.10, while for RREA(basic) it is 76.19). The fact that RREA(basic) solves the issue of noisy information and irrelevant neighbors by utilizing the attention mechanism, without requiring higher number of negative samples, leads to outperforming MTransE, MTransE+RotatE and RDGCN in all new datasets (e.g., RREA(basic)'s Hits@1 for imdb\_tmdb is 27.06, while for MTransE it is 15, RREA(basic)'s Hits@10 for bbc\_db is 65.28, while for MTransE+RotatE it is 41.66 and RREA(basic)'s MR for imdb\_tvdb is 25.86, while for RDGCN it is 377.50).

**RREA(semi)** runs in 5 iterations, where the first iteration is RREA(basic), which means that the 4 extra iterations, increase the performance of RREA(semi) compared to RREA(basic) and consequently to MTransE, MTransE+RotatE and RDGCN in all datasets, in Hits@1, Hits@10 and MRR (e.g., RREA(semi)'s Hits@1 for D\_Y\_15K\_V2 is 97.53, while for RREA(basic) it is 96.56, RREA(semi)'s Hits@10 for imdb\_tmdb is 73.65, while for RDGCN it is 0.73, and RREA(semi)'s MRR for imdb\_tvdb is 0.29, while for MTransE it is 0.23). RREA(semi), like RREA(basic) fails to distinguish relevant and irrelevant neighbors in D\_Y\_15K\_V1, where RDGCN achieves better performance in all metrics (e.g., Hits@1 for RDGCN is 94.10, while for RREA(semi) it is 81.79). The fact that RREA(semi) enriches the seed alignment with new entity pairs leads to increasing the candidate-space (denominator of MR) and also the probability of wrong alignment (compared to RREA(basic)), making RREA(semi) exhibit worse MR in all datasets (e.g., in D\_Y\_15K\_V2, RREA(semi)'s MR is 1.39, while for RREA(basic) it is 1.32).

#### 4.1.2 Attribute-based EA methods

**KDCoE** in addition to one-hop neighborhood and negative sampling, exploits the textual descriptions of entities, when available. Negative sampling also improves the performance of KDCoE in dense versions of the datasets, compared to their sparse versions (e.g., in dataset D\_W\_15K\_V1 Hits@10 is 46.29, while in D\_W\_15K\_V2 Hits@10 is 64.38, and in dataset D\_Y\_15K\_V1 MR is 143.78, while in D\_Y\_15K\_V2 MR is 2.28). In dense entity graphs, negative sampling helps KDCoE to increase the distance of dissimilar entities in an embedding space with limited dimensions. The use of textual descriptions as an additional source of similarity evidence favors KDCoE to outperform the other methods in datasets that are rich in terms of textual descriptions. As shown in Table 8, this is the case of BBC\_DB, which features the highest number of entities with textual descriptions (2,430), compared to the size of the seed alignment (9,396), and where KDCoE exhibits better performance compared to MTransE, MTransE+RotatE and RDGCN. The higher ratio of entities with similar textual descriptions in a dataset boosts the performance of KDCoE. It is worth mentioning that although KDCoE outperforms the aforementioned methods in BBC\_DB, both versions of RREA exhibit higher performance in all metrics, due to the attention mechanism that both utilize and the rule that RREA(semi) uses. In case of lower ratio of entities with similar textual descriptions in a dataset, KDCoE's performance drops as the training set is flooded with falsely similar entity pairs in terms of textual descriptions. This behavior is observed in datasets imdb\_tmdb, imdb\_tvdb, tmdb\_tvdb, where textual descriptions of entities exhibit the minimum similarity among all datasets (e.g., for imdb-tmdb descriptions similarity is 0.35). Statistically meaningful correlations of the method with the specific data characteristics (#Ents\_Descr and Descr\_Sim) are reported in Sect. 4.3. Note that the disfavored effect of negative sampling in these datasets still persists (as in MTransE+RotatE and RDGCN), but with an attenuated strength, due the smaller number of negative samples this method uses.

**MultiKE** exploits multiple sources of similarity evidence of entities (i.e., from attribute names, literals and relations) along with a high number of negative samples. In this respect, MultiKE embeddings are less sensitive to the number of entity relations than in the previous methods. For this reason, MultiKE, in contrary to KDCoE, is not improving in all dense versions of the datasets, even if it uses negative sampling (see Table 8). Specifically, in D\_W\_15K, when comparing the dense to the sparse version, Hits@1 increases from 40.49 to 49.54 as both the number (average attributes per entity is 3 for sparse and 5.58 for dense) and the similarity score of literals (i.e., 81.24 for sparse and 83.79 for dense) increase. In contrary, from the sparse D\_Y\_15K\_V1 to the dense D\_Y\_15K\_V2, there is a small decrease ( $\approx 8\%$ ) in both entity names similarity (from 0.99 to 0.91), causing Hits@1 to drop from 90.18 to 85.64. Finally, in datasets imdb\_tmdb, imdb\_tvdb and tmdb\_tvdb, MultiKE has a lower Hits@1 than KDCoE (i.e., for imdb\_tmdb Hits@1 of KDCoE is 20.02, while for MultiKE it is 15.28), due to the increased number of negative samples. The correlation of MultiKE with the similarity of entity names will be further investigated in Sect. 4.3.3.

**AttrE** is an unsupervised method that consumes fewer negative samples (like KDCoE) and exploits the similarity of literal values. On the other hand, it requires the alignment of predicates in KGs as a pre-processing step. The highest the number of aligned predicates, the better the performance of the method. In datasets without aligned predicates, AttrE is unable to run. Specifically, AttrE is not applicable to datasets D\_W\_15K\_V1 and D\_W\_15K\_V2, while on D\_Y\_15K\_V1 and D\_Y\_15K\_V2, it exhibits a very low performance (e.g., Hits@1 is 0.0003 for D\_Y\_15K\_V1) due to the low predicate similarity (e.g., 0 in D\_W\_15K\_V1 and 47.37% in D\_Y\_15K\_V1). In datasets like BBC\_DB, where both the similarity of predicates (80.93) and the similarity of literals (0.87) are high, AttrE's performance is significantly improved. In datasets imdb-tmdb, imdb-tvdb and tmdb-tvdb, exhibiting the highest similarity of predicates and literals, AttrE outperforms all other methods. The correlation of the method with specific data characteristics (Pred\_Sim and Lit\_Sim) will be presented in Sect. 4.3. Finally, AttrE is the only method that is able to run on the Restaurants datasets, where not all the entities of  $KG_1$  are aligned with entities of  $KG_2$  and vice versa.

**BERT\_INT** is an attribute-based method that exploits multiple sources of similarity evidence between entities (i.e., long-text descriptions, literals, names, and relations). More precisely, BERT\_INT, rather than aggregating the information of multi-hop neighbors that in some cases introduces noise (as in RDGCN), even by the usage of attention mechanisms (as in RREA), it employs interactions between the different views (i.e., description/name-view, neighbor-view and attribute-views), considering multi-hops and consuming very few negative samples compared to RDGCN. Regarding the factual information, it prioritizes textual descriptions over entity names, in order to leverage as much textual information as possible (textual descriptions are typically longer than entity names). BERT\_INT has the advantage, compared to MultiKE, that it also exploits multiple resources but without prioritizing them, since the noise of irrelevant factual information is not incorporated in the embeddings. BERT\_INT also has an advantage over AttrE and KDCoE, that they exclusively exploit either longer textual descriptions or other literals.

Table 8 shows that BERT\_INT outperforms all methods in all datasets and measures, except the two versions of D\_W\_15K. For example, in D\_Y\_15K\_V1, BERT\_INT achieves 99.23 H@1 when RDGCN and RREA(semi) achieve 94.10 and 81.79. In D\_Y\_15K\_V2, we observe similar behavior with BERT\_INT exhibiting 99.37 when RDGCN and RREA(semi) achieve 94.15 and 97.53. This happens, because in both versions of D\_Y\_15K the number of the entities that have textual descriptions, the descriptions similarity and the entity names are high (e.g., in D\_Y\_15K\_V1 3,455 entities have textual descriptions with 0.88 similarity and 0.99 similarity in names). The opposite behavior is observed in both versions of D\_W\_15K, where RREA outperforms BERT\_INT with H@10 90.03 instead of 48.96 that BERT\_INT achieves, since in these datasets the number of the entities that have descriptions is low (only 278 in D\_W\_15K\_V2), with low similarity (0.73) and low similarity in entity names (0.78). In new datasets, the similarity of long-text descriptions and entity names is very low, the literal similarity is very high (0.94 in tmdb-tvdb); as a result, the performance of BERT\_INT is even better, outperforming even AttrE (e.g., in tmdb-tvdb, BERT\_INT finds all the matches correctly, with an impressive

MRR 1, while the corresponding MRR of the second best method, AttrE, is only 0.44).

### 4.1.3 Conventional EA methods

PARIS (Suchanek et al. 2011) is a conventional, probabilistic, holistic approach, i.e., it performs both instance and schema matching, by estimating probabilities of matching (i.e., equivalence), without learning KG embeddings. Such probability estimation relies on the knowledge of quasi-functional relations, i.e., relations that for a given head entity, the expected number of tail entities is close to 1. Given, for example, two relation triples  $\langle h, r, t \rangle$  and  $\langle h', r, t' \rangle$ , with  $r$  being a functional relation (e.g., `born_in`), the equivalence probability  $h \equiv h'$  depends recursively on other equivalence probabilities ( $t \equiv t'$ ). Initial equivalences are computed between literals, using a certain string distance (their similarity is inverse proportional to their edit distance).

Table 9 depicts the performance of PARIS along with the best embedding-based method for each dataset according to Table 8. As in Sun et al. (2020), for PARIS, we report the classification-based metrics (Precision, Recall and F1-score), while for the embedding-based method, we used the H@1. This is because, in the test phase, each source entity gets a list of candidates; as a result, precision, recall and F1-score can be considered equal to Hits@1. Regarding the OpenEA datasets, we observe that in both versions of D\_W\_15K, where the literal similarity is high (for D\_W\_15K\_V1, it is 0.81), PARIS outperforms the best embedding-based method (in this case RREA(semi)). On the other hand, in both versions of D\_Y\_15K, the best method (BERT\_INT) exhibits higher Recall and F1-score (in D\_Y\_15K\_V2 Recall is 99.37, instead of 98.37 of PARIS and F1-score is 99.37, instead of 99.10 of PARIS), while the precision is lower (99.37) but still close to PARIS (99.85). This happens because of the low literal similarity that these datasets exhibit (affecting the Precision) and the high textual descriptions and entity names similarity (improving Recall and F1-score). Finally, in the new datasets, even if the literal similarity is high, the lack of functional relations affect negatively the performance of PARIS, and as a result, BERT\_INT is the clear winner in these datasets (e.g., in `imdb-tmdb` F1-score is 15.53 for PARIS, while BERT\_INT achieves 99.70).

## 4.2 Effectiveness vs efficiency trade-offs

In this section, we investigate whether it is worth paying the time overhead of each method, with respect to its effectiveness (Q3). In this respect, in Sect. 4.2.1, we first rank the methods according to their effectiveness across all datasets of our testbed. For that, we rely on a statistically sound methodology for ranking the different EA methods across all datasets of our testbed, that exploits the non-parametric Friedman test (Demsar 2006) and the post-hoc Nemenyi test (Nemenyi 1963). Then, in Sect. 4.2.2, we report the execution time of each method (Table 14) and analyze their trade-offs in terms of effectiveness and efficiency (Fig. 10), demonstrating also

**Table 9** Performance of PARIS and of the best embedding-based method per dataset

Dataset			PARIS	Best Method
Datasets	D_W_15K_V1	Precision	<b>96.51</b>	71.82
		Recall	<b>74.55</b>	71.82
		F1-score	<b>84.12</b>	71.82
	D_W_15K_V2	Precision	<b>97.59</b>	93.72
		Recall	90.32	<b>93.72</b>
		F1-score	<b>93.81</b>	93.72
OpenEA	D_Y_15K_V1	Precision	<b>99.72</b>	99.23
		Recall	96.24	<b>99.23</b>
		F1-score	97.95	<b>99.23</b>
	D_Y_15K_V2	Precision	<b>99.85</b>	99.37
		Recall	98.37	<b>99.37</b>
		F1-score	99.10	<b>99.37</b>
New datasets	BBC-DB	Precision	78.43	<b>92.50</b>
		Recall	26.38	<b>92.50</b>
		F1-score	38.72	<b>92.50</b>
	imdb-tmdb	Precision	56.84	<b>99.70</b>
		Recall	8.99	<b>99.70</b>
		F1-score	15.53	<b>99.70</b>
	imdb-tvdb	Precision	58.37	<b>99.86</b>
		Recall	9.17	<b>99.86</b>
		F1-score	15.84	<b>99.86</b>
	tmdb-tvdb	Precision	99.76	<b>100</b>
		Recall	97.50	<b>100</b>
		F1-score	98.62	<b>100</b>
	Restaurants	Precision	28.79	<b>94.60</b>
		Recall	69.76	<b>94.60</b>
		F1-score	40.18	<b>94.60</b>

Best scores are given in bold

the training curves of the methods and the time that a method needs to reach 90% of its highest MRR (Fig. 11).

### 4.2.1 Effectiveness-based ranking of EA methods

In this section, for each of the metrics used in Table 8, we first calculate the ranks (rank 1 corresponds to the best performance) of all methods per dataset and average their ranks for each metric. The resulting ranks are reported in Tables 10, 11, 12, and 13.

Then, to infer a statistically significant ranking of the methods, we rely on the non-parametric Friedman test (Demsar 2006). The null hypothesis  $H_0$  of the Friedman test is that “The mean performance for each method is equal”, while the alternative hypothesis ( $H_a$ ) states exactly the opposite. With p-values 0.004, 0.003, 0.001

and 0.003 of Friedman test for Hits@1, Hits@10, MR and MRR, respectively, we can reject the null hypothesis  $H_0$  at a 10% confidence-level ( $\alpha$ ). In the sequel, we conduct the Nemenyi post-hoc test to compare the methods pairwise. This test reports as significance the average ranks of two methods if they differ by a critical distance (CD) given by  $q_\alpha \sqrt{\frac{k(k+1)}{6N}}$ , where  $N$  is the number of the datasets,  $q_\alpha$  is a constant based on  $\alpha$ , and  $k$  is the number of methods in total. For 9 EA methods, 8 datasets and  $\alpha = 0.1$ , the value of  $CD$  is 3.90.

From Fig. 9, we can see that BERT\_INT is a clear winner, since it outperforms the other methods with statistically significant difference (bigger difference than  $CD$ ). Figure 9a, which captures accuracy with 0 error tolerance (Hits@1), shows that BERT\_INT has statistically significant performance difference to MTransE, MTransE+RotatE, AttrE and RDGCN, while RREA(semi) only to MTransE and MTransE+RotatE. In case we accept a 10% error tolerance (Hits@10 in Fig. 9b), none of the examined methods has significant performance difference to the other methods, so the choice should depend on different factors, e.g., the training time (see Sect. 4.2.2) or dataset characteristics (see Sect. 4.3). In addition, BERT\_INT that exploits multiple sources of similarity evidence of entities (e.g., textual descriptions, entity names) and not only structural information, as relation-based methods such as RREA(basic) do, improves the general performance, improving the ranks even of lower ranked entities in the similarities lists (see Sect. 3.3.1). The result is that BERT\_INT has significant performance difference with MTransE, RDGCN, MTransE+RotatE, KDCoE and AttrE, while RREA(basic) has significant performance difference with only RDGCN, MTransE+RotatE, KDCoE and AttrE, for MR. In contrast, MRR (Fig. 9d) which is affected by the first ranks (see Sect. 3.3.2) and is less sensitive to outliers, shows that BERT\_INT is the only method that has statistically significant performance difference with other methods (MTransE and MTransE+RotatE).

#### 4.2.2 Efficiency-based ranking of EA methods

In this subsection, we compare the efficiency of the evaluated EA methods. Particularly, in Table 14, we are measuring the execution time of both training and testing. Figure 10 depicts the trade-off between effectiveness and efficiency of EA methods using BERT\_INT as baseline. In addition, Figure 11 demonstrates the training curves of the methods, as well as the time or the epochs that a method needs to achieve 90% of its own highest MRR.

As MTransE has the most lightweight implementation among all methods (Sun et al. 2020) (only two layers of entity and relation embeddings), it is the fastest method (Table 14). Although MTransE+RotatE also uses only two layers, MTransE is faster because it employs a simpler scoring function than RotatE (see Sect. 2.2.1.1). BERT\_INT, that achieves the best performance in most datasets (excluding the two versions of D\_W\_15K), and MultiKE are competing for the second most efficient method. Specifically, as both BERT\_INT and MultiKE trains and combines multi-view embeddings, they require more training time (long-lasting



**Table 10** Ranking of methods per dataset according to Hits@1

Dataset	MTransE	MTransE+RotatE	RDGCN	RREA(basic)	RREA(semi)	KDCoE	MultiKE	AttrE	BERT_INT
D_W_15K_V1	7	6	3	2	1	8	5	9	4
D_W_15K_V2	8	5	3	2	1	7	4	9	6
D_Y_15K_V1	8	7	2	5	4	6	3	9	1
D_Y_15K_V2	8	5	4	3	2	6	7	9	1
BBC-DB	5	6	9	3	2	4	8	7	1
imdb-tmdb	7	8	9	4	3	5	6	2	1
imdb-tvdb	6	8	9	3	4	5	7	2	1
tmdb-tvdb	7	8	9	5	4	3	56	2	1
Avg rank	7	6.62	6	3.37	2.62	5.5	5.75	6.12	2

**Table 11** Ranking of methods per dataset according to Hits@10

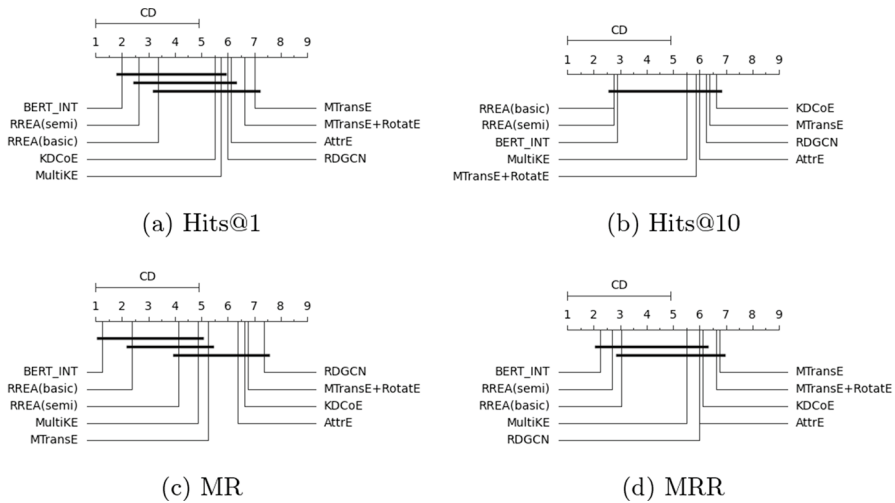
Dataset	MTransE	MTransE+RotatE	RDGCN	RREA(basic)	RREA(semi)	KDCoE	MultiKE	AttrE	BERT_INT
D_W_15K_V1	6	4	3	2	1	8	5	9	7
D_W_15K_V2	7	4	3	2	1	6	5	9	8
D_Y_15K_V1	8	7	2	5	4	6	3	9	1
D_Y_15K_V2	8	4	6	2	1	5	7	9	3
BBC-DB	5	7	9	2	3	4	8	6	1
imdb-tmdb	6	7	9	3	4	8	5	2	1
imdb-tvdb	5	7	9	3	4	8	6	2	1
tmdb-tvdb	6	7	9	3	4	8	5	2	1
Avg rank	6.37	5.87	6.25	2.75	2.75	6.62	5.5	6	2.87

**Table 12** Ranking of methods per dataset according to MR

Dataset	MTransE	MTransE+RotatE	RDGCN	RREA(basic)	RREA(semi)	KDCoE	MultiKE	AttrE	BERT_INT
D_W_15K_V1	5	7	6	2	3	8	4	9	1
D_W_15K_V2	6	5	8	1	2	8	4	9	3
D_Y_15K_V1	7	8	2	4	5	6	3	9	1
D_Y_15K_V2	8	5	7	2	3	4	6	9	1
BBC-DB	5	8	9	2	3	4	7	6	1
imdb-tmdb	4	7	9	3	6	8	5	2	1
imdb-tvdb	4	7	9	3	5	8	6	2	1
tmdb-tvdb	3	7	9	2	6	8	4	5	1
Avg rank	5.25	6.75	7.37	2.37	4.12	6.62	4.87	6.37	1.25

**Table 13** Ranking of methods per dataset according to MRR

Dataset	MTransE	MTransE+RotatE	RDGCN	RREA(basic)	RREA(semi)	KDCoE	MultiKE	AttrE	BERT_INT
D_W_15K_V1	7	6	3	2	1	8	4	9	5
D_W_15K_V2	8	4	3	2	1	6	5	9	7
D_Y_15K_V1	8	7	2	5	4	6	3	9	1
D_Y_15K_V2	8	5	4	2.5	2.5	6	7	9	1
BBC-DB	5	7	9	3	2	4	8	6	1
imdb-tmdb	6	8	9	4	3	7	5	2	1
imdb-tvdb	5	8	9	3	4	6	7	2	1
tmdb-tvdb	7	8	9	3	4	6	5	2	1
Avg rank	6.75	6.62	6	3.06	2.68	6.12	5.50	6	2.25



**Fig. 9** Statistical significance rankings of methods according to all metrics

epochs) than MTransE, despite the fact that the latter needs more epochs according to Table 15 (short-lasting epochs). On the other hand, BERT\_INT and MultiKE are faster than MTransE+RotatE and KDCoE, since BERT\_INT and MultiKE converge faster (see Table 15). In addition, BERT\_INT and MultiKE are faster than RREA(basic) and RREA(semi), since the last two do not utilize early stopping and they need multiple epochs to be trained (1,200 and  $5 \times 1,200 = 6,000$ , respectively). Finally, BERT\_INT and MultiKE are faster than the unsupervised method AttrE, even if the latter needs fewer epochs for training (it does not utilize early stopping). This is due to the nature of AttrE refining embeddings in multiple iterations, by exploiting all the possible literal values and not the seed alignment.

Figure 10 illustrates the trade-off between effectiveness and efficiency of all methods, compared to a baseline method (BERT\_INT). The x-axis represents the relative difference of the methods in effectiveness, measured by MRR, the metric that is the least sensitive to outliers (see Sect. 3.3.2), and the y-axis represents the relative difference in execution time, with respect to the baseline. In particular, the red quadrant (bottom right) indicates methods that both run faster and have a better effectiveness than the baseline (i.e., dominate the baseline) in a specific dataset, and the green quadrant (top left) exactly the opposite (i.e., methods dominated by the baseline). The bottom-left grey quadrant indicates methods that run faster with lower performance than the baseline and the top-right grey quadrant indicates methods that have better performance, but are slower than the baseline method.

The choice of BERT\_INT as baseline is because it is the method that dominates most methods (highest number of methods on the top-left quadrant), and it is not dominated by any other method (0 methods in the bottom-right quadrant). As shown in Fig. 9d, BERT\_INT is the method with the best MRR in both OpenEA and new datasets. Indicatively, the corresponding values for the fastest method, MTransE,

**Table 14** Training + testing time (execution time in seconds) of methods in each dataset

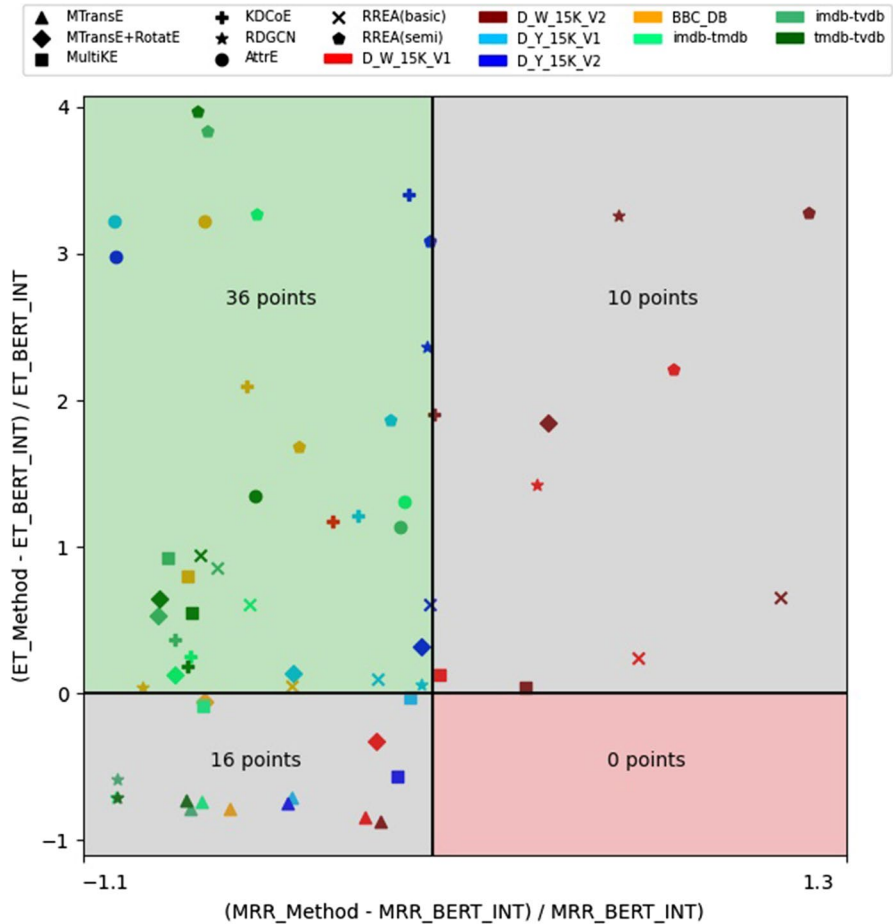
Dataset	MTransE	MTransE+RotatE	RDGCN	RREA(basic)	RREA(semi)	KDCoE	MultiKE	AttrE	BERT_INT
D_W_15K_V1	<b>53.54 + 3.67</b>	236.76 + 6.85	853.80 + 26	442.68 + 8.69	1,158.73 + 8.22	784.75 + 5.67	403.53 + 4.38	-	333.34 + 30.09
D_W_15K_V2	<b>44.41 + 3.62</b>	1,028.26 + 6.94	1,520.21 + 26.86	593.32 + 8.33	1,543.96 + 9	1,052.56 + 4.50	372.29 + 5.37	-	333.61 + 29.80
D_Y_15K_V1	<b>101.76 + 3.63</b>	405.33 + 6.92	356.50 + 27.07	387.43 + 8.12	1,027.93 + 7.56	795.42 + 5.34	348.55 + 4.40	1,472.26 + 54	331.73 + 29.95
D_Y_15K_V2	<b>88.85 + 3.62</b>	470.96 + 6.78	1,185.88 + 27.80	570.38 + 9	1,468 + 8.25	1,585.22 + 4.71	154.12 + 4.52	1,416.78 + 20.59	331.40 + 29.75
BBC-DB	<b>44.88 + 2.21</b>	209.26 + 2.84	221.97 + 10.38	230.33 + 5.27	595.32 + 5.56	690.21 + 2.95	400.93 + 2.05	938.05 + 7.48	206.69 + 17.13
imdb-tmdb	<b>11.94 + 0.30</b>	52.32 + 0.30	12.09 + 1.44	71.41 + 3.74	195.12 + 3.90	57.23 + 1.12	41.90 + 0.64	107.47 + 0.32	42.62 + 4.04
imdb-tvdb	<b>5.16 + 0.24</b>	37.93 + 0.19	9.31 + 1.03	42.77 + 3.48	117.37 + 3.48	32.90 + 1.22	47.57 + 0.55	53.16 + 0.08	23.80 + 1.16
tmdb-tvdb	<b>6.03 + 0.24</b>	37.62 + 0.17	5.61 + 1.02	41.07 + 3.53	110.67 + 3.53	25.93 + 1.22	34.96 + 0.54	53.79 + 0.08	22.06 + 0.91

Best scores are given in bold

The dash '-' means that the method could not run in this dataset

**Table 15** Number of epochs employed for training

Dataset	MTransE rel_emb	MTransE+RotatE rel_emb	RDGCN rel_emb	RREA(basic) rel_emb	RREA(semi) rel_emb	KDCoE rel_emb + desc_emb	MultiKE multi-view_ emb	AtrE rel_emb	BERT_INT BERT + inter- actions
D_W_15K_V1	224	322	220	1200 (fix)	6000 (fix)	256/iter. + 568/ iter.	114	–	5 + 200 (fix)
D_W_15K_V2	144	300	224	1200 (fix)	6000 (fix)	304/iter. + 640/ iter.	154	–	5 + 200 (fix)
D_Y_15K_V1	472	396	70	1200 (fix)	6000 (fix)	204/iter. + 490/ iter.	90	50 (fix)	5 + 200 (fix)
D_Y_15K_V2	318	332	140	1200 (fix)	6000 (fix)	268/iter. + 568/ iter.	86	50 (fix)	5 + 200 (fix)
BBC-DB	236	234	126	1200 (lfix)	6000 (fix)	230/iter. + 344/ iter.	116	50 (fix)	5 + 200 (fix)
imdb-tmdb	196	178	50	1200 (fix)	6000 (fix)	220/iter. + 184/ iter.	74	50 (fix)	5 + 200 (fix)
imdb-tvdb	130	224	50	1200 (fix)	6000 (fix)	122/iter. + 144/ iter.	56	50 (fix)	5 + 200 (fix)
tmdb-tvdb	120	228	50	1200 (fix)	6000 (fix)	74/iter. + 42/iter.	72	50 (fix)	5 + 200 (fix)

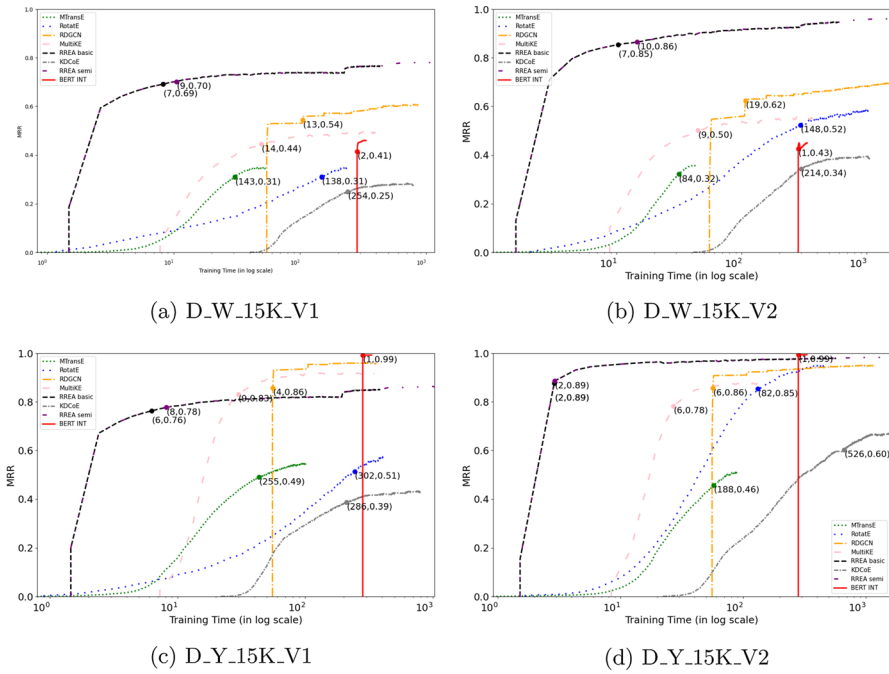


**Fig. 10** Trade-off between efficiency and effectiveness (MRR) of methods across all datasets, using BERT\_INT as baseline. The x-axis is log-scaled using symmetric log function

when chosen as a baseline, are: 16 (top left), 46 (top right), 0 (bottom left), and 0 (bottom right).

As BERT\_INT outperforms all methods in all datasets (except of the two versions of D\_W\_15K) in terms of MRR, with low a execution time, the top-left quadrant consists of 36 out of 62 (58.06%) points. The bottom-left grey quadrant contains 16 out of 62 (25.80%) points, as BERT\_INT is slower than the fastest method MTransE in all datasets, as well as slower than the competing method MultiKE in some datasets. More precisely, MTransE is on average 394.14% faster than BERT\_INT, achieving 138.73% lower MRR, while MultiKE is 43.79% faster than BERT\_INT, at a cost of 10.61% lower MRR in the two versions of D\_Y\_15K. Finally, since BERT\_INT is negatively affected by the low attribute (e.g., textual descriptions and entity names) similarity in dense KGs, in terms of effectiveness, the top-right





**Fig. 11** Methods training time curve for each dataset. Each of the seven points represents the time (measured both in seconds and epochs) that a specific method needs to achieve 90% of its highest MRR. The x-axis is log-scaled

quadrant contains 10 out of 62 (16.12%) points. In KGs with such characteristics, RREA(basic) and RDGCN outperform BERT\_INT, achieving 88.63% and 45.45% better effectiveness with an overhead of 226.62% and 42.53% longer execution time, respectively.

In Fig. 11, we demonstrate the training curves of the methods for each of the OpenEA datasets, while we also point out the training time that a method needs, in order to achieve 90% of its highest MRR (represented as a dot in each curve), similar to the *time-to-accuracy* (TTA) metric in the literature (Coleman et al. 2017, 2019), measured here in seconds (x-axis labels) and epochs (first element of the points coordinates). As suggested in Coleman et al. (2017), Coleman et al. (2019), we run the methods in only one fold of the datasets. We exclude the AttrE method, since it cannot run with D\_W\_15K\_V1 and D\_W\_15K\_V2, while for D\_Y\_15K\_V1 and D\_Y\_15K\_V2 it yields MRR close to zero. In addition, it is worth mentioning that the training curve of RREA(basic) (Fig. 11) overlaps with part of the training curve of RREA(semi), since the former is the first iteration (out of 5) of the latter. Last but not least, in Table 16 we demonstrate the coefficient of variation (CV) of methods training epochs per dataset. More precisely, CV measures the relative variability of the methods epochs and is given by  $\sigma/\mu$ , where  $\sigma$  is the standard deviation and  $\mu$  is the mean. The higher the values, the greater the degree of relative variability.

From Fig. 11, we firstly conclude that the increased training time of RREA(semi) is not worth spending, since the improvement of its MRR is very low, compared to RREA(basic), and it reaches 90% of its highest MRR in its first iteration. In addition, we conclude that in most cases, both versions of RREA, not only achieve the highest MRR, but also reach every threshold faster (steeper slope) in only few seconds or epochs, e.g., in D\_W\_15K\_V1, RREA(basic) needs only 9.43s (7 epochs) for 0.69 MRR. RDGCN achieves quite lower MRR in reasonable time e.g., in D\_W\_15K\_V2 it needs only 107.19s (19 epochs) for 0.62 MRR. Although BERT\_INT is the fastest or the second fastest method (see Table 14), as it needs few short-lasting epochs for training, yet, it takes more time to fine-tune the BERT model and launch the training. This explains why BERT\_INT appears on right extreme of the x-axis in Fig. 11, but it still manages to finish before all other methods. It is worth mentioning that it takes only 1 or 2 epochs to achieve an MRR that is even better than RREA or RDGCN. E.g., in D\_Y\_15K\_V1 it needs 281.97s (1 epoch) for 0.99 MRR. Last but not least, as we can see from Table 16, KDCoE, RDGCN and BERT\_INT have CVs greater than 1, which means that the time of each epoch varies a lot during training. In contrary to those methods, MTransE, MTransE+RotatE, RREA(basic), RREA(semi), and MultiKE exhibit low variability in epochs duration.

### 4.3 Meta-level analysis of EA methods

In this subsection, we report the statistically significant correlations between the various meta-features (Table 5) extracted by the KGs of our datasets and the performance achieved by EA methods (Q4). In this experiment, we exclude the Restaurants dataset, since only AttrE was able to run on this dataset.

The null hypothesis ( $H_0$ ) of our correlation analysis is that “there is no significant correlation between the performance metric  $\mathbf{X}$  and the meta-feature  $\mathbf{Y}$ ”, while the alternative hypothesis ( $H_a$ ) states exactly the opposite. We reject the null hypothesis when the p-values of correlations are less than a significance level  $\alpha$ . We use the Spearman’s correlation, which measures the degree of association between two ranked variables, with a coefficient  $r_s = 1 - \frac{6 \sum d_i^2}{n(n^2-1)}$ , where  $d_i$  is the difference between the ranks of the two variables and  $n$  is the maximum rank. This choice is motivated by the robustness of the Spearman’s correlation to outliers (unlike Pearson’s correlation), as it relies on the ranks of the variables and not on their actual values. The closer  $r_s$  is to zero, the weaker the association between two ranked variables (negative or positive). A positive Spearman’s correlation denotes that when a meta-feature increases the method-specific metric increases too, while negative Spearman’s correlation denotes the opposite direction of the association<sup>17</sup> between the two ranked variables.

Table 17 reports the correlations for a significance level  $\alpha = 0.05$  (otherwise cells are empty) under different categories. To make the results easily comparable across the different performance metrics, we multiply correlations of MR with  $-1$ , since in contrary to the other metrics, the higher the value of MR the worse. Finally, for

<sup>17</sup> Correlation is not causation as the association direction is not known.

**Table 16** Coefficient of variation (CV) of methods training epochs

Dataset	MTransE	MTransE+RotatE	RDGCN	RREA(basic)	RREA(semi)	KDCoE	MultiKE	BERT_INT
D_W_15K_V1	0.031	0.009	1.907	0.022	0.018	5.283	0.019	11.94
D_W_15K_V2	0.031	0.010	1.789	0.019	0.016	4.234	0.014	11.93
D_Y_15K_V1	0.035	0.012	1.955	0.026	0.021	6.788	0.013	11.95
D_Y_15K_V2	0.027	0.008	1.849	0.018	0.016	7.739	0.017	11.94

meta-features that impact only specific EA methods, e.g., description similarity concerns only KDCoE, their correlations with non relevant methods are omitted (cells are filled with dashes).

#### 4.3.1 Seed alignment size

As we can see in Table 17, the size of the seed alignment (#Entity\_Pairs) is positively correlated with the supervised and semi-supervised methods MTransE, MTransE+RotatE, RDGCN, RREA(basic), RREA(semi), KDCoE, MultiKE and BERT\_INT and negatively correlated with the unsupervised method AttrE. Increasing the seed alignment and consequently, increasing the space of candidate matches, penalizes unsupervised methods, as it also increases the probability of wrong alignments. On the other hand, the size of the seed alignment is positively correlated with the performance of supervised and semi-supervised methods, as a large number of entity pairs in seed alignment implies larger KGs<sup>18</sup>, and by extension more training data. Although BERT\_INT is a supervised method, in the new datasets (which are smaller, but exhibit higher factual similarity than the OpenEA datasets), it achieves higher performance than in OpenEA datasets, resulting in a negative correlation with the seed alignment size. Note that we only measure correlations, which are not causal relations, as the association direction is not known.

#### 4.3.2 Density

As we can see in Table 17, the average number of relations per entity is positively correlated with the performance of MTransE, MTransE+RotatE, RDGCN, RREA(basic), RREA(semi), KDCoE and MultiKE, while all methods are negatively correlated with the sole meta-feature and positively correlated with the hyper meta-feature (see Sect. 3.2). In addition, the average number of attributes per entity is negatively correlated with the performance of AttrE and BERT\_INT.

As we discussed in Sects. 2.5.1.2 and 2.5.6, a high number of average relations per entity improves the quality of embeddings. Specifically, entities have the opportunity to minimize their embedding distance with multiple similar entities (neighbors). Also, the more relations we have per entity, the more negative samples we will consider, and as a result, dissimilar entities will be placed farther in the embedding space. Sole and hyper meta-features are essentially alternative measures of KG density. When we have many sole relations, many relation types will never co-occur with others resulting to a low average number of relations per entity (low density). Inversely, when we have many hyper relations, there will be many co-occurring relation types resulting to a high average relations per entity (high density).

---

<sup>18</sup> In supervised and semi-supervised methods, all nodes of aligned KGs should be entities of the employed seed alignment.

**Table 17** Correlation matrix of meta-features with EA methods' performance metrics

Methods	Metrics			Seed alignment size		Density		Heterogeneity						
	H@1	H@10	MR	MRR	#Entity_Pairs of seed alignment	Avg_Rels_per_Entity	Avg_Atrs_per_Entity	Sole_Rels	Hyper_Rels	#Ents_Descr	Descr_Sim	Ent_Name_Sim	Literal_Sim	Pred_Name_Sim
MTransE	H@1	0.94	0.79	-	-0.78	0.78	-	-	-	-	-	-	-	-
	H@10	0.89	0.74	-	-	-	-	-	-	-	-	-	-	-
MTransE + Rot.	MR	0.96	0.79	-	-0.78	0.78	-	-	-	-	-	-	-	-
	MRR	0.94	0.90	-	-0.81	0.81	-	-	-	-	-	-	-	-
	H@1	0.94	0.93	-	-0.81	0.81	-	-	-	-	-	-	-	-
	H@10	0.94	0.93	-	-0.81	0.81	-	-	-	-	-	-	-	-
RDGCN	MR	0.91	0.93	-	-0.81	0.81	-	-	-	-	-	-	-	-
	MRR	0.84	0.81	-	-0.83	0.83	-	-	-	-	-	-	-	-
	H@1	0.84	0.81	-	-0.83	0.83	-	-	-	-	-	-	-	-
	H@10	0.84	0.81	-	-0.83	0.83	-	-	-	-	0.88	-	-	-
RREA(basic)	MR	0.84	0.74	-	-0.78	0.78	-	-	-	-	-	-	-	-
	MRR	0.94	0.9	-	-0.81	0.81	-	-	-	-	-	-	-	-
	H@1	0.89	0.86	-	-	-	-	-	-	-	-	-	-	-
	H@10	0.89	0.86	-	-	-	-	-	-	-	-	-	-	-
RREA(semi)	MR	0.94	0.9	-	-0.81	0.81	-	-	-	-	-	-	-	-
	MRR	0.94	0.9	-	-0.83	0.83	-	-	-	-	-	-	-	-
	H@1	0.89	0.86	-	-	-	-	-	-	-	-	-	-	-
	H@10	0.89	0.86	-	-	-	-	-	-	-	-	-	-	-
MTransE	MR	0.94	0.9	-	-0.81	0.81	-	-	-	-	-	-	-	-
	MRR	0.94	0.9	-	-0.81	0.81	-	-	-	-	-	-	-	-

Table 17 (continued)

Methods	Metrics	Seed alignment size	Density	Heterogeneity									
				#Entity_Pairs of seed alignment	Avg_Rels_Entity	Avg_Attrs_per_Entity	Sole_Rels	Hyper_Rels	#Ents_Descr	Descr_Sim	Ent_Name_Sim	Literal_Sim	Pred_Name_Sim
KDCoE	H@1	0.85	0.93	-0.90	0.90	0.81	0.88	-	-	-	-	-	-
	H@10	0.85	0.93	-0.90	0.90	0.81	0.88	-	-	-	-	-	-
	MR												--
MultiKE	MRR	0.85	0.93	-0.90	0.90	0.81	0.88	-	-	-	-	-	-
	H@1	0.89	0.74					0.71					
	H@10	0.81						0.83					
AttrE	MR							0.9					
	MRR	0.81						0.83					
	H@1	-0.84								0.73		0.9	0.83
BERT_INT	H@10	-0.87		-0.73									
	MR												
	MRR	-0.84									0.73		0.9
BERT_INT	H@1	-0.81		-0.9									
	H@10	-0.78		-0.9									
	MR												
MRR			-0.8										

Empty cells denote statistically not significant correlation, and dashed cells denote that a meta-feature is not relevant to an EA method

The meta-feature measuring the average number of attributes per entity is relevant only to attribute-based methods that exploit similarity of literals for the embeddings, such as MultiKE, AttrE and BERT\_INT. As we can see in Table 17, the performance of AttrE and BERT\_INT is negatively correlated with this meta-feature. This happens because as the average number of attributes per entity increases, the noise (irrelevant attribute values) per entity is increasing too, especially if the literals similarity is low. The correlation of AttrE with the literals similarity is examined in Sect. 4.3.3.

### 4.3.3 Heterogeneity

As the only attribute-based semi-supervised method, KDCoE exploits textual descriptions to enrich the seed alignment with new entity pairs that have very similar descriptions. The presence of many entities with highly similar (Descr\_Sim) textual descriptions (#Ents\_Descr) boosts the performance of KDCoE. For this reason, we observe a positive correlation of these two meta-features with KDCoE. In addition, MultiKE and RDGCN are the only supervised methods that exploit entity names to create and initialize the entity embeddings. Thus, we observe a positive correlation between those methods and entity name similarity. AttrE is the only unsupervised method that is not using pre-trained word embeddings and thus is positively correlated with the meta-feature capturing the similarity of literals (Lit\_Sim). A higher literal similarity is associated with better performance for the method. Finally, as AttrE exploits the Levenshtein distance to align the predicate names of two KGs, its performance is positively correlated with this meta-feature (Pred\_Name\_Sim): the more similar the predicates names of the two KGs, the better for AttrE. Although BERT\_INT, exploits multiple resources of similarity evidence of entities (e.g., textual descriptions, entity names, literal), it uses pre-trained BERT-based word embeddings, so measuring correlations with these meta-features, is irrelevant.

## 4.4 Lessons learned

In this section, we summarize the main conclusions drawn from our experiments that shed light on the four open questions of our empirical study.

**Q1.** *What are the critical factors that affect the effectiveness of relation-based and attribute-based methods and how sensitive are the methods to hyperparameters tuning?*

*Negative sampling* and the range of the entity neighborhood are the two critical factors that affect the performance (see Table 8) of the relation-based methods. Specifically, negative sampling helps the methods to exploit the rich semantic information that dense KGs offer, in order to distantiate dissimilar entities in an embedding space with limited dimensions. However, negative sampling harms the effectiveness of methods when the KGs are very sparse (average relations per entity < 3). This is due to the fact that similar entities lie already far in the embedding space and negative sampling further increases their distance: the more

negative samples per positive sample (see Table 2) of KGs, the more the distance increases.

Moreover, by increasing the *range of the neighborhood* that a method exploits, information from both close and distant neighbors is aggregated to learn entity embeddings. However, an increased number of features for learning embeddings, will also increase the dimensionality of the embedding space under the default neural-network architecture that each method use. Additionally, multi-hop relations entails the risk of modeling noisy information from irrelevant distant neighbors. To alleviate the noise in embeddings, we need either a high number of negative samples or attention mechanisms that weight relevant and important neighbors.

*Interaction-based* methods that compare the pairwise similarity of attributes' values of both close and distant neighbors, are less affected by negative sampling on sparse KGs. They can also cope with the noise introduced by the aggregation of distant neighbors, but they require highly similar attributes, independently of the density of the KGs. As a matter of fact, any additional source of similarity evidence, such as textual descriptions (KDCoE, BERT\_INT), literals (MultiKE, AttrE and BERT\_INT), entity names (MultiKE, BERT\_INT) or predicate names (AttrE), boosts the performance of attribute-based methods.

**Q2.** *What is the improvement in the effectiveness of embedding-based entity alignment methods, if we consider not only the structural relations of entities, but also their attribute values?*

The improvement in the effectiveness (see Table 8) of embedding-based EA methods that exploit the factual information of entities depends on the characteristics of the input KGs. In the dense KGs of the OpenEA datasets, where the similarity of factual information (e.g., textual descriptions and entity names) is low, exploiting only the structural neighborhood of entities seems sufficient for a method to achieve a competitive performance. However, when the similarity of factual information is high, the effectiveness of the methods in the dense KGs is significantly improved, approaching the perfect matching. In the sparse KGs of the new datasets containing a high number of textual descriptions and very similar literals and predicates, the attribute-based methods outperform most relation-based methods, especially BERT\_INT that is the overall winning method.

In particular, among all embedding-based methods, both versions of RREA seem to achieve the best performance in dense KGs, when the similarity of factual information (e.g., textual descriptions and entity names) is low, since RREA exploits multi-hop neighborhoods and utilizes attention mechanism that weights relevant and important neighbors. In contrast, when the similarity of factual information is high, BERT\_INT is the best method, since it seems to properly handle it, achieving very high performance. However, in sparse KGs, RREA is outperformed by BERT\_INT that achieves the best performance among all methods in the new datasets, as in these datasets, the literals similarity is higher than in the OpenEA datasets.



The same improvement is also observed in PARIS (see Table 9) that outperforms the best embedding-based methods when textual descriptions and entity names similarity is low in dense KGs. However when the textual descriptions and entity names similarity is high, PARIS is outperformed by the best embedding-based method. Nonetheless, even in those cases, PARIS's precision remains high, because it uses only functional attributes. In sparse KGs, when literal similarity is high, PARIS is outperformed by the best embedding-based method.

**Q3.** *Is the runtime overhead of each method worth paying, with respect to the achieved effectiveness?*

We conclude that the training time of BERT\_INT is a cost worth spending, since the effectiveness improvement is significant, compared to faster methods, like MTransE (Fig. 10). Moreover, from Fig. 11, we conclude that BERT\_INT reaches 90% of its highest MRR faster than the other methods, employing only 1 or 2 epochs. It is worth noticing that in dense KGs with low similarity in factual information (e.g., textual descriptions and entity names), BERT\_INT is not as effective as RREA(basic). In such KGs RREA(basic) should be preferred, although it is slower.

**Q4.** *To which characteristics of the datasets are supervised, semi-supervised and unsupervised methods sensitive?*

The effectiveness of supervised and semi-supervised methods is positively correlated with the size of the seed alignment and the density of the KGs (see Table 17). Although RDGCN is a supervised method, it is also affected by the degree of similarity between the entity names, as entity embeddings are initialized using the entity name embeddings. Concerning unsupervised methods, the effectiveness of AttrE is negatively correlated to the size of the seed alignment, since it increases the probability of making wrong alignment decisions without improving its learning. In addition, data characteristics such as the similarity of entity names, predicate names, or textual descriptions and literals, are positively correlated with the corresponding relevant method. More precisely, the effectiveness of KDCoE is positively correlated with the number and similarity of textual descriptions, MultiKE with the entity names similarity and AttrE with the predicate names and literals similarity. The unsupervised nature of AttrE justifies the different correlations found w.r.t. other attribute-based methods, in meta-features that refer to the seed alignment size and average relations per entity. In particular, AttrE relies on attribute values (literals) rather than seed alignment, to measure the similarity of entities. Its effectiveness is boosted when the aligned entities have not only a high number of attributes, but also their values are highly similar. In the opposite case, the aligning process is adversely penalized by the noise that is created. BERT\_INT exhibits a similar behavior.

## 5 Conclusions and future work

In this work, we have experimentally evaluated several translational- and Graph Neural Networks-based methods for implementing EA tasks using KG embeddings, as well as, non-learning methods like PARIS. They essentially cover a wide range of supervised, to semi-supervised and unsupervised methods that exploit both relational and factual information of entity descriptions along with different negative sampling and neighborhood range (one- vs multi-hop) strategies. We have measured both the effectiveness and efficiency of all methods over a rich collection of datasets that exhibit different characteristics of entity descriptions provided by KGs.

According to our analysis, *negative sampling* proves to be a critical factor that affects the performance of both relation- and attribute-based KG embeddings methods, as it helps distancing dissimilar entities in the embedding space. Another critical factor that affects the performance of relation-based methods is the *attention mechanism* that weights relevant and important neighbors. Moreover, the performance of relation-based methods depends on the *range of the entity neighborhood (one-hop or multi-hop)*. In particular, we observed that multi-hop methods, require an increased number of negative samples or attention mechanisms to compensate noisy information in the produced embeddings from far distant neighbors. As expected, attribute-based methods (independently of the usage of KG embeddings) are more affected by the different types of factual information exploited in embeddings, such as textual descriptions. In addition, *interaction-based* GNNs seem to better cope with the noise of distant neighbors and improve the utilization of attributes for matching, compared to aggregation-based GNNs. We should notice that the unsupervised attribute-based method AttrE, can run only if the predicates of the different KGs have been sufficiently aligned and it does not impose the 1-to-1 mapping assumption. Supervised and semi-supervised methods do not suffer from this limitation. We conclude that BERT\_INT dominates all methods in terms of effectiveness and efficiency overall, especially when the KGs contain high similar factual information. In dense KGs with low factual similarity, RREA and RDGCN achieve the best and the second best performance, respectively, as also reported in Sun et al. (2020), Zhang et al. (2020), Zhao et al. (2022). MTransE is the fastest method (due to the lack of negative sampling and its simple architecture - see Sects. 2.5.6 and 2.5.7), but it suffers from much lower effectiveness compared to the other methods.

As future work, we would like to benchmark embedding-based EA methods exploiting different Neural Network architectures, e.g., Generative Adversarial (Vretinaris et al. 2021), Self-Adversarial (Sun et al. 2019), Graph Transformer (Cai et al. 2022). Moreover, we plan to extend our testbed with datasets proposed in OAEI<sup>19</sup> that exhibit a higher structural heterogeneity between the KGs to be aligned (e.g., in terms of average degrees). Finally, recent works have revealed that KG embedding methods often exhibit *direct* or *indirect* forms of bias, leading to discrimination (Efthymiou et al. 2022; Fisher et al. 2020). We are particularly interested in assessing the robustness of embedding-based EA methods with respect to

<sup>19</sup> <https://oaei.ontologymatching.org/2022/knowledgegraph/index.html>.

increasing structural diversity (i.e., in terms of connected components) of KGs, that reflect indirect forms of bias (Fanourakis et al. 2023).

**Funding** Open access funding provided by HEAL-Link Greece. This work has received funding from the Hellenic Foundation for Research and Innovation (HFRI) and the General Secretariat for Research and Technology (GSRT), under grant agreement No 969.

## Declarations

**Conflict of interest** The authors have no conflicts of interest to declare that are relevant to the content of this article.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

- Ahmetaj S, Efthymiou V, Fagin R, et al. (2021) Ontology-enriched query answering on relational databases. In: AAAI, pp 15247–15254
- Berrendorf M, Faerman E, Vermue L, et al. (2020) On the ambiguity of rank-based evaluation of entity alignment or link prediction methods. CoRR [arXiv:abs/2002.06914](https://arxiv.org/abs/2002.06914)
- Bordes A, Usunier N, García-Durán A, et al. (2013) Translating embeddings for modeling multi-relational data. In: NeurIPS, pp 2787–2795
- Cai W, Ma W, Zhan J, et al. (2022) Entity alignment with reliable path reasoning and relation-aware heterogeneous graph transformer. In: IJCAI, pp 1930–1937
- Cao Y, Liu Z, Li C, et al. (2019) Multi-channel graph neural network for entity alignment. In: ACL, pp 1452–1461
- Chaurasiya D, Surisetty A, Kumar N, et al. (2022) Entity alignment for knowledge graphs: Progress, challenges, and empirical studies. CoRR [arXiv:abs/2205.08777](https://arxiv.org/abs/2205.08777)
- Chen M, Tian Y, Yang M, et al. (2017) Multilingual knowledge graph embeddings for cross-lingual knowledge alignment. In: IJCAI, pp 1511–1517
- Chen M, Tian Y, Chang K, et al. (2018) Co-training embeddings of knowledge graphs and entity descriptions for cross-lingual entity alignment. In: IJCAI, pp 3998–4004
- Choudhary S, Luthra T, Mittal A, et al. (2021) A survey of knowledge graph embedding and their applications. CoRR [arXiv:abs/2107.07842](https://arxiv.org/abs/2107.07842)
- Christophides V, Efthymiou V, Stefanidis K (2015) Entity Resolution in the Web of Data. Theory and Technology, Morgan & Claypool Publishers, San Rafael, California, Synthesis Lectures on the Semantic Web
- Christophides V, Efthymiou V, Palpanas T, et al. (2021) An overview of end-to-end entity resolution for big data. ACM Comput Surv 53(6):127:1–127:42
- Coleman C, Narayanan D, Kang D et al (2017) Dawnbench: An end-to-end deep learning benchmark and competition. Training 100(101):102
- Coleman C, Kang D, Narayanan D et al (2019) Analysis of dawnbench, a time-to-accuracy machine learning performance benchmark. ACM SIGOPS Oper Syst Rev 53(1):14–25
- Demsar J (2006) Statistical comparisons of classifiers over multiple data sets. J Mach Learn Res 7:1–30

- Devlin J, Chang M, Lee K, et al. (2019) BERT: pre-training of deep bidirectional transformers for language understanding. In: NAACL-HLT, pp 4171–4186
- Dong X, Gabrilovich E, Heitz G, et al. (2014) Knowledge vault: a web-scale approach to probabilistic knowledge fusion. In: SIGKDD, pp 601–610
- Efthymiou V, Stefanidis K, Christophides V (2015) Big data entity resolution: From highly to somehow similar entity descriptions in the web. In: IEEE Big Data, pp 401–410
- Efthymiou V, Stefanidis K, Pitoura E, et al. (2022) FairER: Entity resolution with fairness constraints. In: CIKM, pp 3004–3008
- Fanourakis N, Efthymiou V, Christophides V, et al. (2023) Structural bias in knowledge graphs for the entity alignment task. In: ESWC
- Fisher J, Mittal A, Palfrey D, et al. (2020) Debiasing knowledge graph embeddings. In: EMNLP, pp 7332–7345
- Jiang J, Li M, Gu Z (2021) A survey on translating embedding based entity alignment in knowledge graphs. In: DSC, pp 187–194
- Kamigaito H, Hayashi K (2022) Comprehensive analysis of negative sampling in knowledge graph representation learning. In: ICML, pp 10661–10675
- Kipf TN, Welling M (2017) Semi-supervised classification with graph convolutional networks. In: ICLR
- Kocmi T, Bojar O (2017) An exploration of word embedding initialization in deep-learning tasks. In: ICON, pp 56–64
- Lehmann J, Isele R, Jakob M et al (2015) Dbpedia - A large-scale, multilingual knowledge base extracted from wikipedia. *Semantic Web* 6(2):167–195
- Leone M, Huber S, Arora A et al (2022) A critical re-evaluation of neural methods for entity alignment. *PVLDB* 15(8):1712–1725
- Lin J (1991) Divergence measures based on the shannon entropy. *IEEE Trans Inf Theory* 37(1):145–151
- Mao X, Wang W, Xu H, et al. (2020a) MRAEA: an efficient and robust entity alignment approach for cross-lingual knowledge graph. In: WSDM, pp 420–428
- Mao X, Wang W, Xu H, et al. (2020b) Relational reflection entity alignment. In: CIKM, pp 1095–1104
- Million E (2007) The Hadamard product. *Course Notes* 3(6)
- Nemenyi P (1963) *Distribution-free Multiple Comparisons*. Princeton University, Princeton
- Nickel M, Kiela D (2017) Poincaré embeddings for learning hierarchical representations. In: NeurIPS, pp 6338–6347
- Obraczka D, Schuchart J, Rahm E (2021) EAGER: embedding-assisted entity resolution for knowledge graphs. *CoRR arXiv:abs/2101.06126*
- Parisi L, Neagu D, Ma R et al (2022) Quantum relu activation for convolutional neural networks to improve diagnosis of parkinson's disease and COVID-19. *Expert Syst Appl* 187:115892
- Rebele T, Suchanek FM, Hoffart J, et al. (2016) YAGO: A multilingual knowledge base from wikipedia, wordnet, and geonames. In: ISWC, pp 177–185
- Suchanek FM, Abiteboul S, Senellart P (2011) PARIS: probabilistic alignment of relations, instances, and schema. *PVLDB* 5(3):157–168
- Sun Z, Hu W, Zhang Q, et al. (2018) Bootstrapping entity alignment with knowledge graph embedding. In: IJCAI, pp 4396–4402
- Sun Z, Deng Z, Nie J, et al. (2019) Rotate: Knowledge graph embedding by relational rotation in complex space. In: ICLR
- Sun Z, Zhang Q, Hu W et al (2020) A benchmarking study of embedding-based entity alignment for knowledge graphs. *PVLDB* 13(11):2326–2340
- Tang X, Zhang J, Chen B, et al. (2020) BERT-INT: A bert-based interaction model for knowledge graph alignment. In: IJCAI, pp 3174–3180
- Tarus JK, Niu Z, Mustafa G (2018) Knowledge-based recommendation: a review of ontology-based recommender systems for e-learning. *Artif Intell Rev* 50(1):21–48
- Trisedya BD, Qi J, Zhang R (2019) Entity alignment between knowledge graphs using attribute embeddings. In: AAAL, pp 297–304
- Trouillon T, Welbl J, Riedel S, et al. (2016) Complex embeddings for simple link prediction. In: ICML, pp 2071–2080
- Velickovic P, Cucurull G, Casanova A, et al. (2018) Graph attention networks. In: ICLR
- Vrandečić D, Krötzsch M (2014) Wikidata: a free collaborative knowledgebase. *Commun ACM* 57:78–85. <https://doi.org/10.1145/2629489>
- Vretinari A, Lei C, Efthymiou V, et al. (2021) Medical entity disambiguation using graph neural networks. In: SIGMOD, pp 2310–2318

- Wang Q, Mao Z, Wang B et al (2017) Knowledge graph embedding: a survey of approaches and applications. *IEEE Trans Knowl Data Eng* 29(12):2724–2743
- Wang Z, Yang J, Ye X (2020) Knowledge graph alignment with entity-pair embedding. In: *EMNLP*, pp 1672–1680
- Wang Z, Li M, Gu Z (2021) A review of entity alignment based on graph convolutional neural network. In: *DSC*, pp 144–151
- Wu Y, Liu X, Feng Y, et al. (2019) Relation-aware entity alignment for heterogeneous knowledge graphs. In: *IJCAI*, pp 5278–5284
- Xiong C, Dai Z, Callan J, et al. (2017) End-to-end neural ad-hoc ranking with kernel pooling. In: *SIGIR*, pp 55–64
- Yang B, Yih W, He X, et al. (2015) Embedding entities and relations for learning and inference in knowledge bases. In: *ICLR*
- Zeng K, Li C, Hou L et al (2021) A comprehensive survey of entity alignment for knowledge graphs. *AI Open* 2:1–13
- Zhang C, Zhou M, Han X et al (2017) Knowledge graph embedding for hyper-relational data. *Tsinghua Sci Technol* 22(2):185–197
- Zhang Q, Sun Z, Hu W, et al. (2019) Multi-view knowledge graph embedding for entity alignment. In: *IJCAI*, pp 5429–5435
- Zhang R, Trisedya BD, Li M et al (2022) A benchmark and comprehensive survey on knowledge graph entity alignment via representation learning. *VLDBJ* 31(5):1143–1168
- Zhang Z, Liu H, Chen J, et al. (2020) An industry evaluation of embedding-based entity alignment. In: *COLING*, pp 179–189
- Zhao X, Zeng W, Tang J et al (2022) An experimental study of state-of-the-art entity alignment approaches. *IEEE Trans Knowl Data Eng* 34(6):2610–2625
- Zhu Q, Zhou X, Wu J, et al. (2019) Neighborhood-aware attentional representation for multilingual knowledge graphs. In: *IJCAI*, pp 1943–1949

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

## Authors and Affiliations

Nikolaos Fanourakis<sup>1,2</sup>  · Vasilis Efthymiou<sup>1</sup>  · Dimitris Kotzinos<sup>3</sup> · Vassilis Christophides<sup>3</sup>

✉ Nikolaos Fanourakis  
fanourakis@ics.forth.gr

Vasilis Efthymiou  
vefthym@ics.forth.gr

Dimitris Kotzinos  
Dimitrios.Kotzinos@cyu.fr

Vassilis Christophides  
Vassilis.Christophides@ensea.fr

<sup>1</sup> Institute of Computer Science, FORTH, N. Plastira 100, Heraklion, Greece

<sup>2</sup> Computer Science Department, University of Crete, Voutes Campus, Heraklion, Greece

<sup>3</sup> Lab. ETIS, CNRS, CY Cergy Paris University, ENSEA, 2 av. A. Chauvin, Cergy-Pontoise UMR 8051, France