



LAWRENCE  
LIVERMORE  
NATIONAL  
LABORATORY

# Reconcile: A Coreference Resolution Research Platform

V. Stoyanov, C. Cardie, N. Gilbert, E. Riloff, D.  
Buttler, D. Hysom

November 25, 2009

## **Disclaimer**

---

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

This work performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344.

# Reconcile: A Coreference Resolution Research Platform

**Veselin Stoyanov**  
**Claire Cardie**

Dept. of Computer Science  
Cornell University  
Ithaca, NY

ves@cs.cornell.edu  
cardie@cs.cornell.edu

**Nathan Gilbert**  
**Ellen Riloff**

School of Computing  
University of Utah  
Salt Lake City, UT

ngilbert@cs.utah.edu  
riloff@cs.utah.edu

**David Buttler**  
**David Hysom**

Lawrence Livermore  
National Laboratory  
Livermore, CA

buttler1@llnl.gov  
hysom1@llnl.gov

## Abstract

Despite the availability of standard data sets and metrics, approaches to the problem of noun phrase coreference resolution are hard to compare empirically due to the different evaluation setting stemming, in part, from the lack of comprehensive coreference resolution research platforms. In this tech report we present Reconcile, a coreference resolution research platform that aims to facilitate the implementation of new approaches to coreference resolution as well as the comparison of existing approaches. We discuss Reconcile’s architecture and give results of running Reconcile on six data sets using four evaluation metrics, showing that Reconcile’s performance is comparable to state-of-the-art systems in coreference resolution.

## 1 Introduction

Noun phrase coreference resolution (or simply coreference resolution) is the problem of identifying all noun phrases (NPs) that refer to the same real-world entity in text. Coreference resolution is one of the fundamental problems of natural language processing (NLP) because of its usefulness for other NLP tasks (e.g. Morton (1999), Steinberger et al. (2007)) as well as the theoretical interest in understanding computational mechanisms involved in government, binding and linguistic reference. Figure 1 displays an example of coreferent noun phrases from a sample news article.

Several formal evaluations have been conducted for the coreference resolution task, and the data sets created for this evaluations have become standard benchmarks in the field (e.g., MUC and ACE). However, it is still frustratingly difficult to compare results across different coreference resolution systems. Reported scores for coreference resolution vary wildly across data sets, evaluation metrics, and system configurations.

We believe that one root cause of these disparities is the high cost of implementing an end-

to-end coreference resolution system. Coreference resolution is a complex problem, and successful systems must tackle a variety of non-trivial subproblems, which themselves require substantial implementation efforts. As a result, many researchers exploit gold-standard annotations from available data sets as a substitute for component technologies to solve these subproblems. Many published research results use gold standard annotations to identify the NPs (substituting for mention/markable detection), to distinguish the anaphoric NPs from the non-anaphoric NPs (substituting for anaphoricity determination), to identify named entities (substituting for named entity recognition), and to identify semantic types of nominals (substituting for semantic class identification). The use of gold standard annotations leads to an unrealistic evaluation setting, and makes it impossible to directly compare results against coreference resolvers that solve all of these subproblems from scratch.

The problem is further magnified by the availability of several (non-trivial) evaluation measures, and data sets that have substantially different task definitions and annotation formats. Understandably, research efforts rarely implement all of the evaluation metrics or evaluate their systems on all of the data sets. As a result, even when different approaches use the same system configuration, they may conduct experiments on different data sets or use different evaluation metrics, which again makes it impossible to directly compare results.

In addition, coreference is a pervasive problem in natural language processing, so in principle nearly any NLP application could benefit from having an effective coreference resolver. But the hefty price of creating a coreference resolver from

A state security court suspended a newspaper<sub>1</sub> critical of the government Saturday after convicting it<sub>2</sub> of publishing religiously inflammatory material, the newspaper<sub>3</sub>'s lawyer<sub>1</sub> said.

The Khartoum Monitor<sub>4</sub> may not appear for two months and must pay a fine<sub>1</sub> of 1 million Sudanese pounds (about US\$400), lawyer Ngor Olang Ngor<sub>2</sub> told The Associated Press. Failure to pay the fine<sub>2</sub> would result in a two-month extension of the suspension.

The sentence was the latest in a series of state actions against the Monitor<sub>5</sub>, the only English language daily in Sudan<sub>6</sub> and a leading critic<sub>7</sub> of conditions in the south of the country, where a civil war has been waged for 20 years.

Ngor<sub>3</sub> said the paper<sub>8</sub> would pay the penalty<sub>3</sub> but it<sub>9</sub> would also appeal the verdict.

He<sub>4</sub> said the court ruling stated that the verdict and sentence would be submitted to the National Press Council, a government-appointed body, for approval.

Figure 1: Example noun phrase coreference chains from a sample news article.

scratch often prohibits both researchers and NLP practitioners from incorporating coreference resolution in larger application systems.

To address these issues, we have created a platform for coreference resolution, called Reconcile, that can serve as a software infrastructure to support the creation, experimentation, and evaluation of learning-based coreference resolvers. Our goal in creating Reconcile was to build a coreference resolution system that would allow for different component technologies to be easily swapped in and out (e.g., different parsers, NER systems, or classification algorithms) and for new features to be easily created and plugged in. We also wanted a platform that would provide built-in capabilities to process many of the standard benchmark data sets, respecting their idiosyncratic formats and task definitions, as well as evaluation tools to produce scores that conform to multiple widely used evaluation metrics. This software infrastructure will allow researchers to focus on the new ideas and methodological improvements that they wish to investigate, without having to invest the substantial time and effort required to implement the various subcomponents themselves. Furthermore, this platform will allow researchers to more easily evaluate the impact of their ideas in a complete end-to-end coreference resolution system, with respect to multiple standard data sets and evaluation metrics. We hope that this software infrastructure will help to advance the state-of-the-art in confer-

ence resolution by enabling researchers to easily produce results that are directly comparable, yielding better insights into which ideas and techniques perform the best.

This technical report describes the Reconcile coreference resolution platform. Reconcile was designed with the following 6 desiderata in mind:

- implement the basic underlying software architecture of contemporary state-of-the-art learning-based coreference resolution systems.
- support text analysis on most of the standard coreference resolution data sets.
- implement many widely used coreference resolution scoring metrics.
- exhibit state-of-the-art coreference resolution performance (i.e., it can be configured to create a coreference resolver that achieves performance levels that are close to the best results previously reported in the literature).
- can be extended with new methods and features with little effort.
- is relatively fast and easy to configure and run.

While several other coreference systems are publicly available (e.g., Poesio and Kabadjov (2004), Qiu et al. (2004) and Versley et al. (2008)),

none of them meets all six of these requirements (see Related Work section for more details). Reconcile is a modular software platform that abstracts the basic architecture of the most contemporary supervised learning-based coreference resolution systems (e.g. Soon et al. (2001), Ng and Cardie (2002), McCallum and Wellner (2004), Ng (2008)) and achieves performance comparable to the current state-of-the-art on several benchmark data sets commonly used to evaluate coreference resolution systems. Additionally, due to its modularity, Reconcile can be easily reconfigured to use different algorithms, features, preprocessing elements, evaluation settings and metrics.

This report explains how Reconcile is organized, what it can do, and how it performs on several MUC and ACE data sets. Section 2 first explains how Reconcile differs from other publicly available coreference resolution software systems. Section 3 presents Reconcile's organization and describes its subcomponents. Finally, Section 4.3 shows experimental results for Reconcile on a variety of different data sets and with several evaluation metrics.

## 2 Related Work

Several coreference resolution systems are currently publicly available. JavaRap (Qiu et al., 2004) is an implementation of the Lappin and Leass's (1994) Resolution of Anaphora Procedure (RAP). JavaRap resolves only pronouns and, thus, it is not directly comparable to Reconcile. GUITAR (Poesio and Kabadjov, 2004) and BART (Versley et al., 2008) are both modular systems that target the complete the full coreference resolution task. As such, both of these systems meet the majority of the desiderata set forth in Section 1. BART can be considered an alternative to Reconcile, although we believe that Reconcile's approach is more flexible than BART. The architecture and system components of Reconcile (including a comprehensive set of implemented features that draw on the expertise of state-of-the-art supervised learning approaches, such as Soon et al. (2001) and Ng and Cardie (2002)) result in performance that is closer to the state-of-the-art.

The problem of coreference resolution has received much research attention, resulting in a rich variety of approaches, algorithms and features. Reconcile's architecture has been modeled after

the coreference system of Soon et al. (2001) because of its simplicity and its proven competitive performance. In addition, Reconcile includes in its base distribution the refined and extended feature set of Ng and Cardie (2002).

Reconcile has been modeled after typical supervised learning approaches to coreference resolution because of the popularity and relatively good performance of these systems. However, there have been many other approaches to coreference resolution, including recent unsupervised and semi-supervised approaches (e.g. Ng (2008) and Haghighi and Klein (2007)), structured approaches (e.g. McCallum and Wellner (2004) and Finley and Joachims (2005)), competition approaches (e.g. Yang et al. (2003)) and a bell-tree search approach Luo et al. (2004). Most of the approaches can be directly implemented in Reconcile, while some (such as the latter two) require more serious modifications to the system, but can, nevertheless, be run in the system.

In Section 4, we position Reconcile's performance with respect to the state-of-the-art by comparing Reconcile's results on several standard data sets with the most comparable state-of-the-art results reported in the literature.

## 3 System Description

### 3.1 Overall Architecture

Reconcile was designed to be a research test-bed capable of implementing most current state-of-the-art approaches to coreference resolution. Thus, the basic architecture of the system includes five major steps, which are performed by most coreference resolution systems (e.g., (Soon et al., 2001), (Ng and Cardie, 2002), (Ng, 2008), (Haghighi and Klein, 2007), (McCallum and Wellner, 2004) and (Finley and Joachims, 2005)). Starting with a corpus of text documents together with a manually annotated coreference resolution answer key<sup>1</sup>, Reconcile performs the following 5 steps, in order:

1. **Preprocessing.** All documents are passed through a series of (external) linguistic processors such as tokenizers, part-of-speech

---

<sup>1</sup>When Reconcile is applied to texts to perform coreference resolution, Reconcile does not require the presence of a manually annotated answer key. For the purposes of this report, however, we will assume that Reconcile is being used both to create a coreference resolver and to evaluate its performance on standard data sets, which requires answer keys

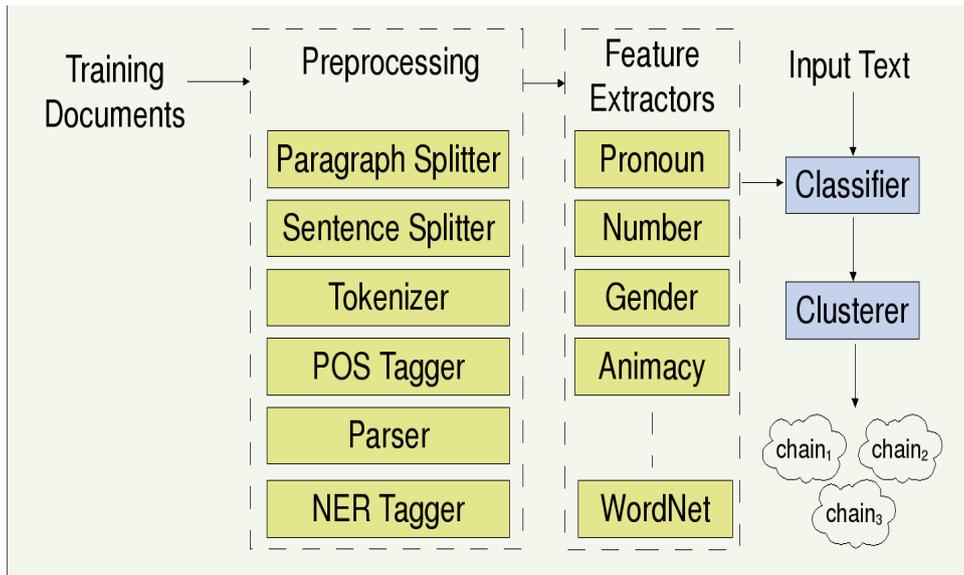


Figure 2: Diagram of the Reconcile tool chain.

taggers, syntactic parsers, etc. These components produce annotations of the text in the documents.

2. **Feature generation.** Using the annotations produced during preprocessing, Reconcile creates a set of features. Each feature is an attribute that characterizes a pair of noun phrases. For example, a feature might denote whether two noun phrases agree in number, or whether two noun phrases share any words in common.
3. **Classification** Since Reconcile adopts a supervised learning framework, the system can be used both to train a classifier (using training texts paired with their answer keys) and to apply the classifier to new text documents. The features are used to create one feature vector for each pair of noun phrases. When performing coreference resolution, the classifier is given the feature vector representing a noun phrase pair and it returns a score indicating how likely it is that the NPs are coreferent.
4. **Clustering.** Based on the scores output by the classifier, Reconcile uses a clustering algorithm to form the final set of coreference clusters (chains).
5. **Scoring.** Finally, Reconcile runs scoring algorithms that compare the coreference chains produced by the system to the gold standard

coreference chains in the answer key. (This is only possible if the documents have corresponding answer keys.)

Note that some structured coreference resolution algorithms (e.g., (McCallum and Wellner, 2004) and (Finley and Joachims, 2005)) combine the classification and clustering steps above. Reconcile can easily accommodate this modification.

Each of the five major steps described above can include multiple substeps and invoke several components. Reconcile allows various subsystems to be introduced or removed according to the desires of the user. Reconcile's standard distribution comes with a comprehensive set of implemented components. Additionally, the user is allowed to implement and introduce her own subsystems/components by using some simple API's. Next, we describe the set of implemented components for each of the five step that come pre-packaged in the standard Reconcile distribution, giving more details of the steps in the process.

### 3.2 Component Subsystems

One of the primary goals in developing Reconcile was to support a plug-and-play architecture where different components could be easily swapped in and out. Toward this end, we have included several different software options for many of the sub-tasks that Reconcile performs. The components that are included in the current release of Reconcile are described below, but it should also be easy

for anyone to plug in their own favorite component to perform any of these tasks.

### 3.2.1 Preprocessing

Reconcile performs a series of preprocessing steps and supports a variety of different preprocessing tools:

**Sentence Splitting:** Reconcile offers two different sentence segmentation options: The University of Illinois Urbana-Champaign sentence segmentation tool, and the OpenNLP project’s sentence splitter (Baldrige, J., 2005). The UIUC segmentation tool is a very fast, heuristic approach. The OpenNLP sentence detector is based on a maximum entropy classifier with a standard English model that can be retrained by the user if desired.

**Tokenizing:** The tokenizer is responsible for identifying word boundaries based on whitespace and punctuation to produce tokens that are used by the downstream components, such as parsers and part of speech taggers. The tokenizer included in the Reconcile distribution is the OpenNLP project’s maximum entropy tokenizer with a standard English model.

**Part of Speech Tagging:** The Reconcile distribution includes the OpenNLP maximum entropy part-of-speech (POS) tagger. Additionally, POS tags are produced by the two parsers listed below, so POS tags could be easily obtained from the parsers.

**Parsing:** Reconcile currently includes two probabilistic context-free grammar (PCFG) full parsers, the Stanford (Klein and Manning, 2003) parser and the Berkeley (Petrov and Klein, 2007) parser. The standard models provided with each of these parsers have been trained on Wall Street Journal data. The Sundance shallow parser (Riloff and Phillips, 2004) is also available for use, although some of the features supported by Reconcile (see the Appendix) require the use of a full parser.

**Named Entity Recognition:** Reconcile supports two named entity recognizers: the OpenNLP Maximum-Entropy NER system, and the Stanford Conditional Random Field NER system (Finkel et al., 2004). All available models for both systems have been trained on Wall Street Journal texts. In addition, in-house trained models are provided to

improve date and time recognitions.

**Coreference Element Identification:** The MUC and ACE evaluations used different criteria to define the set of NPs<sup>2</sup> that can participate in the coreference relation. For the sake of generality, we will use the term **coreference element (CE)** to refer to the set of linguistic expressions that participate in the coreference relation. We define the term CE to be roughly equivalent to (a) the notion of *markable* with respect to the MUC task definition and (b) the structures that can be *mentions* with respect to the ACE task definition.

Reconcile includes CE extractors that conform to both the MUC and ACE specifications. These CE extractors were developed in-house and rely on both parsing and named entity recognition.<sup>3</sup> At a high level, both the MUC and ACE evaluations define NPs as nouns, pronouns, and noun phrases. However, the MUC definition excludes (1) “nested” named entities (NEs) (e.g. “America” in “Bank of America”), (2) relative pronouns, and (3) gerunds, but allows (4) nested nouns (e.g. “union” in “union members”). The ACE definition, on the other hand, includes relative pronouns and gerunds, excludes **all** nested nouns that are not themselves NPs, and allows premodifier NE mentions of geo-political entities and locations, such as “Russian” in “Russian politicians”.

Another major distinction is that ACE restricts CEs to be entities belonging to one of seven semantic classes: person, organization, geo-political entity, location, facility, vehicle, and weapon. MUC has no semantic restrictions.

Figure 3 shows the detailed algorithm that Reconcile uses to collect all *coreference elements* from a document.

### 3.2.2 Classification

Most supervised learning approaches for coreference resolution rely on a pairwise function to determine whether or not two NPs are coreferent. A classifier is trained on NP pairs that are labeled as either *coreferent* or *not coreferent*, and can then be applied to new NP pairs to determine whether that pair is coreferent. Reconcile adopts this pairwise classification framework.

<sup>2</sup>Strictly speaking, both the MUC and ACE evaluations define NP coreference resolution over structures that are not always NPs, but both task definitions are generally described as NP coreference resolution.

<sup>3</sup>Our CE extractors were modeled after the extractors initially developed by Ng and Cardie (2002).

```

procedure ExtractCorefElements(D, NPs, NEs){
  A text document D
  NPs := all noun phrases in D.
  NEs := the set of named entities in D.

  /* Filtering step. */
  foreach noun phrase np ∈ NPs do
  if np contains a preposition or embedded clause:
    Remove np from NPs
  if np is the child of another np:
    Remove np from NPs

  /* Look for overlapping NPs and NEs */
  foreach noun phrase np ∈ NPs do
    foreach named entity ne ∈ NEs do
      if a np bytespan overlaps with a named entity ne bytespan:
        Expand np's bytespan so that it subsumes ne.
      else if ne is a proper substring of np:
        Add ne to NPs

  /* Add in named entities that are not part of the NP set. */
  foreach named entity ne in NEs do
  /* The contains function returns true if the second argument is a substring of the first. */
    if ne ∉ NP and ∃np ∈ NPs, ne ∈ contains(np, ne):
      Add ne to NPs

  /* Add in nested noun phrases that are not already in the set. */
  foreach noun phrase np ∈ NPs do
    if np contains a nested noun phrase nnp ∉ NPs:
      Add nnp to NPs
  return NPs
}

```

Figure 3: The algorithm for generating *coreference elements*.

To train a classifier, Reconcile creates a feature vector for every pair of NPs in a document. Therefore, given a document that contains  $n$  noun phrases, the total number of feature vectors created will be  $\frac{n^2-n}{2}$ . Once the classifier has been trained, it can be given pairwise feature vectors for unseen documents and it will output a probability indicating how likely it is that the two NPs are coreferent.

We developed an extensive set of features, based largely on the feature sets that have been used by other successful coreference resolution systems, most notably the Soon et al. (Soon et al., 2001) and Cardie & Ng (Soon et al., 2001) systems. A comprehensive list and description of the 89 features that are currently available in Reconcile can be found in the Appendix. It should be noted that some of the features rely on properties of NPs pairs, while others rely on individual NPs. For example, the Gender feature has a value of 'Y' if the two NPs in the pair have the same gender, whereas the Quote features simply determine whether an individual NP is surrounded by quotes. Features that look at individual NPs are cached so they are not recomputed every time an NP participates in a pair.

The current Reconcile release supports a wide variety of learning algorithms. Most of these learning algorithms are included via the Weka

toolkit, which is interfaced and included as part of the Reconcile distribution. Additionally, Reconcile includes interfaces for the libSVM (Chang and Lin, 2001) package and SVM<sub>light</sub> (Joachims, 2002), both of which are implementations of Support Vector Machines (SVMs).

### 3.2.3 Clustering Algorithms

The last step of the coreference resolution process is clustering. The goal of the clustering algorithm is to consolidate the pairwise decisions of the classifiers and produce a coherent clustering of the NPs in a document. Three different clustering algorithms are implemented in Reconcile:

1. Single-link Clustering (Transitive Closure): This is a simple clustering algorithm that performs transitive closure of the clustering relation. In other words, the single-link clusterer groups together all NPs that are connected by a path of coreferent links.
2. Best-First: Most classifiers return a value that can be interpreted as a “confidence measure.” For instance, an SVM may return a distance from the hyperplane for each pair. Best-first clustering uses the classifier’s confidence value to cluster each noun phrase with its most confident antecedent, given that the value for the most confident antecedent is above a pre-specified threshold.

3. Most Recent First: This clustering algorithm pairs each noun phrase with the single most recent antecedent that is labeled as coreferent. This algorithm follows the intuition that the coreference relation is local.

### 3.2.4 Scoring

The coreference chains produced by the system can be compared to gold standard answer keys to quantitatively assess their quality. Toward that end, Reconcile implements three widely-used automatic scoring functions, which are described below.

**MUC score.** The MUC scoring algorithm (Vilain et al., 1995) computes the  $F_1$  score (harmonic mean) of precision and recall based on the identification of unique coreference links.

**$B^3$  score.** The  $B^3$  algorithm (Bagga and Baldwin, 1998) computes a precision and recall score for each coreference element ( $ce$ ) using the following formulas:

$$\begin{aligned} \text{precision}(ce) &= |R_{ce} \cap K_{ce}| / |R_{ce}| \\ \text{recall}(ce) &= |R_{ce} \cap K_{ce}| / |K_{ce}|, \end{aligned}$$

where  $R_{ce}$  is the coreference chain to which  $ce$  is assigned in the response (i.e. the system-generated output) and  $K_{ce}$  is the coreference chain that contains  $ce$  in the key (i.e. the gold standard). Precision and recall for a set of documents are computed as the mean over all coreference elements (CEs) in the documents. The final  $B^3$  score is the  $F_1$  (harmonic mean) of precision and recall.

**$B^3$  Implementation Issues.** Unlike the MUC score, which counts links between CEs,  $B^3$  presumes that the gold standard and the system response are clusterings over the same set of NPs. This, of course, is not the case when the system automatically identifies the CEs, so the scoring algorithm requires a mapping between extracted and annotated NPs. We will use the term *twin*( $ce$ ) to refer to the unique annotated/extracted NP to which the extracted/annotated NP is matched. We say that a NP is *twinless* (has no twin) if no corresponding NP is identified. A twinless extracted NP signals that the resolver extracted a spurious NP, while an annotated NP is twinless when the resolver fails to extract it. Unfortunately, it is unclear how the  $B^3$  score should be computed for twinless NPs. Bengtson and Roth (2008) simply discard twinless NPs, but this solution is likely too lenient — it does no punishment for mistakes on

twinless annotated or extracted NPs and it would be tricked, for example, by a system that extracts only the NPs about which it is most confident.

We propose two different ways to deal with twinless NPs for  $B^3$ . One option,  $B^3_{all}$ , retains all twinless extracted NPs. It computes the precision as above when  $ce$  has a twin, and computes the precision as  $1/|R_{ce}|$  if  $ce$  is twinless. (Similarly,  $\text{recall}(ce) = 1/|K_{ce}|$  if  $ce$  is twinless.)

The second option,  $B^3_0$ , discards twinless extracted NPs, but penalizes recall by setting  $\text{recall}(ce) = 0$  for all twinless annotated NPs. Thus,  $B^3_0$  presumes that all twinless extracted NPs are spurious.

**CEAF Score.** Luo’s (2005) CEAF score (for Constrained Entity-Alignment F-Measure) is a coreference resolution evaluation metric resembling the ACE score. Similar to ACE, CEAF relies on a measure of how well a response cluster matches an answer key cluster and computes an optimal mapping between answer key and response clusters. CEAF differs from ACE in that it computes recall by dividing the score of the optimal match by the score for mapping the key to itself (i.e. the maximum is 1) and precision by dividing by the score of matching the response to itself. The reported CEAF score is the  $F_1$  score (harmonic mean of precision and recall).

Luo (2005) suggests several functions to score the goodness of the match of a key cluster  $A$  and response cluster  $B$ . We borrow one of these functions:  $\phi(A, B) = (2 * |A \cap B|) / (|A| + |B|)$ . In other words, the score for the match is the number of items the two clusters have in common proportional to the combined size of the two clusters.

## 4 Evaluation

In this section, we present experimental results to compare the performance of a coreference resolver created with Reconcile against other state-of-the-art supervised learning-based coreference resolvers.

### 4.1 Data Sets

For evaluation we use the six most commonly used coreference resolution data sets. Two of those are from the MUC conferences (MUC-6, 1995; MUC-7, 1997) and four are from the Automatic Content Evaluation (ACE) Program (NIST, 2004). For ACE, we use only the newswire portion because it is closest in composition to the MUC corpora.

		MUC6	MUC7	ACE-2	ACE03	ACE04	ACE05
1. 0.5 THRESHOLD	<i>MUC</i>	70.40	58.20	65.76	66.73	56.75	64.30
	<i>B<sup>3</sup>all</i>	69.91	62.88	77.25	77.56	73.03	72.82
	<i>B<sup>3</sup>0</i>	68.55	62.80	76.59	77.27	72.99	72.43
2. THRESHOLD ESTIMATION	<i>MUC</i>	68.50	62.80	65.99	67.87	62.03	67.41
	<i>B<sup>3</sup>all</i>	70.88	65.86	78.29	79.39	76.50	73.71
	<i>B<sup>3</sup>0</i>	68.43	64.57	76.63	77.88	75.41	72.47
3. OPTIMAL THRESHOLD	<i>MUC</i>	71.20	62.90	66.83	68.35	62.11	67.41
	<i>B<sup>3</sup>all</i>	72.31	66.52	78.50	79.41	76.53	74.25
	<i>B<sup>3</sup>0</i>	69.49	64.64	76.83	78.27	75.51	72.94

Table 1: Reconcile’s performance on six standard data sets.

System		MUC6	MUC7
Soon et al. (2001)	<i>MUC</i>	62.6	60.4
Ng and Cardie (2002)	<i>MUC</i>	70.4	63.4
Yang et al. (2003)	<i>MUC</i>	71.3	60.2

Table 2: Reported results for state-of-the-art coreference systems using comparable evaluation settings.

Statistics for each of the data sets are shown in Table 3. When available, we use the standard test/train split. Otherwise, we randomly split the data into a training and test set following a 70/30 ratio.

dataset	docs	CEs	chains	CEs/ch	tr/tst split
MUC6	60	4232	960	4.4	30/30 (st)
MUC7	50	4297	1081	3.9	30/20 (st)
ACE-2	159	2630	1148	2.3	130/29 (st)
ACE03	105	3106	1340	2.3	74/31
ACE04	128	3037	1332	2.3	90/38
ACE05	81	1991	775	2.6	57/24

Table 3: Dataset characteristics including the number of documents, annotated CEs, coreference chains, annotated CEs per chain (average), and number of documents in the train/test split. We use *st* to indicate a standard train/test split.

## 4.2 The *Reconcile<sub>ACL09</sub>* Configuration

We created a coreference resolution system using Reconcile, which we will call *Reconcile<sub>ACL09</sub>* because this system was used to conduct the experiments presented in (Stoyanov et al., 2009). *Reconcile<sub>ACL09</sub>* is modeled after the state-of-the-art system of Ng and Cardie (2002), using many of the same features. The *Reconcile<sub>ACL09</sub>* system was configured using the following subsystems and options:

### 1. Preprocessing

#### (a) Sentence Segmentation:

*UIUC segmentation tool*

#### (b) Tokenizer: *OpenNLP*

#### (c) POS Tagger: *OpenNLP*

#### (d) Parser: *Berkeley parser*

#### (e) Named Entity Recognizer: *Stanford NER system w/ additional LLNL models*

- Feature Set<sup>4</sup>** - SoonStr, ProStr, ProComp, PNStr, WordsStr, WordOverlap, Modifier, PNSubstr, WordsSubstr, Pronoun1, Pronoun2, Definite1, Definite2, Demonstrative2, Embedded1, Embedded2, InQuote1, InQuote2, BothProperNouns, BothEmbedded, BothInQuotes, BothPronouns, BothSubjects, Subject1, Subject2, Appositive, MaximalNP, Animacy, Gender, Number, SentNum, ParNum, Alias, IAntes, Span, Binding, Contraindices, Syntax, ClosestComp, Indefinite, Indefinite1, Prednom, Pronoun, ContainsPN, Constraints, ProperNoun, Agreement, ProperName, WordNetClass, WordNetDist, WordNetSense, Subclass, RuleResolve, ProResolve, SameSentence, ConsecutiveSentences, AlwaysCompatible, SameParagraph, HeadMatch, PairType, Quantity, WNSynonyms

### 3. Classifier

<sup>4</sup>Full descriptions of the features are found in the Appendix.

4. *Perceptron learning algorithm* - from the Weka toolkit (Witten and Frank, 2005). Run for 15 iterations.
5. **Clustering** - *Single-link Clustering* - The clustering thresholds were tuned by cross validation of the training set as described in Section 4.

### 4.3 Experimental Results

Table 1, box 1 shows the performance of *Reconcile*<sub>ACL09</sub> using a default (0.5) threshold for the coreference classifier. The MUC score is highest for the MUC6 data set, while the four ACE data sets show much higher  $B^3$  scores as compared to the two MUC data sets. The latter occurs because the ACE data sets include singletons.

The classification threshold, however, can be gainfully employed to control the trade-off between precision and recall. This has not traditionally been done in learning-based coreference resolution research — possibly because there is not much training data available to sacrifice as a validation set. Nonetheless, we hypothesized that estimating a threshold *from just the training data* might be effective. Our results (THRESHOLD ESTIMATION box in Table 1) indicate that this indeed works well.<sup>5</sup> With the exception of MUC6, results on all data sets and for all scoring algorithms improve; moreover, the scores approach those for runs using an optimal threshold (box 3) for the experiment as determined by using the **test set**.

One of the goals of *Reconcile* was to provide a platform that can be used to build to assemble a coreference resolution system that exhibits performance that is comparable to the state-of-the-art. Unfortunately, due to the difficulties outlined in Section 1, results presented here can be compared only to a limited number of scores reported in the literature. Table ?? lists some of the best scores that are directly comparable to our system. Based on these comparisons, *Reconcile*<sub>ACL09</sub> does indeed exhibit performance levels that are comparable to the previously reported scores for state-of-the-art supervised learning-based coreference solvers.

---

<sup>5</sup>All experiments sample uniformly from 1000 threshold values.

## Appendix: Feature Set

The following 89 features are currently available in Reconcile:

1. Agreement - If both NPs agree on both number and gender, the value returned is 'C' for Compatible. If they disagree on either number or gender, 'I' or Incompatible is returned. 'NA' is returned if no gender or number information is present for one or both of the NPs in question.
2. Alias - If one NP is an alias of the other, this returns 'C' otherwise 'I'. Here we take Alias to mean a variety of things, such as different ways of representing the same date, monetary value or number. Also, information regarding semantic type is used to determine if both NPs are really two different names of the entity, such as acronyms or abbreviations.
3. AlwaysCompatible - Always returns 'C'.
4. AlwaysIncompatible - Always returns 'I'.
5. AnaMed - This feature returns:  $m - \frac{med(np_1, np_2)}{m}$ , where  $m = len(np_1)$  and  $med$  is the minimum edit distance.
6. Animacy - If both NPs agree of animacy, returns 'C', else returns 'I'. Uses semantic type information.
7. AnteMed - Nearly the same as AnaMed, but focuses on the antecedent. Returns  $m - med(np_1, np_2)/m$  where  $m = len(np_2)$ .
8. Appositive - Return 'C' if both NPs are in an appositive construction, else returns 'I'. Uses semantic type information.
9. Binding - Returns 'C' if both NPs do not violate conditions B and C in Chomsky's binding theory, else 'I'. In short, element  $\alpha$  binds elements  $\beta$  if and only if  $\alpha$  c-commands  $\beta$ . A node  $X$  c-commands node  $Y$  if and only if:
  - $X$  does not dominate  $Y$
  - $Y$  does not dominate  $X$
  - The first branching node dominating  $X$  also dominates  $Y$ .
10. BothDefinites - This feature checks for the appearance of the article 'the' in each NP. If both start with the article, it returns 'C', if only one contains the article, 'NA', else returns 'I'.
11. BothEmbedded - Returns 'C' if both NPs are embedded NPs, 'NA' if only one is, and 'I' if neither are.
12. BothInQuotes - Returns 'C' if both NPs are inside of quotes, 'NA' if only one is, and 'I' if neither are.
13. BothMCSubjects - Returns 'C' if both NPs are in the main clause, 'NA' if only one is, and 'I' if neither are.
14. BothPronouns - Returns 'C' if both NPs are pronouns, 'NA' if only one is, and 'I' if neither are.
15. BothProperNames - Returns 'C' if both NPs are proper names, 'NA' if only one is, and 'I' if neither are.
16. BothProperNouns - Returns 'C' if both NPs are proper nouns, 'NA' if only one is, and 'I' if neither are.
17. BothSubjects - Returns 'C' if both NPs are in the subject position, 'NA' if only one is, and 'I' if neither are.
18. ClosestComp - This feature checks the semantic compatibility between the two NPs. If they are compatible and the closest, it returns 'C', else 'I'. Uses semantic and named entity information.
19. ConsecutiveSentences - Returns 'C' if the NPs are in consecutive sentences, otherwise returns 'I'.
20. Constraints - Checks the compatibility of the GENDER, NUMBER, CONTRAINDICES, ANIMACY, PRONOUN and CONTAINSPN, if all are compatible (or at least not Incompatible types for the last four features), then it returns 'C', else 'I'.
21. ContainsPN - This feature checks that both NPs contain proper names and contain no words in common, if this is true, it returns 'I', else 'C'.
22. Contraindices - The following constraints are implemented for this feature: (1) two NP's

- separated by a preposition cannot be coindexed and (2) two non-pronominal NP's separated by a non-copular verb cannot be coindexed. If the two NP's violate these conditions, then returns 'I', else it is 'C'.
23. Definite1 - Returns 'Y' if the first NP starts with 'the', else 'N'.
  24. Definite2 - Returns 'Y' if the second NP starts with 'the', else 'N'.
  25. Demonstrative2 - Returns 'Y' if the second NP starts with a demonstrative, i.e. this, that, these and those. Returns 'N' otherwise.
  26. DocNo - The number of the document from which the NPs came from. This is only used as a bookkeeping measure. It is not used for training or testing.
  27. Embedded1 - 'Y' if the first NP is an embedded or nested NP, else 'N'.
  28. Embedded2 - 'Y' if the second NP is an embedded or nested NP, else 'N'.
  29. FullString1 - The full string of the first NP.
  30. FullString2 - The full string of the second NP.
  31. Gender - If the two NPs agree in gender, then this feature returns 'C', and 'I' if they disagree. If the gender information is not known, then it returns 'NA'.
  32. GramRole1 - The grammatical role of the first NP.
  33. GramRole2 - The grammatical role of the second NP.
  34. HeadMatch - Checks for matching head noun between the two NPs, if they match, returns 'C', else 'I'.
  35. HeadNoun1 - The text of head noun of the first NP.
  36. HeadNoun2 - The text of head noun of the second NP.
  37. IAntes - Returns 'Y' if the first NP is part of a quoted string.
  38. ID1BeginOffset - The beginning bytespan offset of the first NP.
  39. ID1EndOffset - The ending bytespan offset of the first NP.
  40. ID1 - The identification number of the first NP, this is a bookkeeping feature and is not used in training nor testing.
  41. ID2BeginOffset - The beginning bytespan offset of the second NP.
  42. ID2EndOffset - The ending bytespan offset of the second NP.
  43. ID2 - The identification number of the second NP, this is a bookkeeping feature and is not used in training nor testing.
  44. Indefinite1 - 'I' if the second NP is an indefinite and is not an appositive, 'C' otherwise.
  45. Indefinite -
  46. InQuote1 - Returns 'Y' if the first NP is part of a quoted string.
  47. InQuote2 - Returns 'Y' if the second NP is part of a quoted string.
  48. instClass - Weka bookkeeping feature, used to determine if the current pair is a coreferent instance or not. Must be included in any feature list.
  49. MaximalNP - If both NPs have the same maximal NP projection, then return 'I', else returns 'C'.
  50. MCSubject1 - Returns 'Y' if the first NP is the subject of the main clause.
  51. MCSubject2 - Returns 'Y' if the second NP is the subject of the main clause.
  52. Modifier - If the pronominal modifiers of one np are a subset of the pronominal modifiers of the other nps, then returns 'C', else 'I'.
  53. Number - If the two NPs agree in number, then this feature returns 'C', and 'I' if they disagree. If the gender information is not known, then it returns 'NA'.
  54. PairType - Encodes the type of the np pair, i.e., if the pair are Proper nouns, pronouns, definite or indefinite.
  55. ParNum - The distance between noun phrases in terms of paragraphs.

56. PNStr - If both NPs are proper names and the same string, then 'C', else 'I'.
57. PNSubstr - If both NPs are proper names and one is a substring of the string, then 'C', else 'I'.
58. Prednom - If the NPs form a predicate nominal construction, then 'C', else 'I'. An example: "Barack Obama is the U.S. president.", *the U.S. President* is acting as the predicate nominal.
59. ProComp - If both NPs are what we call "comparable" pronouns, such as *he* and *his*.
60. ProEquiv - If both NPs are pronouns and agree in GENDER and NUMBER and appear in consecutive sentences, then 'I', else 'C'.
61. Pronoun1 - If NP1 is a pronoun, return 'Y', else 'N'.
62. Pronoun2 - If NP2 is a pronoun, return 'Y', else 'N'.
63. Pronoun - If NP1 is a pronoun and NP2 is not, then return 'I', else 'C'.
64. ProperName - If both NPs are proper names and share no words in common, then returns 'I', else 'C'.
65. ProperNoun - If both NPs are proper nouns and share no words in common, then returns 'I', else 'C'.
66. ProResolve - If one NP is a pronoun and the other NP is its antecedent according to a rule-based algorithm, then 'C', else 'I'.
67. ProStr - Return 'C' if both NPs are pronouns and their strings match exactly, otherwise 'I'.
68. Quantity - Returns 'C' if the two NPs form the pattern  $\text{;sum;}$  of  $\text{;money;}$  (e.g. loss of 1 million), else 'I'.
69. RuleResolve - If the two NPs are coreferent according to a rule-based algorithm, then 'C', else 'I'.
70. SameParagraph - The NPs are found in the same paragraph, then return 'Y', else 'N'.
71. SameSentence - The NPs are found in the same sentence, returns 'Y', else 'N'.
72. SentNum - The distance between the two NPs in terms of sentences.
73. SoonStr - If after discarding determiners, the strings two NPs match, then return 'C', else 'I'.
74. SoonStrNonPro - If both NPs are non-pronominal and after discarding determiners the two NPs match, then return 'C', else 'I'.
75. Span - Returns 'I' if one NP spans the other, else 'C'.
76. Subclass - If one NP's WordNet class is a subclass of the other NP return 'Y' else 'N'.
77. Subject1 - Returns 'Y' if NP1 is a subject, otherwise 'N'.
78. Subject2 - Returns 'Y' if NP2 is a subject, otherwise 'N'.
79. Syntax - If the two NP's have incompatible values for BINDING, CONTRAINDICES, SPAN, or MAXIMALNP, then returns 'I', else 'C'.
80. Title - If one or both NPs is a title, returns 'I', else 'C'.
81. unknwnClass - The class of the instance, '+' if the NPs are coreferent, else '-'.
82. WeAntes - The antecedent is a form of 'We'.
83. WNSynonyms - Returns 'C' if the NPs are WordNet synonyms, else 'I'.
84. WordNetClass - Returns 'C' if both NPs have the same WordNet class, else 'I'.
85. WordNetDist - The distance in the WordNet Synset tree between the two NPs.
86. WordNetSense - Returns the first WordNet sense that both NPs share.
87. WordOverlap - If the intersection of the content words of the two nps is not empty, then 'C', else 'I'.
88. WordsStr - If both NPs are non-pronominal and the strings match, then 'C', else 'I'.
89. WordsSubstr - If both NPs are non-pronominal and one np is proper substrings of the other with respect to content words, then returns 'C', else 'I'.

## References

- A. Bagga and B. Baldwin. 1998. Algorithms for scoring coreference chains. In *In Linguistic Coreference Workshop at LREC 1998*.
- Baldrige, J. 2005. *The OpenNLP project*. <http://opennlp.sourceforge.net/>. The OpenNLP Project.
- Eric Bengtson and Dan Roth. 2008. Understanding the value of features for coreference resolution. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, pages 294–303. Association for Computational Linguistics.
- Chih-Chung Chang and Chih-Jen Lin. 2001. Libsvm: a library for support vector machines.
- J. Finkel, S. Dingare, H. Nguyen, M. Nissim, and C. Manning. 2004. Exploiting Context for Biomedical Entity Recognition: From Syntax to the Web. In *Joint Workshop on Natural Language Processing in Biomedicine and its Applications at COLING 2004*.
- T. Finley and T. Joachims. 2005. Supervised clustering with support vector machines. In *International Conference on Machine Learning (ICML)*, pages 217–224.
- A. Haghighi and D. Klein. 2007. Unsupervised Coreference Resolution in a Nonparametric Bayesian Model. In *Proceedings of the 45rd Annual Meeting of the Association for Computational Linguistics*.
- T. Joachims. 2002. SVM light, <http://svmlight.joachims.org>.
- D. Klein and C. Manning. 2003. Fast Exact Inference with a Factored Model for Natural Language Parsing. In *Advances in Neural Information Processing Systems 15*.
- S. Lappin and H. Leass. 1994. An algorithm for pronominal anaphora resolution. *Computational Linguistics*, 20(4):535–561.
- Xiaoqiang Luo, Abe Ittycheriah, Hongyan Jing, Nanda Kambhatla, and Salim Roukos. 2004. A mention-synchronous coreference resolution algorithm based on the bell tree. In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics*.
- X. Luo. 2005. On Coreference Resolution Performance Metrics. In *Proceedings of the 2005 Human Language Technology Conference / Conference on Empirical Methods in Natural Language Processing*.
- A. McCallum and B. Wellner. 2004. Conditional Models of Identity Uncertainty with Application to Noun Coreference. In *18th Annual Conference on Neural Information Processing Systems*.
- Thomas Morton. 1999. Using coreference for question answering. In *The 8th Text REtrieval Conference (TREC-8)*, pages 85–89.
- MUC-6. 1995. Coreference Task Definition. In *Proceedings of the Sixth Message Understanding Conference (MUC-6)*, pages 335–344.
- MUC-7. 1997. Coreference Task Definition. In *Proceedings of the Seventh Message Understanding Conference (MUC-7)*.
- V. Ng and C. Cardie. 2002. Improving Machine Learning Approaches to Coreference Resolution. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*.
- V. Ng. 2008. Unsupervised Models for Coreference Resolution. In *Proceedings of 2008 the Conference on Empirical Methods in Natural Language Processing (EMNLP-2008)*.
- NIST. 2004. *The ACE Evaluation Plan*. NIST.
- S. Petrov and D. Klein. 2007. Improved Inference for Unlexicalized Parsing. In *Proceedings of the Annual Meeting of the North American Chapter of the Association for Computational Linguistics (HLT/NAACL 2007)*.
- M. Poesio and M. Kabadjov. 2004. A general-purpose, off-the-shelf anaphora resolution module: implementation and preliminary evaluation. In *LREC*.
- L. Qiu, M.-Y. Kan, and T.-S. Chua. 2004. A public reference implementation of the rap anaphora resolution algorithm. In *LREC*.
- E. Riloff and W. Phillips. 2004. An Introduction to the Sundance and AutoSlog Systems. Technical Report UUCS-04-015, School of Computing, University of Utah.
- W. Soon, H. Ng, and D. Lim. 2001. A Machine Learning Approach to Coreference of Noun Phrases. *Computational Linguistics*, 27(4):521–541.
- J. Steinberger, M. Poesio, M. Kabadjov, and K. Jezek. 2007. Two uses of anaphora resolution in summarization. *Information Processing and Management, Special Issue on Summarization*, 43:16631680.
- V. Stoyanov, N. Gilbert, C. Cardie, and E. Riloff. 2009. Conundrums in noun phrase coreference resolution: Making sense of the state-of-the-art. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the Association for Computational Linguistics and the 4th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing*.
- Y. Versley, S.P. Ponzetto, M. Poesio, V. Eidelman, J. Jern, A. and Smith, X. Yang, and A. Moschitti. 2008. BART: A modular toolkit for coreference resolution. In *LREC*.

- M. Vilain, J. Burger, J. Aberdeen, D. Connolly, and L. Hirschman. 1995. A Model-Theoretic Coreference Scoring Theme. In *Proceedings of the Sixth Message Understanding Conference (MUC-6)*.
- I. Witten and E. Frank. 2005. *Data Mining: Practical Machine Learning Tools and Techniques, 2nd Edition*. Morgan Kaufmann.
- Xiaofeng Yang, Guodong Zhou, Jian Su, and Chew Lim Tan. 2003. Coreference resolution using competition learning approach. In *ACL '03: Proceedings of the 41st Annual Meeting on Association for Computational Linguistics*, pages 176–183.