

Empirical Methods in Information Extraction

Claire Cardie

■ This article surveys the use of empirical, machine-learning methods for a particular natural language-understanding task—information extraction. The author presents a generic architecture for information-extraction systems and then surveys the learning algorithms that have been developed to address the problems of accuracy, portability, and knowledge acquisition for each component of the architecture.

Most corpus-based methods in natural language processing (NLP) were developed to provide an arbitrary text-understanding application with one or more general-purpose linguistic capabilities, as evidenced by the articles in this issue of *AI Magazine*. Author Eugene Charniak and coauthors Ng Hwee Tou and John Zelle, for example, describe techniques for part-of-speech tagging, parsing, and word-sense disambiguation. These techniques were created with no specific domain or high-level language-processing task in mind. In contrast, my article surveys the use of empirical methods for a particular natural language-understanding task that is inherently domain specific. The task is information extraction. Generally, an *information-extraction system* takes as input an unrestricted text and “summarizes” the text with respect to a pre-specified topic or domain of interest: It finds useful information about the domain and encodes the information in a structured form, suitable for populating databases. In contrast to in-depth natural language-understanding tasks, information-extraction systems effectively skim a text to find relevant sections and then focus only on these sections in subsequent processing. The information-extraction system in figure 1, for example, summarizes stories about natural disasters, extracting for each such event the type of disaster, the date

and time that it occurred, and data on any property damage or human injury caused by the event.

Information extraction has figured prominently in the field of empirical NLP: The first large-scale, head-to-head evaluations of NLP systems on the same text-understanding tasks were the Defense Advanced Research Projects Agency-sponsored Message-Understanding Conference (MUC) performance evaluations of information-extraction systems (Chinchor, Hirschman, and Lewis 1993; Lehnert and Sundheim 1991). Prior to each evaluation, all participating sites receive a corpus of texts from a predefined domain as well as the corresponding answer keys to use for system development. The *answer keys* are manually encoded templates—much like that of figure 1—that capture all information from the corresponding source text that is relevant to the domain, as specified in a set of written guidelines. After a short development phase,¹ the NLP systems are evaluated by comparing the summaries each produces with the summaries generated by human experts for the same test set of previously unseen texts. The comparison is performed using an automated scoring program that rates each system according to measures of recall and precision. *Recall* measures the amount of the relevant information that the NLP system correctly extracts from the test collection; *precision* measures the reliability of the information extracted:

$$\text{recall} = \frac{(\# \text{ correct slot fillers in output templates})}{(\# \text{ slot fillers in answer keys})}$$

$$\text{precision} = \frac{(\# \text{ correct slot fillers in output templates})}{(\# \text{ slot fillers in output templates})}$$

As a result of MUC and other information-ex-

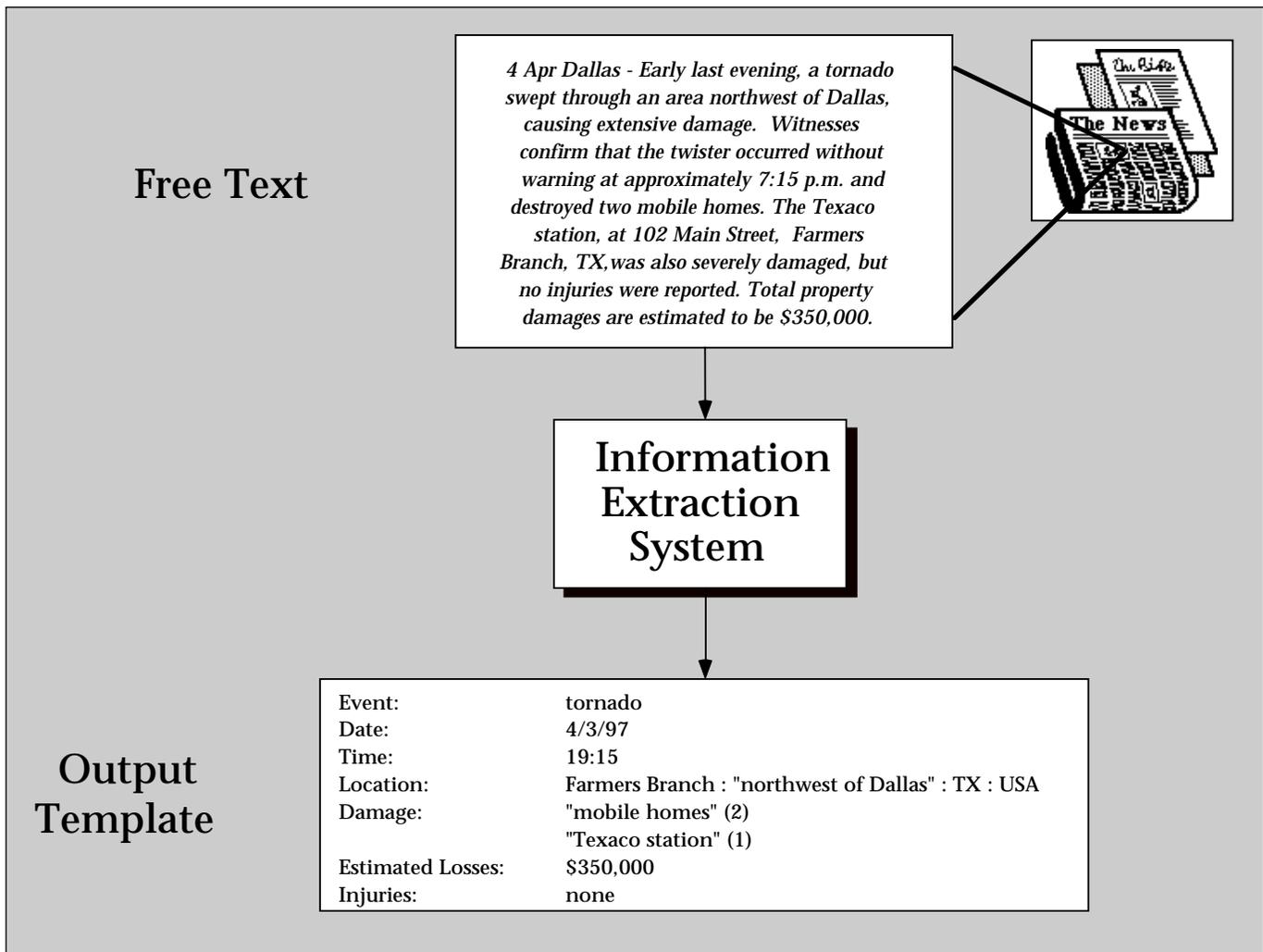


Figure 1. An Information-Extraction System in the Domain of Natural Disasters.

traction efforts, information extraction has become an increasingly viable technology for real-world text-processing applications. For example, there are currently information-extraction systems that (1) support underwriters in analyzing life insurance applications (Glasgow et al. 1997); (2) summarize medical patient records by extracting diagnoses, symptoms, physical findings, test results, and therapeutic treatments to assist health-care providers or support insurance processing (Soderland, Aronow, et al. 1995); (3) analyze news wires and transcripts of radio and television broadcasts to find and summarize descriptions of terrorist activities (MUC-4 1992; MUC-3 1991); (4) monitor technical articles describing microelectronic chip fabrication to capture information on chip sales, manufacturing advances, and the development or use of chip-processing technologies (MUC-5 1994); (5) analyze newspaper articles with the goal of finding and summariz-

ing business joint ventures (MUC-5 1994); and (6) support the automatic classification of legal documents (Holowczak and Adam 1997).

A growing number of internet applications also use information-extraction technologies. Some examples include NLP systems that build knowledge bases directly from web pages (Craven et al. 1997); create job-listing databases from news groups, web sites, and classified advertisements (see www.junglee.com/success/index.html); build news group query systems (Thompson, Mooney, and Tang 1997); and create weather forecast databases from web pages (Soderland 1997).

Although the MUC evaluations have shown that it is possible to rigorously evaluate some aspects of an information-extraction system, it is difficult to state the overall performance levels of today's information-extraction systems: At a minimum, performance depends on the relative complexity of the extraction task, the

quality of the knowledge bases available to the NLP system, the syntactic and semantic complexity of the documents to be processed, and the regularity of the language in the documents. In general, however, the best extraction systems now can achieve levels of about 50-percent recall and 70-percent precision on fairly complex information-extraction tasks and can reach much-higher levels of performance (approximately 90-percent recall and precision) for the easiest tasks. Although these levels of performance might not initially seem impressive, one should realize that information extraction is difficult for people as well as machines. Will's (1993) study, for example, showed that the best machine-extraction systems have an error rate that is only twice that of highly skilled analysts specifically trained in information-extraction tasks.

In spite of this recent progress, today's information-extraction systems still have problems: First, the accuracy and robustness of machine-extraction systems can be improved greatly. In particular, human error during information extraction is generally caused by a lapse of attention, but the errors of an automated extraction system are the result of its relatively shallow understanding of the input text. As a result, the machine-generated errors are more difficult to track down and correct. Second, building an information-extraction system in a new domain is difficult and time consuming, often requiring months of effort by domain specialists and computational linguists familiar with the underlying NLP system. Part of the problem lies in the domain-specific nature of the task: An information-extraction system will work better if its linguistic knowledge sources are tuned to the particular domain, but manually modifying and adding domain-specific linguistic knowledge to an existing NLP system is slow and error prone.

The remainder of the article surveys the empirical methods in NLP that have been developed to address these problems of accuracy, portability, and knowledge acquisition for information-extraction systems. Like the companion articles in this issue, we see that empirical methods for information extraction are corpus-based, machine-learning algorithms. To start, I present a generic architecture for information-extraction systems. Next, I provide examples of the empirical methods designed to increase the accuracy or the portability of each component in the extraction system. Throughout, I focus on the specific needs and constraints that information extraction places on the language-learning tasks.

The Architecture of an Information-Extraction System

In the early days of information extraction, NLP systems varied widely in their approach to the information-extraction task. At one end of the spectrum were systems that processed a text using traditional NLP techniques: (1) a full syntactic analysis of each sentence, (2) a semantic analysis of the resulting syntactic structures, and (3) a discourse-level analysis of the syntactic and semantic representations. At the other extreme lie systems that used keyword-matching techniques and little or no linguistic analysis of the input text. As more information-extraction systems were built and empirically evaluated, however, researchers began to converge on a standard architecture for information-extraction systems. This architecture is shown in figure 2. Although many variations exist from system to system, the figure indicates the main functions performed in an information-extraction system.

Each input text is first divided into sentences and words in a tokenization and tagging step. As indicated in figure 2, many systems also disambiguate, or tag, each word with respect to part of speech and, possibly, semantic class at this point during processing. The *sentence-analysis phase* follows. It comprises one or more stages of syntactic analysis, or parsing, that together identify noun groups, verb groups, prepositional phrases, and other simple constructs. In some systems, the parser also locates surface-level subjects and direct objects and identifies conjunctions, appositives, and other complex phrases. At some point, either before, during, or after the main steps of syntactic analysis, an information-extraction system also finds and labels semantic entities relevant to the extraction topic. In the natural disaster domain, for example, the system might identify locations, company names, person names, time expressions, and money expressions, saving each in a normalized form.

Figure 2 shows the syntactic constituents and semantic entities identified during sentence analysis for the first sentence of the sample text. There are important differences between the sentence-analysis stage of an information-extraction system and traditional parsers. Most importantly, the goal of syntactic analysis in an information-extraction system is not to produce a complete, detailed parse tree for each sentence in the text. Instead, the system need only perform partial parsing; that is, it need only construct as much structure as the information-extraction task requires. Unlike traditional full-sentence parsers, a partial pars-

...
empirical methods for information extraction are corpus-based, machine-learning algorithms.

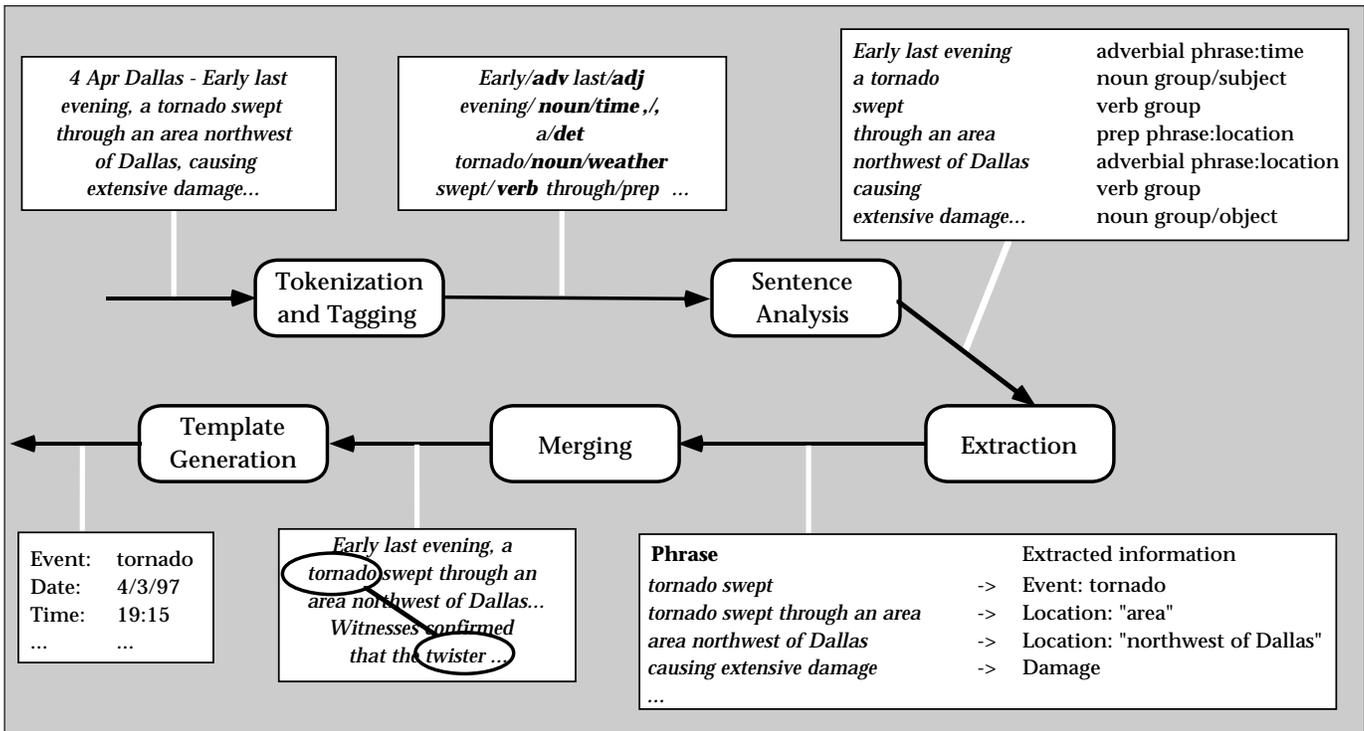


Figure 2. Architecture for an Information-Extraction System.

er looks for fragments of text that can reliably be recognized, for example, noun groups and verb groups. Because of its limited coverage, a partial parser can rely solely on general pattern-matching techniques—often finite-state machines—to identify these fragments deterministically based on purely local syntactic cues. Partial parsing is well suited for information-extraction applications for an additional reason: The ambiguity-resolution decisions that make full-blown parsing difficult can be postponed until later stages of processing where top-down expectations from the information-extraction task can guide the system's actions.

The *extraction phase* is the first entirely domain-specific component of the system. During extraction, the system identifies domain-specific relations among relevant entities in the text. Given the first sentence in our example, this component should identify the type of natural disaster (*tornado*), the location of the event (*area* and *northwest of Dallas, TX*), and the fact that there was some property damage. Figure 2 shows the information extracted from the first sentence and those portions of text responsible for each piece of extracted data. Information for filling the remaining slots of the output template would similarly be extracted from subsequent sentences.

The main job of the *merging phase* is coreference resolution, or anaphora resolution: The system examines each entity encountered in the text and determines whether it refers to an existing entity or whether it is new and must be added to the system's discourse-level representation of the text. In the sample text, for example, the mention of *a tornado* in the first sentence indicates a new entity; *the twister* in sentence two, however, refers to the same entity as the *tornado* in sentence one. Recognizing when two statements refer to the same entity is critical for an information-extraction system because it allows the system to associate the information from both statements with the same object. In some systems, another task of merging is to determine the implicit subjects of all verb phrases. In sentence one, this component would infer that *tornado* is the subject of *causing* (as well as the subject of *swept*), allowing the system to directly associate *damage* with the tornado. The discourse-level inferences made during merging aid the *template-generation phase*, which determines the number of distinct events in the text, maps the individually extracted pieces of information onto each event, and produces output templates. Purely domain-specific inferences can also occur during template generation. In the MUC terrorism domain, for example, terrorist

events involving only military targets were not considered relevant unless civilians were injured, or there was damage to civilian property. The template-generation phase is often the best place to apply this domain-specific constraint. In addition, some slots in the output template must be filled with terms chosen from a set of possibilities rather than a string from the input text. In the sample scenario, the Injuries and Event slots require such *set fills*. Still other slots (for example, Date, Time, Location) might require normalization of their fillers. Both of these subtasks are part of the template-generation phase.

The Role of Corpus-Based– Language Learning Algorithms

With this architecture in mind, we can now return to our original question: How have researchers used empirical methods in NLP to improve the accuracy and portability of information-extraction systems? In general, corpus-based–language learning algorithms have been used to improve individual components of the information-extraction system and, as a result, to improve the end-to-end performance. In theory, empirical methods can be used for each subtask of information extraction: part-of-speech tagging, semantic-class tagging, word-sense disambiguation, named entity identification (for example, company names, person names, locations), partial parsing, extraction-pattern learning, coreference resolution, and each step of template generation. The catch, as is often the case in corpus-based approaches to language learning, is obtaining enough training data. As described in the overview article, supervised language learning algorithms acquire a particular language-processing skill by taking many examples of how to correctly perform the task and then generalizing from the examples to handle unseen cases. The algorithms, therefore, critically depend on the existence of a corpus that has been annotated with the appropriate supervisory information. For language tasks that are primarily domain independent and syntactic in nature, annotated corpora such as the Penn tree bank (Marcus, Marcinkiewicz, and Santorini 1993) already exist and can be used to extract training data for the information-extraction system. Part-of-speech tagging and bracketing text into noun groups, verb groups, clauses, and so on, fall into this category. For these tasks, one can use the tree bank's *Wall Street Journal* corpus, which has been annotated with both word class and syntactic structure, together with any of a number of corpus-based

algorithms: Some examples include using hidden Markov models (HMMs) for part-of-speech tagging and statistical learning techniques for parsing (see Charniak [1993] and Weischedel et al. [1993]), Brill's (1995) transformation-based learning for part-of-speech tagging and bracketing (Ramshaw and Marcus 1995), decision tree models for parsing (Magerman 1995), case-based learning for lexical tagging (Daelemans et al. 1996; Cardie 1993), and inductive logic programming for learning syntactic parsers (Zelle and Mooney 1994). The resulting taggers and bracketers will be effective across information-extraction tasks as long as the input to the information-extraction system uses a writing style and genre that is similar to the training corpus. Otherwise, a new training corpus must be created and used to completely retrain or bootstrap the training of the component. In theory, word-sense-disambiguation algorithms would also be portable across extraction tasks. However, defining standard word senses is difficult, and to date, text collections have been annotated according to these predefined senses only for a small number of selected words. In addition, the importance of word-sense disambiguation for information-extraction tasks remains unclear.

Natural language-learning techniques are more difficult to apply to subsequent stages of information extraction—namely, the learning of extraction patterns, coreference resolution, and template generation. There are a number of problems: First, there are usually no corpora annotated with the appropriate semantic and domain-specific supervisory information. The typical corpus for information-extraction tasks is a collection of texts and their associated answer keys, that is, the output templates that should be produced for each text. Thus, a new corpus must be created for each new information-extraction task. In addition, the corpus simply does not contain the supervisory information needed to train most components of an information-extraction system, including the lexical-tagging, coreference-resolution, and template-generation modules. The output templates are often inadequate even for learning extraction patterns: They indicate which strings should be extracted and how they should be labeled but say nothing about which occurrence of the string is responsible for the extraction when multiple occurrences appear in the text. Furthermore, they provide no direct means for learning patterns to extract set fills, symbols not necessarily appearing anywhere in the text. As a result, researchers create their own training corpora, but because this process is slow, the resulting corpora can be

... corpus-based–language learning algorithms have been used to improve individual components of the information-extraction system and, as a result, to improve the end-to-end performance.

*One of the
earliest
systems
for
acquiring
extraction
patterns
was
AUTOSLOG.*

much smaller than is normally required for statistical approaches to language analysis.

Another problem is that the semantic and domain-specific language-processing skills needed for information extraction often require the output of earlier levels of analysis, for example, tagging and partial parsing. This requirement complicates the generation of training examples for the learning algorithm because there can be no standard corpus from which complete training examples can be read off, as is the case for part-of-speech tagging and parsing. The features that describe the learning problem depend on the information available to the extraction system in which the learning algorithm is embedded, and these features become available only after the training texts have passed through earlier stages of linguistic analysis. Whenever the behavior of these earlier modules changes, new training examples must be generated and the learning algorithms for later stages of the information-extraction system retrained. Furthermore, the learning algorithms must deal effectively with noise caused by errors from earlier components. The cumulative effect of these complications is that the learning algorithms used for low-level tagging or syntactic analysis might not readily apply to the acquisition of these higher-level language skills, and new algorithms often need to be developed.

In spite of the difficulties of applying empirical methods to problems in information extraction, it is precisely the data-driven nature of corpus-based approaches that allows them to simultaneously address both of the major problems for information-extraction systems—accuracy and portability. When the training data are derived from the same type of texts that the information-extraction system is to process, the acquired language skills are automatically tuned to that corpus, increasing the accuracy of the system. In addition, because each natural language-understanding skill is learned automatically rather than manually coded into the system, the skill can be moved quickly from one information-extraction system to another by retraining the appropriate component.

The remaining sections describe the natural language-learning techniques that have been developed for training the domain-dependent and semantic components of an information-extraction system: extraction, merging, and template generation. In each case, I describe how the previously discussed problems are addressed and summarize the state of the art in the field.

Learning Extraction Patterns

As in the sentence-analysis stages, general pattern-matching techniques have also become the technique of choice for the extraction phase of an information-extraction system (MUC-6 1995). The role for empirical methods in the extraction phase, therefore, is one of knowledge acquisition: to automate the acquisition of good extraction patterns, where *good patterns* are those that are general enough to extract the correct information from more than one sentence but specific enough to not apply in inappropriate contexts. A number of researchers have investigated the use of corpus-based methods for learning information-extraction patterns. The learning methods vary along a number of dimensions: the class of patterns learned, the training corpus required, the amount and type of human feedback required, the degree of preprocessing necessary, the background knowledge required, and the biases inherent in the learning algorithm itself.

One of the earliest systems for acquiring extraction patterns was AUTOSLOG (Riloff 1993; Lehnert et al. 1992). AUTOSLOG learns extraction patterns in the form of domain-specific concept-node definitions for use with the CIRCUS parser (Cardie and Lehnert 1991; Lehnert 1990). AUTOSLOG's *concept nodes* can be viewed as domain-specific semantic case frames that contain a maximum of one slot for each frame. Figure 3, for example, shows the concept node for extracting *two mobile homes* as damaged property from sentence two of the sample text. The first field in the concept node specifies the type of concept to be recognized (for example, Damage). The concept type generally corresponds to a specific slot in the output template (for example, the Damage slot of figure 1). The remaining fields in the concept node represent the extraction pattern. The *trigger* is the word that activates the pattern—it acts as the pattern's conceptual anchor point. The *position* denotes the syntactic position where the concept is expected to be found in the input sentence (for example, the direct object, subject, object of a preposition); the *constraints* are selectional restrictions that apply to any potential instance of the concept. In CIRCUS, these semantic constraints can be hard or soft: *Hard constraints* are predicates that must be satisfied before the phrase in the specified position can be extracted as an instance of the concept; *soft constraints* suggest preferences for slot fillers but do not inhibit the extraction of phrases if violated. In all our examples, we assume that the constraints are hard. Finally, the *enabling conditions* are constraints on the linguistic con-

text of the triggering word that must be satisfied before the pattern is activated. The concept node of figure 3, for example, would be triggered by the word *destroyed* when used in the active voice. Once activated, the concept node would then extract the direct object of the clause to fill its Damage slot as long as the phrase denoted a physical object.

Once the concept node is defined, it can be used in conjunction with a partial parser to extract information from novel input sentences. The system parses the sentence; if the trigger word is encountered and the enabling conditions satisfied, then the phrase found in the specified syntactic constituent is extracted, tested for the appropriate semantic constraints, and then labeled as an instance of the designated concept type. The bottom of figure 3 shows the concept extracted from sentence two of the sample text after applying the Damage concept node. Alternatively, given the sentence “the hurricane destroyed two office buildings,” the same Damage concept node would extract *two office buildings* as the damaged entities. The extracted concepts are used during merging and template generation to produce the desired output templates.

AUTOSLOG learns concept-node definitions using a one-shot learning algorithm designed specifically for the information-extraction task. As a training corpus, it requires a set of texts and their answer keys.² The AUTOSLOG learning algorithm is straightforward and depends only on the existence of a partial parser, a small lexicon with semantic-class information, and a small set (approximately 13) of general linguistic patterns that direct the creation of concept nodes. Given a noun phrase from an answer key, AUTOSLOG performs the following steps to derive a concept node for extracting the phrase from the original text:

First, find the sentence from which the noun phrase originated. For example, given the target noun phrase *two mobile homes* that fills the Damage slot, AUTOSLOG would return sentence two from the sample text during this step.

Second, present the sentence to the partial parser for processing. AUTOSLOG’s partial parser must be able to identify the subject, direct object, verb group, and prepositional phrases of each clause. For sentence two of the sample text, the parser should determine, among other things, that *destroyed* occurred as the verb group of the third clause with *two mobile homes* as its direct object.

Third, apply the linguistic patterns in order. AUTOSLOG’s linguistic patterns attempt to identify domain-specific thematic role information

Sentence Two:

“Witnesses confirm that the twister occurred without warning at approximately 7:15 p.m. and *destroyed two mobile homes.*”

Concept-Node Definition:

Concept = Damage
 Trigger = “destroyed”
 Position = direct-object
 Constraints = ((physical-object))
 Enabling Conditions = ((active-voice))

Instantiated Concept Node

Damage = “two mobile homes”

Figure 3. Concept Node for Extracting Damage Information.

for a target noun phrase based on the syntactic position in which the noun phrase appears and the local linguistic context. The first pattern that applies determines the extraction pattern, that is, concept node, for extracting the noun phrase from the training sentence. The linguistic pattern that would apply in the *two mobile homes* example is

<active-voice-verb> followed by <target-np>=
 <direct object> .

This pattern says that the noun phrase to be extracted, that is, the target-np, appeared as the direct object of an active voice verb. Similar patterns exist for the objects of passives and infinitives and for cases where the target noun phrase appears as the subject of a clause or the object of a prepositional phrase. AUTOSLOG’s linguistic patterns are, for the most part, domain independent; they need little or no modification when moving an NLP system from one information-extraction task to another.

Fourth, when a pattern applies, generate a concept-node definition from the matched constituents, their context, the slot type for the target noun phrase, and the predefined semantic class for the filler. For our example, AUTOSLOG would generate a concept node definition of the following form:

Concept = < <slot type> of <target-np> >
 Trigger =
 “< <verb> of <active-voice-verb> >”
 Position = direct-object
 Constraints =
 ((< <semantic class> of <concept> >))
 Enabling Conditions = ((active-voice)) .

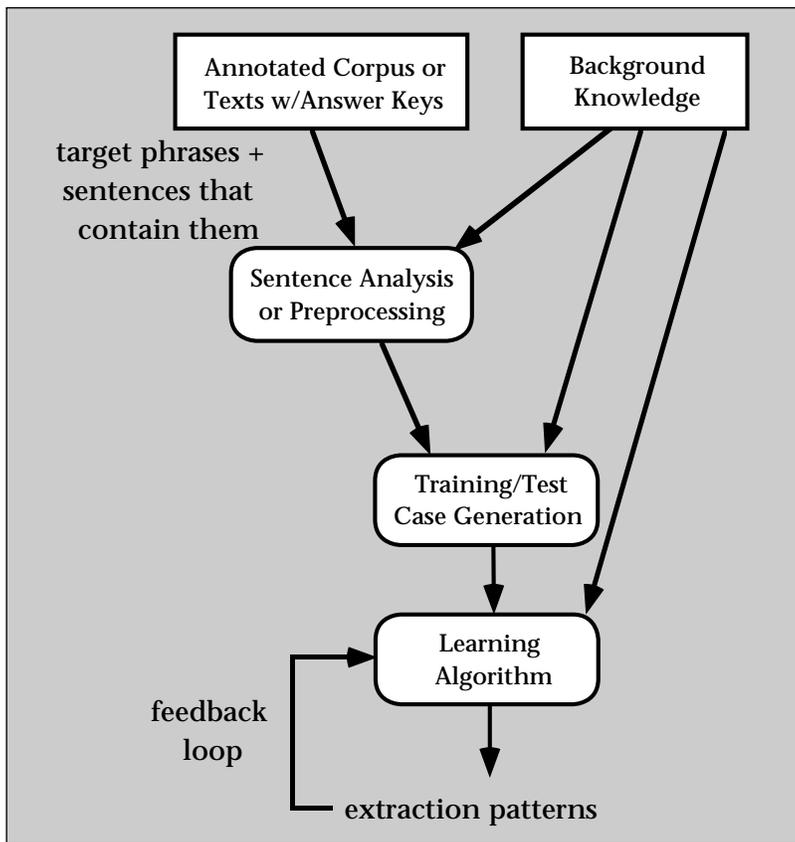


Figure 4. Learning Information-Extraction Patterns.

AUTOSLOG assumes that the semantic class of each concept type is given as part of the domain specification and that the parser has a mechanism for assigning these semantic classes to nouns and noun modifiers during sentence analysis. After substitutions, this concept-node definition will match the Damage concept node of figure 3.

Some examples of extraction patterns learned by AUTOSLOG for the terrorism domain include (in shorthand form): <victim> was **murdered**; <perpetrator> **bombed**; <perpetrator> attempted to **kill**; was **aimed** at <target>. In these examples, the bracketed items denote concept type and the word in boldface is the concept node trigger. Although many of AUTOSLOG's learned patterns are good, some are too general (for example, they are triggered by *is* or *are*); others are too specific; still others are just wrong. These bad extraction patterns are sometimes caused by parsing errors; alternatively, they occur when target noun phrases occur in a prepositional phrase, and AUTOSLOG cannot determine whether the preceding verb or noun phrase should trigger the extraction. As a result, AUTOSLOG requires that a person re-

view the proposed extraction patterns and discard those that seem troublesome.

AUTOSLOG has been used to automatically derive extraction patterns for a number of domains: terrorism, business joint ventures, and advances in microelectronics. In terms of improving the portability of information-extraction systems, AUTOSLOG allowed developers to create extraction patterns for the terrorism domain in 5 hours instead of the approximately 1200 hours required to create extraction patterns for the domain by hand. In terms of accuracy, there was no direct empirical evaluation of the learned patterns, even though some form of cross-validation could have been used. Instead, the learned patterns were evaluated indirectly—by using them to drive the University of Massachusetts–MUC-3 information-extraction system. In short, the AUTOSLOG-generated patterns achieved 98 percent of the performance of the handcrafted patterns. This result is especially impressive because the handcrafted patterns had placed the University of Massachusetts system among the top performers in the MUC-3 performance evaluation (Lehnert et al. 1991). AUTOSLOG offered an additional advantage over the handcrafted rule set: Because domain experts can review the automatically generated extraction patterns with minimal training, building the patterns no longer required the expertise of a computational linguist with a deep understanding of the underlying NLP system. This is a critical step toward building information-extraction systems that are trainable entirely by end users.

Figure 4 shows the general structure of corpus-based approaches to learning information-extraction patterns. AUTOSLOG conforms to this structure except for its human feedback loop, which does not inform the learning algorithm of its findings. Virtually all subsequent attempts to automate the acquisition of extraction patterns also conform to the general structure of figure 4. In the next paragraphs, I describe a handful of these systems.

First, Kim and Moldovan's (1995) PALKA system learns extraction patterns that are similar in form to AUTOSLOG's concept nodes. The approach used to generate the patterns, however, is quite different. The background knowledge is not a set of linguistic patterns to be instantiated but a concept hierarchy, a set of predefined keywords that can be used to trigger each pattern, and a semantic-class lexicon. The concept hierarchy contains generic semantic case-frame definitions for each type of information to be extracted. To learn extraction patterns, PALKA looks for sentences that contain case-

frame slots using semantic-class information. The training corpus is used to choose the correct mapping when more than one is possible, but the concept and semantic-class hierarchies guide PALKA's generalization and specialization of proposed patterns.

Like AUTOSLOG and PALKA, CRYSTAL (Soderland et al. 1995) learns extraction patterns in the form of semantic case frames. CRYSTAL's patterns, however, can be more complicated. Instead of specifying a single trigger word and its local linguistic context, the triggers for CRYSTAL's patterns comprise a much more detailed specification of linguistic context. In particular, the target constituent or any surrounding constituents (for example, the subject, verb, or object of the current clause) can be tested for a specific sequence of words or the presence of heads or modifiers with the appropriate semantic class. CRYSTAL uses a covering algorithm to learn extraction patterns and their relatively complicated triggering constraints. *Covering algorithms* are a class of inductive learning technique that successively generalizes input examples until the generalization produces errors. As a result, CRYSTAL begins by generating the most specific concept node possible for every phrase to be extracted in the training texts. It then progresses through the concept nodes one by one. For each concept node, C , CRYSTAL finds the most similar concept node, C' , and relaxes the constraints of each just enough to unify C and C' . The new extraction pattern, P , is tested against the training corpus. If its error rate is less than some prespecified threshold, P is added to the set, replacing C and C' . The process is repeated on P until the error tolerance is exceeded. At this point, CRYSTAL moves on to the next pattern in the original set. CRYSTAL was initially used to derive extraction patterns for a medical diagnosis domain, where it achieved precision levels ranging from 50 percent to 80 percent and recall levels ranging from 45 percent to 75 percent, depending on how the error-tolerance threshold was set.

Although AUTOSLOG, PALKA, and CRYSTAL learn extraction patterns in the form of semantic case frames, each uses a different learning strategy. AUTOSLOG creates extraction patterns by specializing a small set of general linguistic patterns, CRYSTAL generalizes complex but maximally specific linguistic contexts, and PALKA performs both generalization and specialization of an initial extraction pattern. Where AUTOSLOG makes no attempt to limit the number of extraction patterns created, CRYSTAL's covering algorithm derives the minimum number of patterns that cover the examples in the train-

ing corpus. In addition, CRYSTAL and PALKA use automated feedback for the learning algorithm; AUTOSLOG requires human perusal of proposed patterns. CRYSTAL and PALKA, however, require more background knowledge in the form of a possibly domain-specific *semantic-class hierarchy*, a lexicon that indicates semantic-class information for each word and, in the case of PALKA, a set of trigger words. The parsers of both systems must also be able to accurately assign semantic-class information to words in an incoming text. No semantic hierarchy is needed for AUTOSLOG—a flat semantic feature list will suffice. Also, although AUTOSLOG's patterns perform best when semantic-class information is available, the learning algorithm and the resulting concept nodes can still operate effectively when no semantic-class information can be obtained.

There have been a few additional attempts to learn extraction patterns. Huffman's (1996) LIEP system learns patterns that recognize semantic relationships between two target noun phrases, that is, between two slot fillers of an information-extraction output template. The patterns describe the syntactic context that falls between the target noun phrases as well as the semantic class of the heads of the target phrases and all intervening phrases. I (Cardie 1993) used standard symbolic machine-learning algorithms (decision tree induction and a k nearest-neighbor algorithm) to identify the trigger word for an extraction pattern, the general linguistic context in which the pattern would be applied, and the type of concept that the pattern would identify. Califf and Mooney (1997) have recently applied relational learning techniques to acquire extraction patterns from news group job postings. Like CRYSTAL, their RAPIER system operates by generalizing an initial set of specific patterns. Unlike any of the previously mentioned systems, however, RAPIER learns patterns that specify constraints at the word level rather than the constituent level. As a result, only a part-of-speech tagger is required to process input texts.

Although much progress has been made on learning extraction patterns, many research issues still need to be resolved. Existing methods work well when the information to be extracted is explicitly denoted as a string in the text, but major extensions would be required to handle set fills. Furthermore, existing methods focus on the extraction of noun phrases. It is not clear that the same methods would work well for domains in which the extracted information is another syntactic type or is a component of a constituent rather than a complete constituent (for example, a group of noun

modifiers in a noun phrase). Finally, few of the methods described here have been evaluated on the same information-extraction tasks under the same conditions. Until a direct comparison of techniques is available, it will remain difficult to determine the relative advantages of one technique over another. A related open problem in the area is to determine, a priori, which method for learning extraction patterns will give the best results in a new extraction domain.

Coreference Resolution and Template Generation

In comparison to empirical methods for learning extraction patterns, substantially less research has tackled the problems of coreference resolution and template generation. As mentioned earlier, the goal of the coreference component is to determine when two phrases refer to the same entity. Although this task might not appear difficult, consider the following text from the MUC-6 (1995) corporate management succession domain. In this text, all the bracketed segments are coreferential:

[Motor Vehicles International Corp.] announced a major management shake-up.... [MVI] said the chief executive officer has resigned.... [The Big 10 auto maker] is attempting to regain market share.... [It] will announce significant losses for the fourth quarter.... A [company] spokesman said [they] are moving [their] operations to Mexico in a cost-saving effort.... [MVI, [the first company to announce such a move since the passage of the new international trade agreement],] is facing increasing demands from unionized workers.... [Motor Vehicles International] is [the biggest American auto exporter to Latin America].

The passage shows the wide range of linguistic phenomena that influence coreference resolution, including proper names, aliases, definite noun phrases, definite descriptions, pronouns, predicate nominals, and appositives. Unfortunately, different factors can play a role in handling each type of reference. In fact, discourse processing, and coreference in particular, has been cited as a major weakness of existing information-extraction systems. One problem is that most systems use manually generated heuristics to determine when two phrases describe the same entity, but generating good heuristics that cover all types of reference resolution is challenging. In particular, few discourse theories have been evaluated

empirically, and as a result, it is not clear what information to include in the heuristics. It is also difficult to design heuristics that combine multiple coreference cues effectively, given that the relative importance of each piece of information is unknown. Furthermore, most computational approaches to coreference resolution assume as input fully parsed sentences, often marked with additional linguistic attributes such as grammatical function and thematic role information. Information-extraction systems do not normally have such detailed parse information available: The robust partial parsing algorithms used by most information-extraction systems offer wider coverage in exchange for less syntactic information. A further complication in developing trainable coreference components for an information-extraction system is that discourse analysis is based on information discerned by earlier phases of processing. Thus, any coreference algorithm must take into account the accumulated errors of the earlier phases as well as the fact that some information that would aid the coreference task might be missing. Finally, the coreference component of an information-extraction system must be able to handle the myriad forms of coreference across different domains.

Empirical methods for coreference were designed to address these problems. Unlike the methods for learning extraction patterns, algorithms for building automatically trainable coreference-resolution systems have not required the development of learning algorithms designed specifically for the task. By recasting the coreference problem as a classification task, any of a number of standard inductive-learning algorithms can be used. Given two phrases and the context in which they occur, for example, the coreference-learning algorithm must classify the phrases with respect to whether they refer to the same object. Here, I describe two systems that use inductive-classification techniques to automatically acquire coreference-resolution heuristics: MLR (machine-learning-based resolver) (Aone and Bennett 1995) and RESOLVE (McCarthy and Lehnert 1995).

Both MLR and RESOLVE use the same general approach, which is depicted in figure 5. First, a training corpus is annotated with coreference information; namely, all phrases that refer to the same object are linked using the annotations. Alternatively, just the best (usually the most recent) antecedent for each referent is marked. Training examples for presentation to the machine-learning algorithm are then created from the corpus. There will be one in-

... the goal of the coreference component is to determine when two phrases refer to the same entity.

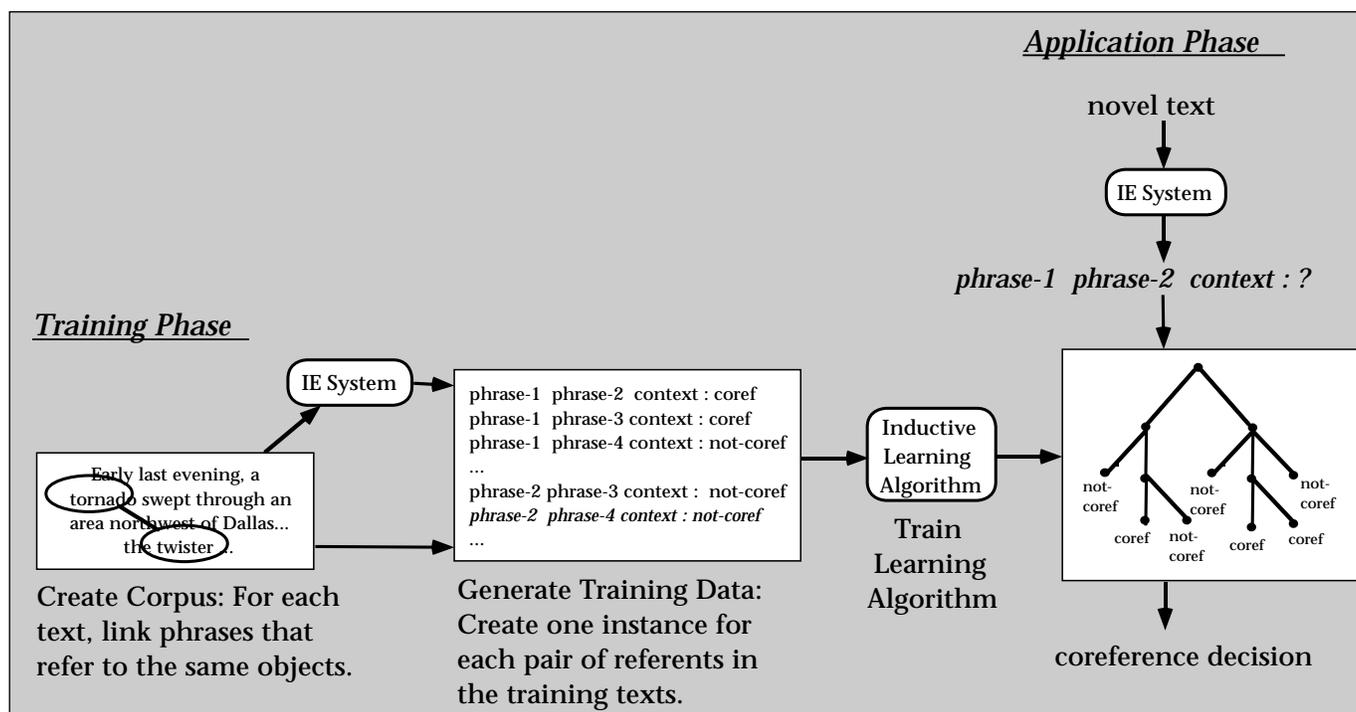


Figure 5. A Machine-Learning Approach to Coreference Resolution.

stance for every possible pairing of referents in the training texts: Some of these are *positive examples* in that they correspond to phrases that are coreferent; others are *negative examples* in that they correspond to phrases that are not referring to the same object. The exact form of the instances depends on the learning algorithm, but for the inductive-learning algorithm used by MLR and RESOLVE, the training examples contain (1) a list of features, or attribute-value pairs, that describe the phrases under consideration and the context in which they occur and (2) supervisory information in the form of a *class value* that indicates whether the two phrases are coreferent. The specific features used depend on the kinds of information available to the information-extraction system when the coreference decision must be made. More details on the creation of training data are given later.

Once the data set has been derived from the corpus, it is presented to the machine-learning algorithm, which uses the examples to derive a concept description for the coreference-resolution task. Figure 5 shows a concept description in the form of a decision tree, but the actual form depends on the particular learning algorithm employed. The idea is that after training, this concept description can be used to decide whether two phrases in an unseen text refer to the same object. This process is

shown as the application phase in figure 5. The information-extraction system processes a new text and reaches a point where coreference decisions must be made. For each such decision, the NLP system creates a test instance. The *test instance* uses the same feature set as the training instances: Its features describe a discourse entity, its possible antecedent, and their shared context. The test instance is given to the learned concept description for classification as either coreferent or not, and the decision is returned to the information-extraction system.

Both MLR and RESOLVE use this general method for automatically constructing coreference components for their information-extraction systems. Both use the widely available C4.5 decision-tree-induction system (Quinlan 1992) as the inductive-learning component. There are, however, a number of differences in how each system instantiated and evaluated the general approach of figure 5. McCarthy tested RESOLVE on the MUC-5 business joint-venture corpus (English version), and Aone and Bennett tested MLR on the Japanese corpus for the same information-extraction domain. The evaluation of MLR focused on anaphors involving entities tagged as organizations (for example, companies, governments) by the sentence-analysis phase of their information-extraction system. The evaluation of RESOLVE focused more specifically on organizations

that had been identified as a party in a joint venture by the extraction component.

Both systems used feature representations that relied only on information that earlier phases of analysis could provide. However, MLR's data set was generated automatically by its information-extraction system, while RESOLVE was evaluated using a manually generated, noise-free data set. In addition, the feature sets of each system varied markedly. MLR's training and test instances were described in terms of 66 features that describe (1) lexical features of each phrase (for example, whether one phrase contains a character subsequence of the other), (2) the grammatical role of the phrases, (3) semantic-class information, and (4) relative positional information. Although all attributes of the MLR representation are domain independent, the values for some attributes can be domain specific.

RESOLVE's instance representation, on the other hand, contains a number of features that are unabashedly domain specific. Its representation includes eight features including whether each phrase contains a proper name (two features), whether one or both phrases refer to the entity formed by a joint venture (three features), whether one phrase contains an alias of the other (one feature), whether the phrases have the same base noun phrase (one feature), and whether the phrases originate from the same sentence (one feature). Note that a number of RESOLVE's features correspond to those used in MLR, for example, the alias feature of RESOLVE versus the character subsequence feature of MLR.

RESOLVE and MLR were evaluated using data sets derived from 50 and 250 texts, respectively. RESOLVE achieved recall and precision levels of 80 percent to 85 percent and 87 percent to 92 percent (depending on whether the decision tree was pruned). A baseline system that always assumed that the candidate phrases were not coreferent would also achieve relatively high scores given that negative examples made up 74 percent of the RESOLVE data set. MLR, however, achieved recall and precision levels of 67 percent to 70 percent and 83 percent to 88 percent (depending on the parameter settings of the training configuration). For both MLR and RESOLVE, recall and precision are measured with respect to the coreference task only, not the full information-extraction task.

Without additional experiments, it is impossible to know whether the differences in results depend on the language (English versus Japanese), the slight variations in training-testing methodology, the degree of noise in the data, or the feature sets used. Interestingly, MLR does

well because a single feature—the character subsequence feature—can reliably predict coreference for phrases that are proper names for organizations, which make up almost half of the instances in the data set. Performance on non-proper-name organization referents was much lower. Definite noun phrases, for example, reached only 44-percent recall and 60-percent precision. Nevertheless, an important result for both MLR and RESOLVE was that each significantly outperformed a coreference system that had been developed manually for their information-extraction systems.

In a subsequent evaluation, the RESOLVE system competed in the MUC-6 coreference competition where it achieved scores of 41-percent to 44-percent recall and 51-percent to 59-percent precision after training on only 25 texts. This result was somewhat below the five best systems, which achieved 51-percent to 63-percent recall and 62-percent to 72-percent precision. All the better-performing systems, however, used manually encoded coreference algorithms. Like some of the manually coded systems, RESOLVE only attempted to resolve references to people and organizations. In fact, it was estimated that a good proper name–alias recognizer would have produced a coreference system with relatively good performance—about 30-percent recall and, possibly, 90-percent precision. One should note, however, that the interannotator agreement for marking coreference in 17 articles was found to be 80-percent recall and 82-percent precision, with definite descriptions (for example, [MVI, [the first company to announce such a move since the passage of the new international trade agreement]]) and bare nominals (for example, “A [company] spokesman”) accounting for most of the discrepancies.

Overall, the results for coreference resolution are promising. They show that it is possible to develop automatically trainable coreference systems that can compete favorably with manually designed systems. In addition, they show that specially designed learning algorithms need not be developed because standard machine-learning algorithms might be up to the challenge. There is an additional advantage to applying symbolic machine-learning techniques to problems in natural language understanding: They offer a mechanism for evaluating the usefulness of different knowledge sources for any task in an NLP system that can be described as a classification problem. Examination of the coreference decision trees created by C4.5, for example, will indicate which knowledge sources are more important for the task: The knowledge source

corresponding to a feature tested at node i in the tree is probably more important than the knowledge sources corresponding to the features tested below it in the tree. Furthermore, once the data set is created, it is a simple task to run multiple variations of the learning algorithm, giving each variation access to a different subset of features. As a result, empirical methods offer data-driven feedback for linguistic theories and system developers alike.

Still, much research remains to be done. The machine-learning approach to coreference should be tested on additional types of anaphor using a variety of feature sets, including feature sets that require no domain-specific information. In addition, if the approach is to offer a general, task-independent solution to the coreference problem, then the role of domain-specific information for coreference resolution must be determined empirically, and the methods must be evaluated outside the context of information extraction. The relative effect of errors from the preceding phases of text analysis on learning algorithm performance must also be investigated.

There have been few attempts to use empirical methods for other discourse-level problems that arise in information extraction. BBN Systems and Technologies has developed a probabilistic method for determining paragraph relevance for its information-extraction system (Weischedel et al. 1993); it then uses the device to control the recall-precision trade-off. I have used symbolic machine-learning techniques to learn relative pronoun disambiguation heuristics (Cardie 1992a, 1992b). Thus, the information-extraction system can process a sentence such as "Castellar was kidnapped by members of the ELN, who attacked the mayor in his office" and infer that *members of the ELN* is the actor of the kidnapping as well as the implicit actor of the attack in the second clause. Two trainable systems that simultaneously tackle merging and template generation have also been developed: TTG (Dolan et al. 1991) and WRAP-UP (Soderland and Lehnert 1994). Both systems generate a series of decision trees, each of which handles some piece of the template-generation or merging tasks, for example, deciding whether to merge two templates into one or deciding when to split an existing template into two or more templates. WRAP-UP used 91 decision trees to make these decisions for the MUC-5 microelectronics domain based on features of the entities extracted from each clause in an input text. Unfortunately, the information-extraction systems that used these trainable discourse components did not perform nearly as well as

systems that used manually generated merging and template-generation subsystems. Additional research is needed to determine the feasibility of an entirely trainable discourse component. Finally, statistical approaches to template merging are also beginning to surface. Kehler (1997), for example, introduced a method for assigning a probability distribution to coreference relationships, as encoded in competing sets of output templates. His initial experiments indicate that the method compares favorably with the greedy approach to template merging that is used in SRI International's FASTUS information-extraction system (Appelt et al. 1995).

Future Directions

Research in information extraction is new. Research in applying learning algorithms to problems in information extraction is even newer: We are only beginning to understand the techniques for automatically acquiring both domain-independent and domain-dependent knowledge for these task-driven systems. As a result, the field can take any number of exciting directions. First, like the trends in statistical language learning, a next step would be to explore unsupervised learning algorithms as a means for sidestepping the lack of large, annotated corpora for information-extraction tasks. In general, there is a dearth of learning algorithms that deal effectively with the relatively small amounts of data available to developers of information-extraction systems. A related but slightly different direction of research is to focus on developing techniques that allow end users to quickly train information-extraction systems for their own needs through interaction with the system over time, completely eliminating the need for intervention by NLP system developers. Many new learning methods will be needed to succeed in this task, not the least of which are techniques that make direct use of the answer keys of an information-extraction training corpus to automatically tune every component of the extraction system for a new domain. Finally, the robustness and generality of current learning algorithms should be investigated and extended by broadening the definition of information extraction to include the extraction of temporal, causal, or other complex relationships among events. The demand for information-extraction systems in industry, government, and education and for personal use is spiraling as more and more text becomes available online. The challenge for empirical methods in NLP is to continue to match this

Research in information extraction is new. Research in applying learning algorithms to problems in information extraction is even newer

demand by developing additional natural language-learning techniques that replace manual coding efforts with automatically trainable components and that make it increasingly faster and easier to build accurate and robust information-extraction systems in new domains.

Acknowledgments

Preparation of this article was supported in part by National Science Foundation CAREER Award IRI-9624639.

Notes

1. The development phase has varied from year to year but has ranged from about one to nine months.
2. A newer version of AUTOSLOG requires only that individual texts are marked as relevant or irrelevant to the domain (Riloff 1996). The learned concept nodes are then labeled according to type by hand.

References

- Aone, C., and Bennett, W. 1995. Evaluating Automated and Manual Acquisition of Anaphora Resolution Strategies. In *Proceedings of the Thirty-Third Annual Meeting of the Association for Computational Linguistics*, 122–129. Somerset, N.J.: Association for Computational Linguistics.
- Appelt, D. E.; Hobbs, J. R.; Bear, J.; Israel, D.; Kameyama, M.; Kehler, A.; Martin, D.; Myers, K.; and Tyson, M. 1995. SRI International FASTUS System: MUC-6 Test Results and Analysis. In *Proceedings of the Sixth Message-Understanding Conference (MUC-6)*, 237–248. San Francisco, Calif.: Morgan Kaufmann.
- Brill, E. 1995. Transformation-Based Error-Driven Learning and Natural Language Processing: A Case Study in Part-of-Speech Tagging. *Computational Linguistics* 21(4): 543–565.
- Califf, M. E., and Mooney, R. J. 1997. Relational Learning of Pattern-Match Rules for Information Extraction. In *Proceedings of the ACL Workshop on Natural Language Learning*, 9–15. Somerset, N.J.: Association for Computational Linguistics.
- Cardie, C. 1993. A Case-Based Approach to Knowledge Acquisition for Domain-Specific Sentence Analysis. In *Proceedings of the Eleventh National Conference on Artificial Intelligence*, 798–803. Menlo Park, Calif.: American Association for Artificial Intelligence.
- Cardie, C. 1992a. Corpus-Based Acquisition of Relative Pronoun Disambiguation Heuristics. In *Proceedings of the Thirtieth Annual Meeting of the ACL*, 216–223. Somerset, N.J.: Association for Computational Linguistics.
- Cardie, C. 1992b. Learning to Disambiguate Relative Pronouns. In *Proceedings of the Tenth National Conference on Artificial Intelligence*, 38–43. Menlo Park, Calif.: American Association for Artificial Intelligence.
- Cardie, C., and Lehnert, W. 1991. A Cognitively Plausible Approach to Understanding Complicated Syntax. In *Proceedings of the Ninth National Conference on Artificial Intelligence*, 117–124. Menlo Park, Calif.: American Association for Artificial Intelligence.
- Charniak, E. 1993. *Statistical Language Learning*. Cambridge, Mass.: MIT Press.
- Chinchor, N.; Hirschman, L.; and Lewis, D. 1993. Evaluating Message-Understanding Systems: An Analysis of the Third Message-Understanding Conference (MUC-3). *Computational Linguistics* 19(3): 409–449.
- Craven, M.; Freitag, D.; McCallum, A.; Mitchell, T.; Nigam, K.; and Quek, C. Y. 1997. Learning to Extract Symbolic Knowledge from the World Wide Web, Internal report, School of Computer Science, Carnegie Mellon University.
- Daelemans, W.; Zavrel, J.; Berck, P.; and Gillis, S. 1996. MBT: A Memory-Based Part-of-Speech Tagger-Generator. In *Proceedings of the Fourth Workshop on Very Large Corpora*, eds. E. Ejerhed and I. Dagan, 14–27. Copenhagen: ACL SIGDAT.
- Dolan, C.; Goldman, S.; Cuda, T.; and Nakamura, A. 1991. Hughes Trainable Text Skimmer: Description of the TTS System as Used for MUC-3. In *Proceedings of the Third Message-Understanding Conference (MUC-3)*, 155–162. San Francisco, Calif.: Morgan Kaufmann.
- Glasgow, B.; Mandell, A.; Binney, D.; Ghemri, L.; and Fisher, D. 1997. MITA: An Information-Extraction Approach to Analysis of Free-Form Text in Life Insurance Applications. In *Proceedings of the Ninth Conference on Innovative Applications of Artificial Intelligence*, 992–999. Menlo Park, Calif.: American Association for Artificial Intelligence.
- Holowczak, R. D., and Adam, N. R. 1997. Information Extraction-Based Multiple-Category Document Classification for the Global Legal Information Network. In *Proceedings of the Ninth Conference on Innovative Applications of Artificial Intelligence*, 1013–1018. Menlo Park, Calif.: American Association for Artificial Intelligence.
- Huffman, S. 1996. Learning Information-Extraction Patterns from Examples. In *Symbolic, Connectionist, and Statistical Approaches to Learning for Natural Language Processing*, eds. S. Wermter, E. Riloff, and G. Scheler, 246–260. Lecture Notes in Artificial Intelligence Series. New York: Springer.
- Kehler, A. 1997. Probabilistic Coreference in Information Extraction. In *Proceedings of the Second Conference on Empirical Methods in Natural Language Processing*, eds. C. Cardie and R. Weischedel, 163–173. Somerset, N.J.: Association for Computational Linguistics.
- Kim, J.-T., and Moldovan, D. I. 1995. Acquisition of Linguistic Patterns for Knowledge-Based Information Extraction. *IEEE Transactions on Knowledge and Data Engineering* 7(5): 713–724.
- Lehnert, W. 1990. Symbolic-Subsymbolic Sentence Analysis: Exploiting the Best of Two Worlds. In *Advances in Connectionist and Neural Computation Theory*, eds. J. Barnden and J. Pollack, 135–164. Norwood, N.J.: Ablex.
- Lehnert, W., and Sundheim, B. 1991. A Performance Evaluation of Text Analysis Technologies. *AI Maga-*

zine 12(3): 81–94.

Lehnert, W.; Cardie, C.; Fisher, D.; Riloff, E.; and Williams, R. 1991. University of Massachusetts: Description of the CIRCUS System as Used in MUC-3. In *Proceedings of the Third Message-Understanding Conference (MUC-3)*, 223–233. San Francisco, Calif.: Morgan Kaufmann.

Lehnert, W.; Cardie, C.; Fisher, D.; McCarthy, J.; Riloff, E.; and Soderland, S. 1992. University of Massachusetts: Description of the CIRCUS System as Used in MUC-4. In *Proceedings of the Fourth Message-Understanding Conference (MUC-4)*, 282–288. San Francisco, Calif.: Morgan Kaufmann.

McCarthy, J. F., and Lehnert, W. G. 1995. Using Decision Trees for Coreference Resolution. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence*, ed. C. Mellish, 1050–1055. Menlo Park, Calif.: International Joint Conferences on Artificial Intelligence.

Magerman, D. M. 1995. Statistical Decision-Tree Models for Parsing. In *Proceedings of the Thirty-Third Annual Meeting of the ACL*, 276–283. Somerset, N.J.: Association for Computational Linguistics.

Marcus, M.; Marcinkiewicz, M.; and Santorini, B. 1993. Building a Large Annotated Corpus of English: The Penn Tree Bank. *Computational Linguistics* 19(2): 313–330.

MUC-3. 1991. *Proceedings of the Third Message-Understanding Conference (MUC-3)*. San Francisco, Calif.: Morgan Kaufmann.

MUC-4. 1992. *Proceedings of the Fourth Message-Understanding Conference (MUC-4)*. San Francisco, Calif.: Morgan Kaufmann.

MUC-5. 1994. *Proceedings of the Fifth Message-Understanding Conference (MUC-5)*. San Francisco, Calif.: Morgan Kaufmann.

MUC-6. 1995. *Proceedings of the Sixth Message-Understanding Conference (MUC-6)*. San Francisco, Calif.: Morgan Kaufmann.

Quinlan, J. R. 1992. *C4.5: Programs for Machine Learning*. San Francisco, Calif.: Morgan Kaufmann.

Ramshaw, L. A., and Marcus, M. P. 1995. Text Chunking Using Transformation-Based Learning. In *Proceedings of the Thirty-Third Annual Meeting of the ACL*, 82–94. Somerset, N.J.: Association for Computational Linguistics.

Riloff, E. 1996. Automatically Generating Extraction Patterns from Untagged Text. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence*, 1044–1049. Menlo Park, Calif.: American Association for Artificial Intelligence.

Riloff, E. 1993. Automatically Constructing a Dictionary for Information-Extraction Tasks. In *Proceedings of the Eleventh National Conference on Artificial Intelligence*, 811–816. Menlo Park, Calif.: American Association for Artificial Intelligence.

Soderland, S. 1997. Learning to Extract Text-Based Information from the World Wide Web. In *Proceedings of the Third International Conference on Knowledge Discovery and Data Mining (KDD-97)*, 251–254. Menlo Park, Calif.: AAAI Press.

Soderland, S., and Lehnert, W. 1994. Corpus-Driven

Knowledge Acquisition for Discourse Analysis. In *Proceedings of the Twelfth National Conference on Artificial Intelligence*, 827–832. Menlo Park, Calif.: American Association for Artificial Intelligence.

Soderland, S.; Fisher, D.; Aseltine, J.; and Lehnert, W. 1995. CRYSTAL: Inducing a Conceptual Dictionary. In *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence*, 1314–1319. Menlo Park, Calif.: International Joint Conferences on Artificial Intelligence.

Soderland, S.; Aronow, D.; Fisher, D.; Aseltine, J.; and Lehnert, W. 1995. Machine Learning of Text-Analysis Rules for Clinical Records, Technical Report, TE-39, Department of Computer Science, University of Massachusetts.

Thompson, C. A.; Mooney, R. J.; and Tang, L. R. 1997. Learning to Parse Natural Language Database Queries into Logical Form. In *Proceedings of the ML-97 Workshop on Automata Induction, Grammatical Inference, and Language Acquisition*. Somerset, N.J.: Association for Computational Linguistics.

Weischedel, R.; Meteer, M.; Schwartz, R.; Ramshaw, L.; and Palmucci, J. 1993. Coping with Ambiguity and Unknown Words through Probabilistic Models. *Computational Linguistics* 19(2): 359–382.

Weischedel, R.; Ayuso, D.; Boisen, S.; Fox, H.; Matsukawa, T.; Papageorgiou, C.; MacLaughlin, D.; Sakai, T.; Abe, H. J. H.; Miyamoto, Y.; and Miller, S. 1993. BBN's PLUM Probabilistic Language-Understanding System. In *Proceedings, TIPSTER Text Program (Phase I)*, 195–208. San Francisco, Calif.: Morgan Kaufmann.

Will, C. A. 1993. Comparing Human and Machine Performance for Natural Language Information Extraction: Results from the TIPSTER Text Evaluation. In *Proceedings, TIPSTER Text Program (Phase I)*, 179–194. San Francisco, Calif.: Morgan Kaufmann.

Zelle, J., and Mooney, R. 1994. Inducing Deterministic Prolog Parsers from Tree Banks: A Machine-Learning Approach. In *Proceedings of the Twelfth National Conference on Artificial Intelligence*, 748–753. Menlo Park, Calif.: American Association for Artificial Intelligence.



Claire Cardie is an assistant professor in the Computer Science Department at Cornell University. She received her Ph.D. in 1994 from the University of Massachusetts at Amherst. Her current research is in natural language learning, case-based learning, and the application of natural language-understanding techniques to problems in information retrieval. Her e-mail address is cardie@cs.cornell.edu.