

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/220675419>

# Comparing methods for the syntactic simplification of sentences in information extraction

Article *in* Literary and Linguistic Computing · October 2011

DOI: 10.1093/lc/fqr034 · Source: DBLP

---

CITATIONS

25

---

READS

154

1 author:



[Richard Evans](#)

University of Wolverhampton

48 PUBLICATIONS 716 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:



Coreference and Pronoun Resolution [View project](#)



Building Research Databases for Researchers (BiRD) [View project](#)

1 **Title:** Comparing Methods for the Syntactic Simplification of Sentences in Information

2 Extraction

3 **Name of Author:** Richard Evans

4 **Address where the work was done:**

5 Research Institute in Information and Language Processing,

6 University of Wolverhampton,

7 Wulfruna Street,

8 Wolverhampton,

9 West Midlands,

10 WV1 1NA.

11 United Kingdom.

12 **Email:** r.j.evans@wlv.ac.uk

13 **Abstract**

14 This article describes research aimed at improving the accuracy of an information  
15 extraction system by treating coordinate structures systematically. Commas, coordinating  
16 conjunctions, and adjacent comma-conjunction pairs are considered to be potential  
17 indicators of coordination in natural language. A recursive algorithm is implemented  
18 which converts sentences containing classified potential coordinators into sequences of  
19 simple sentences. Several approaches to the classification of potential coordinators are  
20 presented, one exploiting memory based learning, another exploiting the publicly  
21 available Stanford parser, and a hybrid approach which classifies commas and  
22 conjunctions using the former system and comma-conjunction pairs using the latter. The  
23 article describes the initial set of features developed for exploitation by the memory based  
24 classifier and presents optimization of that classifier. A baseline system is also described.

25         The sentence simplification module was exploited by an information extraction  
26 system. With regard to the automatic classifiers that form the basis for simplification,  
27 comparative evaluation demonstrated that information extraction can be performed with  
28 greatest accuracy when exploiting the hybrid classifier. It also demonstrated that a simple  
29 baseline classifier induces improved accuracy when compared to systems that ignore the  
30 presence of coordinate structures in input sentences. The article presents an analysis of  
31 the errors made by the different sentence simplification modules and the information  
32 extraction system that exploits them. Directions for future research are suggested.

33 **Comparing Methods for the Syntactic Simplification of Sentences in Information**  
34 **Extraction**

35

36 **1. Introduction**

37 This article presents a method to improve the accuracy of a clinical information  
38 extraction system (IE) by pre-processing syntactically complex input sentences. The  
39 research described investigates the automatic simplification of syntactic complexity in  
40 natural language. It focuses on the relations of *subordination* and *coordination* which  
41 involve the linking of syntactic units of the same rank in a sentence. In subordination, the  
42 linked units form a hierarchy with the subordinate unit being a constituent of the  
43 superordinate unit (1).

44

45 (1) [[For the past 3 days][,] he has had fever, malaise, and headache].

46

47 In coordination, the linked units are constituents at the same level of constituent structure  
48 (2). It is 'a type of linkage whereby the resulting conjoint construction is equivalent,  
49 structurally speaking, to each of its members.' That is, 'if [A] and [B] are conjoins of the  
50 conjoint construction X, then any structural function which may be undertaken  
51 individually by [A] or [B] may also be undertaken by X.' (Quirk *et al.* 1985).

52

53 (2) A 3-day-old boy is brought to the emergency department because of a 6-hour history  
54 of [[rapid breathing] [and] [poor feeding]].

55

56 In general, linked units may comprise a wide range of grammatical categories and levels  
57 of syntactic projection.

58 In the present article, we adopt the terminology used by Quirk *et al.* (1985).  
59 Linked constituents are referred to as *conjoins*. Their linking forms a *coordinated*  
60 *constituent*. Overt linking of conjoins by coordinating conjunctions is referred to as  
61 *syndetic* coordination. Coordination in which the linking is not overtly marked, except by  
62 the occurrence of commas or semicolons in writing or tone unit boundaries in speech, is  
63 termed *asyndetic*.

64 Coordination usually links conjoins which are structurally and grammatically  
65 similar. As noted by Quirk *et al.* (1985), this can include some complex cases in which  
66 combined units such as indirect and direct objects (3), objects and direct complements  
67 (4), and objects and adverbials (5), are coordinated.

68

69 (3) We gave [[[William] [a book on stamps]] [and] [[Mary] [a book on painting]]].

70

71 (4) Jack painted [[[the kitchen] [white]] [and] [[the living room] [blue]]].

72

73 (5) You should serve [[[the coffee] [in a mug]] [and] [[the lemonade] [in a glass]]].

74

75 In the present article, these are considered examples of verb phrase (VP) coordination  
76 with ellipsis of the second head verb.

77 With regard to noun phrases (NPs), Quirk *et al.* describe *segregatory coordination*  
78 (2) and *combinatory coordination* (6). The two can be distinguished by considering the

79 relationship of the coordinated constituent and its conjoins to its predicate. If the  
80 coordinated constituent can be replaced by each of its conjoins in the sentence and the  
81 meanings of the new sentences are consistent with that of the original, then the  
82 coordination is segregatory. If this replacement creates new sentences whose meanings  
83 are not consistent with that of the original (7), then the coordination is combinatory.

84

85 (6) The patient usually complains of [[pins] [and] [needles]] in the deltoid area.

86

87 (7) \*The patient usually complains of [pins] in the deltoid area. The patient usually  
88 complains of [needles] in the deltoid area.

89

90 Due to the scarcity of combinatory coordination in the corpus described in Section 3.1 of  
91 this article, which provides evidence of the occurrence and use of coordination in this  
92 context, all NP coordination is considered segregatory in the research described here.

93 In writing, coordination is indicated by the use of conjunctions and punctuation.  
94 The approach taken in the current article focuses on *potential coordinators* which  
95 comprise the coordinating conjunctions *and*, *but*, and *or*, semicolons, commas, and  
96 adjacent comma-conjunction pairs. By definition, the coordinating conjunctions usually  
97 serve as coordinating links between conjoins. There is far more ambiguity in the use of  
98 commas, which may have either coordinating or subordinating functions.

99 Nunberg *et al.* (2002) describe the use of punctuation in English. They note that  
100 commas, semicolons, and colons normally mark constituent boundaries within sentences.  
101 In addition to coordinated units, commas serve to mark the boundaries of subordinated

102 constituents such as post-modifiers (8), adverbial modifiers (1), and other sub-clausal  
103 constituents which are less central to the main message being conveyed in the sentence  
104 (9). Nunberg *et al.* (2002) note that adjuncts, parentheticals, supplementary relative  
105 clauses, vocatives, and a range of others are all commonly bounded in this way by  
106 delimiting commas.

107

108 (8) [His father has schizophrenia[,] [paranoid type][,] treated with haloperidol and  
109 trihexyphenidyl].

110

111 (9) [Examination[,] [including cardiovascular examination[,] ] shows no abnormalities].

112

113         In the context of our current work, the term *simple sentence* is used to denote  
114 declarative sentences containing no coordinate constituents. The aim of this research is to  
115 improve the performance of an information extraction (IE) system by means of a module  
116 that rewrites input sentences containing coordinate constituents as sequences of simple  
117 sentences. The module is also intended to recognize some types of subordinate  
118 constituent and exploit them in the IE process. One hypothesis tested in this article is that  
119 it is more effective for a system to exploit a small number of rules to extract pertinent  
120 facts from simple sentences than to exploit a larger number of rules in an effort to address  
121 the variation that results from coordination in natural language.

122         The detection of potential coordinators, their classification, and the identification  
123 of their conjoins is a prerequisite to realizing this aim. In the present article, we assume  
124 that coordination in a sentence can be detected by reference to potential coordinators.

125 Given that the function of potential coordinators is often ambiguous, especially in the  
126 case of commas, it is necessary to recognize their use as subordinators and coordinators.  
127 Further, as described in Section 3.1, coordination can hold between a variety of syntactic  
128 categories at various levels of syntactic projection. For this reason, once a system has  
129 identified a coordinator for the purpose of rewriting a complex sentence as a sequence of  
130 simple sentences, it is necessary for it to further identify the particular type of  
131 coordination signaled by the coordinator from the wide range of possibilities that exist.

132 Explicitly, the aim of the sentence rewriting module described in Section 3 is to  
133 convert sentences such as (10) into sequences of sentences such as (11).

134

135 (10) Examination shows [[jaundice][,] [hypothermia][,] [hypotonia][,] [[large [anterior]  
136 [and] [posterior] fontanel]][, and] [a hoarse cry]].

137

138 (11) Examination shows [a hoarse cry]. Examination shows [hypotonia]. Examination  
139 shows [[large] anterior fontanel]. Examination shows [[large] posterior fontanel].  
140 Examination shows [jaundice]. Examination shows [hypothermia].

141

142 In this article, Section 2 motivates research into the rewriting of complex  
143 sentences as sequences of simple sentences for the purpose of an application in natural  
144 language processing (NLP), information extraction (IE). The initial system is described  
145 and several performance issues noted. Section 3 begins with a description of the corpus of  
146 clinical vignettes that serves as the basis for the sentence simplification method presented  
147 in this article. An analysis of this corpus is described and findings regarding the use and



148 the range of types of coordination and subordination that occur in it is presented. This  
149 section also presents a new machine learning classifier for potential coordinators in  
150 natural language sentences. It automatically labels instances as belonging to one of a  
151 wide variety of subordinating or coordinating classes derived from the analysis of the  
152 corpus. A range of baseline classifiers are also described. Finally, Section 3 describes an  
153 algorithm for rewriting complex sentences into sequences of simple sentences that  
154 exploits the classifiers. Section 4 presents related work on coordination, punctuation, its  
155 automatic treatment and exploitation in NLP. Evaluation of the new approaches is  
156 presented in Section 5, which includes a comparison of IE systems exploiting the  
157 classifiers and sentence rewriting module described in Section 3. Section 6 presents plans  
158 for future work while Section 7 discusses the findings of the article and draws some  
159 conclusions. Throughout the article, unless stated otherwise, all linguistic examples are  
160 drawn from the corpus of clinical vignettes presented in Section 3.1. Relevant conjoins,  
161 coordinators, and subordinators are delimited using square brackets.

162

## 163 **2 Motivation: Information Extraction from Clinical Vignettes**

164 The research described in this article was undertaken in the context of a project on  
165 information extraction from vignettes that provide brief clinical descriptions of patients.

166 The discourse structure of these vignettes consists of seven elements:

- 167 1. Basic information
- 168 2. Chief complaint
- 169 3. History
- 170 4. Vital signs

171 5. Physical examination

172 6. Diagnostic study

173 7. Laboratory study

174 Considering each in turn, *Basic information* describes the patient's gender, profession,  
175 ethnicity, and health status. *Chief complaint* presents the main concern that led the patient  
176 to seek therapeutic intervention. *History* is a narrative description of the patient's social,  
177 family, and medical history. *Vital signs* is a description of the patient's pulse and  
178 respiration rates, blood pressure, and temperature. *Physical examination* is a narrative  
179 description of clinical findings observed in the patient. *Diagnostic study* and *Laboratory*  
180 *study* present the results of several different types of clinical test carried out on the  
181 patient.

182 Each element in the discourse structure is represented by a template encoding  
183 related information. For example, the template for physical examinations holds  
184 information on each clinical finding or symptom (*finding*) observed in the examination,  
185 information on the technique used to elicit that finding (*technique*), the bodily location to  
186 which the technique was applied (*location*), the body system that the finding provides  
187 information on (*system*), and any qualifying information about the finding (*qualifier*). In  
188 this article, we focus on automatic extraction of information pertaining to physical  
189 examinations. The goal of the IE system is to identify the phrases used in the clinical  
190 vignette that denote findings and related concepts and add them to its database entry for  
191 the vignette.

192 In the research described in this article, the IE system depends on several NLP  
193 modules:

- 194 1 Sentence tagger;
- 195 2 Concept tagger;
- 196 3 Relation extractor.

197

198 Modules 1 and 2 are arranged in a pipeline, each one adding XML annotation to its input  
199 and passing this on to be exploited by the next module. Both were developed in house.  
200 The concept tagger uses gazetteers to tag references to clinical concepts mentioned in the  
201 vignette. In light of the specificity of the IE task undertaken in this research, the  
202 gazetteers were developed in-house on the basis of corpus analysis. Existing resources  
203 such as SNOMED and UMLS were considered, but their size and scope made them  
204 difficult to exploit in the current research. Hand-crafted finite-state transducers were used  
205 in conjunction with the gazetteers to group sequences of adjacent concepts together.

206 With regard to the third module in the IE pipeline, two relation extraction  
207 modules, BASIC and PATTERNS, were implemented for the purpose of comparison.  
208 Both of them exploit the annotation of sentences and clinical concepts obtained from the  
209 first two modules.

210 BASIC consists of a small number of simple rules. To summarize briefly,  
211 vignettes are processed by considering each sentence in turn. The first clinical finding or  
212 symptom mentioned in a sentence is taken as the basis for a new database entry.  
213 Similarly, the first tagged *technique*, *system*, and *location* within that sentence is  
214 considered to be related to the *finding*. *Qualifiers* (e.g. *bilateral* or *peripheral*) are  
215 extracted in the same way, except in sentences containing the word *no*. In these cases, the  
216 qualifier related to the finding is identified as *none*. Due to their scarcity in the corpus,

217 this rule was not extended to additional negative markers such as *never* or *not*. When  
218 processing sentences generated by the simplification module described in Section 3.3, if  
219 the input sentence contains no tagged **techniques**, then BASIC attempts to extract this  
220 information from any adverbials identified in the sentence.

221 The PATTERNS relation extraction module takes every mention of a finding or  
222 symptom tagged in the input vignette as the basis for a new physical examination entry  
223 in the database. A set of hand-crafted rules is then applied to identify references to related  
224 concepts mentioned in the vignette. By way of illustration, references to the location to  
225 which a technique is applied in order to elicit a clinical finding are identified by selecting,  
226 on the basis of the first applicable rule, any tagged **location** in the pattern:

- 227 1. **technique** {at|over} the \_
- 228 2. **technique** of the \_
- 229 3. \_ **system** is **finding**
- 230 4. **finding** ... {of|at|over} the \_
- 231 5. \_ **technique** {is|are} **finding**
- 232 6. **finding** in the **qualifier** \_
- 233 7. \_ is **finding**
- 234 8. \_ **finding**

235 An underbar is used to indicate the position of the **location** in the pattern. Similar rule  
236 sets are used in the identification of the other concepts related to the finding. For brevity,  
237 they are not presented in this article.

238 The hand-crafted IE rules exploited by the PATTERNS module are implemented

239 using regular expressions. They are applied in order, exploiting lexical and conceptually  
240 tagged elements. Quantitative evaluation of the BASIC and PATTERNS IE systems is  
241 presented in Section 5. In this section, we make some general observations on the outputs  
242 of the latter system.

243         It was noted that many errors were caused by its inability to accurately process  
244 coordinated constituents. Consider (12) and (13).

245

246 (12) Physical examination shows [[enlarged supraclavicular nodes that are stony hard][,  
247 and] [a liver that is [[enlarged] [and] [irregular]]]].

248

249 (13) Examination of [the [[heart][,] [lungs][, and] [abdomen]]] shows normal findings.

250

251 In (12), noun phrase and adjectival coordination means that there is a mismatch in the  
252 number of explicitly mentioned concepts: three findings, one technique, and two  
253 qualifiers. In (13), coordination of the head nouns causes a similar mismatch in the  
254 numbers of explicitly mentioned concepts: one finding, one technique, and three systems.  
255 The rules implemented in the initial IE system cannot detect the ellipsis of elements that  
256 occurs due to this coordination and are unable to reliably identify the relations holding  
257 between explicitly mentioned and elided concepts.

258         The patterns exploited by this initial IE system are too simple to accurately detect  
259 the relations that hold between the concepts tagged in these sentences. While the use of  
260 additional regular expressions would enable more accurate processing of them, they  
261 would be of limited use beyond the specific cases that they were designed to address. The

262 variability of input sentences due to syntactic coordination is so great that it should be  
263 handled systematically rather than heuristically. Attempting to meet this challenge by the  
264 formulation of additional IE patterns would lead only to small improvements and would  
265 be a continual process. This line of reasoning motivates the development of the  
266 systematic approach to coordination presented in Section 3.

267

### 268 **3. An Automatic Treatment of Coordination**

269 This section presents a method to automatically rewrite sentences containing potential  
270 coordinators into sequences of simple sentences. It relies on a corpus in which potential  
271 coordinators have been annotated with information about their specific coordinating or  
272 subordinating function. The annotation is exploited by methods to classify previously  
273 unseen potential coordinators. Finally, a sentence simplification algorithm utilizing the  
274 classifiers is presented.

275

#### 276 *3.1 An Annotated Corpus*

277 A corpus consisting of 138,641 words from 708 clinical vignettes was compiled in order  
278 to support development and evaluation of the IE system described in Section 2. The  
279 vignettes are written in academic US English and are highly consistent in their use of  
280 terminology, punctuation, and grammatical style.

281 Potential coordinators, including conjunctions, commas, and adjacent comma-  
282 conjunction pairs, were manually annotated in this corpus. In this article, seven types of  
283 potential coordinator are considered: *and*, *but*, *or*, *comma*, *comma-and*, *comma-but*, and  
284 *comma-or*. The decision to treat comma-conjunction pairs separately from commas or

285 conjunctions alone was made on the basis that they usually introduce the final conjoin of  
286 coordinated constituents. It is thus likely that they share contexts distinct from those of  
287 the other potential coordinators.

288 The annotated corpus was divided into a training portion and a testing portion.  
289 The characteristics of the two are presented in Table 1.

290

291 **Table 1 Characteristics of the annotated corpus**

292

293 In order to address our aim of implementing a module to automatically rewrite  
294 complex sentences for the purpose of subsequent NLP tasks, it is important to identify the  
295 different roles that may be played by potential coordinators. Instances of potential  
296 coordinators occurring in the corpus of vignettes were manually annotated with labels  
297 indicating their function. Where an instance occurs between two conjoins, its label  
298 conveys information about those conjoins.

299 The different classes of instance are divided into two sets, one for coordinators  
300 and another for subordinators. Table 2 and Table 3 display the different classes of  
301 coordinator and subordinator annotated in the training corpus. The abbreviations used in  
302 the tables consist of a minimum of three components. The first indicates whether the class  
303 has a coordinating (C) or subordinating (S) function. The second component indicates the  
304 projection level of the constituents: morphemic (P), lexical (L), intermediate (I), maximal  
305 (M), or clausal/extended (C). The third element of each acronym is an abbreviation of the  
306 grammatical category of the constituents: nominal (N), verbal (V), adjectival (A),  
307 adverbial (Adv), prepositional (P), quantificational (Q), or unclear (X). A final numerical

308 value is used to differentiate classes that cannot be distinguished on the basis of the  
309 criteria previously listed. To illustrate, CMV2-6 denote coordination of VPs in which the  
310 head of the rightmost VP has been elided and the conjoined VPs have distinct argument  
311 structures, as in sentences (3) to (5). The adoption of such specific classes is expected  
312 firstly to enable automatic classifiers to leverage very specific patterns of PoS tags,  
313 words, and semantic concept labels in their recognition and secondly, to enable each class  
314 to be associated with specific and accurate sentence simplification patterns.

315 **Table 2 Classes of coordinator in the training and testing corpora**

316

317 **Table 3 Classes of subordinator in the training and testing corpora**

318

319 The annotated corpus described here serves as the basis for development and  
320 evaluation of the classification modules for potential coordinators described in Section  
321 3.2.

322

323 *3.2 Automatic Classification of Potential Coordinators*

324 Several methods were implemented for the automatic classification of potential  
325 coordinators. These classifiers are described in Sections 3.2.1 – 3.2.4. Their evaluation is  
326 presented in Section 5.

327

328 *3.2.1 Memory Based Learning (MBL) Classifier*

329 Instances of potential coordinators in the training corpus were processed in order to  
330 represent them as vectors of feature values representing their linguistic properties and



331 context of use. The initial representations exploited sixty-five features. A classifier of  
332 potential coordinators was derived from this training data using the TiMBL memory-  
333 based learner (Daelemans *et al.*, 2010). Feature selection and algorithm optimization  
334 were performed using a simple hill-climbing procedure.

335         The initial set of features encodes different kinds of information about each  
336 potential coordinator. They can be grouped as follows:

- 337         1. Orthographic form of the potential coordinator.
- 338         2. Information on the position of the instance within the document.
- 339         3. Information about items that both precede and follow the potential coordinator  
340         within the same sentence. This includes:
  - 341                 a. words and their parts of speech.
  - 342                 b. clinical concepts.
  - 343                 c. the number of determiners.
  - 344                 d. the distance in words to the next following determiner if a determiner also  
345                 precedes the instance.
  - 346                 e. the parts of speech that immediately precede and follow *other* potential  
347                 coordinators that both precede and follow the instance.
- 348         4. Boolean features asserting various conditions that hold over items occurring in the  
349         same sentence as the potential coordinator:
  - 350                 a. Words with matching parts-of-speech  $p$  precede and follow the instance.  
351                 Here,  $p$  comprises verbs of the past, past participle, and singular present  
352                 tenses, determiners, cardinal numbers, adjectives, pronouns, and nouns.
  - 353                 b. An adverb precedes the instance.

- 354 c. The instance is both preceded and followed by a word with part-of-speech  
355 *q*:
- 356 i. where *q* includes adjectives and past participle verbs.  
357 ii. where *q* includes potentially mismatched singular or plural  
358 common nouns or proper nouns.
- 359 d. The instance is *immediately* preceded and followed by a word with part-  
360 of-speech *q*, where *q* is:
- 361 i. determiner.  
362 ii. cardinal number.
- 363 e. The words *no*, *not*, or *either* precede the instance in the sentence.  
364 f. An adverb or preposition precedes the instance.  
365 g. Textual material that includes a word with part-of-speech *p* followed by a  
366 word with part of speech *q* both precedes and follows the instance where:
- 367 i. *p* is an adjective and *q* is a preposition.  
368 ii. *p* is nominal and *q* is an adverb.  
369 iii. *p* is a cardinal number and *q* is a preposition.  
370 iv. *p* is nominal and *q* is a preposition.
- 371 5. A domain-specific ternary feature indicating whether the potential coordinator is  
372 either preceded or followed by the word *history* in the sentence, or both preceded  
373 and followed by that word.
- 374 6. Features that combine the values of another pair of features into a single feature:  
375 a. Immediately preceding and following part-of-speech tags (built from  
376 features in 3.a).

377           b. Closest preceding and following conceptual tags (built from features in  
378           3.b).

379

380 Table 4 displays the groups of feature selected for the classification of each type of  
381 potential coordinator. The third column shows the proportion of features from the initially  
382 proposed set that are selected for optimal classification accuracy. The most globally  
383 important feature groups appear in bold font in Table 4.

384

385 **Table 4 Features selected for optimal classification of different potential**  
386 **coordinators**

387

388           The optimization revealed that for all potential coordinators, TiMBL worked best  
389 when using the TRIBL2 algorithm. Table 5 presents the optimal settings for other  
390 parameters with respect to each potential coordinator. Instances occurring in input data  
391 are classified using TiMBL with these optimal parameter settings. Section 5 presents an  
392 evaluation of the optimized classifiers described here.

393

394 **Table 5 Optimal parameter settings for TiMBL when classifying potential**  
395 **coordinators**

396

397 *3.2.2 Stanford Parser Classifier*

398 This classifier (STANFORD) exploits the Stanford Lexicalized Parser v1.6.3 (Klein and  
399 Manning, 2003) for the purpose of classifying potential coordinators. The constituent

400 structure returned by the parser can be used to derive a classification for potential  
401 coordinators for most of the classes presented in Table 2. A set of simple conversion rules  
402 exploiting regular expressions was employed for this purpose. The classification of  
403 subordinators is slightly more difficult, and is based on the recognition of patterns in the  
404 upper nodes of the tree output by the parser.

405

### 406 *3.2.3 Hybrid Classifier*

407 Evaluation of the two classifiers over the testing data showed that the MBL and  
408 STANFORD classifiers have somewhat orthogonal performance. This motivated the  
409 development of a hybrid classifier (HYBRID) that uses the MBL classifier when  
410 processing conjunctions and commas and uses the STANFORD classifier when  
411 processing adjacent comma-conjunction pairs.

412

### 413 *3.2.4 Majority Class Baseline Classifier*

414 This baseline classifier (MAJORITY) is based on observation of the frequency with  
415 which different classes of coordination and subordination occur in the training corpus. It  
416 classifies every instance with the most frequently observed class for potential  
417 coordinators of that type. Every instance of *and* and *comma-or* is classified as CMN1,  
418 every instance of *but* and *or* is classified as CMV1, every *comma* is classified as  
419 SMAAdv1, and every instance of *comma-and* and *comma-but* is classified as CCV under  
420 this method.

421

## 422 *3.3 Sentence Rewriting*

423 The syntactic simplification method exploits all the annotations added to input sentences  
424 by the previously described modules. A part of speech tagger (Brill, 1994) is also  
425 exploited. The method is based on a recursive algorithm operating over an array of  
426 sentences (see Fig. 1). In its initial state, this array comprises a single sentence containing  
427 one or more instances belonging to any of the ten classes displayed in Table 6.

428

### 429 **Fig. 1 The sentence simplification algorithm**

430

### 431 **Table 6 Classes of coordinator/subordinator triggering simplification rules**

432

433 The sentence simplification algorithm is presented in Fig. 1. The function *simplify*  
434 consists of an ordered set of quick-fire rules, designed to process different classes of  
435 coordination and subordination indicated by different types of instance. Each rule  
436 identifies a coordinator/subordinator,  $t_i$ , and generates a pair of sentences,  $\int_i^1$  and  $\int_i^2$ . The  
437 former is derived from textual material preceding  $t_i$  in the input sentence, while the latter  
438 is derived from material following it. The function returns any identified adverbials,  $adv$ ,  
439 and a reference,  $\xi_i$ , to the pair of generated sentences. The rules were developed manually  
440 by reference to the test corpus and a key file containing information on the class of each  
441 potential coordinator.

442 To illustrate with two examples of rules:

- 443 • SMAAdv1 triggers a rule that recognizes preceding material as an adverbial  
444 modifier of the input sentence,  $adv$ .  $\int_i^1$  is an empty string and  $\int_i^2$  is the part of the  
445 sentence that follows  $t_i$ . A binary array consisting of  $\int_i^1$  and  $\int_i^2$  is built.

446 • When instances of class CMN1 occur in a context such as A B/vbz C  $t_i$  D in  $s_i$ , a  
447 rule is triggered which constructs an array consisting of the strings  $\int_i^1$ : A B/vbz C  
448 and  $\int_i^2$ : A B/vbz D. If  $s_i$  is (14), this rule derives (15) as  $\int_i^1$  and (16) as  $\int_i^2$ . The  
449 upper case letters A-D are regular expressions matching text intervening between  
450 the strings specified within the pattern. vbz is a part of speech tag denoting  
451 present tense verbs. The fact that the NP conjoints follow the verb implies that they  
452 form a coordinated object. This assumption motivates the form of the rule.

453

454 (14) She has diabetic retinopathy [but] no evidence of renal disease.

455

456 (15) She has diabetic retinopathy.

457

458 (16) She has no evidence of renal disease.

459

460 The output,  $A$  is a set of sentences containing no instances of the types listed in Table 6.

461 The IE function is applied to this set of sentences and any adverbial information derived

462 during the rewriting process.

463 Table 7 displays characteristics of the rewrite rules used by the simplification

464 algorithm. Column 3 of the table displays the possible number of rewrite rules that may

465 be applied by the algorithm on encountering different types of coordinator/subordinator.

466 This information serves as an indirect indicator of the challenge posed by each rewriting

467 task. It can be noted that many more rules are needed to cater for the various functions

468 and contexts of NPs in the clinical vignettes than other types of constituent. The rules

469 include heuristics that exploit PoS tagging and preposition and verb recognition to  
470 identify the syntactic function of coordinated NPs. Column 4 of Table 7 shows the  
471 relative order in which the rules are tested against an input sentence. The success of the  
472 rewriting algorithm depends on both the patterns exploited by the rules and their order of  
473 application. In general, the rules are intended to process the coordination of larger and  
474 more syntactically dominant conjoins first.

475

#### 476 **Table 7 Characteristics of rewrite rules by class**

477

478 The simplified sentences derived by the interaction of this module with each of  
479 the classifiers of potential coordinators described in Section 3.2 are then processed by the  
480 BASIC IE system described in Section 2. The templates produced by this IE system and  
481 the PATTERNS IE system, described in the same section, are evaluated in Section 5.

482

#### 483 **4. Related Work**

484 In conducting the research presented in this article, a review of previous related work was  
485 undertaken. This includes research on the disambiguation of coordinated structures and  
486 the role of punctuation and research describing the exploitation of information about  
487 coordination in syntactic parsing, information extraction, and other NLP applications.

488 Addressing the challenge of disambiguating and processing coordination, Agarwal  
489 and Boggess (1992) present a system to identify the boundaries of conjoins linked by  
490 coordinating conjunctions. Their rule-based algorithm exploits concept tagging, part of  
491 speech tagging, and the use of a “semi-parser” to identify constituents such as NPs, VPs,

492 and PPs. It performs with an accuracy of 81.6%, but is noted to be unable to identify  
493 clausal conjoins and does not recognize coordination indicated by commas.

494 Many of the approaches presented in the literature recognize that there is likely to  
495 be syntactic and semantic similarity between conjoins involved in coordination and  
496 exploit this in order to disambiguate coordinate structures (Kurohashi and Nagao, 1992;  
497 Resnik, 1999; Goldberg, 1999; Chantree *et al.*, 2005).

498 Buyko and Hahn (2008) sought to learn the extent of the contribution made by the  
499 recognition of semantic similarity between conjoins to the processing of coordination.  
500 They found that a system based on conditional random fields exploiting semantic features  
501 was outperformed by one based on output from a syntactic parser. In the present article  
502 (see Section 3.2.1), features encoding semantic information were selected for exploitation  
503 by several classifiers, though the statistical significance of their contribution has not been  
504 assessed. Shimbo and Hara (2007) describe an approach to the disambiguation of  
505 coordinate conjunctions based on methods from sentence-alignment. Their system was  
506 found to outperform state-of-the-art parsers when processing the GENIA Treebank beta  
507 corpus.

508 Kawahara and Kurohashi (2007) present methods to disambiguate coordination in  
509 Japanese. Exploiting verb case frames automatically derived from the web, the method  
510 applies lexical preferences and co-occurrence statistics between potential conjoins to  
511 resolve coordination ambiguities. An updated approach exploiting functional dependency  
512 information was described in Kawahara and Kurohashi (2008).

513 Methods exploiting information about neighbouring syntactic constituents have  
514 been used to disambiguate the role of commas in natural language. This work was



515 described in Bayraktar *et al.* (1998) and Srikumar *et al.* (2008).

516         With regard to IE, Rindflesch *et al.* (2000) used an automatic treatment of  
517 coordination to improve IE of facts about macromolecular binding. In a contrasting  
518 approach, Klebanov *et al.* (2004) present a method to improve performance in IE without  
519 processing coordination. Their approach relies on the identification of 'easy-access  
520 sentences' (EAS) that contain a single finite verb in a 'semantically non-problematic  
521 environment' and a large number of named entities (concepts). In this approach, IE rules  
522 are applied only to EASs. The accuracy with which EASs are identified is reported in this  
523 work, but unfortunately changes in accuracy elicited in their IE system as a result of  
524 applying the method are not presented.

525         Many authors demonstrate that the use of methods to improve the resolution of  
526 coordination ambiguities improves overall performance in syntactic parsing for various  
527 languages (Kim and Lee, 2003; Ratnaparkhi *et al.*, 1994; Rus *et al.*, 2002; Nakov and  
528 Hearst, 2005; Hogan, 2007; Charniak and Johnson, 2005; Kübler *et al.*, 2009). In addition  
529 to this, various papers report on the exploitation of information about coordination for  
530 other tasks in NLP and in industrial contexts. Rindflesch (1995) incorporated a method  
531 for dealing with coordination to improve the mapping of NPs identified in input  
532 documents to concepts in the medical UMLS database. Cederberg and Widdows (2003)  
533 present a method exploiting information about noun coordination to improve automatic  
534 hyponymy extraction. The method is based on well-established lexicosyntactic patterns  
535 modified to allow recognition and exploitation of coordinated structures. The hyponymy  
536 relations identified are then filtered using latent semantic analysis. In their preliminary  
537 work, Tjong and Berry (2008) seek to improve the clarity of industrial requirements

538 specifications by minimizing the ambiguous use of coordination. Their paper describes a  
539 range of semantic relations implied by the use of coordination and urges the adoption of  
540 rules similar to a controlled language specifying a writing policy for coordinated  
541 structures.

542         Despite the amount of work addressing the issue of coordination in natural  
543 language, the contribution brought by these approaches to practical NLP applications has  
544 been little reported. Overall, the work surveyed in this section was useful in guiding  
545 development of the features presented in Section 3.2.1. Authors have drawn differing  
546 conclusions as to the suitability of different types of information in resolving  
547 coordination ambiguities. This observation motivated the approach adopted in the present  
548 article, in which an initial feature set is developed and a feature selection method is  
549 applied in order to derive the optimal subset to be exploited by the classifier.

550

## 551 **5. Evaluation**

552 This section presents an evaluation of the modules described in Section 3. In all cases,  
553 where comparisons are made between different systems in terms of accuracy or F-score,  
554 significance was computed using approximate randomization (Chinchor, 1992). The  
555 significance threshold,  $\alpha = 0.05$ .

556         The production of annotated data for the task of sentence simplification is costly  
557 and complex. For this reason, different settings of the sentence simplification module will  
558 be evaluated extrinsically (Sparck-Jones and Galliers, 1996) via the performance of the  
559 IE system that exploits them.

560         Unfortunately, due to the nature of the sources from which information is to be

561 extracted in this work, no direct comparison can be made with the systems presented in  
562 previous research. Recognizing this problem, the new modules presented in this article  
563 are compared with one based on the publicly accessible Stanford parser.

564

#### 565 *5.1 Evaluation of the Classification of Potential Coordinators*

566 Table 8 presents the accuracy scores of the different classifiers obtained using ten-fold  
567 cross-validation over the training corpus. Given their superiority over the MAJORITY  
568 classifier, the main focus of this section will be in comparing the accuracy of the MBL  
569 and STANFORD classifiers. It can be observed that the MBL classifier classifies all  
570 potential coordinators except *comma-but* with greater accuracy than the STANFORD  
571 classifier. However, the only potential coordinators for which there was a statistically  
572 significant difference in performance were the two most common, *comma* and *and*.

573

574 **Table 8 Classification accuracy obtained via ten-fold cross-validation over the**  
575 **training set**

576

577 A detailed class-by-class examination of the performance of the MBL classifier reveals  
578 that for all potential coordinators, the most common type of error concerns the projection  
579 level of nominal constituents. This finding was also derived from an analysis of inter-  
580 annotator agreement over a sample of the training data. It can be noted that errors made  
581 by the STANFORD classifier are similar in kind, though there is more evidence of the  
582 erroneous assignment of grammatical category as well as projection level to coordinate  
583 constituents. Overall, of the eighty-one combinations of classes and types, the F-score

584 obtained by STANFORD is superior to MBL in thirteen. The most frequent of these  
585 thirteen is CCV signalled by *comma-and*, which accounts for 14.25% of all instances  
586 annotated in the training data. However the margin of difference is slight (F 0.9880 vs.  
587 0.9719), and the contribution of this improvement to the IE system is not envisaged to be  
588 great. A similar description can be made with regard to the greater F-score obtained by  
589 STANFORD with regard to the CCV class signalled by *comma-but*, which accounts for  
590 1.53% of the training data. There are classes for which STANFORD obtains a  
591 significantly higher F-score than MBL, but each of these accounts for less than 1% of the  
592 total training set.

593         One reason for the relatively poor performance of STANFORD is that the labels  
594 returned by the Stanford parser are not as specific as those used in the manual annotation  
595 of the training data exploited by MBL. To illustrate, the label *VP* used by the Stanford  
596 parser subsumes two classes (*CMV1* and *CMV2*) and *NP* subsumes three (*CMN1*, *CIN*,  
597 and *CMV3*). The STANFORD classifier is therefore unable to differentiate between these  
598 classes.

599         Finally, it has been noted that both MBL and STANFORD classifiers fail to  
600 identify instances of class *CIN* (17). The most common type of error involving *CIN* is  
601 misclassification as *CMN1*. The simplification rules applied to these classes are similar in  
602 many ways, relying on identification of nominal and verbal heads in the sentence. It is  
603 therefore expected that such errors will not be too detrimental.

604

605 (17) The sclerae and the skin of the [[head] [and] [upper trunk]] are yellow.

606

607 **Table 9 Classification accuracy over the test set**

608

609 Table 9 presents the classification accuracy of different classifiers when processing a  
610 subsample of the test corpus which consists only of sentences that mention clinical  
611 findings. Over the classes of coordination and subordination occurring in descriptions of  
612 physical examinations, STANFORD classifies potential coordinators consisting of  
613 adjacent *comma-conjunction* pairs more accurately than MBL does. The difference in  
614 classification accuracy between the STANFORD and MBL classifiers is statistically  
615 significant with regard to the potential coordinators *comma-but*, *and*, and *comma*. In  
616 classifying the latter two types, MBL is superior whereas in classifying the first,  
617 STANFORD is superior.

618

619 *5.2 Evaluation of IE Exploiting Classification of Potential Coordinators and Sentence*  
620 *Simplification*

621 Testing data for the IE task was derived from the stems of 70 clinical vignettes. The set  
622 contains 206 clinical findings and related concepts. The IE systems, PATTERNS and  
623 BASIC, described in Section 2 were used to process this data set. Several variants of the  
624 BASIC system were employed, each exploiting one of the different methods for  
625 classification of potential coordinators described in Section 3.2 and the sentence  
626 simplification algorithm presented in Section 3.3. A variant of the PATTERNS relation  
627 extraction module was also implemented that extracts just a single tagged finding from an  
628 input sentence as opposed to all tagged findings.

629 The metrics used in evaluation of the IE systems are based on accuracy. For

630 findings, when the IE system identifies a finding within a particular sentence of a  
631 particular vignette, and the same finding has been marked within the same sentence of the  
632 same vignette in the key, this is considered a true positive. The accuracy score for  
633 findings is the ratio of the number of true positives to the total number of findings marked  
634 in the key. It is computed in a similar way for the concepts related to findings. Due to the  
635 strong semantic typing involved in the IE task and the limited number of candidates for  
636 selection with regard to a particular finding, accuracy was considered a more suitable  
637 metric than F-measure. The evaluation described here is based on exact string matching.  
638 Systems are not awarded rewarded for obtaining partial matches.

639

640 **Table 10 Accuracy of IE systems exploiting different classifiers of potential**  
641 **coordinators (assuming one finding per sentence)**

642

643 **Table 11 Accuracy of IE systems exploiting different classifiers of potential**  
644 **coordinators (assuming multiple findings per sentence)**

645

646 Tables 10 and 11 display the accuracy of different IE systems in identifying clinical  
647 findings and related concepts in descriptions of physical examinations. Table 10 shows  
648 the performance of IE systems implemented only to identify the first tagged finding and  
649 concepts related to that finding in input sentences. Table 11 provides evaluation results  
650 for IE systems that identify all tagged findings and concepts related to those findings in  
651 input sentences.

652 In both tables, the columns IGNORE contain accuracy scores for systems that

653 exploit the sentence rewriting module described in Section 3.3 but do not exploit any  
654 classification of potential coordinators. As a result, for these systems, the sentence  
655 rewriting rules are never activated. The columns MBL, STANFORD, HYBRID, and  
656 MAJORITY present accuracy scores for IE systems that work in the same way, but which  
657 exploit the classification modules for potential coordinators described in Sections 3.2.1,  
658 3.2.2, 3.2.3, and 3.2.4, respectively. The columns PATTERNS present the accuracy scores  
659 obtained by the IE system described in Section 2.

660 The results are broadly in line with expectation. A comparison of the overall  
661 accuracy of the PATTERNS systems with the others supports the hypothesis that it is  
662 more effective to apply a small set of IE rules over simplified input sentences than to  
663 employ a larger set of complex IE rules in an effort to handle the variation exhibited by  
664 sentences containing coordinated constituents. For IE systems identifying single findings  
665 in input sentences, the fact that IGNORE is more accurate than PATTERNS was  
666 unexpected. However, it was found that the PATTERNS approach works significantly  
667 better if multiple findings are extracted from input sentences.

668 In Table 11, the IGNORE system is the most effective one at identifying findings  
669 mentioned in clinical vignettes. This suggests that even after syntactic simplification,  
670 some test sentences still mention multiple findings. Another possibility is that some  
671 coordinated constituents have been erroneously identified as findings in the key file.

672 A significance matrix was computed to plot a pairwise comparison of all systems  
673 presented in this article. With  $\alpha = 0.05$ , the systems can be ranked as follows:

- 674 1. KEY (multiple findings identified per sentence)
- 675 2. KEY (one finding identified per sentence) and HYBRID (multiple findings per

- 676 sentence)
- 677 3. HYBRID and STANFORD (one finding per sentence) and MBL and STANFORD  
678 (multiple findings per sentence)
- 679 4. IGNORE and PATTERNS (multiple findings per sentence) and MBL (one finding  
680 per sentence)
- 681 5. MAJORITY (in both contexts)
- 682 6. IGNORE (one finding per sentence)
- 683 7. PATTERNS (one finding per sentence)

684 This ranking is based on a comparison of the number of systems that a given system  
685 significantly outperforms with the number that significantly outperform it.

686 Although not statistically significant in this setting, the difference in accuracy  
687 between the MAJORITY and IGNORE systems in Table 10 shows that performance in IE  
688 can be improved even when the classification of potential coordinators is quite  
689 inaccurate.

690

### 691 *5.3 Error Analysis*

692 The output of different modules within the IE system was examined in order to  
693 investigate the causes and impact of the errors they make. In this section, categories of  
694 errors are categorized as concerning conceptual tagging, the classification of potential  
695 coordinators, the simplification of sentences, and information extraction.

696 A number of errors arose as a result of the conceptual tagging process. In  
697 particular, there are clinical findings involving clinical procedures that were not included  
698 in our existing gazetteers and could not be recognized (e.g. *requires*



699 *intubation/mechanical ventilation*). Another omission of this type involves general  
700 vocabulary such as the word *moves* in the finding *he moves all extremities to painful*  
701 *stimuli*. Finally, several findings are numerical and their recognition depends on  
702 processing context, as in the example *Deep tendon reflexes are 1+*.

703         One particular weakness of the conceptual tagger is its inability to resolve  
704 ambiguities between adjectives that belong to different concept types according to the  
705 context of use. To illustrate, the modifier *stony-hard* functions as a qualifier or finding  
706 whereas *palpable* functions as a qualifier or technique, depending on the context of use.  
707 One additional challenge in the processing of qualifiers is the decision of whether to tag  
708 them as separate elements or to merge them with adjacent concepts.

709         With regard to the classification of potential coordinators, learning curves were  
710 plotted to show the correlation between training set size and classification accuracy for  
711 each type of potential coordinator. Examination of the learning curves suggests that a  
712 minimum of 200 instances are required in order to obtain a representative sample of the  
713 use of each potential coordinator. The training corpus used in this study contains far  
714 fewer instances than this of the potential coordinators *but, or, comma-but, and comma-or*.  
715 It is suggested that the training sets for these items should be increased considerably  
716 before the accuracy scores of the different classifiers can be regarded as definitive.

717         In the IE task, several errors were caused by a misclassification of sentences  
718 conveying information about the medical history of the patient and those concerning the  
719 physical examination. One example of this type is (18). Such errors arise because the IE  
720 system exploits information on the occurrence of particular verbs in the present tense  
721 when classifying sentences in the vignette as ones which provide information on physical

722 examinations. The verb used in the first clause of sentence (18) is also commonly used in  
723 descriptions of physical examinations.

724

725 (18) [[Needle biopsy shows papillary carcinoma][, and] [he undergoes total  
726 thyroidectomy]].

727

728         There are several instances of errors in the IE key file in which phrases denoting  
729 findings and qualifiers contain potential coordinators. These cases may have some  
730 impact on the accuracy scores obtained by the IGNORE system evaluated in this article.  
731 However, they are infrequent enough that their influence on the evaluation results  
732 reported in Section 5.2 is not expected to be significant. No instances of combinatory  
733 coordination were noted in the test data.

734

## 735 **6. Plans for Future Work**

736 The linguistic studies discussed in Section 1 and the error analysis presented in Section  
737 5.3 motivate five directions in which development of the sentence simplification module  
738 presented in this article may proceed.

739         One non-trivial improvement that could be made to the sentence simplification  
740 module would be to classify coordination as having either a segregatory or a combinatory  
741 interpretation. No assessment has been made of the significance of this issue in the  
742 context of the current IE task, but one possible approach to this challenge would be a  
743 method exploiting very large unannotated corpora. Quirk *et al.* (1985) note that one way  
744 for linguists to distinguish between segregatory and combinatory coordination is to check

745 the acceptability of sentences created by inserting the word *both* before the first conjoin.  
746 This operation would produce sentence (19) from sentence (6).

747

748 (19) The patient usually complains of both [[pins] [and] [needles]] in the deltoid area.

749

750 It may be possible, by examining the frequencies of such constructed sentences in very  
751 large corpora, to recognize combinatory coordination in input sentences. Empirical  
752 approaches comparing the frequency of occurrence of constructions in which the order of  
753 the conjoins is reversed may also be examined.

754 Nunberg *et al.* (2002) present a description of the role and use of other  
755 punctuation symbols besides the comma such as indicators of parenthesis, single and  
756 double dashes, single and double quotation marks, related punctuation indicators, and the  
757 pragmatic implications that arise from the interaction of various punctuation marks. The  
758 modules described in the current article do not address these phenomena. For the current  
759 IE task, this is not problematic, but it is envisaged that IE from sources such as medical  
760 journals, text books, and patient notes may benefit from future work on the simplification  
761 of sentences employing this wider range of punctuation symbols.

762 The MBL classifier of potential coordinators was optimized using a naïve hill-  
763 climbing procedure in which feature selection and algorithm optimization are treated  
764 independently. Methods for joint optimization of the two have been undertaken in  
765 previous work (Daelemans *et al.*, 2003). Such approaches are more computationally  
766 expensive, often exploiting clusters of processors employing genetic algorithms. It has  
767 been shown that joint optimization leads to the derivation of significantly more accurate

768 classifiers by undertaking a more thorough exploration of the possibility space defined by  
769 different parameter settings. It will be interesting to apply such approaches in future work  
770 in order to derive more effective classifiers of potential coordinators.

771         For the scenario described in Section 2, the recognition and use of specific verbs  
772 in the rules used by the IE system is not important. However, this is not true of IE in  
773 alternate scenarios in which pertinent facts are identified by reference to the verbs linking  
774 different concepts. In light of this, it will be beneficial to apply a methodology to ensure  
775 subject-verb concord in the sentences generated by the module described in Section 3.3.  
776 This will ensure that a sentence such as (20) will be rewritten as a sequence such as (21)  
777 rather than (22). This improvement can be made using relatively simple morpho-syntactic  
778 rules.

779

780 (20) [[Pelvic examination] [and] [urinalysis]] show no abnormalities.

781

782 (21) [Pelvic examination] shows no abnormalities. [Urinalysis] shows no abnormalities.

783

784 (22) \*[Pelvic examination] show no abnormalities. [Urinalysis] show no abnormalities.

785

786         In addition to the expansion of the annotated corpus motivated by observations  
787 made in Section 5.3 and with improvement in subject-verb concord, it will be interesting  
788 to assess the contribution of the simplification process in other NLP applications such as  
789 question answering, pronoun resolution, multiple-choice question generation, and IE in  
790 different scenarios.

791 Finally, analysis of documents from alternate domains shows evidence of classes  
792 of subordination absent from the corpus described in Section 3.1. To be effective when  
793 applied to different domains, the annotation scheme for subordinators should be revised  
794 to include classes of comma signalling the left and right boundaries of different types of  
795 subordinated constituent.

796

## 797 **7. Conclusion**

798 Three main conclusions were drawn from the research described in this article. The first  
799 is that the automatic simplification of syntactic complexity can induce significant  
800 improvements in subsequent NLP tasks. A variety of approaches were implemented and  
801 evaluated by reference to the accuracy of an IE system exploiting them. Of the fully  
802 automatic modules tested, the best performing one was a hybrid system combining a  
803 memory-based learning classifier with a classifier derived from a syntactic parser. When  
804 exploiting classifiers based only on a syntactic parser or a memory based learning  
805 method, sentence simplification still significantly improved the accuracy of the IE  
806 system.

807 The second conclusion to be drawn also follows from the comparative evaluation  
808 of variant IE systems. It was found that approaches which bypass a systematic treatment  
809 of coordination and handle coordination and subordination by means of more  
810 sophisticated IE rules perform relatively poorly.

811 The third conclusion to be drawn from this article follows from error analysis. It is  
812 expected that the syntactic simplification method described here will be improved by  
813 pursuing various lines of research. These include increasing the amount of annotated data

814 available for development of some of the classifiers of potential coordinators, introducing  
815 procedures to disambiguate combinatory and segregatory coordination, developing a  
816 module to recognize the functions of a wider range of punctuation symbols in the  
817 simplification model, and introducing methods to ensure subject-verb concord in the  
818 sentences generated by the modules.

819

820 **10. References**

- 821 **Agarwal, R. and Boggess, L.** (1992). A Simple but Useful Approach to Conjunct  
822 Identification, *Proceedings of the 30th annual meeting for Computational Linguistics*,  
823 Newark, Delaware, June 1992.
- 824 **Bayraktar, M., Say, B., and Akman, V.** (1998). An Analysis of English Punctuation: The  
825 Special Case of Comma. In the *International Journal of Corpus Linguistics*, 3 (1), pp.  
826 33–57.
- 827 **Brill, E.** (1994). Some Advances in Transformation-Based Part of Speech Tagging,  
828 *Proceedings of the Twelfth National Conference on Artificial Intelligence*, Seattle,  
829 Washington, July-August 1994.
- 830 **Buyko, E. and Hahn, U.** (2008). Are Morpho-syntactic Features More Predictive for the  
831 Resolution of Noun Phrase Coordination Ambiguity than Lexico-semantic Similarity  
832 Scores?, *Proceedings of the 22nd International Conference on Computational Linguistics*  
833 (*Coling 2008*), Manchester, England, August 2008.
- 834 **Cederberg, S. and Widdows, D.** (2003). Using LSA and Noun Coordination Information  
835 to Improve the Precision and Recall of Automatic Hyponymy Extraction, *Proceedings of*  
836 *the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*, Edmonton,  
837 Canada, May 2003.
- 838 **Chantree, F., Kilgarriff, A., De Roeck, A., and Willis, A.** (2005). Using a Distributional  
839 Thesaurus to Resolve Coordination Ambiguities. Technical Report, The Open University.
- 840 **Charniak, E. and Johnson, M.** (2005). Coarse-to-Fine n-Best Parsing and MaxEnt  
841 Discriminative Reranking, *Proceedings of the 43rd Annual Meeting of the ACL*, Ann  
842 Arbor, Michigan, June 2005.

843 **Chinchor, N.** (1992). The statistical significance of the MUC-4 results, *Proceedings of*  
844 *the Fourth Message Understanding Conference*, McLean, Virginia, June 1992.

845 **Daelemans, W., Hoste, V., De Meulder, F., and Naudts, B.** (2003). Combined  
846 optimization of feature selection and algorithm parameters in machine learning of  
847 language, *Proceedings of the 14th European Conference on Machine Learning (ECML-*  
848 *2003)*, Cavtat-Dubrovnik, Croatia, September 2003.

849 **Daelemans, W., Zavrel, J., Van der Sloot, K., and Van den Bosch, A.** (2010). TiMBL:  
850 Tilburg Memory Based Learner, version 6.3, Reference Guide. Technical Report, ILK  
851 Research Group.

852 **Goldberg, M.** (1999). An Unsupervised Model for Statistically Determining Coordinate  
853 Phrase Attachment, *Proceedings of the 37th annual meeting of the Association for*  
854 *Computational Linguistics on Computational Linguistics*, College Park, Maryland, June  
855 1999.

856 **Hogan, D.** (2007). Coordinate noun phrase disambiguation in a generative parsing model,  
857 *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*,  
858 Prague, Czech Republic, June 2007.

859 **Kawahara, D. and Kurohashi, S.** (2007). Probabilistic Coordination Disambiguation in  
860 a Fully-Lexicalised Japanese Parser, *Proceedings of the 2007 Joint Conference on*  
861 *Empirical Methods in Natural Language Processing and Computational Natural*  
862 *Language Learning*, Prague, Czech Republic, June 2007.

863 **Kawahara, D. and Kurohashi, S.** (2008). Coordination Disambiguation Without any  
864 Similarities, *Proceedings of the 22nd International Conference on Computational*  
865 *Linguistics (Coling 2008)*, Manchester, England, August 2008.



866 **Kim, M.-Y. and Lee, J.-H.** (2003). S-Clause Segmentation for Efficient Syntactic  
867 Analysis Using Decision Trees, *Proceedings of the Australasian Language Technology*  
868 *Workshop*, Melbourne, Australia, December 2003.

869 **Klebanov, B. B., Knight, K., and Marcu, D.** (2004). Text Simplification for  
870 Information-Seeking Applications. In Meersman, R. and Tari, Z. (eds) *On the Move to*  
871 *Meaningful Internet Systems 2004*, Berlin: Springer-Verlag, pp. 735-747.

872 **Klein, D. and Manning, C.D.** (2003). Fast exact inference with a factored model for  
873 natural language parsing. In *Advances in Neural Information Processing Systems (NIPS-*  
874 *15)*, Vancouver, British Columbia, December 2002.

875 **Kübler, S., Hinrichs, E., Maier, W., and Klett, E.** (2009). Parsing Coordinations,  
876 *Proceedings of the 12th Conference of the European Chapter of the ACL*, Athens, Greece,  
877 March 2009.

878 **Kurohashi, S. and Nagao, M.** (1992). Dynamic Programming Method for Analysing  
879 Conjunctive Structures in Japanese, *Proceedings of COLING-92*, Nantes, France, August  
880 1992.

881 **Nakov, P. and Hearst, M.** (2005). Using the Web as an Implicit Training Set:  
882 Application to Structural Ambiguity Resolution, *Proceedings of Human Language*  
883 *Technology Conference and Conference on Empirical Methods in Natural Language*  
884 *Processing (HLT/EMNLP)*, Vancouver, British Columbia, October 2005.

885 **Nunberg, G., Briscoe, T., and Huddleston, R.** (2002). Punctuation. In Huddleston, R.  
886 and Pullum, G. K. (eds) *The Cambridge Grammar of the English Language*. Cambridge:  
887 Cambridge University Press, pp. 1724–1764.

888 **Quirk, R., Greenbaum, S., Leech, G., and Svartvik, J.** (1985). *A Comprehensive*

889 *Grammar of the English Language*. London: Longman.

890 **Ratnaparkhi, A., Roukos, S., and Ward, R. T.** (1994). A Maximum Entropy Model for  
891 Parsing, *Proceedings of the International Conference on Spoken Language Processing*  
892 (*ICSLP*), Yokohama, Japan, September 1994.

893 **Resnik, P.** (1999). Semantic Similarity in a Taxonomy: An Information-Based Measure  
894 and Its Application to Problems of Ambiguity in Natural Language. In *The Journal of*  
895 *Artificial Intelligence Research* (11), pp 95-130.

896 **Rindflesch, T. C.** (1995). Integrating Natural Language Processing and Biomedical  
897 Domain Knowledge for Increased Information Retrieval Effectiveness, *Proceedings of*  
898 *the Fifth Annual Dual-use Technologies and Applications Conference*, Utica/Rome, NY,  
899 May 1995.

900 **Rindflesch, T. C., Rajan, J. V., and Hunter, L.** (2000). Extracting Molecular Binding  
901 Relationships from Biomedical Text, *Proceedings of the Sixth Conference on Applied*  
902 *Natural Language Processing*, Seattle, Washington, April 2000.

903 **Rus, V., Moldovan, D., and Bolohan, O.** (2002). Bracketing Compound Nouns for  
904 Logic Form Derivation, *Proceedings of FLAIRS-2002*, Pensacola Beach, Florida, May  
905 2002.

906 **Shimbo, M. and Hara, K.** (2007). A discriminative learning model for coordinate  
907 Conjunctions, *Proceedings of the 2007 Conference on Empirical Methods in Natural*  
908 *Language Processing and Computational Natural Language Learning*, Prague, Czech  
909 Republic, June 2007.

- 910 **Sparck-Jones, K. and Galliers, J. R.** (1996). *Evaluating Natural Language Processing*  
911 *Systems: An Analysis and Review*. No. 1083 in Lecture Notes on Artificial Intelligence.  
912 Berlin: Springer.
- 913 **Srikumar, V., Reichart, R., Sammons, M., Rappoport, A., and Roth, D.** (2008).  
914 Extraction of Entailed Semantic Relations Through Syntax-based Comma Resolution,  
915 *Proceedings of the Annual Meeting of the Association of Computational Linguistics*  
916 *(ACL)*, Columbus, Ohio, June 2008.
- 917 **Tjong, S. F. and Berry, D. M.** (2008). Can Rules of Inferences Resolve Coordination  
918 Ambiguity in Natural Language Requirements Specification?, *Proceedings of the 11th*  
919 *Workshop on Requirements Engineering (WER'08)*, Barcelona, Spain, September 2008.

920 **Table 1 Characteristics of the annotated corpus**

Characteristic	Training	Testing
#Items	422	286
#Words	107,900	30,741
#Sentences	12,451	3286
#Potential coordinators	4709	1491

922 **Table 2 Classes of coordinator in the training and testing corpora**

Category	Morphemic	Lexical	Intermediate	Phrasal	Clausal
Noun		CLN	CIN	CMN1, CMN2, CMN3	
Verb		CLV		CMV1, CMV2, CMV3, CMV4, CMV5, CMV6	CCV
Adjective	CPA	CLA		CMA1, CMA2	
Adverb				CMAAdv	
Preposition		CLP		CMP	
Quantifier		CLQ			
Miscellaneous				CMM1, CMM2	

924 **Table 3 Classes of subordinator in the training and testing corpora**

Category	Morphemic	Lexical	Intermediate	Phrasal	Clausal
Noun				SMN	
Verb					
Adjective				SMA	
Adverb				SMAAdv1, SMAAdv2	
Preposition					
Quantifier					
Miscellaneous				SMM1, SMM2	SCM

926

927

928 **Table 4 Features selected for optimal classification of different potential**  
 929 **coordinators**

Potential coordinator	Feature group	Proportion of features selected
<i>and</i>	<b>3.a</b> , 3.b, 3.c, <b>3.e</b> , 4.a, 4.c.i, 4.d.ii, 4.f, 4.g.ii, 4.g.iii, 4.g.iv, 5, 6.a, 6.b	0.4923
<i>but</i>	2, <b>3.a</b> , 3.b, 3.c, <b>3.e</b> , 4.a, 4.b, 4.c.i, 4.d, 4.f, 6.a, 6.b	0.4154
<i>or</i>	<b>3.a</b> , <b>3.e</b> , 4.c.ii, 6.a	0.2308
<i>comma</i>	<b>3.a</b> , 3.b, 3.c, <b>3.e</b> , 4.a, 4.c.i, 4.f, 5, 6.a, 6.b	0.4154
<i>comma-and</i>	<b>3.a</b> , 3.b, 3.c, 3.d, <b>3.e</b> , 4.a, 4.c.i, 4.c.ii, 4.f, 4.g.iii, 5, 6.a	0.5231
<i>comma-but</i>	<b>3.a</b> , <b>3.e</b> , 6.a	0.0461
<i>comma-or</i>	3.d, <b>3.e</b> , 6.a	0.0461

931 **Table 5 Optimal parameter settings for TiMBL when classifying potential**  
 932 **coordinators**

Parameter setting	Classifiers
Feature weighting	Gain ratio Shared variance No weighting
Class voting weight	Inverse distance Normal majority voting
Distance metric	Modified value difference Jeffrey divergence Overlap
Neighbours	3 4 5 21

934 **Fig. 1 Sentence simplification algorithm**

**Input:** Sentence containing coordinated constituents,  $s_0$   
**Output:** Array of simple sentences,  $A$ ; Adverbial modifier,  $adv$

```

1  $A \leftarrow \emptyset$ ;
2  $adv \leftarrow$  empty string;
3  $S \leftarrow \{s_0\}$ ;
4 while  $S \neq \emptyset$  do
5    $s_i \leftarrow pop(S)$ ;
6   if  $s_i$  contains a coordinator/subordinator of a type listed in Table 6 then
7      $(adv, \xi_i) \leftarrow simplify(s_i)$ ;
8      $f_i \leftarrow dereference(\xi_i)$ ;
9      $S \leftarrow S \cup \{f_i\}$ ;
10  else
11     $A \leftarrow A \cup \{s_i\}$ 
12  end
13 end

```

936 **Table 6 Classes of coordinator/subordinator triggering simplification rules**

Coordinator/ Subordinator	Classes
<i>and</i>	CCV, CMN1, CIN, CLA, CMA1, CMV1
<i>but</i>	CMN1, CMA1, CMV1
<i>or</i>	CMN1, CIN, CLN, CMV1
<i>comma</i>	SMA <sub>adv</sub> 1, CCV, SMM1, SMM2, CMN1, CMA1, CLA
<i>comma-and</i>	CMV1, CMN1, CLN, CCV
<i>comma-but</i>	CMA1, CCV
<i>comma-or</i>	CMN1

938

939

940

941

942 Table 7 Characteristics of rewrite rules by class

Class	Coordinator/ subordinator	#Rewriting rules	Order of precedence
CCV	<i>and</i>	1	3
	<i>comma</i>	1	4
	<i>comma-and</i>	1	2
	<i>comma-but</i>	1	2
CMN1	<i>and</i>	26	8
	<i>but</i>	1	9
	<i>or</i>	9	10
	<i>comma</i>	7	12
	<i>comma-and</i>	7	11
	<i>comma-or</i>	1	11
CIN	<i>and</i>	2	17
	<i>or</i>	1	18
CLA	<i>and</i>	2	21
	<i>comma</i>	2	22
CMA1	<i>and</i>	1	14
	<i>but</i>	1	14
	<i>comma</i>	2	16
	<i>comma-but</i>	1	15
CMV1	<i>and</i>	2	6
	<i>but</i>	2	6
	<i>or</i>	2	6
	<i>comma-and</i>	1	7
CLN	<i>or</i>	1	19
	<i>comma-and</i>	1	20
SMA <sub>adv</sub> 1	<i>comma</i>	1	1
SMM1	<i>comma</i>	2	5
SMM2	<i>comma</i>	2	5

945 **Table 8 Classification accuracy obtained via ten-fold cross-validation over the**  
 946 **training set**

Potential coordinator	#Instances	MAJORITY	STANFORD	MBL
<i>and</i>	1544	0.3543	0.5971	0.7506
<i>but</i>	100	0.7100	0.8000	0.8700
<i>or</i>	80	0.2625	0.4750	0.6000
<i>comma</i>	1931	0.2952	0.7369	0.8716
<i>comma-and</i>	965	0.6953	0.8788	0.8891
<i>comma-but</i>	75	0.9600	0.9867	0.9733
<i>comma-or</i>	14	0.5714	0.5714	0.7143
ALL	4709	0.4158	0.7205	0.8320

948

949 **Table 9 Classification accuracy over the test set**

Potential coordinator	#Instances	MAJORITY	STANFORD	MBL
<i>and</i>	137	0.3650	0.4453	0.6642
<i>but</i>	13	0.3077	0.5385	0.7692
<i>or</i>	12	0.0833	0.4167	0.4167
<i>comma</i>	91	0.1209	0.5604	0.7363
<i>comma-and</i>	49	0.3673	0.8163	0.7347
<i>comma-but</i>	6	0.6667	0.8333	0.6667
<i>comma-or</i>	2	0.5000	0.5000	0.0000
ALL	310	0.2871	0.5484	0.6871

951

952

953

954



955 **Table 10 Accuracy of IE systems exploiting different classifiers of potential**  
 956 **coordinators (assuming one finding per sentence)**

Template slot	One finding per sentence						
	IGNORE	MAJORITY	PATTERNS	STANFORD	MBL	HYBRID	KEY
finding	0.5845	0.7584	0.5556	0.7971	0.7971	0.8019	0.8696
technique	0.7729	0.7681	0.7778	0.7778	0.7778	0.7874	0.8019
system	0.7536	0.8357	0.7391	0.8309	0.8406	0.8454	0.8744
qualifier	0.6812	0.8164	0.5797	0.7971	0.8261	0.8213	0.8261
location	0.8696	0.8889	0.8985	0.8985	0.9034	0.9082	0.9324
ALL	0.7324	0.8135	0.7101	0.8203	0.8290	0.8328	0.8609

958

959 **Table 11 Accuracy of IE systems exploiting different classifiers of potential**  
 960 **coordinators (assuming multiple findings per sentence)**

Template slot	Multiple findings per sentence						
	IGNORE	MAJORITY	PATTERNS	STANFORD	MBL	HYBRID	KEY
finding	0.9420	0.8068	0.8454	0.8744	0.8551	0.8696	0.9275
technique	0.7729	0.7681	0.8116	0.7778	0.7778	0.7874	0.8019
system	0.7536	0.8357	0.8406	0.8309	0.8406	0.8454	0.8744
qualifier	0.6811	0.8164	0.6135	0.7971	0.8261	0.8213	0.8261
location	0.8696	0.8889	0.9130	0.8985	0.9034	0.9082	0.9324
ALL	0.8039	0.8232	0.8048	0.8357	0.8406	0.8464	0.8725