

# Named Entity Normalization in User Generated Content

## ABSTRACT

Named entity recognition is important for semantically oriented retrieval tasks, such as question answering, entity retrieval, biomedical retrieval, trend detection, and event and entity tracking. In many of these tasks it is important to be able to accurately *normalize* the recognized entities, i.e., to map surface forms to unambiguous references to real world entities. Within the context of structured databases, this task (known as record linkage and data de-duplication) has been a topic of active research for more than five decades. For edited content, such as news articles, the named entity normalization (NEN) task is one that has recently attracted considerable attention. We consider the task in the challenging context of user generated content (UGC), where it forms a key ingredient of tracking and media-analysis systems.

We show that a baseline NEN system from the literature (that normalizes surface forms to Wikipedia pages) performs considerably worse on UGC than on edited news: 75% vs 92% accuracy on an English language data set and 65% vs 81% on a Dutch language data set. We identify several sources of errors: entity recognition errors, multiple ways of referring to the same entity and ambiguous references.

To address these issues we propose five improvements to the baseline NEN algorithm, to arrive at a language independent NEN system that achieves overall accuracy scores of 90% on the English data set and 89% on the Dutch data set. We show that each of the improvements contributes to the overall score of our improved NEN algorithm, and conclude with an error analysis on both Dutch and English language UGC. The NEN system is computationally efficient and runs with very modest computational requirements.

## Categories and Subject Descriptors

H.3.1 [Information Systems]: Information Storage and Retrieval—*Content Analysis and Indexing*

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. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SIGIR'08, July 20–24, 2008, Singapore.

Copyright 2008 ACM 978-1-59593-597-7/07/0007 ...\$5.00.

## General Terms

Algorithms, Experimentation, Evaluation

## Keywords

Named entity normalization, Wikipedia, Language technology

## 1. INTRODUCTION

The task of record linkage (RL) is generally to find entries that refer to the same entity in different data sources. This task has been investigated since the 1950s, usually in the case entries are considered with their attributes (e.g., person with phone, address) [21]. This task received significance because data sources have different ways of referring to the same real-world entity due to, e.g., different naming conventions, misspelling or use of abbreviation. The task of reference normalization is to analyse and detect these different references [6, 7]. When we consider the special case of this problem for natural language texts, we first have to recognize entities in the text, and we resolve these references either to entities exist within the document or to real world entities. This, together, is the problem known as *named entity normalization* (NEN).

We consider the NEN task within the setting of user generated content (UGC), such as blogs, discussion fora, or comments left behind by readers of online documents. For this type of textual data, the NEN task is particularly important within the setting of media analysis, reputation analysis and intelligence gathering. Many strategies deployed in these areas revolve around the idea of determining and tracking the impact of an event, i.e., determining the number, intensity and orientation of responses and identifying the stakeholders and other actors and entities involved.

The specific scenario on which we focus concerns the analysis of data that is increasingly common: online texts decorated with moderated but unedited comments left behind by web users. Examples include news sites (such as BBC news), and discussion and collaboration forum (such as linuxforum.com). These comments contain valuable information that complements the original text that triggered them, but the sheer volume and their (usually) flat organization makes them hard to comprehend. Hence, tools are needed that help organize the list of comments, by clustering them, summarizing them, computing aggregate information, creating hyperlinks between them, etc.

Let's consider an example. In the data set that we use for development purposes in this paper (see Section 4 for de-

**Figure 1: An excerpt from a BBC news article, with the excerpts of the three comments (out of total of 39). Separate named entities are underlined.**

- *News item:* Michael Schumacher wins his sixth victory in eight races—and tightens his grip on another Championship title. Do you think the title race is over? Have Your Say. Michael Schumacher extended his lead to 43 points after Juan Pablo Montoya’s Williams broke down with 12 laps to go. (...)
- *Comments:*
  1. Ferrari and Schumacher are now beyond the point where anyone can stop them (...)
  2. (...) Ralf, Montoya or DC need to win all the remaining seven races without Michael getting any points (...)
  3. (...) Both Williams’ drivers could be giving Schumi more of a challenge if their cars were reliable, as could Coulthard at McLaren (...)

tails), one of the news stories is about racing driver Michael Schumacher<sup>1</sup>; see Figure 1. To normalize the named entities in this news item and the comments it triggers, we need to resolve them to real world entities. We notice that there are two types of reference in the data. One is *within-document reference*, e.g., in the comments in Figure 1, *Michael* and *Schumi* are used to refer to *Michael Schumacher* mentioned earlier. The second kind of references are *references to real-world entities*, e.g., in the second comment, *DC* is used to refer to *David Coulthard*. Notice that resolving references of the latter type involves *named entity disambiguation*: in the context of Figure 1, *DC* is not used to refer to “Daimler-Chrysler,” “direct current” or the number 600 in Roman numerals.<sup>2</sup>

The main challenge in normalizing named entities (NEs) occurring in the comments on a news story is that commenters often do not use the full name of an already mentioned NE, use nicknames, misspell words or creatively pun with them. For example, in one of the examples in our Dutch dataset (a news article with 90 comments), singer *Anneke Grönloh* is referred to in 11 different ways, including variants such as *Anneke Grohnlöh*, *anneke gr ?hnloh*, *Mw. Gronloh*, *Anneke Kreunlo*, *Mevrouw G.*, etc. Other examples of creative language use include *G@@Gle* and *Bu\$h*. Besides, commenters often introduce additional NEs not even mentioned in the triggering news story, and some of the NEs used may actually refer to earlier comments. All in all, this turns NEN on UGC into a challenging problem.

NEN has been considered before, on structured data and on edited content. Of particular relevance to us is the recent work by Cucerzan [4], who compares two methods for NEN on edited content. In the present paper we apply a similar method to user generated content. We find several main sources of errors: NE recognition errors (incorrect boundaries of named entities and missing NEs), multiple ways of referring to the same entity, and ambiguous (out of context) references. We present five improvements to the baseline

<sup>1</sup>URL: [http://news.bbc.co.uk/sport2/low/sports\\_talk/2031010.stm](http://news.bbc.co.uk/sport2/low/sports_talk/2031010.stm).

<sup>2</sup>See <http://en.wikipedia.org/wiki/DC>

NEN algorithm to address these types of errors, namely: trimming, joining and ngramming NEs, approximate name matching, identification of missing references and name disambiguation. We assess the overall performance of the improved system as well as the individual contributions of the improvements.

The main contributions of the paper are: presentation and analysis of the problem of NEN in UGC, an algorithm for addressing the problem, and evaluation and analysis of the algorithm. Our algorithm was developed for Dutch and using Dutch data, but experiments with an English dataset indicate that it is well applicable to other languages. Moreover, the algorithm is computationally efficient.

The remainder of the paper is organized as follows. In Section 2 we discuss related work. In Section 3 we present a baseline algorithm for named entity normalization in user generated content, and describe an improved version based on an error analysis. In Section 4 we present our experimental setup and in Section 5 we present and analyze the results of our evaluation. A set of conclusions in Section 6 completes the paper.

## 2. RELATED WORK

Name matching and disambiguation has been recognized as an important problem in various domains. Borgman and Siegfried [2] present an overview of motivations, applications, common problems and techniques for name matching in the art domain; see [17] for recent experiments with classification for name matching. A dual problem, personal name disambiguation has also attracted a lot of attention, and a number of unsupervised methods were proposed [14, 18].

A similar problem has been known in the database community for over five decades as the *record linkage* or the *record matching* problem [8, 21]. However, there the task is more general: matching arbitrary types of records, not just person names or other types of named entities. Another type of research focuses on identification, disambiguation and matching of text objects other than named entities, specifically, temporal expressions [1].

Related problems occur in a different task: discovering links in text, or *wikifying*. Like NEN, this task involves identifying and disambiguating references to entities, and has also attracted attention of the research community [10, 15].

Research on named entity extraction and normalization has been carried out in both restricted and open domains. For example, for the case of scientific articles on genomics, where gene and protein names can be both synonymous and ambiguous, Cohen [3] normalizes entities using dictionaries automatically extracted from gene databases. For the news domain, Magdy et al. [13] address cross-document Arabic person name normalization using a machine learning approach, a dictionary of person names and frequency information for names in a collection. Cucerzan [4] considers the entity normalization task for news and encyclopedia articles; they use information extracted from Wikipedia combined with machine learning for context-aware name disambiguation; the baseline that we use in this paper (taken from [11]) is a modification (and improved version) of Cucerzan [4]’s baseline. [4] also presents an extensive literature overview on the problem.

Recent research has also examined the impact of normalizing entities in text on specific information access tasks.

Zhou et al. [22] show that appropriate use of domain-specific knowledge base (i.e., synonyms, hypernyms, etc.) yields significant improvement in passage retrieval in the biomedical domain. Similarly, Khalid et al. [11] demonstrate that NEN based on Wikipedia helps text retrieval in the context of Question Answering in the news domain.

Finally, in recent years, there has been a steady increase in the development or adaptation of language technology for UGC. Most of the attention has gone to blogs (see [16] for a recent survey on text analytics for blogs). Online discussion fora are more closely related to the data with which we work; recent research includes work on finding authoritative answers in forum threads [9, 12], as well as attempts to assess the quality of forum posts [20]. To the best of our knowledge, discussion threads as triggered by news stories of the kind considered here have not been studied before.

### 3. AN ALGORITHM FOR NAMED ENTITY NORMALIZATION IN USER GENERATED CONTENT

In this section we present a baseline algorithm for NE-normalization based on [4, 11], perform an error analysis and describe five improvements to the baseline, each accounting for a specific type of error identified.

#### 3.1 Baseline Algorithm

Algorithm 1, our baseline algorithm, takes as input a pair (A, R) where A is the triggering news article and R is the list of comments on A in reverse chronological order. It returns an *entity model*, i.e., a list of triples  $(s, n, p)$ , where  $s$  is a surface form (i.e., a named entity as it occurs in text),  $n$  is the normalized form of  $s$  (e.g., a title of the corresponding Wikipedia article), and  $p$  is the character position of  $s$  in the document. For example, one of the entity triples from the text in Figure 1 is (*Schumi, Michael Schumacher*, 57).

Line 1 of Algorithm 1 performs the NE recognition, i.e., it identifies NEs of types PERSON, LOCATION, ORGANIZATION or MISC (miscellaneous). Lines 2 and 4 do the preprocessing: we remove all noisy NEs, i.e., short (length  $\leq 2$  characters) or stopword-only NEs, with the exception of (capitalized) abbreviations, and remove diacritics (e.g., replacing  $\ddot{o}$  with  $o$ ).

Next, on line 5 we normalize each found NE using the function shown as Algorithm 2. This normalization algorithm treats NEs that are person names differently from other NE types (lines 2 and 3). Specifically, for persons we further remove common titles (such as *Mr*, *Mrs*) and perform *within-document reference resolution*, as detailed in Algorithm 3.

Algorithm 2 continues (line 5) by trying to link the NE to a Wikipedia article, calling the function `findWikiEntity` shown in Algorithm 4. If even after this step the NE is not normalized, we take the string itself as its normalized form (lines 6–7 of Algorithm 1).

The function `ResolveRefInDoc`, described in Algorithm 3, examines the list of entities already found and normalized earlier in the document, and finds matches based on first or last names.

The function `findWikiEntity`, described in Algorithm 4, first tries to match the input reference string with a Wikipedia article title (either exact match or case-insensitive, in

---

**Algorithm 1** Compute the entity model of a document

---

**Require:** a text document  $DOC$

```

1: REFSdoc  $\leftarrow$  NE-Recognition(DOC)
   {REFSdoc: list of (NE, type, position) triples}
2: REFSdoc  $\leftarrow$  Remove-NoisyNEs(REFSdoc)
3: for each (NE, type, position)  $\in$  REFSdoc do
4:   REF  $\leftarrow$  Removing diacritics(NE)
5:   REF-norm  $\leftarrow$  NormalizeNE( (REF, type), REF-
     NORMSdoc)
6:   if REF-norm = NULL then
7:     REF-norm  $\leftarrow$  NE
8:   end if
9:   REF-NORMSdoc  $\leftarrow$  (REF, REF-norm, position)
10: end for
11: return REF-NORMSdoc

```

---

this order).<sup>3</sup> If we find a matched Wikipedia page title, WT, then we further check whether the page is a Wikipedia redirection page (line 2). In case of a redirect, we take the title of the target Wikipedia page instead. Then we check if WT refers to a Wikipedia disambiguation page<sup>4</sup> (i.e., it lists a number of possible candidate pages for a given term). If this is the case, we select one of them, disambiguating between candidates using a heuristic from [11]: we select the candidate which has the highest number of incoming links in Wikipedia.

---

**Algorithm 2** NormalizeNE: named entity normalization

---

**Require:** a pair (NAME, Type), a list REF-NORMS

```

1: if Type = "PERSON" then
2:   NAME  $\leftarrow$  RemoveTitles(NAME)
3:   REF-norm  $\leftarrow$  ResolveRef-InDoc( NAME, REF-
     NORMS) {call Algorithm 3}
4: end if
5: if REF-norm = NULL then
6:   REF-norm  $\leftarrow$  findWikiEntity(REF) {call Algo-
     rithm 4}
7: end if
8: return REF-norm

```

---

#### 3.2 Error Analysis of the Baseline Algorithm

For development purposes, we annotated one news story (about two personalities in Dutch show business, Anneke Grönloh and Paul de Leeuw) and the first 90 accompanying comments. Then, we compared the performance of the baseline algorithm described above to this gold standard, and performed error analysis. Below we list the most common types of errors.

**Recognition errors** The boundaries of some NEs are recognized incorrectly: NEs include noise or are split in the middle, e.g., "`<ne>Gronloh!!!</ne>`", and "`<ne>Paul</ne> de <ne>Leeuw</ne>`".

<sup>3</sup>The order is important, e.g., "MAC" and "Mac" are two different Wiki entities, the former is an abbreviation and the latter may refer to Mac OS X.

<sup>4</sup>Such pages use a Wikipedia Disambiguation template, or having as title a surface form followed by the word "disambiguation", e.g., "Gall\_(disambiguation)", or followed by the word "surname", e.g., "Gray\_(surname)".

---

**Algorithm 3** ResolveRef-InDoc: within-document coreference resolution

---

**Require:** a string NAME, a list REF-NORMS  
{REF-NORMS: list of (REF, REF-norm, position) triples}

- 1: **for** each (REF, REF-norm, position)  $\in$  REF-NORMS **do**
- 2:   **if** NAME = REF or NAME = REF-norm **then**
- 3:     **return** (REF, REF-norm, position)
- 4:   **else if** NAME = firstName(REF) or NAME = lastName(REF) or NAME = firstName(REF-norm) or NAME = lastName(REF-norm) **then**
- 5:     {firstName is the first token of the input string, and the lastName is the rest}
- 6:     **return** (REF, REF-norm, position)
- 7:   **end if**
- 8: **end for**

---

**Algorithm 4** findWikiEntity

---

**Require:** a string NAME

- 1: WT  $\leftarrow$  findWikiTitle(NAME) {If NAME matches with a Wikipedia page’s title}
- 2: **if** isRedirectPage(WT) **then**
- 3:   WT  $\leftarrow$  getTargetPage(WT)
- 4: **end if**
- 5: **if** isDisambiguationPage(WT) **then**
- 6:   **return** Disambiguate(WT)
- 7: **end if**
- 8: **return** WT

---

**Multi-references** Commenters sometimes do not bother with punctuation, and the NER recognizes a group of adjacent names as a single NE. For example, “<ne> Anneke Gronloh Paul de Leeuw</ne>”, actually contains references to two persons. Commenters use capitalization and punctuation in non-standard ways, confusing the NE recognizer, which produces results like <ne>Gronloh ,Maak</ne> or <ne>Gronloh .Die</ne>.

**Variants** Commenters often use partial (first or last) names or nicknames of the entities mentioned in the trigger article. Users also misspell names in comments or use creative puns<sup>5</sup> In our development article with 90 comments we found 12 variants of the name *Paul de Leeuw* and 11 variants of the name *Anneke Grönloh*: “*Anneke G*”, “*Anneke Grohloh*”, “*Mevrouw G*”, “*Mw. Gronloh*”, etc.<sup>6</sup> In the entire annotated dataset (see Section 4.2), 29% of the entities were referred to using more than one surface form and 5% had five or more distinct forms.

**Missing NEs** Often, capitalization is absent in Dutch UGC. Our NE recognizer misses all those variants like “anneke”, “gronloh”, “paul” and “anneke gr’ohnloh”.

**Incomplete entities** Commenters introduce new named entities without ever using full names, e.g., “Pronk” (which

<sup>5</sup>An example of a smart pun: “*Anneke Kreunlo*” was used in one of the comments to refer to the Dutch singer Anneke Grönloh. “Kreun” and “Grön” sound similar in Dutch, and the verb *kreunen* is related to English *crooning* (sentimental singing style) and, moreover, means *to moan* in Dutch.

<sup>6</sup>In Dutch, “*Mevrouw*” and “*Mw.*” mean “Mrs.”

refers to former Dutch minister Jan Pronk).

### 3.3 Improving Named Entity Normalization

Based on the observed causes of errors we suggested five improvements of our baseline algorithm, which we list below.

1. *Pre-processing NEs*: We clean the entities and fix some of the NE boundary recognition errors. Specifically, we perform the following two steps before calling the function `NormalizeNE()`.
  - Clean up: remove non-alphabetical characters from both sides of an NE;
  - Gluing NEs: (Dutch-specific) to improve recognition of Dutch multi-word last names such as “van Gogh” or “ten Cate,” we glue two adjacent person names if they are separated by one of the standard Dutch infixes (*van, ten, de, van der, van ten, van t, vt*).
2. *N-gram NE normalization*: we use ngramming to split NEs that cannot be normalized. This is a modification of part of the baseline algorithm (lines 6–8 of Algorithm 1). Specifically, if an NE cannot be normalized, we split the NE into a set of all overlapping word ngrams and try to normalize each ngram using the function `NormalizeNE`. We proceed from longer to shorter ngrams, and when one of the ngrams is normalized successfully, all other ngrams that overlap with it are ignored.
3. *Person-name matching*: we improve our within-document coreference resolution (Algorithm 3), handling common variants of person names as they are used in UGC. Specifically, we use a set of heuristics (based on first and last names) and inexact string matching (based on Levenshtein edit distance) to identify variants of person names. The details of the improved name matching algorithm are given in the Appendix.
4. *Finding missing NEs*: we add a post-processing step to the baseline algorithm, in order to improve the recall of NEs. Specifically, after the recognition and normalization has been done (line 10 of Algorithm 1), we look for phrases which have not been recognized as NEs, but which are similar to those already recognized and normalized. We consider text segments outside the recognized NEs, and use a procedure similar to improvement 2 above: we split the text segments into word ngrams and try to resolve each ngram within the documents (calling the function `ResolveRef-InDoc`, Algorithm 3). Again, we proceed from longer to shorter ngrams, and ignore ngrams overlapping with those we detect as NEs. For efficiency, in this step we only consider ngrams that contain at most  $n + 2$  tokens, where  $n$  is the maximal number of tokens of a normalized form of NEs in the document.
5. *Normalizing new entities*: Wikipedia disambiguation pages do not cover all possible ambiguous entities, e.g., incomplete entity references such as “Pronk” (a Dutch surname). As noted in [4, 11], anchor texts in Wikipedia often cover these cases. Following [11], if an NE exactly matches the anchor text of some hyperlink in Wikipedia, we take as its normalized form the title of

the Wikipedia article (the target of the hyperlink) with the highest number of incoming links. This is a modification of the function Disambiguate of Algorithm 4.

In the following sections we describe the setup and experiments that we used to test the effect of the improvements listed above.

## 4. EXPERIMENTAL SETUP

### 4.1 Research Questions

In Section 3 we have described the baseline NEN algorithm and five improvements. Now we turn to set up our experiments in order to answer the following research questions.

- RQ1 How does the baseline NEN algorithm from the literature perform on UGC?
- RQ2 What is the overall performance achieved by the combined improvements proposed in Section 3.3?
- RQ3 How important are the proposed improvements, individually? That is, if we leave out one of them, how much do we lose in terms of effectiveness?
- RQ4 To what extent is our NEN algorithm work language dependent? I.e., does it achieve comparable scores across languages? While we do not have the resources to answer this question for all possible language pairs, we are keen to compare the performance of our NEN algorithm on two related languages: Dutch and English.

While our main interest lies with the effectiveness of our NEN algorithm, we will also include some results on efficiency.

### 4.2 Description of the data

As there is no standard data set available for assessing NEN on user generated comments, we decided to create our own data set. We selected 5 BBC<sup>7</sup> articles with comments from the “Have your say” section, and five articles with comments from online versions of two leading Dutch national newspapers.<sup>8</sup> For both languages, we selected articles reflecting the four main types of named entities: persons (one in politics, and one in show business), locations, organizations, and products and chose articles containing relatively many mentions of NEs, both in the article and the comments. On average, each real world entity was mentioned 5.6 and 5.3 times in the Dutch and English datasets, respectively.

Table 1 describes these 10 articles in detail. E.g., the Dutch article about a person in politics had 8,602 words in total, contained 196 comments, with a total of 508 mentions of named entities, which referred to 67 different real world entities.

We asked two assessors to assess each word in the articles with comments in our data set and to judge whether it is a reference to an NE, and in this case what its normalized form is. For this purpose we split the input in such a way that each token of the data occurs on a line; on the same

<sup>7</sup>URL: <http://news.bbc.co.uk/>.

<sup>8</sup>URL: <http://www.telegraaf.nl> and <http://www.ad.nl>

**Table 1: Description of the evaluation data. #Ws is the length in words (of the article plus comments); #Cs is the number of comments; #NEs is the total number of named entites occurring in the article plus comments; and #UNFs is the number of real world entities (thus having Unique Normal Forms) in the article plus comments.**

Category	#Ws	#Cs	#NEs	#UNFs
<i>Dutch language evaluation set</i>				
Persons (politics)	8,602	196	508	67
Persons (show business)	3,705	90	166	22
Locations	7,133	175	214	64
Organizations	2,516	47	190	34
Products	2,642	50	172	32
<i>English language evaluation set</i>				
Persons (politics)	4,265	62	241	43
Persons (show business)	2,230	30	80	16
Locations	4,050	67	257	49
Organizations	2,353	42	137	29
Products	2,324	39	163	26

line, following the token, an NE tag (if applicable) and the normalized form had to be provided by our assessors.

Our assessors had to use two kinds of tags for their annotation effort. For recognizing reference strings, the tags have a format similar to one known from chunking: an I denotes the first or inside term of an entity and an E denotes the last term. We used the four standard types of entities: person names (PER), organizations (ORG), locations (LOC) and miscellaneous names (MISC). For assigning normalized forms, tags have two kinds of format: normalized form (C), and Wikipedia entity (WIKI). If assessor was able to find a Wikipedia article about the entity, the tag WIKI was assigned. If no Wikipedia article could be found, assessors used the tag C, indicating within-document normalization. Here is an example:

```
Schumi   E-PER   WIKI-"Michael Schumacher"
DC       E-PER   WIKI-"David Coulthard"
,        0      0
Montoya  E-PER   WIKI-"Juan Pablo Montoya"
and      0      0
Ralf     E-PER   WIKI-"Ralf Schumacher"
team    0      0
Ross     I-PER
Braun    E-PER   C-"Ross Braun"
deserves 0      0
.        0      0
TRP     E-MISC  C-"XXXX"
```

If an assessor was not able to determine the normalized form, he was asked to use XXXX: this could be the case for obscure entities or incomplete names. The third column of the output of our assessors is taken to be the normalized form, e.g., “Ralf Schumacher” is the normalized form of the reference string “Ralf.”

On average, the assessment took about two hours per article with comments.

### 4.3 Evaluation Measures

The media analysis use case that motivated the work in this paper, suggest two important analysis dimensions for

each (article, comment-thread) pair: 1. who/what is mentioned? and 2. how often? Our evaluation measures reflect this interest. For both questions we define separate notions of recall and precision.

### Measuring real-world entities.

Fix an (article, comment-thread) pair. Let  $UNF_{gt}$  and  $UNF_{nen}$  denote the sets of real-world entities discovered by the assessors and by our NEN algorithm, respectively. We then define

$$entity-recall = \frac{|UNF_{nen} \cap UNF_{gt}|}{|UNF_{gt}|} \quad (1)$$

$$entity-precision = \frac{|UNF_{nen} \cap UNF_{gt}|}{|UNF_{nen}|}. \quad (2)$$

It is natural to order the sets of real world entities  $UNF_{gt}$  and  $UNF_{nen}$  by the number of mentions of each real world entity. Using that order we define  $entity-precision@n$  as the number of correctly identified entities within the first  $n$  returned entities divided by  $n$ .

### Measuring occurrences.

Let  $EM_{gt}$  and  $EM_{nen}$  denote the sets of triples (reference-string, normalized-form, position) discovered by the assessors and by our NEN algorithm, respectively. We define the set of true positives  $TP$  as:  $EM_{nen} \cap EM_{gt}$ . Now define the standard notions of precision and recall:

$$recall = \frac{|TP|}{|EM_{gt}|} \quad (3)$$

$$precision = \frac{|TP|}{|EM_{nen}|}. \quad (4)$$

We also want to measure the quality of normalization separately from the recognition step. Let  $FP_{norm}$  be the set of triples (reference-string, normalized-form, position) which are correctly recognized but incorrectly normalized. Thus  $FP_{norm} = \{(e, n, p) \in EM_{nen} \mid (e, n', p) \in EM_{gt} \text{ and } n \neq n'\}$ . We measure the accuracy of normalization using the following metric:

$$accuracy = \frac{|TP|}{|TP| + |FP_{norm}|}. \quad (5)$$

Because the last three measures work on the individual NE occurrences we use those to evaluate the effect of each of our five improvements.

## 4.4 Significance Testing

We perform significance testing only for accuracy (Eq. 5). There we compare two lists of triples and check for each triple if either they match perfectly with each other or there is no match. Thus we have a list of binary comparison numbers. In this case the McNemar test is appropriate. We look for significance at the  $p = 0.01$  level.

## 5. EXPERIMENTAL EVALUATION

For evaluation purposes, we used the test set specified in Section 4.2. We used two named entity recognition tools in our NEN algorithm, one based on [19] for the Dutch language data set, the other based on [5] for the English data set. For our NEN algorithm we used the English Wikipedia of August 2007, and the Dutch Wikipedia of November 2006.

**Table 2: Precision, recall and accuracy of baseline, full system (FS) and derived versions of the NEN algorithm. Ignore<sub>k</sub> denotes the full system FS but without the  $k$ -th improvement step. Significant improvements over the baseline are marked with <sup>▲</sup> ( $p = 0.01$ )**

	recall%	precision%	accuracy%
<i>Evaluation results on Dutch data set</i>			
Baseline	46	45	65
Ignore <sub>1</sub>	70	60	87 <sup>▲</sup>
Ignore <sub>2</sub>	64	53	79 <sup>▲</sup>
Ignore <sub>3</sub>	80	65	88 <sup>▲</sup>
Ignore <sub>4</sub>	70	<b>66</b>	88 <sup>▲</sup>
Ignore <sub>5</sub>	70	54	78 <sup>▲</sup>
FS	<b>81</b>	62	<b>89<sup>▲</sup></b>
<i>Evaluation results on English data set</i>			
Baseline	70	72	77
Ignore <sub>1</sub>	85	84	<b>90<sup>▲</sup></b>
Ignore <sub>2</sub>	79	75	83 <sup>▲</sup>
Ignore <sub>3</sub>	84	80	87 <sup>▲</sup>
Ignore <sub>4</sub>	84	<b>85</b>	89 <sup>▲</sup>
Ignore <sub>5</sub>	85	81	88 <sup>▲</sup>
FS	<b>89</b>	82	<b>90<sup>▲</sup></b>

**Table 3: Entity-precision, Entity-recall and Entity-precision@10 of the baseline and the full system of our NEN algorithm.**

	e-recall	e-precision	e-precision@10
<i>Evaluation results on Dutch data set</i>			
Baseline	51	21	68
FS	<b>69</b>	<b>34</b>	<b>85</b>
<i>Evaluation results on English data set</i>			
Baseline	73	54	68
FS	<b>78</b>	<b>67</b>	<b>84</b>

## 5.1 Results

The overall accuracy of the baseline NEN algorithm was 80% and 69% for Dutch news article and news articles plus comments, respectively, and 94% and 77% for English news articles and news articles plus comments, respectively. This provided the answer to the first research question, RQ1: the baseline NEN algorithm performed much worse on UGC than on edited text (news).

In order to answer the remaining research questions, we evaluated the baseline system, the extension with all improvements combined, and its derived versions by dropping one of the improvements. Table 2 lists the evaluation results of the variants of our NEN algorithm separated for the Dutch and the BBC data. Table 3 shows the *entity-recall*, *entity-precision* and *entity-precision@10* measurements. In these tables and in Section 5.2 below, **Baseline** denotes the baseline NEN algorithm, **FS** denotes the combined improvements (full system), and **Ignore<sub>k</sub>** denotes the absence of the  $k$ th improvement in the FS NEN algorithm, where  $k = 1, \dots, 5$  presents the improvement as numbered in Section 3.3, i.e.,

1. *Pre-processing NEs*,
2. *N-gram NE Normalization*,
3. *Person-name matching*,
4. *Finding missing NEs*,

### 5. Use Wikipedia link text for disambiguation.

Now we examine the evaluation results shown in the Tables 2 and 3 and compare the full system FS with the baseline algorithm, so as to answer our second research question, RQ2. FS performed much better than the baseline for both the Dutch and the English data set, according to all evaluation metrics. Also, FS is highly recall-oriented (81% vs 46% for Dutch, and 89% vs 70% for BBC) and more accurate than the baseline (89% vs 65% for Dutch, and 90% vs 77% for BBC).

In order to answer our third research question, RQ3, we examine Table 2 thoroughly, and zoom in on Figure 2. We see the impact of each improvement on the performance of the full system, by comparing derived versions of the full system  $\text{Ignore}_k$  with the full system FS itself. When we ignore the first improvement (*Pre-processing NEs*),  $\text{Ignore}_1$ , we find no effect on the accuracy for the English data set, but the accuracy drops a little on the Dutch data set. The reason is that comment threads of in the English data set tend to be cleaner than the comment threads of the Dutch articles. The results in Table 2 also show that this improvement is important in increasing recall for both data sets.

When we ignored the second improvement (*N-gram NE Normalization*), the third improvement (*Person-name matching*), or the fifth improvement (*Use Wikipedia link text for disambiguation*), then the NEN algorithm performed worse than the full system FS, not only in terms of accuracy, but also in terms of recall and precision. These three improvements are really important for the performance of our NEN algorithm.

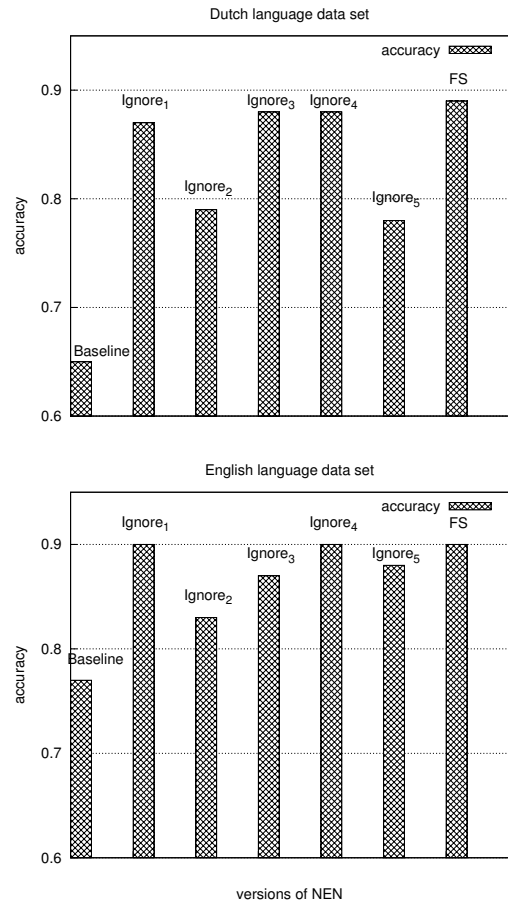
When we ignore the fourth improvement (*Finding missing NEs*), we find the highest precision (bold face numbers in Table 2). The reason for this is that in this improvement we greedily extract all kinds of reference strings. When we split the text into ngrams and use our own *Person-name matching* algorithm (that uses Levenshtein distance and our own heuristics) to handle variants of NEs, a string is sometimes incorrectly matched with an NE. E.g., “grapes” was considered as the last name of “Bill Gates”.

Finally, the NEN algorithm with all combined improvements and each of its derived versions (denoted by  $\text{Ignore}_k$ ) are significantly better than the baseline at  $p = 0.01$ .

## 5.2 Further Analysis

We have answered three of our four research questions from Section 5. Now we provide a further analysis of the experimental results, to answer the fourth research question concerning the performance on different languages. Notice that apart from the language specific rule concerning infixes in Dutch last names (part of the first improvement, in Section 3.3), the only difference between the Dutch and English language version of our NEN method is the corpus against which we normalize, Wikipedia. We examine the difference in accuracy between the baseline NEN algorithm, the NEN algorithm with combined improvements, and its derived algorithms, for Dutch and English. Consider Figure 2. First, while it is hard to draw very firm conclusions given the small number of articles we were able to annotate, we see that, mostly, the different versions of the system perform comparably for the two languages.  $\text{Ignore}_5$  (*Use Wikipedia link text for disambiguation*) seems to be a clear exception: on the Dutch data set the accuracy of the NEN algorithm is affected more strongly than on the English data set. The

**Figure 2: Accuracy of different versions of the NEN algorithm on the Dutch language data set, and English data set.  $\text{Ignore}_k$ : the full system minus the  $k$ th improvement step.**



reason for this seems clear: the English version of Wikipedia constitutes a much larger knowledge base than the Dutch version, with  $\sim 88,000$  vs  $\sim 13,000$  disambiguation pages.

Finally, on standard desktop hardware, the baseline NEN algorithm took on average 4.5s per document (news-article plus comments), and the full system FS and all other versions except  $\text{Ignore}_4$  took 15.2s. The fourth improvement (*Finding missing NEs*) is the most expensive step of the NEN algorithm, accounting for approx. 10s of the running time of the full system. Hence, the fourth improvement, although helpful to increase recall, is expensive in terms of running time.

## 6. CONCLUSION

Our aim in this paper was to create a named entity normalization algorithm that performs well on user generated content (UGC). For this purpose we started with a baseline NEN system from the literature, and found that it performed much worse on UGC than on edited news: 77% vs 94% accuracy on an English language data set and 65% vs 81% on a Dutch language data set.

We identified the following main sources of errors of the baseline system when applied to UGC: NE recognition errors (incorrect boundaries of named entities or missing NEs),

multiple ways of referring to the same entity, and ambiguous (out of context) references. We addressed these issues by proposing five improvements to the baseline NEN algorithm. Our experimental results showed that all improvements are important in increasing recall, precision and accuracy of the algorithm. While helpful in increasing the recall, the improvement that we introduced to cover missing NEs is expensive in terms of running time. The overall system can run on multiple languages, and the main source of differences in performance between languages seems to be the size of the underlying corpus against which named entities are normalized, Wikipedia.

In future work we will attempt to further improve the performance of our NEN algorithm by using context-aware named entity disambiguation, creating small entity-specific language models. In addition, we want to improve the underlying NER tools we use and to consider other measures of string similarity than we have used so far (edit distance) to handle misspellings of the person names.

## 7. REFERENCES

- [1] D. Ahn, J. van Rantwijk, and M. de Rijke. A cascaded machine learning approach to interpreting temporal expressions. In *Human Language Technologies 2007: Proc. ACL 2007*, pages 420–427, 2007.
- [2] C. L. Borgman and S. L. Siegfried. Getty’s synonyme and its cousins: A survey of applications of personal name-matching algorithms. *Journal of the American Society for Information Science*, 43(7):459–476, 1992.
- [3] A. M. Cohen. Unsupervised gene/protein named entity normalization using automatically extracted dictionaries. In *ACL-ISMB Workshop on Linking Biological Literature, Ontologies and Databases*, pages 17–24, 2005.
- [4] S. Cucerzan. Large-scale named entity disambiguation based on wikipedia data. In *EMNLP-CoNLL ’07*, pages 708–716, 2007.
- [5] F. de Meulder and W. Daelemans. Memory-based named entity recognition using unannotated data. In *Proceedings of CoNLL-2003, Edmonton, Canada*, pages 208–211, 2003.
- [6] A. Doan and A. Y. Halevy. Semantic-integration research in the database community. *AI Mag.*, 26(1):83–94, 2005.
- [7] X. Dong, A. Halevy, and J. Madhavan. Reference reconciliation in complex information spaces. In *Proc. SIGMOD ’05*, pages 85–96, 2005.
- [8] A. K. Elmagarmid, P. G. Ipeirotis, and V. S. Verykios. Duplicate record detection: A survey. *Transactions on Knowledge and Data Engineering*, 19(1):1–16, 2007.
- [9] D. Feng, E. Shaw, J. Kim, and E. Hovy. Learning to detect conversation focus of threaded discussions. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association of Computational Linguistics (HLT-NAACL)*, 2006.
- [10] S. Fissaha Adafre and M. de Rijke. Discovering missing links in Wikipedia. In *Proceedings of the Workshop on Link Discovery: Issues, Approaches and Applications (LinkKDD-2005)*, 2005.
- [11] M. Khalid, V. Jijkoun, and M. de Rijke. The impact of named entity normalization on information retrieval for question answering. In *Proc. ECIR 2008*, 2008.
- [12] J. Kim, G. Chern, D. Feng, E. Shaw, and E. Hovy. Mining and assessing discussions on the web through speech act analysis. In *Proceedings of the Workshop on Web Content Mining with Human Language Technologies at the 5th International Semantic Web Conference*, 2006.
- [13] W. Magdy, K. Darwish, O. Emam, and H. Hassan. Arabic cross-document person name normalization. In *CASL Workshop ’07*, pages 25–32, 2007.
- [14] G. S. Mann and D. Yarowsky. Unsupervised personal name disambiguation. In *Proceedings of the 7th Conference on Natural Language Learning at HLT-NAACL 2003*, pages 33–40, 2003.
- [15] R. Mihalcea and A. Csomai. Wikify!: Linking documents to encyclopedic knowledge. In *Proceedings of the 17th ACM Conference on Conference on Information and Knowledge Management*, pages 233–242, 2007.
- [16] G. Mishne. *Applied Text Analytics for Blogs*. PhD thesis, University of Amsterdam, Amsterdam, 2007.
- [17] C. Phua, V. Lee, and K. Smith. The personal name problem and a recommended data mining solution. In *Encyclopedia of Data Warehousing and Mining (2nd Edition)*. 2006.
- [18] Y. Song, J. Huang, I. G. Council, J. Li, and C. L. Giles. Efficient topic-based unsupervised name disambiguation. In *Proceedings of the 2007 Joint Conference on Digital Libraries*, pages 342–351, 2007.
- [19] E. F. Tjong Kim Sang. Memory-based named entity recognition. In *Proceedings of CoNLL-2002, Taipei, Taiwan*, pages 203–206, 2002.
- [20] M. Weimer, I. Gurevych, and M. Mehlhauser. Automatically assessing the post quality in online discussions on software. In *Proceedings of the ACL 2007 Demo and Poster Sessions*, pages 125–128, 2007.
- [21] W. Winkler. The state of record linkage and current research problems. Technical report, Statistical Research Division, U.S. Bureau of the Census, Washington, DC., 1999.
- [22] W. Zhou, C. Yu, N. Smalheiser, V. Torvik, and J. Hong. Knowledge-intensive conceptual retrieval and passage extraction of biomedical literature. In *SIGIR ’07*, pages 655–662, 2007.

## APPENDIX

### Person name matching algorithm

Algorithm 5 takes three arguments: a source string, a target string and a flag indicating whether the target string follows a person title (e.g., Mr, Mrs, etc.). The algorithm checks whether the target string is a variant of the source string.

---

#### Algorithm 5 Person name matching algorithm

---

```

Require: a source string  $S$ , a target string  $T$ , a flag  $L$ 
1: if  $T = S$ , or  $T = \text{firstName}(S)$ , or  $T = \text{lastName}(S)$  then
2:   return matched
3: end if
4: if  $\text{lastName}(S) \neq \text{NULL}$  then
5:   if  $\text{lastName}(T) = \text{NULL}$  then
6:     if  $L = \text{TRUE}$  and  $\text{lastName}(S)$  startsWith  $T$  then
7:       return matched
8:     else if  $\text{firstChar}(T) = \text{firstChar}(\text{lastName}(S))$  then
9:       return  $\text{isSimilar}(\text{lastName}(S), T)$   $\{\text{isSimilar}(\text{str1}, \text{str2})$  returns true if  $\frac{\text{editDist}(\text{str1}, \text{str2})}{\text{length}(\text{str2})} \leq 0.34\}$ 
10:    end if
11:   else if  $\text{firstName}(S)$  startsWith  $\text{firstName}(T)$  then
12:     if  $\text{lastName}(S)$  startsWith  $\text{lastName}(T)$  then
13:       return matched
14:     else
15:       return  $\text{isSimilar}(S, T)$ 
16:     end if
17:   else if  $\text{lastName}(S)$  startsWith  $\text{firstName}(T)$  then
18:     return  $\text{isSimilar}(\text{lastName}(S), T)$ 
19:   end if
20: else if  $\text{lastName}(T) = \text{NULL}$  and  $\text{firstChar}(S) = \text{firstChar}(T)$  then
21:   return  $\text{isSimilar}(S, T)$ 
22: end if

```

---