# How the Sando Search Tool Recommends Queries

Xi Ge*, David Shepherd†, Kostadin Damevski‡, Emerson Murphy-Hill*

*NC State University, Raleigh, NC, USA
xge@ncsu.edu, emerson@csc.ncsu.edu
†ABB Inc, Raleigh, NC, USA
david.shepherd@us.abb.com
‡Virginia State University, Petersburg, VA, USA
kdamevski@vsu.edu

*Abstract*—Developers spend a significant amount of time searching their local codebase. To help them search efficiently, researchers have proposed novel tools that apply state-of-the-art information retrieval algorithms to retrieve relevant code snippets from the local codebase. However, these tools still rely on the developer to craft an effective query, which requires that the developer is familiar with the terms contained in the related code snippets. Our empirical data from a state-of-the-art local code search tool, called Sando, suggests that developers are sometimes unacquainted with their local codebase. In order to bridge the gap between developers and their ever-increasing local codebase, in this paper we demonstrate the recommendation techniques integrated in Sando.

## I. Introduction

Software is hard to maintain. One cause for the difficulty lies in the increasing complexity of software systems. To accomplish a software maintenance task, a developer needs to explore the information space of a software system and comprehend all relevant parts of the codebase. This exploration is both tedious and time-consuming. According to Singer et al., developers spend over 40% of their time in navigating, searching and reading source code [11].

Due to the overwhelming complexity of software systems, developers often start maintenance tasks by searching for a starting point in the local codebase. A recent study conducted by Ko et al. indicates that most (9 out of 12) software maintenance tasks start with local code search [10].

To help developers search their codebase, researchers have proposed several local code search tools that are integrated into popular IDEs. For instance, Sando enhances the Visual Studio IDE and InstaSearch is aimed for Eclipse developers [6], [5]. Both of these search tools significantly improve upon the built-in search tools available in these IDEs in the following aspects: (1) they retrieve different levels of software entities, such as fields, methods, and classes, as search results, instead of retrieving lines of text; (2) they rank the search results according to their relevance to the given search query; and (3) these tools support multi-word searches, that is, querying by space-separated terms and retrieving software entities containing any of these terms.

Regardless of their sophistication, all search tools require developers issuing both specific and relevant queries to produce useful results. However, the size of the codebase under search as well as the vocabulary mismatch problem, where a domain concept is expressed differently in the code, often make this requirement unreasonable. Empirical data collected from Sando usage in the field further substantiates this conjecture. Specifically, about 20% of all Sando queries fail, returning no results at all, which suggests developers in the field are facing difficulties in writing effective queries.

To solve this problem, we improve Sando by a set of query recommendation techniques. These techniques fall into two categories: pre-search recommendations and post-search recommendations. The former recommends queries before a developer starts searching and the later recommends queries after a developer's search fails, that is, when the search returns no results. Below, we present how these recommendation techniques assist developers to compose queries in a semi-automatic fashion that are better at finding the relevant code snippets.

## II. Approach

Before detailing the recommendation techniques, we first briefly introduce the local code search tool called Sando, the platform on which our query recommendation techniques are implemented and evaluated. Sando is a state-of-the-art code search tool for Visual Studio developers [5], which up to this point has been downloaded over three thousand times. According to our usage data collected from Sando users, roughly forty developers use Sando on a daily basis. Figure 1 presents the user interface of Sando. Sando supports multiple programming languages including C#, C++, C and XML.

After installing Sando, a developer can issue queries in the search box at A, press the button on the right, and be presented by search results at B. Next, the developer can examine the search results in one of two ways: by single clicking on a search result or by double clicking the search result. A single click triggers a pop-up menu that summarizes the search result, as illustrated by C and D in Figure 1, where C shows the entire software entity, which can be the declaration of a class, a method, or a field; and D shows the lines of code containing the search query; the developer can also double click the search result to open the containing file in the Visual Studio editor.

### A. Components

To recommend queries, our recommendation techniques rely on five different data source components, namely the
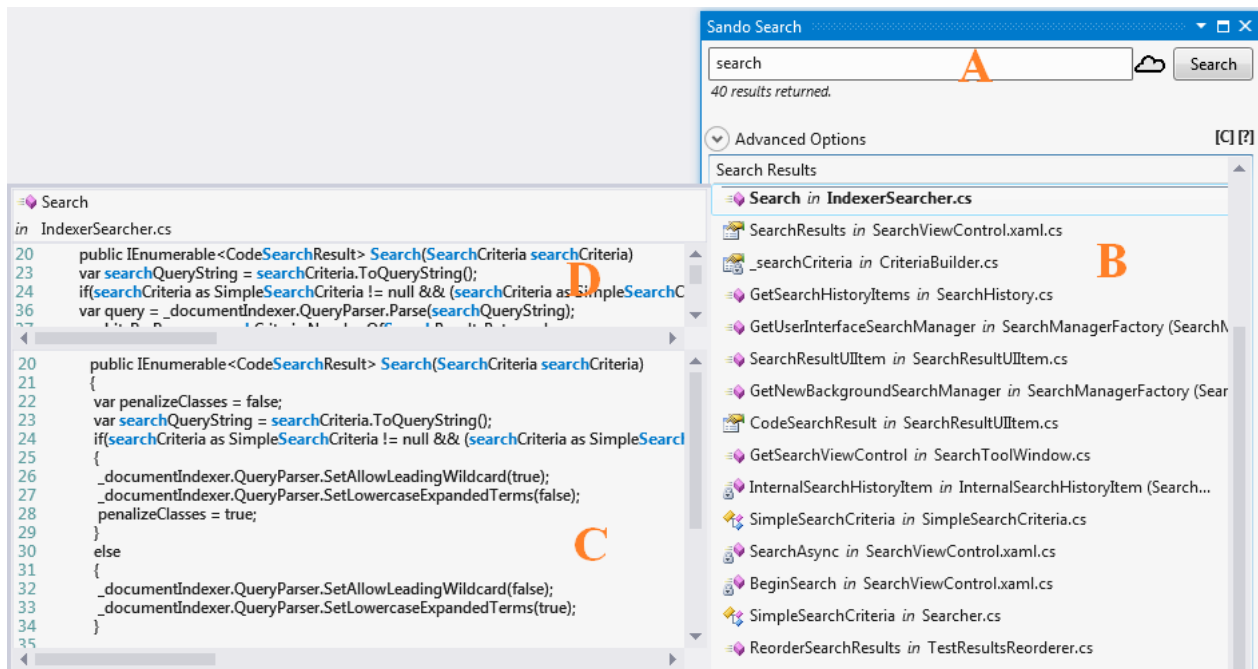
Fig. 1: Sando screenshot after retrieving search results.

local dictionary, the term co-occurrence matrix, the verb-direct-object pairs, the software engineering thesaurus, and the general English thesaurus, which we will explain next.

*1) Local Dictionary:* The first component used by our recommendation techniques is the local dictionary of the codebase under search. Specifically, the local dictionary contains terms that appear in the codebase at least once. To construct the dictionary, we reuse the indexing performed by the Sando search engine. The indexing process breaks each software entity from the codebase into terms; later, after a developer issues a search query, the search engine retrieves the software entities whose indexed terms match with the given query as search results. Taking the following code snippet as an example, the Sando search engine indexes the method `Perform` to the terms of "perform", "output", "func", "invoke", "input", "finished", "event", and "finishedevent". In addition to indexing raw identifiers, Sando indexes the terms that result from splitting the identifiers that are in camel case or underscore delimited format.

```
void Perform()
{
    var output = func.Invoke(input);
    if(FinishedEvent != null)
        FinishedEvent(this, output);
}
```

Sando performs either a entire indexing or an incremental indexing: entire indexing, performed when a developer opens a new project that Sando has no cached index, traverses and collects terms for the entire project; incremental indexing, on the other hand, monitors the changed part of a cached project and re-indexes that part. Indexing the entire project of $10K$ LOC takes about 30 seconds to finish; incrementally indexing

an updated C# file takes about 30 milliseconds.

The local dictionary consists of all the terms collected from indexing without redundancy. To facilitate searching this dictionary, we construct a binary search tree on these terms. Taking the codebase of Sando as an example, the local dictionary contains about two thousand different terms. Finding a given term in the local dictionary costs trivial time to finish. Using the local dictionary, our recommendation technique can ensure that the recommended queries actually appear in the local codebase under search.

*2) Term Co-occurrence Matrix:* In addition to the local dictionary, we maintain the collected terms from the local codebase through a co-occurrence matrix. The matrix has an equal number of rows and columns; each column or row represents one term appearing in the local dictionary. Each element in the matrix saves the count of two terms, as represented by the row and the column that occur together in the codebase. For instance, the element at [red, blue] is the count of occurrences of "red" and "blue" together. By appearing together, we mean these two terms appear in the same software entity. Again taking the method `Perform` in Section II-A1as an example, because the method is an independent software entity, any two of the indexed eight terms appear together.

Literally keeping the full co-occurrence matrix in memory is inefficient due to the observation that the matrix is usually sparse. To improve the memory-efficiency, our technique keep the matrix by using Yale format, a data structure that maintains a sparse matrix using three rows of integers [12].

*3) Verb-direct-object Pairs:* Many code snippets conceptually correspond to performing certain actions on certain objects, such as "open file", "close stream" and "create instance". Based on this observation, Fry et al. proposed a

text mining technique that collects these concepts, or verb-direct-object pairs, from a given codebase [7]. Applying this technique, Sando mines the verb-direct-object pairs from the code base under search, and caches them for recommending either the verb or the object when the developer queries the other part. For instance, supposing Sando mined the verb-direct-object pair "open file" from the codebase under search. If the developer queries "open", Sando recommends "file" to the search query. Similarly, if the developer queries "file", Sando recommends "open" to the original query.

*4) Software Engineering Thesaurus:* Sando's recommendation techniques also tries to solve the vocabulary mismatch problem, which happens when a developer queries terms that do not appear in the codebase, however semantically relate to some terms that do. For instance, if the codebase under search uses "retrieval" consistently, then the developer's querying of "search" are unlikely to return useful results. To solve this problem, we applied a thesaurus-based technique. Taking a term as input, the thesauri return synonyms that appear in the codebase as recommended queries. Only using general English thesaurus is not enough because software development has developed many field-specific synonyms, such as "instantiate" to "create" and "update" to "refresh". Therefore, we apply the work of Gupta et al. that mined the source code of open software projects, generating 1724 pairs of related terms. Among these pairs, about 91% are field-specific [8].
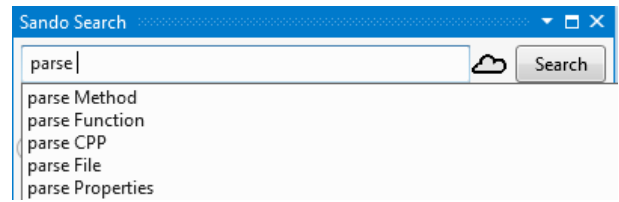
*5) General English Thesaurus:* In addition to the field-specific thesaurus, we also include a general English thesaurus to help the developer when the field-specific thesaurus fails. We derive the general English thesaurus from WordNet, which is database of English words with the relationship between them [9]. Keeping the entire WordNet in memory is costly; hence we only include the top $100k$ most frequent terms in the English language.
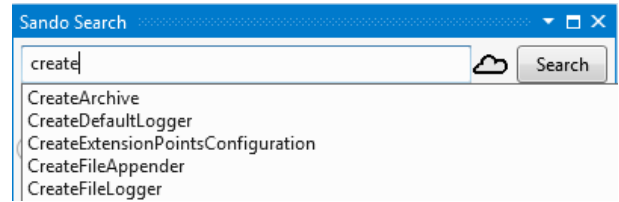
### B. Pre-search Recommendations

Based on the aforementioned components, we next detail how Sando computes and presents the recommended queries. The first category of our recommendation techniques assists developers before they query. These pre-search recommendations allow developers to select terms that are either frequent in the local codebase or closely related to the input. More precisely, the recommended queries originate from the following three sources: (1) the verb-direct-object pairs, (2) the identifiers, and (3) the co-occurring terms.

When the developer inputs a verb or a noun in the search box, a drop-down menu presents the verb-direct-object pairs which contain the term given in the search box. For instance, as illustrate in Figure 2a, when the developer inputs "parse" to the search box, the drop-down menu lists pairs such as "parse file" and "parse method".
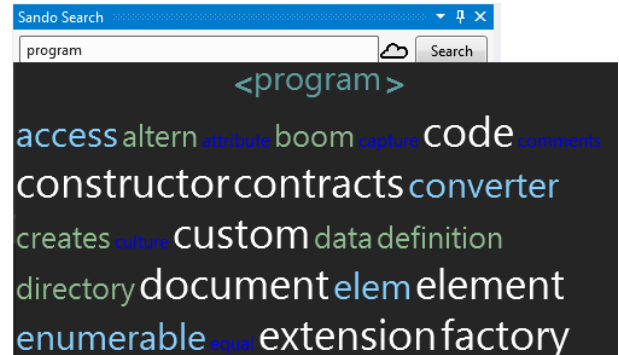
When the developer inputs the prefix of cached identifiers, the drop-down menu presents the identifiers starting with the given prefix. For instance, the drop-down menu in Figure 2b shows identifiers starting with "create", which is the developer's input in the search box. These identifiers could be method names, class names, and field names. Sando retrieves these identifiers from the local dictionary.


(a) Recommending verb-direct-object pairs.


(b) Recommending identifiers.


(c) Frequently co-occurring terms.

Fig. 2: UI of pre-search recommendations.

The third source of the recommended terms is the co-occurrence matrix. After the developer inputs a term in the search box, she can view the terms that co-occur with the term. Different from the aforementioned two sources, Sando presents the co-occurring terms through a tag cloud. The size of a term in the tag cloud indicates the comparative co-occurrence count. The bigger the font size of a term, the more frequent the term co-occurs with the term in the search box.

For instance, after the developer inputs "program" in the search box and clicks the cloud button near the search box, a tag cloud appears as illustrated in Figure 2c. From the tag cloud, the developer can infer that "code" appears with "program" more frequently than "data" does. Furthermore, each term in the tag cloud is a hyper link, the click on which adds the term to the original query. The reason for using the tag cloud instead of the drop-down menu is that the co-occurring terms significantly outnumber the identifiers and the verb-direct-object pairs; presenting the co-occurring terms through the drop-down menu may be inconvenient for the developer to browse.

### C. Post-search Recommendations

The second category of the recommendation techniques helps developers when their manual queries fail, which happens when the queries contain terms not appearing in the local
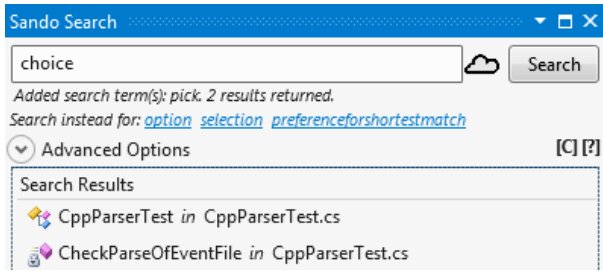
Fig. 3: Post-search recommendations.

codebase. The post-search recommendation assumes that the failure is either due to the vocabulary mismatch problem or to misspelling. Hence, the recommended terms are either the synonyms to the original terms or the original terms with corrected typos. Figure 3 depicts the UI of the post-search recommendations after searching "choice" fails.

For a query leading to no search results, the recommendation technique first pre-processes the query by splitting it. Using the terms in the local dictionary, we greedily extract the terms of the query that appear in the local codebase. Our recommendation technique does not recommend replacements for these terms. For the remaining terms that do not appear in the local dictionary, the post-search recommendation goes through the following steps to compute the recommended replacements:

**Step 1:** For a term that does not appear in the local dictionary, Sando first queries the Software Engineering thesaurus to find the synonyms of the term. After finding several synonyms, Sando recommends them to the developer. Sando filters the the recommended sysnonyms with the local dictionary to ensure that every recommended term appears in the local codebase at least once.

**Step 2:** When finding no field-specific synonyms, our technique queries the general English dictionary to find synonyms. After finding several synonyms, Sando recommends them to the developer. Similar with Step 1, Sando also filters these sysnonmys.

**Step 3:** If neither thesauri contains the synonyms of the given term, Sando considers the term a typo. Therefore, Sando corrects the typo by using the terms in the local dictionary.

## III. RELATED WORK

A large body of existing research relates to ours. We briefly summarize them according to two aspects: search tools and recommendation techniques.

**Search Tools.** Similar to Sando, researchers have proposed multiple tools allowing developers to efficiently access the enormous information space. For instance, Grechanik et al. proposed a search engine called Exemplar that retrieves function-level code elements from the web and also visualizes them [1]. Bazrafshan et al. presented a search tool that can retrieve code across the boundary of versions and branches in the source control system [2]. Differing from these tools, Sando searches code snippets from the local codebase.

**Recommender Systems.** Recommender systems play an important role in the software engineering research. For instance, Thummalapenta and Xie presented PARSEWeb that recommends method invocations sequence that bridges the starting object with the destination object [3]. Ankolekar et al. proposed Dhruv to recommend reusable artifacts to open source developers [4]. Differing from these techniques, the recommendation techniques adopted by Sando assist developers in composing promising queries.

## IV. CONCLUSION

Facing an ever-increasing codebase, developers often start tasks with searching their local code. However, even the state-of-the-art code search tools require developers' recalling the text of the intended code snippets. To further bridge the gap between developers and their intended code snippets, we demonstrate a set of recommendation techniques integrated into Sando. Issuing recommendations both before and after a developer's search, our technique potentially mitigates developers' cognitive burden when using local code search tools.

## REFERENCES

[1] C. McMillan, M. Grechanik, D. Poshyvanyk, C. Fu, and Q. Xie. *Exemplar: A Source Code Search Engine for Finding Highly Relevant Applications*. In TSE, 38(5):1069–1087, 2012.

[2] S. Bazrafshan, R. Koschke, and N. Gode. *Approximate code search in program histories*. In Proc. WCRE, pages 109-118, 2011.

[3] S. Thummalapenta and T. Xie. *Parseweb: a programmer assistant for reusing open source code on the web*. In Proc. ASE, pages 204-213, 2007.

[4] A. Ankolekar, K. Sycara, J. Herbsleb, R. Kraut, and C. Welty. *Supporting online problem-solving communities with the semantic web*. In Proc. WWW, pages 575-584, 2006.

[5] D. Shepherd, K. Damevski, B. Ropski, and T. Fritz. *Sando: an extensible local code search framework*. In Proc. FSE, pages 15:1-15:2, 2012.

[6] InstaSearch: Eclipse plug-in for quick code search. https://code.google.com/a/eclipselabs.org/p/instasearch/, 2013.

[7] Z.P. Fry, D. Shepherd, E. Hill, L. Pollock, and K. Vijay-Shanker. *Analysing source code: looking for useful verb-direct object pairs in all the right places*. IET Software, 2(1):27-36, 2008.

[8] S. Gupta, S. Malik, L. Pollock, and K. Vijay-Shanker. *Part-of-speech tagging of program identifiers for improved text-based software engineering tool*. In Proc. ICPC, 2013.

[9] G. A. Miller. *Wordnet: A lexical database for english*. Communications of the ACM, 38:39-41, 1995.

[10] A. J. Ko, B. A. Myers, M. J. Coblenz, and H. H. Aung. *An exploratory study of how developers seek, relate, and collect relevant information during software maintenance tasks*. IEEE TSE, 32(12):971-987, 2006.

[11] J. Singer, T. Lethbridge, N. Vinson, and N. Anquetil. *An examination of software engineering work practices*. In Proc. CCASCR, pages 21-36, 1997.

[12] M. H. Schultz S. C. Eisenstat, M. C. Gursky and A. H. Sherman. *Yale sparse matrix package i: The symmetric codes*. Journal of Numerical Methods in Engineering, 18:1145-1151, 1982