

# Re-ranking for Joint Named-Entity Recognition and Linking

Avirup Sil  
Temple University  
Computer and Information Sciences  
Philadelphia, PA  
avi@temple.edu

Alexander Yates  
Temple University  
Computer and Information Sciences  
Philadelphia, PA  
yates@temple.edu

## ABSTRACT

Recognizing names and linking them to structured data is a fundamental task in text analysis. Existing approaches typically perform these two steps using a pipeline architecture: they use a Named-Entity Recognition (NER) system to find the boundaries of mentions in text, and an Entity Linking (EL) system to connect the mentions to entries in structured or semi-structured repositories like Wikipedia. However, the two tasks are tightly coupled, and each type of system can benefit significantly from the kind of information provided by the other. We present a joint model for NER and EL, called NEREL, that takes a large set of candidate mentions from typical NER systems and a large set of candidate entity links from EL systems, and ranks the candidate mention-entity pairs together to make joint predictions. In NER and EL experiments across three datasets, NEREL significantly outperforms or comes close to the performance of two state-of-the-art NER systems, and it outperforms 6 competing EL systems. On the benchmark MSNBC dataset, NEREL provides a 60% reduction in error over the next-best NER system and a 68% reduction in error over the next-best EL system.

## Categories and Subject Descriptors

I.2.7 [Natural Language Processing]: Text analysis; I.3.1 [Content Analysis and Indexing]: Linguistic processing

## Keywords

Named Entity Recognition; Entity Linking; Entity Disambiguation

## 1. INTRODUCTION

Entity Linking (EL) [1, 3] is the task of identifying when a name that appears in text refers to a known entity in a reference set of named entities, such as a relational database or the set of articles in Wikipedia. Existing EL systems depend on knowing where the set of names in a text appear. Nearly all EL systems therefore use a pipeline architecture: they run a Named-Entity Recognition (NER) [6, 17] system on the text to identify the boundaries of names. They then use the names that were identified by the NER system, typi-

cally called *mentions*, as the input set of names for the EL task<sup>1</sup>. Unfortunately, as is common with pipeline architectures for probabilistic models, errors from the NER system propagate to the EL system. Any mention that goes undetected by the NER system clearly cannot be linked correctly by the downstream EL system. Other common errors from NER systems include incorrectly splitting a multi-word name into multiple mentions; joining together two mentions that happen to appear near one another; or leaving off one or more words from a multi-word mention. Each of these errors can cause the downstream EL system to generate too many, too few, or wrong links.

Consider the following phrase that has been incorrectly processed by a pipeline of a state-of-the-art NER system [17] and our own baseline EL system:

... creator of [ABC] | *ABC's The [View] | Pay\_Per\_View*,  
said ...

The word “the” commonly precedes certain mentions, as in “the [Rocky Mountains],” so it is reasonable but incorrect for the NER system to place the mention boundaries around just “View” in this example. Likewise, it is reasonable but incorrect in this context for an EL system to predict that “View” refers to *Pay\_Per\_View*, a common method of watching certain television programs, since most of the surrounding text is related to television programming.

We develop a joint NER and EL model that we use to re-rank candidate mentions and entity links produced by base NER and EL models. Our base NER techniques would produce candidate mentions like ABC, View, and The View from the above phrase. Our base EL system would then propose several promising entity links for each candidate mention, including *Pay\_Per\_View* and *The\_View* (an American television program on a broadcast channel called ABC). The reranking model then chooses among the set of all possible mention and entity link labelings for the whole phrase to determine the best choice. Since the reranking model considers labels for “ABC” and “(The) View” simultaneously, it can use features for known relationships between the television channel ABC and the television program *The\_View* to encourage these as outputs. For efficiency, such features are not available to the pipeline models, which consider only one mention at a time. We use the pipeline models to prune the set of all possible candidate mentions and entity links to a manageable size while maintaining high recall. The reranking model can then use more sophisticated features for collective classification over this pruned set.

Previous work [3, 15, 9, 11] has considered collective classification for entities, but never for the combination of NER and EL. Ratnov *et al.* [18] argue that collective classification for entity linking provides only a small benefit over a purely-local model. However, we find that by combining NER and EL, the set of correct mentions that are available to be linked improves drastically, and EL performance improves as a result.

<sup>1</sup>For this paper we ignore type labels like Person, Location, or Organization during NER, since we will link to a relational database with a finer-grained ontology.

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 the author(s) 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.

CIKM '13, Oct. 27–Nov. 1, 2013, San Francisco, CA, USA.

Copyright is held by the owner/author(s). Publication rights licensed to ACM.

ACM 978-1-4503-2263-8/13/10 ...\$15.00.

<http://dx.doi.org/10.1145/2505515.2505601>.

The remainder of this paper is organized as follows: The next section discusses related work. Section 3 describes the Freebase database, and why it is suited to the EL task. Section 4 presents our joint model for NER and EL. Section 5 describes our methods for finding candidate mentions and entities, including a pipeline EL system for Freebase. Section 6 gives our experimental results, and Section 7 concludes.

## 2. RELATED WORK

All but one previous entity linking system that we are aware of uses a pipeline architecture. Ratnov *et al.* [18] point out that finding mentions using NER is a non-trivial task, and hence their experiments report linking performance using gold-standard mentions, as do several other recent EL systems [20, 4, 8]. Nevertheless, in practice NER performance limits the performance of these EL systems. Our work relies on pipeline NER and EL techniques [1, 14, 22, 5, 18, 12, 13, 10, 20] to provide high-quality candidate mentions and entity links. We then subject these candidates to further processing by a joint NER and EL model that outperforms state-of-the-art NER systems and state-of-the-art EL systems on their respective tasks.

Guo *et al.* [7] are the only authors we are aware of that consider joint mention detection and entity disambiguation, but their techniques were developed for, and are best-suited for, microblog text. Because microblog texts are so short, Guo *et al.* can use computationally expensive machinery; they use a structural SVM algorithm which requires NP-hard inference. Their technique can consider nearly every token as part of a candidate mention, again because of the short length of microblog texts. In contrast, our techniques are better suited for longer documents. We use linear maximum-entropy models to re-rank a set of candidate mentions and entities provided by efficient NER and EL base models. A more minor difference is that Guo *et al.* link to Wikipedia; our technique links to both Wikipedia and Freebase—a large, user-contributed, relational database. Also, Guo *et al.*’s techniques cannot identify mention boundaries that have no corresponding Wikipedia entries, whereas our techniques can identify mentions with no corresponding entity in our reference set; we follow the Text Analysis Conference’s (TAC) guidelines in linking such mentions to the special symbol *NIL*.

## 3. ENTITY LINKING WITH FREEBASE

Most existing large-scale entity linkers use Wikipedia titles as a reference set. We chose to use a combination of Freebase and Wikipedia as our reference set because it provides us all of the information in the Wikipedia articles, as well as a well-developed ontology of types and binary relations available in Freebase. Freebase is a freely-available, user-contributed repository of structured data containing almost 40 million unique entities. In comparison, Wikipedia contains 4.2 million entities and DBpedia contains 3.77 million entities (all numbers are for the English versions). The Freebase schema is divided into *domains*, each of which contains its own set of *types* (such as people, books, bridges, *etc.*) and *properties*, like date of birth for a person or latitude and longitude for a location. Freebase commons, which we use in our experiments, contains 86 different domains, including such diverse areas as schools, geography, sports, and astronomy. Every entity in Freebase is given a unique *ID*, and separately a set of possible names that can refer to the entity, through the built-in name and alias relations. Many Freebase entities also have links to their corresponding Wikipedia pages, when available. Furthermore, Freebase includes a built-in, nontrivial API for finding candidate entity links for a given name, which is based on the user-contributed name and alias relations.

## 4. JOINT NAMED-ENTITY RECOGNITION AND ENTITY LINKING

### 4.1 Examples

While EL systems can go wrong even when provided correct mentions, we have found that in practice a large number of EL errors are caused by poor mention boundaries. Table 1 shows examples of three common errors by NER systems that propagate to EL systems. Other types of errors include completely-missed mentions, and correct mentions but incorrect links. Joint NER and EL can help with many of these problems by allowing for features that, for instance, favor joining two mentions that link to the same entity over having the same entity repeated twice.

### 4.2 Problem Formulation

Labeled data for joint NER and EL consists of a set of documents  $D$ , and for each document  $d \in D$ , a set  $Entities(d)$  of triples of the form  $(l, r, e)$ . Here,  $l$  and  $r$  indicate the left and right boundaries of the mention of some named entity, and  $e$  indicates the unique identifier of the named entity in a reference set. Given a labeled training set  $T$  of this form, the task is to learn a model that can predict  $Entities(d)$  for new documents  $d$ .

### 4.3 Decomposition into Connected Components

Classification over all possible sets of (mention, entity)-pairs for a document is generally intractable. Typical approaches to EL have addressed this by decomposing the task into a pipeline: first, identify the set of mentions using a standard NER model. Second, determine a local model for  $P(E|d, l, r)$ , where  $E$  is a random variable over entity ids  $e$ . Many recent EL approaches have investigated a third step, which is to select up to  $K$  tuples of entity ids  $(e_1, \dots, e_n)$ , one  $e_i$  for each mention in the document found by the NER model, and then build a model to rank these tuples of entity ids [1, 3]. However, Ratnov *et al.* [18] argue that this kind of global model provides a relatively small improvement over the purely-local approach, and Sil *et al.* [20] dispense with the global classification step entirely.

We instead adopt a different approach to decomposing the task, which offers a way to jointly classify mention boundaries and entity links. Rather than assume that the first-stage NER model produces the correct mention boundaries for linking, we assume that it provides a high-recall set of potentially-overlapping candidate mentions. One way to produce this set is to develop a probabilistic NER model  $P(M|d, l, r)$ , where  $M$  is a binary random variable indicating whether  $l$  and  $r$  are the exact boundaries of a mention. By selecting a relatively low threshold for this distribution, the model can produce a high-recall set of candidate mentions for  $d$  that have some reasonable probability of being correct. In practice, we have found that an alternative approach of getting candidate mentions (described in Sec. 5.1), which uses an NP chunker and an existing NER system, to be efficient and effective. Let  $A$  be the set of all sequences  $a_i = (m_1, \dots, m_{n_i})$  of mentions  $m$  such that each  $m$  is taken from the candidate set, and no two mentions in the same sequence have overlapping boundaries. As in the pipeline model, we also develop a local model  $P(E|d, l, r)$ , and select several promising candidates from the local model for each candidate mention. The joint classification task is over all tuples  $B$  of the form  $b_i = (m_1, e_1, \dots, m_{n_i}, e_{n_i})$ , where each  $e_j$  is taken from the local EL model’s candidate set for  $m_j$ . We call these *entity-mention tuples*.

The cost of high recall in the set of candidate mentions is that the set of candidate mentions is large, and as a result  $B$  contains far too many candidate tuples for tractable classification. To decompose the problem, we observe that most words in (non-microblog) text are not part of named-entities. We adopt a simple heuristic of separating all candidate mentions that are three or more words apart. In all three of the corpora in our experiments, this simple heuristic results in a partitioning of the mentions into small sets of mentions

Error Type	Bad Example	Correct Example
Under-segmentation	[Orange County, Calif.]→Orange County	[Orange County]→Orange County [Calif.]→California
Under-segmentation	[Diaz and Timberlake]→Justin Timberlake	[Diaz]→Cameron Diaz and [Timberlake]→Justin Timberlake
Over-segmentation	[Blue Cross]→BCBS of Alabama and [Blue Shield of Alabama]→BCBS of Alabama	[Blue Cross and Blue Shield of Alabama]→ BCBS of Alabama
Mis-segmentation	[Florida Supreme]→(FL Supreme Court) Court	[Florida Supreme Court]→FL Supreme Court

**Table 1: Examples of three kinds of NER errors that cause a pipeline EL system to miss links or add extra links. Mention boundaries are shown in [square brackets], and entity links are indicated by →.**

that appear near one another. We refer to these sets as the *connected components* of  $d$ , or  $CC(d)$ . We perform classification over the set of entity-mention tuples  $B(cc)$  that are formed using candidate mentions and entities within the same connected component  $cc \in CC(d)$ . As an example, the phrase “ABC’s The View” would constitute a connected component, since several candidate mentions (ABC and View) are separated by three or fewer tokens. Two of the entity-mention tuples for this connected component would be ([ABC], ABC, [View], Pay-Per-View) and ([ABC], ABC, [The View], The View).

In over 98% of our training documents, all of the connected components contained tuples  $b_i$  with  $n_i \leq 10$  mentions. The few remaining large connected components were from documents describing the results of cricket matches, where long lists of player names were separated only by one or two statistics each. For such long sequences, we heuristically divide them into chunks of a length that is the maximum length that our implementation can handle, which is 15. With relatively low thresholds for selecting candidate mentions and entities, the number of tuples in each connected component in our training data was 31.02 on average, while still allowing a maximum possible recall of 96%. In contrast, the maximum possible recall for a pipeline EL system using a standard NER system in our datasets is 82%. The joint model adds an extra layer of processing on top of the pipeline model, but in experiments the re-ranking step has nearly identical running time as the combination of our pipeline NER and EL models. Overall, this simple decomposition leads to a tractable joint NER and EL classification task.

#### 4.4 Re-ranking Model

We use a maximum-entropy model to estimate  $P(b|d, cc)$ , the probability of an entity-mention tuple  $b$  for a given connected component  $cc \in CC(d)$ . The model involves a vector of real-valued feature functions  $\mathbf{f}(b, d, cc)$  and a vector of real weights  $\mathbf{w}$ , one weight per feature function. The probability is given by

$$P(b|d, cc, \mathbf{w}) = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(b, d, cc))}{\sum_{b' \in B(cc)} \exp(\mathbf{w} \cdot \mathbf{f}(b', d, cc))} \quad (1)$$

We use L2-regularized conditional log likelihood (CLL) as the objective function for training:

$$CLL(T, \mathbf{w}) = \sum_{(b, d, cc) \in T} \log P(b|d, cc, \mathbf{w}) + \sigma \|\mathbf{w}\|_2^2$$

where  $(b, d, cc) \in T$  indicates that  $b$  is the correct tuple of entities and mentions for connected component  $cc$  in document  $d$  in training set  $T$ , and  $\sigma$  is a regularization parameter. We used LBFG-S for gradient-based convex optimization.

#### 4.5 Features

The primary advantage of the discriminative re-ranking approach is that it provides the flexibility to include many novel features into NER and EL. For instance, we include features that penalize NER boundaries that oversegment a mention by counting how often an entity-mention tuple  $b$  links to the same entity multiple times in

a row. We also include features that help to learn and encourage common sequences of entity types, like (job title, employee name) and (city, country or state). These features help both to prevent undersegmentation in NER and to improve linking accuracy. Details of our features are below.

**Token-level features:** Given  $b$ ,  $d$ , and  $cc$ , we define features that find the minimum (**MIN**), maximum (**MAX**) and the average (**AVG**) number of tokens inside the mentions in  $b$ , and a **NUM** feature for the total number of mentions in  $b$ . We also use features that measure how many tokens are not part of any mention in  $b$ , but are part of a mention in some other  $b' \in B(cc)$ . Let  $Non\_Mention\_Tokens(b, cc)$  be the set of all such tokens, and  $All\_Tokens(cc)$  be the set of all tokens for  $cc$ . We construct the **NON-MENTION-COUNT** feature as  $|Non\_Mention\_Tokens(b, cc)|$ , and the **NON-MENTION-RATIO** feature as  $\frac{|Non\_Mention\_Tokens(b, cc)|}{|All\_Tokens(cc)|}$ .

**Capitalization:** Similar to [7], we include one binary feature to check if all candidate mentions in  $b$  are capitalized, one binary feature to check if all words in all mentions in  $b$  are capitalized, and one binary feature to check if there is any mention that has capitalized, non-mention tokens to its immediate left or right.

**Features from entity links:** We include a number of features that make use of the entity links in  $b$ , as well as the types and known relations between entities, as stored in Freebase. The most basic versions of these features include: **COUNT-NIL**, which counts the number of entities that link to *NIL*, to help penalize tuples  $b$  that are poor matches with Freebase; **COUNT-EXACT-MATCH**, which counts the number of mentions whose surface form matches exactly with one of the names for the linked entity stored in Freebase; **ALL-EXACT-MATCH**, which is true if all mentions in  $b$  match a Freebase name exactly; and **ACRONYM-MATCH**, if the mention’s surface form is an acronym for a name or alias of the linked entity in Freebase.

The **DISTINCT-LINKS-PER-MENTION** measures the ratio of distinct entities in  $b$  to mentions in  $b$ . This feature helps to penalize over-segmented phrases. This is because two mentions that should have been joined together often have high probability for links to the same entity, as in a bracketing of “[Home] [Depot]” that splits this name for a business into two parts, each of which has a high probability of linking to the company Home Depot. Typically, any value of less than 1 for this feature is an indication that  $b$  is oversegmented, as it is uncommon to refer to the same entity twice in the short span of text that makes up one of our connected components, although appositives are an obvious exception.

To further establish the relation between mention candidates and their links, we define two features that mimic the Normalized Google Distance (NGD) as in [7, 18, 2]. For every consecutive pair of entities  $(e, e')$  that belongs to mentions in  $b$ , the **NGD-FREEBASE** feature computes

$$\sum_{(e, e') \in b} \frac{|types(e) \cap types(e')|}{|types(e) \cup types(e')|}. \text{ The Wikipedia version, } \mathbf{NGD-WIKIPEDIA}, \text{ computes } \sum_{(e, e') \in b} \frac{|Wiki(e) \cap Wiki(e')|}{|Wiki(e) \cup Wiki(e')|},$$

where  $Wiki(e)$  denotes the set of tokens in the Wikipedia page of  $e$ , excluding stop-words.

Two features, **ENTITY-TYPE-PMI** and **ENTITY-TYPE-PRODUCT-PMI**, make use of Freebase’s type system to find patterns of entities that commonly appear next to one another. Let  $T(e)$  be the set of Freebase types for entity  $e$ . We remove common Freebase types which are associated with almost every entity in text, like `/location/location` or `/people/person`, since they have lower discriminating power. From the training data, the system first computes pointwise mutual information (PMI) [21] scores for the Freebase types of consecutive pairs of entities,  $(e_1, e_2)$ :

$$PMI(T(e_1), T(e_2)) = \frac{\sum_{(e, e') \in T} \mathbf{1}[T(e_1) = T(e) \wedge T(e_2) = T(e')]}{\sum_{e \in T} \mathbf{1}[T(e_1) = T(e)] \times \sum_{e \in T} \mathbf{1}[T(e_2) = T(e)]}$$

where the sum in the numerator is taken over consecutive pairs of entities  $(e, e')$  in training data. The feature **ENTITY-TYPE-PMI** adds these scores up for every consecutive  $(e_1, e_2)$  pair in  $b$ . We also include another feature **ENTITY-TYPE-PRODUCT-PMI** which does the same, but uses an alternative product variant of the PMI score.

The **BINARY-RELATION-COUNT** feature encourages entity tuples  $b$  where the entities have a known relationship with one another. For instance, for the text “Home Depot CEO Robert Nardelli”, we wish to encourage the tuple that links “Home Depot” to the entity id of the company Home Depot, and “Robert Nardelli” to the entity id of the person Nardelli. Freebase contains a relation called `organization_board_memberships` that indicates that Nardelli is a board member of Home Depot. **BINARY-RELATION-COUNT** counts the number of such relations between every consecutive pair of entities in  $b$ . This feature helps not just to improve entity linking accuracy; it also helps to prevent under-segmented phrases. For instance, an entity-mention tuple  $b$  that contains just a single mention for the whole phrase “Home Depot CEO Robert Nardelli” would have a lower count for this feature than the correctly-segmented version, even if the incorrect mention was linked to Nardelli.

Finally, we include a similar version of this feature that uses Wikipedia rather than Freebase. The **WIKIPEDIA-COOCURRENCE-COUNT** feature computes, for every pair of consecutive entities  $(e, e') \in b$ , the number of times that  $e'$  appears in the Wikipedia page for  $e$ , and *vice versa*. It adds these counts up to get a single number for  $b$ .

## 5. BASE NER AND EL MODELS

### 5.1 Overgenerating candidate mentions

Previous work on EL systems [18, 7] have argued that EL systems often suffer because of the mention detection phase. The impact is mostly because of false negative mentions: any missed mention by the NER system is also a missed entity for EL. To improve on the pipeline model, our joint system requires *overgeneration* of candidate mentions: that is, a high-recall set of mentions with just enough precision to prevent the number of false positives from overwhelming subsequent processing.

Rather than modify an existing NER system to overgenerate mentions, we adopt an approach of combining multiple existing systems and heuristic techniques. We use the mentions found by the state-of-the-art UIUC NER system. In addition, we included strings matching the following heuristic tests as candidate mentions:

- All Noun-phrase (NP) chunks from the Illinois shallow parser [16]. We also retain all sub-expressions of up to 7 tokens of the NP chunks.
- All possible n-grams (for  $n \leq 8$ ) which match names and aliases of entities in Freebase.

- For any string selected by any of the above techniques, we add in one new mention that includes the word to the left, and another new mention that includes the word to the right.

### 5.2 A Base Entity Linker using Freebase

Let mention  $m = (d, l, r)$  consist of a document  $d$  and the left and right boundaries of the mention, and let  $S$  be a reference set containing entity ids and the special symbol *NIL*, which indicates an entity that does not match any id in the reference set. The task of a base (or local) EL model is to predict  $P(E|m)$ , where random variable  $E$  ranges over  $S$ . We call our base EL model BASEEL.

BASEEL, like our re-ranking model, is a maximum entropy model. Using Freebase and Wikipedia, we construct a vector of real-valued feature functions  $\mathbf{f}(m, s)$  describing the degree to which  $m$  and  $s \in S$  appear to match one another. The model includes a vector of real weights  $\mathbf{w}$ , one weight per feature function.

$$P(E = s|m, \mathbf{w}) = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(m, s))}{\sum_{s' \in S} \exp(\mathbf{w} \cdot \mathbf{f}(m, s'))}$$

As with the joint model, we use L2-regularized CLL as the objective function, and LBFG-S for gradient-based convex optimization. We use only gold-standard mentions for training the base EL model. Our manually annotated data contains only positive examples for entity ids. To generate negative examples, we use Freebase’s candidate generation technique (described in Section 3) to identify a set of potentially confusing bad matches.

BASEEL uses a combination of information from Wikipedia and Freebase to judge the similarity between a mention  $m$  and Freebase entity  $s$ . As a starting point, we leverage the four domain-independent features previously described by Sil *et al.* [20], adapted to Freebase. Intuitively, these features measure the number of attribute values and attribute names for entity  $s$  in Freebase that appear in the context of  $m$ .<sup>2</sup> We add in Freebase features that take advantage of its type system. One such feature, **TF**( $m, s$ ), counts, for each type  $t$  that entity  $s$  belongs to, the number of mentions in the context of  $m$  that have any candidate entities with the same type. The feature then adds up these counts for all the types of  $s$ . Additional features normalize these counts in various ways, and still other features make use of cosine similarity between the Wikipedia page for  $s$  and the context of  $m$ .

## 6. EXPERIMENTS

We evaluate the performance of NEREL first on the NER task against 2 state-of-the-art NER systems, and second on the EL task against 6 state-of-the-art EL systems.

**Datasets:** We use 3 standard datasets for EL as our test sets: the ACE dataset, as annotated by Ratnov *et al.* [18] for Wikipedia links; the MSNBC dataset, as collected and annotated by Cucerzan [3] for mention boundaries and Wikipedia links; and the CoNLL-2003 NER dataset (testb), with YAGO2 and Freebase IDs added by Hoffart *et al.* [10]. Table 2 provides key statistics on these datasets. For datasets annotated with Wikipedia links, we automatically extracted the corresponding Freebase IDs using Freebase’s built-in links to Wikipedia. In less than 1% of the examples for all the 3 test sets, Wikipedia contained a correct entity, but Freebase did not. We left the Wikipedia page as the correct label for these examples, and our BASEEL and NEREL systems always gets these examples wrong. In 10% of the MSNBC examples, Wikipedia had no entry, and the examples were marked *NIL* (ACE had no *NIL* entries). We left these *NIL* cases even when Freebase does contain a correct entity, to provide a fairer comparison with wikifiers. The MSNBC data consists of news articles covering 10 domains, 2 articles from each domain: Business, Entertainment, Politics, Science, etc, and is thus more diverse than either the CoNLL or ACE datasets.

The ACE dataset includes hand-labeled entity links for only the first nominal mention of each annotated coreference chain, and

<sup>2</sup>For a full list of features see [19].

domain	$ M $	$E \phi_{FB}(m) $	$NIL$	In Freebase
ACE	257	43.65	0%	100%
MSNBC	747	28.15	10%	89%
CoNLL	5616	27.85	20%	79%
Wikipedia	158715	12.62	0%	100%

**Table 2: Number of mentions, average number of candidate referents, % of mentions that are  $NIL$ , and % of mentions that are in Freebase (and Wikipedia) in our datasets. We train our system on the Wikipedia dataset.**

there are many correct mentions and entities which have not been labeled. This gold standard allows us to estimate the recall of an NER system, but not the precision, since many of the system’s guesses that don’t match the gold standard may still be correct. To estimate precision for NER on this dataset, we took a random sample of 100 examples of each system’s guesses, and manually judged whether they were correct or not. For EL, we measured precision and recall using the gold-standard mentions in the dataset.

Our training dataset consists of 10,000 random Wikipedia pages, where all of the phrases that link to other Wikipedia articles are treated as mentions, and the target Wikipedia page is the label. The dataset is made available by Ratinov *et al.* For training, we disregard mentions for which there is no Freebase entity that links to the target Wikipedia page; for the remaining mentions, we use the corresponding Freebase entity as the label.

**Evaluation Metric:** For NER evaluation, we measure precision, recall and F1 scores. To evaluate EL accuracy, we report on a Bag-of-Freebase ids (BOF) evaluation analogous to the Bag-of-Titles (BOT) F1 evaluation as introduced by [15, 18]. In BOF-F1, we compare the set of IDs output for a document with the gold set of IDs for that document (ignoring duplicates), and utilize standard precision, recall, and F1 measures. We separately computed exact-match accuracy, but the results were very similar to the BOF-F1 score, and we chose to report BOF-F1 as it allows for comparison with previous systems. For more details on BOT-F1, see [18]. Note that by construction, BOT-F1 and BOF-F1 provide identical numbers for our experiment.

**Competitors:** We compare NEREL with 2 state-of-the-art NER systems: the current state-of-the-art NER system for the benchmark CoNLL 2003 test set, the UIUC NER system [17]; and the Stanford NER system [6]. The former uses a regularized averaged perceptron model and external gazetteers for strong performance. The latter uses Conditional Random Fields and Gibbs sampling to incorporate long-distance dependencies into its NER model; the constraints ensure that the same phrase should be labeled consistently within the same document. We downloaded these two systems and ran the configurations that were reported to perform best on the CoNLL dataset. We found the results to be almost the same as reported by their authors.

The six EL competitors for NEREL include: 1) Cuc07 [3], which was the first system to use collective classification of entity links. 2) MW08 [15], which uses the set of unambiguous mentions in the text surrounding a mention to define the mention’s context. MW08 uses Normalized Google Distance to compute the similarity between this context and the candidate Wikipedia entry. 3) Kul09 [11] is an extension of both Cuc07 and MW08. It uses a hill-climbing approach to compute the joint probability distribution of all entity link assignments for a document. 4) Rat11 [18], the current state-of-the-art wikifier, has both a local and a global component. The local model uses techniques similar to traditional EL systems [14, 1] and similar to the Wikipedia-features used in our BASEEL model. The global component uses the predictions of the local model for all mentions in a document as the context of a mention, and they use both NGD and PMI [21] to compute the similarity between this context and a candidate Wikipedia page. 5)

Dataset	System	Prec	Rec	F1
ACE	UIUC	92.0	89.3	90.7
	Stanford	<b>97.0</b>	84.8	90.5
	NEREL	92.0	<b>92.4</b>	<b>92.2</b>
MSNBC	UIUC	69.8	75.4	72.5
	Stanford	<b>86.0</b>	72.2	78.5
	NEREL	83.7	<b>91.0</b>	<b>87.2</b>
CoNLL	UIUC	91.2	<b>90.5</b>	<b>90.9</b>
	Stanford	<b>95.1</b>	78.3	85.9
	NEREL	86.8	89.5	88.2

**Table 3: NEREL outperforms two state-of-the-art NER systems on two out of three datasets, and outperforms one of them on the third dataset.**

Hof11 [10], a state-of-the-art database EL system, links mentions to the YAGO2 knowledge base. Their technique uses a dense coherence graph algorithm to measure mention-entity similarity. 6) Sil12 represents our re-implementation of Sil *et al.*’s [20] domain-independent database EL system, which we used as a starting point for our BASEEL model. The four wikifiers and Hof11 use a collective classification framework, but only for collectively classifying sets of entity links, not for combining NER and EL. Sil12 and BASEEL are purely local models.

**Parameter Settings:** NEREL includes two important parameters:  $\sigma$ , the regularization weight; and the number of candidate entities we select from BASEEL. We set the value of  $\sigma$  by trying five possible values in the range [0.1, 10] on held-out data. We chose 5 random documents as development data from the CoNLL-2003 training set. We found  $\sigma = 0.5$  to work best for our experiments. We chose to select 3 candidate entities from BASEEL for each candidate mention (or fewer if BASEEL had fewer than 3 links with nonzero probability), as higher thresholds led to a computationally intractable number of entity-mention tuples. As mentioned previously, using 3 candidate links per mention resulted in a maximum possible recall for NEREL of 0.96 across our test sets, which is significantly higher than the maximum possible recall of 0.82 for pipeline models.

**Results:** Table 3 shows the performance of NEREL on the NER task. NEREL achieves the highest recall and F1 scores on two of the three datasets, and second-highest on the third dataset (CoNLL). The Stanford NER system has the best precision on all three datasets, and the UIUC system has the best overall performance on CoNLL. On MSNBC, NEREL outperforms both the state-of-the-art UIUC and Stanford NER systems by 0.14 and 0.07 gains in F1 score respectively. On the other datasets, differences are smaller, but NEREL consistently performs very well. We suspect that the gazetteers in the UIUC system help significantly on the CoNLL dataset. Another problem for NEREL on CoNLL was the presence of several long chains of 20 or more entities, which occurred in several articles reporting the results of cricket matches. NEREL tended to have much lower accuracy on these long chains of entities.

Table 4 compares NEREL with previously reported results by MW08 and Rat11 on the ACE and MSNBC datasets, as well as the results of our implementations of Sil12 and BASEEL. NEREL achieves an F1 score of 85.9 on ACE and 84.6 on MSNBC, clearly outperforming all competitors. Compared with the state-of-the-art Rat11 wikifier, NEREL improves by 0.086 and 0.097 F1 on ACE and MSNBC respectively. Part of this is due to better linking accuracy, probably resulting from our use of Freebase’s attributes, relations, and types for constructing features. Even the BASEEL and the relatively simple Sil12 system (ported to Freebase) outperform the Rat11 wikifier, despite the fact that the Rat11 wikifier is given gold-standard mentions as input, and BASEEL is given the output of the UIUC NER system.

Systems	ACE			MSNBC		
	Prec	Rec	F1	Prec	Rec	F1
MW08*	-	-	72.8	-	-	68.5
Rat11*	-	-	77.3	-	-	74.9
Sil12*	82.1	74.7	78.2	87.6	67.1	75.9
BASEEL	84.6	77.0	80.7	<b>89.3</b>	68.4	77.5
NEREL	<b>85.5</b>	<b>86.4</b>	<b>85.9</b>	85.8	<b>83.4</b>	<b>84.6</b>

**Table 4: NEREL outperforms all competitors on the ACE and MSNBC datasets. We obtain the numbers for MW08 and Rat11 from [18], which does not report precision and recall. \* indicates that the EL results are based on tests where gold-standard mentions are given as input.**

	Systems			
	NEREL	Hof11	Kul09	Cuc07
Prec@Rec=max	<b>84.22</b>	81.91	76.74	43.74

**Table 5: NEREL outperforms all previously-tested competitors on the CoNLL-2003 (testb) test set. Precision is shown at the highest recall achievable by each system. All competitors to NEREL use only the correct mentions detected by the Stanford NER system as input, and thus their recall is limited by the recall of the Stanford NER system.**

However, another important aspect of NEREL’s performance is its ability to correct poor mention boundaries that propagate to errors in linking. Compared with BASEEL, NEREL achieves significantly higher recall and F1, similar to the differences between NEREL’s NER performance and the performance of the UIUC NER system that BASEEL uses in its pipeline. Overall, NEREL improves by 0.052 and 0.071 F1 over BASEEL on these two benchmark datasets.

Table 5 illustrates the EL performance of NEREL on the CoNLL dataset. In order to compare with Hoff11 [10], we replicate the experimental evaluation that Hoff11 consider most important: systems are evaluated according to their precision at their maximum possible recall level. Furthermore, no mention candidates with *NIL* links are considered; in Hoff11’s experiments, this resulted in the removal of 20% of the mention candidates. In Hoff11’s evaluation, each system was supplied with the mentions from the Stanford NER system, and thus their recall was limited by the recall of the Stanford NER system. Hoff11 do not report this recall, but in our experiments we found that the Stanford system had a recall of 0.783 on the CoNLL test set, or 0.801 when excluding NIL mentions. As before, NEREL uses its overgeneration techniques to generate mention candidates. Its recall on this test set is 0.826; thus NEREL outperforms the next-best system, Hoff11, in both precision and recall.

## 7. CONCLUSION

Much of the previous research on entity linking has gone into improving linking accuracy over gold-standard mentions, but we observe that many of the common errors made by entity linkers in practice have to do with the pipeline architecture, which propagates errors from named-entity recognition systems to the entity linkers. We introduce a re-ranking model that performs joint named entity recognition and entity linking. The discriminative re-ranking framework allows us to introduce features into the model that capture the dependency between entity linking decisions and mention boundary decisions, which existing models do not handle. Furthermore, the model can handle collective classification of entity links, at least for nearby groups of entities. The joint NER and EL model has strong empirical results, outperforming a number of state-of-

the-art NER and EL systems on several benchmark datasets while remaining computationally inexpensive.

## Acknowledgements

This material is based upon work supported by the National Science Foundation under Grant No. IIS-1218692.

## 8. REFERENCES

- [1] R. Bunescu and M. Pasca. Using encyclopedic knowledge for named entity disambiguation. In *EACL*, 2006.
- [2] R. Cilibrasi and P. Vitanyi. The google similarity distance. *IEEE Transactions on Knowledge and Data Engineering*, 19(3):370–383, 2007.
- [3] S. Cucerzan. Large-scale named entity disambiguation based on wikipedia data. In *EMNLP-CoNLL*, pages 708–716, 2007.
- [4] A. Davis, A. Veloso, A. S. da Silva, W. Meira Jr, and A. H. Laender. Named entity disambiguation in streaming data. In *ACL*, 2012.
- [5] P. Ferragina and U. Scaiella. Tagme: on-the-fly annotation of short text fragments (by wikipedia entities). In *CIKM*, 2010.
- [6] J. R. Finkel, T. Grenager, and C. D. Manning. Incorporating non-local information into information extraction systems by gibbs sampling. In *ACL*, 2005.
- [7] S. Guo, M.-W. Chang, and E. Kıcıman. To link or not to link? a study on end-to-end tweet entity linking. In *NAACL*, 2013.
- [8] X. Han, L. Sun, and J. Zhao. Collective entity linking in web text: a graph-based method. In *SIGIR*, 2011.
- [9] X. Han and J. Zhao. Named entity disambiguation by leveraging Wikipedia semantic knowledge. In *CIKM*, pages 215–224, 2009.
- [10] J. Hoffart, M. A. Yosef, I. Bordino, H. Furstenu, M. Pinkal, M. Spaniol, B. Taneva, S. Thater, and G. Weikum. Robust Disambiguation of Named Entities in Text. In *EMNLP*, pages 782–792, 2011.
- [11] S. Kulkarni, A. Singh, G. Ramakrishnan, and S. Chakrabarti. Collective annotation of wikipedia entities in web text. In *KDD*, pages 457–466, 2009.
- [12] T. Lin, Mausam, and O. Etzioni. Entity Linking at Web Scale. In *AKBC-WEKEX*, 2012.
- [13] P. N. Mendes, M. Jakob, and C. Bizer. Evaluating DBpedia Spotlight for the TAC-KBP Entity Linking Task. In *TAC*, 2011.
- [14] R. Mihalcea and A. Csomai. Wikify!: Linking documents to encyclopedic knowledge. In *CIKM*, pages 233–242, 2007.
- [15] D. Milne and I. H. Witten. Learning to link with wikipedia. In *CIKM*, 2008.
- [16] V. Punyakanok and D. Roth. The use of classifiers in sequential inference. *arXiv preprint cs/0111003*, 2001.
- [17] L. Ratnov and D. Roth. Design challenges and misconceptions in named entity recognition. In *CoNLL*, 2009.
- [18] L. Ratnov, D. Roth, D. Downey, and M. Anderson. Local and global algorithms for disambiguation to wikipedia. In *Proc. of the Annual Meeting of the Association of Computational Linguistics (ACL)*, 2011.
- [19] A. Sil. A Re-ranking Framework for Joint Named-Entity Recognition and Linking. In *ACM PIKM*, 2013.
- [20] A. Sil, E. Cronin, P. Nie, Y. Yang, A.-M. Popescu, and A. Yates. Linking Named Entities to Any Database. In *EMNLP-CoNLL*, 2012.
- [21] P. D. Turney. Thumbs up or thumbs down? semantic orientation applied to unsupervised classification of reviews. In *Procs. of ACL*, pages 417–424, 2002.
- [22] Y. Zhou, L. Nie, O. Rouhani-Kalleh, F. Vasile, and S. Gaffney. Resolving surface forms to wikipedia topics. In *Coling*, pages 1335–1343, 2010.