

Generation and Visualization of Static Function Call Graph for Large C Codebases

Sourabh S Badhya, Shobha G



Abstract: As software systems evolve, there is a growing concern on how to manage and maintain a large codebase and fully understand all the modules present in it. Developers spend a significant amount of time analyzing dependencies before making any changes into codebases. Therefore, there is a growing need for applications which can easily make developers comprehend dependencies in large codebases. These applications must be able to analyze large codebases and must have the ability to identify all the dependencies, so that new developers can easily analyze the codebase and start making changes in short periods of time. Static analysis provides a means of analyzing dependencies in large codebases and is an important part of software development lifecycle. Static analysis has been proven to be extremely useful over the years in their ability to comprehend large codebases. Out of the many static analysis methods, this paper focuses on static function call graph (SFCG) which represents dependencies between functions in the form of a graph. This paper illustrates the feasibility of many tools which generate SFCG and locks in on Doxygen which is extremely reliant for large codebases. The paper also discusses the optimizations, issues and its corresponding solutions for Doxygen. Finally, this paper presents a way of representing SFCG which is easier to comprehend for developers.

Keywords: Static function call graph, Static analysis, Duplicate functions, Doxygen, Cytoscape.js

I. INTRODUCTION

Large software companies write quality software on a daily basis and track individual components very well. Codebases written in C language have been in use in system development, network development and many applications mainly because C language is easier to interface with machine hardware, consumes less memory and has faster runtimes. It also provides great control to the programmer to create efficient programs.

Over the years, it has been found that large codebases are difficult to analyze due to its sheer scale and complexity involved and the different ways in which files can be linked with each other. There could be several links of a function present in a file with several other functions present within the file or other files which makes it hard for the developer to comprehend these links while making changes to the codebase. There exists a study by T. D. LaToza et al. [1]

which indicates ease of comprehending complex codebases using call graphs and it mentions that developers who use visual tools for call graphs were more likely to complete a task much faster than developers who do not use visual tools for call graphs. Hence, there is a need for a visual representation which can be used for easily visualizing all function call dependencies amongst files in the codebases.

The early mention of using Call Graphs for static analysis was done in B. G. Ryder et al. [2] which explicitly defines a call graph as a representation in which “the nodes of the graph are the procedures of the program; each edge represents one or more invocations of a procedure P_j by a procedure P_i ”. The first efforts in this direction were in the Fortran language in which call graphs were represented as a directed acyclic graph (since Fortran 77 is a non-recursive language). However, from then onwards several implementations of call-graphs have arisen in different languages and are still under active research. Most of the implementations make use of the abstract syntax tree (AST) of the underlying language in order to get the call dependencies amongst different functions. However, most of these languages have a single compiler hence making it simple to extract the required information. This is unlike the C language wherein multiple compilers exist and all these compilers were written for some use case such as Intel’s C/C++ compiler was written for Intel processors to achieve best runtimes.

There have been several efforts in making developers comfortable with analyzing large C codebases. Compiler specific solutions for static analysis have been present for several years, such as Clang Static Analyzer [3]. However, these tools are compiler specific and since there are many compilers, using any of these tools might result in compatibility problems. Hence the required call graph tool must not be dependent on a specific compiler. This paper discusses the usage of Doxygen [4] as one such tool which is scalable, almost accurate and compiler independent. It further discusses how the functions and function calls can be visualized by using Cytoscape.js [5] which is a state-of-the-art graph theory visualization JavaScript framework to enable a better viewing experience for developers.

II. RELATED WORK

Several implementations have been developed in this field in order to accurately represent function call graphs. The early formal research into the extraction of function calls was achieved by D. Callahan et al. [6] which proposes an algorithm to extract function calls from Fortran 8x language which handles recursion and has a time complexity

Manuscript received on June 05, 2021.

Revised Manuscript received on June 12, 2021.

Manuscript published on July 30, 2021.

* Correspondence Author

Sourabh S Badhya*, Department of Computer Science, R.V. College of Engineering, Bengaluru (Karnataka), India. Email: sourabhsbadhya.cs17@rvce.edu.in

Shobha G, Professor, Department of Computer Science, R.V. College of Engineering, Bengaluru (Karnataka), India. Email: shobhag@rvce.edu.in

© The Authors. Published by Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP). This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

Retrieval Number: 100.1/ijsce.F35070510521

DOI: 10.35940/ijsce.F3507.0710621

Journal Website: www.ijsce.org

Published By:
Blue Eyes Intelligence Engineering
and Sciences Publication
© Copyright: All rights reserved.



dependent exponentially on the number of procedures. It also asserts that polynomial time complexity can be achieved by setting an upper bound on the number of procedures. This was one of the earliest algorithms which could identify both normal function calls as well as recursion correctly, which is used in modern programming languages.

A more rigorous overview of call-graph as a graph theory problem was done in T. Reps et al. [7] whose focus was to find solutions to a large set of interprocedural dataflow problems in polynomial time. He also proposes that instead of calculating the worst-time complexity of each algorithm, it is better to bound the total cost of the operations performed at each aspect.

An in-depth dissertation on the effect of function pointers on call graphs was given in G. Antoniol et al. [8] wherein it proposes an algorithm to identify function pointers called the 'points-to' algorithm. It also provides a quantitative evaluation of function pointers and the key role it plays in call graph construction.

An algorithm to specifically identify virtual function calls as well as interfaces was also proposed in X. Zhuo et al [9], wherein it makes use of type flow analysis to get the call dependencies. This algorithm also takes less time and space usage when compared to points-to algorithm. However, this algorithm only applies to object-oriented languages.

An elaborate framework for call graph construction was proposed by D. Grave et al. [10]-[11]. These papers propose a general parameterized algorithm which provides a detailed vocabulary for depicting call graph algorithms, illustrates the differences and similarities of different algorithms and investigates the design space of call graph algorithms. It also assesses call graph algorithms with respect to an optimizing compiler (Vortex compiler) and the algorithms can be applied to any functional language.

An alternative method for call-graph construction was proposed in Y. Terashima et al. [12] wherein a tool by the name 'dgg' was proposed which made use of DWARF2 debugging information. This paper illustrates a method which combines both binary analysis and debugging information in order to extract function call dependencies. This method could extract inline functions in C code as well virtual function calls in C++ apart from the default functions present in code.

One of the earliest tools in order to generate call graphs is described in G. Antoniol et al. [13] by the name 'XOgastan'. It was developed as a static analysis tool which makes use of gcc/g++ compiler. It exploits the internal representation of gcc/g++ compiler which is the abstract syntax tree and translates it into graph exchange language representation. The final output is in the XML format which can be easily parsed with XML parsers. However, since it is compiler dependent it cannot be used for codebases which are not built on gcc/g++ compiler.

A framework was also proposed in H. Hoogendorp [14] which describes the complete steps right from data extraction to visualization of call graphs. This was one of the first frameworks which accounted for scalability of the system. It also explains about the different ways in which visualization can be done. The problem with this framework is that compiler wrapping must be done for every compiler. This makes the system compiler dependent.

Another alternative method for call graph generation was done in F. Zhang et al. [15] which proposes a static analysis

method that analyses the LLVM IR (Internal representation) generated by compilation of source programs. This method was used in order to analyze the parent-child relationship in between threads which are created using the pthread library. Compiler dependency of LLVM is the drawback of this method. A recent and similar tool which was developed for static analysis which is described in P. D. Schubert et al. [16] by the name 'Phasar'. Phasar is built as an extension to the LLVM compiler infrastructure. The analysis is done on LLVM IR, since solving data-flow problems on IR is easier than source code itself. However, this makes the implementation compiler dependent in nature.

Another tool for static analysis of source code was proposed in M. L. Collard et al. [17] which is named as srcML. It is a highly scalable tool and robust tool for source code analysis. It converts source code into XML format by making use of Clang AST, hence making the tool compiler dependent. A similar approach for C/C++ source code is also stated in A. M. Bogar et al. [18] in their implementation MLSA wherein it is designed to provide support for multiple languages. MLSA is lightweight, scalable and is written as an island grammar (a technique used to support multiple languages). The C/C++ code is parsed using Clang AST which is again compiler dependent.

From all the papers featured above, we see that most of the tools are compiler dependent. Some tools are not scalable and hence do not work for large codebases. In order to address these issues, we are making use of Doxygen which is an effort to document large codebases. However, the parser that is used in Doxygen can be exploited to extract function call dependencies. The proposed system also discusses a unique way of representing function calls using Cytoscape.js which is interactive, user friendly and scalable by the number of nodes.

III. PROPOSED SYSTEM

The proposed system makes use of Doxygen as its tool for static analysis. It can be broken down into 2 phases - Preprocessing phase and visualization phase.

A. Preprocessing phase

This phase mainly involves extraction of call dependencies by using Doxygen. Doxygen requires a config file called Doxyfile which specifies the input files that must be analyzed, along with the output format that should be generated. Doxygen supports various output formats like HTML, XML, Latex, SQLite etc. In this system, we will use SQLite [19] database since it is lightweight, reliable and easily portable. The steps involved in preprocessing phase are as follows –

- **Identification of C files that are used during build / compilation** - The C files which are used to create final executables / images are being used as input to the Doxyfile. The identification can be either as simple as looking into all C files in the codebase or looking into the Make files / past execution logs and deriving information about the C files. Once the C files are identified which are used for a specific build, it is inserted into the Doxyfile in the INPUT tag.



```

INPUT= main1.c main2.c ... mainN.c
EXTRACT_ALL = YES
EXTRACT_PRIVATE = YES
EXTRACT_STATIC = YES
REFERENCES_RELATION = YES
REFERENCED_BY_RELATION = YES
LOOKUP_CACHE_SIZE = 2
GENERATE_XML = NO
GENERATE_HTML = NO
GENERATE_LATEX = NO
GENERATE_SQLITE3 = YES
SQLITE3_OUTPUT = sqlite
SQLITE3_RECREATE_DB = YES
    
```

Fig. 1. A sample Doxyfile

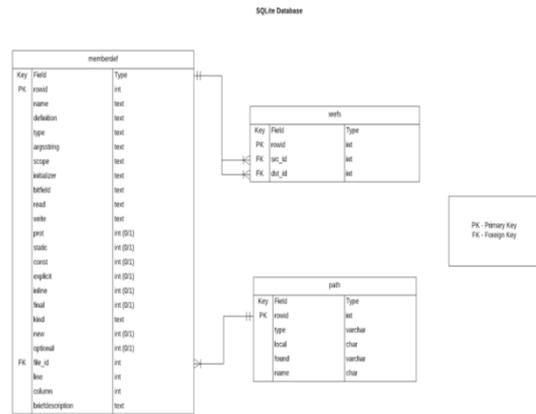


Fig. 2. Doxygen SQLite Database

- Use Doxygen to extract call dependencies** - Run Doxygen by using the Doxyfile. A sample Doxyfile would have the following tags as specified in Fig. 1. Some additional tags such as LOOKUP_CACHE_SIZE, NUM_PROC_THREADS can be set for faster execution time for very large codebases. This will be discussed in detail in the experimental results section.
- Create support for duplicate functions** - Duplicate functions are functions with the same name and same arguments but are implemented more than once. This is usually found in codebases wherein a functionality has several implementations and all of them are equally important. However, they are compiled separately in different images. These functions are usually seen in networking codebases wherein a single functionality has multiple implementations and depending upon the use case, one function is used over another. However, such dependencies are not handled in static analysis tools and hence support needs to be created. This can be done as follows -

 - Identify all the functions which have the same name, return type and arguments.
 - Get all the caller dependencies of all the duplicate functions involved.
 - Insert a link between the caller dependencies with all the duplicate functions.

In other words, the caller dependencies are shared amongst all the duplicate functions. This will make sure that the call graph generated will reflect all the dependencies of the duplicate functions. The data extracted by Doxygen and the SQLite database output is described in Fig. 2. Similarly, other fields are defined which indicate the properties of the member indexed. The table 'xrefs' has the information pertaining to call dependencies. The fields 'src_rowid' and 'dst_rowid' are foreign keys to 'memberdef' table's 'rowid'. A record in 'xrefs' table signifies that there is a call dependency from member with ID 'src_rowid' to member with ID 'dst_rowid' provided that both these members are functions.

This is a unique way in which nodes and edges are stored and this can be easily used in the visualization phase. As discussed earlier, the links to the duplicate functions are inserted to the 'xrefs' table and no other table is modified. The 'path' table has information about the location of the members indexed in 'memberdef'. In summary, the entire preprocessing phase can be shown in Fig. 3.

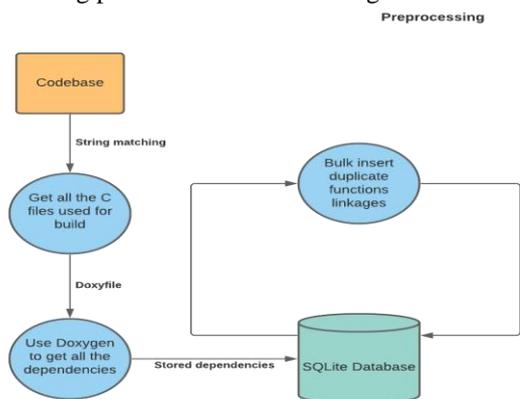


Fig. 3. Preprocessing phase

B. Visualization phase

The visualization phase involves using the SQLite database in a Node.js application. Cytoscape.js is used for graph visualization since it is highly optimized, scalable and user-friendly. The function call graph is represented by making use of a compound graph wherein the file in which the function is present is symbolized as the parent node and function present in the file is represented as the child node. This creates a sense of inclusivity of functions within files and is easily understandable to the developer.

The caller and callee edges are present in between different functions and no edges are present in between files. Cytoscape.js expects JSON formatted data in order to include edges and nodes. Each node and edge are uniquely identified by its ID, hence this ID is set to the row ID in 'memberdef' table for nodes and row ID of 'xrefs' table for edges to make it unique. In order to facilitate the usage of compound graph, Cola extension [20] is also used along with Cytoscape.js. It is also used to preassign the position of nodes, without manual setting of positions.



Every node and edge are randomly positioned with a predefined distance amongst them in order to avoid overlapping.

Interactive functionalities can be implemented by creating buttons on the nodes of the graph. Since Cytoscape.js uses canvas element to plot the nodes and edges of a graph, all these buttons must be plotted on top of the canvas. The buttons are created using the context-menus extension [21].

The architecture of the visualization phase can be shown in Fig. 4. When an input such as a function name and file name is taken from the sidebar window, appropriate results from the SQLite database are extracted and given to the webview which uses Cytoscape.js and plots the nodes and directed edges on the graph. All the information transfer between the extension backend and the webview happens in the form of messages sent from application backend to webview and vice-versa.

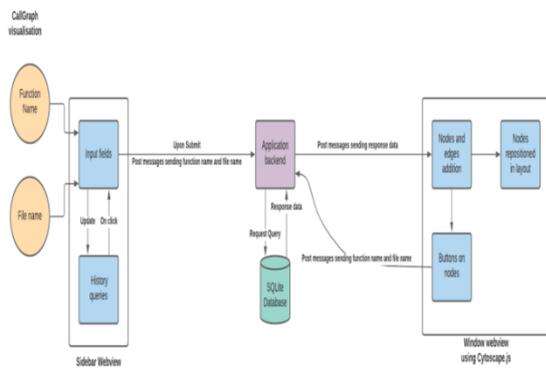


Fig. 4. Visualization phase

IV. EXPERIMENT RESULTS

The experiments were conducted in a RHEL 8 machine with system specifications of 8 cores, 8 GB RAM, 512 GB hard disk. Several validations were done in order to identify duplicate functions correctly. The source of errors was seen in unused function pointers. However, unused function pointers are usually considered a bad practice, but they are valid. An example for duplicate function is shown in Fig. 5. The graph represents a small codebase with 3 files – main.c, a.c and b.c. The file main.c has the driver function ‘main’ which calls a function called ‘type1’ which is present in both a.c and b.c. The definition of ‘type1’ is different in both files and hence they call different functions. However, when the code is compiled, only one file is passed into compilation thereby having only one definition of ‘type1’ during compilation.

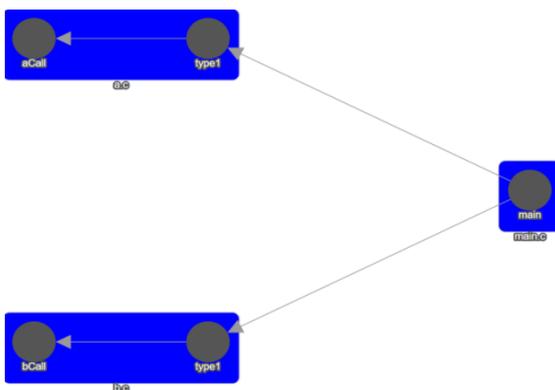


Fig. 5. Example of duplicate functions

Table - I: Comparison of codebases on different metrics

Codebases	Size of codebase (kLOC)	Time taken to analyze (s)	Output SQLite database (MB)
Redis	181.118	9.861	4.4
OpenSSL	522.121	34.607	14
PostgreSQL	1469.177	62.677	25
Linux	20550.393	10255.646	732

A comparison was done amongst 4 codebases - Redis [22], OpenSSL [23], PostgreSQL [24] and Linux [25]. Table. I illustrate the results obtained when Doxygen was run on these codebases. This table indicates that smaller codebases having smaller number of functions generate less amount of data when compared to its larger counterpart which is the obvious norm. It also indicates that if function call dependencies are more in a codebase, then a lot of data is generated, hence taking more time. This is especially seen in Linux codebase since the order at which memory increases is higher than order at which size of the codebases increases. Higher variable usage is also a potential reason for large database output as well since Doxygen generates data for variables as well.

Time required vs. Lookup Cache size

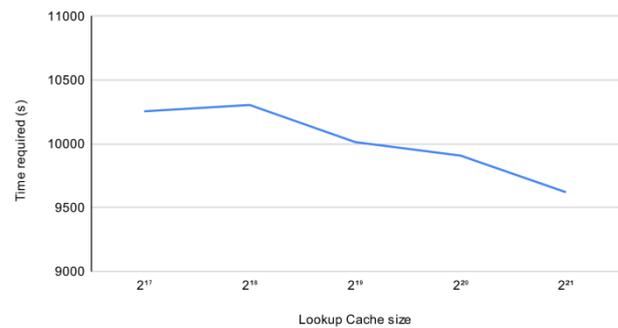


Fig. 6. Time required v/s. Lookup cache size

Cache misses vs. Lookup Cache size

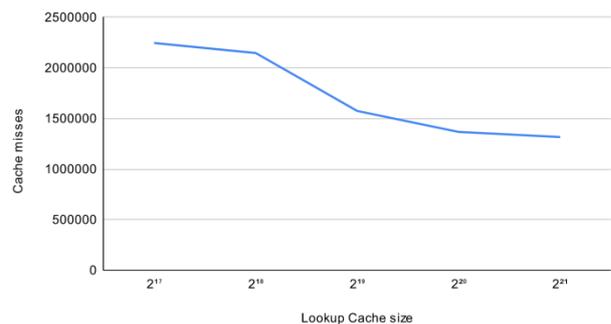


Fig. 7. Cache misses v/s. Lookup cache size

Optimization on execution performance was also considered. Doxygen was not originally coded for providing the best runtime behavior, however efforts were made to bring in efficiency. Doxygen makes use of a symbol lookup cache while parsing the code whose sizes can be manipulated to bring better performance. It is controlled by the following tag - LOOKUP_CACHE_SIZE.



AUTHORS PROFILE



Natural Language Processing domain.

Sourabh S Badhya, is an undergraduate student who is pursuing Computer Science & Engineering in R.V. College of Engineering, Bengaluru. His area of interests lies in Machine Learning and Natural Language Processing. He has 3 publications in his name in the



Dr. **Shobha G**, is a professor in R. V. College of Engineering, Bengaluru. Her area of interests lies in data mining, image processing and networking. She has a teaching experience of 25 years and research experience of 14 years. She has 145 publications in international journals and conferences. She has guided 30 UG projects, 40 PG projects and 7 Ph.D. students. She also has 4 patents and has been a reviewer of several books.