

ClusterEA: Scalable Entity Alignment with Stochastic Training and Normalized Mini-batch Similarities

Yunjun Gao*
Zhejiang University
Hangzhou, China
gaoyj@zju.edu.cn

Xiaoze Liu*
Zhejiang University
Hangzhou, China
xiaoze@zju.edu.cn

Junyang Wu
Zhejiang University
Hangzhou, China
wujunyang@zju.edu.cn

Tianyi Li
Aalborg University
Aalborg, Denmark
tianyi@cs.aau.dk

Pengfei Wang
Zhejiang University
Ningbo, China
wangpf@zju.edu.cn

Lu Chen
Zhejiang University
Hangzhou, China
luchen@zju.edu.cn

ABSTRACT

Entity alignment (EA) aims at finding equivalent entities in different knowledge graphs (KGs). Embedding-based approaches have dominated the EA task in recent years. Those methods face problems that come from the geometric properties of embedding vectors, including hubness and isolation. To solve these geometric problems, many normalization approaches have been adopted for EA. However, the increasing scale of KGs renders it hard for EA models to adopt the normalization processes, thus limiting their usage in real-world applications. To tackle this challenge, we present ClusterEA, a general framework that is capable of scaling up EA models and enhancing their results by leveraging normalization methods on mini-batches with a high entity equivalent rate. ClusterEA contains three components to align entities between large-scale KGs, including stochastic training, ClusterSampler, and SparseFusion. It first trains a large-scale Siamese GNN for EA in a stochastic fashion to produce entity embeddings. Based on the embeddings, a novel ClusterSampler strategy is proposed for sampling highly overlapped mini-batches. Finally, ClusterEA incorporates SparseFusion, which normalizes local and global similarity and then fuses all similarity matrices to obtain the final similarity matrix. Extensive experiments with real-life datasets on EA benchmarks offer insight into the proposed framework, and suggest that it is capable of outperforming the state-of-the-art scalable EA framework by up to 8 times in terms of *Hits@1*.

CCS CONCEPTS

• **Computing methodologies** → **Knowledge representation and reasoning**; *Semantic networks*.

*Both authors contributed equally to this research.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

KDD '22, August 14–18, 2022, Washington, DC, USA

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9385-0/22/08...\$15.00

<https://doi.org/10.1145/3534678.3539331>

KEYWORDS

Entity Alignment, Knowledge Graph, Graph Neural Network

ACM Reference Format:

Yunjun Gao, Xiaoze Liu, Junyang Wu, Tianyi Li, Pengfei Wang, and Lu Chen. 2022. ClusterEA: Scalable Entity Alignment with Stochastic Training and Normalized Mini-batch Similarities. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '22)*, August 14–18, 2022, Washington, DC, USA. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3534678.3539331>

1 INTRODUCTION

Knowledge graphs (KGs) represent collections of relations between real-world objects, which facilitate many downstream applications, such as semantic search [39] and recommendation systems [44]. Although various KGs have been constructed in recent years, they are still highly incomplete. To be more specific, KGs built from different data sources hold unique information individually while having overlapped entities. This motivates us to integrate the knowledge of different KGs with the overlapped entities to complete KGs. Entity alignment (EA) [6], a fundamental strategy for knowledge integration, has been widely studied. EA aims to align entities from different KGs that refer to the same real-world objects, and thus, it facilitates the completion of KGs.

Embedding-based EA has been proposed [6], and has been witnessed rapid development in recent years [20, 23, 26, 29, 33] thanks to the use of Graph Neural Networks (GNNs) [14, 18, 36]. They assume that the neighbors of two equivalent entities in KGs are also equivalent [22]. Based on this, they align entities by applying representation learning to KGs. We summarize the process of Embedding-based EA as the following three steps: (i) taking two input KGs and collecting *seed alignment* as training data; (ii) training an EA model with the isomorphic graph structure of two KGs to encode entities into embedding vectors; and (iii) aligning the equivalent entities between the two input KGs based on a specific similarity measurement (e.g., cosine similarity) of their corresponding embeddings.

The size of real-world KGs is much larger than that of conventional datasets used in evaluating EA tasks. For instance, a real-world KG YAGO3 includes 17 million entities [12]. Thus, EA methods should be scaled up to massive data in order to adapt to real-world applications. However, a recent proposal [12] lost

too much graph structure information, trading the quality of results for scalability. Worse still, as the input KGs become larger, using greedy search [34] to find corresponding entities from one KG to another with top-1-nearest neighbor becomes more challenging. Specifically, the geometric properties of high-dimensional embedding vector spaces lead to problems for embedding-based EA, namely, *geometric problems*, mainly including the hubness and isolation problems [34]. The hubness problem indicates that some points (known as hubs) frequently appear as the top-1 nearest neighbors of many other points in the vector space. The isolation problem implies that some outliers would be isolated from any point clusters. As the scale of input KGs grows, these problems become even more severe due to the increasing number of candidates of nearest neighbor search. EA is generally assumed to follow 1-to-1 mapping [6]. Many existing methods have been proposed to solve the geometric problems by making the similarity matrix better satisfy this assumption, i.e., *normalization* methods. A list of widely used normalization methods for EA includes (i) *Cross-domain Similarity Local Scaling (CSLS)* [19] that is adopted from word translation [11, 21, 26, 32]; (ii) *Gale-Shapley algorithm* [27] that treats EA as stable matching [41, 42]; (iii) *Hungarian algorithm* [15] and (iv) *Sinkhorn iteration* [9] that both transform EA as the assignment problem [10, 11, 24]. Specifically, Sinkhorn iteration [9] could significantly improve the accuracy of matching entities with embeddings [11, 24]. Moreover, Sinkhorn iteration suits the normalization approach of EA best as it can be easily parallelized on the GPU. Nonetheless, the existing normalization approaches [9, 15, 19, 27] are at least of quadratic complexity. This prohibits them from being applied to large-scale data, thus limiting their real-world applications.

In this work, we aim to scale up the normalization process of the similarity matrix to achieve higher EA performance. We adopt a standard machine learning technique, sampling mini-batches, to perform EA in linear time. Specifically, we first train a GNN model to obtain the global embeddings of two KGs. Then, we generate mini-batches for two KGs by placing entities that could find their equivalent together. Next, we calculate and normalize a local similarity matrix between two sets of entities selected for each mini-batch by Sinkhorn iteration. Finally, we merge the local similarities into a unified and sparse similarity matrix. With this strategy, the final similarity matrix is normalized with Sinkhorn iteration, thus achieving higher accuracy, and the time and space complexities can also be significantly reduced. However, its materialization is non-trivial due to the two major challenges:

- *How to sample mini-batches with high entity equivalence rate to ensure 1-to-1 mapping?* Splitting mini-batches on two KGs is quite different from conventional tasks. To transfer EA within mini-batches into the assignment problem, we should place possibly equivalent entities into the same batch to meet the 1-to-1 mapping assumption. Nevertheless, the mapping is only partially known (as the training set) for the EA task, making it hard for two entities to be aligned and hence to be placed in the same batch. Intuitively, randomly splitting two KGs will make most mini-batches fail to correspond. An existing study [12] proposes a rule-based method to split batches with a higher rate of entity equivalence. Nonetheless, its results are still not enough to satisfy the 1-to-1 mapping assumption.

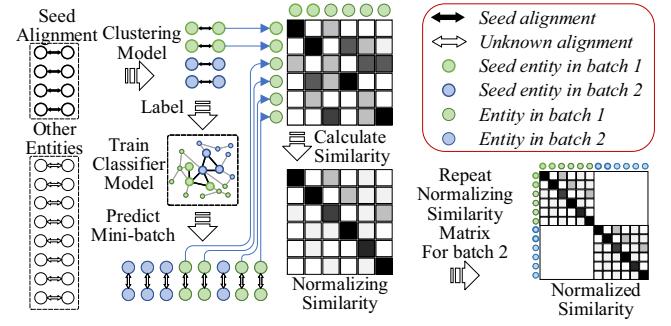


Figure 1: An example of normalizing mini-batch similarity matrix with ClusterSampler strategy

- *How to fuse the mini-batch similarity matrices to ensure high accuracy?* The similarity matrix of each batch focuses only on its local information. Thus, the results generally deviate from those obtained without sampling, even if the batch sampler is good enough. This motivates us to study two problems. First, how to combine the similarities of the mini-batches to get a better global optimal match, after getting the mini-batch. Second, how to improve the accuracy of results obtained by batch samplers.

To tackle the aforementioned two challenges, we present a general EA framework, ClusterEA, that can be applied to arbitrary GNN-based EA models while can be able to achieve high accuracy and scalability. To scale up GNN training over large-scale KGs, ClusterEA first utilizes neighborhood sampling [14] to train a large-scale GNN for EA in a stochastic fashion. Then, with the embedding obtained from the GNN model, ClusterEA proposes a novel *strategy* for generating high-quality mini-batches in EA, called ClusterSampler. As depicted in Figure 1, ClusterSampler first labels the training pairs into different batches with a *clustering* method. Then, it performs supervised *classification* using previously labeled training pairs to generate mini-batches for all the entities. By changing the *clustering* and *classification* models, the strategy can sample mini-batches by capturing multiple aspects of the graph structure information. We propose two sampling methods within ClusterSampler, including (i) *Intra-KG Structure-based ClusterSampler (ISCS)* that retains the intra-KG graph structure information such as the neighborhood of nodes; and (ii) *Cross-KG Mapping-based ClusterSampler (CMCS)* which retains the inter-KG matching information provided by the learned embedding. These methods sample two KGs into multiple mini-batches with a high entity equivalent rate to satisfy the 1-to-1 mapping assumption better. Finally, ClusterEA uses our proposed SparseFusion, which fuses the similarity matrices calculated separately using Sinkhorn iteration with normalized global similarity. The SparseFusion produces a matrix with high sparsity while keeping as much valuable information as possible. Our contributions are summarized as follows:

- *Scalable EA framework.* We develop ClusterEA¹, a scalable framework for EA, which reduces the complexity of normalizing similarity matrices with mini-batch. To the best of our knowledge, this is the first framework that utilizes stochastic GNN training for large-scale EA. Any GNN model can be easily integrated into

¹<https://github.com/joker-xii/ClusterEA>

ClusterEA to deal with large-scale EA (Section 4) with better scalability and higher accuracy.

- *Fast and accurate batch samplers.* We present ClusterSampler, a novel *strategy* that samples batches by learning the latent information from embeddings of the EA model. In the strategy, we implement two samplers to capture different aspects of KG information, including (i) ISCS that aims at retaining the intra-KG graph structure information, and (ii) CMCS that aims at retaining the inter-KG alignment information. Unlike the previous rule-based method, the two samplers are learning-based, and thus can be parallelized and produce high-quality mini-batches.
- *Fused local and global similarity matrix.* We propose SparseFusion for normalizing and fusing not only local similarity matrices but also global similarity matrices into one unified sparse matrix, which enhances the expressive power of EA models.
- *Extensive experiments.* We conduct comprehensive experimental evaluation on EA tasks compared against state-of-the-art approaches over the existing EA benchmarks. Considerable experimental results demonstrate that ClusterEA successfully achieves satisfactory accuracy on both conventional datasets and large-scale datasets (Section 5).

2 RELATED WORK

Most existing EA proposals find equivalent entities by measuring the similarity between the embeddings of entities. Structures of KGs are the basis for the embedding-based EA methods. Representative EA approaches that rely purely on KGs' structures can be divided into two categories, namely, *KGE-based EA* [6, 31, 32, 47] and *GNN-based EA* [4, 20, 25, 29, 33, 37]. The former incorporates the KG embedding models (e.g., TransE [3]) to learn entity embeddings. The latter learns the entity embeddings using GNNs [18], which aggregates the neighbors' information of entities.

In recent years, GNN-based models have demonstrated their outstanding performance [23]. This is contributed by the strong modeling capability on the non-Euclidean structure of GNNs with anisotropic attention mechanism [36]. Nonetheless, they suffer from poor scalability [12] due to the difficulty in sampling mini-batches with neighborhood information on KGs. LargeEA [12], the first study focusing on the scalability of EA, proposes to train GNN models on relatively small batches of two KGs independently. The small batches are generated with a rule-based partition strategy called METIS-CPS. However, massive information of both graph structures and seed alignment is lost during the process, resulting in poor structure-based accuracy. In contrast, ClusterEA utilizes neighborhood sampling [14]. Specifically, it trains one unified GNN model on the two KGs, during which the loss of structure information is neglectable. Moreover, ClusterEA creates mini-batches using multiple aspects of graph information, resulting in better entity equivalent rate than METIS-CPS.

In addition to structure information, many existing proposals facilitate the EA performance by employing *side information* of KGs, including *entity names* [2, 10–12, 22, 24, 30, 35, 43, 45, 46], *descriptions* [35, 45], *images* [21], and *attributes* [22, 30, 35, 37, 38, 40, 45]. Such proposals are able to mitigate the geometric problems [34]. Nonetheless, the models using side information mainly have two main limitations. First, side information may not be available due to privacy concerns, especially for industrial applications [23, 25,

26]. Second, models that incorporating machine translation or pre-aligned word embeddings may be overestimated due to the name bias issue [5, 21, 22, 28]. Thus, compared with the models employing side information, the structure-only methods are more general and not affected by bias of benchmarks. To this end, we do not incorporate side information in ClusterEA.

In order to solve the geometric problems of embedding-based EA approaches, CSLS [19] have been widely adopted, which normalizes the similarity matrix in recent studies [11, 21, 26, 32]. However, CSLS does not perform full normalization of the similarity matrix. As a result, its improvement over greedy search of top-1-nearest neighbor is limited. CEA [41, 42] adopts the Gale-Shapley algorithm [27] to find the stable matching between entities of two KGs, which produces higher-quality results than CSLS. Nevertheless, the Gale-Shapley algorithm is hard to be parallelized, and hence, it is almost infeasible to perform large scale EA. Recent studies have transformed EA into the assignment problem [11, 24]. They adopt Hungarian algorithm [15] or Sinkhorn iteration [9] to normalize the similarity matrix. However, the computational cost of such algorithms is high, prohibiting them from being applied to large scale EA. ClusterEA also utilizes Sinkhorn iteration [9] which performs full normalization on the similarity matrix with GPU acceleration. Moreover, ClusterEA develops novel batch-sampling methods for adopting the normalization to large-scale datasets with the loss of information being minimized.

3 PRELIMINARIES

We proceed to introduce preliminary definitions. Based on these, we formalize the problem of entity alignment.

DEFINITION 1. A **knowledge graph** (KG) can be denoted as $G = (E, R, T)$, where E is the set of entities, R is the set of relations, and $T = \{(h, r, t) \mid h, t \in E, r \in R\}$ is the set of triples, each of which represents an edge from the head entity h to the tail entity t with the relation r .

DEFINITION 2. **Entity alignment** (EA) [34] aims to find the 1-to-1 mapping of entities ϕ from a source KG $G_s = (E_s, R_s, T_s)$ to a target KG $G_t = (E_t, R_t, T_t)$. Formally, $\phi = \{(e_s, e_t) \in E_s \times E_t \mid e_s \equiv e_t\}$, where $e_s \in E_s$, $e_t \in E_t$, and \equiv is an equivalence relation between e_s and e_t . In most cases, a small set of equivalent entities $\phi' \subset \phi$ is known beforehand and used as seed alignment.

DEFINITION 3. **Embedding-based EA** aims to learn a set of embeddings for all entities E_s and E_t , denoted as $\mathbf{f} \in \mathcal{R}^{(|E_s|+|E_t|) \times D}$, and then, it tries to maximize the similarity (e.g. cosine similarity) of entities that are equivalent in ϕ , where D is the size of embedding vectors.

4 OUR FRAMEWORK

In this section, we present our proposed framework ClusterEA, a novel scalable EA framework. We start with the overall framework, followed by details on each component of our framework.

4.1 Overall Framework of ClusterEA

As shown in Figure 2, to scale up GNN training over large-scale KGs, ClusterEA first utilizes neighborhood sampling [14] to train a large-scale GNN for EA in a stochastic fashion. Thereafter, ClusterEA proposes a novel learning-based batch sampling strategy ClusterSampler. It uses the embedding vectors obtained from stochastic

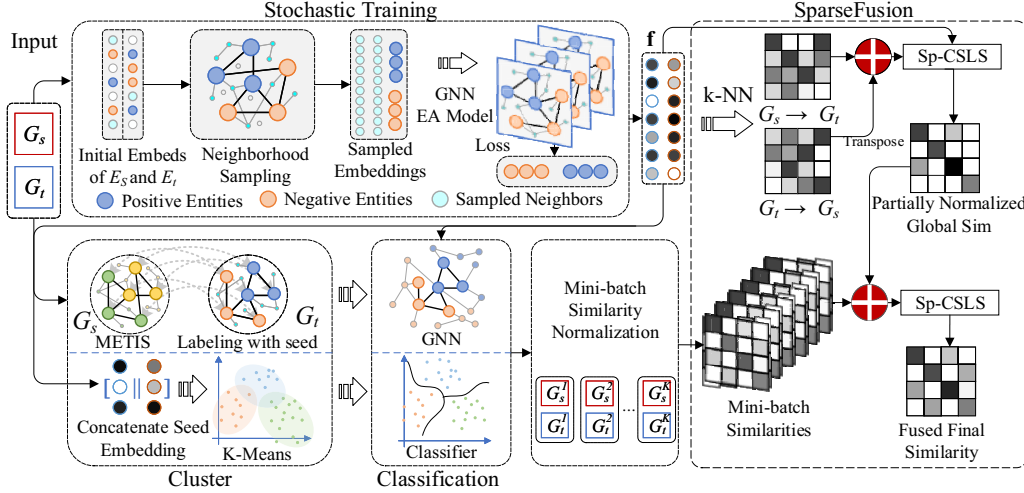


Figure 2: The overall ClusterEA framework

training. In the strategy, two batch samplers that learn multiple aspects of the input KGs are applied to the input KGs, including (i) ISCS that aims at retaining the intra-KG graph structure information such as the neighborhood of nodes, and (ii) CMCS that aims at retaining the inter-KG matching information provided by the learned embedding. These methods sample two KGs into multiple mini-batches with a high entity equivalent rate that better satisfies the 1-to-1 mapping assumption. Finally, ClusterEA proposes SparseFusion, which fuses the normalized local similarity matrices with partially normalized global similarity. The SparseFusion produces the fused final matrix with high sparsity while keeping as much valuable information as possible.

4.2 Stochastic Training of GNNs for EA

GNN-based methods [20–23, 33] have dominated the EA tasks with promising performances by propagating the information of seed alignments to their neighbors. Inspired by this, we propose to incorporate GNN-based models into ClusterEA. ClusterEA provides a general framework for training a Siamese GNN on both G_s and G_t that all GNN-based EA models follow [23]. As a result, any existing GNN model can be used in ClusterEA to produce the structural feature embeddings of entities.

To scale up the existing GNN-based EA models, ClusterEA trains the models with neighborhood sampling. We sample mini-batches based on the seed alignment. Specifically, following the negative sampling process, we first randomly select a mini-batch of size N_p containing source and target entities ϕ'_s and ϕ'_t in the seed alignment ϕ' that are equivalent. Then, we randomly sample source and target entities size of N_n with seed alignment in their whole entity sets that does not overlap with selected seed entities, denoted as $\theta_s = \{e_s | e_s \in E_s \cap e_s \notin \phi'_s\}$ and $\theta_t = \{e_t | e_t \in E_t \cap e_t \notin \phi'_t\}$. Finally, a mini-batch $B = \{B_s, B_t\} = \{(\phi'_s \cup \theta_s), (\phi'_t \cup \theta_t)\}$ is generated waiting for the GNN model to produce its embeddings.

Generally, the GNN-based models train the entity's embedding by propagating the neighborhood information [18, 26, 36]. Formally, the embedding of an entity $v \in B$ in the k_{th} layer of GNN h_v^k is obtained by aggregating localized information via

$$\begin{aligned} a_v^{(k)} &= \text{Aggregate}^{(k)}(\{h_u^{(k-1)} \mid u \in \mathcal{N}(v)\}) \\ h_v^{(k)} &= \text{Update}^{(k)}(a_v^{(k)}, h_v^{(k-1)}) \end{aligned} \quad (1)$$

where $h_v^0 \in \mathcal{R}^D$ is a learnable embedding vector initialized with Glorot initialization, and $\mathcal{N}(v)$ represents the set of neighboring entities around v . The model's final output on entity e is denoted as \mathbf{f}_e . By applying neighborhood sampling, the size of neighborhood in each GNN layer of each entity is limited that no more than a fan out hyperparameter F , formally $|\mathcal{N}(v)| \leq F$. We place the graph information on CPU memory. When computing one layer of GNN in each batch, the neighbor of entities in current batch is sampled to form a block of graph, the graph information of this block will be transferred to GPU memory, along with the computation graph, making the final loss backward propagated with GPU.

To maximize the similarities of equivalent entities in each mini-batch, GNN-based EA models often use triplet loss along with negative sampling [20, 25, 26, 37]. In this paper, the *Normalized Hard Sample Mining (NHSM)* loss [23], an extension for negative sampling that could significantly reduce training epochs, is adopted by us for training large-scale GNNs. We detail how we apply the NHSM loss in Appendix A.

Discussions. Due to the neighborhood sampling process, our stochastic version of EA training will have a certain graph information loss, which is minimized with the randomness of sampling. Recent studies [8] have proposed to limit the sampling in small fixed sub-graphs for better training speed, which will further decrease the accuracy. LargeEA [12], on the other hand, trains multiple GNNs on small batches generated with a rule-based method. Such an approach could fasten the training process. Nonetheless, much of the structure information is lost during partitioning, incurring poor performance. Although LargeEA can scale up the EA models to deal with large-scale datasets, it has too much trade-off on the accuracy, which is unreasonable.

4.3 Learning-based Mini-batch Samplers

After obtaining the KG embeddings, we aim to build batch samplers that utilize the features learned by EA models for generating mini-batches with high entity equivalent rates. To this end, we present the ClusterSampler strategy, which first *clusters* the training set for obtaining batch labels, and then fits a *classification* model with train labels to put all the entities into the right batch. The two models must satisfy two rules: (i) *scalability* that both the classification and

clustering method should be able to apply on large-scale data; and (ii) *distinguishability* that the classifier must be able to distinguish entities into different labels the clustering method provides. If the *scalability* rule is not satisfied, the model will crash due to limited computing resources. If the *distinguishability* rule is not satisfied, the model cannot produce reasonable output. For example, if we randomly split the train set, there is no way for any model to classify the entities with such a label. Following the two rules, by changing different clustering and classification methods, we propose two batch samplers capturing different aspects of information in the two KGs. Each of them produces a set containing K batches, where K is a hyperparameter. Intuitively, larger K would result in smaller batches, making the normalization process consumes less memory. However, it also makes the ClusterSampler process more difficult. We detail the effect of different K on the accuracy of ClusterSampler in Section 5.4. To distinguish the batches from batches in Section 4.2, the batches are denoted as $\mathcal{B} = \{(\mathcal{B}_s, \mathcal{B}_t)\}$, where $\mathcal{B}_s \subset E_s$ and $\mathcal{B}_t \subset E_t$ are the sets of source and target entities respectively for each batch.

Mini-batch sampling with Intra-KG information. In the ClusterSampler strategy, we first present Intra-KG Structure-based ClusterSampler (ISCS) for sampling batches based on learning neighborhood information of KGs. Following the *distinguishability* rule, for retaining the neighborhood information, the clustering method should put nodes that are neighbors into the same batch. This brings us to minimize the edge cut. A previous study proposes METIS-CPS [12], which clusters two KGs with METIS [17], a classic algorithm to partition large graphs for minimizing the edge cut. METIS-CPS is designed to minimize both the edge cuts of two KGs and the decrease of entity equivalent rate. It first clusters source KG with METIS, and then clusters target KG with higher weights set for train nodes guiding the METIS algorithm. However, the METIS algorithm on the target KG is also a clustering algorithm, thus does not necessarily follow the guidance of train nodes. Following the ClusterSampler strategy, ISCS also adopts METIS for *clustering* source KG, but trains a GNN [18] for *classifying* target KG. The labels provided by METIS will keep nodes that neighbor together as much as possible, and thus can be learned with neighborhood propagation of GNNs. The training and inference on target KG could follow standard node classification task settings. For scalability consideration, we adopt the classic GCN as the classification model. We use the learned embeddings \mathbf{f} of the EA model as the input feature and cross-entropy as the learning loss to train a two-layer GCN. Since GCN does not recognize different relation types in KGs, we adopt the computation of weights in the adjacency matrix from GCNAlign [37] to convert triples with different relation types into different edge weights. We detail the adjacency matrix construction in Appendix B.

Mini-batch sampling with Inter-KG information. Recall that the cross-KG mapping information is learned into two sets of embeddings \mathbf{f}_s and \mathbf{f}_t , partitioning based on these embeddings could preserve the mapping information as much as possible. We propose Cross-KG Mapping-based ClusterSampler (CMCS) for partitioning directly based on the embedding vectors. For obtaining labels of training sets, in *clustering* process of CMCS, we adopt the K-Means algorithm, a widely-used and scalable approach to cluster embeddings in high-dimensional vector spaces. K-Means may lead to the

entities unevenly distributed in different batches. To reduce this effect and to obtain more balanced mini-batches, we normalize the entity features with the standard score normalization. We concatenate the normalized embeddings of the training set into one unified set of embeddings. Then, we cluster the embeddings to obtain the labeled batch number for the training set C' . Formally, $C'_{e_s} = C'_{e_t} = \text{k-Means}(\mathbf{f}_n(e_s, e_t), K)$, $(e_s, e_t) \in \phi'$, $\mathbf{f}_n(e_s, e_t) = [\text{z-score}(\mathbf{f}_{e_s}) || \text{z-score}(\mathbf{f}_{e_t})]$, where $\text{z-score}(X) = \frac{X - \mu(X)}{\sigma(X)}$ is the standard score normalization.

Next, we use the label to train two classifiers for both E_s and E_t , and predict on all the embeddings. We describe how to obtain the batch number of each entity as follows: $C_{E_s} = \{\text{clf}(\mathbf{f}_{\phi'_s}, C'_{\phi'_s}, \mathbf{f}_{e_s}) \mid e_s \in E_s\}$ and $C_{E_t} = \{\text{clf}(\mathbf{f}_{\phi'_t}, C'_{\phi'_t}, \mathbf{f}_{e_t}) \mid e_t \in E_t\}$, where $\text{clf}(\mathbf{f}_{\text{train}}, C_{\text{train}}, \mathbf{f})$ denotes the classification model. It trains based on the training set embedding $\mathbf{f}_{\text{train}}$ and label C_{train} , and then, it predicts the class for embedding \mathbf{f} . In this paper, we use XGBoost Classifier [7], a scalable classifier that has been a golden standard for various data science tasks. After classification, we can easily obtain the batches with the label C .

4.4 Fusing Local and Global Similarities

Since the ClusterSampler strategy utilizes different aspects of information to learn the mini-batches, the mini-batch similarity matrices generated may be biased by the corresponding batch sampler. For example, CMCS only relies on the embeddings, tending to put entities with similar embeddings together. This information bias may have a negative effect on the final accuracy. To avoid such bias as much as possible, we propose SparseFusion. It first applies Sinkhorn iteration on mini-batch similarity matrices generated by multiple batch samplers. Then, SparseFusion sums all the similarity matrices of generated batches to obtain a fused local similarity matrix. Finally, it further fuses the local similarity matrix with a partially normalized global similarity based on a newly proposed sparse version of CSLS [19], namely, Sp-CSLS.

Local Similarity Matrix Normalization. Previous section describes how to sample the input KGs into multiple batches. For each batch generated from a batch-sampler $\mathcal{B}^i \in \mathcal{B} = \{\mathcal{B}_s^i, \mathcal{B}_t^i\}$, $i \in K$, we assume that there exists 1-to-1 mapping between the source and target entities. We first obtain the local similarity matrix $\mathcal{M}^i \in \mathcal{R}^{|E_s| \times |E_t|}$ of current batch. Formally,

$$\mathcal{M}_{e_s, e_t}^i = \begin{cases} \text{sim}(e_s, e_t) & \text{if } e_s \in \mathcal{B}_s^i \text{ and } e_t \in \mathcal{B}_t^i \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

where $\text{sim}(e_s, e_t) = \mathbf{h}_{e_s} \cdot \mathbf{h}_{e_t}$ is the similarity of two entities obtained with the GNN output feature.

Then, we follow [9] to implement the Sinkhorn iteration. We iteratively normalize the similarity matrix K_s rounds in each batch, converting the similarity matrix into a doubly stochastic matrix. The entities of mini-batches in one graph do not have overlap with each other, meaning that there will be no overlapped values in all the mini-batch similarities. Therefore, to obtain the locally normalized similarity for the whole dataset, we directly sum up all the mini-batch similarities, denoted as $\mathcal{M} = \sum_{i \in K} \text{Sinkhorn}(\mathcal{M}^i, K_s) \in [0, 1]^{|E_s| \times |E_t|}$.

Fusing multi-aspect local similarities. To avoid the bias from one batch sampler, we calculate multiple similarity matrices as described above with different batch samplers. We obtain the cross-KG information-based similarity \mathcal{M}_C with CMCS, and intra-KG information-based similarity \mathcal{M}_I with ISCS. Since the ISCS process is unidirectional, indicating that this process on $G_s \rightarrow G_t$ produces different result with $G_t \rightarrow G_s$. According to this characteristic, we apply ISCS on both direction, resulting into two matrices $\mathcal{M}_{I,G_s \rightarrow G_t}$ and $\mathcal{M}_{I,G_t \rightarrow G_s}$. Following [41], we sum up all the similarity matrices without setting any weight to obtain the final local similarity matrix. Formally, $\mathcal{M}_L = \mathcal{M}_C + \mathcal{M}_{I,G_s \rightarrow G_t} + \mathcal{M}_{I,G_t \rightarrow G_s}^T$. This simple approach of fusing multi-view similarity matrices is proved to be useful in various previous studies [11, 12, 41, 42].

Normalize global similarity with Sp-CSLS. To fuse the normalized local similarity matrix \mathcal{M}_L with the global similarity, we first obtain the global similarity, and normalize it partially. A widely-used normalization approach for solving geometric problems is to apply CSLS [19] on the similarity matrix. Formally, for two entities e_s and e_t , $\text{CSLS}(e_s, e_t) = 2 \text{sim}(e_s, e_t) - r_S(e_t) - r_T(e_s)$, where r_S and r_T are the nearest neighborhood similarity mean, which can be obtained by k-NN search with K_n as the neighborhood size.

However, CSLS also lack scalability, which motivates us to propose a sparse version of CSLS, i.e., Sp-CSLS. Recent studies [12, 28] apply FAISS [16] to compute K-nearest neighbours of E_s on the embedding space \mathbf{h}_t . Similar to the dense version of CSLS, Sp-CSLS normalizes a sparse similarity matrix, resulting in a partially normalized similarity matrix. It first uses FAISS to calculate mean neighborhood similarities $r_T(x_s)$ and $r_S(y_t)$. Then, for a sparse matrix \mathcal{M} , the Sp-CSLS only subtracts nonzero values of $2\mathcal{M}$ with $r_T(x_s)$ and $r_S(y_t)$, the result is denoted as \mathcal{M}' . To keep the nonzero values to be useful, the final output is normalized with min-max normalization. Formally, $\text{Sp-CSLS}(\mathcal{M}) = \frac{\mathcal{M}' - \min(\mathcal{M}')}{\max(\mathcal{M}') - \min(\mathcal{M}')} \cdot$

To obtain the normalized global similarity. We first utilize FAISS to obtain an initial global similarity matrix by fusing k-NN similarity matrices on both $\mathbf{f}_s \rightarrow \mathbf{f}_t$ and $\mathbf{f}_t \rightarrow \mathbf{f}_s$ directions. Next, we normalize it with Sp-CSLS. Formally, $\mathcal{M}_G = \text{Sp-CSLS}(\text{k-NN}(\mathbf{h}_s, \mathbf{h}_t, K_r) + \text{k-NN}(\mathbf{h}_t, \mathbf{h}_s, K_r)^T, K_n)$, where $\text{k-NN}(\mathbf{X}, \mathbf{Y}, K_r)$ returns the similarity matrix of $\mathbf{X} \rightarrow \mathbf{Y}$ remaining only top- K_r values, and K_n is the neighborhood size for CSLS. Note that the sparse similarity of two directions may have some values overlapped, this overlapped values will become twice the value before. However, it will hardly have negative impact on accuracy since entity pairs that both have high ranking on the another KG are more possible to be correctly aligned [25].

Fusing normalized similarities. To obtain the fused final similarity matrix, we first fuse the global and local similarities, and then apply the Sp-CSLS for further normalization on the final matrix. Formally, $\mathcal{M}_F = \text{Sp-CSLS}(\mathcal{M}_L + \mathcal{M}_G, K_n)$.

5 EXPERIMENTS

In this section, we report on extensive experiments aimed at evaluating the performance of ClusterEA.

5.1 Experimental Settings

Datasets. We conduct experiments on datasets with different sizes from two cross-lingual EA benchmarks, i.e., IDS [34] and DBP1M [12].

- *IDS* contains four cross-lingual datasets, i.e., English and French (IDS15K_{EN-FR} and IDS100K_{EN-FR}), and English and German (IDS15K_{EN-DE} and IDS100K_{EN-DE}). These benchmarks are sampled with consideration of keeping the properties (e.g., degree distribution) consistent with their source KGs. We use the latest 2.0 version of IDS, where the URIs of entities are encoded to avoid possible name bias.
- *DBP1M* is the largest cross-lingual EA benchmark. It contains two large-scale datasets extracted from DBpedia [1], i.e., English and French (DBP1M_{EN-FR}), and English and German (DBP1M_{EN-DE}). However, DBP1M is biased with name information. Specifically, part of the entities in inter-language links (ILLs) does not occur in the two KGs. Thus, we remove those ILLs to solve the name bias issue while retaining all the triples.

Following previous studies, we use 30% of each dataset as seed alignment, and use 70% of it to test the EA performance. As can be seen, we consider both degree distribution issue [13, 34] and the name bias issue [22] when selecting benchmarks, which meets the requirements of real-world applications. Table 4 in Appendix C lists the detailed information of the datasets used in our experiments.

Evaluation metrics. We use the widely-adopted Hits@N (H@N) and Mean Reciprocal Rank (MRR) to verify the accuracy of ClusterEA [6, 11, 20, 23, 26, 37, 47]. Here, for H@N, $N=1, 10$. Higher H@N and MRR indicate better performance. Also, we use running time and maximum GPU Memory cost (*Mem.*) to evaluate the scalability of ClusterEA. Specifically, running time is measured in seconds, and *Mem.* is measured in Gigabytes.

Baselines. We compare ClusterEA with structure-only based methods. If a baseline includes side information components, we remove them in order to guarantee a fair comparison [12, 21, 23, 26, 29, 34]. All the baselines are enhanced with CSLS (Sp-CSLS for large-scale datasets) before evaluation if possible. The implementation details and parameter settings of ClusterEA and all baselines are presented in Appendix D. Considering the scalability of models, we divide the compared baselines into two major categories, as listed below:

- *Non-scalable baselines* that includes (i) *GCNAlign* [37], the first GNN-based EA model; (ii) *RREA* [26], a GNN-based EA model that utilizes relational reflection transformation to obtain relation-specific embeddings for each entity, and is used as the default of LargeEA [12]; and (iii) *Dual-AMN* [23], a SOTA EA model that contains Simplified Relational Attention Layer and Proxy Matching Attention Layer for modeling both intra-graph and cross-graph relations.
- *Scalable baselines* that includes (i) *LargeEA* [12], the first EA framework that focuses on scalability by training multiple EA models on mini-batches generated by a rule-based strategy, and excepting for the two variants presented in [12], we provide a new variant *LargeEA-D* that incorporates recently proposed EA model Dual-AMN; and (ii) *Stochastic training (Section 4.2) variant of GNN models*, which incorporates EA models with Stochastic training, including *GCNAlign-S* for GCNAlign [37], *RREA-S* for RREA [26], and *Dual-AMN-S* for Dual-AMN [23].

Variants of ClusterEA. Since ClusterEA is designed to be integrated with GNN-based EA models, we present three versions of ClusterEA, i.e., ClusterEA-G that includes GCNAlign, ClusterEA-R that incorporates RREA, and ClusterEA-D that incorporates Dual-AMN. Specifically, we treat ClusterEA-D as the default setting.

Table 1: Overall EA results on IDS15K and IDS100K

Methods	IDS15K _{EN-FR}					IDS15K _{EN-DE}					IDS100K _{EN-FR}					IDS100K _{EN-DE}				
	H@1	H@10	MRR	Time	Mem.	H@1	H@10	MRR	Time	Mem.	H@1	H@10	MRR	Time	Mem.	H@1	H@10	MRR	Time	Mem.
GCNAlign	38.2	78.5	0.51	10.90	0.13	58.7	85.5	0.67	12.27	0.13	29.9	61.7	0.40	71.37	1.00	41.0	66.1	0.49	79.52	1.00
RREA	63.3	91.4	0.73	136.32	4.07	75.5	94.9	0.82	156.85	4.07	–	–	–	–	–	–	–	–	–	–
Dual-AMN	64.6	91.5	0.74	12.50	4.05	76.5	95.2	0.83	13.72	4.05	49.3	77.5	0.59	413.73	21.91	59.3	81.8	0.67	456.87	22.56
LargeEA-G	30.0	63.5	0.39	9.76	0.13	40.3	68.7	0.50	9.94	0.13	19.5	45.7	0.28	36.65	0.50	21.4	39.7	0.27	39.29	0.50
LargeEA-R	47.2	74.0	0.56	41.01	1.01	58.0	77.6	0.65	42.12	1.01	32.5	55.5	0.40	163.92	4.04	27.3	42.5	0.32	157.15	4.04
LargeEA-D	45.7	69.6	0.54	13.51	0.75	58.7	76.2	0.65	13.96	0.75	33.0	54.8	0.40	134.5	3.41	28.6	42.7	0.65	121.9	3.65
GCN-Align-S	37.1	73.8	0.49	24.10	1.11	53.6	83.4	0.63	23.73	1.11	25.3	45.5	0.35	184.43	1.72	35.5	61.5	0.44	189.12	1.72
RREA-S	62.7	90.3	0.72	34.09	4.89	76.3	95.0	0.83	34.23	5.01	46.4	75.4	0.56	250.80	7.16	57.3	80.6	0.65	256.35	8.50
Dual-AMN-S	60.9	88.9	0.71	15.59	5.30	75.0	94.3	0.82	15.64	5.10	48.2	76.6	0.57	122.78	7.79	58.8	81.4	0.66	124.49	8.13
ClusterEA-G	46.6	77.8	0.57	40.99	4.43	62.0	86.3	0.70	40.47	4.43	30.6	57.9	0.40	236.37	2.77	41.4	64.3	0.49	246.90	2.77
ClusterEA-R	67.9	90.0	0.76	52.29	4.89	79.4	94.5	0.85	51.89	5.01	52.0	76.3	0.60	329.78	7.16	62.2	81.7	0.69	339.54	8.50
ClusterEA-D	67.4	89.5	0.75	35.76	5.30	79.5	94.6	0.85	35.82	5.10	54.2	78.1	0.62	210.50	8.52	63.7	82.8	0.70	212.80	8.31

¹ The symbol “–” indicates that the model fails to perform EA on IDS100K dataset due to extensive GPU memory usage.

Table 2: Overall EA results on DBP1M

Methods	DBP1M _{EN-FR}					DBP1M _{EN-DE}				
	H@1	H@10	MRR	Time	Mem.	H@1	H@10	MRR	Time	Mem.
LargeEA-G	5.1	13.4	0.08	463	4.00	3.4	9.5	0.05	378	4.00
LargeEA-R	9.4	21.5	0.13	1681	20.05	6.4	15.0	0.09	1309	20.49
LargeEA-D	10.5	21.9	0.15	2546	19.72	6.6	14.7	0.09	1692	17.71
GCN-Align-S	6.0	19.1	0.10	1817	9.25	4.5	14.4	0.08	1491	7.78
RREA-S	21.1	42.2	0.28	2080	16.86	20.5	40.2	0.27	1754	15.65
Dual-AMN-S	22.4	43.4	0.29	588	15.85	22.0	41.7	0.28	490	14.61
ClusterEA-G	10.0	24.5	0.15	2501	17.43	6.9	17.7	0.11	2027	17.76
ClusterEA-R	26.0	45.6	0.32	3025	21.11	25.0	45.0	0.32	2526	19.68
ClusterEA-D	28.1	47.4	0.35	1647	20.10	28.8	48.8	0.35	1360	18.20

5.2 Overall Results

Performance on IDS. Table 1 summarizes the EA performance on IDS15K and IDS100K. First, ClusterEA improves H@1 by 3.3% – 29.7% compared with the non-scalable methods (viz., GCNAlign, RREA, and Dual-AMN), and by 3.2% – 42.3% compared against the scalable ones. They validate the accuracy of ClusterEA. Moreover, all variants of ClusterEA perform better than the structure-only based models in terms of H@1. It confirms the superiority of the way how we enhance the output, compared with CSLS that is the default setting for most recent EA models [11, 21, 23, 25, 26, 34] to enhance the output similarity matrix. Second, it is observed that the accuracy of LargeEA variants is significantly reduced compared with their corresponding original models. This is because LargeEA discards structure information during its mini-batch training process. Unlike the LargeEA variants, the accuracy of Stochastic Training variants only drops slightly compared to the original non-scalable models. This shows that the Stochastic training process can minimize the structure information loss (cf. Section 4.2). Third, as observed, H@10 of some non-scalable models (RREA, Dual-AMN) is higher than that of ClusterEA. This is mainly due to (i) the information loss during ClusterEA’s Stochastic training and (ii) the incompleteness of global similarity normalization (to be detailed in Section 5.3). Nevertheless, H@1 is the most representative indicator for evaluating the accuracy of EA since higher H@1 directly indicates more correct proportion of aligned entities. Thus, achieving the highest H@1 among all the competitors is sufficiently to demonstrate the high performance of ClusterEA. Finally, the experimental results show that the training of all LargeEA variants is faster than that of other variants of the corresponding models. The reason is that LargeEA omits most information of graph structure, where

Table 3: The result of ablation study

Methods	DBP1M _{EN-FR}			DBP1M _{EN-DE}		
	H@1	H@10	MRR	H@1	H@10	MRR
ClusterEA	28.1	47.4	0.35	28.8	48.8	0.35
ClusterEA - Sp-CSLS	27.9	46.2	0.34	28.4	46.0	0.34
ClusterEA - Global Sim	27.6	48.3	0.34	28.3	48.7	0.35
ClusterEA - Sinkhorn	20.0	46.1	0.29	22.7	46.2	0.30
ClusterEA - CMCS	26.5	46.2	0.33	26.8	47.1	0.33
ClusterEA - ISCS	25.3	45.9	0.32	25.5	46.5	0.32
ClusterEA - Dual-AMN	10.0	24.5	0.15	6.9	17.7	0.11

less data is subjected to the training process. However, in this case, the running time of ClusterEA is still comparable.

Performance on DBP1M. Table 2 reports the EA performance of ClusterEA and its competitors on DBP1M. Note that we do not report the results of the non-scalable models because it is infeasible to perform their training phases on DBP1M due to large GPU memory usage. We observe that the accuracy of all the variants of ClusterEA on DBP1M is higher than those of LargeEA. Specifically, H@1 is improved by 5.7% – 23.0% and 6.8% – 25.4% on DBP1M_{EN-FR} and DBP1M_{EN-DE}, respectively. Next, compared with each Stochastic training variants of the corresponding incorporated model of ClusterEA, ClusterEA brings about 4.0% – 5.7% and 2.4% – 6.8% absolute improvement in H@1 on DBP1M_{EN-FR} and DBP1M_{EN-DE} datasets, respectively. As the expressiveness of the model improves, ClusterEA also offers more improvement. Last, all the scalable variants incorporating GCNAlign variants produce poor EA results. Specifically, GCNAlign’s model cannot be directly applied to large-scale datasets due to its insufficient expressive ability. Note that ClusterEA not only outperforms baselines in terms of accuracy but also achieves comparable performance with them in terms of both Mem. and running time. Overall, ClusterEA is able to scale-up the GNN-based EA models while enhancing H@1 by at most 8× compared with the state-of-the-arts. More experimental results of scalability are provided in Appendix E.

5.3 Ablation Study

In ablation study, we remove each component of ClusterEA, and report H@1, H@10, and MRR in Table 3. First, after removing the Sp-CSLS component, the accuracy of ClusterEA drops. This shows that the Sp-CSLS normalization indeed address the geometric problems of global similarity on large-scale EA. Second, after removing the global similarity, H@1 of ClusterEA drops but H@10 grows on

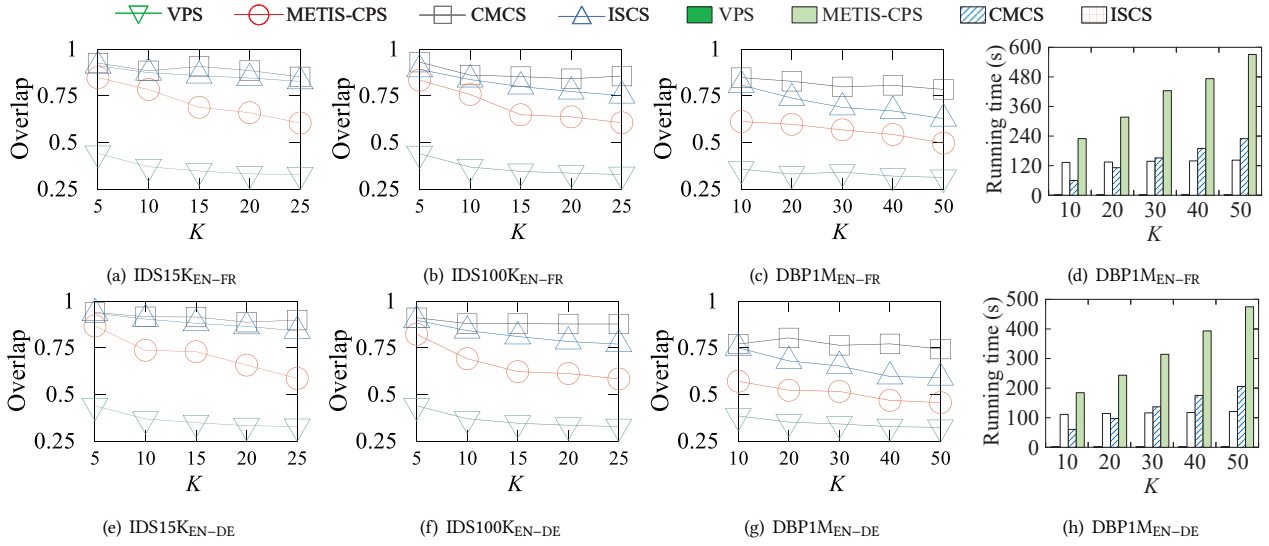


Figure 3: Comparison results of different batch sampling strategies

DBP1M_{EN-FR}. The fluctuation on H@10 may be due to the incompleteness of global similarity normalization that disturbs the final similarity matrix. Specifically, the local similarity is normalized into nearly permutation matrices with Sinkhorn iteration, where most values of one row are close to zero. When fusing local and global similarity matrices, elements that are not top-1 in the local matrix will be biased towards the value of the global matrix, which is only partially normalized. This causes the disturbance. However, the incompleteness of normalization does not degrade H@1. This is because the value of one row in M_L that is correctly aligned will be normalized into a higher value, providing resistance to the global matrix disturbance. Third, after removing the Sinkhorn iteration, the accuracy of ClusterEA drops significantly. This confirms the importance of normalizing the similarity matrices. Finally, after removing each sampler of ClusterSampler, the accuracy of ClusterEA drops on all metrics. This validates the importance of fusing information from multi-aspects, including cross-KG information and intra-KG information. In addition, we observe that CMCS has less influence compared with ISCS. This is mainly due to the following two reasons. First, CMCS generally clusters entities with similar embedding vectors, where each batch still suffers from geometric problems. On the contrary, ISCS samples batches based on the graph neighborhood information, which is a strong constriction on the learning model. Second, we apply ISCS in both directions. In this case, it can capture more information than CMCS. By replacing its Dual-AMN model with the GCNAlign model, the accuracy of ClusterEA drops significantly on all metrics. This demonstrates the importance of the ability of the incorporated EA model in ClusterEA.

5.4 Case Study: ClusterSampler Analysis

The ClusterSampler is a vital component of ClusterEA. To ensure scalability, we set the batch number K such that the space cost of the normalization process does not exceed the GPU memory. We also need to guarantee that our batch sampling method can produce acceptable mini-batches under different K settings. Thus, we provide a detailed analysis on varying the batch number K for

different batch samplers from ClusterSampler (i.e., ISCS and CMCS). Specifically, we study how much are the mini-batches generated by one sampler acceptable as the percentage of equivalent entities that are placed into the same mini-batches (denoted as *Overlap*). Next, we report the *Overlap* metric and running time of the proposed sampler, where the mini-batch number K of the proposed sampler is varied from 5 to 25 on IDS and is varied from 10 to 50 on DBP1M. We compare the proposed sampler with two rule-based baselines, VPS and METIS-CPS. VPS randomly partitions seed alignments and all other entities into different mini-batches. METIS-CPS sets the training entities with higher nodes to sample better mini-batches. Note that both METIS-CPS and ISCS are uni-directional. Thus, we apply these methods in both directions, and present their average performance of the two directions. We report the *Overlap* of all sampling methods on all benchmarks in Fig 3, where Figures 3(a), 3(b), 3(c), 3(e), 3(f), and 3(g) are the *Overlap* of different datasets. The results show that CMCS generally outperforms two baselines on all the datasets, and its performance is stable when varying K . However, although it is better overlapped, it may incur hubness and isolation problems in mini-batches (cf. Section 5.3). Thus, fusing multi-aspect is essential for ClusterEA. ISCS results in less overlapped mini-batches while still much better than METIS-CPS. This is because METIS does not necessarily follow the guidance provided by METIS-CPS. The GNN node classification model in ISCS, which considers the cross-entropy loss as a penalty, is forced to learn mini-batches more effectively. Since we set the training ratio to 30%, all samplers have an *Overlap* over 30%, including VPS that splits mini-batches randomly. Finally, we report the running time on DBP1M in Figures 3(d) and 3(h). Since all samplers achieve sufficiently high efficiency on IDS datasets, we do not present the running time of all samplers on IDS due to space limitation. It is observed that, although the result is unacceptable, VPS is the fastest sampling method. Both the proposed ISCS and CMCS are always about 2× faster than the rule-based METIS-CPS when changing K . The reason is that both of them utilize machine learning models that can be easily accelerated with GPU.

6 CONCLUSIONS

We present ClusterEA to align entities between large-scale knowledge graphs with stochastic training and normalized similarities. ClusterEA contains three components, including stochastic training, ClusterSampler, and SparseFusion, to perform large-scale EA task based solely on structure information. We first train a large-scale Siamese GNN for EA in a stochastic fashion to produce entity embeddings. Next, we propose a new ClusterSampler strategy for sampling highly overlapped mini-batches taking advantage of the trained embeddings. Finally, we present a SparseFusion process, which first normalizes local and global similarities and then fuses them to obtain the final similarity matrix. The whole process of ClusterEA guarantees both high accuracy and comparable scalability. Considerable experimental results on EA benchmarks with different scales demonstrate that ClusterEA significantly outperforms previous large-scale EA study. In future, it is of interest to explore dangling settings [28] of EA on large scale datasets.

7 ACKNOWLEDGEMENTS

This work was supported in part by the National Key Research and Development Program of China under Grant No. 2021YFC3300303, the NSFC under Grants No. (62025206, 61972338, and 62102351), and the Zhejiang Provincial Natural Science Foundation under Grant No. LR21F020005. Lu Chen is the corresponding author of the work.

REFERENCES

- [1] Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary G. Ives. 2007. DBpedia: A Nucleus for a Web of Open Data. In *ISWC*. 722–735.
- [2] Fabio Azzalini, Songle Jin, Marco Renzi, and Letizia Tanca. 2021. Blocking Techniques for Entity Linkage: A Semantics-Based Approach. *Data Sci. Eng.* 6, 1 (2021), 20–38.
- [3] Antoine Bordes, Nicolas Usunier, Alberto García-Durán, Jason Weston, and Oksana Yakhnenko. 2013. Translating Embeddings for Modeling Multi-relational Data. In *NIPS*. 2787–2795.
- [4] Yixin Cao, Zhiyuan Liu, Chengjiang Li, Zhiyuan Liu, Juanzi Li, and Tat-Seng Chua. 2019. Multi-Channel Graph Neural Network for Entity Alignment. In *ACL*. 1452–1461.
- [5] Muhao Chen, Weijia Shi, Ben Zhou, and Dan Roth. 2021. Cross-lingual entity alignment with incidental supervision. (2021), 645–658.
- [6] Muhao Chen, Yingtao Tian, Mohan Yang, and Carlo Zaniolo. 2017. Multilingual Knowledge Graph Embeddings for Cross-lingual Knowledge Alignment. In *IJCAI*. 1511–1517.
- [7] Tianqi Chen and Carlos Guestrin. 2016. XGBoost: A Scalable Tree Boosting System. In *KDD*. 785–794.
- [8] Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. 2019. Cluster-GCN: An Efficient Algorithm for Training Deep and Large Graph Convolutional Networks. In *KDD*. 257–266.
- [9] Marco Cuturi. 2013. Sinkhorn Distances: Lightspeed Computation of Optimal Transport. In *NIPS*. 2292–2300.
- [10] Matthias Fey, Jan Eric Lenssen, Christopher Morris, Jonathan Masci, and Nils M. Kriege. 2020. Deep Graph Matching Consensus. In *ICLR*.
- [11] Congcong Ge, Xiaozhe Liu, Lu Chen, Baihua Zheng, and Yunjun Gao. 2021. Make It Easy: An Effective End-to-End Entity Alignment Framework. In *SIGIR*. 777–786.
- [12] Congcong Ge, Xiaozhe Liu, Lu Chen, Baihua Zheng, and Yunjun Gao. 2022. LargeEA: Aligning Entities for Large-scale Knowledge Graphs. *PVLDB* 15, 2 (2022), 237–245.
- [13] Lingbing Guo, Zequn Sun, and Wei Hu. 2019. Learning to Exploit Long-term Relational Dependencies in Knowledge Graphs. In *ICML*. 2505–2514.
- [14] William L Hamilton, Rex Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *NIPS*. 1025–1035.
- [15] Harold W. Kuhn. 1955. The Hungarian method for the assignment problem. *Naval Research Logistics Quarterly* 2, 1, 83–97.
- [16] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2017. Billion-scale similarity search with GPUs. *arXiv preprint arXiv:1702.08734* (2017).
- [17] George Karypis and Vipin Kumar. 1998. A Fast and High Quality Multilevel Scheme for Partitioning Irregular Graphs. *SIAM J. Sci. Comput.* 20, 1 (1998), 359–392.
- [18] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *ICLR*.
- [19] Guillaume Lample, Alexis Conneau, Marc'Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2018. Word translation without parallel data. In *ICLR*.
- [20] Chengjiang Li, Yixin Cao, Lei Hou, Jiaxin Shi, Juanzi Li, and Tat-Seng Chua. 2019. Semi-supervised Entity Alignment via Joint Knowledge Embedding Model and Cross-graph Model. In *EMNLP*. 2723–2732.
- [21] Fangyu Liu, Muhao Chen, Dan Roth, and Nigel Collier. 2020. Visual Pivoting for (Unsupervised) Entity Alignment. *arXiv preprint arXiv:2009.13603* (2020).
- [22] Zhiyuan Liu, Yixin Cao, Liangming Pan, Juanzi Li, and Tat-Seng Chua. 2020. Exploring and Evaluating Attributes, Values, and Structures for Entity Alignment. In *EMNLP*. 6355–6364.
- [23] Xin Mao, Wenting Wang, Yuanbin Wu, and Man Lan. 2021. Boosting the Speed of Entity Alignment 10x: Dual Attention Matching Network with Normalized Hard Sample Mining. In *WWW*. 821–832.
- [24] Xin Mao, Wenting Wang, Yuanbin Wu, and Man Lan. 2021. From Alignment to Assignment: Frustratingly Simple Unsupervised Entity Alignment. In *EMNLP*. 2843–2853.
- [25] Xin Mao, Wenting Wang, Huimin Xu, Man Lan, and Yuanbin Wu. 2020. MRAEA: An Efficient and Robust Entity Alignment Approach for Cross-lingual Knowledge Graph. In *WSDM*. 420–428.
- [26] Xin Mao, Wenting Wang, Huimin Xu, Yuanbin Wu, and Man Lan. 2020. Relational Reflection Entity Alignment. In *CIKM*. 1095–1104.
- [27] Alvin E Roth. 2008. Deferred acceptance algorithms: History, theory, practice, and open questions. *international journal of game Theory* 36, 3 (2008), 537–569.
- [28] Zequn Sun, Muhao Chen, and Wei Hu. 2021. Knowing the No-match: Entity Alignment with Dangling Cases. In *ACL*. 3582–3593.
- [29] Zequn Sun, Muhao Chen, Wei Hu, Chengming Wang, Jian Dai, and Wei Zhang. 2020. Knowledge Association with Hyperbolic Knowledge Graph Embeddings. In *EMNLP*. 5704–5716.
- [30] Zequn Sun, Wei Hu, and Chengkai Li. 2017. Cross-Lingual Entity Alignment via Joint Attribute-Preserving Embedding. In *ISWC*. 628–644.
- [31] Zequn Sun, Wei Hu, Qingheng Zhang, and Yuzhong Qu. 2018. Bootstrapping Entity Alignment with Knowledge Graph Embedding. In *IJCAI*. 4396–4402.
- [32] Zequn Sun, JiaCheng Huang, Wei Hu, Muhao Chen, Lingbing Guo, and Yuzhong Qu. 2019. TransEdge: Translating Relation-Contextualized Embeddings for Knowledge Graphs. In *ISWC*. 612–629.
- [33] Zequn Sun, Chengming Wang, Wei Hu, Muhao Chen, Jian Dai, Wei Zhang, and Yuzhong Qu. 2020. Knowledge Graph Alignment Network with Gated Multi-Hop Neighborhood Aggregation. In *AAAI*. 222–229.
- [34] Zequn Sun, Qingheng Zhang, Wei Hu, Chengming Wang, Muhao Chen, Farahnaz Akrami, and Chengkai Li. 2020. A Benchmarking Study of Embedding-based Entity Alignment for Knowledge Graphs. *PVLDB* 13, 11 (2020), 2326–2340.
- [35] Xiaobin Tang, Jing Zhang, Bo Chen, Yang Yang, Hong Chen, and Cuiping Li. 2020. BERT-INT: A BERT-based Interaction Model For Knowledge Graph Alignment. In *IJCAI*. 3174–3180.
- [36] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. In *ICLR*.
- [37] Zhichun Wang, Qingsong Lv, Xiaohan Lan, and Yu Zhang. 2018. Cross-lingual Knowledge Graph Alignment via Graph Convolutional Networks. In *EMNLP*. 349–357.
- [38] Zhichun Wang, Jinjian Yang, and Xiaojun Ye. 2020. Knowledge Graph Alignment with Entity-Pair Embedding. In *EMNLP*. 1672–1680.
- [39] Chenyan Xiong, Russell Power, and Jamie Callan. 2017. Explicit Semantic Ranking for Academic Search via Knowledge Graph Embedding. In *WWW*. 1271–1279.
- [40] Kai Yang, Shaoqin Liu, Junfeng Zhao, Yasha Wang, and Bing Xie. 2020. COTSAE: CO-Training of Structure and Attribute Embeddings for Entity Alignment. In *AAAI*. 3025–3032.
- [41] Weixin Zeng, Xiang Zhao, Jiuyang Tang, and Xuemin Lin. 2020. Collective Entity Alignment via Adaptive Features. In *ICDE*. 1870–1873.
- [42] Weixin Zeng, Xiang Zhao, Jiuyang Tang, Xuemin Lin, and Paul Groth. 2021. Reinforcement Learning-based Collective Entity Alignment with Adaptive Features. *TOIS* 39, 3 (2021), 26:1–26:31.
- [43] Weixin Zeng, Xiang Zhao, Wei Wang, Jiuyang Tang, and Zhen Tan. 2020. Degree-Aware Alignment for Entities in Tail. In *SIGIR*. 811–820.
- [44] Fuzheng Zhang, Nicholas Jing Yuan, Defu Lian, Xing Xie, and Wei-Ying Ma. 2016. Collaborative Knowledge Base Embedding for Recommender Systems. In *KDD*. 353–362.
- [45] Qingheng Zhang, Zequn Sun, Wei Hu, Muhao Chen, Lingbing Guo, and Yuzhong Qu. 2019. Multi-view Knowledge Graph Embedding for Entity Alignment. In *IJCAI*. 5429–5435.
- [46] Xiang Zhao, Weixin Zeng, Jiuyang Tang, Xinyi Li, Minnan Luo, and Qinghua Zheng. 2022. Toward Entity Alignment in the Open World: An Unsupervised Approach with Confidence Modeling. *Data Sci. Eng.* 7, 1 (2022), 16–29.
- [47] Hao Zhu, Ruobing Xie, Zhiyuan Liu, and Maosong Sun. 2017. Iterative Entity Alignment via Joint Knowledge Embeddings. In *IJCAI*. 4258–4264.

Table 4: Statistics of IDS15K, IDS100K, and DBP1M

Datasets		#Entities	#Relations	#Triples
IDS15K	EN-FR	15,000-15,000	267-210	47,334-40,864
	EN-DE	15,000-15,000	215-131	47,676-50,419
IDS100K	EN-FR	100,000-100,000	400-300	309,607-258,285
	EN-DE	100,000-100,000	381-196	335,359-336,240
DBP1M	EN-FR	1,877,793-1,365,118	603-380	7,031,172-2,997,457
	EN-DE	1,625,999-1,112,970	597-241	6,213,639-1,994,876

A NORMALIZED HARD SAMPLE MINING LOSS

We detail the NHSM loss used in our work here. Formally, for the current mini-batch, the training loss is defined as

$$L = \text{LogSumExp}(\lambda z(e_s, e_t)) + \text{LogSumExp}(\lambda z(e_t, e_s)) \quad (3)$$

, where $e_s \in B_s$, $e_t \in B_t$, $(e_s, e_t) \in \phi'$, and $\text{LogSumExp}(X) = \log(\sum_{x \in X} e^x)$ is an operator to smoothly generate hard negative samples, λ is the smooth factor of LogSumExp , and $z \in \mathcal{R}^{|B_t|}$ is the normalized triple loss, defined as

$$z(e_s, e_t) = z\text{-score}(\{y + \text{sim}(e_s, e_t) - \text{sim}(e_s, e'_t) | e'_t \in B_t\}), \quad (4)$$

where $z\text{-score}(X) = \frac{X - \mu(X)}{\sigma(X)}$ is the standard score normalization, and $\text{sim}(e_s, e_t) = \mathbf{h}_{e_s} \cdot \mathbf{h}_{e_t}$ is the similarity of two entities obtained with the GNN output feature.

B ADJACENCY MATRIX CONSTRUCTION FOR ISCS

Following [37], to construct the adjacency matrix $A \in \mathbb{R}^{|E| \times |E|}$, we compute two metrics, which are called functionality and inverse functionality, for each relation:

$$\begin{aligned} \text{fun}(r) &= \frac{\#Head_Entities_of_r}{\#Triples_of_r} \\ \text{ifun}(r) &= \frac{\#Tail_Entities_of_r}{\#Triples_of_r} \end{aligned} \quad (5)$$

where $\#Triples_of_r$ is the number of triples of relation r , $\#Head_Entities_of_r$ is the number of head entities of r , and $\#Tail_Entities_of_r$ is the number of tail entities of r . To measure the influence of the i -th entity over the j -th entity, we set $a_{ij} \in A$ as:

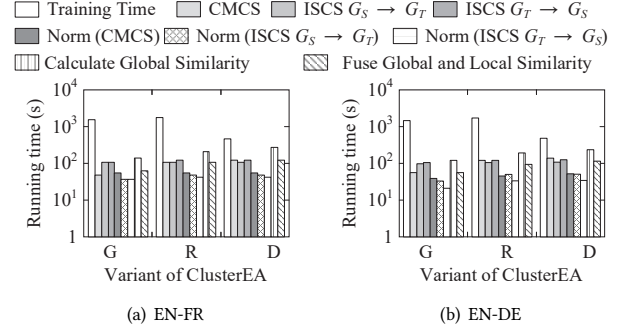
$$a_{ij} = \sum_{\langle e_i, r, e_j \rangle \in G} \text{ifun}(r) + \sum_{\langle e_j, r, e_i \rangle \in G} \text{fun}(r) \quad (6)$$

C STATISTICS OF DATASETS

D IMPLEMENTATION DETAILS

The results of all the baselines are obtained by our re-implementation with their publicly available source codes. All experiments were conducted on a personal computer with an Intel Core i9-10900K CPU, an NVIDIA GeForce RTX3090 GPU, and 128GB memory. The programs were all implemented in Python. Since the deep learning frameworks used vary on different implementations of EA models, we use the NVIDIA Nsight Systems² to record the GPU memory

²<https://developer.nvidia.com/nsight-systems>

**Figure 4: Scalability analysis vs. variants of ClusterEA**

usage of all approaches uniformly. We detail the hyper-parameters used in the experiment as follows. All the hyper-parameters are set without special instructions.

D.1 Non-scalable methods

For non-scalable methods, including GCNAlign, RREA, and Dual-AMN, we follow the parameter settings reported in their original papers [23, 26, 37].

D.1.1 GCNAlign. Note that GCNAlign contains an attribute encoder as the side information component. We re-implement the GCNAlign model by removing attribute features from it.

D.2 LargeEA

We reproduce LargeEA by removing entirely the name channel. It is worth noting that LargeEA incorporates a name-based data augmentation process for generating alignment seeds other than training data. For fair comparison, we also remove this augmentation process, making the LargeEA framework not use any alignment signals apart from the training data. To be more specific, we only compare the structure-channel of LargeEA. The running time and maximum GPU memory usage are also recorded only for the structure-channel of LargeEA, including METIS-CPS partitioning and mini-batch training. We use the default settings reported in [12] of the two versions, LargeEA-G and LargeEA-R. For the newly proposed variant LargeEA-D, we keep all the hyper-parameters unchanged except the structure-based EA model-related parameters switched to the default settings of Dual-AMN [23].

D.3 Stochastic training variant of non-scalable methods

Contributed by the NHSM loss [23], the training epoch number could be greatly reduced. We replace the training epoch for GCNAlign-S and RREA-S from 2000 and 1200 to 50. We follow [23] to train 20 epochs for Dual-AMN-S on IDS datasets. We further notice that the training loss of Dual-AMN-S is stable after 10 epochs for the large-scale dataset DBP1M, thus setting the training epoch of Dual-AMN-S to 10 for DBP1M. We set the fan out number in neighborhood sampling $F = 8$, the batch size $N_p = 2000$, $N_n = 4000$, and Adam as the training optimizer for all variants. We follow the setup of [23] for NHSM loss. For other hyper-parameter settings of

the model, including the embedding dimension D , we follow their original papers [23, 26, 37].

D.4 ClusterEA

In the *stochastic training process*, we adopt the aforementioned settings for each ClusterEA variant. In the *ClusterSampler process*, we set the mini-batch number $K = 5$ for IDS15K, $K = 10$ for IDS100K, and $K = 30$ for DBP1M. Moreover, we utilize a cache in disk for pre-processed edge information to quickly load the constructed matrix of relation-aware GNN variants. For CMCS, we set the max iteration number to 300, the tolerance to 10^{-4} , and distance metric to euclidean distance for the KMeans algorithm. We adopt the default setting for the XGBoost Classifier [7]. We utilize GPU acceleration for those methods. For ISCS, we adopt the default setting of METIS [17], and train a two layer GCN with Adam. The learning rate is set to 0.01, and the training epoch of GCN set to 800 for IDS15K, 1500 for IDS100K, and 3000 for DBP1M. In the *SparseFusion process*, we set $K_s = 100$ for the iteration round of Sinkhorn, $K_r = 50$ for top-K similarity serach, and $K_n = 10$ for CSLs.

E SCALABILITY EVALUATION

To further investigate the scalability of our proposed ClusterEA framework, we verify the running time of each components in different variants of ClusterEA on the DBP1M dataset, including (1) *Training time* of the EA model; (2) *CMCS batch sampling time*; (3) *ISCS batch sampling time* on both directions; (4) *Local Similarity Normalization time* of the batches generated by each samplers, denoted as Norm(Sampler); (5) *Calculating Global Similarity* time in SparseFusion; and (6) *Fusing Global and Local Similarity* time in SparseFusion. The experimental results are shown in Figure 4, where G, R, and D denote ClusterEA-G, ClusterEA-R, and ClusterEA-D, respectively. As depicted in Figure 4, we observe that the running time of each component except for Training time does not exceed 10^3 seconds on the large-scale dataset. This confirms the scalability of ClusterEA. Note that the training time of the Dual-AMN model is significantly less than that of other variants. This is contributed by the model design of Dual-AMN that could dramatically reduce the required training epochs.