



ELSEVIER

Information and Software Technology xx (0000) 1–17

**INFORMATION
AND
SOFTWARE
TECHNOLOGY**
www.elsevier.com/locate/infosof

An experiment in software component retrieval

Hafedh Mili*, Estelle Ah-Ki, Robert Godin, Hamid Mcheick

Département d'Informatique, Université du Québec à Montréal, Case Postale 8888 (A), Montréal, Que., PQ, Canada H3C 3P8

Received 29 September 2002; revised 11 November 2002; accepted 30 December 2002

Abstract

Our research centers around exploring methodologies for developing reusable software, and developing methods and tools for building inter-enterprise information systems with reusable components. In this paper, we focus on an experiment in which different component indexing and retrieval methods were tested. The results are surprising. Earlier work had often shown that controlled vocabulary indexing and retrieval performed better than full-text indexing and retrieval [IEEE Trans. Software Engng (1994) 1, IEEE Trans. Software Engng (1991) 800], but the differences in performance were often so small that some questioned whether those differences were worth the much greater cost of controlled vocabulary indexing and retrieval [Commun. Assoc. Comput. Mach. 28 (1985) 289, Commun. Assoc. Comput. Mach. 29 (1986) 648]. In our experiment, we found that full-text indexing and retrieval of software components provided comparable precision but much better recall than controlled vocabulary indexing and retrieval of components. There are a number of explanations for this somewhat counter-intuitive result, including the nature of software artifacts, and the notion of relevance that was used in our experiment. We bring to the fore some fundamental questions related to reuse repositories.

© 2003 Elsevier Science B.V. All rights reserved.

Keywords: Software reuse; Multi-faceted classification; Boolean retrieval; Plain-text retrieval; Retrieval evaluation; Approximate retrieval

1. Introduction

1.1. Component retrieval: do we still care?

Software reuse is seen by many as an important factor in improving software development productivity and software products quality [2,13]. It is customary in the software reuse literature to make the distinction between the *generative* approach whereby developers reuse development *processors* such as code generators or high-level specification language interpreters, and the *building blocks approach*, whereby developers reuse the *product* of previous software development efforts in the process of building new ones. The building blocks approach modifies the traditional, analytical, divide and conquer approach to system specification and design by introducing three reuse tasks that must be performed before one falls back on analytical methods:

(1) searching and retrieving reusable components based on partial specifications, (2) assessing the reuse worth of the retrieved components, and, possibly, (3) tailoring the reusable components to the specifics of the problem at hand [22]. In this paper, we focus on computer support for software component search and retrieval.

The problem of component retrieval has been widely addressed in the software reuse literature. A number of developments have rendered this problem somewhat uninteresting. From a technical point of view, research in the area has hit the formal methods cost barrier: the investment needed to get the next level of performance—to get beyond signature matching or multi-faceted classification—overshadowed the anticipated productivity gains. Second, there was a widespread recognition in the object-oriented reuse community that classes are too small units of reuse, for two reasons. First, classes cannot be reused in isolation. Second, considering that more is gained by reusing designs than by reusing code, individual classes embody mostly code, but little design. Finally, empirical evidence from reuse repositories had shown that small components may account for a good fraction of reuse instances, but in the end, account for little reuse volume,¹ and thus little benefit [12]. The underlying lesson was ‘focus

* Corresponding author. Tel.: +1-514-987-3943; fax: +1-514-987-8477.
E-mail address: hafedh.mili@uqam.ca (H. Mili).

¹ Isoda reported on an experimental reuse program at NTT where they found that components of 50 lines or less accounted for 48% of the reuse instances and 6% of the reuse volume, while modules 1000 lines or larger accounted for only 6% of the reuse instances, but of 56% of the reuse volume [12].

113 on a small number of large components embodying design
114 as well as code', i.e. application frameworks.

115 Interestingly, the Internet has brought repository issues
116 back to the forefront. First, it has enabled a virtual market
117 for software components: developers have been searching
118 the web for software components, both free and for-fee, for
119 the past decade. Second, inter-enterprise (B2B) electronic
120 commerce relies on enterprises ability to 'plug-in' each
121 other's systems to be able to complete transactions, end to
122 end. The ability to plug systems together has become a
123 major factor in entering into business relationships [2],
124 some times the overriding one [26]. The pluggability of
125 information systems for the purposes of entering into
126 electronic commerce starts with the lookup of industry-wide
127 registries of APIs exported by potential partners. Standards
128 are emerging to represent such APIs in a technology
129 independent way (see e.g. ebXML [26]), but the issue of
130 conceptual appropriateness remains whole. Notwithstand-
131 ing things such as ebXML registries or software vendor-
132 specific web sites, it seems that much reuse is taking place in
133 the unstructured world of the world wide web, as opposed to
134 a corporate managed reuse repository with dedicated
135 personnel and strict quality control. This paper explores
136 component classification and retrieval methods with an
137 overriding concern for automation.

139 1.2. The component retrieval problem

140
141 A wide range of component categorization and searching
142 methods have been proposed in the literature, from the
143 simple string search (see e.g. Ref. [21]), to faceted
144 classification and retrieval (e.g. Refs. [27,28]) to signature
145 matching (see e.g. Ref. [37]) to behavioral matching (see
146 e.g. Refs. [10,17,38]). Different methods rely on more or
147 less complex descriptions for both software components and
148 search queries, and strike different trade-offs between
149 performance and cost of implementation [22]; the cost of
150 implementation involves both initial set-up costs, and the
151 cost associated with formulating, executing and refining
152 queries. In the context of our research, we developed four
153 classes of retrieval algorithms (1) retrieval using full-text
154 search on software documents and program files, (2) multi-
155 faceted classification and retrieval of components, (3)
156 navigation through the structure of components, and (4)
157 signature matching. The first two use the documentation or
158 the *meta-data* that accompanies software components, and
159 thus rely on its existence, its quality, and some pre-
160 processing. The last two focus on the *structure* of the
161 software components themselves, and thus depend on the
162 availability of that structure in some form—source code,
163 interface—and the availability of (computer) language
164 processors.

165 An age-old debate, first in the information retrieval
166 literature [4,31], and later in the context of reuse repositories
167 [6,8,16,23], has opposed the free-text classification and
168 retrieval of components to the so-called controlled vocabu-

lary, multi-faceted classification and retrieval of com- 169
ponents. The conventional wisdom is that free-text 170
retrieval costs nothing—no manual labour—but produces 171
many false positives (matches words taken out of context) 172
and false negatives (misses out relevant components 173
because of the use of a non-standard terminology). 174
Controlled-vocabulary indexing and retrieval is supposed 175
to solve both problems by providing a common vocabulary 176
for classification and retrieval, and by having actual human 177
beings classify documents/components. However, it 178
involves a major cost in building and maintaining such 179
vocabularies and in classifying/indexing components. 180
Research in the area has traditionally attempted to bridge 181
the gap between the two approaches in terms of cost and 182
performance. From the free-text end, research has aimed at 183
making the matching more intelligent and less dependent on 184
surface-level similarity, but keeping humans out of the 185
loop—e.g. using associations between terms instead of term 186
matching or identity, as in latent semantic analysis methods 187
[6,11,16]. From the controlled vocabulary end, research has 188
aimed at automating or assisting the manual steps, but 189
hopefully without losing much in terms of quality of 190
retrieval. Our own work has covered both approaches, and 191
this paper reports on a number of experiments trying out 192
different ideas and comparing approaches. 193

194 Our first experiment dealt with the construction of
195 domain vocabularies. Much of the earlier work on
196 automated indexing of textual documents had relied on
197 the statistics of the occurrences (and co-occurrences) of key
198 terms or phrases within document collections to infer
199 content indicators for documents and relations between key
200 terms [15,30]. Our work furthers these ideas to build
201 concept hierarchies based on statistics of (co)occurrences
202 alone. A technique that worked well in previous exper-
203 iments was less successful with software documentation.
204 The experiment is described, and the results are analyzed in
205 Section 3. The second experiment dealt with the automatic
206 indexing of software components (their documentation)
207 using a controlled vocabulary: the basic idea is that an index
208 term (say 'Database Management Systems') is assigned to a
209 component if 'most' of its constituent words appear 'close'
210 to each other within the documentation of the component;
211 most and close are both tunable parameters of the method.
212 In principle, the automatic assignment of index terms suffers
213 from the same problems as free text search: matching words
214 out of context (false positives), and missing out on relevant
215 components because of choice of terminology (false
216 negatives). However, we felt that the use of compound
217 terms would reduce the chances of false positives, and the
218 use of inexact matches (most, close) would reduce the
219 chances of false negatives. The results bear this out, and are
220 discussed in Section 4.

221 Our third experiment consisted of comparing an all-
222 manual controlled vocabulary indexing and retrieval
223 method with an all-automatic free-text indexing and
224 retrieval method, using a variant of the traditional

information retrieval measures, *recall* and *precision*. Instead of computing recall and precision based on some abstract measure of ‘relevance’, as is done in information retrieval and in most reuse library experiments, we adapted the measure to take into account the true utility of the retrieved components to solve the problem at hand. Further, we used a realistic experimental protocol, one that is closer to the way such tools would be used in practice. Here the results were surprising. Full-text retrieval yielded *significantly* better recall and somewhat better precision—although the difference is statistically insignificant. The experiment is described in Section 5. We analyze the results in light of new evidence about the behavior of users in an information retrieval setting. We conjecture that multi-faceted retrieval requires more information than the user is able to provide in the early stages of problem solving, and fails to capture a faithful expression of users’ needs at the later stages.

Section 2 provides a brief introduction to our tool set. We conclude in Section 6.

2. ClassServer: an experimental component repository

2.1. Overview

This work is part of ongoing research at the University of Québec at Montréal aiming at developing methods and tools for developing reusable software, and for developing with reusable software. The work described in this paper centers around a tool kit called *ClassServer* that consists of various tools for classifying, retrieving, navigating, and presenting reusable components (see Fig. 1). Reusable components consist essentially of object-oriented source code components, occasionally with the accompanying textual documentation. Raw input source files are put through

various tools—called *extractors*—which extract the relevant pieces of information, and package them into ClassServer’s internal representation format for the purposes of supporting the various reuse tasks. So far, we have developed extractors for Smalltalk and C++. The information extracted by these tools consists of built-in language structures, such as *classes*, *variables*, *functions*, and *function parameters*. To these, we added a representation for *object frameworks*, which are class-like object aggregates that are used to represent application frameworks and design patterns [24]; unlike the built-in language structures, which are extracted by parsers, object frameworks need to be manually encoded. Fig. 1 shows a very schematic view of the ClassServer tool set. The tool set may be seen as consisting of three subsystems. The first subsystem, labeled ‘Full-text retrieval’ supports the required functionalities for full-text retrieval of source code *files*, namely, the ‘Full-text indexer’, and the ‘Full-text search tool’. Their functionalities are explained in Section 2.3.1.

The component browser and the keyword retrieval subsystems use the structured representation of the components that is extracted by the tool referred to as ‘semantic/structural parser’ in Fig. 1. Typically, the parsing produces a trace of the traversal of the abstract syntax tree. The trace consists of a batch of component creation commands (in Smalltalk), which are executed when we ‘load’ the trace; that is the *structured component loader*. Each kind of component is defined by a descriptive template that includes: (1) structural information describing the kind of subcomponents a component can or must have (e.g. a *class* has *variables* and *methods*, a *framework* has *participants*, *message sequences*, etc.), (2) code, which is a string containing the definition or declaration of the component in the implementing language, and (3) descriptive attributes, which are used for search purposes; for

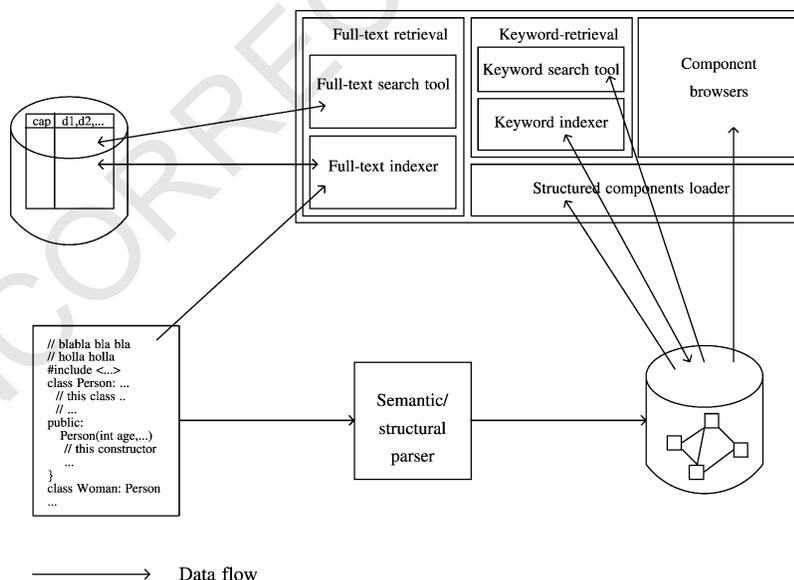


Fig. 1. Overall architecture of ClassServer.

337 example, a *class* has an *author* and an *application domain*, a
 338 *method* has a *purpose*, etc. Descriptive attributes, or simply,
 339 attributes, represent non-structural, non-intrinsic properties
 340 of software components, and are often derived from non-
 341 code information such as documentation, or entered
 342 explicitly by the person(s) responsible for managing the
 343 component library. Attributes will be described in more
 344 detail in Section 2.2.

345
 346 2.2. A multi-faceted classification of components

347
 348 Attributes are used in ClassServer to represent categor-
 349 ization/classification facets, as in Prieto-Diaz’s multi-
 350 faceted categorization of components [28]. Attributes are
 351 themselves objects with two properties of their own: (1) *text*,
 352 which is a (natural language) textual description, and (2)
 353 *values*, which is a collection of key words or phrases, taken
 354 from a predefined set referred to as the *vocabulary of the*
 355 *attribute*. The text is used mainly for human consumption
 356 and for documentation generation [21]. Filling in the *values*
 357 property is referred to as *classification*, *categorization* or
 358 *indexing*. When human experts assign those key words or
 359 phrases from a predefined list, we talk about *manual*
 360 *controlled-vocabulary indexing* [30]. In our case, we used
 361 *automatic controlled-vocabulary indexing* whereby a key
 362 word or phrase is assigned to an attribute if it occurs within
 363 the text field. More on this in Section 4.

364 For a given attribute multiple values are considered to be
 365 *alternative* values (ORed), rather than *partial* values
 366 (ANDed). For example, for the attribute ‘Purpose’ of a
 367 component, several values mean that the component has
 368 *many* purposes, and *not* a single purpose defined by the
 369 conjunction of several terms. For a given vocabulary, the
 370 terms of the vocabulary (key words and phrases) may be
 371 organized along a conceptual hierarchy. Fig. 2 shows
 372 excerpts of the conceptual hierarchies of key phrases for the
 373 attributes ‘Application Domain’ (Fig. 2a) and Purpose
 374 (Fig. 2b). Notice that the Application Domain hierarchy of
 375 key phrases is inspired from the (ACM) *Computing*
 376 *Reviews*’s classification structure [1]. The hierarchical

relationship between key phrases is a loose form of 393
 generalization, commonly referred to in information 394
 retrieval as ‘Broader-Term’ [30]. Attribute values (key 395
 words and phrases) are used in boolean retrieval whereby 396
 component attribute values are matched against required 397
 attribute values (queries, see below). The hierarchical 398
 relationships within an indexing vocabulary are used to 399
 extend the basic retrieval algorithms, as explained in 400
 Section 2.3.2. 401

402
 403 2.3. Software component retrieval in ClassServer

404
 405 As mentioned earlier, ClassServer provides two methods
 406 of classifying (and retrieving) software components,
 407 namely, free-text indexing and search of software com-
 408 ponents (source code and documentation), and multi-
 409 faceted classification and retrieval of components. We
 410 describe them both briefly below. 411

412 2.3.1. Free text indexing and search

413 By free-text indexing, we refer to the class of methods
 414 whereby the contents of a document are described by a
 415 weighted set of words or lexical units occurring in the
 416 document. Different methods use different selection mech-
 417 anisms to restrict the set of eligible content indicators, and
 418 different weighting schemes [30]; the algorithm we used
 419 does not use a weighting scheme. Let us assume for the
 420 moment that *all* the words found in a document are used as
 421 potential content indicators. Given a natural language query
 422 Q , the free-text retrieval algorithm returns the set of
 423 components S computed as follows: 424

- 425 (0) Break the query Q into its component words
 426 w_1, \dots, w_n ,
- 427 (1) $S \leftarrow$ set of components whose documentation
 428 included w_1 ,
- 429 (2) **For** $i = 2$ **To** n **Do**
- 430 (2.1) $S_i \leftarrow$ set of components whose documentation
 431 included w_i ,
- 432 (2.2) $S \leftarrow S \cap S_i$ 433

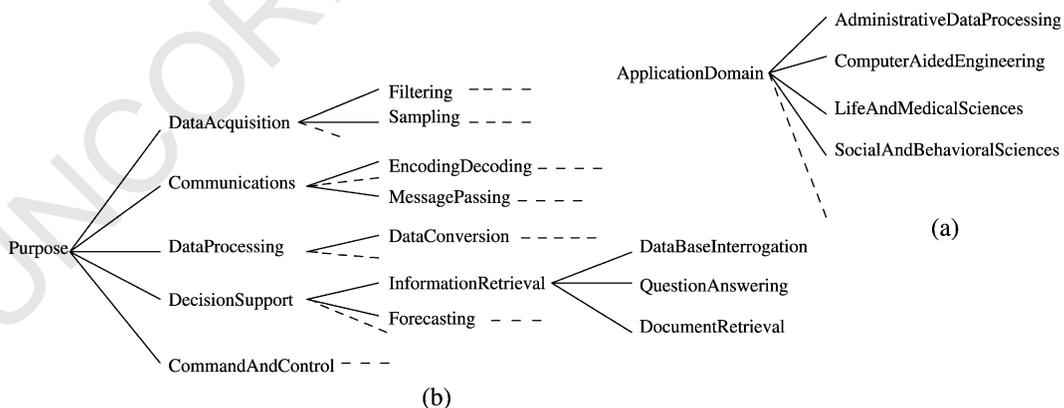


Fig. 2. Hierarchies of key phrases for the attributes Application Domain and Purpose.

When we break a query into its component words, we exclude all the words that are not significant in the application domain. This includes common language words such as ‘the’, ‘an’, ‘before’, and so forth. It also includes domain specific words that are likely to be found in every document—software component documentation in this case. For example, we would expect the word ‘computer’ to appear everywhere in a computer science collection. These are called *stop* words. In order to account for lexical variations when matching words of the query to words of the documents, we reduce both to their roots, as in mapping ‘Managing’ and ‘Management’ to ‘Manag’. Algorithms to perform this mapping are called *word stemmers*, and we used one published in Ref. [9]. Finally, to speed search, we pre-process the entire document collection by creating an inverted list which is a table whose keys are unique words stems such as Manag, and whose values are lists of the documents in which the word occurred in one lexical form or another (e.g. as Managing or Management). This reduces the step (2.1) above to a simple table look-up.

2.3.2. Multi-faceted controlled-vocabulary retrieval

Our choice for the representation of queries involved a trade-off between flexibility and expressiveness, on the one hand, and allowing users to specify the most common queries most easily and most efficiently, on the other. The simplest form of a query is a list of so called *attribute query terms* (AQTs), considered to be ANDed. In its simplest form, an AQT consists of an attribute, and a list of key phrases, considered to be ORed. In the actual implementation, each AQT is assigned a weight and cut-off point, used for weighted boolean retrieval and conceptual distance, respectively (see below). Symbolically:

- * Query :: = AQT|AQT AND Query
- * AQT :: = Attribute Weight CutOff ListOfKey-Phrases
- * ListOfKeyPhrases :: = KeyPhrase|KeyPhrase OR ListOfKeyPhrases

A single AQT retrieves the components whose attribute ⟨Attribute⟩ has at least one value in common with ⟨ListOfKeyPhrases⟩. Viewing attributes as functions, an AQT denoted by the four-tuple ⟨Attribute, Weight, Cut Off, ListOfKeyPhrases⟩ retrieves the components C such that $\text{Attribute}(C) \cap \text{ListOfKeyPhrases} \neq \Phi$. The query denoted by the tuple $(\text{AQT}_1, \dots, \text{AQT}_k)$, returns the intersection of sets of components that would have been returned by the individual AQTs.

With weighted boolean retrieval, components are assigned numerical scores that measure the extent to which they satisfy the query, instead of being either ‘in’ or ‘out’. Let Q be a query with terms $(\text{AQT}_1, \dots, \text{AQT}_k)$, where $\text{AQT}_i = \langle \text{Attribute}_i, \text{Weight}_i, \text{CutOff}_i, \text{ListOfKeyPhrases}_i \rangle$. The score of a component C is

computed as follows:

$$\text{Score}(Q, C) \equiv \frac{\sum_{i=1}^k \text{Weight}_i \times \text{Score}(\text{AQT}_i, C)}{\sum_{i=1}^k \text{Weight}_i} \quad (1)$$

where $\text{Score}(\text{AQT}_i, C)$ equals 1.0 if $\text{ListOfKeyPhrases}_i \cap \text{Attribute}_i(C) \neq \Phi$, and 0 otherwise.

Another extension meant to handle approximate matches is based on the number of edges separating the key terms of the query from the key terms of the attribute of the component in the conceptual hierarchies that enclose them (as in Fig. 2). If, for some i , $\text{ListOfKeyPhrases}_i \cap \text{Attribute}_i(C) \neq \Phi$, we look at some aggregate of the path lengths that separate elements of $\text{ListOfKeyPhrases}_i$ from elements of $\text{Attribute}_i(C)$ and use that to assign a score between 0 and 1 for the query term; the higher the average distance, the lower the score. The mathematical properties of the resulting similarity metric—called DISTANCE—and its effectiveness at emulating human relevance judgements have been thoroughly documented in Ref. [29]. In ClassServer, the cut-off value puts an upper limit on the path lengths to be considered in the computation; key phrases that are separated by more than ‘cut-off’ edges are considered totally unrelated.² A third extension uses the hierarchical relationships between key terms to ‘classify’ the query within a virtual classification structure of components that is based on the relationships between their attribute values, returning the most ‘specific’ components that are more ‘general’ than the query. The ‘specialization’ relationship has a formal meaning in this case [17]. Neither of the last two extensions was used in the experiments of Section 5, and will not be discussed further.

2.4. The component library

For the purposes of the experiment described in Section 5, we loaded the ClassServer repository with the OSE library [7] which contained some 200 classes and 2000 methods distributed across some 230 *.h files with, typically, one class per file. For the purposes of supporting plain-text indexing and retrieval, the 230 files were put through the plain text indexing tool, which generated an inverted list of unique word stems (see Section 2.3.1). Further, a shell script put the files through a C++ pre-processor before they were input into the C++ extractor (see Section 2.1). Because of the good quality and format consistency of the in-line documentation (comparable to Javadoc), we were able to automatically assign C++ comments as text values for the ‘Description’ attribute of various components (classes, methods, variables). Overall, we classified components using two attributes Application-Domain, and Description. ApplicationDomain was indexed manually, but in a fairly systematic fashion, using the on-

² This ‘sunsetting’ is used to fix some singularities in the otherwise well-behaved similarity metric [29].

561 line documentation of the library. In fact, the section
 562 headers of the documentation were themselves used as
 563 index terms (see Ref. [20] for a justification). The
 564 documentation grouped the various classes by application
 565 area. Further, each class was first described by a general
 566 statement about what the class does, followed by a more
 567 detailed description of its services, which mapped closely to
 568 methods. Some utility methods were not documented, and
 569 we could not assign those an ApplicationDomain; however,
 570 all classes were properly classified.

571 For the Description attribute (telling what a component
 572 does and how it does it, rather than ‘what it is used for’), we
 573 did not have a ready-made indexing vocabulary. We
 574 considered using available classification structures that
 575 include computer science concepts, including the 1200 +
 576 terms Computing Reviews classification structure [1].
 577 However, the classification terms were too general to be
 578 of any use to our library of components. For example,
 579 whereas we needed terms that corresponded to the different
 580 sorting algorithms (‘MergeSort’, ‘RadixSort’), the term
 581 ‘Sorting’ was a leaf node of the ACM hierarchy. Accord-
 582 ingly, we decided to develop our own vocabulary by
 583 analyzing the available software documentation; the process
 584 of building the vocabulary is described next. Further, we
 585 decided to perform the actual indexing of the attribute (the
 586 assignment of key terms to attribute values) automatically.
 587 The algorithm and the results are discussed in Section 4.
 588

590 3. Constructing domain vocabulary

592 A hierarchy of the important concepts in a domain has
 593 many uses in the context of software component retrieval. In
 594 addition to the advantages of having a standard vocabulary,
 595 its hierarchical structure helps ‘librarians’ locate the most
 596 appropriate term to describe a component, and ‘re-users’
 597 find the closest term to their need to use in a search query.
 598 Those same relations may also be used to extend boolean
 599 retrieval methods to account for ‘close’ matches, as shown
 600 in Section 2.3.2 (see also Refs. [11,27]). Constructing a
 601 hierarchy of the important concepts in a domain (or
 602 thesaurus) involves identifying those important concepts
 603 and their preferred terminology (Section 3.1), and organiz-
 604 ing them into a hierarchy (Section 3.2). We discuss these
 605 issues in turn.
 606

608 3.1. Extracting a set of concepts

610 A good place to look for the important computer science
 611 concepts that are germane to a library of reusable
 612 components is the documentation of the library itself. By
 613 looking *only* at the documentation, we run the risk of getting
 614 a partial and narrow view of the underlying domain, and of
 615 depending too much on the terminology used by the
 616 documenter. At the same time, we are assured that we will

617 not miss any concepts that are important to the particular
 618 library (or libraries) at hand.

619 The next question is one of identifying the right lexical
 620 unit that corresponds to key concepts, and extracting such
 621 units from the text. Computer science being a relatively new
 622 field, most of the important concepts are described by noun
 623 phrases, as in ‘Software Engineering’ ‘Bubble Sort’,
 624 ‘Printing Monitor’, and so forth, rather than single words
 625 as is the case for more mature fields such as medicine.³ In
 626 order to extract those higher level lexical units, to which we
 627 will abusively refer as noun phrases, we used *Xerox Part Of*
 628 *Speech Tagger* (XPost) [5]. XPost is a program that takes as
 629 input a natural language text and produces the syntactic
 630 (‘part of speech’) category or *tag* for each word or *token* for
 631 the text. For example, it assigns to the phrase ‘The common
 632 memory pool’ the tag sequence ‘at jj nn nn’, where ‘at’
 633 stands for *article*, ‘jj’ for *adjective*, and ‘nn’ for *noun*. XPost
 634 uses two major sources of information to assign tags to
 635 words of a sentence: (1) a ‘tag table’, giving the set of tags
 636 that correspond to a given token, and (2) a probabilistic
 637 (markovian) model of the allowable sequences of tags. For
 638 example, the word ‘book’ can be a noun (‘nn’) or a verb
 639 (‘vb’). If we also know that ‘the’ is an article and that only
 640 nouns can follow articles, we know that ‘book’ in the phrase
 641 ‘the book’ is a noun. XPost falls within the category of parts
 642 of speech taggers that derive the probabilistic model using
 643 unsupervised learning [5].
 644

645 One of the typical uses of XPost is to extract phrases that
 646 follow a given pattern. We used XPost to extract ‘noun
 647 phrases’ that are likely to represent important domain
 648 concepts. To this end, we ran XPost on a training sample, we
 649 identified the tag sequences for the noun phrases in which
 650 we were interested, and then looked for a set of regular
 651 expressions that would have extracted those phrases. Those
 652 regular expression were then used to filter the output of
 653 XPost to extract noun phrases. The first set of regular
 654 expressions (a grammar) accepted far too many phrases, and
 655 we had to refine the grammar through trial and error, with a
 656 bias towards minimizing false positive phrases, at the
 657 expense of missing out some valid phrases. Fig. 3 shows the
 658 regular expressions in an awk-like format.

659 In case a sentence matched several expressions, we take
 660 the longest running expression. For example, if we analyze
 661 the sentence ‘Memory management of event based
 662 systems’, we produce a single noun phrase consisting of
 663 the entire sentence, rather than the two phrases ‘Memory
 664 management’ and ‘event based systems’, both of which
 665 matching the pattern BASIC.

666 We used this approach on the on-line documentation of
 667 the library. The documentation consisted of 13 html files,
 668 one of which giving an overview of the library, and the
 669 remaining 12 describing specific subsets of the library. The
 670 13 files contained a total of 37,777 words (244 Kbytes). The

671 ³ See Ref. [39] for a discussion on the evolution of languages and
 672 terminology.

PREFIX	≡ (JJ VBG VBN)(NN NNS NP NPS) (NN NNS NP NPS)	729
#	Example "Small System", or "system", but not "small"	730
BASIC	≡ PREFIX (JJ NN NNS NP NPS VBG VBN)*	731
#	Example "Event based systems"	732
X_OF_Y	≡ BASIC IN BASIC	733
#	Example "Memory management of event based systems"	734
X_OF_A_Y	≡ BASIC IN AT BASIC	735
#	Example "storage requirements of an event based system"	736

Fig. 3. Grammar for noun phrases.

extraction process identified 2616 unique noun phrases, with overall occurrences ranging from 163 (for the word ‘function’) to 1, with 1,765 phrases occurring just once, including phrases such as ‘command line options’ or ‘Conversion operator to a standard pointer’. Typically, phrases that occur too often are not good discriminators [30]. Further, phrases that occur rarely may not be important for the domain at hand. We found 8 ‘phrases’ that occurred more than a 100 times, and discarded them: OTC (the name of the library, 267 times), ‘Function’ (163), ‘String’ (134), ‘Member function’ (114), ‘Object’ (114), ‘Class’ (113), ‘Example’ (105), and ‘Program’ (104). We also discarded the phrases that occurred less than five times. Overall, we used 229 phrases. These include C++ identifiers that may have appeared in the code examples, and possibly referred to thereafter in the running text. We could have removed them from the vocabulary that was fed into the hierarchy builder, but we chose to exclude any manual processing or decisions that cannot be systematized or automated.

3.2. Constructing a hierarchy of important domain concepts

Having identified a set of the important concepts in a domain, we need to organize those concepts in a conceptual hierarchy. We present a simple algorithm that does just that based on statistics of occurrences of these concepts in documents. Next, we describe an earlier experiment with the algorithm that provided encouraging results. We conclude with the results of the algorithm on the set of concepts extracted with the method described in Section 3.1.

3.2.1. Principles

Given a set of terms $T = \{t_1, \dots, t_m\}$, a set of documents $D = \{d_1, \dots, d_m\}$ with manually assigned indices $\text{Idx}(d_i) = \{t_{i_1}, t_{i_2}, \dots\}$, we argued that [18]:

- H₁ Terms that co-occurred often in document indices were related in a way that is important to the domain of discourse,
- H₂ The more frequently occurring a term, the more general its conceptual scope, and
- H₃ If two terms co-occur often in document indices (and thus are related, according to H₁), and if one has a more general scope than the other, than there

is a good chance that the relationship between them is a generalization/specialization-like relationship.

The H₁ hypothesis is based on fact that documents tend to exhibit conceptual cohesion and logic, and because index terms reflect the important concepts within a document, they tend to be related. The second hypothesis is based on observations made about both terms occurring in free-format natural language [14] as well as index terms [34].

We developed an algorithm that generates an acyclic graph with a single node with in-degree 0 (root) based on the above hypotheses [18]. Given m index terms t_1, \dots, t_m , the algorithm operates as follows:

- (1) Rank the index terms by decreasing order of frequency,
- (2) Build a matrix of co-occurrences (call it M) where the i th row (column) corresponds to the i th most frequent term,
- (3) Normalize the elements of the matrix M by dividing $M(i, j)$ by the square root of $M(i, i) \times M(j, j)$; note that after this normalization, $M(i, j) \leq 1$,
- (4) Choose terms to include in the first level of the hierarchy; assume that the terms t_1 through t_{l_1} were chosen to be included in the first level,
- (5) For $i = l_1 + 1$ through m
 - 5.1 Find the maximum of the elements $M(i, 1)$ through $M(i, l_1 - 1)$. Note that because of the ordering of rows and columns (step 2), these are the frequencies of co-occurrences of t_i with the terms whose occurrences are higher than that of t_i ,
 - 5.2 Create a link between the term t_i and all the terms t_j such that $j < i$ and $M(i, j) = \text{maximum found in 5.1}$.

The choice of the first level nodes is quite arbitrary although, ultimately, it has a very little impact on the overall hierarchy.

3.2.2. A case-study: the Genbank experiment

In one experiment, we used the GenBank genetic sequences database (databank). The GenBank Genetic Sequence DataBank serves as a repository for genetic sequences [3]. The entry for each sequence includes, among

785 other things, the article that reported the discovery of the
 786 sequence, and a set of keywords that describe the sequence,
 787 including names of components or processes that are
 788 involved either in the composition and transformation of
 789 the sequence, or it is discovery. For our purposes, each entry
 790 corresponded to a document. We ran the experiment on
 791 5700 such ‘documents’. The co-occurrences matrix was
 792 limited to those keywords that occurred more than 10 times,
 793 and there were 274 of those. The resulting hierarchy was
 794 evaluated both qualitatively and quantitatively. The quali-
 795 tative evaluation had to do with whether the parent-child
 796 links that were created were meaningful in general
 797 (hypothesis H₂), and whether they were generalization/
 798 specialization-like in particular (hypothesis H₃). Experts
 799 found that 50% of the links were indeed ‘generalization/
 800 specialization’ (G/S), as in the link between ‘Heavy Chain
 801 Immunoglobulins’ and ‘Immunoglobulins’. Another 15% of
 802 the links were deemed meaningful as when the two terms
 803 represent a chemical component, and the process that
 804 creates it. The remaining 35% could not be characterized.
 805 Clearly, the resulting hierarchy was by far not as ‘coherent’
 806 or ‘enlightening’ as manually built hierarchies such as the
 807 Computing Reviews Classification Structure, for example.
 808

809 The quantitative evaluation had to do with the extent to
 810 which the resulting hierarchy supported extended boolean
 811 retrieval (DISTANCE-based, see Section 2.3.2) of docu-
 812 ments any better or worse than a manually built hierarchy⁴
 813 that contained the same terms. For a given hierarchy H, the
 814 evaluation consists of: (1) using DISTANCE on H to rank a
 815 set of documents by order of relevance with respect to a set
 816 of queries, (2) asking human subjects to do the same, and (3)
 817 computing the correlation between the two rankings; the
 818 higher the correlation, the more faithful is distance to human
 819 evaluation, and the more useful is the hierarchy. Our
 820 experiments showed that the automatically constructed
 821 hierarchy performed as well, if not better than the manually
 822 built one [18].

823 Overall, the experiments showed that while the hierarchy
 824 may not be ‘user-friendly’ or make as much sense as a
 825 manually built one, it can perform useful retrieval tasks
 826 equally well. We had observed that the keywords did not
 827 belong to a single conceptual domain, and that across-
 828 domain relationships could dominate within-domain ones.
 829 An algorithm that focuses on the strongest relationships
 830 would miss potential generalization relationships. For
 831 example, we had chemicals as well as chemical processes,
 832 and we had hypothesized (but not tested) that, had we
 833 separated them and applied the algorithm to the separate
 834 sets, we might have gotten more consistent hierarchies [18].
 835 In other words, we felt that there was room for
 836 improvement.
 837

838 ⁴ The Medical Subject Headings hierarchy, maintained by the National
 839 Library of Medicine, and used to support its on-line bibliographic retrieval
 840 system MEDLINE [32].

3.2.3. Constructing the graph based on OSEs on-line documentation 841 842

843 The construction of the hierarchy requires co-occurrence
 844 data between phrases within relatively coherent text units.
 845 We can break the documentation different ways, where a
 846 ‘document’ may be either, an entire file, a major section
 847 within a file, a subsection within a file, or even a paragraph.
 848 Whatever the document, we have to make sure that: (1) the
 849 phrases are good content indicators for that document, and
 850 (2) the co-occurrence of two phrases within the same
 851 document is not fortuitous and does reflect a significant
 852 relationship. The first constraint may suggest that we use
 853 documents that are big enough that phrase occurrence
 854 statistics become significant. The second constraint suggests
 855 that we use documents that are small enough that phrase co-
 856 occurrence be confined to a coherent textual unit. We
 857 decided to use subsections in files (an average of 10
 858 subsections per file) as documents. Further, for each
 859 document, if a phrase P₁ occurred *m* times and a phrase P₂
 860 occurred *n*, we consider that the phrases co-occurred
 861 minimum(*m*, *n*) times.

862 The first run of the algorithm generated a hierarchy with
 863 291 relations between 291 phrases, including the dummy
 864 root node. Because we had no other hierarchy to which to
 865 compare it on a specific task, as was the case for the
 866 experiment described in Section 4.1, we could only evaluate
 867 the hierarchy *qualitatively*. To this end, we presented six
 868 subjects with the hierarchy and asked them to mark, for each
 869 node, whether the node represented a valid concept from the
 870 domain of discourse, and in case it did, to label the node’s
 871 relationship to its parent as one of (a) *has broader-term* [33],
 872 which is a loose form of generalization, (b) *related*, to
 873 indicate any relationship other than has broader term, and
 874 (c) *unrelated*. *Unrelated* was used when there was no
 875 apparent relationship between a node and its parent. We
 876 show below excerpts from the hierarchy to illustrate the
 877 three kinds of relations. The relationship between LENGTH
 878 OF THE STRING and LENGTH is *has-broader-term*. That
 879 between RANGE and LENGTH is *related*.
 880

```

881 ... 881
882 0.2.1.1.2.1 LENGTH 882
883 0.2.1.1.2.1.1 LENGTH OF THE STRING 883
884 0.2.1.1.2.1.2 CAPACITY 884
885 0.2.1.1.2.1.2.1 CAPACITY OF THE STRING 885
886 0.2.1.1.2.1.3 RANGE 886
887 ... 887
888 0.2.1.1.2.3.2 B 888
889 0.2.1.1.2.3.2.1 CONVERSION 889
890 0.2.1.1.2.3.2.1.1 SUBJECT 890
891 0.2.1.1.2.3.2.1.2 CONVERSION OPERATOR 891
892
```

893 We note the ‘term’ B, which is a C++ identifier that was
 894 tagged by XPost as a noun, because it is not a known verb or
 895 noun, and because it occurred in the text where a
 896 subject/object was expected. B occurred enough times to

make it into the vocabulary. As mentioned earlier, we decided to leave such terms in to get an idea about what the hierarchy would look like without any manual filtering. In this case, not only B should not have been there, but all of the relationships between B and its children (CONVERSION) are non-significant. Such relationships are labeled as unrelated. The relationship between CONVERSION and SUBJECT is an interesting one. SUBJECT is the name of the class representing strings. This class supports several conversion operations, and hence the association. Somebody thinking of CONVERSION in general, would not think of strings. However, in the context of this library, the association is important and useful. This is similar to the kind of indirect associations between keywords exploited by the CODEFINDER system [11], which reflect the structure of the library as much as it reflects the structure of the semantic domain.

The evaluation of the six subjects are summarized in Table 1. The second line shows the results obtained by rederiving the hierarchy after we have removed the invalid terms (26 of them). Notice that because not all 26 terms were leaf nodes, by removing them we needed to reassign parents to 18 valid terms.

These results are disappointing compared to those obtained in the GenBank experiment [18], even after we remove manually the invalid terms from the input. The reasons are easy to identify. In the GenBank experiment the terms of the hierarchy did indeed describe important concepts in the domain, as opposed to the indiscriminate noun phrases extracted from our software documentation. We attempted a number of refinements using statistical measures to eliminate ‘spurious’ terms. Our first attempt was to eliminate the terms with the lowest frequency (5). This reduced the number of terms from 291 to 194, but ironically, only one non-applicable term was eliminated, and the distribution of the remaining relationships (has-broader-term, related, and unrelated) remained about the same. We used another measure of the *information value* carried by a given term, i.e. the extent to which it differentiates a specific and relatively small subgroup of the document set. Let T be a term, and d a document, we define $FREQ(T, d)$ as the number of occurrences of T in d , and $FREQ(T)$ as the total number of occurrences of T . The

Table 1
Evaluating the individual links created by the statistical algorithm

Hierarchy	Percentage of invalid terms	Percentage of has-broader-term	Percentage of related	Percentage of unrelated
With invalid terms	9	20	37	34
Without invalid terms	0	27	39	34
Links removed	26	8	19	28
Links added	0	17	13	18

Table 2
Evaluating the hierarchy after filtering the terms that occurred more than 20 times, and whose entropy is more than half of the maximum possible entropy

Percentage of non-app. term	Percentage of has-broader-term	Percentage of related	Percentage of unrelated
7	19	36	38

entropy of a term T is defined as follows:

$$ENTROPY(T) = \sum_{d \in \text{Documents}} \frac{FREQ(T, d)}{FREQ(T)} \times \log \frac{FREQ(T)}{FREQ(T, d)}$$

For a given number of occurrences $FREQ(T) = N$, the entropy is maximal if N are evenly spread across the document collection. If there are N documents, that entropy is $\log(N)$, and it correspond to T occurring exactly once in each of N documents. Let $MAXENTROPY(T)$ be that maximum. Generally speaking, good terms are the ones with the smallest spread possible, i.e. whose entropy is closest to zero. Accordingly, we filtered the terms based on the ratio $ENTROPY(T)/MAXENTROPY(T)$: among the terms that occurred more than a threshold frequency⁵ F_0 , we rejected the ones for which the above ratio is above a certain threshold ρ . We tried several values of F_0 , and several values of ρ . For $F_0 = 20$, and $\rho = 0.75, 0.666$, and 0.5 , we eliminated 14, 22, and 37 terms, respectively, from the initial set of 194 terms. Table 2 shows an evaluation of the relationships within the generated hierarchy when $F_0 = 20$, and $\rho = 0.5$.

By looking at the remaining list of terms, a considerable number remain that should not be there. Hence, this test is not very effective at filtering invalid terms.

The second explanation for these results is related to the size of the document set. The GenBank experiment used 5700 documents, while this one used 120 documents. This makes statistical inferences unreliable. Finally, because we are dealing with software documentation, the terms tend to be rather specific, and their common ancestors are less likely to appear within the document set. We hypothesized that the higher level relationships cut across branches of a ‘virtual hierarchy’. This is consistent with the earlier observation that, from a conceptual scope point of view, the concepts we need to describe software components tend to be at the lowest levels of the ACM classification structure, or even lower. This means that, potentially, most of the second level relationships are invalid since the software documentation is not likely to contain general computer science terms, or if it does, those will appear infrequently. Table 3 shows a level by level breakdown of relationships. The overall degra-

⁵ If a term occurred only a handful of times, its ENTROPY will be close to the maximum even in those cases where it identifies a narrow subset of documents. F_0 is the overall frequency over which we start ‘demanding’ focussed occurrences. We used thresholds that were close to 1/5th the size of the document set, which here means around 24.

1009 Table 3
1010 Distribution of links across the levels of the hierarchy

	No. of terms	Percentage of invalid terms	Percentage of has-broader-term	Percentage of related	Percentage of unrelated
1015 Level 2	10	10.0	20.0	60.0	10.0
1016 Level 3	24	0.0	37.5	50.0	12.5
1017 Level 4	36	13.0	22.0	47.0	16.67
1018 Level 5	53	9.0	21.0	41.0	28.0
1019 Level 6	65	6.0	20.0	35.0	38.0
1020 Level 7	47	6.0	17.0	30.0	46.0
1021 Level 8	28	11.0	21.0	29.0	39.0
1022 Level 9	14	14.0	0.0	36.0	50.0
1023 Level 10	6	0	16.67	16.67	66.67

1024 dation of the quality of the links within the hierarchy as we
1025 go down is consistent with the unreliability of the results for
1026 the less frequent terms; we cannot make much of the fact
1027 that the level links are of a lesser quality than the level 3
1028 terms, but the above hypothesis is worth exploring.

1029 We considered merging the resulting hierarchy with the
1030 ACM hierarchy (see e.g. Ref. [19]) whereby, if a term T
1031 appears in both ACM and the automatically generated
1032 hierarchy, we carry over the subtree from the automatically
1033 generated tree to the ACM subtree. We found only eight
1034 such common terms between the ACM tree (1200 + nodes)
1035 and the automatically generated hierarchy (190 + nodes).

1036 We made several other refinements that improved the
1037 quality of the hierarchy only marginally, if at all. Ways to
1038 improve the results include using larger data sets in general,
1039 but also using a document collection that covers a broad
1040 spectrum of conceptual depth and precision. For the
1041 purposes of the retrieval experiment, the automatically
1042 generated hierarchy was used as a flat set of terms, since we
1043 could not rely on the quality of relationships.

1046 4. Automatic indexing from controlled vocabulary

1048 4.1. The algorithm

1050 Traditionally, controlled-vocabulary indexing is done
1051 manually, which is a labor-intensive task. We attempted to
1052 automate it, at the cost of losing *some*, but hopefully not *all*
1053 of the advantages of controlled vocabulary indexing. Simply
1054 put, our approach works as follows: a document D is
1055 assigned a term $T = w_1w_2 \dots w_n$ if it contains (most of) its
1056 component words, consecutively ('... $w_1w_2 \dots w_n$...'), or in
1057 close proximity ('... $w_1n_1n_2w_2w_3 \dots w_n$...'). In our
1058 implementation, we reduced the words of both the terms
1059 of the vocabulary and the documents to their word stem by
1060 removing suffixes and word endings. Also, we used two
1061 tunable parameters for indexing, (1) proximity, and (2)
1062 threshold for the fraction of the number of words found in a
1063 document, to the total number of words of a term; a term
1064 was assigned if that fraction is above the threshold. Assume

1065 that the vocabulary contains the term (key phrase) Database
1066 Management Systems. A threshold of 2/3 would assign the
1067 term to any document that contained two or more words out
1068 of three. The proximity parameter indicates how many
1069 words apart should words appear to be considered part of the
1070 same noun phrase (term). Maarek et al. had found that five
1071 worked well for two-word phrases in English [16]. It has
1072 been our experience that indexing works best when both
1073 parameters depend on the size of the term. A threshold that
1074 is an increasing function of the number of words in a term
1075 seems to yield a balanced mix of short and long terms, with
1076 reasonably few false-positive assignments. Similarly, what
1077 seems to work best for proximity is to use an m -word
1078 distance between any two neighboring words, but a smaller
1079 overall spread than $n \times m$, where n is the number of words in
1080 the term.

1081 At first glance, this approach seems to suffer from similar
1082 problems to automatic plain-text indexing because of its
1083 potential for false positives—still matching words regard-
1084 less of semantic context—and false negatives—still relying
1085 on the terminology used by technical writers or developers.
1086 We felt, however, that because we are dealing mostly with
1087 compound terms, the proximity and threshold parameters
1088 provide both some context for the matching, thereby
1089 reducing the chances of false positives, and some flexibility
1090 in matching, thereby reducing the chances of false
1091 negatives. Further, notwithstanding the quality of indexing,
1092 the fact that searchers are constrained to use the same
1093 vocabulary that was used for indexing can eliminate a good
1094 many sources of retrieval errors.

1097 4.2. Results

1099 The indexing algorithm was used to index the Descrip-
1100 tion attribute of the library components. In particular, we
1101 used a threshold of 2/3 and a proximity of 5, meaning that
1102 we assign a term when at least two thirds of the words of the
1103 term occurred in the textual part of the attribute, with no two
1104 words more than five words apart. The results of the
1105 indexing were somewhat difficult to analyze directly
1106 because the quality of indexing is related as much to the
1107 quality of the vocabulary as it is to the indexing algorithm.
1108 For example, we know that names of classes, methods, or
1109 variables should not have been included in the indexing
1110 vocabulary in the first place—about 26 terms. Another
1111 factor came into play: terms that make sense in the context
1112 of other terms, make little sense when taken alone. For
1113 example, the hierarchy contained the path 'Size' → '
1114 Allocation' → 'Block of Memory', and one intuitively
1115 reads Size as Size of Block of Memory or Size of Allocation
1116 of Block of Memory. Suppose, however, that the term Size
1117 alone were assigned to the description of a component; it
1118 means very little in this context. This problem is not unique
1119 to the automatically generated hierarchy: the ACM
1120 Computing Reviews classification structure has several

instances of nodes which should be ‘read’ in conjunction with their ancestors to be meaningful.⁶

In order to separate the issue of vocabulary control from the performance of automatic indexing per se, we indexed the in-line textual documentation of classes with the ApplicationDomain vocabulary. While we did not expect to find the same term assignment as the manual indexing, we wanted to get an idea about ‘how often’ terminology issues miss some important term assignments, and about the appropriateness of the indexing parameters (threshold and maximum word distance). Our evaluation takes into account what is in the vocabulary, and what is in the text, and the question was, given the same limited vocabulary and limited textual description, would a human being have done it any differently?

We studied 80 textual descriptions ranging in size from a single sentence such as ‘Do not define an implementation for this’, to half a page of text. The results are summarized in the table below.

	Exact	Related	Extraneous terms	Missing because termin. differences	Missing because words missing
Number of terms	42	3	6	11	24
Percentage among assigned	82	6	12		
Coverage	52	4		14	30

The extraneous terms are terms that should not have been assigned (false-positive). Examples include the indexer mistaking the verb [this method] ‘sets’ for the word ‘Sets’ (as in collections). Some of these cases can be resolved if we combine word matching with part-of-speech tag matching so that names match names, and verbs match verbs. Other examples of extraneous terms include a case where the indexer assigned the term ‘Copying Strings’ to the sentence ‘This class does not make copies of the character strings it is given...’.

The missing terms are terms that a human indexer would have assigned if they had the same text, and fall into two categories, (a) a synonym for the actual word(s) was used instead of the actual words, or (b) the concept does not appear ‘verbally’ altogether, but is implicit. An example of (a) is the use of the word ‘Array’ in the text, and the word ‘Vector’ in the on-line documentation.⁷ Examples of (b)

⁶ Such terms are sometimes called *minor descriptors*, i.e. property names attached to their parent concepts; obviously the property name alone does not mean much as several concepts may share the same property.

⁷ This is an interesting discrepancy because it illustrates a fundamental difficulty in software component retrieval. The on-line documentation rightly focuses on abstractions, and hence used the word ‘Vector’. The in-line documentation (within program code, like javadoc comments) describes the *implementation*. Developers will be querying based on abstractions, and not on implementations. Actually, ideally, we should let them query based on *problems*, altogether, but that is another story.

include the sentence ‘matches upper case character’ missing the term ‘Pattern Matching’ or ‘String comparison’. It also includes a number of cases where a term is a conjunction as in ‘Strings and Symbols’, and only one of the two words appearing in the text, coming short of the 2/3 threshold. This happened quite a few times, and can be easily resolved by tagging conjunctive terms to tell the indexer to assign the whole term if it matches one or the other. This may involve, among other things, rewriting terms such as ‘Information Storage and Retrieval’ as ‘(Information Storage) and (Information Retrieval)’.

In summary, only 6% of the assigned terms were wrong, which should only minimally affect retrieval precision. However, the indexer seems to have missed a significant number of terms (44%), although that number can be reduced using minor refinements. We cannot estimate what the effect of these ‘false-negative’ term assignments will be on retrieval recall. For instance, on any given document or component, the effect of removing an index term on the retrievability of the document or component will depend on the other terms already assigned (are there any, are they related to the removed term), and on the retrieval algorithm used (does it use exact retrieval, does it measure ‘conceptual distance’ between related terms, etc).

5. Retrieval experiments

5.1. Experimental design

We were as concerned with establishing the usefulness of the library tool in a production setting as we were with performing comparisons between the various retrieval methods. It is our belief that such comparisons do not mean much if a developer will not use ANY of the methods in a real production setting. The decision for a developer to use or not use a tool has to do with, (1) his/her estimate of the effort it takes to build the components from scratch [35], (2) the cost of using the library tool, including formulating the queries and looking at the results, and (3) the perceived track record of the tool and the library in terms of either finding the right components, or quickly ‘convincing’ the developer that none could be found that satisfy the query. By contrast, comparative studies between the retrieval methods focus on the retrieval performance, regardless of the cost factors. Further, to obtain a fair and finely detailed comparison, the format of the queries is often restricted in those experiments to reduce the number of variables, to the point that they no longer reflect normal usage of the library.

With these considerations in mind, we made the following choices:

- (1) We only controlled the search method that the users could use to answer each of the queries, without giving a time limit on each query, or a limit on

1233 the number of trials made for each query; we assumed
 1234 that users will stop when they are convinced that they
 1235 have found all that is relevant,
 1236 (2) We logged the actions of the subjects with the
 1237 tool. This provided us with finer experimental data
 1238 without interfering with the subjects' workflow.
 1239

1240 By giving users this much freedom, we run the risk that
 1241 user bias will skew the data in one direction, preventing us
 1242 from performing reliable analyses. For example, with
 1243 boolean retrieval, subjects could search on two search
 1244 attributes, separately or in combination. Recall that one
 1245 attribute, Application Domain, was indexed manually with a
 1246 manually built vocabulary, while the other, Description, was
 1247 indexed automatically with the automatically generated
 1248 hierarchy (see Section 4). We did not ask the subjects to use
 1249 one or the other, or both in combination. When we studied
 1250 the traces, it turned out that the Description attribute was
 1251 used only twice out of a possible 43 keyword queries, and
 1252 neither query returned a relevant document, which makes
 1253 any formal comparison of the two attributes impossible.
 1254 However the fact that the Description attribute was used
 1255 only twice tells us that subjects did not feel it provided
 1256 useful information, and that, in and of itself, is a valuable
 1257 data.
 1258

1259 The experimental data set consisted of about 200 classes
 1260 and 2000 methods from the OSE library. We used 11
 1261 queries, whose format is discussed Section 5.2. Seven
 1262 subjects participated in the experiment, although only the
 1263 data from 5 subjects was usable. All subjects were
 1264 experienced C++ programmers. They included two
 1265 professors, three graduate students, and two professional
 1266 developers working for the industrial partners of the project.
 1267 The subjects were given a questionnaire which included the
 1268 statements of the queries, and blank spaces to enter the
 1269 answer as a list of component names. For each of the initial
 1270 77 (subject,query) pairs, we randomly assigned a search
 1271 method (keyword-based versus plain text). For each
 1272 (subject,query,search method) triplet, the subject could
 1273 issue as many search statements as s/he wishes using the
 1274 designated search, with no limitation on the time or on the
 1275 number of search statements. The experiment started with a
 1276 general presentation of the functionality of the tool set
 1277 (about 45 mn), followed by a hands-on tutorial with the tool
 1278 set (about 1 h), providing the subjects with an understanding
 1279 of the theoretical underpinnings of the functionalities, as
 1280 well as some practical know-how. Before leaving, the
 1281 subjects were asked to fill out a questionnaire to collect their
 1282 qualitative appreciation of the tool set.

1283 In order to analyze the results, we used the query
 1284 questionnaires to compare the subjects' answers to ours,
 1285 which were based on a thorough study of the library's user
 1286 manual and some code inspection, where warranted. The log
 1287 traces provided more detailed information and were used to
 1288 support finer analyses.

5.2. Queries

1289 Information retrieval systems suffer from the difficulty
 1290 users have in translating their needs into searchable queries.
 1291 The issue is one of translating the description of a problem
 1292 (their needs) into a description of the solution (relevant
 1293 documents). With *document* retrieval systems, problems
 1294 may be stated as 'I need to know more about ⟨X⟩', and
 1295 solutions as 'A document that talks about ⟨Y⟩'. For a given
 1296 problem, the challenge is one of making sure that ⟨X⟩ and
 1297 ⟨Y⟩ are the same, and in systems that use controlled
 1298 vocabulary indexing, trained librarians interact with naive
 1299 users to help them use the proper search terms.
 1300

1301 With *software component* retrieval, the gap between
 1302 problem statement (a requirement) and solution description
 1303 (a specification) is not only terminological, but also
 1304 conceptual. In an effort to minimize the effect of the
 1305 expertise of subjects in an application, and their familiarity
 1306 with a given library, controlled experiments in component
 1307 retrieval usually use queries that correspond closely to
 1308 component specifications. This does not reflect normal
 1309 usage for a reusable components library tool. For instance,
 1310 users typically do not know how the solution to their
 1311 problem is structured, and for the case of a C++
 1312 component library, e.g. the answer could be a class, a
 1313 method, a function, or any combination thereof. It has
 1314 generally been observed that developers need to know the
 1315 underlying structure or architecture of a library to search for
 1316 components effectively [22]. Accordingly, in an effort to get
 1317 a realistic experiment, we formulated our queries as
 1318 problems to be solved. Each query was preceded by a
 1319 problem description setting up the context, followed by a
 1320 statement 'Find a way of ⟨performing a given task⟩'. The
 1321 problem description is also used to familiarize the subjects
 1322 with the terminology of the application domain using
 1323 textbook-like language.
 1324

5.3. Component relevance: a performance-based evaluation

1325 The difference between traditional bibliographic docu-
 1326 ment retrieval and reusable component retrieval manifests
 1327 itself in the retrieval evaluation process as well. The concept
 1328 of relevance, which serves as the basis for recall and
 1329 precision measures, is notoriously difficult to define. With
 1330 bibliographic document retrieval, a search query for a
 1331 concept X is understood as meaning 'I want documents that
 1332 talk about X', and hence, a document is relevant if it 'talks
 1333 about' X. This definition is different from pertinence which
 1334 reflects a document's usefulness to the user [30]. The
 1335 usefulness of a document to the user depends, among other
 1336 things, on the user's prior knowledge, or on the pertinence of
 1337 the other documents shown to them. Recall, which measures
 1338 the number of relevant documents returned by a query to the
 1339 total number of relevant documents in the document set,
 1340 implicitly assumes that all the relevant documents are
 1341 equally pertinent and irreplaceable: the user needs all of
 1342

1345 them. In other words, with traditional document retrieval,
 1346 assuming that a query Q has N relevant documents, and
 1347 retrieved a set of documents $S = \{D_1, \dots, D_m\}$, we can define
 1348 *pertinence*, and recall as follows:

$$1349 \text{PERT}(D_i) = \begin{cases} \frac{1}{N}, & \text{if } D_i \text{ is relevant} \\ 0, & \text{if } D_i \text{ is not relevant} \end{cases} \quad \text{and } \text{PERT}(S) \\ 1350 \\ 1351 \\ 1352 \\ 1353 \\ 1354 = \text{RECALL}(S) = \sum_{j=1}^m \text{PERT}(D_j) \\ 1355$$

1356 With software component retrieval, the notions of pertinence
 1357 (usefulness) and substitutability are much easier to define as
 1358 both relate to a developer’s ability to solve a problem with
 1359 the components at hand. Symbolically, we view query as a
 1360 requirement Q , which may be satisfied by several, possibly
 1361 overlapping, sets of components S_1, \dots, S_k , where $S_i =$
 1362 $\{D_{i_1}, D_{i_2}, \dots, D_{i_k}\}$. As a first approximation, we define as
 1363 follows:
 1364

$$1365 \text{PERT}(S_i) = \text{PERT}(D_{i_1}, D_{i_2}, \dots, D_{i_k}) = \sum_{j=1}^{k_i} \text{PERT}(D_{i_j}/S_i) = 1 \\ 1366 \\ 1367 \\ 1368 \quad (2)$$

1369 where $\text{PERT}(D/S_i)$ is the usefulness or pertinence of the
 1370 component D in the context of the solution set S_i . This
 1371 illustrates the fact that a retrieved component D is useful
 1372 ‘only if’ the other components required to build a solution
 1373 are retrieved with it. Further, this definition of PERT means
 1374 that total user satisfaction can be achieved with a subset of
 1375 the set of relevant components, which is not the case for
 1376 recall. We illustrate the properties of PERT through an
 1377 example.

1378 Consider two solutions sets $S_1 = \{D_1, D_2\}$ and $S_2 =$
 1379 $\{D_1, D_3, D_4\}$ and assume that D_1, D_2 and D_3 have the sizes
 1380 30, 20, 40, and 30, respectively, giving S_1 and S_2 the sizes
 1381 50, and 100, respectively. We can use the relative sizes of
 1382 the components with respect to the enclosing solution as
 1383 their contextual/conditional pertinence, i.e. $\text{PERT}(D_i/S_j) =$
 1384 $\text{size}(D_i)/\text{size}(S_j)$. In this case $\text{PERT}(D_1/S_1) = 0.6$, PERT
 1385 $(D_2/S_1) = 0.4$, $\text{PERT}(D_1/S_2) = 0.3$, $\text{PERT}(D_3/S_2) = 0.4$,
 1386 and $\text{PERT}(D_4/S_2) = 0.3$. Assume that a query retrieves the
 1387 component D_1 . In this case, $\text{PERT}(D_1) = \text{Max}(\text{PERT}$
 1388 $(D_1/S_1), \text{PERT}(D_1/S_2)) = 0.6$. If the query retrieved D_1
 1389 and D_3 , instead, $\text{PERT}(\{D_1, D_3\}) = \text{Max}(\text{PERT}(\{D_1, D_3\}/$
 1390 $S_1), \text{PERT}(\{D_1, D_3\}/S_2)) = \text{Max}(\text{PERT}(D_1/S_1) + \text{PERT}$
 1391 $(D_3/S_1), \text{PERT}(D_1/S_2) + \text{PERT}(D_3/S_2)) = \text{Max}(0.6 + 0.0,$
 1392 $0.3 + 0.4) = 0.7$. This illustrates the fact that when several
 1393 partial solutions are returned by the system, we take into
 1394 account the one that is most complete, and the value of
 1395 individual components is relative to that solution. Symbolically,
 1396 given the solution sets S_i, \dots, S_k , a query that returns
 1397 a set of components S has the pertinence:
 1398

$$1399 \text{PERT}(S) = \text{Max}_{j=1, \dots, k} \text{PERT}(S \cap S_j/S_j) \quad (3) \\ 1400$$

1401 Finally, we add another refinement which takes into account
 1402 the overlap of two components within the same solution set.
 1403 Consider the solution S_1 above, and assume that the system
 1404 retrieves D_1 and D'_2 , where D'_2 is a *superclass* of D_2 that
 1405 implements only part of the functionality required of D_2 . In
 1406 this case, we could take $\text{PERT}(D_1/D'_2) = 0.6 + 0.3 = 0.9$. If
 1407 the query retrieved D'_2 AND D_2 , then we discard the weaker
 1408 component. This is similar to viewing solutions sets as role
 1409 fillers and, for each role, take the component that most
 1410 closely matches the role. Within the context of reusable OO
 1411 components, *roles* may be seen as class interfaces, and role
 1412 fillers as class implementations.

1413 For our experiments, some of the 11 queries were
 1414 straightforward in the sense that there was a single
 1415 component (a method or a class) that answered the query,
 1416 and both component relevance and recall were straightfor-
 1417 ward to compute. Queries whose answers involved several
 1418 classes collaborating together (e.g. an object framework)
 1419 were more complex to evaluate and involved all of the
 1420 refinements discussed above.

1421 For the case of precision, we used the traditional
 1422 measure, i.e. the ratio of the retrieved components that
 1423 were relevant (i.e. had a non-zero $\text{PERT}(\cdot)$) to the total
 1424 number of retrieved components. We can also imagine
 1425 refining the definition of precision to take into account the
 1426 effective usefulness of the individual components, and
 1427 factor that in with the cost of retrieving and examining a
 1428 useless component. The cost of examining a useless
 1429 component is a function of its complexity, and size could
 1430 be used as a very first approximation of that complexity.
 1431

1432 5.4. Performance results

1433 Table 4 shows recall and precision for the 11 queries.

1434 For each query, we randomly selected three subjects out
 1435 of the initial seven to perform the query using full-text
 1436 retrieval, and the remaining four subjects to perform
 1437 keyword retrieval, or vice versa, while making sure that
 1438

1439 Table 4
 1440 Summary of retrieval results

Query	Full-text retrieval			Keyword retrieval		
	Subjects	% Recall	% Precision	Subjects	% Recall	% Precision
1	3	100	88.666	2	50	50
2	4	50	100	1	50	100
3	1	100	100	4	100	100
4	1	100	80	4	50	100
5	4	25	12.5	1	0	0
6	3	33.333	33.333	2	12.5	25
7	2	65	75	3	66.333	50
8	2	30	75	3	30	83.333
9	3	53.333	100	2	30	78
10	3	78.333	80.333	1	35	100
Average (26)		63.49	74.47	(23)	42.41	68.33

1457 each subject had a balanced load of full-text and keyword
 1458 queries (6 and 5, respectively, or vice-versa). Because the
 1459 results of two subjects could not be used, we ended up with
 1460 some queries answered by four subjects using full-text
 1461 retrieval, say, and only once using keyword retrieval (see
 1462 e.g. query 2). The 11th query was rejected because the three
 1463 keyword-based answers were all rejected for one reason or
 1464 another. Hence, comparisons between the two methods for
 1465 the individual queries are not reliable.

1466 At first glance, it appears that plain-text retrieval yielded
 1467 significantly better recall and somewhat better precision. It
 1468 also appears that it has done consistently so for the 10
 1469 queries, with a couple of exception. In order to validate
 1470 these two results statistically, we have to ascertain that none
 1471 of this happened by chance. We performed a number of
 1472 ANOVA tests, to check whether recall and precision were
 1473 random variables of the pair (query, search method), and
 1474 both tests were rejected. Next, we isolated the effect of the
 1475 search type to see if the difference in recall and precision
 1476 performance is significant. The results are shown in Table 5.

1477 The ‘Pr > F’ shows the probability that such a difference
 1478 in performance could have been obtained by chance. It is
 1479 generally accepted that a threshold of 5 percent is required
 1480 to affirm that the differences are significant. Thus, we
 1481 conclude that:

- 1483 • Full-text retrieval yields provably/significantly better
 1484 recall than controlled vocabulary-based retrieval
- 1485 • Full-text retrieval yields comparable precision per-
 1486 formance to that of controlled vocabulary-based
 1487 retrieval.

1489 Our results seem to run counter to the available
 1490 experimental evidence. Document retrieval experiments
 1491 have consistently shown that controlled vocabulary-based
 1492 indexing and retrieval yielded better recall and precision
 1493 than plain-text search [4,30,31], although the difference was
 1494 judged by many as being too small to justify the extra costs
 1495 involved in controlled vocabulary-based indexing and
 1496 retrieval [31]. Similarly, a comparative retrieval experiment
 1497 for reusable components conducted by Frakes and Pole⁸ at
 1498 the SPC showed that recall values were comparable, and a
 1499 superior precision for controlled vocabulary-based retrieval
 1500 [8]. Most surprising in our results is the significant
 1501 difference is recall performance. We analyze these results
 1502 in more detail below.

1503 To explain these results, we formulated and tested a
 1504 number of hypotheses. We first note that out of the 11
 1505 queries, some were supposed to retrieve single components

1507 ⁸ Frakes and Pole compared four methods, and their test of statistical
 1508 significance was based on variance analysis of the precision averages for
 1509 the four methods, which was inconclusive [8]. However, we are quasi-
 1510 certain that by performing pairwise comparison between plain-text search
 1511 (50%) and controlled vocabulary search (what appears to be 100% on the
 1512 plot [8]), they would have established, statistically, the superiority of
 controlled-vocabulary retrieval.

Table 5
Significance of differences between plain-text retrieval and keyword retrieval

Effect of search method	Recall	Precision
F value	4.1	0.93
Pr > F	0.0500	0.3404

(often methods), as in Query 7, formulated as ‘getting the length of a string’, and the others were supposed to retrieve a collection of components with complex interactions, often a mix of classes and methods. With full-text search, queries retrieve indiscriminately methods and classes. With controlled-vocabulary search, users have to instantiate different query templates, depending on the kind of components they are seeking (a class or a method). We hypothesize that this makes the search more tedious and users may give up search easily, yielding lower recall. For this hypothesis to hold, there has to be a marked difference between the performance for the single-component queries (queries 1, 7, 8, 9) and the queries whose answers consisted of collections of components (queries 2, 3, 4, 5, 6, 10). Table 6 compares the two kinds of queries.

Our hypothesis that plain-text retrieval favors component collection queries is not validated. Along the same lines, we hypothesized that plain-text retrieval favored queries whose answers involved a mix of methods and classes, or just classes, since the same query would retrieve both kinds of components. Table 7 shows recall and precision values for the two retrieval methods, separated into the two kinds of queries.

This hypothesis is not validated: in both cases, plain-text retrieval is markedly superior to controlled-vocabulary retrieval with regard to recall—and somewhat with regard to precision for the case of queries whose answers included both classes and methods. Note, however, that there is a marked difference in performance between the two groups of queries.

Could the quality of indexing be to blame for the lower performance of controlled-vocabulary based retrieval? Recall that we indexed two attributes, Application Domain and Description. The Application Domain attribute was indexed manually and fairly systematically, thanks to the quality of on-line documentation. There are two potential weakness of this indexing, but none can account for the observed difference in performance

Table 6
Comparing the two sets of queries

Query set	Full-text retrieval		Keyword retrieval	
	% Recall	% Precision	% Recall	% Precision
Single comp. queries	62.08	84.67	44.08	65.333
Comp. coll. queries	64.43	68.17	41.30	70.33

1569 Table 7
1570 Comparing the two sets of queries depending on whether they retrieve
1571 methods or not

1572 Query set	1573 Full-text retrieval		1574 Keyword retrieval	
	1575 % Recall	1576 % Precision	1577 % Recall	1578 % Precision
1579 Answer = classes and 1580 methods	41.333	59.17	27.77	47.27
1581 Answer = classes only	85.65	89.77	57.05	89.40

1582 between the two retrieval methods. First, some methods
1583 were left unindexed because the on-line documentation
1584 said nothing about these methods—such as constructors.
1585 The results of Table 7 bear this out: keyword retrieval
1586 performed better on queries whose answers involved only
1587 classes than on queries whose answers involved a
1588 combination of classes of methods. However, this does
1589 not explain the fact that free-text retrieval performed
1590 better than keyword retrieval *for both types of queries*.
1591 The second potential weakness of the indexing of the
1592 ApplicationDomain is the fact that index terms are
1593 sometimes perceived as too general. Indexing that is too
1594 general results in poor precision, but is known to produce
1595 *better* recall, which is not what we observed.

1596 The attribute Description was indexed automatically (see
1597 Section 4) with the vocabulary that was generated
1598 automatically (see Section 3). Notwithstanding the quality
1599 of indexing of this attribute, the experiment logs showed
1600 that this attribute was actually used only three times, and in
1601 all three cases, it was used in conjunction with ApplicationDomain,
1602 but failed to match any component. Accordingly, even in those cases
1603 where it was used, it did not affect the ranking of components
1604 returned by weighted boolean retrieval (see Section 2.3.2). The fact
1605 that the attribute Description was not used as often as ApplicationDomain
1606 could be explained by the nature of queries: the queries were
1607 presented as programming problems or tasks to solve, rather than
1608 a look-up for components given a set of specifications. Because
1609 ApplicationDomain talks about problems that components help solve
1610 whereas Description talks about how these components are implemented,
1611 it makes sense that the former be used more often than the latter
1612 in the queries.

1613 We continue to analyze the results of this experiment, as
1614 the logs provide us with a wealth of information and hypotheses
1615 that we could validate. We do not expect this experiment to reverse
1616 the long-held consensus that controlled vocabulary performs better
1617 than free-text retrieval; more experiments that target narrower
1618 retrieval tasks, and that involve fewer operational parameters
1619 would be needed for that. We can view it in light of another
1620 emerging consensus according to which, whichever performance
1621 benefits controlled vocabulary indexing and retrieval might have—
1622 in our case none, quite the contrary—they hardly justify the
1623 added cost. Perhaps more importantly, four subjects out of five
1624 preferred plain-text search.

1625 More importantly, we believe that this experiment
1626 contributes to a needed rethinking of reusable component
1627 retrieval paradigms and tools. Such implications are
1628 discussed next.

1630 6. Conclusion and directions

1631 We set out to develop, evaluate, and compare two classes
1632 of component retrieval methods which, supposedly, strike
1633 different balances along the costs/benefits spectrum,
1634 namely, the (quasi-) zero-investment free text classification
1635 and retrieval versus the ‘up-front investment-laden’ but
1636 presumably superior controlled vocabulary faceted indexing
1637 and retrieval. Recent experiments with software component
1638 repositories have put into question the *cost-effectiveness* of
1639 the controlled vocabulary approach, but not its *superior* or
1640 *at least equally good* retrieval performance [8]. We
1641 attempted to bring the two kinds of methods to a level-
1642 playing field by: (1) addressing the costs issue by
1643 automating as much as possible of the pre-processing
1644 involved in controlled vocabulary-based methods, and (2)
1645 using a realistic experimental setting and realistic evaluation
1646 measures. Our experiments showed that: (1) those aspects of
1647 the pre-processing involved in controlled vocabulary
1648 methods that we automated were of poor enough quality
1649 that they were not used (the Description attribute), and (2)
1650 the fully automatic free text search performed better than the
1651 fully manual controlled-vocabulary based indexing and
1652 retrieval of components.

1653 Because these results are somewhat counter-intuitive, we
1654 continue to analyze them, along with the log data, and to
1655 design new experiments that are better targeted towards
1656 validating the various hypotheses discussed in Section 5.4.
1657 However, they give legitimacy and some urgency to some of
1658 the questions we and others have raised about the retrieval
1659 of reusable software components [22,23,36].

1660 From an organizational issues point of view, there was
1661 wide recognition in the late eighties that reuse will not
1662 happen at a large scale within organizations without the
1663 proper structuring and management. It was possible, in that
1664 context, to conceive of centralized reuse repositories with
1665 well-defined roles and quality control criteria and mechanisms
1666 [25]. Nowadays, a lot more reuse happens in the
1667 unstructured and decentralized world of the Internet and
1668 open source software, and any ‘virtual reuse repository’ can
1669 only rely on *automated indexing and retrieval methods*,
1670 *regardless of differences in performance*.

1671 Reuse repositories are also facing a number of *paradigmatic*
1672 issues. First, there exist qualitative differences
1673 between bibliographic document retrieval and software
1674 component retrieval [22], which make some of the
1675 document retrieval analogies inappropriate. Document
1676 library users who do not find the documents they are
1677 looking for will look even harder because they cannot
1678 perform the tasks for which they needed the information
1679

1681 otherwise. A software developer will more easily give up
 1682 and get on with developing the software component from
 1683 scratch. As reuse repository designers, we need to account
 1684 for the fact that software developers are not our captive
 1685 users, which puts more pressure on us to provide more
 1686 useful and less intrusive tools. It is important that the use of
 1687 the repository integrates well into the workflow of
 1688 developers; this has led some people to suggest that reuse
 1689 repositories should be active in the sense of presenting
 1690 potentially relevant information to users before they ask for
 1691 it [36]. It also means that issues of usability are paramount;
 1692 if developers prefer a particular search method, then that is
 1693 the one we should focus on. Our tool set does not address
 1694 this issue specifically, but we take seriously the fact that four
 1695 out of five users preferred free text search, which confirms
 1696 earlier studies. In our case, it even performed better.

1697 Surely, our experiments suggest that there is ample room
 1698 for improvement in several areas (see Sections 3.2.3, 4.2,
 1699 and 5.4). However, we believe that there is something more
 1700 fundamental at play. We believe that multi-faceted
 1701 classification and retrieval of reusable components to be at
 1702 the wrong level of formality for the typical workflow of
 1703 developers using a library of reusable components. We
 1704 identify two very distinct search stages. The first stage
 1705 coincides with analysis, and is fairly exploratory, as
 1706 developers do not yet know which form (specification?)
 1707 the solution to their problem will take. During this stage, a
 1708 free-format search technique such as plain-text search is
 1709 appropriate, as multi-faceted search may be too rigid and
 1710 constraining. After contemplating several designs, a devel-
 1711 oper may then start searching for components that would
 1712 play a given role within a design, and multi-faceted
 1713 classification may be too poor for this stage. The format
 1714 of our queries (problems to be solved), and the fact that
 1715 experimental subjects used mostly the ApplicationDomain
 1716 attribute, setting aside the more implementation-oriented
 1717 Description attribute seem to point in this direction. A
 1718 combination of free-text search and active reuse repositories
 1719 [36] may be worth exploring.

1722 Acknowledgements

1724 This work was supported by grants from Canada's
 1725 Natural Sciences and Engineering Research Council
 1726 (NSERC), TANDEM Computers, Québec's *Fonds pour la*
 1727 *Création et l'Aide à la Recherche* (FCAR), and Québec's
 1728 *Ministère de l'Enseignement Supérieur et de la Science*
 1729 (MESS) under the IGLOO project organized by the Centre
 1730 de Recherche Informatique de Montréal.

1731 Bertrand Fournier, a statistician with the *Service de*
 1732 *Consultation en Analyse de Données* (SCAD, <http://www.scad.uqam.ca>), and Professor Manzour Ahmad, director of
 1733 SCAD, provided us with invaluable assistance in measuring
 1734 and interpreting the results.
 1735
 1736

References

- 1737
 1738
 1739
 1740
 1741
 1742
 1743
 1744
 1745
 1746
 1747
 1748
 1749
 1750
 1751
 1752
 1753
 1754
 1755
 1756
 1757
 1758
 1759
 1760
 1761
 1762
 1763
 1764
 1765
 1766
 1767
 1768
 1769
 1770
 1771
 1772
 1773
 1774
 1775
 1776
 1777
 1778
 1779
 1780
 1781
 1782
 1783
 1784
 1785
 1786
 1787
 1788
 1789
 1790
 1791
 1792
- [1] ACM, An introduction to the CR classification system, *Computing Reviews* January (1985) 45–57.
 - [2] P. Allen, Reuse in the component marketplace, *Component Development Strategies* 11 (8) (2001).
 - [3] H. Bilofsky, C. Burks, J.W. Fickett, W.B. Goad, F.I. Lewitter, W.P. Rindone, C.D. Swindell, C. Tung, The GenBank genetic sequence databank, *Nucleic Acids Research* 14 (1986) 1–4.
 - [4] D. Blair, M.E. Maron, An evaluation of retrieval effectiveness for a full-text document-retrieval system, *Communications of the Association for Computing Machinery* 28 (3) (1985) 289–299.
 - [5] D. Cutting, J. Kupiec, J. Pedersen, P. Sibun, A practical part-of-speech tagger, *Proceedings of the Applied Natural Language Processing Conference* (1992).
 - [6] E. Damiani, M.G. Fugini, C. Belletini, A hierarchy-aware approach to faceted classification of object-oriented components, *ACM Transactions on Software Engineering and Methodology* 8 (3) (1999) 215–262.
 - [7] G. Dumpleton, *OSE—C++ Library User Guide*, Dumpleton Software Consulting Pty Limited, Parramatta, 2124, New South Wales, Australia, 1994, 124 pp.
 - [8] W.B. Frakes, T. Pole, An empirical study of representation methods for reusable software components, *IEEE Transactions on Software Engineering* August (1994) 1–23.
 - [9] W.B. Frakes, R. Baeza-Yates, *Information Retrieval: Data Structures and Algorithms*, Prentice-Hall, Englewood Cliffs, NJ, 1994.
 - [10] R.J. Hall, Generalized Behavior-based Retrieval, *Proceedings of the 15th International Conference on Software Engineering*, ACM Press, Baltimore, MD, 1993, pp. 371–380.
 - [11] S. Henninger, Using iterative refinement to find reusable software, *IEEE Software* 11 (5) (1994) 48–59.
 - [12] S. Isoda, Experience report on a software reuse project: its structure, activities, and statistical results, *Proceedings of the 14th International Conference on Software Engineering*, Melbourne, Australia May (1992) 320–326.
 - [13] I. Jacobson, M. Griss, P. Jonsson, *Software Reuse: Architecture, Process and Organization for Business Success*, Addison-Wesley, Reading, MA, 1997.
 - [14] K.S. Jones, A statistical interpretation of term specificity and its application in retrieval, in: B.C. Griffith (Ed.), *Key Papers in Information Science*, Knowledge Industry Publications, Inc, White Plains, NY, 1980, pp. 305–315.
 - [15] M.E. Lesk, Word–word associations in document retrieval systems, *American Documentation* 20 (1) (1969) 27–38.
 - [16] Y.S. Maarek, D.M. Berry, G.E. Kaiser, An information retrieval approach for automatically constructing software libraries, *IEEE Transactions on Software Engineering* 17 (8) (1991) 800–813.
 - [17] A. Mili, R. Mili, R. Mittermeir, Storing and retrieving software components: a refinement-based approach, *Proceedings of the 16th International Conference on Software Engineering*, Sorrento, Italy May (1994).
 - [18] H. Mili, R. Rada, Building a knowledge base for information retrieval, *Proceedings of the Third Annual Expert Systems in Government Conference* October (1987) 12–18.
 - [19] H. Mili, R. Rada, Merging Thesauri: principles and evaluation, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 10 (2) (1988) 204–220.
 - [20] H. Mili, R. Rada, Medical experttext as regularity in semantic nets, *Artificial Intelligence in Medicine* 2 (1990) 217–229. Elsevier Science Publishers.
 - [21] H. Mili, R. Rada, W. Wang, K. Strickland, C. Boldyreff, L. Olsen, J. Witt, J. Heger, W. Scherr, P. Elzer, Practitioner and SoftClass: a comparative study of two software reuse research projects, *Journal of Systems and Software* 27 (1994).

- 1793 [22] H. Mili, F. Mili, A. Mili, Reusing software: issues and research
1794 directions, IEEE Transactions on Software Engineering 21 (6) (1995)
1795 528–562.
- 1796 [23] H. Mili, E. Ah-Ki, R. Godin, H. McHeick, Another nail to the coffin of
1797 faceted controlled-vocabulary component classification and retrieval,
1798 Proceedings of the '97 Symposium on Software Reuse, Boston, MA
1799 May (1997) 89–98.
- 1800 [24] H. Mili, H. Sahraoui, Describing and using frameworks, in: R.E.
1801 Johnson (Ed.), Building Application Frameworks: Object-oriented
1802 Foundations of Framework Design, Wiley, New York, 1999, pp.
1803 523–561.
- 1804 [25] H. Mili, A. Mili, S. Yacoub, E. Addy, Reuse-based Software
1805 Engineering: Techniques, Organization, and Control, Wiley, New
1806 York, 2002, ISBN 0-471-39819-5.
- 1807 [26] OASIS, Business Process, Business Information Analysis Overview
1808 (ebXML), Organization for the Advancement of Structured Infor-
1809 mation Standards, May 11, 2001, [http://www.ebxml.org/specs/
1810 bpOVER.pdf](http://www.ebxml.org/specs/bpOVER.pdf).
- 1811 [27] E. Ostertag, J. Hendler, R. Prieto-Diaz, C. Braun, Computing
1812 similarity in a reuse library system: an AI-based approach, ACM
1813 Transactions on Software Engineering and Methodology 1 (3) (1992)
1814 205–228.
- 1815 [28] R. Prieto-Diaz, P. Freeman, Classifying software for reusability, IEEE
1816 Software January (1987) 6–16.
- 1817 [29] R. Rada, H. Mili, E. Bicknell, M. Blettner, Development and
1818 application of a metric on semantic nets, IEEE Transactions on
1819 Systems, Man, and Cybernetics 19 (1) (1989) 17–30.
- 1820 [30] G. Salton, M. McGill, Introduction to Modern Information Retrieval, 1849
1850 McGraw-Hill, New York, 1983.
- 1851 [31] G. Salton, Another look at automatic text-retrieval systems,
1852 Communications of the Association of Computing Machinery 29
1853 (7) (1986) 648–656.
- 1854 [32] C. Smith, MEDLINE Queries and Distances in MeSH, Internal
1855 Report, National Library of Medicine, 1985.
- 1856 [33] D. Soergel, Organizing Information: Principles of Data Base and
1857 Retrieval Systems, Academic Press, Orlando, FL, 1985.
- 1858 [34] B.H. Weinberg, J.A. Cunningham, The relationship between term
1859 specificity in MeSH and online postings in MEDLINE, Bulletin
1860 Medical Library Association 73 (4) (1985) 365–372.
- 1861 [35] S.N. Woodfield, D.W. Embley, D.T. Scott, Can programmers reuse
1862 software, IEEE Software July (1987) 52–59.
- 1863 [36] Y. Ye, G. Fischer, Promoting Reuse with Active Reuse Repository
1864 Systems, Proceedings of the Sixth International Conference on
1865 Software Reuse, Lecture Notes in Computer Science, vol. 1844,
1866 Springer, Berlin, 2000, pp. 302–317.
- 1867 [37] A.M. Zaremski, J.M. Wing, Signature matching: a key to reuse,
1868 Software Engineering Notes 18 (5) (1993) 182–190. First ACM
1869 SIGSOFT Symposium on the Foundations of Software Engineering.
- 1870 [38] A.M. Zaremski, J.M. Wing, Specification matching: a key to reuse,
1871 Software Engineering Notes 21 (5) (1995) Third ACM SIGSOFT
1872 Symposium on the Foundations of Software Engineering.
- 1873 [39] G.K. Zipf, The Psycho-Biology of Language, MIT Press, Cambridge,
1874 MA, 1965.
- 1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904