

1 Article

2 **Collective List-only Entity Linking: A**
3 **Graph-based Approach**4 Weixin Zeng ¹ , Xiang Zhao ^{1,2,*}  and Jiuyang Tang ^{1,2}5 ¹ National University of Defense Technology, China6 ² Collaborative Innovation Center of Geospatial Technology, China

7 * Correspondence: xiangzhao@nudt.edu.cn; Tel.: +86-731-84576563

8 **Abstract:** List-only entity linking is the task of mapping ambiguous mentions in texts to target
9 entities in a group of entity lists. Different from traditional entity linking task, which leverages rich
10 semantic relatedness in knowledge bases to improve linking accuracy, list-only entity linking can
11 merely take advantage of co-occurrences information in entity lists. State-of-the-art work utilizes
12 co-occurrences information to enrich entity descriptions, which are further used to calculate local
13 compatibility between mentions and entities to determine results. Nonetheless, entity coherence is also
14 deemed to play an important part in entity linking, which is yet currently neglected. In this work, in
15 addition to local compatibility, we take into account global coherence among entities. Specifically, we
16 propose to harness co-occurrences in entity lists for mining both explicit and implicit entity relations.
17 The relations are then integrated into an entity graph, on which Personalized PageRank is incorporated
18 to compute entity coherence. The final results are derived by combining local mention-entity similarity
19 and global entity coherence. The experimental studies validate the superiority of our method. Our
20 proposal not only improves the performance of list-only entity linking, but also opens up the bridge
21 between list-only entity linking and conventional entity linking solutions.

22 **Keywords:** list-only entity linking; named entity disambiguation; graph-based approach23 **1. Introduction**

24 Entity Linking (EL) is the task of detecting corresponding named *entities* for ambiguous *mentions*
25 in text. Mention refers to character string, such as *Jackson* in the example in Figure 1, the true
26 meaning of which needs to be determined by being linked to an entity, such as the basketball coach *Phil*
27 *Jackson*. Traditional EL methods leverage *knowledge bases* (KBs), which offer rich semantic information
28 of entities, for robust and accurate disambiguation process. Nevertheless, despite the effectiveness of
29 knowledge-based EL, it might not be applicable in situations where there is insufficient information of
30 entities, such as *entity lists*.

31 Entity list, as is often the case, consists of a group of closely-related entities, and it exists in
32 various information sources [1]. In contrast to KBs, where complete structure of entities facilitates
33 almost all entity-related tasks, entity list minimizes necessary information to mere co-occurrences of
34 interrelated entities, thus serving as a light-weight alternative in terms of describing entity correlations.

35 Entity lists can be found useful, for instance, in the scenario concerning detection of emerging
36 stock names. When investors search new stock names in Wikipedia ¹, a frequently updated KB, chances
37 are that there are no corresponding items. In fact, as is shown in [2], for a dataset including 2,468
38 stock names, merely 340 of them can be found in Wikipedia. Nevertheless, those stocks can be found
39 co-occurring with others in stock lists on financial websites. Thus, the stock lists will be of great use if
40 people have doubts concerning new stocks. There are much more similar situations, such as searching
41 for specific car brand or collecting information about bars in a small town, where the knowledge about
42 target entities are sparse.

¹ <https://www.wikipedia.org/>

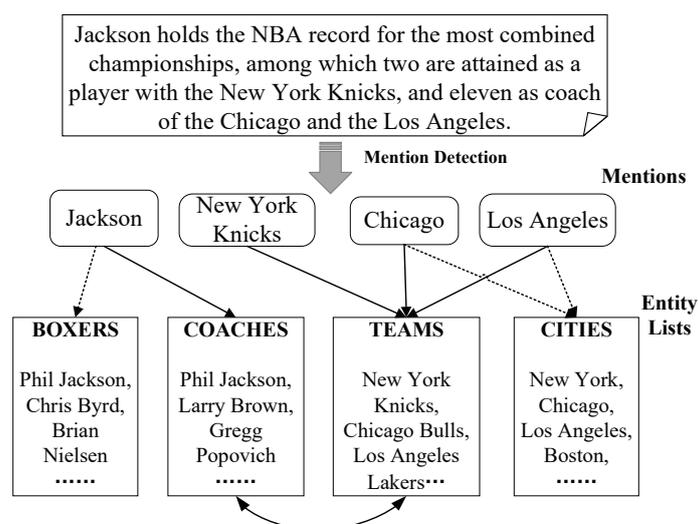


Figure 1. Example of list-only EL

Consequently, the demand for list-only EL emerges [1], which targets at solving the problem of mapping ambiguous mentions to entity lists (rather than KBs); Figure 1 describes an example of list-only EL problem. State-of-the-art method [1] addresses the challenge by merely considering the local compatibilities between mentions and entities to determine matching pairs, whereas neglecting the global coherence among entities.

Example 1. As shown in Figure 1, there is a piece of text with mentions Jackson, New York Nicks, Chicago and Los Angeles; and there are 4 sample entity lists to be linked to, namely *Boxers*, *Coaches*, *Teams* and *Cities*. The task of list-only EL is to link mentions to correct entities in the entity lists. It can be seen that entity Chicago and entity Los Angeles in entity list *Cities* have the same name strings with mention Chicago and mention Los Angeles in the text. Because of the high mention-entity compatibility, existing method tends to map mentions Chicago and Los Angeles to entities Chicago and Los Angeles in the entity list featured *Cities*. However, the true entities for them are Chicago Bulls and Los Angeles Lakers in entity list *Teams*. Furthermore, it is hard for current method to decide which entity that mention Jackson should be linked to, since there are two possible candidate entities with the same name Phil Jackson and they are in different entity lists.

Moreover, we observe from the dataset currently used for empirical study that each document only contains one mention for disambiguation, which may not reflect the reality well. A pragmatic scenario may look like the example in Figure 1, where there are four mentions to be disambiguated. Therefore, we contend that the existing evaluations could be inappropriate and need a redesign.

In short, the shortcomings of the existing list-only EL solution is two-fold:

- Entity coherence within or across entity lists were overlooked and not leveraged; and
- Results were supportless for lack of appropriate dataset and deliberate experiment design.

We close the gap and address the deficiencies in this article. In particular, we propose to solve list-only EL task by taking account of the correlations in entities and converting the disambiguation problem to a graph problem. We show the merits of graph-based list-only EL by referring to the example in Figure 1. It is easy to map mention *New York Knicks* to entity *New York Knicks* in the entity list featured *Teams*. Then by considering the interdependence of entities in the same list *Teams*, mention *Chicago* will be mapped to entity *Chicago Bulls*, and mention *Los Angeles* will be mapped to entity *Los Angeles Lakers*. Additionally, by further taking into account cross-dependence of entities across different entity lists, entity *Phil Jackson* in the *Coaches* entity list, rather than entity *Phil Jackson* in the *Boxers* entity list, will be chosen as the target entity for mention *Jackson*.

To implement graph-based list-only EL, we mainly carry through the following three steps. (1) Pre-processing—including an optional *named entity recognition* process and the candidate entity generation process. This step formalizes raw texts and produces mentions and candidate entities as inputs for later steps. (2) Entity information enrichment. The descriptions of entities are enriched by collecting representative texts from the inputs, which in turn enable the establishment of coherence among entities. (3) Graph-based entity disambiguation. An entity graph is constructed by integrating outputs from earlier steps. We propose a graph-based algorithm Gloel, which implements Personalized PageRank to determine how likely an entity is the target entity by taking into consideration both coherences among entities, and compatibilities between mentions and entities. The output is a list of pairs comprising mentions and their most possible entities.

Furthermore, we put forward a new procedure to construct datasets applicable to evaluating list-only EL. The experimental results in this new dataset validate the effectiveness of graph-based linking, a popular method of collective linking, and the in-depth analysis shows that compared with existing list-only linking method, our graph-based solution achieves better performance in list-only EL task.

Contributions. The main contributions of this article can be summarized into three ingredients:

- We motivate to revise list-only EL by taking into account relations between entity lists, i.e., global coherence, in addition to local compatibilities between mentions and entities.
- We tackle the problem by a graph-based method and offer a new algorithm Gloel, where Personalized PageRank is adopted to capture global coherence among candidate entities.
- A new dataset construction procedure is presented to cater to the redefined task, and Gloel is experimentally evaluated on top of it, and shown to outperform state-of-the-art method.

Organization. This paper is organized as follows. In Section 2, the new definition of list-only EL problem and the methodology, which contains three steps, are elaborated. New dataset construction and experiment results are detailed in Section 3. Section 4 summarizes related work and bridges list-only EL with conventional KB-oriented EL, followed by conclusion in Section 5.

2. Methodology

We start with defining the proposed problem. Existing work defined list-only EL as mapping a *single* mention m_i in document d_i to the corresponding entity $e_{i,j} \in E_j$ in the entity lists. Nonetheless, on the one hand, in most real-life documents, there are more than one mention, differentiating this definition from reality. On the other hand, the ambiguity between entity lists is not stressed, which can turn the problem of mapping mentions to a group of highly ambiguous entities into determining whether the mentions have corresponding entities in the entity lists. And the latter also deviates from the original motivation of EL task, which centres on disambiguating mentions from several possible meanings. We will further elaborate the definition of ambiguity between entity lists via mathematical equations in Section 3.

As a consequence, it is vital to extend the definition of this task so as to cater to broader scenarios. Specifically, we formalize list-only EL problem as follows.

Definition 1 (List-only entity linking). *Given a set of documents $D = \{d_1, \dots, d_n\}$, each of which contains a set of mentions $M_i = \{m_{i1}, \dots, m_{is}\}$, an ambiguous set of entity lists $\mathcal{E} = \{E_1, \dots, E_l\}$, the task is to determine the most possible entity $e_{ij,k} \in E_k$ for each mention m_{ij} , or return NIL if there is no corresponding entity.*

Note that the set of entity lists has to be ambiguous. In other words, for the majority of entities, there ought to be at least one more ambiguous entity in the entity lists.

We take the example in Figure 1 to explain the definition. There are four mentions to be disambiguated in the document. By utilizing list-only EL, mentions *New York Knicks*, *Chicago*, *Los Angeles* should be mapped to entities *New York Knicks*, *Chicago Bulls*, *Los Angeles Lakers* in the entity

list featured **Teams** respectively, instead of *New York*, *Chicago*, *Los Angeles* in the entity list featured **Cities**. And mention *Jackson* should be linked to entity *Phil Jackson* in the **Coaches** entity list, rather than entity *Phil Jackson* in the **Boxers** entity list.

The specific procedure for graph-based list-only entity linking includes three steps, namely, pre-processing, entity information enrichment and graph disambiguation. The former two steps generate inputs, based on which the entity graph is constructed and Gloel is performed to determine results.

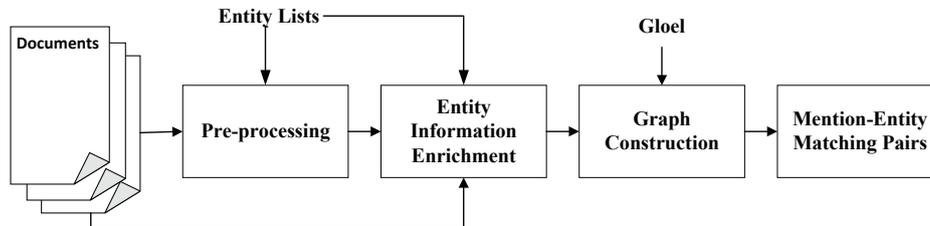


Figure 2. The flow diagram of graph-based list-only entity linking

2.1. Preparation for Graph Input

This subsection presents treatment of raw text data and generation of inputs for graph construction.

2.1.1. Pre-processing

In the pre-processing step, mentions in the text are detected and the candidate entities are also generated.

Specifically, the initial input for EL is a set of raw documents, either with specified mentions to be disambiguated or without. Under the circumstance where mentions are not pointed out, Named Entity Recognition (NER) should be harnessed to finish the mention detection task. State-of-the-art NER methods utilize Neural Networks and Deep Learning techniques to achieve better performances, whereas they have not been widely used yet on account of the freshness and complexity. Instead, Stanford NER Tagger, a NER tool which is less accurate but maturer, embraces higher popularity in tasks involving but not focusing on NER. In our experiment, we have already extracted the mentions during dataset construction process.

After obtaining mentions, the following step is retrieving possible candidate entities for each mention. Take Figure 1 for instance, for mention *Chicago*, both entities *Chicago Bulls* and *Chicago* should be generated as candidates. In order to improve recall and generate more candidate entities, most KB-oriented EL methods tend to take advantage of name dictionaries embedded in KBs, or use alias dictionaries built from collecting Wikipedia redirecting and disambiguation pages. However, considering the limited number of target entities and sparse information of entity lists, we design a set of simple but efficient string matching rules for entity generation, as is shown in Table 1. In the examples, the left are mentions while the right are candidate entities.

Table 1. String matching rules

Rules	Examples
Containment	Chicago → Chicago Bulls
Partial Matching	President Trump → Donald Trump LA → Los Angeles
Alternative Names	National Capital → Washington, D.C. Smiley → Miley Cyrus

The generated candidate entities for mention m_{ij} are represented by $Can(m_{ij})$. Noteworthy, we adopt candidate-pruning policy to ensure that a mention will not have two or more candidates from the same list, since entity list is utilized to help candidate entity within it to compete with entities

151 from other lists, and choosing among candidate entities from the same list will render coherence within
152 entity list useless.

153 2.1.2. Entity Information Enrichment

154 Solely relying on co-occurrences between entities is not enough to establish relations among entities,
155 let alone semantically bridge mentions with candidate entities. Therefore, we enrich information on the
156 entity side by selecting representatives derived from input documents.

157 Given input documents $D = \{d_1, \dots, d_n\}$, the mentions $M_i = \{m_{i1}, \dots, m_{is}\}$ in each d_i , a
158 set of entity lists $\mathcal{E} = \{E_1, \dots, E_l\}$, the enrichment process should collect *a set of highly relevant*
159 *and representative texts* $T^r = \{t_1^r, \dots, t_h^r\}$ around mentions for E_r , which is achieved by harnessing
160 co-occurrences of entities in the same entity list.

161 Specifically, the idea is that, since a document is not only composed of mentions, but also a lot
162 of other irrelevant information, we merely extract the texts around all mentions in all documents as
163 *candidate representatives* τ to avoid noisy information. If a candidate representative $t_p \in \tau$ contains
164 many entity names from the same entity list E_r , chances are that it indeed shares the same category or
165 topic with entity list E_r , and the mention m_p in candidate representative t_p is thus much more likely
166 to refer to the candidate entity from E_r . Consequently, t_p is a representative of E_r and the text in t_p
167 can be used to enrich the textual descriptions of entities in E_r .

168 We further illustrate the method in Figure 3. Note that in each candidate representative, the
169 bold text represents a mention, and the rest texts are its surroundings. Given an entity list **Cities**
170 and the entity *Chicago*, the goal is to collect relevant representatives for *Chicago* from documents,
171 which are then used to enrich representatives of entity list **Cities**. In *Document: United Paramount*
172 *Network*, there are three candidate representatives, two of them contain name string *Chicago*. However,
173 Candidate Representative 1 includes no extra name strings of other entities from the entity list, thus
174 might not refer to Entity List **Cities**. In contrary, both *New York* and *Los Angeles* co-occur with
175 *Chicago* in Candidate Representative 3, indicating the high possibility that it is a true representative
176 for Entity List **Cities**. Switching to *Document: Gotham City*, both Candidate Representatives contain
177 name string *Chicago*. Despite that Candidate Representative 4 is derived from mention *Chicago*, we
178 cannot consider it as a representative due to lack of co-occurrences information. Conversely, containing
179 several name strings from entity list **Cities**, Candidate Representative 5 is chosen as a representative,
180 even though it is built surrounding mention *Detroit*.

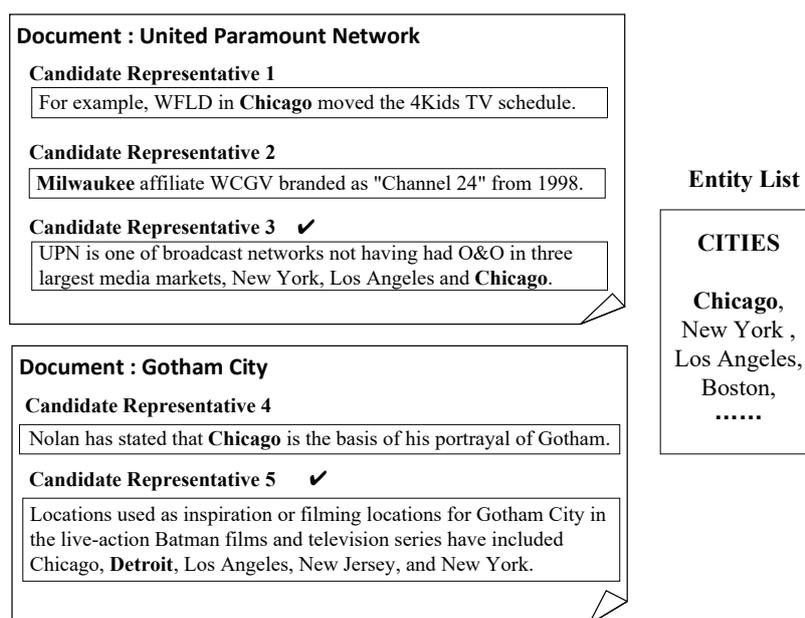


Figure 3. Example of entity information enrichment.

181 2.2. Graph Construction and Disambiguation

182 In this subsection, we illustrate the construction of candidate entity graph, followed by the
183 description of our proposed algorithm Gloel, which takes advantage of Personalized PageRank so as to
184 determine target entities.

185 2.2.1. Graph Construction

186 Through the pre-processing step, mentions and their candidate entities are obtained. Then after
187 enriching textual descriptions in the entity side, the *compatibility score* between each mention and
188 corresponding candidate entity can be calculated in terms of text similarity. Previous list-only EL ranks
189 the candidate entities for each mention merely based on mention-entity compatibility scores, thereby
190 producing the results accordingly. We argue that the judgement simply depending on compatibility score
191 is not convincing enough because the coherence among entities is ignored, which plays an indispensable
192 role in the linking process. For instance, as is shown in Figure 1, it is easy to map mentions *Chicago*
193 and *Los Angeles* to the *Cities* entities *Chicago* and *Los Angeles* due to the short text information
194 and high name string similarity. Provided that the candidate entity coherence is considered, the high
195 interdependence among *Teams* entities *New York Knicks*, *Chicago Bulls*, *Los Angeles Lakers* will lead
196 to the correct answers for mentions *Chicago* and *Los Angeles*.

197 To better capture the correlations among entities, similar to many existing KB-based EL methods,
198 we construct an entity graph, which is depicted in Figure 4. The definition of entity graph is defined as
199 follows:

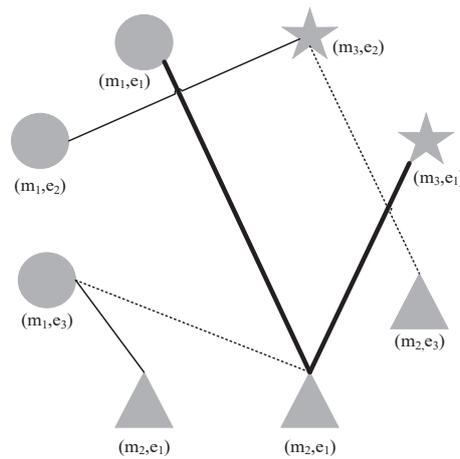


Figure 4. Entity graph.

200 **Definition 2** (Entity graph). An entity graph $G = \{V, E\}$ is a weighed graph, in which the nodes
 201 V represent all candidate entities, with their source mentions specified, and edges E include relations
 202 between entities.

203 It is noteworthy that we differentiate the mentions with identical name strings even though they
 204 might appear in the same document, and similarly, by specifying the source mention of nodes, the
 205 candidate entities with the same name but generated from different mentions are also treated differently.
 206 In this way, the situations where there are duplicate nodes, either caused by mentions or entities, can be
 207 avoided. In the mathematical form, we represent the r -th candidate entity for mention m_{ij} in document
 208 d_i as $e_{ij,r}$, which clearly shows the source mention m_{ij} of candidate entity $e_{ij,r}$.

209 With reference to edges, following the tradition in KB-based EL and adapting it to list-only
 210 problem, we connect two nodes with an edge under three circumstances: 1) The name strings of the
 211 two entities are in the same entity list $E \in \mathcal{E}$, and in this case, the edge weight is defined as 1; 2) The
 212 name strings of the two entities simultaneously appear in at least one candidate representative $t \in \tau$
 213 ; 3) The name strings of the other entities in the entity lists these two entities separately belong to,
 214 simultaneously appear in at least one candidate representative $t \in \tau$. The first two kinds of relations
 215 are termed as *explicit relations*, while the third method of adding edges among entities, named *implicit*
 216 *relations* mining, leverages the unique characteristic of entity list — that the rest entities $E'_i = E_i \setminus \{e_j\}$
 217 in the same entity list E_i can help mine more correlations for entity e_j even if e_j is in the long tail. As
 218 for edge weight, which is defined below, takes into account both explicit and implicit relations between
 219 entities. Furthermore, the edges among candidate entities with the same source mention are pruned so
 220 as to eliminate the influence generated by competitors themselves.

221 We further assign *initial node weight* $ini(v)$ and *edge weight* on the graph. The *initial node weight*
 222 $ini(v)$ is defined as the *compatibility score* between candidate entity and its source mention, while *edge*
 223 *weight* is determined by *relation score* between the entities on the two sides of the edge. The specific
 224 approaches to calculate *compatibility score* and *relation score* are:

225 **Compatibility score.** Given a document d_i , and m_{ij} , a mention contained in d_i , suppose $e_{ij,r} \in E_r$
 226 is a candidate entity for m_{ij} and $T^r = \{t_1^r, \dots, t_h^r\}$ is the set of representative texts for E_r . The
 227 compatibility score $\phi(m_{ij}, e_{ij,r})$ can be measured by the following equation

$$ini(m_{ij}, e_{ij,r}) = \phi(m_{ij}, e_{ij,r}) = \frac{1}{|T^r|} \sum_{p=1}^{|T^r|} Sim(m_{ij}, t_p^r).$$

228 Since entities in the same entity list share the same representative texts, which are collected
 229 according to the method proposed in Section 2.1.2, calculating compatibility between a pair of mention

230 m_{ij} and candidate entity $e_{ij,r} \in E_r$ can be converted to computing the average text similarity Sim
 231 between texts surrounding mention m_{ij} and all the text representatives T^r of candidate entity $e_{ij,r}$.

There are many ways to measure text similarity Sim and in this paper, we choose to compute the similarity between *embedding vectors* of two texts, which is represented as $E(m_{ij}, t_p^r)$. Additionally, we also regard the name string similarity between mention and candidate entity as an appropriate indicator, and it is denoted as $N(m_{ij}, e_{ij,r})$. Thus, the *Compatibility score* equation is converted to

$$\phi(m_{ij}, e_{ij,r}) = \alpha N(m_{ij}, e_{ij,r}) + \beta \frac{1}{|T^r|} \sum_{p=1}^{|T^r|} E(m_{ij}, t_p^r).$$

232 In the equation above, α and β are the weight coefficients balancing the importance of text similarity
 233 and name string similarity.

Relation score. Given two entities $e_i^p \in E_p, e_j^q \in E_q$ (We merely consider relationships among entities when calculating Relation Score, which is mention-irrelevant, thus we neglect the mention here), the *Relation Score* is denoted in the following equation

$$Rel(e_i^p, e_j^q) = \begin{cases} \eta O(e_i^p, e_j^q) + \frac{\theta}{M} \sum_u^{E_p-i} \sum_v^{E_q-j} O(e_u^p, e_v^q), & p \neq q; \\ 1, & p = q, \end{cases}$$

234 where $O(e_i, e_j) = \frac{|Occur(e_i) \cap Occur(e_j)|}{|Occur(e_i) \cup Occur(e_j)|}$, $Occur(e) = \{t | e \in t, t \in \tau\}$, and $M = (|E_p| - 1)(|E_q| - 1)$.

235 We illustrate equations above as follows: $Occur(e)$ denotes the occurrences of entity e in *All*
 236 candidate representatives τ , since compared with noisy textual information contained in the whole
 237 documents, merely considering texts around mentions (candidate representatives) can improve the
 238 accuracy. The Co-occurrence Frequency $O(e_i, e_j)$ of two entities e_i and e_j is defined as the number of
 239 candidate representatives they both occur in, divided by all the candidate representatives they either
 240 occur in together, or separately. As for the *Relation Score* $Rel(e_i^p, e_j^q)$ of two entities $e_i^p \in E_p, e_j^q \in E_q$, if
 241 p equals q , which means e_i^p and e_j^q are from the same entity list, we set the relation score as 1. Otherwise,
 242 the score is composed of two parts. The first component is the direct Co-occurrence Frequency of these
 243 two entities, multiplied by a weight factor η , which indicates explicit relations. While the implicit
 244 relations are represented by indirect Co-occurrence Frequency, the second component with a coefficient
 245 θ , which takes into account the co-occurrences of the rest entities in E_p and E_q in a pair-wise fashion.

246 Furthermore, as is shown in Figure 4, there are three kinds of lines. The bold line represents that
 247 entities on the two sides are in the same entity list, and the *Relation Score* is 1. The dotted line denotes
 248 that two entities merely have implicit relations, while the normal line requires that there are explicit
 249 relations between entities.

250 It is noteworthy that, different from traditional KB-oriented EL problem which merely considers
 251 the direct relations between two entities, we extend the definition by taking into account the contribution
 252 made by relations between two entity lists as well, and represent them as implicit relations of two
 253 entities. The detailed approach to quantitatively describe the implicit relations is embodied in the
 254 equations above.

255 2.2.2. Ranking Mention-entity Pairs

256 Given a weighed entity graph G_i of document d_i , the target is to find for each mention m_{ij} in
 257 document d_i the most likely entity $e_{ij,k}$ from a group of entities. In line with popular methods proposed
 258 in KB-oriented EL [3], we propose graph-based list-only entity linking algorithm, namely Gloel, which
 259 utilizes Personalized PageRank to depict the coherence among candidate entities.

Specifically, we assign a vector $p(v_s)$ with length n to each node v_s to represent the results of a PageRank process starting from v_s . To better capture the coherence among entities within the same document, instead of regarding the similarity between the vectors of nodes as the *coherence score*, we

define it as how a candidate entity fits in the document. To enable the definition, a n -length vector $p(d_i)$ is also assigned to document d_i , representing the results of the PageRank process initiating from a group of unambiguous nodes. Consequently, the *coherence score* of a candidate entity $e_{ij,r}$ for mention m_{ij} in document d_i is defined as

$$\psi(e_{ij,r}, d_i) = \frac{p(v_{ij,r})p(d_i)}{|p(v_{ij,r})||p(d_i)|}.$$

260 We first elaborate the random walk process initiating from a single node, then extend it to
 261 calculating document PageRank vector. The PageRank algorithm, based on random walk theory, is
 262 firstly proposed to measure the importance of web pages by counting the number and quality of links to
 263 this page. It has been applied to EL problems in recent years, and has achieved great performance [3–6].
 264 The basic elements of PageRank include initial vector r^0 , transition matrix A , and preference vector s .
 265 Note that in our method, $r^0 = s$.

Transition Matrix A is the same in both individual and collective processes, the value at i th row and j th column is defined as

$$A_{ij} = \frac{Rel(e_i, e_j)}{\sum_{e_k \in Edges(e_i)} Rel(e_i, e_k)}.$$

266 where $Edges(e_i)$ represents the edges connected to entity e_i .

267 When computing the vector $p(v_t)$ for a single node v_t , $r^0 = s = (0 \dots 0, 1(tth), 0 \dots)_n$, which
 268 means that r^0 and s are identical n -length vectors, the position t of the vector is assigned with 1 and
 269 the rest are endowed with 0.
 270

271 The situation is slightly more complicated as for document PageRank vector $p(d_i)$. Firstly, we
 272 regard a candidate entity $e_{ij,r}$ as a unambiguous one if it satisfies one of the following conditions:

- 273 1. $e_{ij,r}$ is the only candidate entity of mention m_{ij} and $ini(e_{ij,r})$ is above threshold μ . The
 274 unambiguous entities of this kind is endowed with initial weight λ .
- 275 2. When there are more than one candidate entities and $e_{ij,r}$ is the candidate entity with the largest
 276 initial value, suppose $e'_{ij,r}$ is the candidate entity with the second largest initial value. It suffices
 277 that $ini(e_{ij,r}) - ini(e'_{ij,r}) \geq \nu$. The initial weight of this kind is κ .
- 278 3. If there are no candidate entities meeting the conditions, all the candidate entities will be added
 279 to the unambiguous entities set, with the same weight endowments.

280 After obtaining unambiguous entities set, the actual weight can be assigned via normalization of
 281 initial weight values. Note that in graph, unambiguous entities are presented as equivalent nodes, and
 282 by placing the actual weight of unambiguous nodes in the corresponding positions of the n -length vector,
 283 we can attain r^0 and s for document accordingly. Furthermore, we adopt an iterative disambiguation
 284 approach. In other words, after erasing ambiguity for each mention, the chosen result entity will be
 285 regarded as unambiguous and added to the unambiguous entities set, with initial weight of ι . Afterwards,
 286 the document PageRank vector will be re-computed by utilizing the new unambiguous entities set.

With initial vector r^0 , transition matrix A , and preference vector s defined as above, the Personalized PageRank is presented as following

$$r^{t+1} = (1 - \rho) \times A \times r^t + \rho \times s.$$

287 In the equation above, t represents the t th iteration, and ρ denotes the probability that the random
 288 walk process jumps out of the original iteration and starts from a new vector, which is usually set at
 289 0.15. Normally, the restarting nodes are all nodes in the graph, and the weights in vector s are the
 290 same, which equal to $\frac{1}{|V|}$. Nonetheless, in this work, vector s is personalized and set as the same with
 291 initial vector, which means that the random walk merely restarts from the initial nodes, eliminating the
 292 effect from other nodes. When the iterative calculation reaches to a stage where r^k does not change
 293 any more or the variation is within a minimal range, we consider that it converges and $p(v_s)$, $p(d_i)$ are
 294 thereby attained. At last, we formalize the list-only EL problem in a mathematical way:

Definition 3 (List-only entity linking in mathematical form). Given a set of documents $D = \{d_1, \dots, d_n\}$, each of which contains a set of mentions $M_i = \{m_{i1}, \dots, m_{is}\}$, an ambiguous set of entity lists $\mathcal{E} = \{E_1, \dots, E_l\}$, the task is to determine the most possible entity $e_{ij,k} \in E_k$ for each mention m_{ij} , and

$$e_{ij,k} = \arg \max_{e_{ij,r} \in \text{Can}(m_{ij})} (\gamma \phi(m_{ij}, e_{ij,r}) + \delta \psi(e_{ij,r}, d_i)).$$

295

where γ and δ are two weight coefficients balancing the weight between mention-entity compatibility score and entity coherence score. NIL will be returned if there is no corresponding entity.

296

297

3. Experiments and Results

298

Considering the deficiency in current list-only EL dataset, we propose a similar but more comprehensive approach for dataset construction. Then our method is validated via experiments on this dataset and the merits are highlighted through comparison with state-of-the-art method.

299

300

301

3.1. Dataset

302

There are two shortcomings in current dataset [1]. For one thing, each document merely contains a single mention to be disambiguated, which does not fit in most real-life occasions. For another, the target entity lists are not ambiguous enough, giving rise to the situation that most mentions merely have one candidate entity, and the disambiguation problem is converted to judging whether this sole candidate entity is true or not. Take entity *Apple* in **Company** entity list shown in the dataset of [1], there is no other similar entities in the set of entity lists. As a result, when given a mention *Apple*, the candidate entity for it will only be *Apple* in **Company** entity list, and the problem is transformed into deciding whether the mention can be mapped to entity lists or not.

303

304

305

306

307

308

309

310

In order to overcome the deficiencies, we propose to mine target entity lists and collect documents. The entity lists can be constructed both manually and automatically, but the ambiguity must be ensured. Given two entity lists $E_m = \{e_{1,m}, \dots, e_{i,m}\}$ and $E_n = \{e_{1,n}, \dots, e_{j,n}\}$, for E_m , the ambiguity caused by the existence of E_n is defined as

$$\text{Amb}(E_m, E_n) = \frac{1}{|E_m|} \sum_{e_{i,m} \in E_m} \operatorname{argmax}_{e_{j,n} \in E_n} \text{amb}(e_{i,m}, e_{j,n}).$$

311

Note that $\text{amb}(e_{i,m}, e_{j,n})$ represents the ambiguity between two entities in different entity lists. Many approaches can be utilized to measure it, and in this paper, we harness the matching rules defined in the candidate entities retrieval section. If matching rules are satisfied, we endow 1 to $\text{amb}(e_{i,m}, e_{j,n})$. Otherwise, the value is determined by name string similarity. Furthermore, the reason why only the highest ambiguity value for $e_{i,m}$ is chosen lies in the fact that we merely need to assure $e_{i,m}$ has one ambiguous competitor to avoid the situation as the example above.

312

313

314

315

316

317

For the whole entity lists set $\mathcal{E} = \{E_1, \dots, E_l\}$, the ambiguity is denoted as

$$A(\mathcal{E}) = \frac{1}{|\mathcal{E}|} \sum_{E_p \in \mathcal{E}} \operatorname{argmax}_{E_q \in \mathcal{E} \setminus E_p} \text{Amb}(E_p, E_q).$$

318

Again, for each entity list E_p , we only consider the highest ambiguity it has with the rest entity lists in \mathcal{E} , since constructing a entity lists set with high ambiguity between each pair of entity lists is nearly impossible.

319

320

321

322 Referring to [1], we generated raw entity lists by utilizing NeedleSeek ², which were then filtered
 323 and processed according to the definition of ambiguity. At last, seven entity lists with 70 entities in
 324 total were generated. The ambiguity of newly-constructed entity lists set is 0.965, calculated according
 325 to the equations given above, while the value ³ for entity lists set in [1] is 0.267. Part of the newly
 326 constructed entity lists are presented in Table 2.

Table 2. Part of entity lists.

E_i	Entities
1	Atlanta, Chicago, Boston, Houston, New York, Detroit, Cleveland, Atlanta, Milwaukee
2	Atlanta Hawks, Chicago Bulls, Boston Celtics, Houston Rockets, Detroit Pistons...
3	Atlanta Braves Chicago Cubs, Boston Red Sox, Houston Astros, Detroit Tigers...
4	Toyota Camry, Ford Ikon, Tata Indica, Honda Accord, Hyundai Accent...
5	Toyota, Ford, Tata Motors, Honda, Hyundai, Chevrolet, Porsche, Volkswagen, Buick...
6	Cambridge, Oxford, St Andrews, Warwick, London, Edinburgh, Glasgow, Manchester...
7	University of Cambridge, University of Oxford, University of St Andrews...

327 As for building the documents dataset, we emphasize that there have to be at least two mentions
 328 in the same document to enable the construction of candidate entity graph. Otherwise there will be no
 329 difference between independent linking method and the proposed collective linking method based on
 330 graph.

331 To be specific, we utilized *wikilinks* in Wikipedia to obtain the documents. For each entity $e_{i,k}$ in
 332 entity list E_k , its referent Wikipedia page was determined in the first place. For instance, the Wikipedia
 333 page of entity *Atlanta* is en.wikipedia.org/wiki/Atlanta. Then we randomly retrieved 1,000 Wikipedia
 334 pages directing at $e_{i,k}$ via the *WhatLinksHere* page. As for *Atlanta*, the url of its *WhatLinksHere* page
 335 is en.wikipedia.org/wiki/Special:WhatLinksHere/Atlanta. After conducting the same operation for
 336 all the entities in entity list E_k , the links appearing in at least three entities' 1,000 Wikipedia pages
 337 were selected and the web pages texts they refer to were considered as documents. In this way, we can
 338 affirm that each document involves at least three mentions. Table 3 describes the specific information
 339 of documents and mentions.

340 3.2. Results and Analyses

341 We compare Gloel and the method utilized in [1] (denoted as *Independent*) on the dataset
 342 we create. The results are shown in Table 4 and the settings of parameters are listed as follows:
 343 $\alpha = 0.4, \beta = 0.6, \eta = 0.7, \theta = 0.3, \gamma = 0.5, \delta = 0.5, \lambda = 0.5, \kappa = 0.4, \iota = 0.3$.

344 The measurements we adopt are the same with the metrics in [7], namely *Precision*, *Recall* and
 345 *F1*. *Precision* takes into account all entity mentions that are linked by the system and determines the
 346 correctness. *Recall* on the other hand, considers all the mentions should be linked, and reflects the
 347 fraction of correctly linked mentions. *F1* is a balanced indicator of *Precision* and *Recall*.

² <http://needleseek.msra.cn>

³ Since the author did not offer the full entity lists information, we compute the ambiguity of the segmental entity lists presented in the previous work.

Table 3. Dataset statistics

Target E_i	#documents	#mentions
1	156	731
2	535	3,450
3	528	3,054
4	41	108
5	151	742
6	97	472
7	115	564
Total	1,623	9,121

Table 4. Experimental results on original dataset

Method		E_1	E_2	E_3	E_4	E_5	E_6	E_7	Overall
Independent	P	0.914	0.996	0.998	0.818	0.998	0.667	0.997	0.962
	R	0.990	0.984	0.988	0.998	0.959	0.989	0.585	0.959
	F1	0.950	0.990	0.993	0.900	0.979	0.797	0.734	0.961
Gloel	P	0.997	1.000	0.998	0.931	1.000	0.675	0.997	0.974
	R	0.997	0.998	0.997	1.000	0.981	0.992	0.598	0.971
	F1	0.997	0.999	0.998	0.964	0.990	0.803	0.748	0.972

Table 5. Experimental results on 50% corrupted dataset

Method		E_1	E_2	E_3	E_4	E_5	E_6	E_7	Overall
Independent	P	0.968	0.735	0.987	0.167	0.968	0.867	0.639	0.766
	R	0.856	0.994	0.624	0.944	0.287	0.318	0.961	0.764
	F1	0.909	0.845	0.765	0.284	0.443	0.465	0.768	0.765
Gloel	P	0.641	1.000	0.998	0.943	0.967	0.555	0.995	0.910
	R	0.997	0.966	0.899	0.769	0.985	0.992	0.332	0.907
	F1	0.781	0.982	0.946	0.847	0.976	0.712	0.497	0.909

348 We first report the results on original dataset. As is depicted in Table 4, Gloel outperforms
349 independent EL method in all occasions, with a overall F1 gain at 1.1%. Nevertheless, it is evident that
350 both methods achieve high Precision, Recall and F1 scores. This can be justified that most mentions
351 in the documents appear in the same name string form as the entity name strings. For instance,
352 in documents containing mention referring to entity *University of Cambridge*, the name form of the
353 mention is also *University of Cambridge*, thus the high name string similarity basically guarantees
354 the correct matching and rules out the possibility of other candidate entities. Plus, this does not fit
355 in situations of most text sources other than Wikipedia. In news reports concerning *University of*
356 *Cambridge*, it constantly goes by the name *Cambridge* as in sentence *Cambridge beats Oxford in terms*
357 *of computer science*. In these cases, the probability of generating result entity *Cambridge* is enhanced
358 significantly and the disambiguation difficulty also rises up.

359 As a consequence, we corrupted the dataset to observe the corresponding results produced by
360 these two methods. To achieve corruption and increase ambiguity, we replaced mention names of the
361 entities in lists 2,3,5,7 to the corresponding ambiguous names in lists 1,4,6. For instance, the mention
362 names *Atlanta Hawks* and *Atlanta Braves* were substituted by *Atlanta*. Considering the fact that
363 after corruption, the name string similarity (in [1] the NER results) might be of no use and possibly
364 lead to negative contributions, which was unfair for *Independent* results, we merely took into account
365 the embedding vectors similarity in terms of mention-entity similarity calculation and altered the
366 corresponding parameter setting.

367 The results of 50% corruption in Table 5 are generated after half of the entities' corresponding
368 mention names in lists 2,3,5,7 get replaced. For fair comparison, the parameters are optimized for

369 separate methods. As can be seen, the gap between the results of Independent and Gloel widens. Gloel
 370 achieves better outcomes with overall F1 score at 90.9%, while the overall F1 value of previous method
 371 is 76.5%, hence validating the superiority of the proposed method.

372 We further conducted 25% and 75% corruption on the dataset and Figure 5 dynamically depicts
 373 the F1 scores of these two methods under corruption. With input texts getting more difficult, the
 374 results of EL based solely on mention-entity compatibility decline rapidly, while Gloel, a method based
 375 on graph, still yields robust results with smaller decreases.

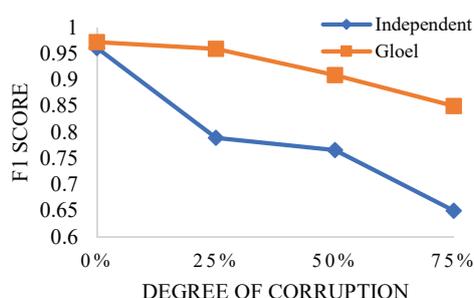


Figure 5. F1 score of independent and Gloel over corrupted dataset

376 4. Related Work

377 In this section, we brief related work, and discuss the differences and connections between list-only
 378 EL and traditional EL.

379 4.1. List-only Entity Linking

380 Over recent years, in accordance with the emergence of various text sources, EL tasks in new forms
 381 have been put forward. List-only EL task, first formally defined by Lin et al. [1], is the task of mapping
 382 mentions to a group of entity lists, rather than complete KBs. Lin et al. selected seed mentions for each
 383 entity list to bridge the gap between mentions and non-informative target entities, and then conducted
 384 the independent linking process to determine final results. Noticing that they merely harnessed entity
 385 lists co-occurrences information for generating entity descriptions, in this work, we further utilize the
 386 co-occurrences information to model entity relatedness and integrate it in the entity graph, which yields
 387 a more robust and accurate EL framework when confronting difficult input texts.

388 There are other new forms of EL problems which are similar to the list-only task. One is the
 389 Target Entity Disambiguation problem [2,8]. The main disparity is that the focus of Target Entity
 390 Disambiguation task lies in finding documents related to the entities given a entity list, whereas the
 391 starting point of list-only EL task is to eliminate the ambiguity in documents by using entity lists.
 392 Another similar task is the Named Entity Disambiguation with Linkless KBs [9]. Different from the
 393 mere entity lists in our task, there are still textual descriptions for entities in Linkless KBs.

394 4.2. Knowledge Base Oriented Entity Linking

395 Earlier works on EL focus on the situation where abundant information exists on the entity
 396 side. Specifically, KBs such as YAGO, Freebase and Wikipedia, offer rich semantic structures among
 397 entities as well as detailed textual descriptions, thus resulting in robust and accurate linking procedure.
 398 KB-oriented EL work can generally be divided into independent and collective methods.

399 In the former approach, mentions are disambiguated merely according the similarity between
 400 mentions and entities, and the problem is transformed into candidate entities ranking so as to obtain
 401 the most possible result. The similarity is mainly measured by lexical features such as bag-of-words of
 402 surrounding texts and statistical features such as prior popularities of entities. Then as for ranking
 403 process, unsupervised methods [10] calculate cosine similarities of feature vectors and output the results,
 404 whereas supervised approaches [11,12] construct classifiers by training on annotated dataset, and the

405 linking process is in the charge of classifiers when inputs are given. Although methods of this kind can
406 achieve good results, semantic coherences within entities are neglected, which prove to be essential in
407 improving overall performances.

408 With respect to collective linking methods in conventional EL task, most of them assume mentions
409 in the same document are semantically coherent, which also should fit in the textual topic of the whole
410 document. Therefore, the resulting entities also are expected to have high relatedness and the problem
411 is in turn converted to find matching pairs maximizing the coherence. Cucerzan [13] proposed to
412 harness Wikipedia categories to model coherence among entities, while Milne and Witten [14] reckoned
413 normalized Google Distance as another useful tool for measurement, which was utilized by Kulkarni et
414 al. [15] to form integer linear programming problem so as to collectively obtain results. Hoffart et al. [16]
415 defined keyphrase relatedness to capture entity coherence, and proposed to construct a mention-entity
416 graph, on which dense sub-graph generation algorithm was put forward to determine the sub-graph
417 containing one-to-one mention-entity matches. The method of re-formalizing the linking problem by
418 constructing mention-entity or entity-only graph distinguished itself among other works due to its
419 capability to integrate both local similarity information between mentions and entities, along with the
420 coherence information among entities. Based on this, several works [3–6] proposed and applied modified
421 graph algorithm on the graph, which improved the disambiguation accuracy and the adaptability to
422 difficult texts. Overall, the collective linking methods generally perform better than the independent
423 counterparts in terms of conventional KB-oriented EL.

424 4.3. Discussion on Differences and Connections

425 There are indeed many similarities between these two lines of works, despite of the evident
426 differences. The disparity mainly lies in the information on the entity side. Regarding conventional
427 KB oriented EL, entities have rich and well-structured descriptions offered by KBs, in terms of both
428 text description and internal links among entities [14–16]. Thereafter, researchers merely need to filter
429 valuable information to improve linking results. In stark contrast, with respect to list-only scenarios,
430 the mere information existing on the entity side is the co-occurrences among entity name strings in the
431 same entity list, which in turn requires information mining and enrichment. In this paper, to avoid
432 help from structured or semi-structured knowledge source, the dataset itself is leveraged to harvest the
433 relevant relations among entities, thus fulfilling the entity information enrichment task.

434 Nevertheless, aside from information mining process, the methods utilized in conventional researches
435 can be applied to this newly-defined problem and will achieve promising results. For instance, with
436 disambiguation problem taking the form of graph, [3–6] can all be implemented.

437 Above all, the techniques developed in traditional EL also apply in list-only EL problem, and the
438 extra work for the latter is to mine information on the entity side.

439 5. Conclusion

440 List-only entity linking task, as a new form of traditional EL problem, distinguishes itself by the
441 sparse information on the entity side. In this work, on the one hand, we propose to utilize entity
442 co-occurrences information to mine both textual description of entities and relations among entities, so
443 as to enrich entity information. On the other hand, inspired by conventional EL methods, we construct
444 an entity graph to capture relations among entities, on which the newly proposed algorithm Gloel is
445 applied to obtain results. Similar to the situation in traditional EL, our approach, a collective EL
446 method based on graph, outperforms independent EL on the dataset we create for fair comparison.

447 For future works, we plan to investigate two aspects. One is to consider the situation where an
448 entity appears in more than one entity list. For instance, *Washington, D.C.* can appear in entity lists
449 featured `American Cities` and `Country Capitals`. Another is to leverage word embedding techniques
450 and deep neural networks to better model mention-entity compatibility and entity coherence.

451 **Acknowledgments:** This work was in part supported by NSFC under grants No. 61402494, 71690233 and
452 71331008, and NSF of Hunan Province under grant No. 2015JJ4009.

453 **Author Contributions:** W. Zeng and X. Zhao conceived the problem and the method, conducted the
454 experiments and wrote the paper; J. Tang managed the project and revised the paper.

455 **Conflicts of Interest:** The authors declare no conflict of interest. The founding sponsors had no role in the
456 design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, and
457 in the decision to publish the results.

458 References

- 459 1. Lin, Y.; Lin, C.; Ji, H. List-only Entity Linking. Proceedings of the 55th Annual Meeting of the
460 Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume
461 2: Short Papers, 2017, pp. 536–541.
- 462 2. Cao, Y.; Li, J.; Guo, X.; Bai, S.; Ji, H.; Tang, J. Name List Only? Target Entity Disambiguation in
463 Short Texts. Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing,
464 EMNLP 2015, Lisbon, Portugal, September 17-21, 2015, 2015, pp. 654–664.
- 465 3. Guo, Z.; Barbosa, D. Robust Entity Linking via Random Walks. Proceedings of the 23rd ACM
466 International Conference on Conference on Information and Knowledge Management, CIKM 2014,
467 Shanghai, China, November 3-7, 2014, 2014, pp. 499–508.
- 468 4. Han, X.; Sun, L.; Zhao, J. Collective entity linking in web text: a graph-based method. Proceeding of
469 the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval,
470 SIGIR 2011, Beijing, China, July 25-29, 2011, 2011, pp. 765–774.
- 471 5. Alhelbawy, A.; Gaizauskas, R.J. Graph Ranking for Collective Named Entity Disambiguation.
472 Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL
473 2014, June 22-27, 2014, Baltimore, MD, USA, Volume 2: Short Papers, 2014, pp. 75–80.
- 474 6. Pershina, M.; He, Y.; Grishman, R. Personalized Page Rank for Named Entity Disambiguation. NAACL
475 HLT 2015, The 2015 Conference of the North American Chapter of the Association for Computational
476 Linguistics: Human Language Technologies, Denver, Colorado, USA, May 31 - June 5, 2015, 2015, pp.
477 238–243.
- 478 7. Shen, W.; Wang, J.; Han, J. Entity Linking with a Knowledge Base: Issues, Techniques, and Solutions.
479 *IEEE Trans. Knowl. Data Eng.* **2015**, *27*, 443–460.
- 480 8. Wang, C.; Chakrabarti, K.; Cheng, T.; Chaudhuri, S. Targeted disambiguation of ad-hoc, homogeneous
481 sets of named entities. Proceedings of the 21st World Wide Web Conference 2012, WWW 2012, Lyon,
482 France, April 16-20, 2012, 2012, pp. 719–728.
- 483 9. Li, Y.; Tan, S.; Sun, H.; Han, J.; Roth, D.; Yan, X. Entity Disambiguation with Linkless Knowledge
484 Bases. Proceedings of the 25th International Conference on World Wide Web, WWW 2016, Montreal,
485 Canada, April 11 - 15, 2016, 2016, pp. 1261–1270.
- 486 10. Bunescu, R.C.; Pasca, M. Using Encyclopedic Knowledge for Named entity Disambiguation. EACL 2006,
487 11st Conference of the European Chapter of the Association for Computational Linguistics, Proceedings
488 of the Conference, April 3-7, 2006, Trento, Italy, 2006.
- 489 11. Dredze, M.; McNamee, P.; Rao, D.; Gerber, A.; Finin, T. Entity Disambiguation for Knowledge Base
490 Population. COLING 2010, 23rd International Conference on Computational Linguistics, Proceedings
491 of the Conference, 23-27 August 2010, Beijing, China, 2010, pp. 277–285.
- 492 12. Mihalcea, R.; Csomai, A. Wikify!: linking documents to encyclopedic knowledge. Proceedings of the
493 Sixteenth ACM Conference on Information and Knowledge Management, CIKM 2007, Lisbon, Portugal,
494 November 6-10, 2007, 2007, pp. 233–242.
- 495 13. Cucerzan, S. Large-Scale Named Entity Disambiguation Based on Wikipedia Data. EMNLP-CoNLL
496 2007, Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing
497 and Computational Natural Language Learning, June 28-30, 2007, Prague, Czech Republic, 2007, pp.
498 708–716.
- 499 14. Milne, D.N.; Witten, I.H. Learning to link with wikipedia. Proceedings of the 17th ACM Conference on
500 Information and Knowledge Management, CIKM 2008, Napa Valley, California, USA, October 26-30,
501 2008, 2008, pp. 509–518.

- 502 15. Kulkarni, S.; Singh, A.; Ramakrishnan, G.; Chakrabarti, S. Collective annotation of Wikipedia entities
503 in web text. Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery
504 and Data Mining, Paris, France, June 28 - July 1, 2009, 2009, pp. 457–466.
- 505 16. Hoffart, J.; Yosef, M.A.; Bordino, I.; Fürstenau, H.; Pinkal, M.; Spaniol, M.; Taneva, B.; Thater, S.;
506 Weikum, G. Robust Disambiguation of Named Entities in Text. Proceedings of the 2011 Conference on
507 Empirical Methods in Natural Language Processing, EMNLP 2011, 27-31 July 2011, John McIntyre
508 Conference Centre, Edinburgh, UK, A meeting of SIGDAT, a Special Interest Group of the ACL, 2011,
509 pp. 782–792.