

A Logical Framework for Template Creation and Information Extraction

David Corney¹, Emma Byrne², Bernard Buxton¹, and David Jones¹

¹ Department of Computer Science, University College London,
Gower Street, London WC1E 6BT, UK

D.Corney@ucl.ac.uk, B.Buxton@cs.ucl.ac.uk, D.Jones@cs.ucl.ac.uk

² School of Primary Care and Population Sciences, University College London,
Highgate Hill, London N19 5LW, UK
emma.byrne@ucl.ac.uk

(Final submitted version)

Abstract. Information extraction is the process of automatically identifying facts of interest from pieces of text, and so transforming free text into a structured database. Past work has often been successful but ad hoc, and in this paper we propose a more formal basis from which to discuss information extraction. We introduce a framework which will allow researchers to compare their methods as well as their results, and will help to reveal new insights into information extraction and text mining practices.

One problem in many information extraction applications is the creation of templates, which are textual patterns used to identify information of interest. Our framework describes formally what a template is and covers other typical information extraction tasks. We show how common search algorithms can be used to create and optimise templates automatically, using sequences of overlapping templates, and we develop heuristics that make this search feasible. Finally we demonstrate a successful implementation of the framework and apply it to a typical biological information extraction task.

Keywords: Information extraction; biological text mining.

1 Introduction

Information extraction (IE) [1] has developed over recent decades with applications analysing text from news sources [2], financial sources [3], and biological research papers [4–6]. Competitions such as MUC and TREC have been promoted as using real text sources to highlight problems in the real world, and more recently TREC has included a genomics track [7], again highlighting biology and medicine as growing areas of IE research. It has long been recognised that there is a need to share resources between research groups in order to allow a fair comparison of their different systems and to motivate and direct further research. We strongly feel that there is also a need to provide a *theoretical* framework within which these information extraction systems can be described,

compared and developed, by identifying key issues explicitly. The framework we present here will allow researchers to compare their methods as well as their results, and also provides new methods for template creation. This will aid the identification of important issues within the field, allowing us to identify the questions to ask as well as to formulate some answers.

The terms information extraction and text mining are often used interchangeably. Some authors use the term text mining to suggest the detection of novelty, such as combining information from several sources to generate and test new hypotheses [8]. In contrast, IE extracts only that which is explicitly stated. This framework focuses on IE and template creation, but it also applies to text mining.

Information extraction is a large and diverse research area. One approach widely used is to develop modular systems such as GATE [9], so that each component can be optimised individually. One of the most challenging of these is the template component. A template is a textual pattern designed to identify “interesting” information to be extracted from documents, where “interesting” is relative to the user’s intentions. An ideal template can be used to extract a large proportion of the interesting information available with only a little uninteresting information.

Different types of templates exist, but in general, they can be thought of as regular expressions over words and the features of those words. Standard regular expressions match sequences of characters, but IE templates can also match features of words. To guide our discussion, consider the sentence “The cat sat on the mat”. Informally, one template that would match that sentence is “The ANIMAL VERB on the FLOOR_COVERING”, where ‘ANIMAL’ and ‘FLOOR_COVERING’ are pre-defined semantic categories, and ‘VERB’ is a part-of-speech label. A different template that would match the same sentence is “DETERMINER * * PREPOSITION DETERMINER *”, where each ‘*’ is a wildcard, matching any single word, and the other symbols are part-of-speech labels. Any sentence can be matched with a large number of templates, and many templates match a large number of sentences. This makes template creation a challenging problem.

Although it covers several key areas, this paper focuses on template creation. Currently, templates are typically designed by hand, which can be laborious and limits the rapid application of IE to new domains. There have been several attempts at automatic template creation [10–12], and there are likely to be more in the future. To the best of our knowledge, no such system has demonstrated widespread applicability, but tend to be successful only within narrow domains. Some systems are effective, but require extensive annotation of a training set [13], which is also laborious. One way to view the automatic creation of useful templates is as a search problem of a kind familiar to the artificial intelligence community [14, ch. 3–4]. To formulate it this way, we need to define the space of candidate solutions (i.e. templates); a means of evaluating and comparing these candidate solutions; a means of generating new candidate solutions; and an algorithm for guiding the search (including starting and stopping). Any useful framework describing IE must provide a way to define and create templates, and our framework proposes using these AI search methods,

an idea we expand in Sect. 6.2, where we “grow” useful templates from given seed phrases.

One alternative to using templates is *co-occurrence analysis* [15]. This identifies pieces of text (typically sentences, abstracts or entire documents) that mention two entities, and assumes that this implies that the two entities are in some way related. Within our framework, this can be seen as a special case of a template, albeit a very simple one, as we show in Sect. 2.3.

The framework itself is presented in Sections 2–5, with the subsequent sections discussing various implementation issues.

Section 2 defines various concepts formally, moving from words and documents to templates and information extraction. Section 3 describes how templates can be ordered according to how specific or general they are, as a precursor to template creation and optimisation. Section 4 discusses how to modify a template to make it more general. Section 5 gives formal definitions of recall and precision within our framework and discusses how they might be estimated in practice. Section 6 discusses heuristic search algorithms and their implementation and includes a detailed example, before a concluding discussion.

A shorter form of this work is published in [16].

2 Basic Definitions

In this section, we define several terms culminating in a formal definition of information extraction templates.

Definition 1. A literal λ is a word in the form of an ordered list of characters. We assume implicitly a fixed alphabet of characters.

Examples: “cat”, “jumped”, “2,5-dihydroxybenzoic”.

Definition 2. A document d is a tuple (ordered list) of literals: $d = \langle \lambda_1, \lambda_2, \dots, \lambda_{|d|} \rangle$.

Examples: $d_1 = \langle \text{the, cat, sat, on, the, mat} \rangle$, $d_2 = \langle \text{a, mouse, ran, up, the, clock} \rangle$.

Definition 3. A corpus D is a set of documents: $D = \{d_1, d_2, \dots, d_{|D|}\}$.

Example: $D_1 = \{d_1, d_2\}$.

Definition 4. A lexicon A is the set of all literals found in all documents in a corpus: $A_D = \{\lambda \mid \lambda \in d \text{ and } d \in D\}$.

Example: $A_{D_1} = \{\text{the, cat, sat, on, mat, a, mouse, ran, up, clock}\}$.

Every word has a set of attributes, such as its part-of-speech or its membership of a semantic class, which we now discuss. Although particular attributes are not a formal part of the framework, they are used in various illustrative examples throughout this paper.

Words that share a common stem, or root, typically share a common meaning, such as the words “sit”, “sitting” and “sits”. It is therefore common practice in information retrieval to index words according to their stem to improve the performance [17]. Similarly in information extraction, it is often helpful to identify words that share a common stem. The most common approach is to remove suffixes to produce a single stem for each word [17], although in principle, each word could have multiple stems, such as if prefixes were removed independently of suffixes.

Words may also belong to pre-defined semantic categories, such as “business”, “country” or “protein”. One common way to define these semantic categories is by using gazetteers. A gazetteer is a named list of words and phrases that belong to the same category. Rather than simple lists, some ontologies are based on hierarchies or directed acyclic graphs, such as MeSH³ and GO⁴ respectively. In this framework, we are not concerned with the nature of such categories, but assume only that there exists some method for assigning such attributes to individual words.

The role of each word in a sentence is defined by its *part of speech*, or lexical category. Common examples are noun, verb and adjective, although these are often subdivided into more precise categories such as “singular common noun”, “plural common noun”, “past tense verb” and so on. The part of speech can usually only be ascertained for a word in a given context. For example, compare “He cut the bread” to “The cut was deep”. In practice, an implementation may limit this to exactly one label per word, based on the context of that word. Following the Penn Treebank tags [18], in some examples we use the symbol “DT” to represent determiners such as “the”, “a” and “this”; “VB” to represent verbs in their base form, such as “sit” and “walk”; “VBD” to represent past-tense verbs, such as “sat” and “walked”; “NN” to represent common singular nouns, such as “cat” and “shed” and so on.

We also introduce wildcards as an extension to the idea of word attributes. In regular expressions, a wildcard can “stand in” for a range of characters, and we use the same notion here to represent ranges of words. For example, we use the symbol “*” as the universal wildcard which can be replaced by any word in the lexicon. Then every word has the attribute “*”. We also use the symbol “?” to represent any word *or no word at all*. We discuss these wildcards further in Sect. 4.3.

Other categories may be introduced to capture other attributes, such as orthography (e.g. upper case, lower case or mixed case), word length, language and so on. A parser could be used to label words as belonging to different types of phrases, such as verb phrases and noun phrases. We could also treat punctuation symbols as literals if required, or as a separate category. However, the categories described above are sufficient to allow us to develop and demonstrate the framework.

Definition 5. *A category κ is set of attributes of words of the same type. Common categories include “parts of speech” and “stems”.*

³MeSH, Medical Subject Headings, <http://www.nlm.nih.gov/mesh>

⁴Gene Ontology, <http://www.geneontology.org>

For convenience, we will label certain categories in these and subsequent examples. This is not part of the framework but reflects categories likely to be used in a practical implementation. In particular, we use A to label the category “literals”; Π for “parts of speech”; Γ for “gazetteers”; Σ for “stems”; and Ω for “wildcards”.

Example:

$\kappa_A = \{\text{the, cat, sat, on, mat, mouse, ...}\}$

$\kappa_\Sigma = \{\text{the_stem, cat_stem, sit_stem, on_stem, mat_stem, mouse_stem, ...}\}$

$\kappa_\Pi = \{\text{DT, NN, VBD, IN, ...}\}$

$\kappa_\Gamma = \{\text{FELINE, RODENT, ANIMAL, FLOOR_COVERING, ...}\}$

$\kappa_\Omega = \{*, ?\}$

We use the suffix ‘_stem’ in stem labels to avoid confusing them with the corresponding literal.

Definition 6. *Let K be a set of categories of attributes. Each element κ of K is a single category of word attributes.*

Example: $K_1 = \{\kappa_A, \kappa_\Sigma, \kappa_\Pi, \kappa_\Gamma, \kappa_\Omega\}$.

Definition 7. *A term t is a value that an attribute may take, i.e. an element of a category of word attributes.*

Examples: $t_1 = \text{cat}$, $t_2 = \text{NN}$, $t_3 = \text{FELINE}$, where $t_1 \in \kappa_A$, $t_2 \in \kappa_\Pi$, $t_3 \in \kappa_\Gamma$.

Definition 8. *We define a template element T to be a set of terms belonging to a single category. Let $T = \{t_1, t_2, \dots, t_n\}$, such that $t_i \in T$. Then $t_i \in \kappa \iff t_j \in \kappa, \forall t_j \in T$.*

Examples:

$T_1 = \{\text{NN, VBD}\}$

$T_2 = \{\text{FELINE, RODENT, FLOOR_COVERING}\}$

The set $\{\text{NN, FELINE}\}$ is not a template element because “NN” and “FELINE” belong to different categories, namely κ_Π and κ_Γ respectively. The name “template element” refers to templates as defined in Definition 13 below.

Definition 9. *The attributes of a literal are the set of template elements defining the values of the literal in each category. We first define the set of attributes of a literal λ for a particular category κ as $\alpha(\lambda, \kappa) = \{T \mid \forall t \in T, t \in \kappa \text{ and } \lambda \text{ has attribute } t\}$. The set of all attributes of a literal is the union of the attributes in each category: $\alpha(\lambda) = \bigcup_{\kappa \in K} \alpha(\lambda, \kappa)$. If a literal has no value for a particular category, then the category is omitted from the set α .*

When we say “ λ has attribute t ”, we assume that this relationship is defined outside of the framework. For example, there maybe functions to assign a stem attribute to a word, or to assign a particular semantic category to any of a given list of words.

For convenience, we label the attributes using the category label as a subscript in these examples.

Examples:

$\alpha_A(\text{cat}) = \{\text{cat}\}$
 $\alpha_R(\text{cat}) = \{\text{FELINE, ANIMAL}\}$
 $\alpha_{II}(\text{cat}) = \{\text{NN}\}$
 $\alpha_\Sigma(\text{cat}) = \{\text{cat_stem}\}$
 $\alpha(\text{cat}) = \{\{\text{cat}\}, \{\text{NN}\}, \{\text{FELINE, ANIMAL}\}, \{\text{cat_stem}\}\}$

In the case of $\alpha(\text{the})$, the word “the” has a part-of-speech tag “DT” (determiner) and the obvious literal, but no gazetteer or stem entries. So $\alpha_{II}(\text{the}) = \{\text{DT}\}$, and $\alpha_R(\text{the})$ is undefined, and so $\alpha(\text{the}) = \{\{\text{the}\}, \{\text{DT}\}\}$. The set of literal attributes A has the special property that every word has exactly one literal. Other categories in K may contain terms such that one literal may correspond to one or more terms, or to no term at all. For example, one literal may belong to more than one gazetteer, while another literal may belong to none. Therefore for any λ , $|\alpha_A(\lambda)| = 1$. As a consequence, $|\alpha(\lambda)| \geq 1$.

2.1 Membership of Terms

We now define the concept of membership to refer to the set of literals that share a particular attribute.

Definition 10. We define the members μ of a term t as being the set of all literals that share the attribute value defined by that term. $\mu(t) = \{\lambda | t \in \alpha(\lambda)\}$. Also, we define the members of a set of terms (such as a template element) as the union of the members of each term in the set: $\mu(\{t_1, t_2, \dots, t_n\}) = \bigcup_{i=1}^n \mu(t_i)$

Examples:

$\mu(\text{sit_stem}) = \{\text{sit, sits, sat, sitting}\}$.
 $\mu(\text{FELINE}) = \{\text{cat, lion, tiger, \dots}\}$.
 $\mu(\text{RODENT}) = \{\text{mouse, rat, hamster, \dots}\}$.
 $\mu(\{\text{FELINE, RODENT}\}) = \{\text{cat, lion, tiger, mouse, rat, hamster, \dots}\}$.
 $\mu(\alpha_A(\text{cat})) = \{\text{cat}\}$. I.e. the membership of the literal category of a literal is the set containing only the literal itself.

Definition 11. We also define membership for a term limited to a particular document or set of documents: $\mu(t, d) = \{\lambda | \lambda \in d, \lambda \in \mu(t)\}$, and $\mu(t, D) = \bigcup_{d \in D} \mu(t, d)$

Examples:

$\mu(\text{NN}, d_1) = \{\text{cat, mat}\}$.
 $\mu(\text{ANIMAL}, D_1) = \{\text{cat, mouse}\}$.

2.2 Templates and Document Fragments

Definition 12. We define a fragment of a document as being a tuple of successive literals taken from some document d . If

$$d = \langle \lambda_1, \lambda_2, \dots, \lambda_a, \lambda_{a+1}, \dots, \lambda_{a+b-1}, \dots, \lambda_{|d|} \rangle,$$

then

$$f(d, a, b) = \langle \lambda_a, \lambda_{a+1}, \dots, \lambda_{a+b-1} \rangle.$$

I.e. the function $f(d, a, b)$ returns a tuple of b words in order, from d , starting with the a^{th} word. f_1 is the first word of the fragment, i.e. λ_a , f_2 is the second word, λ_{a+1} , and so on. Note that $|f| = b$.

Example: If $d_1 = \langle \text{the, cat, sat, on, the, mat} \rangle$, then $f(d_1, 4, 3) = \langle \text{on, the, mat} \rangle$.

Definition 13. A template τ is a tuple of one or more template elements, $\langle T_1, T_2, \dots, T_n \rangle$, where $T_1 = \{t_{1,1}, t_{1,2}, \dots\}$, $T_2 = \{t_{2,1}, t_{2,2}, \dots\}$ and so on. $|\tau|$ is the number of template elements in template τ , and is always greater than zero. Each template element T_i within a template consists of one or more terms of the same type.

Example:

$\tau_1 = \langle \{\text{the}\}, \{\text{FELINE, RODENT}\}, \{\text{VB, VBN}\} \rangle$.

Definition 14. A template matches a fragment of a document d if each successive template element in the template contains a term whose membership includes each successive word in the fragment. Let $f = \langle f_1, \dots, f_n \rangle$ be a fragment. Given a template $\tau = \langle T_1, T_2, \dots, T_n \rangle$, we extend the membership function thus:

$$\mu(\tau, d) = \{f \mid f \in d \text{ and } \forall i \in 1 \dots |\tau|, f_i \in \mu(T_i)\}$$

For a corpus D , the template matches the union of the template membership for each document: $\mu(\tau, D) = \bigcup_{d \in D} \mu(\tau, d)$.

This function returns a set of fragments, each of which consists of a tuple of literals that matches each successive element of the template τ , and each of which is found in a document in the corpus. Matching terms in this way is the core of information extraction. The words that are matched define the information that is to be extracted.

Example: Given a template $\tau_1 = \langle \{\text{DT}\}, \{\text{ANIMAL}\}, \{\text{VBD}\} \rangle$ and a corpus D_1 (as defined after Definition 3), then $\mu(\tau_1, D_1) = \{\langle \text{the, cat, sat} \rangle, \langle \text{a, mouse, ran} \rangle\}$.

2.3 Co-occurrence Analysis

Co-occurrence analysis assumes that two entities in the same piece of text are related, without attempting a more sophisticated linguistic analysis of the text. In our framework, this can be represented by a template (or set of templates) that defines the two entities with a series of wildcards between them.

For example, suppose we use co-occurrence analysis to discover every sentence in a corpus that mentions two entities, as matched by template elements T_i and T_j . Let us assume that all sentences to be considered are finite with a maximum length of Q words. Then we could define two template $\tau_1 = \langle T_1, \dots, T_i, \dots, T_j, \dots, T_Q \rangle$ and $\tau_2 = \langle T_1, \dots, T_j, \dots, T_i, \dots, T_Q \rangle$. Two templates are required if we wish to allow for sentences with the two entities in different orders. We replace every template element except for T_i and T_j with the wildcard element '?' so as to match any sentence containing words that match our terms. With three or more entities, larger sets of templates may be required.

3 Template Ordering

One motivation for creating this framework is to enable the use of common search algorithms for template creation. To do this effectively, we must define an *ordering* over the templates, which we can then use to develop practical search heuristics.

For any given document, each template matches a certain number of fragments. A template that matches every possible fragment is useless, as is one that matches no fragments at all. Somewhere between these two extremes of generic templates and specific templates, lie useful templates that match the interesting fragments only, so the aim of template creation is to find a suitable trade-off between the generic and the specific. We therefore suggest that a useful ordering is one based on the number of fragments that a template is likely to match. We can use such an order to search across a range of templates and explore the trade-off. For unseen text, it is impossible to predict the amount of information to be extracted in advance, so instead, we develop a heuristic ordering that approximates it.

In this section, we define possible orderings of terms and templates. In the next section, we define algorithms that use these orderings to modify templates. In Section 6 we discuss some search algorithms that might be used to locate optimal templates.

Before we consider an ordering over templates, we consider ordering over the terms that make up a template. We want to define the relations $t_1 >_s t_2$ to mean that term t_1 is a more specific attribute than term t_2 , and $t_1 \geq_s t_2$ to mean that t_1 is at least as specific as t_2 . Similarly, $t_1 <_s t_2$ and $t_1 \leq_s t_2$ mean “less specific than” and “no more specific than” respectively. We now define these by introducing “superset ordering”.

3.1 Superset Ordering of Terms and Template Elements

We start by defining a specificity ordering over terms such that each term matches every word that its antecedent matches, along with zero or more extra words. We call this superset ordering.

Definition 15. *Let \geq_s be the ordering over superset specificity. Let t_1, t_2 be terms. If $\mu(t_1) \subseteq \mu(t_2)$ then $t_1 \geq_s t_2$, and we say that t_1 is at least as specific as t_2 . If $\mu(t_1) \subset \mu(t_2)$ then $t_1 >_s t_2$, and we say that t_1 is more specific than t_2 .*

Examples:

Let FELINE and ANIMAL be two terms (specifically, gazetteers), such that

$\mu(\text{FELINE}) = \{\text{cat, lion, tiger, } \dots\}$ and $\mu(\text{ANIMAL}) = \{\text{antelope, dog, cat, } \dots, \text{lion, } \dots, \text{tiger, } \dots, \text{zebra}\}$. Then $\mu(\text{FELINE}) \subset \mu(\text{ANIMAL})$ and so $\text{FELINE} >_s \text{ANIMAL}$.

$\mu(\alpha_\Lambda(\text{cat})) = \{\text{cat}\}$.

$\alpha_\Gamma(\text{cat}) = \{\text{FELINE, ANIMAL}\}$

$\mu(\{\text{FELINE, ANIMAL}\}) = \{\text{cat, lion, tiger, antelope, dog, } \dots\}$

Therefore $\mu(\alpha_\Lambda(\text{cat})) \subset \mu(\alpha_\Gamma(\text{cat}))$.

Some specificity orderings are dependent on the categories of the two terms. For example, by definitions 9 and 10, it is clear that each literal is a member of all the terms that are attributes of the literal. Therefore $\mu(\lambda) \subseteq \mu(\alpha_\lambda) \forall \kappa \in K$ and therefore if $t_1 \in \kappa_A$, then $\forall t_2 \in K, t_1 \geq_s t_2$. In other words, terms that represent literals are at least as specific as any other terms. At the other extreme, the wildcards ‘*’ and ‘?’ match every word, so we can say that wildcard terms are no more specific than any other term. I.e. if $t_1 \in \kappa_\Omega$, then $\forall t_2 \in K, t_2 \geq_s t_1$.

Having defined an ordering over terms, we now consider sets of terms, and template elements in particular. Let T_1 and T_2 be two template elements. We say that T_1 is at least as specific as T_2 if and only if every literal matched by any term in T_1 is also matched by some term in T_2 :

$$T_1 \geq_s T_2 \iff \forall \lambda \in \mu(T_1), \lambda \in \mu(T_2)$$

Similarly, $T_1 >_s T_2 \iff \forall \lambda \in \mu(T_1), \lambda \in \mu(T_2)$ and $\exists \lambda \in \mu(T_2)$ such that $\lambda \notin \mu(T_1)$. I.e. T_1 is more specific than T_2 if every literal matched by a term in T_1 is also matched by a term in T_2 , and at least one literal is matched by a term in T_2 and not by a term in T_1 .

We can trivially extend these ‘more specific than’ relations to define ‘less specific than’ and ‘equally specific as’ relations:

$$\begin{aligned} t_1 >_s t_2 &\iff t_2 <_s t_1 \\ t_1 \geq_s t_2 &\iff t_2 \leq_s t_1 \\ t_1 =_s t_2 &\iff t_1 \geq_s t_2 \text{ and } t_1 \leq_s t_2 \end{aligned}$$

In this last case, $t_1 =_s t_2 \iff \mu(t_1) = \mu(t_2)$. Note that this is a narrower definition than requiring that $|\mu(t_1)| = |\mu(t_2)|$.

3.2 Ordering of Templates

Above, we have discussed ordering of terms and of template elements. Here we generalise this to discuss entire templates. We want to be able to compare two templates, τ_1 and τ_2 , and say which is more specific, i.e. which one matches fewer fragments. If τ_1 and τ_2 are related through superset generalisation, then for a given set of documents D , this is testing whether $\mu(\tau_1, D) \supset \mu(\tau_2, D)$ and therefore whether $|\mu(\tau_1, D)| > |\mu(\tau_2, D)|$, or vice versa, or they are equal. E.g. $\tau_1 >_s \tau_2 \iff \mu(\tau_1, D) \subset \mu(\tau_2, D)$. This depends on the corpus D and is potentially slow to evaluate, especially if D contains a large number of documents. We would rather have an estimate of the relative specificity which is independent of D , which will be useful when developing search heuristics.

Suppose τ_1 and τ_2 are almost identical, and differ only in one template element:

$$\begin{aligned} \tau_1 &= \langle T_1, T_2, \dots, T_i, \dots, T_n \rangle \\ \tau_2 &= \langle T_1, T_2, \dots, T'_i, \dots, T_n \rangle \end{aligned}$$

where $T_i \neq T'_i$. Then we can say that $\tau_1 >_s \tau_2 \iff T_i >_s T'_i$.

More generally, if τ_1 and τ_2 contain the same number of sets of terms, then

$$\tau_1 \geq_s \tau_2 \iff T_{1,i} \geq_s T_{2,i} \quad i = 1 \dots |\tau_1|, \text{ and}$$

$$\tau_1 >_s \tau_2 \iff T_{1,i} \geq_s T_{2,i} \quad i = 1 \dots |\tau_1| \text{ and } T_{1,j} >_s T_{2,j} \text{ for some } j.$$

In a slight extension to previous notation, we use $T_{n,i}$ to refer to the i^{th} element of template τ_n .

Also, if two templates are identical except that one is “missing” the first or last template element of the other, then the shorter of the two is less specific. I.e. if $\tau_1 = \langle T_1, T_2, \dots, T_{n-1}, T_n \rangle$, $\tau_2 = \langle T_1, T_2, \dots, T_{n-1} \rangle$ and $\tau_3 = \langle T_2, \dots, T_{n-1}, T_n \rangle$ then $\tau_1 \geq_s \tau_2$ and $\tau_1 \geq_s \tau_3$.

Although these relationships do not provide a complete ordering over all templates, they do allow us to compare similar templates, and this is sufficient to allow us to create and modify templates, and to develop useful search heuristics. We use this ordering to develop functions that create and modify templates in Sect. 4, and to develop methods to search efficiently for *good* templates in Sect. 6.

4 Template Generalisation

Whether created manually or automatically, templates are usually based on examples of “interesting” phrases. These phrases may be identified by hand (e.g. by a domain expert) or automatically (e.g. by information retrieval methods). These phrases are then used as “seeds” to help define more general templates. In this framework we will follow this “seed phrase” approach, and assume that we have some suitable phrases. We show how a wide range of templates can be created from each seed phrase. We first discuss how terms can be created and generalised, and then expand this to template creation and generalisation (Sect. 4.2).

4.1 Creating and Modifying Single Template Elements

We now bring together several concepts discussed above, and define functions that create and generalise template elements. This leads onto a discussion of creating and modifying entire templates.

Definition 16. *Given a literal, we want to create a new template element, which is simply a set containing the literal. We define the trivial function *initialise* for this purpose: $initialise(\lambda) = \{\lambda\}$.*

Have created a template element, we can then modify it. We now define a function that modifies any given template element to produce a new set of template elements that is at least as general as the element given. This is based on the notion of superset ordering (Sect. 3.1) in that the new template elements match a superset of the literals matched by the original template element. Furthermore, the new set of elements belongs to a specified category which is different from the category of the source element.

Definition 17. *We define a function to create a more general set of template elements from a given template element, such that all of the terms in each new template element are members of a specified category. Given a template element $T = \{t_1, t_2, \dots, t_{|T|}\}$ of category κ and*

given a target category $\kappa' \neq \kappa$, we create a set of template elements $\{T'_1, T'_2, \dots, T'_{|T'|}\}$:

$$\begin{aligned} \mathbf{generalise}(T, \kappa') &= \{T' | T' = \{t'_1, t'_2, \dots, t'_m\}, \text{ and } t'_1, t'_2, \dots, t'_m \in \kappa', \text{ and} \\ &\quad \forall t'_i \in T', |\mu(t'_i) \cap \bigcup_{t \in T} \mu(t)| \geq 1, \text{ and} \\ &\quad \bigcup_{t \in T} \mu(t) \subseteq \bigcup_{t' \in T'} \mu(t'), \text{ and there is no set } \{t'_p \dots t'_q\} \\ &\quad \text{such that } \{t'_p \dots t'_q\} \subset T' \text{ and } \bigcup_{t \in T} \mu(t) \subseteq \bigcup_{j=p}^q \mu(t'_j)\}. \end{aligned}$$

I.e. For each template element produced by $\mathbf{generalise}(T, \kappa')$, each term within that element belongs to category κ' ; and each term within that element matches at least one literal matched by a term in T ; and every literal matched by a term or terms in T is matched by at least one term in the template element; and that no subset of terms exists that meets these two requirements. Note that $\bigcup_{t \in T} \mu(t)$ is the set of all literals that are members of terms in the original template element T , and that $\bigcup_{t' \in T'} \mu(t')$ is the set of all literals that are members of terms in the new template elements T' . Each new template element has to be different from the original template element, as they belong to different categories. This ensures that the new element is *more* general than the original, and not just *as* general.

Example: Let category $\kappa_T = \{\text{FELINE}, \text{CANINE}, \text{ANIMAL}\}$ contain three sets of literals:

$$\mu(\text{FELINE}) = \{\text{cat}\},$$

$$\mu(\text{CANINE}) = \{\text{dog}\}, \text{ and}$$

$$\mu(\text{ANIMAL}) = \{\text{cat}, \text{dog}, \text{mouse}, \text{horse}\}.$$

Let $T = \{t_1, \dots, t_n\}$ such that $\bigcup_{i=1}^m \mu(t_i) = \{\text{cat}, \text{dog}\}$. Then $|\mu(\text{FELINE}) \cap \bigcup_{i=1}^m \mu(t_i)| = 1$, therefore the set FELINE can be considered for inclusion in the sets in $\mathbf{generalise}(T, \kappa_T)$. $|\mu(\text{CANINE}) \cap \bigcup_{i=1}^m \mu(t_i)| = 1$, therefore the set CANINE can be considered for inclusion in the sets in $\mathbf{generalise}(T, \kappa_T)$. $|\mu(\text{ANIMAL}) \cap \bigcup_{i=1}^m \mu(t_i)| = 2$, therefore the set ANIMAL can be considered for inclusion in the sets in $\mathbf{generalise}(T, \kappa_T)$. The sets $\{\text{FELINE}\}$ and $\{\text{CANINE}\}$ are individually insufficient to represent all the literal members of $\mathbf{generalise}(T, \kappa_T)$ and so will not be in $\mathbf{generalise}(T, \kappa_T)$. The remaining candidate sets are

$$\{\{\text{ANIMAL}\},$$

$$\{\text{FELINE}, \text{CANINE}\}, \{\text{FELINE}, \text{ANIMAL}\}, \{\text{CANINE}, \text{ANIMAL}\},$$

$$\{\text{FELINE}, \text{CANINE}, \text{ANIMAL}\}\}.$$

Note that $\{\text{ANIMAL}\} \subset \{\text{FELINE}, \text{ANIMAL}\}$, $\{\text{ANIMAL}\} \subset \{\text{CANINE}, \text{ANIMAL}\}$, and $\{\text{ANIMAL}\} \subset \{\text{FELINE}, \text{CANINE}, \text{ANIMAL}\}$. Likewise $\{\text{FELINE}, \text{CANINE}\} \subset \{\text{FELINE}, \text{CANINE}, \text{ANIMAL}\}$. As a result, $\{\text{FELINE}, \text{ANIMAL}\}$, $\{\text{CANINE}, \text{ANIMAL}\}$ and $\{\text{FELINE}, \text{CANINE}, \text{ANIMAL}\}$ are excluded from $\mathbf{generalise}(T, \kappa_T)$, because there are subsets of these sets that meet the conditions of $\mathbf{generalise}$. Therefore $\mathbf{generalise}(T, \kappa_T) = \{\{\text{FELINE}, \text{CANINE}\}, \{\text{ANIMAL}\}\}$. This is

a set of two template elements, the first consisting of two terms and the second consisting of one term.

Note that the cardinality of the sets returned by $\mathbf{generalise}(T, \kappa_\Gamma)$ is not necessarily an indication of their specificity. For example, the set $\{\text{ANIMAL}\}$ matches more literals but represents a single semantic category, while the set $\{\text{FELINE}, \text{CANINE}\}$ matches fewer literals but combines two semantic categories. In an implementation we may decide to use the set $\{\text{FELINE}, \text{CANINE}\}$ which is a more cumbersome, but more specific set of terms than $\{\text{ANIMAL}\}$, depending on our requirements. The decision made regarding which set to use in a template would depend on the application and on the details of the heuristic searches, as we discuss in Sect. 6.

Note also that $\mathbf{generalise}$ will return an empty set if no generalisation exists that matches all the required literals. Thus if the input set contains a literal that is not contained in any member of κ_x , then $\mathbf{generalise}(T, \kappa_x) = \emptyset$.

Example: $\mathbf{generalise}(\{\text{the}\}, \kappa_\Gamma) = \emptyset$, because the literal “the” is not in any gazetteer.

4.2 Creating and Modifying Entire Templates

In the previous section, we defined the creation and generalisation of template elements. Templates are ordered lists of template elements (Definition 13), and we now apply the above concepts to create and modify templates. Given a seed phrase, in the form of a tuple of literals (i.e. a fragment), we can easily define a very specialised template that matches only that fragment. We can then modify this to increase its generality.

Definition 18. *We extend the $\mathbf{initialise}$ function to create a specialised template from a fragment.*

$$\mathbf{initialise}(\langle \lambda_1, \lambda_2, \dots, \lambda_n \rangle) = \langle \{\lambda_1\}, \{\lambda_2\}, \dots, \{\lambda_n\} \rangle.$$

This can also be written as: $\mathbf{initialise}(f) = f$, where f is a text fragment.

We define a new function that generalises any given template to create a new set of templates by modifying a single element of the template using the element generalisation function defined in Sect. 4.1. One template will be created for each possible generalisation of the specified template element.

Definition 19. *Given the template, $\tau = \langle T_1, T_2, \dots, T_i, \dots, T_n \rangle$. Then*

$$\mathbf{generalise}(\tau, \kappa, i) = \{\tau' \mid T'_i \in \mathbf{generalise}(T_i, \kappa) \text{ and} \\ \tau' = \langle T_1, T_2, \dots, T'_i, \dots, T_n \rangle\}$$

i.e. we replace the i^{th} template element with the result of its own generalisation.

Example: Let $\tau_1 = \langle \text{the}, \text{cat}, \text{sat} \rangle$. Then $\mathbf{generalise}(\tau_1, \kappa_\Gamma, 2) = \{\langle \text{the}, \text{FELINE}, \text{sat} \rangle, \langle \text{the}, \text{ANIMAL}, \text{sat} \rangle\}$. In this case, $\mathbf{generalise}(\tau_1, \kappa_\Gamma, 2)$ returns two templates because the second literal “cat” belongs to two gazetteers. In contrast, $\mathbf{generalise}(\tau_1, \kappa_\Gamma, 1) = \emptyset$, because the first literal “the” does not belong to any gazetteer in the category κ_Γ .

4.3 Matching Fragments of Differing Lengths

So far, we have assumed that all templates are generated from a seed phrase by replacing literals in that phrase with other attributes. This restricts the templates created to be the same length as the seed phrase, and so will only match fragments of this fixed length. This is clearly undesirable, so we will now extend the framework to allow for templates that match fragments of various lengths, firstly to shorter fragments and then to longer ones.

Suppose we want to be able to define a template that can match a particular class of noun phrase, in the form of a determiner, zero or more words (e.g. adjectives), and an animal. To do this, we use the wildcard attribute '?' which matches any single word *or no word at all*. This can be used to define templates such as $\tau = \langle \text{DT}, ?, ?, \text{ANIMAL} \rangle$.

To use the '?' wildcard in our examples, we assume that '?' is in the wildcard category κ_Ω and use the existing $\text{generalise}(T, \kappa_\Omega)$ function. So for example, we take a fragment and apply a series of **generalise** functions after initialisation like this:

$\tau_1 = \langle \text{the, quick, brown, fox} \rangle$	$= \text{initialise}(\text{the, quick, brown, fox})$
$\tau_2 = \langle \text{the, ?, brown, fox} \rangle$	$\in \text{generalise}(\tau_1, \kappa_\Omega, 2)$
$\tau_3 = \langle \text{the, ?, ?, fox} \rangle$	$\in \text{generalise}(\tau_2, \kappa_\Omega, 3)$
$\tau_4 = \langle \text{DT, ?, ?, fox} \rangle$	$\in \text{generalise}(\tau_3, \kappa_\Omega, 1)$
$\tau_5 = \langle \text{DT, ?, ?, ANIMAL} \rangle$	$\in \text{generalise}(\tau_4, \kappa_\Omega, 4)$

Consider a document containing these fragments:

$f_1 = \langle \text{the, cat} \rangle$
$f_2 = \langle \text{the, lazy, dog} \rangle$
$f_3 = \langle \text{the, quick, brown, fox} \rangle$

Of the above list of templates, τ_5 matches all three fragments, whereas $\tau' = \langle \text{DT}, *, *, \text{ANIMAL} \rangle$ would not match f_1 or f_2 .

Although this approach allows templates to match fragments of varying length as desired, it is restricted to matching fragments that are no longer than the original seed. The seed phrase must therefore be chosen with this in mind, or else modified before the search begins by inserting wildcards. One solution is to extend the **generalise** function so that it inserts a wildcard into the template, making the template longer and allowing it to match phrases longer than the seed phrase. For example, we could define $\text{generalise}(\tau, \kappa_{\Omega'}, i)$ as a function that creates new templates from template τ by inserting wildcards from $\kappa_{\Omega'}$ between template elements i and $i + 1$. The category $\kappa_{\Omega'}$ contains the same wildcards as κ_Ω , with the distinction that they are inserted *after* existing template elements, rather than used to replace template elements. Introducing extra wildcards in this way complicates the search algorithms substantially, so restrictions would have to be introduced, either limiting how many wildcards could be added to a template, or limiting the maximum length of a template. These heuristics could be chosen for a particular implementation after experimentation.

Note that there can be no gain from inserting ‘?’ wildcards before the first non-wildcard template element, or after the last non-wildcard element, as the resulting template would match exactly the same fragments with or without these optional ‘external’ wildcards.

5 Measuring Template Quality

In this section, we define recall, precision and related terms formally within our framework. We then discuss how these can be estimated in practice. These measures are needed to guide the automatic search for templates discussed in the next section.

5.1 Defining Recall and Precision

We have already defined a function $\mu(\tau, D)$ that retrieves the set of document fragments that are matched by template τ from document set D (Definition 14).

Definition 20. We define $I(D)$ as the set of interesting, relevant fragments contained in corpus D .

The ideal template would match these, and only these, fragments. This set is generally not known⁵, but allows us to define concepts such as “true positive” and the precision rate. Figure 1 shows the relationships diagrammatically. We use the notation $\wp_f(D)$ to show the set of all fragments in D , independent of any template. This is the set of all tuples of all lengths obtainable as fragments using $f(d, a, b)$ (Definition 12).

We now define a series of sets with respect to a template τ and a corpus D .

Definition 21. True-positives $TP(\tau, D) = \mu(\tau, D) \cap I(D)$

Definition 22. False-positives $FP(\tau, D) = \mu(\tau, D) \setminus I(D)$

Definition 23. False-negatives $FN(\tau, D) = I(D) \setminus \mu(\tau, D)$

Definition 24. True-negatives $TN(\tau, D) = D \setminus \{I(D) \cup \mu(\tau, D)\}$

Definition 25. Recall $r(\tau, D) = \frac{|\mu(\tau, D) \cap I(D)|}{|I(D)|} = \frac{|TP|}{|I(D)|}$

Definition 26. Precision $p(\tau, D) = \frac{|\mu(\tau, D) \cap I(D)|}{|\mu(\tau, D)|} = \frac{|TP|}{|\mu(\tau, D)|}$

We can conceive of a “perfect” template, τ^* , which matches all the positive fragments, and nothing else. I.e. $\mu(\tau^*, D) = I(D)$. This has a recall and a precision of one. Although such a perfect template may not be known in general, we can use these definitions of true positive and false positive counts to guide a search for templates.

⁵This set is task-dependent and user-dependent, as different people will pick different fragments as being interesting, even if attempting the same task [19].

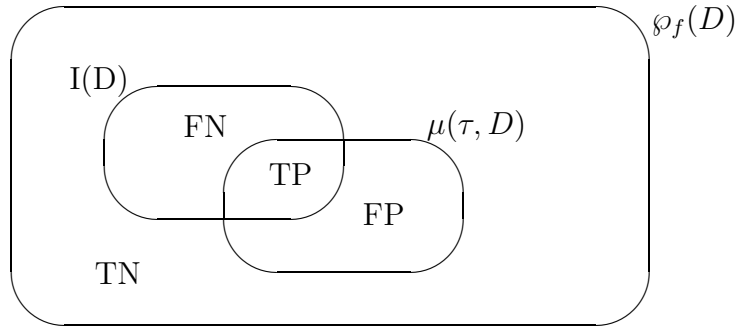


Fig. 1. Venn diagram showing the relationship between: the set of all possible fragments $\varphi_f(D)$; the results of applying template τ to it — $\mu(\tau, D)$; the ideal set of results $I(D)$; and the false negative (FN), true positive (TP), false positive (FP) and true negative (TN) regions.

An ideal template would have $|TP(\tau, D)| = |I(D)|$ and $|FP(\tau, D)| = 0$. We therefore wish to maximise $|TP(\tau, D)|$ while minimising $|FP(\tau, D)|$, but there is typically a trade off between the two. This is an example of multi-objective optimisation. If we knew the relative value of true positives and the cost of false positives, then we could combine these into a single objective function, such as maximising $|TP(\tau, D)| - k \cdot |FP(\tau, D)|$. In practice, such a weighting is not usually available, but a number of evolutionary approaches have been successfully applied to similar problems [20], as we discuss further in Sect. 7.

6 Searching for Good Templates

In 6.2 we will discuss the development of search algorithms, but first we need a practical way to estimate recall and precision.

6.1 Estimating Recall and Precision

As noted earlier, we do not know which fragments are interesting *a priori*, and so the above definitions cannot be used directly in calculating the recall or precision of a template. Instead, we need something that we can measure in practice, and which should approximate the “ideal” values above. We now consider several options.

Suppose that we had a set of documents such that every fragment was labelled as either “interesting” or “not interesting”. Then we could use standard supervised learning algorithms to construct useful templates and directly measure the number of true positives, false positives and so on, to find an optimal template. This could then be used to find further information in the same field. However, while a small number of such

labelled corpora do exist (e.g. [21]), they are for a few very precisely defined application areas, and so of little general use, as they cannot be used to aid IE in other application areas. Annotating documents in this way is very time consuming for a domain expert, and one aim of information extraction is to reduce the time and effort required to find relevant information.

Suppose instead that we have one set of documents where for each sentence, the probability that it contains a relevant piece of information is above some threshold, and a second set of documents, where the probability is below a threshold. We could then treat this as a classification problem, albeit with noisy labels on the data. However, such a set of irrelevant documents is hard to define, and furthermore, even interesting documents are likely to contain irrelevant facts such as background information, although this approach has been used successfully [22].

Suppose that instead of having irrelevant (negative) documents, we have a set of “neutral” documents, each of which may or may not contain relevant information. I.e. we have no prior knowledge about relevant information in neutral documents. We can then compare the proportion of information retrieved from neutral and from positive documents to evaluate a template. We assume that a “good” template will retrieve more information from positive documents than from neutral documents, even if we don’t know in advance which pieces of information are useful, or how much useful information exists in any particular document.

Let D be a corpus containing $|D|$ documents. We define the set of positive documents as D^+ and neutral documents as D^N , such that $D = D^+ \cup D^N$ and $D^+ \cap D^N = \emptyset$. A “positive” document is one that the user believes is likely to contain information of interest. A “neutral” document is one where the user has no reason to believe that the document does or does not contain information of interest.

We now use these two sets of documents to define estimates of the numbers of true-positive fragments and false-positive fragments matched by a template τ .

Definition 27. *We define an estimated true-positive set $\widehat{TP}(\tau, D) = \mu(\tau, D^+)$, for a template τ and a set of positive documents, $D^+ \subseteq D$.*

Definition 28. *We define an estimated false-positive set $\widehat{FP}(\tau, D) = \mu(\tau, D^N)$, for a template τ and a set of neutral documents, $D^N \subseteq D$.*

These are very crude estimates, as they assume that every fragment matched by τ in D^+ contains information of interest, and that every fragment matched by τ in D^N contains *no* information of interest. Thus they cannot be used to estimate the precision or recall scores⁶, but they are sufficient to guide the search of useful templates. In fact, as a successful search progresses and the quality of the template improves, then these estimates will become more accurate, although they are unreliable, especially at the start of the search process.

⁶Substitution into definitions 25 and 26 gives recall \equiv precision \equiv 1 for every template, which is clearly optimistic.

Let us consider some other properties of templates generated using superset generalisation. Suppose we have two templates, τ_1 and τ_2 , and that $\tau_1 >_s \tau_2$. Then by definition, $|\mu(\tau_1, D)| < |\mu(\tau_2, D)|$. If τ_2 is a generalisation of τ_1 derived using superset generalisation, then $\mu(\tau_1, D) \subset \mu(\tau_2, D)$. This relative specificity relation holds for any set of documents, so if one template matches fewer fragments than another in one corpus, then it will in any other corpus as well. Thus given two corpora D^a and D^b :

$$|\mu(\tau_1, D^a)| < |\mu(\tau_2, D^a)| \iff |\mu(\tau_1, D^b)| < |\mu(\tau_2, D^b)|.$$

This property will be useful in developing heuristic search methods, because it allows us to rationally prune search graphs, as we discuss in Sect. 6.2.

We defined terms such as “true positive” and “false positive” above. Now we can say that if template τ_1 is at least as specific as template τ_2 , then the number of true positives returned by τ_1 is no more than the number returned by τ_2 , and equivalently for other scores. I.e. if $\tau_1 \geq_s \tau_2$ then:

$$\begin{aligned} |\text{TP}(\tau_1, D)| &\leq |\text{TP}(\tau_2, D)|. \\ |\text{FP}(\tau_1, D)| &\leq |\text{FP}(\tau_2, D)|. \\ |\text{TN}(\tau_1, D)| &\geq |\text{TN}(\tau_2, D)|. \\ |\text{FN}(\tau_1, D)| &\geq |\text{FN}(\tau_2, D)|. \end{aligned}$$

The equivalent inequalities also hold for the estimates defined above:

$$\begin{aligned} |\widehat{\text{TP}}(\tau_1, D)| &\leq |\widehat{\text{TP}}(\tau_2, D)|. \\ |\widehat{\text{FP}}(\tau_1, D)| &\leq |\widehat{\text{FP}}(\tau_2, D)|. \end{aligned}$$

As these relationships hold for any set of documents D , we can predict some properties of templates without fully evaluating them. We can use these properties to guide heuristic searches.

If our assumptions here are correct, then the probability of finding an interesting fragment is higher in positive documents than in neutral documents. We can write this assumption as $p(f \in I(D) | f \in \mu(\tau, D^+)) > p(f \in I(D) | f \in \mu(\tau, D^N))$.

When comparing two templates, we can say that one is certainly better than the other if either: a) it matches more true positive fragments and no more false positives; or b) it matches fewer false positive fragments and no fewer true positives. Unfortunately, we cannot be certain of whether it is better when c) it matches more true positives *and* more false positives. We need some way of exploring this trade-off, an issue we return to later. We can now generate a series of templates of varying generality and compare them with each other in order to guide our search for useful templates. We discuss this further in the following section.

6.2 Search Algorithms

As stated in the introduction, heuristic searching requires: definitions of candidate solutions; a means to generate and evaluate candidate solutions; and a suitable search algorithm. We have now defined the candidate solutions (i.e. templates, Definition 13) and means to generate (Sect.

4) and evaluate them (using estimates of true positive and false positive scores given in Sect. 5, Definitions 27–28). We now turn to the search algorithms themselves. First, we define an exhaustive search over all possible templates, given a particular seed phrase. We will then show that in general, this is computationally infeasible, owing to the combinatorial growth of the search space with respect to the length of the template. We then introduce heuristics to make the search feasible while still (we hope) producing good templates, and define and discuss a simple best-first search algorithm. We also discuss possible merits of evolutionary algorithms.

In all cases, we assume that we are given a seed fragment f . The root node of the search corresponds to a template consisting of a tuple of template elements, each containing a single literal: $\tau_{root} = \text{initialise}(f)$ (Definition 18). From this, we can modify each element in the template by a single application of the $\text{generalise}(\tau, \kappa, x)$ function (Definition 19). We can estimate the number of true positives and false positives for each of these new templates, and then decide which template to explore and modify next. The exact number of templates produced at each stage depends on the words themselves, because each $\text{generalise}(\tau, \kappa, x)$ function will return 0, 1 or more templates (see Definition 17 and the accompanying discussion). We must also decide when to terminate the search, as we do not know *a priori* the quality of the best possible template, as we are assuming that we do not have a fully labelled corpus.

A simple exhaustive search method would be to start with a literal template created from the seed phrase using the **initialise** function. For each element in this template, we then apply the generalisation function, using each category in turn, so that each application generates a new template. We need some fixed order over the categories, but the actual order is not critical here⁷. For simplicity, let us assume that we have an implementation with five categories, where every word has one literal, one stem, one gazetteer, one part of speech, and one wildcard. We can then list all possible templates derived from a two-term phrase, starting with two sets of literals, and ending with two sets of wildcards:

$\langle \lambda, \lambda \rangle \langle \sigma, \lambda \rangle \langle \gamma, \lambda \rangle \langle \pi, \lambda \rangle \langle \omega, \lambda \rangle$
 $\langle \lambda, \sigma \rangle \langle \sigma, \sigma \rangle \langle \gamma, \sigma \rangle \langle \pi, \sigma \rangle \langle \omega, \sigma \rangle$
 $\langle \lambda, \gamma \rangle \langle \sigma, \gamma \rangle \langle \gamma, \gamma \rangle \langle \pi, \gamma \rangle \langle \omega, \gamma \rangle$
 $\langle \lambda, \pi \rangle \langle \sigma, \pi \rangle \langle \gamma, \pi \rangle \langle \pi, \pi \rangle \langle \omega, \pi \rangle$
 $\langle \lambda, \omega \rangle \langle \sigma, \omega \rangle \langle \gamma, \omega \rangle \langle \pi, \omega \rangle \langle \omega, \omega \rangle$

If every literal has exactly $|\alpha|$ attributes, then a seed phrase of $|f|$ literals has $|\alpha|^{|f|}$ possible templates. Thus a 20-word seed phrase consisting of literals having 5 attributes will be matched by $5^{20} \approx 10^{14}$ templates. In practice, some words will have fewer attributes (e.g. words that don't appear in any gazetteers), some will have more, and all may have more than one wildcard. We therefore need to introduce heuristics to make the search feasible, unless the seed fragment is very short.

⁷One example would be to use the order: literals, stems, gazetteers, parts-of-speech and wildcards. This order reflects the typical size of the membership functions. E.g. each stem has only a few members; gazetteers often contain hundreds or thousands; parts-of-speech such as verb or noun contain millions of members and so on.

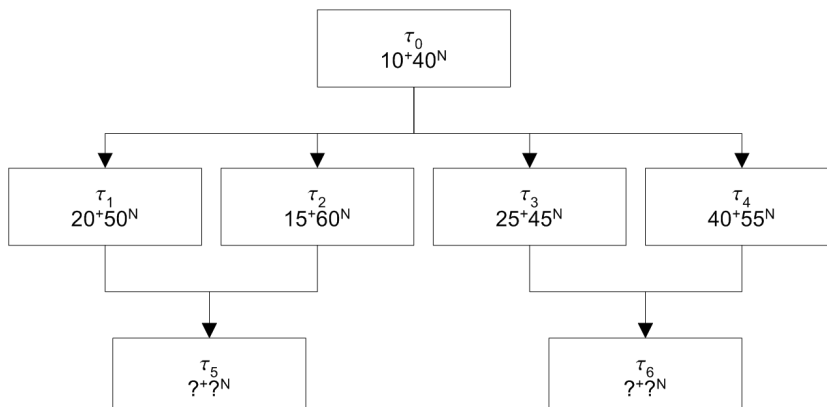


Fig. 2. Part of a search tree. Each node is a template with the number of true positives (x^+) and false positives (y^N) shown, with a “?” for unknown values. Each child node can be created by superset generalisation from any of its parents.

6.3 Feeding Knowledge Forward

After estimating the numbers of true positives and false positives for any template, we can place a lower bound on those values for all templates that are derived using the **generalise** function. This is because we know that the derived templates will match a superset of the fragments matched by the ancestor template (Sect. 6.1).

For example, suppose we have a template τ_1 and we evaluate this and find it matches 20 positive fragments and 50 neutral fragments, i.e. $|\mu(\tau_1, D^+)| = 20$ and $|\mu(\tau_1, D^N)| = 50$. If we then modify τ_1 using superset generalisation to create τ_2 , then we know that $\mu(\tau_1, D) \subseteq \mu(\tau_2, D)$ for any corpus D . We therefore know that τ_2 matches *at least* 20 positive fragments and *at least* 50 neutral fragments, from D^+ and D^N respectively. Furthermore, any other templates derived using superset generalisation from τ_1 or from τ_2 will also match at least those numbers. Therefore, as we carry out a search, we can calculate lower bounds on the estimates of true positives and false positives for a wide range of templates without the computational expense of evaluating each template. Instead, we can just choose the best template from the range available, if we use a best-first search. By “best” here, we might mean the template with the (estimated) most true positives, which will tend to produce templates with a high recall, though possibly with a low precision. If instead we choose the template that matches the least false-positive fragments then we will tend to produce templates with a high precision, though possibly with a low recall. Which of these options is more appropriate depends on the task at hand.

Every template has at least one parent (except for the single root template), because they are created using superset generalisation; but most templates can be created in more than one way, from several parent templates. Therefore, many templates will have more than one parent. The

whole space may be considered as a lattice. Consider the partial search graph shown in Fig. 2. Here, templates $\tau_0 - \tau_4$ have been evaluated, and the number of true positives and false positives are shown for each. For example, τ_1 matches 20 fragments from D^+ and 50 from D^N , and is therefore assumed to match 20 true positives and 50 false positives. At this stage of the search, the decision to be made is which node to evaluate next: τ_5 or τ_6 ? We can feed forward the facts that the two parent templates of τ_5 , τ_1 and τ_2 have 20 and 15 true positives respectively, and that therefore τ_5 must have at least 20 true positives. Similarly, it must have at least 60 false positives. On the other hand, τ_6 must have at least 40 true positives, and at least 55 false positives. We would therefore decide to evaluate τ_6 in preference to τ_5 at this point, because it has a higher lower-bound on the number of true positives, and a lower lower-bound on the number of false positives.

6.4 A Best-First Algorithm

We now present a formal definition of a best-first algorithm suitable for identifying good templates, followed by an example. We start by defining two sets of templates. Set O (“open”) contains templates that have been created (via **initialise** or **generalise**), but not yet evaluated. Set C (“closed”) contains templates that have been evaluated, i.e. had values of TP and FP calculated. Note that $O \cap C \equiv \emptyset$.

Figure 3 gives the algorithm. After initialisation, we take the best template that has not yet been evaluated, and evaluate it, i.e. calculate its true positive and false positive scores. We then generalise it in every way possible to create child templates, and update the lower bounds on the true and false positive scores for these children. Because there are several ways that each template could be created, some children will have more than one parent. In these cases, the lower bounds are the *maximum* of the lower bounds of all the parents.

1. **begin**
2. $C \leftarrow \emptyset, O \leftarrow \emptyset$
3. **initialise**(f) = $\tau_{root} \rightarrow O$
4. **while not finished do**
 - (a) find estimated best template τ in O
 - (b) evaluate τ
 - (c) delete τ from O
 - (d) add τ to C
 - (e) expand τ by adding to O all templates that can be created by superset generalisation of τ
 - (f) update lower bounds on TP and FP for all templates in O
5. return best template from C .
6. **end**

Fig. 3. A best-first algorithm. C is the “closed” set of evaluated templates and O is the “open” set of unevaluated templates. See Sect. 6.4 for further details.

By “best template” (Fig. 3, Steps 4a and 5), we mean select the template with the highest lower-bound on the number of true positives, as inherited from each template’s ancestors. If more than one template has the same maximum true positive lower-bound, then we choose between them by selecting the template with the smallest lower-bound on false positives. If this still selects more than one template, we can either pick one randomly; use them all successively; or use all the selected templates together in the subsequent steps (i.e. evaluate and generalise more than one template in one pass through the main loop).

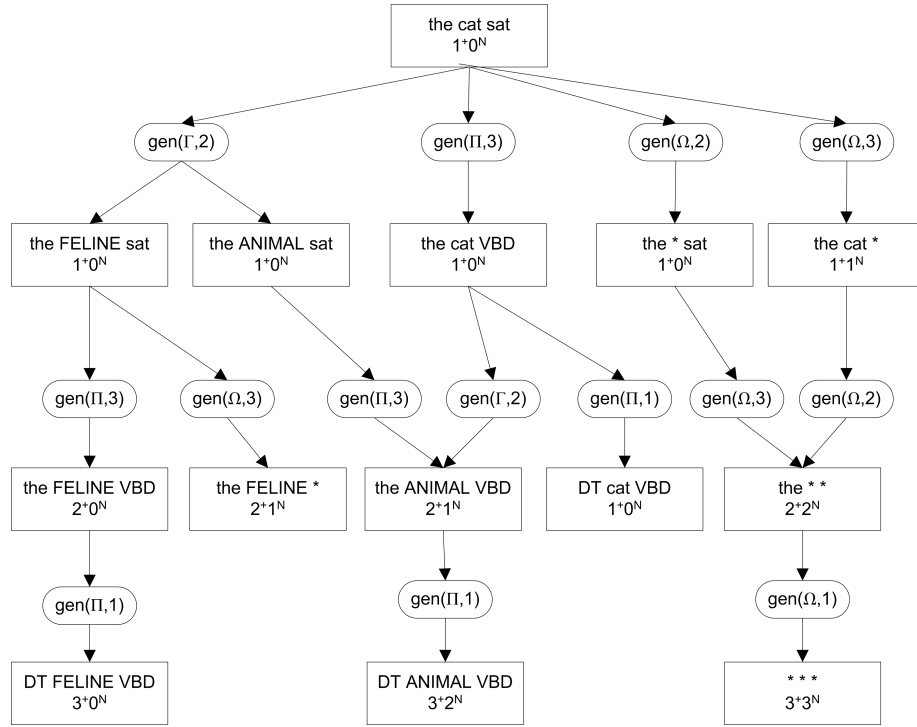
In Step 4f we could choose to either update the bounds of just the descendants of τ in O , or else update the bounds of the descendants of every template in C . The latter would be more computationally expensive, but would lead to better estimates of the lower bounds. For each new template created, it would require a search through C to find all the other possible ancestors, besides τ , in order to calculate the new lower bounds.

Finally, we terminate the search when some pre-specified criterion is satisfied. Possible criteria include terminating: when O is empty (in which case every possible template has been evaluated); or when a limit on the number of evaluations is reached (e.g. stop after 1000 templates have been evaluated); or when a certain number of true positives (or false positives) are matched by the current best template. It would be quite possible for the user to stop the search and consider the current best template, before re-starting the search if necessary.

As a more concrete example, Fig. 4 shows part of a search graph containing various templates created from the seed fragment “the cat sat”, and evaluated with respect to the two small corpora shown.

Consider the right-hand portion of the graph. Near the top are two templates: $\langle \text{the}, *, \text{sat} \rangle$ and $\langle \text{the}, \text{cat}, * \rangle$, both created using the $\text{generalise}(\tau, \kappa_{\Omega}, x)$ function. The first matches one true positive and no false positives (shown as $1+0^N$ in the figure). The second matches one true positive and one false positive. By the definition of superset generalisation, we know that every template derived from this second template will have at least one true positive and one false positive. So if, at one stage of a search, we only want to consider templates with no false positives, then we need not consider *any* descendent of this template. The two descendants shown ($\langle \text{the}, *, * \rangle$ and $\langle *, *, * \rangle$) need not be evaluated explicitly therefore; they can be annotated as having at least one false positive, although in the figure, the fully evaluated scores are shown.

Now consider the third row of templates in Fig. 4, i.e. those created after two generalisation functions. From the left of the graph, three of these have two true positives ($\langle \text{the}, \text{FELINE}, \text{VBD} \rangle$, $\langle \text{the}, \text{FELINE}, * \rangle$ and $\langle \text{the}, \text{ANIMAL}, \text{VBD} \rangle$) and one has only one true positive, so we focus the search on the first three. These have zero, one and one false positives respectively, making the first one look most promising: $\langle \text{the}, \text{FELINE}, \text{VBD} \rangle$. This is the best unevaluated template so far. Several further generalisations are possible from this template, including applying $\text{generalise}(\tau, \kappa_{\Omega}, 3)$ to form $\langle \text{the}, \text{FELINE}, * \rangle$. But this has already been evaluated, so need not be considered again. Applying



Legend:

D^+	D^N	Template definition TP*FP ^N	generalise(category, position)
{the cat sat}	{the mouse ran}		
{the lion roared}	{a cow jumped}		
{a tiger slept}	{the cat plays}		

Fig. 4. Part of a search graph for a three-term template. Each rectangle shows a template, starting with three sets of literals at the top initialised from the fragment “the cat sat”. Various generalisation functions are then applied to it. For example, $\text{gen}(\Gamma, 2)$ means apply the $\text{generalise}(\tau, \kappa_\Gamma, 2)$ function to the upper template, τ , to produce the lower template(s). The graph includes part-of-speech labels “VBD” meaning past-tense verb and “DT” meaning determiner. The numbers in each box below the template represent the estimated numbers of true-positive and false-positive matches for each template, with respect to the positive and neutral document sets D^+ and D^N shown. Note that not every node or edge is shown.

generalise($\tau, \kappa_H, 1$) forms <DT, FELINE, VBD> which has three true positives, and still no false positives. Given the two very small document sets, this is an optimal template, so we stop the search here. In a more realistic application, the search would continue until a stopping criterion was reached, but would not find any superior template.

One modification to the algorithm would require more memory but should lead to a faster convergence to a (possibly) superior solution. This is to expand the unevaluated template set O by repeated applications of the **generalise** function until it contains every template that could possibly be derived from the seed phrase, or as many as is practically possible. We would then proceed with a best-first search, but with the advantage that after each evaluation, a large number of unevaluated templates will have the lower bounds of their true- and false-positive estimates updated, because their ancestry would be explicitly represented on the search graph. This should improve the selection of the best template at each stage.

The algorithm above does not prune any part of the search graph, but merely tends to search promising areas first. In practice, this is likely to lead to very large memory requirements, which can be avoided through pruning. Pruning search *graphs* is a lot more difficult than pruning *trees* because each node can have multiple parents, and so after a node is pruned, it may reappear as part of a different path. Although algorithms such as A* are inappropriate here⁸, recent variations may provide useful pruning methods, such as SMA* [14, p. 104] and Sweep A* [23].

6.5 Sample Results

The algorithm outlined in Fig. 3 has been implemented as an extension to the BioRAT tool [6]. As a proof of principle of the methods described here, we used the implementation to derive a template which could be used to identify protein-protein interactions from biological literature.

We start with the seed phrase “Rad53p protein binds to Dbf4p”, where “Rad53p” and “Dbf4p” are proteins, and “binds to” suggests a direct physical interaction. We create a positive corpus by extracting 500 abstracts that are listed in the Database of Interacting Proteins (DIP [24]). We assume that most of these abstracts mention at least one protein-protein interaction, although the information will be expressed in many different ways. We also create a neutral corpus by selecting the first 500 abstracts dated September 2005 retrieved from PubMed. This is essentially a random set of biomedical texts. It may contain some references to particular protein-protein interactions, but we assume that most of these abstracts will not.

To select the “best template” during the search (Fig. 3, Steps 4a and 5), we define a quality score thus:

score = (number of positive matches) – $w \times$ (number of neutral matches)

The weight w is used to control the trade-off between true- and false-positives and hence balance the search for high recall and high precision

⁸We have no notion of a path cost, and are only interested in the final solution rather than the path to it.

templates. (As explained in Sect. 6.1, we are assuming that each match in the neutral corpus is a false-positive result.) We start with a small value of $w = 0.5$ and slowly increase this as the search progresses. This encourages early exploration of the search space while later penalising false positive matches heavily. We select the unevaluated template with the highest score at each step. In cases where two or more templates have the same score, candidate templates are selected if they contain more gazetteer elements and fewer literal elements than the others. This deliberately introduces a slight bias favouring templates that contain references to gazetteers, as these represent useful domain knowledge. Table 1 summarises the search as it progresses for the first 50 iterations.

Table 1. Each row shows the template evaluated at one iteration, starting with the literal template derived from the given seed phrase at iteration zero. The number of matches in the positive (TP) and neutral (FP) corpora are shown for the current template (“measured matches”) and for its parent (“inherited matches”). Each template element is shown delimited by square brackets in the form [CATEGORY: value], with categories literal (A), gazetteer (T) and wildcard (Ω). A wildcard element with n query symbols matches between 0 and n words. The sequence of templates generated tends to generalise from literals through gazetteers to wildcards.

Iteration	Measured matches		Inherited matches		Template pattern
	TP	FP	TP	FP	
0	1	0	-	-	[A: Rad53p] [A: protein] [A: binds] [A: to] [A: Dbf4p]
10	2	0	1	0	[T: protein, sp] [A: protein] [T: binding] [Ω : *] [T: protein, sp]
20	86	1	25	1	[T: protein] [Ω : ?] [T: binding] [Ω : *] [T: protein, sp]
30	286	4	247	4	[T: protein] [Ω : ?????] [T: binding] [Ω : ???] [T: protein, sp]
40	1202	34	1202	34	[Ω : ??] [Ω : ?????] [T: binding] [Ω : ?????] [T: protein]
50	2621	151	2621	151	[Ω : ???] [Ω : ?????] [T: binding] [Ω : ?????] [Ω : *]

On looking at Table 1, we see that the template evaluated on the 30th iteration matches 286 fragments from the positive corpus, and 4 from the neutral corpus. The template pattern is: [T: protein] [Ω : ?????] [T: binding] [Ω : ???] [T: protein, sp]

This matches a protein followed by between 0 and 5 words, followed by a protein-binding term, followed by between 0 and 3 other words, followed by another protein (of sub-type “sp”; this refers to terms in a gazetteer derived from the SwissProt database). The five sentences listed below were among those found in the positive corpus. The *italicised* portions show the fragments matched by the template; the rest is given for context only.

- *Protein kinase C delta associates with and phosphorylates Stat3* in an interleukin-6-dependent manner.
- Furthermore, *Stat3 was phosphorylated by PKC delta* in vivo on Ser-727...

- We report here that *Grb2 also interacts with tyrosine-phosphorylated IRS-1* in response to gastrin.
- ...RII alpha fused to *endonexin II formed dimers but did not bind MAP2*.
- *A protein interaction map* for cell polarity development.

While the first three results correctly show protein-protein interactions, the last two highlight imperfections in this particular template. The phrase “did not bind” indicates that no interaction was observed, but the template still treats this as a positive statement. In the last example shown, the word “map” is not a protein, although the same word is used to refer to a protein in other contexts and so it is in the protein gazetteer. Table 1 also shows that the template evaluated in the 50th iteration matches 2621 fragments in the positive corpus, and 151 in the neutral, with the template pattern: [Ω : ???] [Ω : ?????] [I : binding] [Ω : ?????] [Ω : *]

This might be seen as an example of over-generalisation, in that the template consists almost entirely of wildcards. However, the quality of any template depends on the requirements of the user, and for some tasks, even this very general template may still be useful. Note that the total runtime on a stand-alone PC was 20 minutes for these 50 iterations.

7 Discussion and Extensions

We now briefly discuss a few of the possible extensions to the framework and its implementation.

We could introduce other wildcards, such as a wildcard which matches an entire phrase, which could itself be defined as a series of terms, much like a template. This would allow optional subclauses, such as subordinate clauses, to be matched. Let $?^{\tau_1}$ designate an optional wildcard that matches a sequence of literals defined by template τ_1 or matches nothing at all. Then if $\tau_1 = \langle \{\text{which}\}, \{*\}, \{*\} \rangle$, the template $\tau_2 = \langle \{\text{DT}\}, \{\text{ANIMAL}\}, \{?^{\tau_1}\}, \{\text{sat}\} \rangle$ would match the fragment $\langle \text{the, cat, which, was, black, sat} \rangle$ as well as $\langle \text{the, cat, sat} \rangle$.

An additional approach to finding good templates is to repeatedly *merge* useful templates to produce more general templates [25]. Our framework could easily be extended to allow this, by ensuring that the product of merging two templates matches every fragment that either template matches. This could be achieved by considering each pair of template elements in turn, and either performing a simple set union if they belong to the same category, or else generalising them both to the same category before such a union. In either case, the new template would match the union of the true-positives matched by the two parents, and the union of the false-positives, allowing the lower-bounds on each to be calculated.

So far, we have considered template that exist in isolation, whereas in practical systems, it is more common to apply a set of templates together. Our framework can be extended to include this by using a sequential covering algorithm. Suppose we have a template τ that matches some true positives and some false positives. We could reduce the number of false positives by creating a second template τ' that is optimised to match

just the false positive fragments matched by τ . This could be achieved by defining two new versions of D^+ and D^N based on the fragments matched by τ , and using these to guide the search for τ' . We could then apply τ and τ' together, predicting interesting fragments as $\mu(\tau, D) \setminus \mu(\tau', D)$ (i.e. fragments matched by τ but not by τ'). In many practical applications, more than one template will be applied to a set of documents, each designed to match a different piece of information, or a different way of expressing that information.

We have assumed that we do not have a set of annotated examples, i.e. fragments known in advance to be positive or negative. Creating and annotating large sets of examples is extremely time consuming for a user, although giving a yes/no response to automatic annotations is simpler [26]. One enhancement to our system therefore would be to start with the estimates of true positive and false positive as outlined above, and search for a good template, and then use this template to annotate a number of fragments and to present these to the user. The user then marks each fragment as interesting or not interesting, and this could then be used to improve the quality of the function used to estimate the numbers of true and false positives. This improved function could be used to guide a new template search.

The best-first search described above may miss out on good templates because of its greedy decision making. This would be true even if the estimates of the numbers of true and false positives were perfect, owing to the structure of the graph: we can't guarantee that the best parents will have the best children. One likely improvement therefore would be a population based search, such as a simple beam search or an evolutionary algorithm. Evolutionary algorithms have been successfully used to solve a wide range of multi-objective optimisation problems [20], including problems where evaluations must be limited due to time or financial constraints [27]. Extracting information from a large corpus takes considerable computing effort, so this is an aspect worth considering. Multi-objective evolutionary algorithms can efficiently generate a range of Pareto-optimal solutions, and so explore the trade-off between the different objectives. In our case, this means that for each number of true positives, we find the template with the fewest false positives, and for each number of false positives, we find the template with the most true positives. This produces a range of solutions from which the user can then select whichever template or templates are most sui for their particular problem. This is more flexible than the simple weighting suggested in Sect. 6.5.

Finally, rather than starting with a seed fragment and a template consisting solely of literals, we could start the search using a hand-written template. This would *not* have to be optimised in advance, and in some cases, would be easy to create. The search could then start from a point chosen to be useful and optimised further through similar search processes to those outlined above.

8 Conclusion

We have presented a formal framework to describe information extraction, focusing on the definition of the template patterns used to convert free text into a structured database. The framework has allowed us to explicitly identify some of the fundamental issues underlying information extraction and to formulate possible solutions. We have shown that the framework allows computationally feasible heuristic search methods to be developed for automatic template creation. We have shown that a practical implementation of this framework is feasible and allows automatic template creation. We also hope that the framework will allow other researchers to gain further insights into the theory and practice of information extraction and text mining.

Acknowledgements This work is partly funded by the BBSRC grant BB/C507253/1, “Biological Information Extraction for Genome and Superfamily Annotation.”

References

1. Cowie, J., Lehnert, W.: Information extraction. *Communications of the ACM* **39**(1) (1996) 80–91
2. Cowie, J., Wilks, Y.: Information extraction. In Dale, R., Moisl, H., Somers, H., eds.: *Handbook of Natural Language Processing*. Marcel Dekker, New York (2000)
3. Costantino, M.: Financial Information Extraction using pre-defined and user-definable templates in the LOLITA System. PhD thesis, University of Durham, Department of Computer Science (1997)
4. Blaschke, C., Valencia, A.: The frame-based module of the SUISEKI information extraction system. *IEEE Intelligent Systems* **17**(2) (2002) 14–20
5. Hirschman, L., Yeh, A., Blaschke, C., Valencia, A.: Overview of BioCreAtIvE: critical assessment of information extraction for biology. *BMC Bioinformatics* **6**(Suppl 1) (2005)
6. Corney, D.P.A., Buxton, B.F., Langdon, W.B., Jones, D.T.: BioRAT: Extracting biological information from full-length papers. *Bioinformatics* **20**(17) (2004) 3206–13
7. Hersh, W., Bhuptiraju, R., Ross, L., Cohen, A., Kraemer, D.: TREC 2004 genomics track overview. In: *The Thirteenth Text Retrieval Conference (TREC 2004)* NIST Special Publication SP 500-261. (2004)
8. Sehgal, A.: Text mining: The search for novelty in text. Ph.D comprehensive examination report, Dept. of Computer Science, The University of Iowa (2004)
9. Cunningham, H.: GATE, a General Architecture for Text Engineering. *Computers and the Humanities* **36**(2) (2002) 223–254
10. Collier, R.: Automatic template creation for information extraction. PhD thesis, Department of Computer Science, University of Sheffield (1998)

11. Huang, M., Zhu, X., Hao, Y., Payan, D.G., Qu, K., Li, M.: Discovering patterns to extract protein–protein interactions from full texts. *Bioinformatics* **20**(18) (2005) 3604–3612
12. Pierce, D., Cardie, C.: Limitations of co-training for natural language learning from large datasets. In: Proceedings of the 2001 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics Research (2001)
13. Riloff, E.: Automatically constructing a dictionary for information extraction tasks. In: National Conference on Artificial Intelligence. (1993) 811–816
14. Russell, S., Norvig, P.: *Artificial Intelligence: A Modern Approach*. 2 edn. Prentice Hall (2003)
15. Koike, A., Niwa, Y., Takagi, T.: Automatic extraction of gene/protein biological functions from biomedical text. *Bioinformatics* **21**(7) (2005) 1227–1236
16. Corney, D., Byrne, E., Buxton, B., Jones, D.: A logical framework for template creation and information extraction. In: Foundations of Semantic Oriented Data and Web Mining workshop, part of ICDM2005 (the Fifth IEEE International Conference on Data Mining). (2005)
17. Porter, M.F.: An algorithm for suffix stripping. *Program* **14**(3) (1980) 130–137
18. Marcus, M., Santorini, B., Marcinkiewicz, M.A.: Building a large annotated corpus in English: the Penn Treebank. *Computational Linguistics* **19** (1993) 313–330
19. Colosimo, M.E., Morgan, A.A., Yeh, A.S., Colombe, J.B., Hirschman, L.: Data preparation and interannotator agreement. *BMC Bioinformatics* **6**(Suppl 1) (2005)
20. Fonseca, C.M., Fleming, P.J.: An overview of evolutionary algorithms in multiobjective optimization. *Evolutionary Computation* **3**(1) (1995) 1–16
21. Kim, J., Ohta, T., Tateisi, Y., Tsujii, J.: GENIA corpus–semantically annotated corpus for bio-textmining. *Bioinformatics* **19** **Suppl 1** (2003) 180–182
22. Riloff, E.: Automatically generating extraction patterns from untagged text. In: Proceedings of the Thirteenth National Conference on Artificial Intelligence (AAAI-96). (1996) 1044–1049
23. Zhou, R., Hansen, E.: Sweep A*: Space-efficient heuristic search in partially-ordered graphs. In: Fifteenth IEEE International Conference on Tools with Artificial Intelligence, Sacramento, CA (2003)
24. Xenarios, I., Salwinski, L., Duan, X., Higney, P., Kim, S., Eisenberg, D.: DIP: The Database of Interacting Proteins. A research tool for studying cellular networks of protein interactions. *Nucleic Acids Research* **30**(1) (2002) 303–305
25. Nobata, C., Sekine, S.: Towards automatic acquisition of patterns for information extraction. In: International Conference of Computer Processing of Oriental Languages. (1999)
26. Pierce, D., Cardie, C.: User-oriented machine learning strategies for information extraction: Putting the human back in the loop. In: Working Notes of the IJCAI-2001 Workshop on Adaptive Text Extraction and Mining. (2001) 80–81

27. Knowles, J., Hughes, E.J.: Multiobjective optimization on a budget of 250 evaluations. In: *Evolutionary Multi-Criterion Optimization (EMO 2005)*, LNCS 3410, Springer-Verlag (2005) 176–190