

A Benchmarking Study of Embedding-based Entity Alignment for Knowledge Graphs

Zequn Sun[†], Qingheng Zhang[†], Wei Hu^{†*}, Chengming Wang[†],
Muhaao Chen[‡], Farahnaz Akrami[§], Chengkai Li[§]

[†] State Key Laboratory for Novel Software Technology, Nanjing University, China

[‡] Department of Computer Science, University of California, Los Angeles, USA

[§] Department of Computer Science and Engineering, University of Texas at Arlington, USA

{zqsun, qhzhang, cmwang}.nju@gmail.com, whu@nju.edu.cn,
muhaochen@ucla.edu, farahnaz.akrami@mavs.uta.edu, cli@uta.edu

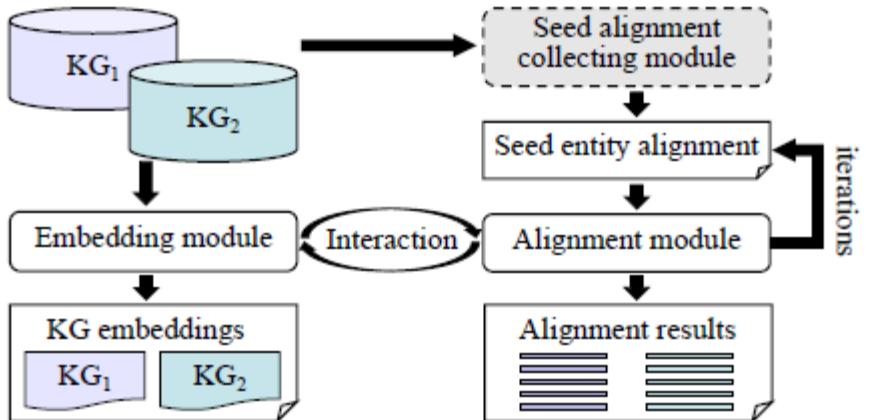


Figure 1: Framework of embedding-based entity alignment

Knowledge graphs (KGs) store facts as triples in the form of (subject entity, relation, object entity) or (subject entity, attribute, literal value).
 relation embedding and attribute embedding

The goal is to identify entities from different KGs that refer to the same entity. There are two typical combination paradigms for module interaction:

- (i) the embedding module encodes the two KGs in two independent embedding spaces, meanwhile the alignment module uses seed alignment to learn a mapping between them; or (ii) the alignment module guides the embedding module to represent the two KGs into one unified space by forcing the aligned entities in seed alignment to hold very similar embeddings. Finally, entity similarities are measured by the learned embeddings.

Besides, to overcome the shortage of seed entity alignment, several approaches deploy semi-supervised learning to iteratively augment new alignment.

Overview

Sect. 1 INTRODUCTION

Sect. 2 PRELIMINARIES: A comprehensive survey on approaches for embedding-based entity alignment

Sect. 3 Benchmark datasets: propose a new sampling algorithm

Sect. 4 Open-source library: develop an open-source library OpenEA using Python and TensorFlow.

Sect. 5 Comprehensive comparison and analysis

Sect. 6 Exploratory experiments: 1. geometric properties of entity embeddings 2. KG embedding models which have not been exploited for entity alignment 3. Comparison to Conventional Approaches

Sect. 7 Future research directions

Literature Review

Existing KG embedding models can be generally divided into three categories:

- (i) translational models, e.g. TransE, TransH, TransR and TransD
- (ii) semantic matching models, e.g. DistMult , ComplEx , Hole , SimplE , RotatE and TuckER
- (iii) deep models, e.g. ProjE,ConvE, R-GCN, KBGAN and DSKG

Embedding-based Entity Alignment :the translational models & graph convolutional networks (GCNs)

Datasets & evaluation metrics. DBP15K and DWY100K are two benchmark datasets for entity alignment in KGs.

Three metrics are widely used in evaluation: (i) proportion of correct links in the top-m ranked results (called Hits@ m , for example, $m = 1$), (ii) mean rank (MR) of correct links, and (iii) mean reciprocal rank (MRR).

Current datasets contain much more high-degree entities (i.e., entities connected with many other entities, which are relatively easy for entity alignment) than real world KGs do. As a result, many approaches may exhibit good performance on these biased datasets.

Additionally, these datasets only focus on one aspect of heterogeneity, e.g., multilingualism, while overlook other aspects, e.g, different schemata and scales. This brings difficulties in understanding the generalization and robustness of embedding-based entity alignment.

Categorization of Techniques

1. Embedding Module
2. Alignment Module
3. Interaction Module

Embedding Module:

1. Triple-based embedding
2. Path-based embedding
3. Neighborhood-based embedding
4. Attribute correlation embedding
5. Literal embedding

Relation embedding is employed by all existing approaches.

- Triple-based embedding captures the local semantics of relation triples. $\phi(e_1, r_1, e_2) = \| \mathbf{e}_1 + \mathbf{r}_1 - \mathbf{e}_2 \|,$
- Path-based embedding: exploit the long-term dependency of relations spanning over relation paths
IPTTransE models relation paths by inferring the equivalence between a direct relation and a multi-hop path.
Assume that there is a direct relation r_3 from e_1 to e_3 . IPTTransE expects the embedding of r_3 to be similar to the path embedding, which is encoded as a combination of its constituent relation embeddings:

$$\mathbf{r}^* = \text{comb}(\mathbf{r}_1, \mathbf{r}_2), \quad \text{minimized } \| \mathbf{r}^* - \mathbf{r}_3 \|$$

- Neighborhood-based embedding $\mathbf{H}^{(i+1)} = \sigma(\hat{\mathbf{D}}^{-\frac{1}{2}} \hat{\mathbf{A}} \hat{\mathbf{D}}^{-\frac{1}{2}} \mathbf{H}^{(i)} \mathbf{W}),$

Attribute embedding

There are two ways for attribute embedding:

Attribute correlation embedding: consider the correlations among attributes

Attributes are regarded as correlated if they are frequently used together to describe an entity. For example, longitude is highly correlated with latitude as they often form a coordinate.

JAPE exploits such correlations for entity alignment, based on the assumption that similar entities should have similar correlated attributes.

$\Pr(a_1, a_2) = \text{sigmoid}(a_1 \cdot a_2)$, Attribute embeddings can be learned by maximizing the probability over all attribute pairs

Literal embedding: introduces literal values to attribute embedding.

Let $v = (c_1, c_2, \dots, c_n)$ be a literal with n characters, where c_i ($1 \leq i \leq n$) is the i^{th} character. AttrE embeds v as
 $\mathbf{v} = \text{comb}(\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_n)$.

With this representation, literals are treated as entities and the relation embedding models like TransE can be used to learn from attribute triples. However, the character-based literal embedding may fail in cross-lingual settings.

Alignment Module

1. Pick a distance metric
2. Design alignment inference strategy.

Distance metrics: Cosine, Euclidean and Manhattan distances are three widely-used metrics. In high-dimensional spaces, a few vectors (called hubs) may repeatedly occur as the k-nearest neighbors of others, the so-called hubness problem.

Alignment inference strategies:

1. Greedy search
2. Collective search: find a global optimal alignment that minimizes $\sum_{(e_1, e_2) \in \mathcal{S}_{\mathcal{K}\mathcal{G}_1, \mathcal{K}\mathcal{G}_2}} \pi(e_1, e_2)$ $O(N^3)$

Another solution is the stable marriage algorithm. The alignment between E_1 and E_2 satisfies a stable marriage if there does not exist a pair of entities that both prefer each other than their current aligned ones. Its solution takes $O(N^2)$ time

Interaction Mode

1. Combination modes.

- Embedding space transformation: embeds two KGs in different embedding spaces and learns a transformation matrix M between the two spaces. $Me_1 \approx e_2$
- Embedding space calibration: minimizes $\| e_1 - e_2 \|$ for each $(e_1, e_2) \in \mathcal{S}'_{KG_1, KG_2}$
- Parameter sharing: directly configures $e_1 = e_2$
- Parameter swapping: swap seed entities in their triples to generate extra triples as supervision.

2. Learning strategies.

- Supervised learning
- Semi-supervised learning:
 - self-training: iteratively propose new alignment to augment seed alignment.
 - co-training: combine two models learned from disjoint entity features and alternately enhances the alignment learning of each other.
- Unsupervised learning

Sect. 3 Benchmark datasets

Iterative Degree-based Sampling

Iterative degree-based sampling (IDS) algorithm, which simultaneously deletes entities in two source KGs with reference alignment until achieving the desired size, meanwhile keeping a similar degree distribution of each sampled dataset as the source KG.

Algorithm 1: Iterative degree-based sampling (IDS)

Input: $\mathcal{KG}_1, \mathcal{KG}_2$, reference alignment \mathcal{S}_{ref} , entity size N ,
hyper-parameters μ, ϵ

// only retain entities in reference alignment

1 Filter $\mathcal{KG}_1, \mathcal{KG}_2$ by \mathcal{S}_{ref} ;

2 Get degree distributions Q_1, Q_2 for $\mathcal{KG}_1, \mathcal{KG}_2$, resp.;

3 do // if fails, run it again

4 Initialize datasets $\mathcal{DS}_1, \mathcal{DS}_2$ from $\mathcal{KG}_1, \mathcal{KG}_2$, resp.;

5 while $|\mathcal{DS}_1| > N \ \&\& \ |\mathcal{DS}_2| > N$ do

6 for \mathcal{DS}_j ($j = 1, 2$) do

7 Get $dsize_j(x, \mu)$ for each degree x ;

8 Get entity deletion probability by PageRank;

9 Delete $dsize_j(x, \mu)$ entities w.r.t. probabilities;

10 Filter $\mathcal{DS}_1, \mathcal{DS}_2$ by \mathcal{S}_{ref} ; update \mathcal{S}_{ref} accordingly;

11 Get degree distributions P_1, P_2 for $\mathcal{DS}_1, \mathcal{DS}_2$, resp.;

12 while $JS(Q_1, P_1) > \epsilon \ || \ JS(Q_2, P_2) > \epsilon$;

13 return $\mathcal{DS}_1, \mathcal{DS}_2, \mathcal{S}_{ref}$;

pling procedure. During iterations, the proportion of entities having degree x in the current dataset, denoted by $P(x)$, cannot always equal the original proportion $Q(x)$. We adjust the entity size to be deleted by $dsize(x, \mu) = \mu(1 + P(x) - Q(x))$, where μ is the base step size (see Line 7). Moreover, we prefer not to delete entities having a big influence on the overall degree distribution, such as the ones of high degree. To achieve this, we leverage the PageRank value for measuring the probability of an entity to be deleted (Line 8).

The difference of two degree distributions

$$JS(Q, P) = \frac{1}{2} \sum_{x=1}^n \left(Q(x) \log \frac{Q(x)}{M(x)} + P(x) \log \frac{P(x)}{M(x)} \right)$$

where $Q(x)$ and $P(x)$ denote the proportions of entities with degree x ($x = 1 \dots n$) in Q, P , respectively, and $M = \frac{Q+P}{2}$.

A small JS divergence between Q and P reveals that they have similar degree distributions.

Dataset Evaluation

We generate two versions of datasets for each pair of source KGs. V1 is gained by directly using the IDS algorithm. For V2, we first randomly delete entities with low degrees ($d \leq 5$) in the source KG to make the average degree doubled, and then execute IDS to get the new KG. As a result, V2 is twice denser than V1 and more similar to existing datasets

- *Random alignment sampling (RAS)* first randomly selects a fixed size (e.g., 15K) of entity alignment between two KGs, and then extracts the relation triples whose head and tail entities are both in the sampled entities.
- *PageRank-based sampling (PRS)* first samples entities from one KG based on the PageRank scores (entities not involved in any alignment are discarded), and then extracts these entities' counterparts from the other KG.

Table 2: Dataset statistics

Datasets	KGs	15K (V1)				15K (V2)				100K (V1)				100K (V2)			
		#Rel.	#Att.	#Rel tr.	#Att tr.	#Rel.	#Att.	#Rel tr.	#Att tr.	#Rel.	#Att.	#Rel tr.	#Att tr.	#Rel.	#Att.	#Rel tr.	#Att tr.
EN-FR	EN	267	308	47,334	73,121	193	189	96,318	66,899	400	466	309,607	497,729	379	364	649,902	503,922
	FR	210	404	40,864	67,167	166	221	80,112	68,779	300	519	258,285	426,672	287	468	561,391	431,379
EN-DE	EN	215	286	47,676	83,755	169	171	84,867	81,988	381	451	335,359	552,750	323	326	622,588	560,247
	DE	131	194	50,419	156,150	96	116	92,632	186,335	196	252	336,240	716,615	170	189	629,395	793,710
D-W	DB	248	342	38,265	68,258	167	175	73,983	66,813	413	493	293,990	451,011	318	328	616,457	467,103
	WD	169	649	42,746	138,246	121	457	83,365	175,686	261	874	251,708	687,860	239	760	588,203	878,219
D-Y	DB	165	257	30,291	71,716	72	90	68,063	65,100	287	379	294,188	523,062	230	277	576,547	547,026
	YG	28	35	26,638	132,114	21	20	60,970	131,151	32	38	400,518	749,787	31	36	865,265	855,161

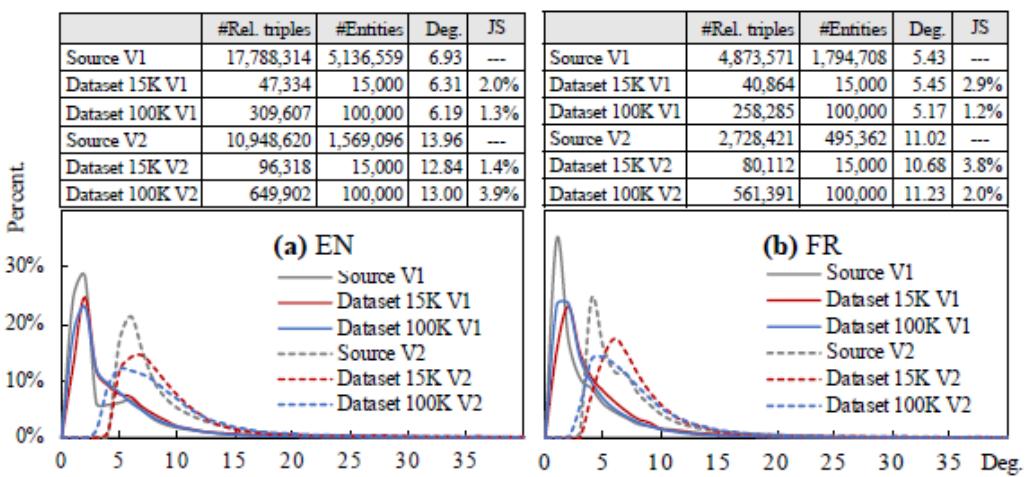


Figure 3: Degree distributions and average degrees of our sampled datasets EN-FR-15K (V1, V2) and EN-FR-100K (V1, V2), compared with DBpedia (the source KG)

Table 3: Comparison of the EN-FR-15K (V1) datasets generated by RAS, PRS and IDS

Datasets	KGs	#Alignment	Deg.	JS	Isolates	Cluster coef.
DBpedia	EN	525,807	6.39	—	0	0.342
	FR		5.43	—	0	0.080
RAS	EN	15,000	0.27	14.5%	85.5%	0.002
	FR		0.17	12.1%	90.1%	0.001
PRS	EN	15,000	1.20	7.3%	68.9%	0.025
	FR		0.63	9.3%	69.4%	0.015
IDS	EN	15,000	6.31	2.0%	0	0.233
	FR		5.45	2.9%	0	0.190

Sect. 4 Open-source library

Sect. 5 Comprehensive comparison and analysis

Table 5: Cross-validation results of current representative approaches on the 15K and 100K datasets

		15K (V1)			15K (V2)			100K (V1)			100K (V2)		
		Hits@1	Hits@5	MRR									
EN-FR	MTransE	.247 <small>± .006</small>	.467 <small>± .009</small>	.351 <small>± .007</small>	.240 <small>± .005</small>	.436 <small>± .007</small>	.336 <small>± .005</small>	.138 <small>± .002</small>	.261 <small>± .004</small>	.202 <small>± .002</small>	.090 <small>± .003</small>	.174 <small>± .003</small>	.135 <small>± .003</small>
	IPTransE	.169 <small>± .013</small>	.320 <small>± .025</small>	.243 <small>± .019</small>	.236 <small>± .012</small>	.449 <small>± .021</small>	.339 <small>± .016</small>	.158 <small>± .004</small>	.277 <small>± .008</small>	.219 <small>± .006</small>	.234 <small>± .007</small>	.431 <small>± .015</small>	.329 <small>± .010</small>
	JAPE	.262 <small>± .006</small>	.497 <small>± .010</small>	.372 <small>± .007</small>	.292 <small>± .009</small>	.524 <small>± .006</small>	.402 <small>± .007</small>	.165 <small>± .002</small>	.310 <small>± .002</small>	.240 <small>± .002</small>	.125 <small>± .003</small>	.239 <small>± .005</small>	.183 <small>± .004</small>
	KDCoE	.581 <small>± .004</small>	.680 <small>± .004</small>	.628 <small>± .003</small>	.730 <small>± .007</small>	.837 <small>± .006</small>	.778 <small>± .005</small>	.482 <small>± .005</small>	.515 <small>± .006</small>	.490 <small>± .005</small>	.611 <small>± .012</small>	.653 <small>± .015</small>	.632 <small>± .014</small>
	BootEA	.507 <small>± .010</small>	.718 <small>± .012</small>	.603 <small>± .011</small>	.660 <small>± .008</small>	.850 <small>± .005</small>	.745 <small>± .005</small>	.389 <small>± .004</small>	.561 <small>± .004</small>	.474 <small>± .004</small>	.640 <small>± .001</small>	.806 <small>± .001</small>	.716 <small>± .000</small>
	GCNAlign	.338 <small>± .002</small>	.589 <small>± .009</small>	.451 <small>± .005</small>	.414 <small>± .005</small>	.698 <small>± .007</small>	.542 <small>± .005</small>	.230 <small>± .002</small>	.412 <small>± .004</small>	.319 <small>± .003</small>	.257 <small>± .002</small>	.455 <small>± .003</small>	.351 <small>± .002</small>
	AttrE	.481 <small>± .010</small>	.671 <small>± .009</small>	.569 <small>± .010</small>	.535 <small>± .015</small>	.746 <small>± .014</small>	.631 <small>± .014</small>	.403 <small>± .019</small>	.572 <small>± .019</small>	.483 <small>± .019</small>	.466 <small>± .011</small>	.644 <small>± .012</small>	.549 <small>± .011</small>
	IMUSE	.569 <small>± .006</small>	.717 <small>± .010</small>	.638 <small>± .008</small>	.607 <small>± .013</small>	.760 <small>± .014</small>	.678 <small>± .013</small>	.439 <small>± .002</small>	.546 <small>± .004</small>	.492 <small>± .003</small>	.461 <small>± .003</small>	.605 <small>± .005</small>	.529 <small>± .004</small>
	SEA	.280 <small>± .015</small>	.530 <small>± .026</small>	.397 <small>± .019</small>	.360 <small>± .018</small>	.651 <small>± .018</small>	.494 <small>± .017</small>	.225 <small>± .011</small>	.399 <small>± .013</small>	.314 <small>± .012</small>	.297 <small>± .002</small>	.500 <small>± .002</small>	.395 <small>± .002</small>
	RSN4EA	.393 <small>± .007</small>	.595 <small>± .012</small>	.487 <small>± .009</small>	.579 <small>± .006</small>	.759 <small>± .006</small>	.662 <small>± .006</small>	.293 <small>± .004</small>	.452 <small>± .006</small>	.371 <small>± .004</small>	.495 <small>± .003</small>	.672 <small>± .005</small>	.578 <small>± .004</small>
EN-DE	MTransE	.307 <small>± .007</small>	.518 <small>± .004</small>	.407 <small>± .006</small>	.193 <small>± .016</small>	.352 <small>± .023</small>	.274 <small>± .018</small>	.140 <small>± .003</small>	.264 <small>± .004</small>	.204 <small>± .004</small>	.115 <small>± .003</small>	.215 <small>± .004</small>	.168 <small>± .003</small>
	IPTransE	.350 <small>± .009</small>	.515 <small>± .012</small>	.43 <small>± .011</small>	.476 <small>± .012</small>	.678 <small>± .011</small>	.571 <small>± .010</small>	.226 <small>± .014</small>	.357 <small>± .019</small>	.292 <small>± .017</small>	.346 <small>± .013</small>	.535 <small>± .016</small>	.437 <small>± .014</small>
	JAPE	.288 <small>± .016</small>	.512 <small>± .018</small>	.394 <small>± .016</small>	.167 <small>± .011</small>	.329 <small>± .015</small>	.250 <small>± .013</small>	.152 <small>± .006</small>	.291 <small>± .009</small>	.223 <small>± .007</small>	.11 <small>± .004</small>	.218 <small>± .006</small>	.167 <small>± .005</small>
	KDCoE	.529 <small>± .014</small>	.629 <small>± .015</small>	.580 <small>± .014</small>	.649 <small>± .017</small>	.788 <small>± .017</small>	.715 <small>± .016</small>	.506 <small>± .014</small>	.591 <small>± .019</small>	.549 <small>± .016</small>	.651 <small>± .011</small>	.756 <small>± .010</small>	.701 <small>± .011</small>
	BootEA	.675 <small>± .004</small>	.820 <small>± .004</small>	.740 <small>± .004</small>	.833 <small>± .015</small>	.912 <small>± .008</small>	.869 <small>± .012</small>	.518 <small>± .003</small>	.673 <small>± .003</small>	.592 <small>± .003</small>	.739 <small>± .004</small>	.851 <small>± .003</small>	.791 <small>± .004</small>
	GCNAlign	.481 <small>± .003</small>	.670 <small>± .005</small>	.571 <small>± .003</small>	.534 <small>± .005</small>	.717 <small>± .005</small>	.618 <small>± .005</small>	.317 <small>± .007</small>	.485 <small>± .008</small>	.399 <small>± .011</small>	.375 <small>± .005</small>	.549 <small>± .006</small>	.457 <small>± .005</small>
	AttrE	.517 <small>± .011</small>	.687 <small>± .013</small>	.597 <small>± .011</small>	.650 <small>± .015</small>	.816 <small>± .008</small>	.726 <small>± .012</small>	.399 <small>± .010</small>	.554 <small>± .012</small>	.473 <small>± .011</small>	.464 <small>± .011</small>	.637 <small>± .010</small>	.546 <small>± .011</small>
	IMUSE	.580 <small>± .017</small>	.720 <small>± .014</small>	.647 <small>± .015</small>	.674 <small>± .011</small>	.803 <small>± .008</small>	.734 <small>± .010</small>	.421 <small>± .005</small>	.516 <small>± .005</small>	.469 <small>± .005</small>	.457 <small>± .005</small>	.588 <small>± .007</small>	.521 <small>± .006</small>
	SEA	.530 <small>± .027</small>	.718 <small>± .026</small>	.617 <small>± .025</small>	.606 <small>± .024</small>	.779 <small>± .018</small>	.687 <small>± .020</small>	.341 <small>± .016</small>	.502 <small>± .017</small>	.421 <small>± .016</small>	.447 <small>± .008</small>	.625 <small>± .006</small>	.532 <small>± .006</small>
	RSN4EA	.587 <small>± .001</small>	.752 <small>± .003</small>	.662 <small>± .001</small>	.791 <small>± .009</small>	.890 <small>± .006</small>	.837 <small>± .008</small>	.430 <small>± .002</small>	.57 <small>± .001</small>	.497 <small>± .001</small>	.639 <small>± .001</small>	.763 <small>± .001</small>	.697 <small>± .001</small>
D-W	MTransE	.259 <small>± .008</small>	.461 <small>± .012</small>	.354 <small>± .008</small>	.271 <small>± .013</small>	.49 <small>± .014</small>	.376 <small>± .013</small>	.210 <small>± .003</small>	.358 <small>± .003</small>	.282 <small>± .003</small>	.148 <small>± .004</small>	.268 <small>± .005</small>	.209 <small>± .005</small>
	IPTransE	.232 <small>± .012</small>	.38 <small>± .016</small>	.303 <small>± .014</small>	.412 <small>± .007</small>	.623 <small>± .010</small>	.511 <small>± .007</small>	.221 <small>± .004</small>	.352 <small>± .008</small>	.285 <small>± .006</small>	.319 <small>± .017</small>	.516 <small>± .024</small>	.413 <small>± .020</small>
	JAPE	.250 <small>± .007</small>	.457 <small>± .010</small>	.348 <small>± .007</small>	.262 <small>± .013</small>	.484 <small>± .019</small>	.368 <small>± .015</small>	.211 <small>± .004</small>	.369 <small>± .004</small>	.287 <small>± .004</small>	.154 <small>± .004</small>	.287 <small>± .005</small>	.221 <small>± .005</small>
	KDCoE	.247 <small>± .020</small>	.412 <small>± .029</small>	.325 <small>± .023</small>	.405 <small>± .020</small>	.640 <small>± .019</small>	.515 <small>± .020</small>	.157 <small>± .003</small>	.243 <small>± .007</small>	.199 <small>± .005</small>	.373 <small>± .010</small>	.550 <small>± .014</small>	.458 <small>± .012</small>
	BootEA	.572 <small>± .008</small>	.744 <small>± .007</small>	.649 <small>± .008</small>	.821 <small>± .004</small>	.926 <small>± .003</small>	.867 <small>± .003</small>	.516 <small>± .006</small>	.685 <small>± .006</small>	.594 <small>± .005</small>	.766 <small>± .007</small>	.892 <small>± .005</small>	.822 <small>± .006</small>
	GCNAlign	.364 <small>± .009</small>	.580 <small>± .010</small>	.461 <small>± .008</small>	.506 <small>± .006</small>	.743 <small>± .005</small>	.612 <small>± .005</small>	.324 <small>± .002</small>	.507 <small>± .004</small>	.409 <small>± .003</small>	.353 <small>± .004</small>	.559 <small>± .006</small>	.449 <small>± .004</small>
	AttrE	.299 <small>± .004</small>	.467 <small>± .003</small>	.381 <small>± .003</small>	.489 <small>± .016</small>	.695 <small>± .018</small>	.585 <small>± .015</small>	.209 <small>± .008</small>	.335 <small>± .011</small>	.273 <small>± .009</small>	.301 <small>± .015</small>	.475 <small>± .018</small>	.386 <small>± .016</small>
	IMUSE	.327 <small>± .016</small>	.523 <small>± .024</small>	.419 <small>± .019</small>	.581 <small>± .016</small>	.778 <small>± .011</small>	.671 <small>± .014</small>	.276 <small>± .010</small>	.437 <small>± .016</small>	.355 <small>± .013</small>	.431 <small>± .011</small>	.631 <small>± .013</small>	.525 <small>± .012</small>
	SEA	.360 <small>± .012</small>	.572 <small>± .015</small>	.458 <small>± .013</small>	.567 <small>± .008</small>	.770 <small>± .007</small>	.660 <small>± .008</small>	.291 <small>± .012</small>	.470 <small>± .014</small>	.378 <small>± .013</small>	.382 <small>± .003</small>	.585 <small>± .003</small>	.479 <small>± .002</small>
	RSN4EA	.441 <small>± .008</small>	.615 <small>± .007</small>	.521 <small>± .007</small>	.723 <small>± .007</small>	.854 <small>± .006</small>	.782 <small>± .006</small>	.384 <small>± .004</small>	.533 <small>± .006</small>	.454 <small>± .005</small>	.634 <small>± .004</small>	.776 <small>± .002</small>	.699 <small>± .003</small>
D-Y	MTransE	.463 <small>± .013</small>	.675 <small>± .011</small>	.559 <small>± .012</small>	.443 <small>± .017</small>	.635 <small>± .013</small>	.533 <small>± .015</small>	.244 <small>± .004</small>	.414 <small>± .006</small>	.328 <small>± .005</small>	.100 <small>± .003</small>	.195 <small>± .005</small>	.152 <small>± .004</small>
	IPTransE	.313 <small>± .009</small>	.456 <small>± .015</small>	.378 <small>± .011</small>	.752 <small>± .018</small>	.873 <small>± .013</small>	.808 <small>± .015</small>	.396 <small>± .014</small>	.558 <small>± .018</small>	.474 <small>± .015</small>	.456 <small>± .016</small>	.620 <small>± .017</small>	.534 <small>± .017</small>
	JAPE	.469 <small>± .009</small>	.687 <small>± .011</small>	.567 <small>± .009</small>	.345 <small>± .010</small>	.546 <small>± .013</small>	.440 <small>± .011</small>	.287 <small>± .007</small>	.474 <small>± .008</small>	.379 <small>± .007</small>	.127 <small>± .004</small>	.244 <small>± .006</small>	.189 <small>± .005</small>
	KDCoE	.661 <small>± .013</small>	.764 <small>± .036</small>	.710 <small>± .021</small>	.895 <small>± .013</small>	.974 <small>± .003</small>	.932 <small>± .008</small>	.565 <small>± .001</small>	.646 <small>± .002</small>	.605 <small>± .002</small>	.540 <small>± .001</small>	.621 <small>± .002</small>	.581 <small>± .001</small>
	BootEA	.739 <small>± .014</small>	.849 <small>± .010</small>	.788 <small>± .012</small>	.958 <small>± .001</small>	.984 <small>± .001</small>	.969 <small>± .001</small>	.703 <small>± .004</small>	.827 <small>± .003</small>	.761 <small>± .003</small>	.886 <small>± .003</small>	.944 <small>± .002</small>	.912 <small>± .002</small>
	GCNAlign	.465 <small>± .012</small>	.626 <small>± .011</small>	.536 <small>± .011</small>	.875 <small>± .005</small>	.948 <small>± .004</small>	.907 <small>± .004</small>	.528 <small>± .003</small>	.695 <small>± .005</small>	.605 <small>± .004</small>	.620 <small>± .006</small>	.779 <small>± .006</small>	.693 <small>± .006</small>
	AttrE	.668 <small>± .012</small>	.803 <small>± .009</small>	.731 <small>± .010</small>	.914 <small>± .015</small>	.97 <small>± .007</small>	.939 <small>± .012</small>	.678 <small>± .017</small>	.81 <small>± .012</small>	.739 <small>± .015</small>	.720 <small>± .01</small>	.846 <small>± .007</small>	.778 <small>± .008</small>
	IMUSE	.392 <small>± .013</small>	.571 <small>± .023</small>	.473 <small>± .017</small>	.899 <small>± .011</small>	.949 <small>± .007</small>	.922 <small>± .009</small>	.536 <small>± .018</small>	.700 <small>± .019</small>	.613 <small>± .018</small>	.629 <small>± .011</small>	.774 <small>± .010</small>	.696 <small>± .010</small>
	SEA	.500 <small>± .011</small>	.706 <small>± .012</small>	.591 <small>± .012</small>	.899 <small>± .005</small>	.950 <small>± .003</small>	.923 <small>± .004</small>	.490 <small>± .042</small>	.677 <small>± .040</small>	.578 <small>± .040</small>	.526 <small>± .021</small>	.687 <small>± .021</small>	.603 <small>± .021</small>
	RSN4EA	.514 <small>± .003</small>	.655 <small>± .004</small>	.580 <small>± .003</small>	.933 <small>± .003</small>	.974 <small>± .001</small>	.951 <small>± .002</small>	.620 <small>± .002</small>	.769 <small>± .003</small>	.688 <small>± .002</small>	.841 <small>± .003</small>	.922 <small>± .001</small>	.877 <small>± .002</small>
MultiKE	.903 <small>± .004</small>	.939 <small>± .003</small>	.920 <small>± .003</small>	.856 <small>± .004</small>	.908 <small>± .002</small>	.881 <small>± .003</small>	.884 <small>± .005</small>	.920 <small>± .004</small>	.901 <small>± .005</small>	.853 <small>± .003</small>	.896 <small>± .003</small>	.874 <small>± .003</small>	
	RDGCN	.931 <small>± .004</small>	.969 <small>± .003</small>	.949 <small>± .003</small>	.936 <small>± .003</small>	.966 <small>± .001</small>	.950 <small>± .002</small>	.897 <small>± .001</small>	.950 <small>± .001</small>	.921 <small>± .001</small>	.911 <small>± .002</small>	.949 <small>± .002</small>	.928 <small>± .002</small>

Means \pm stds. are shown. Top-3 results on each dataset are marked in red, blue and cyan, respectively. The same to the following.

		15K (V1)			15K (V2)			100K (V1)			100K (V2)		
		Hits@1	Hits@5	MRR									
EN-FR	MTransE	.247 ± .006	.467 ± .009	.351 ± .007	.240 ± .005	.436 ± .007	.336 ± .005	.138 ± .002	.261 ± .004	.202 ± .002	.090 ± .003	.174 ± .003	.135 ± .003
	IPTTransE	.169 ± .013	.320 ± .025	.243 ± .019	.236 ± .012	.449 ± .021	.339 ± .016	.158 ± .004	.277 ± .008	.219 ± .006	.234 ± .007	.431 ± .015	.329 ± .010
	JAPE	.262 ± .006	.497 ± .010	.372 ± .007	.292 ± .009	.524 ± .006	.402 ± .007	.165 ± .002	.310 ± .002	.240 ± .002	.125 ± .003	.239 ± .005	.183 ± .004
	KDCoE	.581 ± .004	.680 ± .004	.628 ± .003	.730 ± .007	.837 ± .006	.778 ± .005	.482 ± .005	.515 ± .006	.499 ± .005	.611 ± .012	.653 ± .015	.632 ± .014
	BootEA	.507 ± .010	.718 ± .012	.603 ± .011	.660 ± .006	.850 ± .005	.745 ± .005	.389 ± .004	.561 ± .004	.474 ± .004	.640 ± .001	.806 ± .001	.716 ± .000
	GCNAlign	.338 ± .002	.589 ± .009	.451 ± .005	.414 ± .005	.698 ± .007	.542 ± .005	.230 ± .002	.412 ± .004	.319 ± .003	.257 ± .002	.455 ± .003	.351 ± .002
	AttrE	.481 ± .010	.671 ± .009	.569 ± .010	.535 ± .015	.746 ± .014	.631 ± .014	.403 ± .019	.572 ± .019	.483 ± .019	.466 ± .011	.644 ± .012	.549 ± .011
	IMUSE	.569 ± .006	.717 ± .010	.638 ± .008	.607 ± .013	.760 ± .014	.678 ± .013	.439 ± .002	.546 ± .004	.492 ± .003	.461 ± .003	.605 ± .005	.529 ± .004
	SEA	.280 ± .015	.530 ± .026	.397 ± .019	.360 ± .018	.651 ± .018	.494 ± .017	.225 ± .011	.399 ± .013	.314 ± .012	.297 ± .002	.500 ± .002	.395 ± .002
	RSN4EA	.393 ± .007	.595 ± .012	.487 ± .009	.579 ± .006	.759 ± .006	.662 ± .006	.293 ± .004	.452 ± .006	.371 ± .004	.495 ± .003	.672 ± .005	.578 ± .004
MultiKE	.749 ± .004	.819 ± .005	.782 ± .004	.864 ± .007	.909 ± .005	.885 ± .006	.629 ± .002	.680 ± .002	.655 ± .002	.642 ± .003	.696 ± .003	.670 ± .003	
	RDGCN	.755 ± .004	.854 ± .003	.800 ± .003	.847 ± .006	.919 ± .004	.880 ± .005	.640 ± .004	.732 ± .004	.683 ± .004	.715 ± .003	.787 ± .002	.748 ± .002

Sparse datasets (V1) vs. dense datasets (V2).

Most relation-based approaches perform better on the dense datasets than on the sparse ones, e.g., IPTTransE, BootEA, SEA and RSN4EA.

For the approaches considering attribute triples, KDCoE, GCNAlign, AttrE, IMUSE and RDGCN also perform better on the dense datasets, indicating that the relation embeddings still make contributions.

Differently, MultiKE relies on multiple "views" of features, which make it relatively insensitive to the relation changes. Interestingly, we also see that the performance of two relation-based approaches, MTransE and JAPE, drops on some dense datasets. We believe that this is because they are based on TransE, which has deficiency in handling multi-mapping relations in the dense datasets.

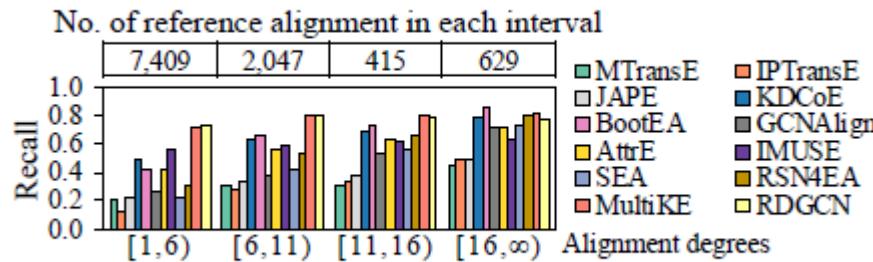


Figure 5: Recall w.r.t. alignment deg. on EN-FR-15K (V1)

We find that all the relation-based approaches run better in aligning entities with rich relation triples while their results decline on long-tail entities,

15K datasets vs. 100K datasets.

We observe that all the approaches perform better on the 15K datasets than on the 100K datasets, except D-Y, because the 100K datasets have more complex structures, causing more difficulties for embedding-based approaches to capture entity proximity.

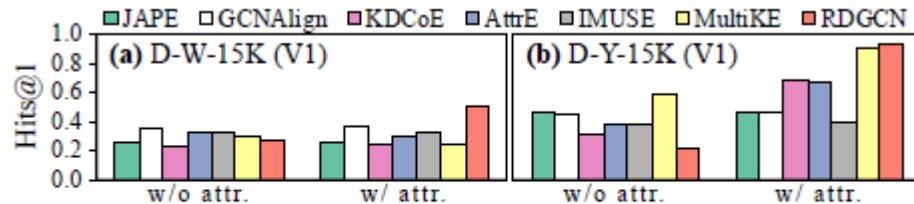


Figure 6: Hits@1 results of JAPE, GCNAlign, KDCoE, AttrE, IMUSE, MultiKE, RDGCN and their degraded variants without attribute embedding

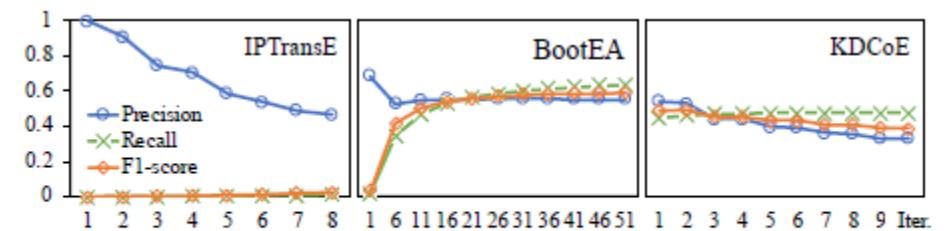


Figure 7: Precision, recall and F1-score of augmented alignment during iterations on EN-FR-100K (V1)

Sect. 6 EXPLORATORY EXPERIMENTS

Similarity Distribution

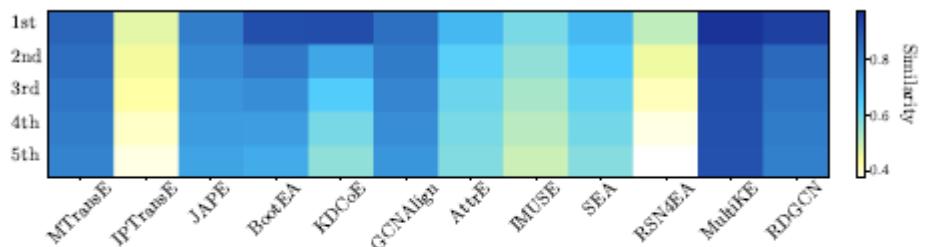


Figure 9: Visualization of the similarities between entities and their top-5 nearest cross-KG neighbors on the D-Y-15K (V1) dataset. The five rows from top to bottom correspond to the similarities from the first to the fifth nearest neighbors, respectively. Darker color indicates larger similarity.

BootEA, KDCoE, MultiKE and RDGCN yield a very high top-1 similarity, while IPTransE and RSN4EA show the opposite. Intuitively, a high top-1 similarity indicates a better quality because it can reflect how confidently the entity embeddings capture the alignment information between two KGs. Most approaches with a high top-1 similarity, such as BootEA, MultiKE and RDGCN, also achieve good performance for entity alignment

The similarity variances between the top-5 nearest neighbors also differ greatly

The ideal similarity distribution for entity alignment is to hold a high top-1 similarity and a large similarity variance.

		15K (V1)			15K (V2)			100K (V1)			100K (V2)		
		Hits@1	Hits@5	MRR									
EN-FR	MTransE	.247 <small>± .006</small>	.467 <small>± .009</small>	.351 <small>± .007</small>	.240 <small>± .005</small>	.436 <small>± .007</small>	.336 <small>± .005</small>	.138 <small>± .002</small>	.261 <small>± .004</small>	.202 <small>± .002</small>	.090 <small>± .003</small>	.174 <small>± .003</small>	.135 <small>± .003</small>
	IPTransE	.169 <small>± .013</small>	.320 <small>± .025</small>	.243 <small>± .019</small>	.236 <small>± .012</small>	.449 <small>± .021</small>	.339 <small>± .016</small>	.158 <small>± .004</small>	.277 <small>± .008</small>	.219 <small>± .006</small>	.234 <small>± .007</small>	.431 <small>± .015</small>	.329 <small>± .010</small>
	JAPE	.262 <small>± .006</small>	.497 <small>± .010</small>	.372 <small>± .007</small>	.292 <small>± .009</small>	.524 <small>± .006</small>	.402 <small>± .007</small>	.165 <small>± .002</small>	.310 <small>± .002</small>	.240 <small>± .002</small>	.125 <small>± .003</small>	.239 <small>± .005</small>	.183 <small>± .004</small>
	KDCoE	.581 <small>± .004</small>	.680 <small>± .004</small>	.628 <small>± .003</small>	.730 <small>± .007</small>	.837 <small>± .006</small>	.778 <small>± .005</small>	.482 <small>± .005</small>	.515 <small>± .006</small>	.499 <small>± .005</small>	.611 <small>± .012</small>	.653 <small>± .015</small>	.632 <small>± .014</small>
	BootEA	.507 <small>± .010</small>	.718 <small>± .012</small>	.603 <small>± .011</small>	.660 <small>± .006</small>	.850 <small>± .005</small>	.745 <small>± .005</small>	.389 <small>± .004</small>	.561 <small>± .004</small>	.474 <small>± .004</small>	.640 <small>± .001</small>	.806 <small>± .001</small>	.716 <small>± .000</small>
	GCNAlign	.338 <small>± .002</small>	.589 <small>± .009</small>	.451 <small>± .005</small>	.414 <small>± .005</small>	.698 <small>± .007</small>	.542 <small>± .005</small>	.230 <small>± .002</small>	.412 <small>± .004</small>	.319 <small>± .003</small>	.257 <small>± .002</small>	.455 <small>± .003</small>	.351 <small>± .002</small>
	AttrE	.481 <small>± .010</small>	.671 <small>± .009</small>	.569 <small>± .010</small>	.535 <small>± .015</small>	.746 <small>± .014</small>	.631 <small>± .014</small>	.403 <small>± .019</small>	.572 <small>± .019</small>	.483 <small>± .019</small>	.466 <small>± .011</small>	.644 <small>± .012</small>	.549 <small>± .011</small>
	IMUSE	.569 <small>± .006</small>	.717 <small>± .010</small>	.638 <small>± .008</small>	.607 <small>± .013</small>	.760 <small>± .014</small>	.678 <small>± .013</small>	.439 <small>± .002</small>	.546 <small>± .004</small>	.492 <small>± .003</small>	.461 <small>± .003</small>	.605 <small>± .005</small>	.529 <small>± .004</small>
	SEA	.280 <small>± .015</small>	.530 <small>± .026</small>	.397 <small>± .019</small>	.360 <small>± .018</small>	.651 <small>± .018</small>	.494 <small>± .017</small>	.225 <small>± .011</small>	.399 <small>± .013</small>	.314 <small>± .012</small>	.297 <small>± .002</small>	.500 <small>± .002</small>	.395 <small>± .002</small>
	RSN4EA	.393 <small>± .007</small>	.595 <small>± .012</small>	.487 <small>± .009</small>	.579 <small>± .006</small>	.759 <small>± .006</small>	.662 <small>± .006</small>	.293 <small>± .004</small>	.452 <small>± .006</small>	.371 <small>± .004</small>	.495 <small>± .003</small>	.672 <small>± .005</small>	.578 <small>± .004</small>
	MultiKE	.749 <small>± .004</small>	.819 <small>± .005</small>	.782 <small>± .004</small>	.864 <small>± .007</small>	.909 <small>± .005</small>	.885 <small>± .006</small>	.629 <small>± .002</small>	.680 <small>± .002</small>	.655 <small>± .002</small>	.642 <small>± .003</small>	.696 <small>± .003</small>	.670 <small>± .003</small>
	RDGCN	.755 <small>± .004</small>	.854 <small>± .003</small>	.800 <small>± .003</small>	.847 <small>± .006</small>	.919 <small>± .004</small>	.880 <small>± .005</small>	.640 <small>± .004</small>	.732 <small>± .004</small>	.683 <small>± .004</small>	.715 <small>± .003</small>	.787 <small>± .002</small>	.748 <small>± .002</small>

Hubness and Isolation

We observe that the approaches which yield fewer isolated and hub entities, such as MultiKE and RDGCN, achieve the leading performance of entity alignment, and vice versa.

To resolve the hubness and isolation problem, we explore cross-domain similarity local scaling (CSLS)

$$\text{CSLS}(\mathbf{x}_s, \mathbf{x}_t) = 2 \cos(\mathbf{x}_s, \mathbf{x}_t) - \psi_t(\mathbf{x}_s) - \psi_s(\mathbf{x}_t),$$

$\psi_t(\mathbf{x}_s)$ denotes the average similarity between the source entity \mathbf{x}_s and its top- k nearest neighbors in the target KG.

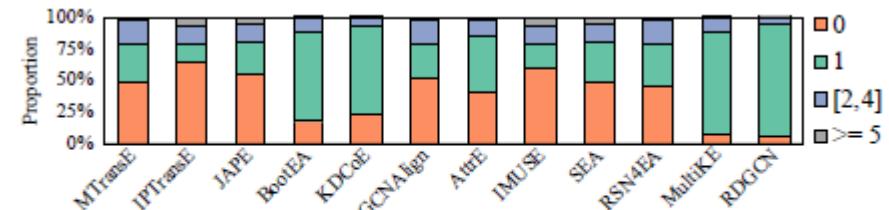


Figure 10: Proportions of target entities appearing 0, 1 times as the nearest neighbors on D-Y-15K (V1)

		15K (V1)			15K (V2)			100K (V1)			100K (V2)		
		Hits@1	Hits@5	MRR									
EN-FR	MTransE	.247 ± .006	.467 ± .009	.351 ± .007	.240 ± .005	.436 ± .007	.336 ± .005	.138 ± .002	.261 ± .004	.202 ± .002	.090 ± .003	.174 ± .003	.135 ± .003
	IPTransE	.169 ± .013	.320 ± .025	.243 ± .019	.236 ± .012	.449 ± .021	.339 ± .016	.158 ± .004	.277 ± .008	.219 ± .006	.234 ± .007	.431 ± .015	.329 ± .010
	JAPE	.262 ± .006	.497 ± .010	.372 ± .007	.292 ± .009	.524 ± .006	.402 ± .007	.165 ± .002	.310 ± .002	.240 ± .002	.125 ± .003	.239 ± .005	.183 ± .004
	KDCoE	.581 ± .004	.680 ± .004	.628 ± .003	.730 ± .007	.837 ± .006	.778 ± .005	.482 ± .005	.515 ± .006	.499 ± .005	.611 ± .012	.653 ± .015	.632 ± .014
	BootEA	.507 ± .010	.718 ± .012	.603 ± .011	.660 ± .006	.850 ± .005	.745 ± .005	.389 ± .004	.561 ± .004	.474 ± .004	.640 ± .001	.806 ± .001	.716 ± .000
	GCNAign	.338 ± .002	.589 ± .009	.451 ± .005	.414 ± .005	.698 ± .007	.542 ± .005	.230 ± .002	.412 ± .004	.319 ± .003	.257 ± .002	.455 ± .003	.351 ± .002
	AttrE	.481 ± .010	.671 ± .009	.569 ± .010	.535 ± .015	.746 ± .014	.631 ± .014	.403 ± .019	.572 ± .019	.483 ± .019	.466 ± .011	.644 ± .012	.549 ± .011
	IMUSE	.569 ± .006	.717 ± .010	.638 ± .008	.607 ± .013	.760 ± .014	.678 ± .013	.439 ± .002	.546 ± .004	.492 ± .003	.461 ± .003	.605 ± .005	.529 ± .004
	SEA	.280 ± .015	.530 ± .026	.397 ± .019	.360 ± .018	.651 ± .018	.494 ± .017	.225 ± .011	.399 ± .013	.314 ± .012	.297 ± .002	.500 ± .002	.395 ± .002
	RSN4EA	.393 ± .007	.595 ± .012	.487 ± .009	.579 ± .006	.759 ± .006	.662 ± .006	.293 ± .004	.452 ± .006	.371 ± .004	.495 ± .003	.672 ± .005	.578 ± .004
	MultiKE	.749 ± .004	.819 ± .005	.782 ± .004	.864 ± .007	.909 ± .005	.885 ± .006	.629 ± .002	.680 ± .002	.655 ± .002	.642 ± .003	.696 ± .003	.670 ± .003
	RDGCN	.755 ± .004	.854 ± .003	.800 ± .003	.847 ± .006	.919 ± .004	.880 ± .005	.640 ± .004	.732 ± .004	.683 ± .004	.715 ± .003	.787 ± .002	.748 ± .002

Stable matching (a.k.a. stable marriage)

The entity alignment between two KGs is stable when there does not exist another predicted aligned pair $(e_1; e_2)$ of higher preference than those of e_1 and e_2 to their current matches. The preference can be calculated based on a similarity metric such as CSLS.

Table 6: Hits@1 w.r.t. distance metrics and alignment inference strategies on D-Y-15K (V1)

	Greedy	Greedy w/ CSLS	SM	SM w/ CSLS
MTransE	.463 \pm .013	.550 \pm .009	.694 \pm .006	.697 \pm .010
IPTransE	.313 \pm .009	.339 \pm .013	.370 \pm .018	.369 \pm .018
JAPE	.469 \pm .009	.549 \pm .009	.692 \pm .015	.691 \pm .015
KDCoE	.661 \pm .013	.679 \pm .000	.840 \pm .024	.815 \pm .031
BootEA	.739 \pm .014	.741 \pm .009	.783 \pm .007	.782 \pm .006
GCNAlign	.465 \pm .012	.531 \pm .008	.613 \pm .008	.582 \pm .010
AttrE	.668 \pm .012	.778 \pm .012	.845 \pm .012	.857 \pm .012
IMUSE	.392 \pm .013	.448 \pm .018	.520 \pm .028	.518 \pm .030
SEA	.500 \pm .011	.557 \pm .017	.647 \pm .012	.650 \pm .012
RSN4EA	.514 \pm .003	.548 \pm .003	.571 \pm .002	.575 \pm .004
MultiKE	.903 \pm .004	.925 \pm .003	.951 \pm .003	.956 \pm .002
RDGCN	.931 \pm .004	.956 \pm .002	.962 \pm .002	.979 \pm .001

In summary, existing approaches concentrate on developing more powerful embedding and interaction methods, but some methods for the alignment module can also improve performance.

Sect. 7 SUMMARY AND FUTURE DIRECTIONS

- (i) RDGCN, BootEA and MultiKE achieve the most competitive performance. This suggests that incorporating both literal information and carefully-designed bootstrapping can help entity alignment.
- (ii) For the embedding models designed for link prediction, we find that not all of them are suitable for entity alignment.
- (iii) Currently, the alignment inference strategy receives little attention. Our preliminary results show that the CSLS distance metric and the stable matching strategy can bring performance improvement to all the approaches.
- (iv) We also find that embedding-based and conventional entity alignment approaches are complementary to each other.
- (v) For choosing appropriate approaches based on the available resources in real-world scenarios, Table 9 summarizes the required information of embedding-based and conventional entity alignment approaches in our experimental analysis.

Table 9: Summary of the required information of embedding-based and conventional entity alignment approaches

	MTransE	IPTransE	JAPE	KDCoE	BootEA	GCNAlign	AttrE	IMUSE	SEA	RSN4EA	MultiKE	RDGCN	LogMap	PARIS
Relation/attribute triples	* /	* /	* / o	o / o	* /	* / o	o / o	o / o	* /	* / o	o / o	* / o	o / *	o / *
Pre-aligned ent./prop.	* / o	* / o	* / o	* /	* / o	* /	* /	* /	* /	* / o	* /	* /	o / o	o / o
Word embed./Google trans.				o /		o /			o /	o /	o /		/ Δ	/ Δ

"*" means "mandatory", "o" means "optional", "Δ" means "mandatory for cross-lingual entity alignment", and blank means "not applicable".

Future Directions

- Unsupervised entity alignment
- Long-tail entity alignment
- Large-scale entity alignment
- Entity alignment in non-Euclidean spaces