

**A Hybrid Approach for Selecting and Optimizing Graph Traversal Strategy for
Analyzing BigCode**

by

Ramanathan Ramu

A thesis submitted to the graduate faculty
in partial fulfillment of the requirements for the degree of
MASTER OF COMPUTER SCIENCE

Major: Computer Science

Program of Study Committee:
Dr. Hridesh Rajan, Major Professor
Dr. Andrew Miner
Dr. Wei Le

Iowa State University

Ames, Iowa

2017

Copyright © Ramanathan Ramu, 2017. All rights reserved.

TABLE OF CONTENTS

LIST OF TABLES	v
LIST OF FIGURES	vi
ACKNOWLEDGEMENTS	viii
ABSTRACT	ix
CHAPTER 1. Introduction	1
CHAPTER 2. Contributions	6
2.1 Traversal Declaration and Traverse Expression	6
2.2 Data-Flow and Loop Sensitivity Analyses for Traversals	6
2.3 Graph Cyclicity	7
2.4 Decision Tree for Traversal Strategy Selection	7
2.5 Evaluation Contribution	7
CHAPTER 3. Background	9
3.1 Graph	9
3.2 Graph traversal	9
3.3 Graph traversal strategies	10
3.3.1 Depth-first search	10
3.3.2 Breadth-first search	10
3.4 Program analysis	11
3.5 Control-flow and Data flow analysis	11
3.6 Graph traversal for Program analysis	11
3.6.1 Random order	11

3.6.2	Postorder	12
3.6.3	Reverse postorder	12
CHAPTER 4.	Hybrid Traversal Selection for Efficient Source Code Analysis	13
4.1	A System For Expressing Source Code Analysis As Traversals	14
4.2	Static and Runtime Properties	19
4.2.1	Data-Flow Sensitivity	20
4.2.2	Computing Data-Flow Sensitivity	21
4.2.3	Loop Sensitivity	22
4.2.4	Computing Loop Sensitivity	23
4.2.5	Graph Cyclicity	26
4.3	Traversal Strategies - Candidates	26
4.4	Decision Tree for Traversal Strategy Selection	27
4.4.1	An Example	31
4.5	Optimizing the Selected Traversal Strategy	31
CHAPTER 5.	Implementation on Boa framework	34
5.1	Boa language and infrastructure	34
5.2	Source code analysis using Traversal construct	34
5.3	Putting it all together	36
CHAPTER 6.	Empirical Evaluation	42
6.1	Analyses, Datasets and Experiment Setting	42
6.1.1	Analyses	42
6.1.2	Datasets	43
6.1.3	Setting	44
CHAPTER 7.	Running Time and Time Reduction	46
7.1	Running Time	46
7.2	Time Reduction	47
7.3	Time reduction against hand optimized analysis	48
CHAPTER 8.	Correctness of Analysis Results	50

CHAPTER 9. Traversal Strategy Selection Precision	51
CHAPTER 10. Analysis on the Decision Tree Distribution	55
CHAPTER 11. Analysis on Traversal Optimization	58
CHAPTER 12. Case Studies	60
12.1 API Precondition Mining (APM).	60
12.2 API Usage Mining (AUM).	61
CHAPTER 13. Threats to Validity	63
CHAPTER 14. Related Work	64
14.1 Mixing static and dynamic information.	64
14.2 Optimizing program analysis.	64
14.3 Ultra-large-scale source code mining.	65
14.4 Graph traversal optimization.	66
CHAPTER 15. Conclusion and Future work	67
15.1 Conclusion	67
15.2 Future Work	68
CHAPTER 16. Appendix	74

LIST OF TABLES

4.1	Syntax reference.	16
4.2	Operations on collections.	16
6.1	List of source code analyses and the properties of their involved traversals.	43
6.2	Statistics of the generated control flow graphs from two datasets. . . .	43
6.3	Time contribution of each phase (in miliseconds).	44
9.1	Traversal strategy prediction precision.	51

LIST OF FIGURES

1.1	Running times (ms) of the three analyses on graph A using different traversal strategies.	2
1.2	Running times of three analyses using different traversal strategies on a large codebase.	4
4.1	Overview of the hybrid approach for selecting and optimizing graph traversal strategy.	13
4.2	Running example of applying the post dominator analysis on an input graph containing branch and loop.	19
4.3	Traversal strategy selection decision tree.	28
7.1	Reduction in running times.	47
7.2	Reduction in running times against hand optimized analysis.	49
9.1	Scatter charts for analyses that have loop sensitive traversals.	51
9.2	Hybrid approach's performance against best approaches for mis-predicted graphs.	54
10.1	Distribution of decisions over the paths of the decision tree.	56
10.2	Distribution of decisions over the paths of the decision tree for the DaCapo Dataset.	57
11.1	Reduction in execution time of the hybrid approach due to traversal optimization.	58
12.1	Running time (minutes) of the case studies on GitHub data.	60

12.2	First and second most mined pre-condition for some of the methods in String java API.	61
12.3	Top 10 API Usage pattern for java.util API.	61

ACKNOWLEDGEMENTS

I would like to take this opportunity to express my thanks to those who helped me with various aspects of conducting research and the writing of this thesis. I would like to thank Dr. Hridesh Rajan, Dr.Hoan A Nguyen and Ganesha Upadhaya for their guidance, patience and support throughout this research and the writing of this thesis. Thanks are due to the US National Science Foundation for financially supporting this project. I would like to thank my committee members Dr. Wei Le and Dr. Andrew Miner for their efforts and contributions to this work. Also, I would like to thank the reviewers of ECOOP 2017 and OOPSLA 2017 conference for their insightful feedback. I would like to extend my thanks to all the members of Laboratory of Software Design for offering constructive criticism and timely suggestions during research.

I am very grateful to my parents Ramu and Meenal and my friends for their moral support and encouragement throughout the duration of my studies.

ABSTRACT

Our newfound ability to analyze source code in massive software repositories such as GitHub has led to an uptick in data-driven solutions to software engineering problems. Source code analysis is often realized as traversals over source code artifacts represented as graphs. Since the number of artifacts that are analyzed is huge, in millions, the efficiency of the source code analysis technique is very important. The performance of source code analysis techniques heavily depends on the order of nodes visited during the traversals: the traversal strategy. For instance, selecting the best traversal strategy and optimizing it for a software engineering task, that infers the temporal specification between pairs of API method calls, could reduce the running time on a large codebase from 64% to 96%. While, there exists several choices for traversal strategy, like depth-first, post-order, reverse post-order, etc., there exists no technique to choose the most time-efficient strategy for traversals. In this paper, we show that a single traversal strategy does not fit all source code analysis scenarios. Somewhat more surprisingly, we demonstrate that given the source code expressing the analysis task (in a declarative form) one can compute static characteristics of the task, which together with the runtime characteristics of the input, can help predict the most time-efficient traversal strategy for that (analysis task, input) pair. We also demonstrate that these strategies can be realized in a manner that is effective in accelerating ultra-large-scale source code analysis. Our evaluation shows that our technique successfully selected the most time-efficient traversal strategy for 99.99%-100% of the time and using the selected traversal strategy and optimizing it, the running times of a representative collection of source code analysis in our evaluation were considerably reduced by 1%-28% (13 minutes to 72 minutes in absolute time) when compared against the best performing traversal strategy. The overhead imposed by collecting additional information for our approach is less than 0.2% of the total running time for a large dataset that contains 287K Control Flow Graphs (CFGs) and less than 0.01% for an ultra-large dataset that contains 162M CFGs.

CHAPTER 1. Introduction

The availability of open source repositories like GitHub is driving data-driven solutions to software engineering problems, e.g. specification inference by Nguyen et al. (2014), discovering programming patterns by Thummalapenta and Xie (2009), suggesting bug fixes by Livshits and Zimmermann (2005); Li et al. (2006), etc. These software engineering tasks analyze different source code representations, such as text, abstract syntax trees (ASTs), control flow graphs (CFGs), at massive scale. The performance of source code analysis over graphs heavily depends on the order of the nodes visited during the traversals: *the traversal strategy*. While graph traversal is a well-studied problem, and various traversal strategies exists; e.g., depth-first, post-order, reverse post-order, etc, no single strategy works best for different kinds of analyses and different kinds of graphs. Our contribution is *hybrid traversal selection*, a novel source code analysis optimization technique for source code analyses expressed as graph traversals. This work was done collaboratively with Dr. Hridesh Rajan, Dr.Hoan A Nguyen and Ganesha Upadhaya. Consequently, chapters 1-4 in this thesis are based on that collaborative work.

Motivation and Key Observations. To motivate, consider a software engineering task that infers the temporal specifications between pairs of API method calls, i.e., a call to a must be followed by a call to b by Engler et al. (2001); Ramanathan et al. (2007); Weimer and Necula (2005); Yang et al. (2006). A data-driven approach for inference is to look for pairs of API calls that frequently go in pairs in the same order at API call sites in the client methods' code. Such an approach contains (at least) three source code analyses on the control flow graph (CFG) of each client method: 1) identifying references of the API classes and call sites of the API methods which can be done using *reaching definition* analysis(Nielson et al. (2010)); 2) identifying the pairs of API calls (a, b) where b follows a in the client code which can be done using *post-dominator* analysis(Aho et al. (2006)); and 3) collecting pairs of temporal API calls

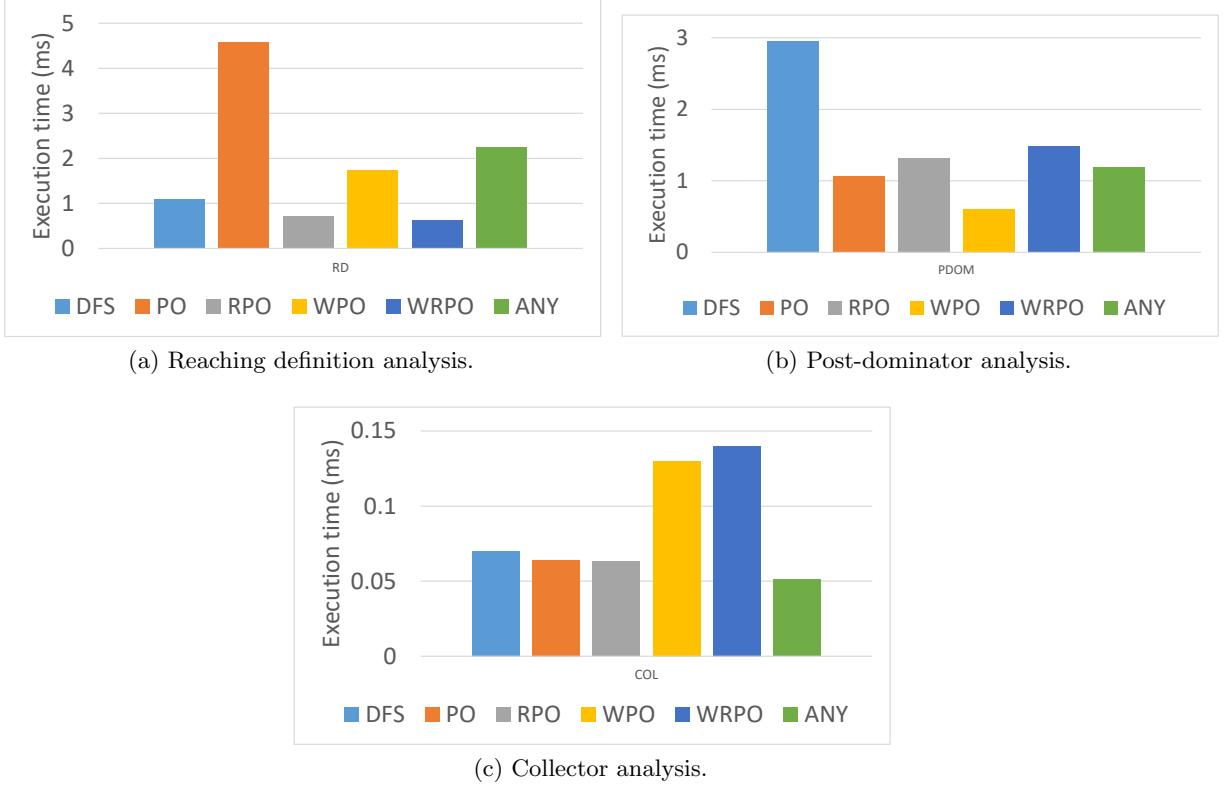


Figure 1.1: Running times (ms) of the three analyses on graph A using different traversal strategies.

by traversing all nodes in the CFG—let us call this *collector* analysis. These analyses need to be run on a large number of client methods to produce temporal specifications with high confidence.

Implementing each of these analyses involves traversing the CFG of each client method. The traversal strategy could be chosen from a set of standard strategies e.g depth-first search (DFS), post-order (PO), reverse post-order (RPO), worklist with post-ordering (WPO), worklist with reverse post-ordering (WRPO) and any order (ANY). In ANY order traversal strategy, nodes can be visited in any order. In post-order traversal, the successors of any node are visited before visiting the node, while in reverse post-order traversal, the predecessors of any node are visited before visiting the node. In Worklist with Post-Order and Worklist with reverse post-order, the nodes are visited in the order they appear in the worklist. A worklist is a data structure used to keep track of nodes to be visited. In WPO, worklist is initialized with post-ordering of nodes, while in WRPO, The worklist is initialized with nodes in the reverse post-order.

Figure 1.1 shows the performance of each of these three analyses when using standard traversal strategies. These runs are analyzing the CFG of a method in the DaCapo benchmark(Blackburn et al. (2006)). Actual implementation of this method is not important, but it suffices to know that the CFG, which we shall refer to as Graph A, has 50 nodes and has branches but no loops. Figure 1.1 shows that, for graph A, the WRPO performs better than other strategies for the reaching definition analysis while the WPO outperforms the others for the post-dominator analysis and the ANY traversal works best for the collector analysis.

No Traversal Strategy Fits All. The performance results are somewhat expected, but require understanding the subtleties of the analyses. Reaching definition analysis is a forward data-flow analysis where the output at each node in the graph is dependent on the outputs of their predecessor nodes. So, DFS, RPO and WRPO by nature are the most suitable. However, worklist is the most efficient strategy here because it visits only the nodes that are yet to reach fixpoint unlike other strategies that also visit notes that have already reached fixpoint. Post-dominator analysis, on the other hand, is a backward analysis meaning that the output at each node in the graph is dependent on the outputs of their successor nodes. Therefore, the worklist with post-ordering is the most efficient traversal. For the collector analysis, any order traversal works better than other traversal strategies for graph A. This is because for this analysis the output at each node is not dependent on the output of any other nodes and hence it is independent of the order of nodes visited. The any order traversal strategy does not have the overhead of visiting the nodes in any particular order like DFS, PO, RPO nor the overhead to maintain the worklist. Therefore any order traversal performs better than other traversal strategies.

Properties of Input Graph Determine Strategy. For the illustrative example discussed above, DFS and RPO were worse than WRPO for the reaching definition analysis and PO was worse than WPO for post-dominator because they require one extra iteration of analysis to be performed and realize that fixpoint has been reached. However, since graph A does not have any loops, if the graph A's nodes are visited in such a way that each node is visited after its predecessors for reaching definition analysis and after its successors for post-dominator analysis, then the additional iteration is actually redundant. Given that graph A has no loops, one could

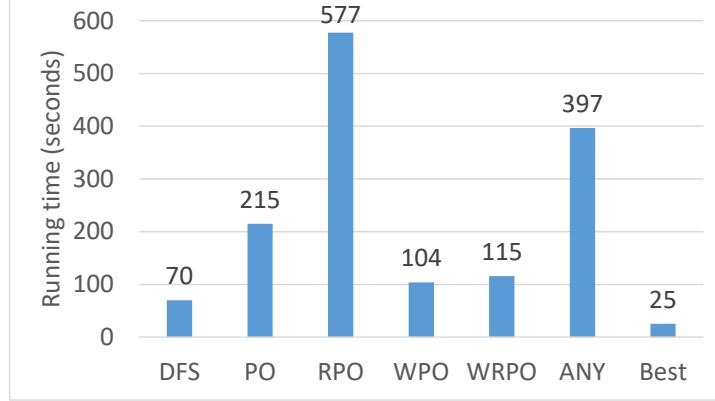


Figure 1.2: Running times of three analyses using different traversal strategies on a large codebase.

optimize RPS or PO to bypass the extra iteration and fixpoint checking. Thus, the optimized RPS or PO would run the same number of iterations as the respective worklist-based ones and finish faster than them because the overhead of maintaining the list is eliminated.

The potential gains of selecting a suitable traversal strategy can be significant. To illustrate, consider Figure 1.2 that shows the performance of our entire illustrative example (inferring temporal specifications) on a large corpus of 287,000 CFGs extracted from the DaCapo benchmark dataset(Blackburn et al. (2006)). Figure 1.2 shows the bar chart for the total running times of the three analyses. The **Best** strategy is an ideal one where we can always choose the most efficient with all necessary optimizations. The bar chart confirms that fixing a traversal strategies for running different analyses on different types of graphs is not efficient. Selecting and optimizing traversal strategy for each analysis on each graph is desirable. Such a strategy could reduce the running time on a large dataset from 64% (against DFS) to 96% (against RPO).

Our approach relies on our observations that a suitable traversal strategy is dependent on both the source code analyses and input graphs' properties. The former are static properties and the latter are runtime properties. More importantly, depending on the properties of analyses and graphs, existing traversals could be optimized to improve their performance. We have evaluated our technique using a set of 21 source code analysis that includes control and data-flow analysis, and analysis to find bugs. The evaluation is performed on two datasets: a dataset containing

well-maintained projects from DaCapo benchmark (contains a total of 287K graphs), and a ultra-large dataset containing more than 380K projects from GitHub (contains a total of 162M graphs). Our evaluation shows that our technique successfully selected the most time-efficient traversal strategy for 99.99%–100% of the time and using the selected traversal strategy and optimizing it, the running times of a representative collection of source code analysis in our evaluation were considerably reduced by 1%-28% (13 minutes to 72 minutes in absolute time) when compared against the best performing traversal strategy. The overhead imposed by our approach is negligible (less than 0.2% of the total running time for a large dataset and less than 0.01% for an ultra-large dataset).

The rest of the thesis is organized as follows. Chapter 2 lists the contribution of our work. Chapter 3 gives a background of many ideas that this thesis builds upon, such as Graph and its traversal, various traversal strategies, Program analysis, Control and data flow analysis. Chapter 4.1 describes our system For expressing source code analysis as traversals, while Chapter 4.2 presents the factors that influence the selection of the optimal traversal strategy for any given traversal. Chapter 4.3 lists all the candidate traversal strategies while Chapter 4.4 describes our decision tree to select traversal strategy and also presents an example analysis and graph and walks through the decision tree. Chapter 5 gives the implementation detail on how we extended Boa language to have hybrid capability. Chapter 6 presents our experimental setup and the dataset used. Chapter 7 provides the evaluation and the anlaysis that we did on running time comparison between hybrid and various other strategies. Chapter 8 presents the correctness of analysis using Hybrid approach. Chapter 9 presents the Hybrid’s traversal strategy prediction precision. Chapter 10 analyses our decision tree’s distribution while Chapter 11 analyses the importance of our traversal optimization. Chapter 12 discusses the two case studies that we implemented using our formalism for source code analysis and their results. Chapter 13 discusses about the threat to validity while Chapter 14 discusses related work. Chapter 15 concludes and gives some possible ideas for future work.

CHAPTER 2. Contributions

In this work, we develop *hybrid traversal selection*, a novel program analysis optimization technique for BigCode analyses expressed as graph traversals. Our approach relies on our observations that a suitable traversal strategy is dependent on both the program analyses and input graphs' properties. The former are static properties and the latter are dynamic properties. More importantly, depending on the properties of analyses and graphs, existing traversals could be optimized to improve their performance. Hybrid traversal selection relies on several technical underpinnings:

2.1 Traversal Declaration and Traverse Expression

Programmers can declare their program analyses as one or more **traversal** declarations and run them using **traverse** expression. The runtime implementation of the **traverse** expression selects a suitable traversal strategy based on the **traversal** declaration and the input graph. Main benefit of these linguistic abstractions is that they abstract away traversal related code so that the traversal strategy can be replaced as needed by the analysis runtime.

2.2 Data-Flow and Loop Sensitivity Analyses for Traversals

We show that traversal strategy selection depends on three critical properties of the traversal: *data-flow sensitivity*, *loop sensitivity*, and *traversal direction*. We propose algorithms for computing these properties. Our analysis system implements these algorithms. These properties are computed statically and their values are stored as metadata to be utilized by the traversal selection at runtime.

2.3 Graph Cyclicity

We have observed that the traversal strategy selection depends on one dynamic property of the input which is *graph cyclicity*. This property partitions the set of graphs into three categories: those that are sequential, those with branches but no cycles, and those with cycles. Our system computes this property at graph construction time and stores it as an attribute in the runtime graph representation.

2.4 Decision Tree for Traversal Strategy Selection

We have devised a *decision tree* for traversal strategy selection that given data-flow sensitivity, loop sensitivity, and traversal direction properties of the analysis and the cyclicity property of the input graph produces a selection for traversal strategy. While the tree is utilized by our automated system, it could also be used by a programmer for manual traversal selection.

Hybrid traversal selection has two direct benefits. First, it improves the efficiency of BigCode analysis thus speeding up data-driven science in this important area. Second, it frees up programmers from having to write traversal related code and then optimizing it based on the analysis and the graph at hand.

2.5 Evaluation Contribution

We have evaluated our technique using a set of 21 source code analysis that includes control and data-flow analysis, and analysis to find bugs. The evaluation is performed on two datasets: a dataset containing well-maintained projects from DaCapo benchmark (contains a total of 287K graphs), and a ultra-large dataset containing more than 380K projects from GitHub (contains a total of 162M graphs). Our evaluation shows that our technique successfully selected the most time-efficient traversal strategy for 99.99%–100% of the time and using the selected traversal strategy and optimizing it, the running times of a representative collection of source code analysis in our evaluation were considerably reduced by 1%–28% (13 minutes to 72 minutes in absolute time) when compared against the best performing traversal strategy. The case studies show that hybrid traversal reduces 80–175 minutes in running times for three software

engineering tasks. The overhead imposed by our approach is negligible (less than 0.2% of the total running time for a large dataset and less than 0.01% for an ultra-large dataset).

In summary, this paper makes the following contributions:

- It describes a system for expressing source code analysis as traversals. The constructs and operations in the system allows different source code analysis to be expressed in a manner that allows automatic selection of the best traversal strategies.
- It defines a set of novel properties about the traversal expressed in our system. It also describes algorithms for analysing traversals for inferring these properties.
- It describes a novel decision tree for selecting the most suitable traversal strategy. The static and runtime properties also allows certain optimizations to be performed on the selected traversal strategy to further improve the performance.
- It demonstrates the potentially broad range of applications of hybrid traversal selection for optimizing source code analysis such as available expressions, local may alias, live variable, nullness analysis, post dominator, reaching definitions, resource status, very busy expression, etc.

CHAPTER 3. Background

In this chapter, we will have a look at the background of many ideas that this thesis builds upon.

3.1 Graph

A graph is a structure amounting to a set of objects in which some pairs of the objects are in some sense "related". The objects correspond to mathematical abstractions called vertices (also called nodes) and each of the related pairs of vertices is called an edge. Typically, a graph is depicted in diagrammatic form as a set of dots for the vertices, joined by lines or curves for the edges. The edges may be directed or undirected. For example, if the vertices represent people at a party, and there is an edge between two people if they shake hands, then this graph is undirected because any person A can shake hands with a person B only if B also shakes hands with A. In contrast, if any edge from a person A to a person B corresponds to A's admiring B, then this graph is directed, because admiration is not necessarily reciprocated. The former type of graph is called an undirected graph and the edges are called undirected edges while the latter type of graph is called a directed graph and the edges are called directed edges.

3.2 Graph traversal

In computer science, graph traversal refers to the process of visiting each vertex in a graph. Such traversals are classified by the order in which the vertices are visited. Graph traversal may require that some vertices be visited more than once, since it is not necessarily known before transitioning to a vertex that it has already been explored. As graphs become more dense, this redundancy becomes more prevalent, causing computation time to increase; as graphs become

more sparse, the opposite holds true. Thus, it is usually necessary to remember which vertices have already been explored by the algorithm, so that vertices are revisited as infrequently as possible. This may be accomplished by associating each vertex of the graph with a "visitation" state during the traversal, which is then checked and updated as the algorithm visits each vertex. If the vertex has already been visited, it is ignored and the path is pursued no further; otherwise, the algorithm checks/updates the vertex and continues down its current path.

3.3 Graph traversal strategies

The two most common traversal patterns are breadth-first traversal and depth-first traversal.

3.3.1 Depth-first search

A depth-first search (DFS) is an algorithm for traversing a finite graph. DFS visits the child vertices before visiting the sibling vertices; that is, it traverses the depth of any particular path before exploring its breadth. A stack is generally used when implementing the algorithm. The algorithm begins with a chosen "root" vertex; it then iteratively transitions from the current vertex to an adjacent, unvisited vertex, until it can no longer find an unexplored vertex to transition to from its current location. The algorithm then backtracks along previously visited vertices, until it finds a vertex connected to yet more uncharted territory. It will then proceed down the new path as it had before, backtracking as it encounters dead-ends, and ending only when the algorithm has backtracked past the original "root" vertex from the very first step

3.3.2 Breadth-first search

A breadth-first search (BFS) is another technique for traversing a finite graph. BFS visits the neighbor vertices before visiting the child vertices, and a queue is used in the search process. This algorithm is often used to find the shortest path from one vertex to another

3.4 Program analysis

In computer science, program analysis is the process of automatically analyzing the behavior of computer programs regarding a property such as correctness, robustness, safety and liveness. Program analysis focuses on two major areas: program optimization and program correctness. The first focuses on improving the programs performance while reducing the resource usage while the latter focuses on ensuring that the program does what it is supposed to do. Program analysis can be performed without executing the program (static program analysis), during runtime (dynamic program analysis) or in a combination of both.

3.5 Control-flow and Data flow analysis

The purpose of control-flow analysis is to obtain information about which functions can be called at various points during the execution of a program. The collected information is represented by a control flow graph (CFG) where the nodes are instructions of the program and the edges represent the flow of control. By identifying code blocks and loops CFG becomes a starting point for compiler made optimizations.

Data-flow analysis is a technique designed to gather information about the values at each point of the program and how they change over time. This technique is often used by compilers to optimize the code. One of the most known examples of data-flow analysis is taint checking which consists of considering all variables which contain user supplied data which is considered "tainted", i.e. insecure and preventing those variables from being used until they have been sanitized. This technique is often used to prevent SQL injection attacks.

3.6 Graph traversal for Program analysis

3.6.1 Random order

This iteration order is not aware whether the data-flow equations solve a forward or backward data-flow problem. Therefore, the performance is relatively poor compared to specialized iteration orders.

3.6.2 Postorder

This is a typical iteration order for backward data-flow problems. In postorder iteration, a node is visited after all its successor nodes have been visited. Typically, the postorder iteration is implemented with the depth-first strategy.

3.6.3 Reverse postorder

This is a typical iteration order for forward data-flow problems. In reverse-postorder iteration, a node is visited before any of its successor nodes has been visited, except when the successor is reached by a back edge.

CHAPTER 4. Hybrid Traversal Selection for Efficient Source Code Analysis

In this chapter we first provide a brief overview of our technique, followed by an overview of the constructs used for expressing source code analyses. We then describe properties, analyses, and a decision tree that are the technical underpinnings of our selection technique.

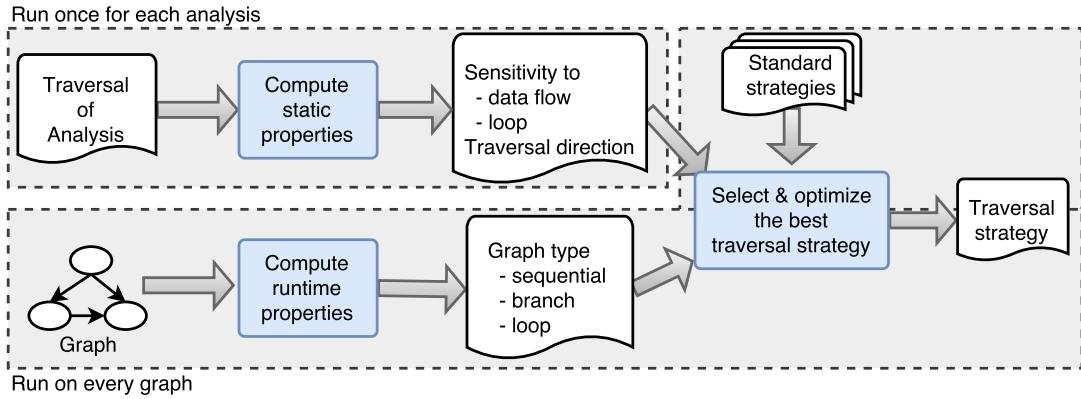


Figure 4.1: Overview of the hybrid approach for selecting and optimizing graph traversal strategy.

Figure 4.1 provides an overview of our approach and its key components. Inputs to our approach are source code analysis that contains one or more traversals (§4.1), and a graph. Output of our technique is an optimal traversal strategy for every traversal in the analysis. For selecting an optimal traversal strategy for a traversal, our technique computes a set of static properties of the analysis (§4.2), such as data-flow sensitivity, loop sensitivity, and extracts a runtime property about the graph that defines the cyclicity in the graph (sequential/branch/loop). Upon computing the static and runtime properties, our approach selects a traversal strategy from a set of candidate strategies (§4.3) for each traversal in the analysis (§4.4) and optimizes it (§4.5). The static properties of the traversals are computed only once for each analysis, whereas

graph cyclicity is determined for every input graph.

4.1 A System For Expressing Source Code Analysis As Traversals

A source code analysis is performed on various source code artifacts such as source code text, intermediate representations like abstract syntax trees (ASTs), graph-based representations like control flow graphs (CFGs) and call graphs (CGs), etc. In our system, source code analysis such as control- and data-flow analysis are expressed as traversals over CFGs.

Definition 1 *A control flow graph (CFG) is a directed graph $G = (N, E, n_{start}, N_{end})$ with a set of nodes N representing the program statements and a set of edges $E \subseteq N \times N$ representing the control flow relation between the program statements. A CFG has a single start node, n_{start} , and a set of end nodes, N_{end} .*

For any node $n \in N$, $n.\text{preds}$ is a set of immediate predecessors, $n.\text{succs}$ is a set of immediate successors, $n.\text{stmt}$ provides the program statement at the node, and $n.\text{id}$ is a unique identifier of the node. Here on, we use graph to refer to CFG.

A source code analysis over a graph visits nodes in the graph in certain order and collects information at nodes (aka, analysis facts or outputs). For instance, the *reaching definition* analysis over a CFG, visits every node in the CFG and collects the variable definitions at nodes as analysis facts. An analysis may require multiple traversals over a graph and each traversal may visit nodes multiple times (for fixpoint). For instance, the *reaching definition* analysis requires two traversal of the CFG: an *initialization* traversal for collecting the variable definitions at nodes as analysis facts, and a *propagation* traversal for propagating the analysis facts along the graph. The *initialization* traversal visits every node exactly once, whereas the *propagation* traversal may visit the nodes multiple times until a fixpoint is reached and the analysis facts at nodes does not change further.

In our system, a source code analysis over a graph is expressed by defining and invoking one or more traversals. A traversal is defined using a special `traversal` block:

```
t := traversal(n : Node) : T { tbody }
```

In this traversal block definition, `t` is the name of the traversal that takes a single parameter `n` representing the graph node that is being visited. A traversal may define a return type `T` representing the output type. The output type can be a primitive or a collection data type. A block of code that generates the traversal output at a graph node is given by `tbody`. The `tbody` may contain common statements and expressions, such as variable declarations, assignments, conditional statements, loop statements, and method calls, along with some special expressions discussed in this chapter.

A traversal can be invoked using a special `traverse` expression:

```
traverse(g, t, d, df, ls, fp)
```

A `traverse` expression takes six parameters: `g` is the graph to be traversed, `t` is the traversal to be invoked, `d` is the traversal direction and `df`, `ls`, `fp` are optional parameters. `df` is of boolean type which indicates whether the analysis is data flow sensitive or not. `ls` is also an boolean variable, indicating whether the analysis is loop sensitive or not. If `df` is not provided, Algorithm 1 in Chapter 3.2.1 will be used to compute this property. Similarly, if `ls` is not provided, Algorithm 2 in Chapter 3.2.3 will be used to compute this property. `fp` is a variable name of the user defined fixpoint function. A traversal direction is a value from the set `{FORWARD, BACKWARD, ITERATIVE}`, where `FORWARD` is used to represent a forward analysis (predecessors of a node are processed before the node), `BACKWARD` is used to represent a backward analysis (successors of a node are processed before the node), and `ITERATIVE` is used to represent a sequential analysis (visits nodes as they appear in the nodes collection). A user defined fixpoint function can be defined using the `fixp` block:

```
fp := fixp(...) : bool { fbody }
```

In this `fixp` block, `fixp` is a keyword for defining a fixpoint function. A fixpoint function can take any number of parameters, and it must always return a boolean. The body of the fixpoint function is defined in the `fbody` block. A fixpoint function can be assigned a name, which can be passed in the `traverse` expression.

Accessing Facts of Other Nodes. We also provide a special expression `output(n, t)` for querying the traversal output associated with a graph node `n`, in the traversal `t`.

Table 4.1: Syntax reference.

Construct	Syntax	Description
Traversal	<code>t := traversal(n : Node): T { tbody }</code>	<code>t</code> is the name of the traversal that takes a single parameter <code>n</code> representing the graph node that is being visited. A traversal may define a return type <code>T</code> representing the output type. A block of code that generates the traversal output at a graph node is given by <code>tbody</code> .
Traverse	<code>traverse(g, t, d, df, ls, fp)</code>	<code>g</code> is the graph to be traversed, <code>t</code> is the traversal to be invoked, <code>d</code> is the traversal direction and <code>df</code> , <code>ls</code> , <code>fp</code> are optional parameters. <code>df</code> is of boolean type which indicates whether the analysis is data flow sensitive or not. <code>ls</code> is also an boolean variable, indicating whether the analysis is loop sensitive or not. <code>fp</code> is a variable name of the user defined fixpoint function. A traversal direction is a value from the set {FORWARD, BACKWARD, ITERATIVE}
Fixpoint	<code>fp := fixp(...) : bool { fbody }</code>	<code>fixp</code> is a keyword for defining a fixpoint function. A fixpoint function can take any number of parameters, and it must always return a boolean. The body of the fixpoint function is defined in the <code>fbody</code> block.
Output	<code>output(n, t)</code>	<code>output</code> is used for querying the traversal output associated with a graph node <code>n</code> , in the traversal <code>t</code>

Table 4.2: Operations on collections.

Operation	Description
<code>add(C, e)</code>	Adding an element <code>e</code> to collection <code>C</code>
<code>addAll(C1, C2)</code>	Adding all elements from collection <code>C2</code> to collection <code>C1</code>
<code>remove(C, e)</code>	Removing an element <code>e</code> from collection <code>C</code>
<code>removeAll(C1, C2)</code>	Removing all elements from collection <code>C1</code> that are also present in collection <code>C2</code>
<code>get(C, i)</code>	Element at index <code>i</code> from collection <code>C</code> is accessed
<code>has(C, e)</code>	Checking if collection <code>C</code> has element <code>e</code>
<code>equals(C1, C2)</code>	Checking if collection <code>C1</code> and collection <code>C2</code> has the same elements
<code>C1 = C2</code>	Assigning collection <code>C2</code> to collection <code>C1</code>
<code>union(C1, C2)</code>	Returns the union of the elements in collection <code>C1</code> and collection <code>C2</code>
<code>intersection(C1, C2)</code>	Returns the intersection of the elements in collection <code>C1</code> and collection <code>C2</code>

Data Types and Collections. Our system for expressing source code analysis as traversals provides primitive and collection data types. Primitive types include: `bool`, `int`, `string` and collection types include: `Set` and `Seq`, where `Set` is a collection with distinct and unordered elements, whereas, `Seq` is a collection with distinct and ordered elements. A set of operations that can be performed on collection types is described in Table 4.2.

To summarize, we described a system for expressing source code analysis as traversals over graphs using two special constructs: `traversal` for defining a traversal, and `traverse` for invoking a defined traversal. A traversal may visit graph nodes multiple times (in case of fixpoint) and it can be invoked using several parameters specifying the direction of the traversal, a user defined fixpoint function, etc. A traversal output associated with graph nodes can be queried using a

special expression `output()`. To be able to express a variety of source code analysis, our system provides primitive and collection datatypes with well-defined operations. Later in this chapter we demonstrate how the constructs and operations of the system enables determining properties of the source code analysis expressed in our system, such that optimal traversal strategies can be automatically selected.

An Example: Post dominator analysis. We now describe how to use our system to express source code analysis as traversals using an example source code analysis. Post dominator analysis is a backward control flow analysis that collects node ids of all nodes that post dominates every node in the CFG (Aho et al. (2006)). This analysis can be expressed using our system as shown in Listing 4.1.

Listing 4.1: Post dominator analysis: an example source code analysis expressed using our system.

```

1 allNodes: Set<int>;
2 initT := traversal(n: Node) {
3     add(allNodes, n.id);
4 }
5 domT := traversal(n: Node): Set<int> {
6     Set<int> dom;
7     if (output(n, domT) != null) {
8         dom = output(n, domT);
9     } else {
10        if (node.id == exitNodeId) {
11            dom = {};
12        } else {
13            dom = allNodes;
14        }
15    }
16    foreach (s : n.succs)
17        dom = intersection(dom, output(s, domT))
18    add(dom, n.id);
19    return dom;
20 }
```

```

21 fp := fixp(Set<int> curr, Set<int> prev): bool {
22   if>equals(curr, prev))
23     return true;
24   return false;
25 }
26 traverse(g, initT, ITERATIVE);
27 traverse(g, domT, BACKWARD, fp);

```

Listing 4.1 mainly defines two traversals `initT` (lines 2-4) and `domT` (lines 5-20), and invokes them using `traverse` expressions (lines 26 and 27). Line 21-25 defines a fixpoint function using `fixp` block, which is used in the `traverse` expression in line 27. Line 1 defines a variable `allNodes` of collection type `Set`, where `Set<int>` defines a collection type `Set` with elements of type `int`. Line 3 uses an operation `add` (defined in Table 4.2) on collection `allNodes`. The common statements and expressions used in the language to express the analysis are not described in our system, however all standard statements and expressions are allowed. For instance, `if-else` statements are used in lines 7-15, `foreach` iteration is used in lines 16-17, and so on. Lines 26 and 27 provides two flavors of invoking traversals using `traverse` expressions: one without a fixpoint and other with a user-defined fixpoint function. A usage of special expression `output(n, domT)` can be seen in line 8. The traversal `initT` does not define any output for CFG nodes, whereas, the traversal `domT` defines an output of type `Set<int>` for every node in the CFG. For managing the analysis output of nodes, `domT` traversal maintains an internal map that contains analysis output for every node, which can be queried using `output(n, domT)`. A pre-defined variable `g` that represents the CFG is used in the `traverse` expressions in lines 26 and 27.

Figure 4.2 takes an example graph, and shows the results of `initT` and `domT` traversals. Our example graph is a CFG containing seven nodes with a branch and a loop. The `initT` traversal visits nodes sequentially and adds node id to the collection `allNodes`. The `domT` traversal visits nodes in the post-order¹ and computes a set of nodes that post dominate every visited node (as indicated by the set of node ids). For instance, node 7 is post dominated by itself, hence the output at node 7 is $\{7\}$. In Figure 4.2, under `domT` traversal, for each node visited, we show the key intermediate steps indicated by @ line number. These line numbers correspond to the line

¹The traversal strategies chosen for `initT` and `domT` traversals is explained in §4.4.1.

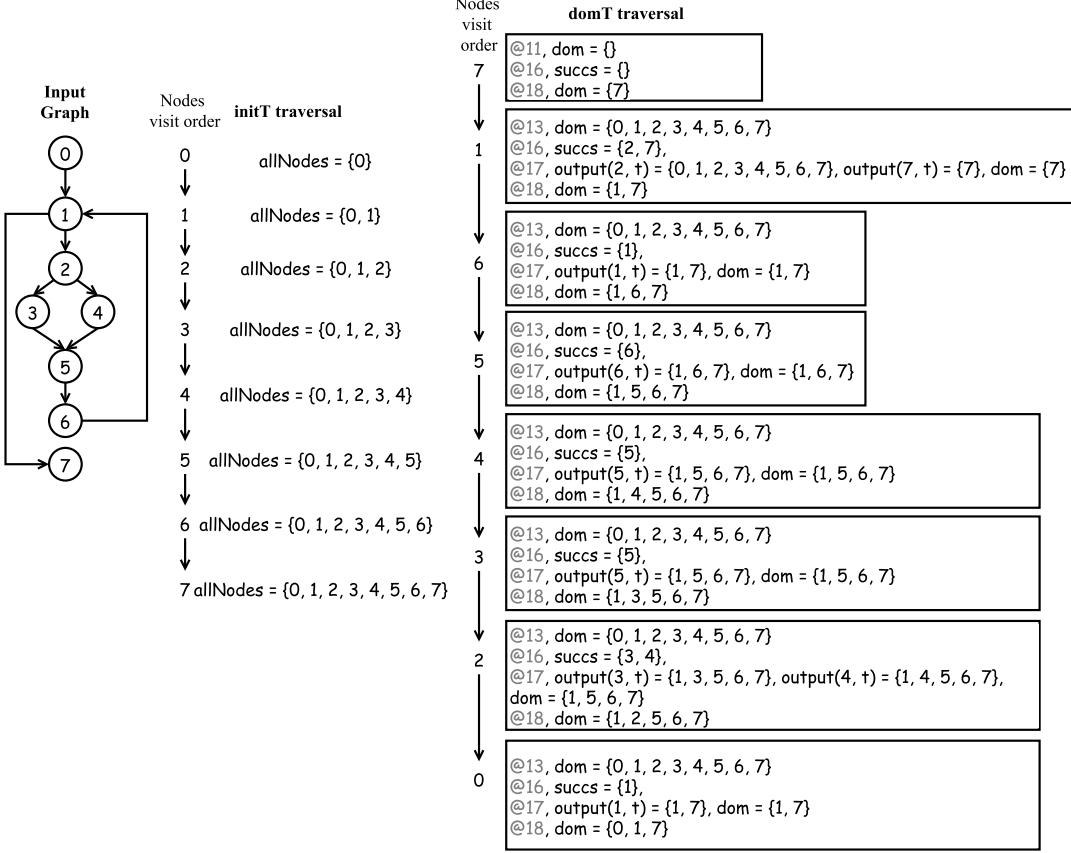


Figure 4.2: Running example of applying the post dominator analysis on an input graph containing branch and loop.

numbers shown in Listing 4.1. We will explain the intermediate results while visiting node 2. In the `domT traversal`, at line 13, the output set `dom` is initialized using `allNodes`, hence `dom` = {0, 1, 2, 3, 4, 5, 6, 7}. At line 16, node 2 has two successors: {3, 4}. At line 17, the set `dom` is updated by performing an `intersection` operation using the outputs of successors 3 and 4. The output of 3 and 4 are {1, 3, 5, 6, 7} and {1, 4, 5, 6, 7} respectively. By performing the intersection of these two sets, the `dom` set becomes {1, 5, 6, 7}. At line 18, the id of the visited node is added to the `dom` set and it becomes {1, 2, 5, 6, 7}. Hence, the post dominator set for node 2 is {1, 2, 5, 6, 7}. Similarly, the post dominator set for other nodes can be calculated.

4.2 Static and Runtime Properties

While it is known in the literature that choosing a right traversal strategy for the source code analysis can significantly improve the performance. For example, Atkinson and Griswold

(2001) have developed a hybridized iterative-worklist algorithm that processes fewer blocks than the traditional algorithm. Despite these advances, how to choose a right traversal strategy, what factors influence the selection of the right traversal strategy, what properties of the analysis and graph are important, and how to determine them, were not known.

In this section, we describe the factors that influence the selection of the optimal traversal strategy for any given traversal. These factors include: the static properties of the analysis and the runtime properties of the graph. We also describe how the challenge of computing these properties is solved with the help of the constructs and operations proposed in our system of expressing source code analysis as traversals (§4.1).

4.2.1 Data-Flow Sensitivity

The *data-flow sensitivity* property of a traversal models the dependence of the traversal outputs of nodes in the input graph. A traversal is data-flow sensitive if the output of a node is computed using the outputs of other nodes. For instance, in the reaching definition data-flow analysis, the outputs of nodes are computed using the outputs of predecessors and in live variable data-flow analysis, the outputs of nodes are computed using the outputs of successors. Hence, both reaching definition and live variable analysis are data-flow sensitive. The term data-flow sensitive is different from the traditional term flow-sensitive, that is used in the literature (Choi et al. (1993)). In literature, a problem is "flow-sensitive" if information about control internal to subroutines is used to compute the final set of data flow facts, while our definition of data-flow sensitive is given below.

Definition 2 $P_{DataFlow}$ (**Data-flow sensitivity**). *Given a traversal t with body $tbody$, a map O that collects and maintains the traversal output of nodes (O is indexed using node ids), and F , a function representing the computation of $O[.]$ in the traversal body $tbody$, if for any node n , its output $O[n]$ is computed by applying F over one or more of $O[n']$, where $n' \neq n$, then t is data-flow sensitive. That is $P_{DataFlow}$ is true otherwise false.*

Algorithm 1: Algorithm to detect data-flow sensitivity

Input: $t := \text{traversal}(n : \text{Node}) : T \{ \text{tbody} \}$

Output: *true/false*

```

1  $A \leftarrow \text{getAliases}(\text{tbody}, n);$ 
2 foreach  $\text{stmt} \in \text{tbody}$  do
3   if  $\text{stmt} = \text{output}(n', t')$  then
4     if  $t' == t$  and  $n' \notin A$  then
5       return true;
6 return false;

```

4.2.2 Computing Data-Flow Sensitivity

To determine the data-flow sensitivity property of a traversal, the operations performed in the traversal needs to be analyzed to check if the output of a node is computed using the outputs of other nodes. In our system of expressing analysis as traversals, the only way to access the output of a node is via `output()` expression, hence given a traversal $t := \text{traversal}(n : \text{Node}) : T \{ \text{tbody} \}$ as input, Algorithm 1 parses the statements in the traversal body `tbody` to identify method calls of the form `output(n', t')` that fetches the output of a node n' in the traversal t' . If such method calls exists, they are further investigated to determine if n' does not point to n and t' points to t . If such method calls exists, it means that the traversal output for the current node n is computed using the traversal outputs of other nodes (n') and hence the traversal is data-flow sensitive. For performing the points to check, Algorithm 1 assumes that an alias environment is computed by using must alias analysis(Jagannathan et al. (1998)). Algorithm 1 requires that the must alias analysis computes all names in the `tbody` that must alias each other at any program point. The must alias information ensures that Algorithm 1 never classifies a data-flow sensitive traversal as data-flow insensitive. A `tbody` may contain more than one `output()` statement, however Algorithm 1 requires only one `output()` statement that fetches the output of other nodes than the current node, to classify the traversal as data-flow sensitive. The control and loop statements in the `tbody` do not have any impact on Algorithm 1 for computing the data-flow sensitivity property.

Worked Out Example . For instance, consider the `domT` traversal shown in Listing 4.1.

Lets apply 1 to determine if `domT` traversal is data-flow sensitive or not. Here `domT` is the `t`. We first compute the aliases of `n`. Only `n` is an alias of `n` (every variable is an alias of itself). Then for each statement in `domT` traversal, we apply line 3-5 in 1. Consider the line 17 of the `domT` traversal. Here, there is a statement `output (s, domT)` which if of the form `output (n', t')`. Here `domT` which is `t'` is equal to `t` and `s` which is `n'` is not an alias of `n`. Hence we return true, indicating `domT` is data-flow sensitive.

Intuitively, in line 16 and 17, a variable `dom` holds the traversal output of a node `n` and it is computed by applying an `intersection` operation on the traversal outputs of successors of `n`. Here, the function `F` is `intersection` over the outputs of all successors of a node.

4.2.3 Loop Sensitivity

The loop sensitivity property models the effect of the loops in the input graph. If an input graph contains loops and if the traversal is affected by the loop, the traversal may require multiple iterations to compute the output of nodes. In the multiple iterations, the traversal outputs of nodes either shrinks or expands to reach a fixpoint. Hence, we define a traversal as loop sensitive, if the traversal output of nodes in subsequent iterations shrinks or expands. The term loop-sensitive is different from the iteration dependences, that is used in the literature (Blume et al. (1996)). In literature, Iteration dependences arise when two different iterations access the same memory location and one or both accesses are a memory write, while our definition of loop-sensitive is given below.

Definition 3 P_{Loop} (**Loop sensitivity**). Given a traversal t , a map O that collects and maintains the traversal output of nodes (O is indexed using node ids), and $O^i[n]$ represents the output of node n in the i^{th} iteration, if $O^{i+1}[n] \lll O^i[n]$ or $O^{i+1}[n] \ggg O^i[n]$, then t is loop sensitive, i.e. P_{Loop} is true otherwise false. The relation \lll represents `shrink` and it is given by, $O^{i+1}[n] \lll O^i[n]$, if $|O^{i+1}[n]| < |O^i[n]|$, and the relation \ggg represents `expand` and it is given by, $O^{i+1}[n] \ggg O^i[n]$, if $|O^{i+1}[n]| > |O^i[n]|$, where $|C|$ is the cardinality of the output collection C .

Since the loop sensitivity property is defined only for data-flow sensitive traversals, we know that the traversal output of nodes in each iteration is computed using the traversal output of other nodes (possible neighbors), we have $O^i[n] = F(O^i[n'])$ and $O^{i+1}[n] = F(O^{i+1}[n'])$, where n, n' are any two nodes such that $n' \neq n$. By substituting these in the `shrink` relation $O^{i+1}[n] \ll O^i[n]$, we get, $F(O^{i+1}[n']) \ll F(O^i[n'])$. For this relation to be `true`, 1) the output of any node n' in any two subsequent iterations i and $i + 1$ must shrink and 2) the function F has the shrink property. Similarly, for expand relation, 3) the output of any node n' in any two iterations i and $i + 1$ must expand and 4) the function F has the expand property. As we know F represents the function in the traversal body that computes the outputs of nodes, if F has the property of shrink or expand, then the traversal can be classified as loop sensitive.

To give an example, consider the `domT` traversal shown in Listing 4.1. Since `domT` is data-flow sensitive, we can check the loop sensitivity property. There are two functions that contributes to the traversal output of any node n in `domT` traversal body. These are `intersection` (line 16) and `add` (line 17). For `domT` to be loop sensitive, we require that both `intersection` and `add` have either shrink or expand property. However, `intersection` has the shrink property and `add` has the expand property, hence we cannot classify `domT` to be loop sensitive.

4.2.4 Computing Loop Sensitivity

In general, computing the loop sensitivity property statically is challenging in the absence of an input graph, however the constructs and operations of our system enables static inference of this property.

A traversal is loop sensitive, if the output of any node in any two subsequent iterations either shrinks or expands. To determine if the traversal output expands or shrinks in the subsequent iterations, the operations performed in the traversal needs to be analyzed. Table 4.2 provides several operations that can be performed on the traversal outputs. The operations `add`, `addAll`, and `union` always expands the output and the operations `remove`, `removeAll`, and `intersection` always shrinks the output.

Given a traversal $t := \text{traversal}(n : \text{Node}) : T \{ \text{tbody} \}$, Algorithm 2 determines the loop sensitivity of t . Algorithm 2 investigates the statements in the `tbody` to determine if the traversal

Algorithm 2: Algorithm to detect loop sensitivity

```

Input: t := traversal(n: Node): T { tbody } 15 foreach stmt  $\in$  tbody do
  Output: true/false 16
  1  $V \leftarrow \{\}$  // a set of output variables related to n; 17
  2  $V' \leftarrow \{\}$  // a set of output variables not related to n; 18
  3 expand  $\leftarrow$  false; 19
  4 shrink  $\leftarrow$  false; 20
  5 gen  $\leftarrow$  false; 21
  6 kill  $\leftarrow$  false; 22
  7  $A \leftarrow getAliases(n);$  23
  8 foreach stmt  $\in$  tbody do 24
    9   if stmt is  $v = output(n', t')$  then 25
      10     if  $t' == t$  then 26
        11       if  $n' \in A$  then 27
          12          $V \leftarrow V \cup v;$  28
        13       else 29
          14          $V' \leftarrow V' \cup v;$  30
    15   if stmt = union( $c_1, c_2$ ) then 31
      16     if ( $c_1 \in V$  and  $c_2 \in V'$ ) || ( $c_1 \in V'$  and  $c_2 \in V$ ) then
        17       expand  $\leftarrow$  true;
    18   if stmt = intersection( $c_1, c_2$ ) then
      19     if ( $c_1 \in V$  and  $c_2 \in V'$ ) || ( $c_1 \in V'$  and  $c_2 \in V$ ) then
        20       shrink  $\leftarrow$  true;
    21   if stmt = add( $c_1, e$ ) || addAll( $c_1, c_2$ ) then
      22     if  $c_1 \in V$  then
        23       gen  $\leftarrow$  true;
    24   if stmt = remove( $c_1, e$ ) || removeAll( $c_1, c_2$ ) then
      25     if  $c_1 \in V$  then
        26       kill  $\leftarrow$  true;
    27   if (expand and gen) || (shrink and kill) then
      28     return true;
    29   else
      30     return false;
  
```

outputs of nodes in multiple iterations either expands or shrinks. For doing that, first it parses the statements to collect all output variables related and not related to input node `n` using the `must alias` information as in Algorithm 1. This is determined in lines 8-14, where all output variables are collected (output variables are variables that gets assigned by the `output` operation) and added to two sets V (a set of output variables related to `n`) and V' (a set of output variables not related to `n`). Upon collecting all output variables, Algorithm 2 makes another pass over all statements in the `tbody` to identify six kinds of operations: `union`, `intersection`, `add`, `addAll`, `remove`, and `removeAll`. These operations are defined in Table 4.2². In lines 16-18, the algorithm looks for `union` operation, where one of the variables involved is an output variables related to `n` and the other variable involved is not related to `n`. These conditions are simply the true conditions for the data-flow sensitivity, where the output of the current node is computed using the outputs of other nodes (neighbors). Similarly, in lines 19-21, the algorithm looks for `intersection` operation. The lines 22-27, identifies add and remove operations that adds or removes elements from the output related to node `n`. Finally, if there exists `union` and `add` operations, the output of a node always expands, and if there exists `intersection` and `remove` operations, the output of a node always shrinks. For a data-flow traversal to be loop sensitive, the output of nodes must either expand or shrink, not both (lines 28-29).

Worked Out Example. Consider the `domT` traversal shown in Listing 4.1. Lets apply 2 to determine if `domT` traversal is loop-sensitive or not. Here `domT` is the `t`. We first compute the aliases of `n`. Only `n` is an alias of `n` (every variable is an alias of itself). Then for each statement in `domT` traversal, we apply line 9-14 in 2 to compute `v` and `v'`. There is only one statement of form `v = output(n', t')` in 4.1 (line 8). Here `dom` maps to `v`, `domT` maps to `t'` and `n` maps to `n`. Since `domT` is also `t`, line 12 in 2 gets executed. At the end of line 14, `v` contains `dom` and `v'` is empty. We once again analyze every statement in `domT` traversal to search for `union`, `intersection`, `add`, `addAll`, `remove`, `removeAll` statement (line 15-27). We have `intersection` (line 17) and `add` (line 18) statement in `domT` traversal, which sets variable `shrink` and `gen` to true. Line 28 in 2 checks if either both `expand` and `kill` are true or `shrink` and `kill` are true. Since both these

²The operations not listed here do not expand or shrink the output.

conditions are not met, the algorithm return `domT` as loop-insensitive.

4.2.5 Graph Cyclicity

So far we have described the two static properties of the analysis that influences the traversal strategy selection. A property of the input graph also influences the selection. This property is the cyclicity in the graph. Based on the cyclicity, we classify graphs into four categories: {sequential, branch only, loop w/o branch, loop w/ branch}. In case of sequential graphs, all nodes in the graph have no more than one successor and predecessor. In case of graphs with branches, nodes may have more than one successor and predecessor. In case of graphs with loops, there exists cycles in the graph. The graph cyclicity is determined during the construction of the graph.

In a source code analysis, traversal output of nodes may depend on each other. For instance, in forward data-flow analysis, output of a node is computed using the outputs of its predecessors. Similarly, in the backward data-flow analysis, output of the successors is required. Graph cyclicity plays an important role in the selection of the appropriate traversal strategy. In case of graphs with branches and loops, the outputs of all dependent nodes of a node (predecessors or successors) may not be available at the time of visiting the node, hence a traversal strategy must be selected that guarantees that the outputs of all dependent nodes of a node are available prior to computing the node's output.

4.3 Traversal Strategies - Candidates

We have picked seven traversal strategies as candidates for choosing an optimal traversal strategy for given a traversal and an input graph. The selected candidate strategies were arrived at by carefully reviewing compilers textbooks, implementations, and source code analysis frameworks. We also made sure that the selected candidate strategies are applicable to any graphs and analysis. We did not consider strategies like chaotic iteration based on Weak Topological Ordering because they are effective only for computing fixed points of continuous function over lattices of infinite height (Bourdoncle (1993)). The selected traversal strategies are describe below:

- **Any order (*ANY*):** In this traversal strategy, nodes can be visited in any order. In our implementation, we visit the nodes in the order they appear in the nodes list N (Definition 1).
- **Increasing order of node ids (*INC*):** In this traversal strategy, the nodes are visited in the increasing order of their node ids. The node ids are assigned during the construction of the graph. For instance, while constructing a CFG, the node ids are assigned in the control flow order.
- **Decreasing order of node ids (*DEC*):** In this traversal strategy, the nodes are visited in the reverse order of their node ids (decreasing order of node ids).
- **Post-Order (*PO*):** In this traversal, the successors of any node are visited before visiting the node.
- **Reverse Post-Order (*RPO*):** In this traversal, the predecessors of any node are visited before visiting the node.
- **Worklist with Post-Order (*WPO*):** In this traversal, the nodes are visited in the order they appear in the worklist. A worklist is a data structure used to keep track of nodes to be visited. In WPO, worklist is initialized with post-ordering of nodes. The worklist is maintained as follows: whenever a node from the worklist is removed and visited, all its successors (for forward traversals) or predecessors (for backward traversals) are added to the worklist as done in Atkinson and Griswold (2001).
- **Worklist with Reverse Post-Order (*WRPO*):** In this traversal, the nodes are visited in the order they appear in the worklist. The worklist is initialized with nodes in the reverse post-order.

4.4 Decision Tree for Traversal Strategy Selection

At this point, we know the factors that influence the traversal strategy selection: the static properties of the analysis, and the runtime property of the graph. Our goal is to check these

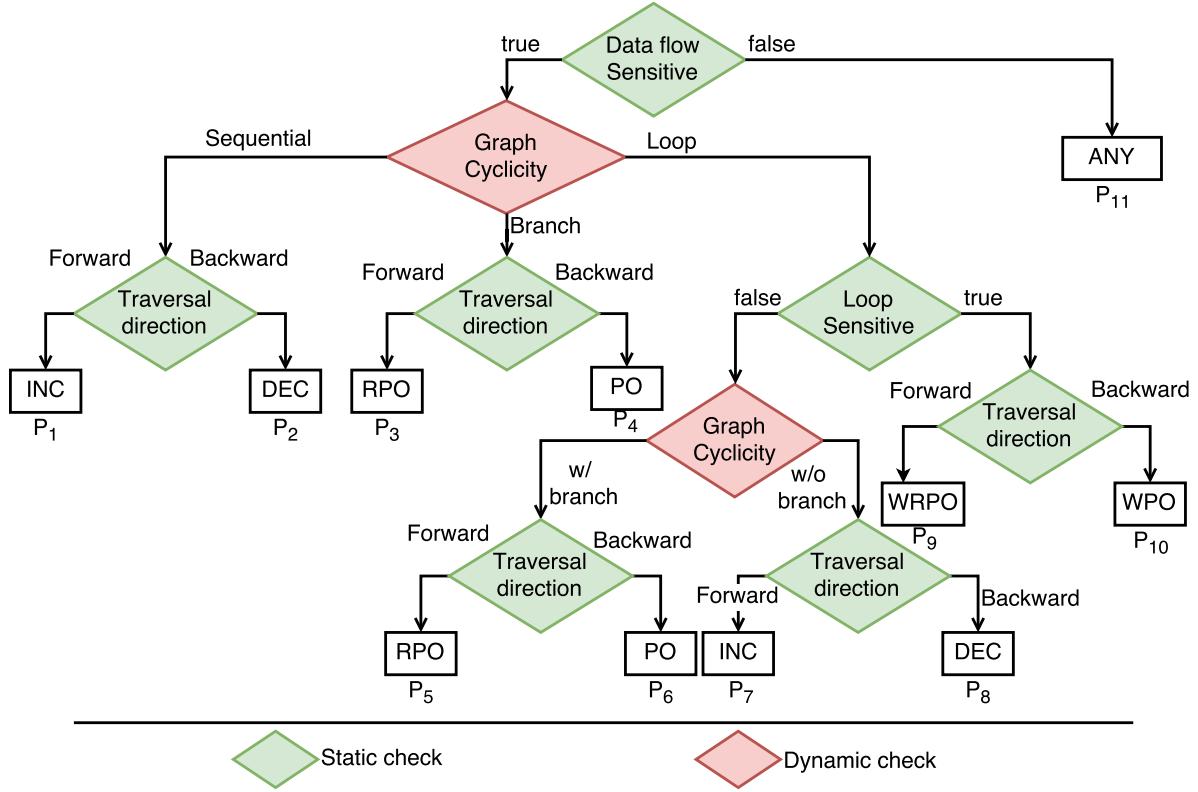


Figure 4.3: Traversal strategy selection decision tree.

properties in certain order to quickly decide the best traversal strategy for a given analysis and a graph, such that only relevant properties are checked and the overhead of static/runtime check is minimized³. To that end, we carefully devised a decision tree as shown in Figure 4.3 for traversal strategy selection.

The leaf nodes of the tree are one of the seven traversal strategies and non-leaf nodes are static/runtime checks. The decision tree has eleven paths marked P_1 through P_{11} . Given a traversal and an input graph, one of the eleven paths will be taken to decide the best traversal strategy. The longest paths P_5 , P_6 , P_7 , and P_8 requires five checks and the shortest path P_{11} requires only one check. The static checks are marked green and the runtime checks are marked red. The static properties that are checked are: data-flow sensitivity ($P_{DataFlow}$), loop sensitivity (P_{Loop}), and traversal direction. The runtime property that is checked is the graph cyclicity: sequential, branch, loop w/ branch, and loop w/o branch. We now provide rationale

³Our evaluation shows that the overhead is less than 0.2% of the total running time for a large dataset and less than 0.01% for an ultra-large dataset.

for arranging the decision tree as shown in Figure 4.3.

The first property that is checked is the data-flow sensitivity of the traversal. This is a static property determined by analyzing the traversal that indicates whether the traversal output of any node is dependent on the traversal output of its neighbors (successors or predecessors). This property is defined in Definition 2 and an algorithm to compute this property is given in Algorithm 1. The rationale for checking this property first is that, if a traversal is data-flow insensitive ($P_{DataFlow}$ is *false*), irrespective of the type of the input graph, the traversal can finish in a single iteration (no fixpoint computation is necessary). In such cases, visiting nodes in any order will be efficient, hence we assign any order (*ANY*) traversal strategy (path P_{11}).

For traversals that are data-flow sensitive ($P_{DataFlow}$ is *true*), further checks are performed to determine the best traversal strategy. Next property that is checked is the input graph cyclicity. This is because the loop sensitive property is applicable to only graphs with loops.

- *Sequential Graphs (paths P_1 and P_2):* In this type of graphs, no branches or loops exists, and all nodes have a single successor and predecessor. At this point, we know that the traversal is data-flow sensitive and it requires output of the neighbors to compute the output for any node. As sequential graphs have only one neighbor (successor or predecessor), a traversal strategy that visits the neighbor prior to visiting any node is sufficient to produce an optimal traversal order. To determine which neighbor (successor or predecessor), we check the traversal direction property. For **FORWARD** traversal direction, predecessor of the node must be visited before the node and for **BACKWARD** traversal direction, successor of the node must be visited before the node. These two traversal orders are provided by our *INC* and *DEC* traversal strategies. The corresponding paths in the decision tree are: P_1 and P_2 .
- *Graphs with branches (paths P_3 and P_4):* In this type of graphs, branches exists, however loops don't exists, which means that a node may have more than one successor or predecessor. At this point, we know that our traversal is data-flow sensitive and it requires output of all neighbors (successors or predecessors) to compute the output for any node, we need a traversal order that ensures that all successors and predecessors are visited

prior to visiting any node. This traversal order is given by the post-order (*PO*) and reverse post-order (*RPO*) traversal strategies. To pick between *PO* and *RPO*, we check the traversal direction. For **FORWARD** traversal direction, we need to visit all predecessors of any node prior to visiting the node, hence we pick the *RPO* traversal strategy. For **BACKWARD** traversal direction, we need to visit all successors of any node prior to visiting the node, hence we pick the *PO* traversal strategy.

- *Graphs with loops (paths P₅ to P₁₀)*: In this type of graphs, loops exists and in addition branches may also exists. We first need to check if the traversal is sensitive to the loop (the loop sensitive property). At this point, we know that our analysis is data-flow sensitive and the input graph has loop based control flow.
 - *Loop sensitive (paths P₉ and P₁₀)*: When the traversal is loop sensitive, for correctly propagating the output, the traversal visits nodes multiple times until a fix point condition is satisfied (user may provide a fix point function). No iterative traversal strategy can guarantee that fix point will be reached in a single traversal of the nodes, hence we adopt a worklist based traversal strategy that visits only required nodes (property of the worklist strategy). The worklist traversal strategy requires that the worklist (a data structure) is initialized with nodes. For picking the best order of nodes for initialization, we further investigate the traversal direction. We know that, for **FORWARD** traversal direction, *RPO* traversal strategy gives the best order of nodes and for **BACKWARD** traversal direction, *PO* traversal strategy gives the best order of nodes, we pick worklist strategy with reverse post-order *WRPO* for **FORWARD** and worklist strategy with post-order *WPO* for **BACKWARD** traversal directions.
 - *Loop insensitive (paths P₅ to P₈)*: When the traversal is loop insensitive, the selection problem reduces to the sequential and branch case that is discussed previously, because for loop insensitive traversal, the loops in the input graph is irrelevant and what remains relevant is the existence of the branches.

4.4.1 An Example

In this section we explain the decision tree using an example source code analysis and a graph. The example analysis that we choose is the *Post Dominator Analysis* shown in Listing 4.1 and the graph that we choose is shown in Figure 4.2. The *Post Dominator Analysis* contains two traversals: `initT` and `domT`. The `initT` traversal is data-flow insensitive ($P_{DataFlow}$ is *false*) and loop insensitive (P_{Loop} is *false*). The `domT` traversal is data-flow sensitive ($P_{DataFlow}$ is *true*) and loop insensitive (P_{Loop} is *false*). The traversal directions of `initT` and `domT` are `ITERATIVE` and `BACKWARD` respectively. Our example graph shown in Figure 4.2 has branches and loops, meaning the graph has nodes with more than one successor or predecessor and has cycles.

For selecting the best traversal strategies for `initT` and the graph with loops and branches, we check the data-flow sensitivity property of the traversal. As `initT` is data-flow insensitive, *ANY* traversal strategy is picked as shown by the path P_{11} in Figure 4.3 and no further checks are required. The traversal strategy *ANY* represents nodes visited in any order.

For selecting the best traversal strategies for `domT` and the graph with branch and loop, we check the data-flow sensitivity property of the traversal. As `domT` is data-flow sensitive, the next property to be checked is the graph cyclicity. As our input graph has loops, the next property to be checked is whether `domT` is loop sensitive ie sensitive to the loops present in the graph. Since `domT` is loop-insensitive, we ignore the loops present in the graph and investigate the rest of the graph structure. As our graph contains branches, the next property to be checked is the traversal direction. The traversal direction for `domT` is `BACKWARD`, we pick *PO* traversal strategy for `domT` traversal, as shown by the path P_6 in Figure 4.3. The traversal strategy post-order (*PO*) visits all successors of a node before visiting that node. This is most suitable for backward analysis like `domT`, because backward analysis analyzes successors of a node prior to analyzing the node.

4.5 Optimizing the Selected Traversal Strategy

Checking the static and dynamic properties not only determines the best traversal strategies, it also helps to perform several optimizations to the selected traversal strategies. In the traversal-

based program analysis with a fixpoint function, in addition to the analysis traversal, two additional traversals of all nodes is required: one traversal to re-compute the analysis results at graph nodes, and another traversal to perform the fixpoint check to ensure that results at nodes have stabilized. The fixpoint check traversal compares the outputs of two traversals of each nodes. The two additional traversals can be eliminated and we formulate them as two optimizations, as described below:

[Opt1] Eliminating the result re-computation traversal: A traversal that re-computes the results at graph nodes for enabling a fixpoint traversal to compare the two results (the analysis traversal result and the re-computation result) can be eliminated, if it is known that results have stabilized and not going to change.

[Opt2] Eliminating the fixpoint check traversal: A fixpoint check traversal that checks the results of the two traversals at graph nodes can be eliminated, if it is known that results are stabilized and not going to change.

The same static and dynamic checks that are performed to determine a traversal strategies also helps to determine if the optimizations can be performed. For instance, consider path P_1 in our decision tree shown in Figure 4.3. This path selects INC traversal strategy that visits nodes in the increasing order of the node ids. While selecting this strategy, we came to know that the analysis is data-flow sensitive (which means the analysis output of neighbors (predecessors or successors) is required for computing the output of a node), the graph is sequential, (which means there is only one successor or predecessor), and the analysis is a forward analysis (which means the output of the predecessor is required to compute the output of a node). Together, we know that if a traversal strategy ensures that the predecessor is visited before visiting any node, the analysis output can be computed in one traversal and no fixpoint check is necessary. The traversal strategy selected for this path is INC and it ensures this property. Hence, both Opt1 and Opt2 can be applied to the selected traversal strategy to further improve the performance. As we show in our evaluation (11), upto 60of the analysis time can be saved by performing these optimizations. In our decision tree shown in Figure 4.3, all paths from P_1 to P_8 are eligible for both the optimizations. The paths P_9 and P_{10} uses a worklist-based traversal strategy

that visits only relevant nodes (relevant nodes are the nodes whose outputs have changed from the last visit) and must perform a fixpoint check every time a node is visited, hence the optimizations cannot be performed. Finally, for path P_{11} , the optimizations are not applicable, because data-flow insensitive traversals requires only one traversal to compute the results and no fixpoint check is performed.

CHAPTER 5. Implementation on Boa framework

5.1 Boa language and infrastructure

Boa is a domain-specific language and infrastructure that eases mining software repositories. Boa's infrastructure leverages distributed computing techniques to execute queries against hundreds of thousands of software projects very efficiently. Boa provide users a domain-specific language tailored for software repository mining and allow them to submit queries via their web-based interface. These queries are then automatically parallelized and executed on a cluster, analyzing a dataset containing almost 700,000 projects, history information from millions of revisions, millions of Java source files, and billions of AST nodes. The language also provides an easy to comprehend visitor syntax to ease writing source code mining queries. The underlying infrastructure contains several optimizations, including query optimizations to make single queries faster as well as a fusion optimization to group queries from multiple users into a single query.

5.2 Source code analysis using Traversal construct

Users must also be able to easily express source code analysis tasks. For users who are intimately familiar with compilers and interpreters, the visitor pattern is well understood. Our traversal pattern is quite similar to it, except that our traversal pattern is for graphs. Without traversal pattern, it generally requires writing a significant amount of boiler-plate code whose length is proportional to the complexity of the source code analysis being written. In this section, we describe how we extended Boa language to support graph traversal and describe our proposed syntax for writing source code analysis. The new syntax is shown in [5.1](#). This syntax is exactly in line with what we described in [4.1](#). The top-level syntax for a analysis task

is a traversal type. During traversal of the graph, the body of the traversal is executed when visiting a node of the graph type. The result of this code is a collection of all node ids in a global variable.

Listing 5.1: Example traversal construct that traverses a CFG and collects each nodes id in a global variable

```
1 allnodeIds := traversal(node: CFGNode) : string {
2     add(cfgnode_ids, string(node.id));
3     return string(node.id);
4 }
```

To begin a analysis task, users write a traverse statement:

```
1 traverse(cfg, TraversalDirection.BACKWARD, TraversalKind.DFS, allnodeIds);
```

that has four parts: the graph to traverse, the traversal direction, the traversal strategy type and a traversal. When this statement executes, a traversal starts at the start node of the cfg using traversal `allnodeIds`.

[5.2](#) shows the implementation code that shows how the traverse calls each node of the cfg using the specified traversal strategy and stores the output of each node. Line 13 - 17 initializes the visit status for all the nodes in the given cfg. Line 18 - 28 applies the specified traversal strategy. At line 2, you can see the variable `outputMap`, which is a map that maintains the output of each node after applying the `allNodeIds` traversal. The traversal can also take a user defined fixpoint function as another parameter and the traversal of the cfg will continue till all its node attain the fixpoint. [5.5](#) shows a traversal that computes post dominators of a node. This traversal requires fixpoint and hence the traverse that calls this traversal contains a user defined fixpoint function `fixp1` as show in [5.3](#).

[5.4](#) shows the user defined fixpoint function.

[5.6](#) lists the implementation code that shows how the traverse calls each node of the cfg using the specified traversal strategy and how the traversal continues till fixpoint is reached. Line 2-3 shows two variables, `outputMapObj` and `prevOutputMapObj`. `outputMapObj` maintains the current output of the nodes ie output of the nodes in the current traversal iteration while `prevOutputMapObj` maintains the previous output of the nodes ie output of the nodes in the

previous traversal iteration. At line 23, you can see that `prevOutputMapObj` is initialized with the values of `outputMapObj`. At this point of time, where the traversal have not yet started, both `outputMapObj` and `prevOutputMapObj` have the same values. Then at line 29, the traversal starts and at line 15, `outputMapObj` gets updated with the current traversal iteration's output. After a iteration of a traversal is completed, line 31-41 applies the user defined fixpoint, which basically checks if both `outputMapObj` and `prevOutputMapObj` have the same values for every node. If not, line 22-41 repeats till the fixpoint is reached. Line 5-7 shows the `getValue` function that returns the current traversal output of the specified node. This function is used in line 11 in [5.5](#).

5.3 Putting it all together

[5.7](#) lists the complete code for post dominator analysis. Lines 15-18 contains the `allNodeIds` traversal that collects all cfg node ids. Line 20 - 44 contains the `cfgPdom` traversal that computes the post dominator. Lines 46-50 contains the user defined fixpoint function `fixp1`. Lines 52-61 is the boa visitor. Line 53 - 59 contains a set of operation that is done before visiting each method in the input dataset. At line 54, `getcfg` method is used which returns the cfg od the given method. Line 57 contains the traverse call to `allNodeIds` with cfg, traversal direction and traversal strategy. Line 58 contains the traverse call to `cfgPdom` with additional fixpoint parameter.

After these two traversal is complete, the `cfgPdom` contains the post dominator set for each node. You can query it using `getValue` function by supplying `cfgPdom` and the node you want to get the post-dominator set. [5.8](#) shows the code snippet to get each node's output of `cfgPdom` traversal.

Listing 5.2: Implementation code of traversal without fixpoint.

```

1 public final void dfsBackward(final CFGNode node, java.util.HashMap<Integer,
2     String> nodeVisitStatus) throws Exception {
3     outputMap.put(node.id, allnodeIds(node));
4     nodeVisitStatus.put(node.getId(), "visited");
5     for (CFGNode pred : node.getPredecessorsList()) {
6         if (nodeVisitStatus.get(pred.getId()).equals("unvisited"))
7             dfsBackward(pred, nodeVisitStatus);
8     }
9 }
10 public final void traverse(final boa.graphs.cfg.CFG cfg, final Traversal.
11     TraversalDirection direction, final Traversal.TraversalKind kind) throws
12     Exception {
13     try {
14         if (preTraverse(cfg)) {
15             java.util.HashMap<Integer, String> nodeVisitStatus=new java.util.
16                 HashMap<Integer, String>();
17             CFGNode[] nl = cfg.nodes();
18             for(int i=0;i<nl.length;i++) {
19                 nodeVisitStatus.put(nl[i].getId(), "unvisited");
20             }
21             switch(kind) {
22                 case TraversalKind.DFS:
23                     switch(direction) {
24                         case TraversalDirection.BACKWARD :
25                             dfsBackward(cfg.getExitNode(), nodeVisitStatus);
26                             break;
27                         case TraversalDirection.FORWARD :
28                             ....
29                         }
30                     }
31             } catch(Exception e) {return;}
32     }

```

Listing 5.3: Traverse call with fixpoint function.

```

1 traverse(cfg, TraversalDirection.BACKWARD, TraversalKind.POSTORDER, cfgPdom,
2     fixp1);

```

Listing 5.4: User defined fixpoint function.

```

1 fixp1 := fixp(curr, prev: T) : bool {
2     if (len(difference(curr, prev)) == 0)
3         return true;
4     return false;
5 };

```

Listing 5.5: Example traversal construct that computes post dominator.

```

1 cfgPdom := traversal(node: CFGNode): T {
2     cur_value : T;
3     if(node.id==exitId) {
4         self_dom: set of string;
5         cur_value = self_dom;
6     }
7     else
8         cur_value = cfgnode_ids;
9
10    if(def(getvalue(node))) {
11        cur_value = getvalue(node);
12    }
13
14    preds:=node.successors;
15    foreach(i:int;def(preds[i])) {
16        pred_value := getvalue(preds[i]);
17        if(def(pred_value)) {
18            cur_value = intersect(cur_value,pred_value);
19        }
20    }
21
22    gen_kill := getvalue(node, allnodeIds);
23    if(def(gen_kill)) {
24        add(cur_value, gen_kill);
25    }
26
27    return cur_value;
28 };

```

Listing 5.6: Implementation code of traversal with fixpoint.

```

1  public abstract class BoaAbstractTraversal<T1> {
2      public java.util.HashMap<Integer,T1> outputMapObj;
3      public java.util.HashMap<Integer,T1> prevOutputMapObj;
4
5      public T1 getValue(final CFGNode node) throws Exception {
6          return (T1)outputMapObj.get(node.getId());
7      }
8
9      public final void postorderBackward(final CFGNode node, java.util.HashMap<
10         Integer,String> nodeVisitStatus) throws Exception {
11         nodeVisitStatus.put(node.getId(),"visited");
12         for (CFGNode succ : node.getSuccessorsList()) {
13             if (nodeVisitStatus.get(succ.getId()).equals("unvisited"))
14                 postorderBackward(succ, nodeVisitStatus);
15             outputMapObj.put(node.id, cfgPdom(node));
16         }
17     }
18     public final void traverse(final boa.graphs.cfg.CFG cfg, final Traversal.
19         TraversalDirection direction, final Traversal.TraversalKind kind, final
20         BoaAbstractFixP fixp) {
21         switch(kind.getNumber()) {
22             case 12 :
23                 boolean fixp_flag;
24                 do {
25                     prevOutputMapObj = new java.util.HashMap<Integer,T1>(
26                         outputMapObj);
27                     java.util.HashMap<Integer,String> nodeVisitStatus=new java.util.
28                         .HashMap<Integer,String>();
29                     CFGNode[] nl = cfg.nodes();
30                     for(int i=0;i<nl.length;i++) {
31                         nodeVisitStatus.put(nl[i].getId(),"unvisited");
32                     }
33                     postorderBackward(cfg.getEntryNode(), nodeVisitStatus);
34                     fixp_flag=true;
35                     for(CFGNode node : nl) {
36                         boolean curFlag=outputMapObj.containsKey(node.getId());
37                         boolean prevFlag=prevOutputMapObj.containsKey(node.getId())
38                             ;
39                         if(curFlag) {
40                             if(outputMapObj.containsKey(node.getId()) &&
41                                 prevOutputMapObj.containsKey(node.getId())) {
42                                 fixp_flag=fixp_flag && fixp.invoke((T1)outputMapObj
43                                     .get(node.getId()),(T1)prevOutputMapObj.get(
44                                         node.getId()));
45                             } else {
46                                 fixp_flag = false; break;
47                             }
48                         }
49                     }
50                 }while(!fixp_flag);
51                 break;
52             default : break;
53         }
54     }
55 }
```

Listing 5.7: Code to compute post dominator using Boa language.

```

1 p: Project = input;
2 cfgnode_ids:set of string;
3 type T = set of string;
4 exitId : int;
5 cfg: CFG;
6 m: output collection[string] of string;
7
8 #clones set of string
9 clone := function(data: T) : T {
10     s: set of string;
11     s = union(s, data);
12     return s;
13 };
14
15 allnodeIds := traversal(node: CFGNode) : string {
16     add(cfgnode_ids, string(node.id));
17     return string(node.id);
18 };
19
20 cfgPdom := traversal(node: CFGNode): T {
21     cur_value : T;
22     if(node.id==exitId) {
23         self_dom:set of string;
24         cur_value = self_dom;
25     }
26     else
27         cur_value = cfgnode_ids;
28     if(def(getvalue(node))) {
29         cur_val1 := getvalue(node);
30         cur_value = clone(cur_val1);
31     }
32
33     preds:=node.successors;
34     foreach(i:int;def(preds[i])) {
35         pred_value := getvalue(preds[i]);
36         if(def(pred_value))
37             cur_value = intersect(cur_value,pred_value);
38     }
39     gen_kill := getvalue(node, allnodeIds);
40     if(def(gen_kill))
41         add(cur_value, gen_kill);
42
43     return cur_value;
44 };
45
46 fixp1 := fixp(curr, prev: T) : bool {
47     if (len(difference(curr, prev)) == 0)
48         return true;
49     return false;
50 };
51
52 controlDependenceVisitor:= visitor {
53     before node: Method -> {
54         cfg = getcfg(node);
55         exitId = len(cfg.nodes) - 1;
56         clear(cfgPdom); clear(cfgnode_ids); clear(allnodeIds);
57         traverse(cfg, TraversalDirection.BACKWARD, TraversalKind.Hybrid,
58                  allnodeIds);
59         traverse(cfg, TraversalDirection.BACKWARD, TraversalKind.Hybrid,
60                  cfgPdom, fixp1);
61     }
62 };
63 visit(p, controlDependenceVisitor);

```

Listing 5.8: Code to query post dominator set of a node.

```
1 result := traversal(node: CFGNode)  {
2     cur_val := getvalue(node, cfgPdom);
3     if(def(cur_val))
4         m[string(node.id)] << string(cur_val);
5 }
```

CHAPTER 6. Empirical Evaluation

We conducted an empirical evaluation on a set of 21 basic source code analyses and 2 public massive code datasets to evaluate several factors about our hybrid approach for selecting and optimizing traversal strategies. First, we show the benefit of using our optimized selected traversal over standard ones by evaluating the *reduction in running times* of our hybrid approach over the standards ones (§7). Then, we evaluate the correctness of the analysis results using our hybrid approach to show that the decision analyses and optimizations in our approach does not affect the correctness of the source code analyses (§8). We also evaluate the precision of our selection algorithm by measuring how often the hybrid approach selects the most time-efficient traversal (§9). Finally, we evaluate how the different components in the approach and different kinds of static and runtime properties impact the overall performance? This is done in §10 and §11, and via various insight analysis of our results.

6.1 Analyses, Datasets and Experiment Setting

6.1.1 Analyses.

We collected source code analyses that traverse control flow graphs from textbooks and source code analysis tools. We also made sure that the analyses list covers all the static properties discussed in §4.2, i.e., data-flow sensitivity, loop sensitivity and traversal direction. We ended up with 21 source code analyses as shown in Table 6.1. They include 10 basic ones (analyses 1, 2, 8, 9, 10, 11, 12, 14, 15 and 19) from textbooks (Aho et al. (2006); Nielson et al. (2010)) and 11 tothers for detecting source code bugs, and code smells from the Soot framework (Vallée-Rai et al. (1999)) (analyses 3, 4, 5, 13, 17 and 18), and FindBugs tool (Ayewah et al. (2007)) (analyses 6, 7, 16, 20 and 21). Table 6.1 also shows the number of traversals each

Table 6.1: List of source code analyses and the properties of their involved traversals.

Analysis	Ts	t_1			t_2			t_3		
		Flw	Lp	Dir	Flw	Lp	Dir	Flw	Lp	Dir
1 Copy propagation (CP)	3	\times	\times	—	✓	✓	→	\times	\times	—
2 Common sub-expression detection (CSD)	3	\times	\times	—	✓	✓	→	\times	\times	—
3 Dead code (DC)	3	\times	\times	—	✓	✓	←	\times	\times	—
4 Loop invariant code (LIC)	3	\times	\times	—	✓	✓	→	\times	\times	—
5 Upsafety analysis (USA)	3	\times	\times	—	✓	✓	→	\times	\times	—
6 Valid FileReader (VFR)	3	\times	\times	—	✓	✓	→	\times	\times	—
7 Mismatched wait/notify (MWN)	3	\times	\times	—	✓	✓	→	\times	\times	—
8 Available expression (AE)	2	\times	\times	—	✓	✓	→			
9 Dominator (DOM)	2	\times	\times	—	✓	\times	→			
10 Local may alias (LMA)	2	\times	\times	—	✓	✓	→			
11 Local must not alias (LMNA)	2	\times	\times	—	✓	✓	→			
12 Live variable (LV)	2	\times	\times	—	✓	✓	←			
13 Nullness analysis (NA)	2	\times	\times	—	✓	✓	→			
14 Post dominator (PDOM)	2	\times	\times	—	✓	\times	←			
15 Reaching definition (RD)	2	\times	\times	—	✓	✓	→			
16 Resource status (RS)	2	\times	\times	—	✓	✓	→			
17 Very busy expression (VBE)	2	\times	\times	—	✓	✓	←			
18 Safe Synchronization (SS)	2	\times	\times	—	✓	✓	→			
19 Used and defined variable (UDV)	1	\times	\times	—						
20 Useless increment in return (UIR)	1	\times	\times	—						
21 Wait not in loop (WNIL)	1	\times	\times	—						

Table 6.2: Statistics of the generated control flow graphs from two datasets.

Dataset	All graphs	Sequential	Branches	Loops		
				All graphs	Branches	No branches
DaCapo	287K	186K (65%)	73K (25%)	28K (10%)	21K (7%)	7K (2%)
GitHub	161,523K	111,583K (69%)	33,324K (21%)	16,617K (10%)	11,674K (7%)	4,943K (3%)

analysis contains and their static properties as described in §4.2. The sets of traversals cover all types of static properties for flow-sensitivity, loop-sensitivity and direction (forward, backward and iterative). All analyses are intra-procedural. We implemented all twenty one of these analysis using Boa language(Dyer et al. (2013a)). In Table 6.1, Ts denotes total number of traversals, t_i denotes properties of traversal i -th, Flw denotes data-flow sensitive, Lp denotes loop sensitive, Dir denotes traversal direction where —, \rightarrow and \leftarrow mean iterative, forward and backward, respectively.

6.1.2 Datasets.

We ran the analyses on two datasets: DaCapo 9.12 benchmark (Blackburn et al. (2006)), DaCapo for short, and a ultra-large-scale dataset containing projects from GitHub from Boa. DaCapo dataset contains the source code of 10 open source Java projects: Apache Batik, Apache

Table 6.3: Time contribution of each phase (in miliseconds).

Analysis	Avg. Time		Static	Runtime				
	DaCapo	GitHub		DaCapo		GitHub		
				Avg.	Total	Avg.	Total	
CP	0.21	0.008	53	0.21	62,469	0.008	1359K	
CSD	0.19	0.012	60	0.19	56,840	0.012	1991K	
DC	0.19	0.010	45	0.19	54,822	0.010	1663K	
LIC	0.21	0.006	69	0.20	60,223	0.006	992K	
USA	0.19	0.006	90	0.19	54,268	0.009	1444K	
VFR	0.18	0.007	42	0.18	52,483	0.007	1142K	
MWN	0.18	0.006	36	0.18	52,165	0.006	1109K	
AE	0.18	0.007	43	0.18	53,290	0.007	1169K	
DOM	0.21	0.008	35	0.21	62,416	0.008	1307K	
LMA	0.18	0.008	76	0.18	52,483	0.008	1346K	
LMNA	0.18	0.008	80	0.18	53,182	0.008	1407K	
LV	0.17	0.007	32	0.17	49,231	0.007	1273K	
NA	0.16	0.008	64	0.16	46,589	0.008	1398K	
PDOM	0.20	0.012	34	0.20	57,203	0.012	2040K	
RD	0.20	0.007	48	0.20	57,359	0.007	1155K	
RS	0.16	0.006	28	0.16	46,367	0.006	996K	
VBE	0.17	0.006	44	0.17	49,138	0.006	1062K	
SS	0.17	0.006	32	0.17	48,990	0.006	1009K	
UDV	0.14	0.005	10	0.14	41,617	0.005	928K	
UIR	0.14	0.006	14	0.14	41,146	0.006	1020K	
WNIL	0.14	0.007	15	0.14	41,808	0.007	1210K	

FOP, Apache Aurora, Apache Tomcat, Jython, Xalan-Java, PMD, H2 database, Sunflow and Daytrader. GitHub dataset contains the source code of more than 380K Java projects collected from GitHub.com. Each method in the datasets was used to generate a control flow graph (CFG) on which the analyses would be run. The statistics of the two datasets are shown in Table 6.2. DaCapo dataset contains 287K non-empty CFGs while GitHub dataset contains more than 162M. Both have similar distributions of CFGs over graph cyclicity. Most CFGs are sequential and only 10% have loops.

6.1.3 Setting.

We compared our *hybrid* approach against the six standard traversal strategies in §4.3: DFS, PO, RPO, WPO, WRPO and ANY. The running time for each analysis is measured from the start to the end of the analysis which includes constructing CFGs and traversing CFGs. The running time for our hybrid approach also includes the time for computing the static and runtime properties, making the traversal strategy decision, optimizing it and then using the optimized traversal strategy to traverse the CFG and run the analysis. The analyses on DaCapo dataset were run on a single machine with 24 GB of memory and 24-cores, running on Linux

3.5.6-1.fc17 kernel. Running analyses on GitHub dataset on a single machine would take weeks to finish, so analyses were executed on a Hadoop [Apache Software Foundation 2015a] install with 1 name node, 1 job tracker node, and 10 compute nodes. The compute nodes have a total of 148 CPU cores and 2GB memory per core. All machines run Ubuntu 12.04LTS. The cluster has been tuned for performance, including setting the maximum number of map tasks for each compute node equal to the number of cores on that node, increasing the VM heap size to 1GB per task, and enabling short-circuit local reads in the distributed filesystem.

CHAPTER 7. Running Time and Time Reduction

We first report the running times and then study the achieved reductions against standard traversal strategies.

7.1 Running Time

Table 6.3 shows the running times for 21 analyses on the two datasets. On average (column **Avg. Time**), each analysis took 0.14–0.21 ms and 0.005–0.012 ms to analyze a graph in Dacapo and GitHub datasets, respectively. The variation in the average analysis time is mainly due to the difference in the machines used to run the analysis for DaCapo and GitHub datasets, as described in Chapter 4.1. Also, the average sizes of graphs in DaCapo are much larger than the average sizes of the graphs in the GitHub. Columns **Static** and **Runtime** show the time contributions for different components of the hybrid approach: the time for determining the static properties of each analysis which is done once for each analysis, and the time for constructing the CFG of each method and traversing the CFG which is done once for every constructed CFG. We can see that the time for collecting static information is negligible, less than 0.2% for DaCapo dataset and less than 0.01% for GitHub dataset, when compared to the total runtime information collection time, as it is performed only once per traversal. When compared to the average runtime information collection time, the static time is quite significant. However, the overhead introduced by static information collection phase diminishes as the number of CFGs increases and becomes insignificant when running on those two large datasets. This result shows the benefit of our hybrid approach when applying on ultra-large-scale analysis.

Analysis	DaCapo						GitHub					
	DFS	PO	RPO	WPO	WRPO	ANY	DFS	PO	RPO	WPO	WRPO	ANY
CP	17%	83%	9%	66%	11%	72%	17%	88%	12%	80%	5%	82%
CSD	41%	93%	39%	74%	4%	89%	31%	—	24%	—	12%	—
DC	41%	30%	89%	7%	64%	81%	25%	22%	—	7%	—	—
LIC	17%	84%	8%	67%	7%	73%	19%	89%	15%	81%	19%	88%
USA	36%	92%	34%	72%	9%	87%	22%	—	17%	—	9%	—
VFR	20%	41%	18%	51%	15%	62%	15%	40%	10%	44%	9%	53%
MWN	21%	35%	16%	35%	22%	49%	17%	31%	12%	33%	11%	46%
AE	40%	14%	39%	73%	14%	87%	16%	—	16%	—	11%	—
DOM	54%	97%	48%	70%	6%	95%	27%	—	32%	—	6%	—
LMA	35%	46%	28%	74%	6%	46%	22%	—	13%	—	6%	—
LMNA	29%	39%	22%	68%	9%	41%	21%	—	15%	—	7%	—
LV	38%	30%	84%	11%	56%	75%	25%	21%	68%	11%	69%	80%
NA	26%	88%	30%	50%	10%	80%	13%	87%	12%	71%	10%	85%
PDOM	51%	41%	95%	10%	72%	95%	24%	20%	—	24%	—	—
RD	15%	80%	7%	62%	9%	68%	19%	91%	10%	79%	5%	86%
RS	31%	31%	30%	31%	28%	30%	16%	40%	9%	31%	7%	49%
VBE	40%	36%	88%	13%	76%	81%	28%	24%	—	10%	—	—
SS	26%	39%	22%	37%	25%	57%	20%	35%	13%	34%	10%	50%
UDV	6%	5%	6%	10%	9%	3%	3%	4%	2%	7%	6%	0%
UIR	2%	2%	1%	3%	3%	0%	2%	5%	4%	7%	7%	0%
WNIL	3%	4%	5%	6%	8%	2%	3%	6%	5%	5%	6%	0%
Overall	31%	83%	70%	55%	35%	81%	—	—	—	—	—	—

(a) Time reduction for each analysis.

Property	DaCapo						GitHub					
	DFS	PO	RPO	WPO	WRPO	ANY	DFS	PO	RPO	WPO	WRPO	ANY
Data-flow	32%	84%	72%	57%	36%	83%	20%	74%	63%	55%	28%	72%
¬Data-flow	4%	4%	4%	6%	6%	2%	31%	81%	66%	58%	40%	92%

(b) Overall reduction over analysis properties.

Property	DaCapo					
	DFS	PO	RPO	WPO	WRPO	ANY
Sequential	20%	74%	63%	55%	28%	72%
Branch	31%	81%	66%	58%	40%	92%
Loop	53%	88%	75%	62%	37%	95%

(c) Overall reduction over graph properties.

Figure 7.1: Reduction in running times.

7.2 Time Reduction

To evaluate the efficiency in running time of the hybrid approach over other traversal strategies, we ran the 21 analyses on DaCapo and GitHub datasets using hybrid approach and other candidate traversals. When comparing the hybrid approach to a standard strategy S , we computed the reduction rate $R = (T_S - T_H)/T_S$ where T_S and T_H are the running times using the standard and the hybrid strategy, respectively. Some analyses have some worst case traversal strategies which might not be feasible to run on dataset at the scale of 162 million graphs as in GitHub dataset. For example, using post-order for forward data-flow analysis will visit the CFG in the direction which is opposite to the natural direction of the analysis and hence takes a lot of time to complete the analysis. For such combinations of analyses and

traversal strategies, the map and the reducer tasks time out in the cluster setting and, thus, we did not provide the running times. The corresponding cells in Figure 7.1a are denoted with symbol –.

The result in Figure 7.1a shows that the hybrid approach helps reduce the running times in almost all cases. The values indicates the reduction in running time by adopting hybrid approach compared against the standard strategies. Most of positive reductions are from 10% (light yellow cells) or even from 50% (light green cells). More importantly, the most time-efficient and the worst traversal strategies vary across the analyses which supports the need of our hybrid traversal strategy. Over all analyses, the reduction was highest against any order and post-order (PO and WPO) strategies. The reduction was lowest against the strategy using depth-first search (DFS) and worklist with reverse post-ordering (WRPO). When compared with the next best performing traversal strategy for each analysis, hybrid approach reduces the overall execution time by about 13 minutes to 72 minutes on GitHub dataset. We do not report the overall numbers for GitHub dataset due to the presence of failed runs.

Figure 7.1b shows time reductions for different types of analyses. For *data-flow sensitive* ones, the reduction rates are high ranging from 32% to 84%. The running time was not improved much for *non data-flow sensitive* traversals, which correspond to the last three rows in Figure 7.1a with mostly one digit reductions (light orange cells). We actually perform almost as same as Any order traversal strategy for analyses in this category. This is because Any order traversal strategy is the best strategy for all the CFGs in these analyses. Hybrid approach also chooses any order traversal strategy and hence there is no scope for performance gain.

Figure 7.1c shows time reduction for different cyclicity types of input graphs. We can see that reductions over graphs with loops is highest and those over any graphs is lowest.

7.3 Time reduction against hand optimized analysis

Another way to extract more performance is to hand optimize the analysis. Figure 7.2 compares Hybrid approach against hand optimized analysis. Hand optimized analysis has single best optimized traversal strategy applied for each analysis. For data-flow analysis, hand optimized analysis uses WPO/WRPO guideline while for non data flow analysis, it uses ANY

Analysis	DaCapo		GitHub
CP	9%		5%
CSD	4%		12%
DC	7%		7%
LIC	7%		19%
USA	9%		9%
VFR	15%		9%
MWN	15%		9%
AE	14%		11%
DOM	6%		6%
LMA	6%		6%
LMNA	9%		7%
LV	11%		11%
NA	10%		10%
PDOM	10%		24%
RD	9%		5%
RS	28%		7%
VBE	13%		10%
SS	22%		10%
UDV	3%		0%
UIR	0%		0%
WNIL	2%		0%

Figure 7.2: Reduction in running times against hand optimized analysis.

traversal strategy, as it is the best traversal strategy for non data-flow analysis. WPO/WRPO guideline suggests that if the direction of the traversal is backward, use WPO else use WRPO. This guideline simplifies the hybrid approach, where the decision is based on only traversal direction while Hybrid approach uses the decision tree in Figure 4.3. Figure 7.2 shows the comparison of hybrid approach against hand optimized analysis. We can see that for about half of the data-flow analysis, we gain at least 10% reduction. For analysis like RS, SS, VFR and MWN, the gain is much higher since these analyses are selective in terms of the program statements that they analyze. WPO and WRPO would not work for these analyses since in case of both WPO and WRPO, there exists a fixed cost of creating and maintaining a worklist of size equals to the number of CFG nodes. When analyses selectively analyzes a small subset of all nodes in the CFGs, this overhead becomes substantial. Our hybrid strategy is not only able to select an alternative strategy instead of WPO/WRPO for such selective analyses, whenever possible (when the input graphs do not contain loops), but also optimize WPO/WRPO (in case of input graphs with loops), such that the overheads are minimized. For non-data flow analysis, there is no gain against hand optimized analysis since ANY is the best traversal strategy for such analysis and both Hybrid approach and hand optimized analysis applied ANY traversal strategy for all the graphs for non-data flow analysis.

CHAPTER 8. Correctness of Analysis Results

To evaluate the correctness of analysis results, we first chose worklist as standard strategy to run analyses on DaCapo dataset to create the groundtruth of the results. We then ran analyses using our hybrid approach and compared the results with the groundtruth. In all analyses on all input graphs from the dataset, the results from our hybrid approach always exactly matched the corresponding ones in the groundtruth.

CHAPTER 9. Traversal Strategy Selection Precision

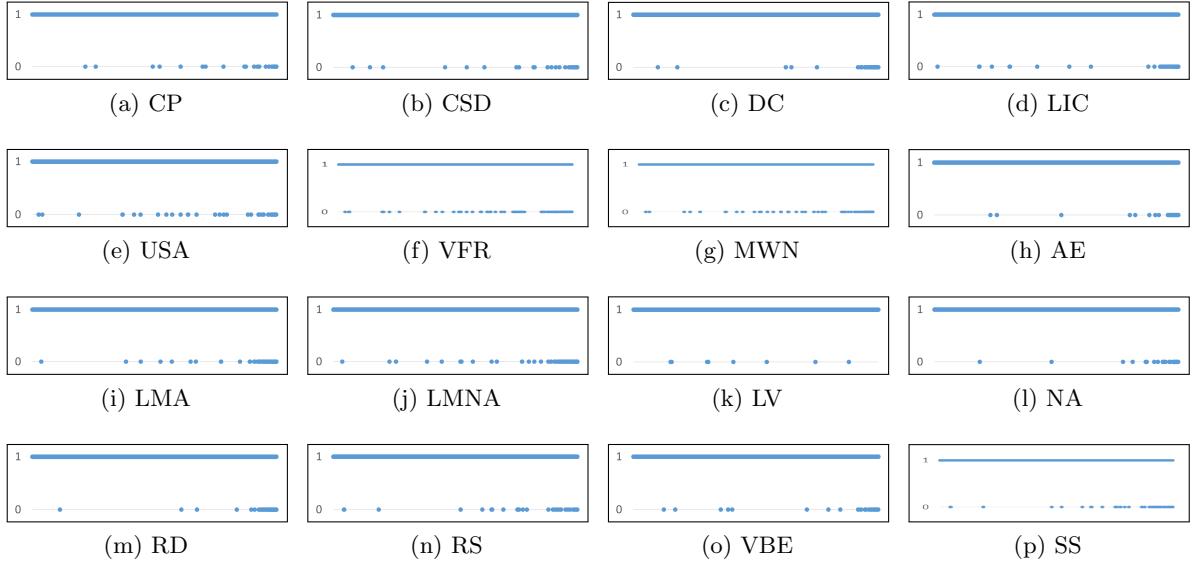


Figure 9.1: Scatter charts for analyses that have loop sensitive traversals.

In this experiment, we evaluated how well the hybrid approach picks the most time-efficient strategy. We ran the 21 analyses on the DaCapo dataset using all the candidate traversals and the one selected by the hybrid approach. One selection is counted for each pair of a traversal and an input graph where the hybrid approach selects a traversal strategy based on the properties of the analysis and input graph. A selection is considered correct if its running time is at least as good as the running time of the fastest among all candidates. The precision is computed as the ratio between the number of correct selections over the total number of all selections. As

Table 9.1: Traversal strategy prediction precision.

Analysis	Precision
DOM, PDOM, WNIL, UDV, UIR	100.00%
CP, CSD, DC, LIC, USA, VFR, MWN, AE, LMA, LMNA, LV, NA, RD, RS, VBE, SS	99.99%

shown in Table 9.1, the selection precision is 100% for all analyses that are not *loop sensitive*. For analyses that involve *loop sensitive* traversals, the prediction precision is 99.99%.

We further analyzed the result to see what contributed to these mispredictions. Let us break the CFGs in the DaCapo dataset by the graph cyclicity: sequential CFGs, CFGs with branches and no loops, and CFGs with loops, and discuss the selection precision.

For sequential CFGs & CFGs with branches and no loops, the selection precision is 100%—the hybrid approach always picks the most time-efficient traversal strategy.

For CFGs with loops, the selection precision is 100% for *loop insensitive* traversals. The mispredictions occur with *loop sensitive* traversals on CFGs with loops. Figure 9.1 shows scatter charts for the traversal selection results for 16 analyses that are *loop sensitive*. In the chart, 1 indicates a correct selection and 0 indicates a misprediction. CFGs are organized along the x -axis in the increasing order of their sizes measured as the numbers of nodes. The scatter charts show that the mispredictions tend to happen with larger CFGs. The reason is that, for *loop sensitive* traversals, the hybrid approach picks worklist as the best strategy. The worklist approach was picked because it visits only as many nodes as needed when compared to other traversal strategies which visit redundant nodes. However using worklist imposes an overhead of creating and maintaining a worklist containing all nodes in the CFG. This overhead is negligible for small CFGs. However, when running analyses on large CFGs, this overhead could become higher than the cost for visiting redundant nodes. Therefore, selecting worklist for *loop sensitive* traversals on large CFGs might not always result in the best running times.

Figure 9.2 shows the Hybrid approach’s performance against best approaches for mispredicted graphs for 16 analyses that has *loop sensitive* traversals. We can see that for majority of the mis-predicted graphs, Hybrid approach’s performance is comparable to the best approaches. For analyses like CP, CSD, LIC, USA, DC, LMA, AE, LMNA, VBE and RD, Hybrid running time are almost similar to the Best traversal strategy for smaller graphs. For graphs of size greater than 200, Hybrid starts to perform worse (1.5x - 2x) than the best traversal strategy. But for analyses like MWN, VFR, LV, NA, RS and SS, Hybrid performs worse even for smaller graphs that were mis-predicted. This is because of the complexity of the analyses. Since these analyses are lesser in complexity, even for small graphs that are mis-predicted, hybrid’s

performance is worse. This is because hybrid chooses Worklist for these mis-predicted analyses and in these analyses, time to maintain the Worklist is much higher than the analyses itself.

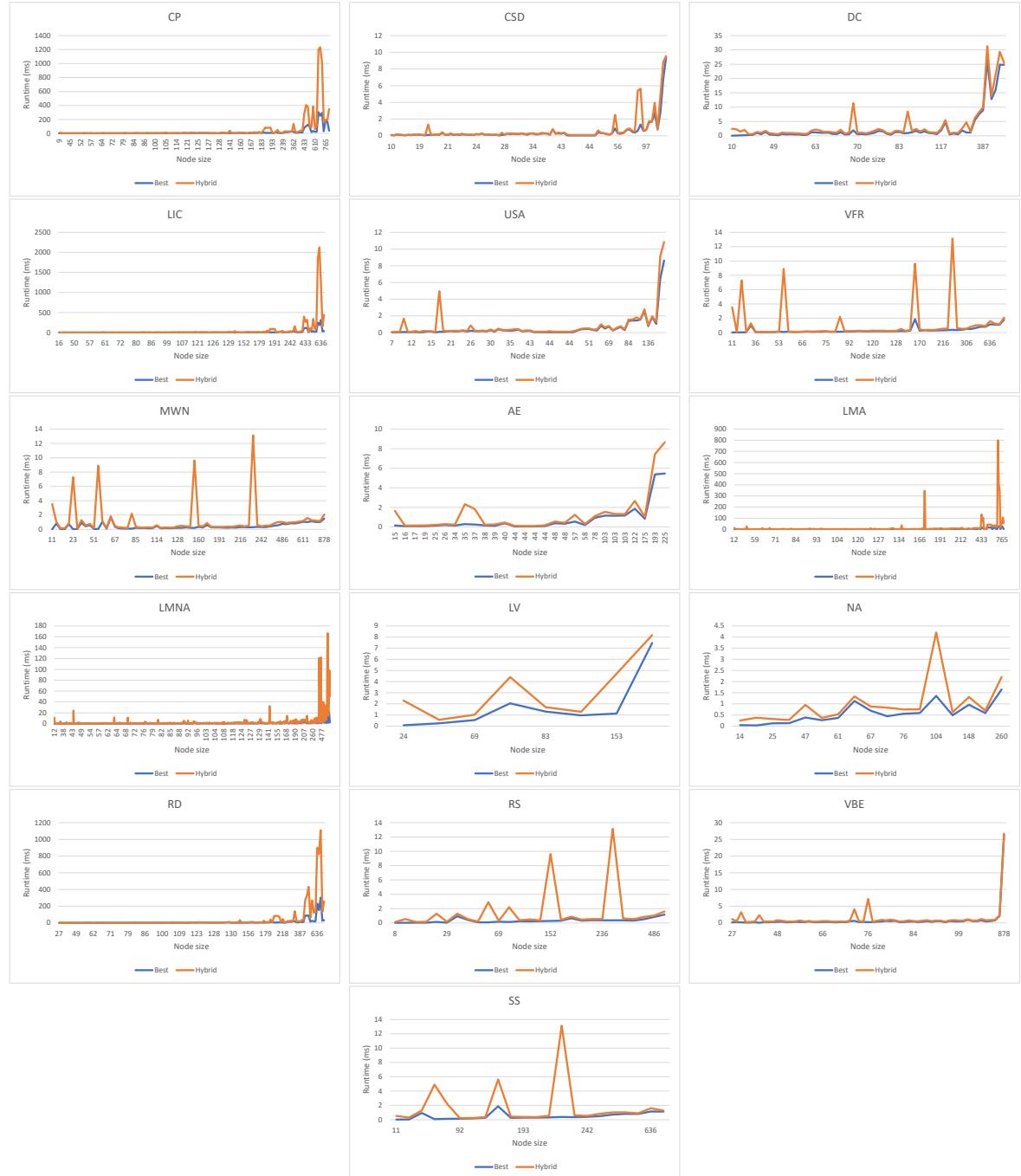


Figure 9.2: Hybrid approach's performance against best approaches for mis-predicted graphs.

CHAPTER 10. Analysis on the Decision Tree Distribution

Decision tree is the key component in our hybrid approach. Given an analysis traversal and an input graph, a path along the check points in the tree will be used to determine the traversal strategy at the corresponding leaf node. There are such 11 paths leading to 11 leaf nodes as shown in Figure 4.3. In this experiment, we want to study the contribution of each path in determining strategies for CFGs from the two datasets. Two tables in Figure 10.1 show the result for 21 analyses. Background colors indicate the ranges of values: 0%, (0%, 1%), [1%, 10%) and [10%, 100%]. The result shows a trend which is consistent between two datasets.

Path P11 is visited most often because path P11 leads to ANY traversal strategy. This path is taken when we encounter data flow insensitive traversal and since atleast once such traversal is present in every analyses and all graphs irrespective of the cyclicity, take this path, we see about 36% of the graphs takes Path P11.

Path P1 is also visited about 35% of the time. This path P1 is taken by sequential graphs and forward data flow sensitive traversal. Since we have about 70% sequential graphs in our dataset and about 14 forward data flow sensitive traversal, Path P1 is taken 35% of the time.

P2 was visited 10% of the time even though there are only 4 backward analyses. This is because there are about 70% sequential CFGs and that boosts up the number to 10%. P3 is visited much often than P4, since P3 is for forward analyses and P4 is for backward analyses.

Paths P9 and P10 are less frequently used since only 10% of the CFGs has loops but since there are about 16 loop sensitive traversal, they were visited 5% and 2% of the times respectively.

Paths P5, P6, P7 and P8 were taken less often than the others because they are only used for CFGs with loops which are only 10% of the CFGs in the datatsets. In addition, these paths are taken when the traversal is data-flow sensitive and loop insensitive. Only two of our analyses contains such traversal. It is also worth to note that, from Figure 4.3, those four paths

	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10	P11
CP	35%	0%	10%	0%	0%	0%	0%	5%	0%	50%	
CSD	35%	0%	10%	0%	0%	0%	0%	5%	0%	50%	
DC	0%	35%	0%	10%	0%	0%	0%	0%	5%	50%	
LIC	35%	0%	10%	0%	0%	0%	0%	5%	0%	50%	
USA	35%	0%	10%	0%	0%	0%	0%	5%	0%	50%	
VFR	35%	0%	10%	0%	0%	0%	0%	5%	0%	50%	
MWN	35%	0%	10%	0%	0%	0%	0%	5%	0%	50%	
AE	69%	0%	21%	0%	0%	0%	0%	10%	0%	0%	
DOM	69%	0%	21%	0%	7%	0%	3%	0%	0%	0%	
LMA	69%	0%	21%	0%	0%	0%	0%	10%	0%	0%	
LMNA	69%	0%	21%	0%	0%	0%	0%	10%	0%	0%	
LV	0%	69%	0%	21%	0%	0%	0%	0%	10%	0%	
NA	69%	0%	21%	0%	0%	0%	0%	10%	0%	0%	
PDOM	0%	69%	0%	21%	0%	7%	0%	3%	0%	0%	
RD	69%	0%	21%	0%	0%	0%	0%	10%	0%	0%	
RS	69%	0%	21%	0%	0%	0%	0%	10%	0%	0%	
VBE	0%	69%	0%	21%	0%	0%	0%	0%	10%	0%	
SS	69%	0%	21%	0%	0%	0%	0%	10%	0%	0%	
UDV	0%	0%	0%	0%	0%	0%	0%	0%	0%	100%	
UIR	0%	0%	0%	0%	0%	0%	0%	0%	0%	100%	
WNIL	0%	0%	0%	0%	0%	0%	0%	0%	0%	100%	
Overall	34.54%	9.87%	10.32%	2.95%	0.26%	0.26%	0.11%	0.11%	4.78%	1.10%	35.71%

Figure 10.1: Distribution of decisions over the paths of the decision tree.

(P5–P8) are the longest paths in the tree. The fact that these longest paths are rare (less than 1% for both DaCapo and GitHub datasets) shows that most analyses and graphs are classified by our technique using fewer dynamic checks.

	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10	P11
CP	32%	0%	13%	0%	0%	0%	0%	0%	5%	0%	50%
CSD	32%	0%	13%	0%	0%	0%	0%	0%	5%	0%	50%
DC	0%	32%	0%	13%	0%	0%	0%	0%	0%	5%	50%
LIC	32%	0%	13%	0%	0%	0%	0%	0%	5%	0%	50%
USA	32%	0%	13%	0%	0%	0%	0%	0%	5%	0%	50%
VFR	32%	0%	13%	0%	0%	0%	0%	0%	5%	0%	50%
MWN	32%	0%	13%	0%	0%	0%	0%	0%	5%	0%	50%
AE	65%	0%	25%	0%	0%	0%	0%	0%	10%	0%	0%
DOM	65%	0%	25%	0%	7%	0%	2%	0%	0%	0%	0%
LMA	65%	0%	25%	0%	0%	0%	0%	0%	10%	0%	0%
LMNA	65%	0%	25%	0%	0%	0%	0%	0%	10%	0%	0%
LV	0%	65%	0%	25%	0%	0%	0%	0%	0%	10%	0%
NA	65%	0%	25%	0%	0%	0%	0%	0%	10%	0%	0%
PDOM	0%	65%	0%	25%	0%	7%	0%	2%	0%	0%	0%
RD	65%	0%	25%	0%	0%	0%	0%	0%	10%	0%	0%
RS	65%	0%	25%	0%	0%	0%	0%	0%	10%	0%	0%
VBE	0%	65%	0%	25%	0%	0%	0%	0%	0%	10%	0%
SS	65%	0%	25%	0%	0%	0%	0%	0%	10%	0%	0%
UDV	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	100%
UIR	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	100%
WNIL	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	100%
Overall	32.46%	9.27%	12.69%	3.62%	0.26%	0.26%	0.10%	0.10%	4.50%	1.04%	35.70%

Figure 10.2: Distribution of decisions over the paths of the decision tree for the DaCapo Dataset.

CHAPTER 11. Analysis on Traversal Optimization

We evaluated the importance of optimizing the chosen traversal strategy by comparing the hybrid approach with the non-optimized version. We computed the reduction rate on the running times for the 21 analyses. Figure 11.1 shows the reduction in execution time due to traversal strategy optimization. For analyses that involve at least one *data-flow sensitive* traversal, the optimization helps to reduce at least 60% of running time. This is because optimizations in such traversals reduce the number of iterations of traversals over the graphs by eliminating the redundant result re-computation traversal and the unnecessary fixpoint condition checking traversal. Thus in effect, we are ignoring two traversals per graph. Since the dataset contains about 90% of sequential graph, we only traverse the graph once due to optimization. If no optimization, we would have traversed the graph thrice and therefore we reduced two-third of the running time. This is the reason why we see about 60% reduction in time due to optimization.

For analyses involving only *data-flow insensitive* traversal, there is no reduction in execution time, as hybrid approach does not attempt to optimize. This is because, Hybrid mainly optimizes by removing the redundant result re-computation traversal and the unnecessary fixpoint condition

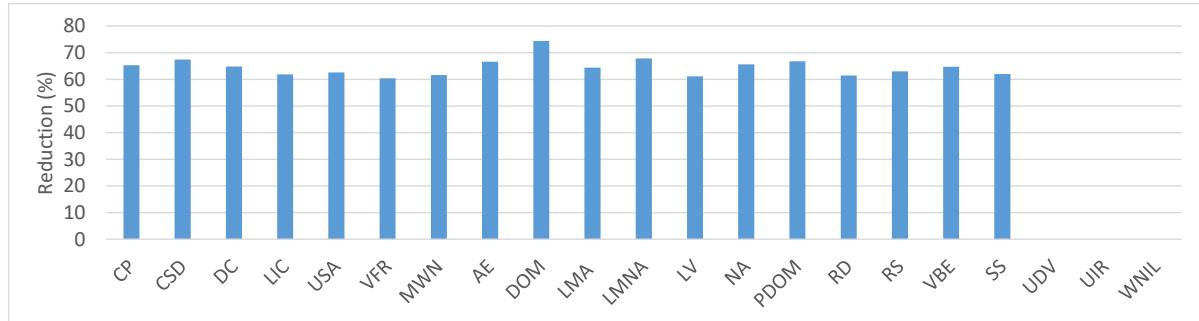


Figure 11.1: Reduction in execution time of the hybrid approach due to traversal optimization.

checking traversal. Both these traversals are not present in *data-flow insensitive* traversal, and hence Hybrid cannot optimize further.

CHAPTER 12. Case Studies

We implemented two case studies using our formalism for source code analysis and evaluated using hybrid and WRPO traversal. We are comparing against only WRPO since it is the next best performing traversal.

12.1 API Precondition Mining (APM).

This case study mines a large corpus of API usages to derive potential preconditions for API methods(Nguyen et al. (2014)). The key idea of this work is that API preconditions would be checked frequently in a corpus with a large number of API usages, while project-specific conditions would be less frequent. This case study analysis mined the preconditions for all methods of `java.lang.String`. Boa source code for this case study can be found in <http://boa.cs.iastate.edu/boa/index.php?q=boa/job/61006> .

Result analysis. Figure 12.2 lists the first and second most mined pre-condition for some of the methods in String java API. For all the methods, the most mined pre-condition is `var!=null`, which is expected as null condition check is done on the parameters before using these parameters in the API methods. Methods like `equals`, `startsWith`, `contains`, `matches` have `var!=null && arg!=null` as the second most mined precondition, as these API methods take a parameter and is applied on a variable. Hence null check is performed on both of these variables. `charAt` takes an index as parameter and a check on whether this index is lesser than the variable's

Case	Hybrid	WRPO	Reduce
APM	1527	1702	10%
AUM	883	963	8%
SVT	1417	1501	6%

Figure 12.1: Running time (minutes) of the case studies on GitHub data.

API Method	First	Second
equals	var!=null	var!=null && arg!=null
startsWith	var!=null	var!=null && arg!=null
contains	var!=null	var!=null && arg!=null
matches	var!=null	var!=null && arg!=null
split	var!=null	!isEmpty(var)
charAt	var!=null	arg var.length
substring	var!=null	arg!=-1

Figure 12.2: First and second most mined pre-condition for some of the methods in String java API.

#	API Pattern
1	iter=delegates.iterator(); LOOP (iter.hasNext()) { o=iter.next(); IF(o instanceof Object && !iter.hasNext()) { result=(Object) o IF(result != null) { return result; } } }
2	iter=delegates.iterator(); LOOP (iter.hasNext()) { o=iter.next(); IF(o instanceof Object) { store=(Object) o; } }
3	return list.size()
4	return iterator.hasNext()
5	CATCH(Exception) { log.error(e.getMessage(),e) }
6	CATCH(Exception) { log.error(e.getMessage(),e); ErrorDialog.open(e) }
7	CATCH(IOException) { log.log(Level.FINE,e.toString(),e) }
8	return this.map.size()
9	this.myArrayList.add(value:Object)
10	TRY { LOG.info("STARTING "+getName()) } CATCH(Exception) {}

Figure 12.3: Top 10 API Usage pattern for java.util API.

length is made (`arg < var.length`). Hence this is the second most mined precondition for `charAt` method.

12.2 API Usage Mining (AUM).

This case study analyzes API usage code and mines API usage patterns(Zhong et al. (2009)). The mined patterns help developers understand and write API usages more effectively with less errors. Our analysis mined usage patterns for `java.util` APIs. Boa source code for this case study can be found in <http://boa.cs.iastate.edu/boa/?q=boa/job/61007> .

Result analysis. Figure 12.3 lists the top 10 API Usage pattern for `java.util` API. Top two usage patterns are looping a collection using an iterator. Returning the size of the list and return an boolean indicating whether the collection contains more element are the next most mined usage patterns. The other most mined API usage patterns are related to `java.util.log` class (logging error, info, message).

Figure 12.1 shows that hybrid traversal helps reduce running times significantly by 80–175

minutes, which is from 6%–10% relatively.

CHAPTER 13. Threats to Validity

Our first threat to validity is our selection of source code analysis used in our evaluation. While there exists no source for a standard set of analysis, we relied mainly on text books and source code analysis tools to select analysis. We have selected basic control and data-flow analyses, and analyses to find bugs or code smells. We made sure to include analysis that covers all the properties of interest. For instance, our analysis set includes: both forward/backward analysis, data-flow sensitive and insensitive analysis, loop sensitive and insensitive analysis.

Our next threat to validity is our selection of ultra-large-scale datasets that provide graphs for running the analyses. The datasets do not contain a balanced distribution of different graph cyclicity (sequential, branch and loop). Both DaCapo and GitHub datasets contains majority of sequential graphs (65% and 69%, respectively) and only 10% are graphs with loops. The impact of this threat can be seen in our evaluation of the importance of paths and decisions in our decision tree. Paths and decisions along sequential graphs are taken more often. This threat is not easy to mitigate, as it is hard to find and difficult to expect a real-world code dataset to contain a balanced distribution of graphs of various types. Nonetheless, our evaluation shows that the selection and optimization of the best traversal strategy for these 35% of the graphs (graphs with branches and loops) plays an important role in improving the overall performance of the analysis over a large dataset of graphs.

CHAPTER 14. Related Work

To the best of our knowledge, our proposal to leverage information about the program analysis code, and the nature of the data on which analysis is applied to select appropriate traversal strategies has not been explored previously. Below we discuss works that are related to various aspects of our proposal.

14.1 Mixing static and dynamic information.

The general philosophy of mixing static and dynamic information has a long history(Ernst (2003)) in both the software engineering and the programming languages communities, with examples such as DSD-Crasher(Csallner et al. (2008)), Palus(Zhang et al. (2011)), segmented symbolic execution(Le (2013)), guided dynamic symbolic execution(Christakis et al. (2016)), gradual typing(Siek and Taha (2007)) , hybrid type checking (Flanagan (2006)), intensional polymorphism(Harper and Morrisett (1995), Crary et al. (1998)), etc. While our proposal also mixes static information about program analysis with dynamic information about the data, none of the previous works have proposed utilizing this information for selecting appropriate traversal strategies for realizing the program analysis.

14.2 Optimizing program analysis.

Atkinson Atkinson and Griswold (2001) presented techniques that reduce the time and space required to perform data-flow analysis of large programs. While their techniques proposed modifications to the underlying data-flow analyses that yield improvement in performance and also proposed reclamation of the data-flow sets during data-flow analysis that result in saving space, hybrid approach gives performance gain by analyzing the user written algorithm and the

input graph received.

Kildall Kildall (1973) presented an algorithm which, in conjunction with various optimizing functions, provides global program optimization. Optimizing functions have been described which provide constant propagation, common sub-expression elimination, and a degree of register optimization. While their approach provides unified approach to global program optimization, we concentrate on optimizing the process that does program optimization using the program's structure and the algorithms characteristics.

14.3 Ultra-large-scale source code mining.

In terms of ultra large scale processing, Boa(Dyer et al. (2013a)) is a language and infrastructure for analyzing ultra-large-scale software repositories. Boa provides a different kind of performance gain through its infrastructure and eases testing MSR-related hypotheses, it is not suitable for graph processing algorithms and does not leverage information from algorithms written in Boa. Upadhyaya and Rajan (2017) also tried to accelerate Ultra large scale score code mining. Their key idea is to analyze the interaction pattern between the mining task and the artifact to cluster artifacts such that running the mining task on one candidate artifact from each cluster is sufficient to produce results for other artifacts in the same cluster. Their artifact clustering criteria go beyond syntactic, semantic, and functional similarities to mining-task-specific similarity, where the interaction pattern between the mining task and the artifact is used for clustering. While their approach does task-specific clustering and extrapolates results, we try to analyze the analysis and come up with the best way to traverse the graph so that we can finish the analysis by visiting lesser number of nodes and lesser operations. Dyer et al. (2013b) developed domain-specific language abstractions for easily writing source code mining tasks on billions of AST nodes. While their language abstractions were for AST nodes, our language features were for traversing CFGs. Another important difference is that we can specify a user defined fixpoint for the traversals and the traversal will run till the fixpoint is reached. We can also provide the direction of traversal and the traversal can return outputs while the visitor construct in Dyer et al. (2013b) does not.

14.4 Graph traversal optimization.

There have been many works that targeted graph traversal optimization through various ways. Green-Marl(Hong et al. (2012)) provides performance benefits by using domain specific knowledge in applying optimizations. It uses high-level algorithmic description written in Green-Marl to exploiting the exposed data level parallelism. While green marl provides performance benefits by taking algorithm written into consideration, hybrid approach takes both algorithm and graph structure into account. And in ultra large scale dataset, containing millions of graphs with different structures, the gain that we can incur by taking graph structure into account is significant.

Pregel(Malewicz et al. (2010)) is a MapReduce like framework that aims to bring distributed processing to graph algorithms. While Pregel's performance gain is through parallelism and handles large graphs processing through vertex centric approach, our approach achieves performance gain by traversing the graph efficient suitable to the algorithm.

There have also been few libraries that support parallel or distributed graph analysis: Parallel BGL(Gregor and Lumsdaine (2005)) is a distributed version of BGL while SNAP(Bader and Madduri (2008)) is a stand-alone parallel graph analysis package.

CHAPTER 15. Conclusion and Future work

15.1 Conclusion

Improving the performance of source code analyses that runs on massive code bases is an ongoing challenge. One way to improve the performance of source code analysis expressed as traversals over graphs like CFGs, is by picking the optimal traversal strategy that defines the order of nodes visited. The selection of the best traversal strategy depends both on the properties of the analysis and the input graph on which the analysis is run. We proposed a hybrid technique for selecting and optimizing graph traversal strategies for source code analysis expressed as traversals over graphs. Our solution includes a system for expressing source code analysis as traversals, a set of static properties of the analysis and algorithms to compute them, a decision tree that checks static properties along with graph properties to select the most time-efficient traversal strategy. Our evaluation shows that the hybrid technique successfully selected the most time-efficient traversal strategy for 99.99%–100% of the time and using the selected traversal strategy and optimizing it, the running times of a representative collection of source code analysis in our evaluation were considerably reduced by 1%–28% (13 minutes to 72 minutes in absolute time) when compared against the best performing traversal strategy. The case studies show that hybrid traversal reduces 80–175 minutes in running times for three software engineering tasks. The overhead imposed by collecting additional information for our approach is less than 0.2% of the total running time for a large dataset and less than 0.01% for an ultra-large dataset.

15.2 Future Work

One possible future work is to understand how the complexity of analysis affects the decision making of traversal strategies for graphs with loops. Right now, We have mis-predictions for graphs with loops. We have identified that as graph size increases, mis-prediction increase. As part of future work, we would like to explore other factors that affect and expand the decision tree to predict correct strategies for large graphs with loops.

Another possible future work will be to build an agent and train it using the decision tree so that it understands and builds a model that gives much more accurate decision tree with lesser mis-predictions.

We could also expand our framework to inter-procedural analysis as we support only intra-procedural analysis right now. Expanding this will be a challenge as we will be dealing with multiple cfgs and the current algorithm for data flow sensitivity and loop sensitivity should be carefully revised. It also requires investigating if any new factors play a role in traversal strategy decision making for inter-procedural analysis.

Also there are analysis whose output changes with the traversal used. They are called traversal sensitive analysis and we need to come up with how to handle such analysis and investigate if there is any way we can recommend a traversal to the user.

Another potential future work is to expand the decision tree to general graphs as we deal with only CFGs now. It requires investigating the factors that affect the traversal strategy decision for general graph analysis and if there is any room for improvement.

We could also expand our framework such that after running the check once to select the optimal traversal, one would not have to rerun this unless the analysis technique under scrutiny has fundamentally changed. That is, one could store a cache of solutions which could be consulted prior to searching for the optimal traversal algorithm again.

Bibliography

- Aho, A. V., Lam, M. S., Sethi, R., and Ullman, J. D. (2006). *Compilers: Principles, Techniques, and Tools (2Nd Edition)*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA.
- Atkinson, D. C. and Griswold, W. G. (2001). Implementation techniques for efficient data-flow analysis of large programs. In *Proceedings of the IEEE International Conference on Software Maintenance (ICSM'01)*, ICSM '01, pages 52–, Washington, DC, USA. IEEE Computer Society.
- Ayewah, N., Pugh, W., Morgenthaler, J. D., Penix, J., and Zhou, Y. (2007). Evaluating static analysis defect warnings on production software. In *Proceedings of the 7th ACM SIGPLAN-SIGSOFT Workshop on Program Analysis for Software Tools and Engineering*, PASTE '07, pages 1–8, New York, NY, USA. ACM.
- Bader, D. A. and Madduri, K. (2008). SNAP, small-world network analysis and partitioning: An open-source parallel graph framework for the exploration of large-scale networks. In *2008 IEEE International Symposium on Parallel and Distributed Processing*, pages 1–12.
- Blackburn, S. M., Garner, R., Hoffmann, C., Khang, A. M., McKinley, K. S., Bentzur, R., Diwan, A., Feinberg, D., Frampton, D., Guyer, S. Z., Hirzel, M., Hosking, A., Jump, M., Lee, H., Moss, J. E. B., Phansalkar, A., Stefanović, D., VanDrunen, T., von Dincklage, D., and Wiedermann, B. (2006). The DaCapo benchmarks: Java benchmarking development and analysis. In *Proceedings of the 21st Annual ACM SIGPLAN Conference on Object-oriented Programming Systems, Languages, and Applications*, OOPSLA '06, pages 169–190, New York, NY, USA. ACM.

- Blume, W., Doallo, R., Eigenmann, R., Grout, J., Hoeflinger, J., and Lawrence, T. (1996). Parallel programming with polaris. *Computer*, 29(12):78–82.
- Bourdoncle, F. (1993). Efficient Chaotic Iteration Strategies With Widenings. In *Proceedings of the International Conference on Formal Methods in Programming and their Applications*, pages 128–141. Springer-Verlag.
- Choi, J.-D., Burke, M., and Carini, P. (1993). Efficient flow-sensitive interprocedural computation of pointer-induced aliases and side effects. In *Proceedings of the 20th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 232–245. ACM.
- Christakis, M., Müller, P., and Wüstholtz, V. (2016). Guiding dynamic symbolic execution toward unverified program executions. In *Proceedings of the 38th International Conference on Software Engineering*, ICSE ’16, pages 144–155, New York, NY, USA. ACM.
- Crary, K., Weirich, S., and Morrisett, G. (1998). Intensional polymorphism in type-erasure semantics. In *Proceedings of the International Conference on Functional Programming*.
- Csallner, C., Smaragdakis, Y., and Xie, T. (2008). Dsd-crasher : A hybrid analysis tool for bug finding. *ACM Trans. Softw. Eng. Methodol.*, 17(2):8:1–8:37.
- Dyer, R., Nguyen, H. A., Rajan, H., and Nguyen, T. N. (2013a). Boa: A language and infrastructure for analyzing ultra-large-scale software repositories. In *Proceedings of the 2013 International Conference on Software Engineering*, ICSE ’13, pages 422–431, Piscataway, NJ, USA. IEEE Press.
- Dyer, R., Rajan, H., and Nguyen, T. N. (2013b). Declarative visitors to ease fine-grained source code mining with full history on billions of ast nodes. In *ACM SIGPLAN Notices*, volume 49, pages 23–32. ACM.
- Engler, D., Chen, D. Y., Hallem, S., Chou, A., and Chelf, B. (2001). Bugs as deviant behavior: A general approach to inferring errors in systems code. In *Proceedings of the Eighteenth ACM Symposium on Operating Systems Principles*, SOSP ’01, pages 57–72, New York, NY, USA. ACM.

- Ernst, M. D. (2003). Static and dynamic analysis: Synergy and duality. In *WODA 2003: ICSE Workshop on Dynamic Analysis*, pages 24–27.
- Flanagan, C. (2006). Hybrid type checking. In *Conference Record of the 33rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL '06, pages 245–256, New York, NY, USA. ACM.
- Gregor, D. and Lumsdaine, A. (2005). The parallel BGL: A generic library for distributed graph computations. *Parallel Object-Oriented Scientific Computing (POOSC)*, 2:1–18.
- Harper, R. and Morrisett, G. (1995). Compiling polymorphism using intensional type analysis. In *Proceedings of the Symposium on Principles of Programming Languages*, POPL '95, pages 130–141, New York, NY, USA. ACM.
- Hong, S., Chafi, H., Sedlar, E., and Olukotun, K. (2012). Green-marl: A DSL for easy and efficient graph analysis. In *Proceedings of the Seventeenth International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS XVII, pages 349–362, New York, NY, USA. ACM.
- Jagannathan, S., Thiemann, P., Weeks, S., and Wright, A. (1998). Single and Loving It: Must-alias Analysis for Higher-order Languages. In *Proceedings of the 25th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL '98, pages 329–341, New York, NY, USA. ACM.
- Kildall, G. A. (1973). A unified approach to global program optimization. In *Proceedings of the 1st Annual ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages*, POPL '73, pages 194–206, New York, NY, USA. ACM.
- Le, W. (2013). Segmented symbolic analysis. In *Proceedings of the 2013 International Conference on Software Engineering*, ICSE '13, pages 212–221, Piscataway, NJ, USA. IEEE Press.
- Li, Z., Lu, S., and Myagmar, S. (2006). CP-Miner: Finding Copy-Paste and Related Bugs in Large-Scale Software Code. *IEEE Trans. Softw. Eng.*, 32(3):176–192.

- Livshits, B. and Zimmermann, T. (2005). Dynamine: finding common error patterns by mining software revision histories. In *ACM SIGSOFT Software Engineering Notes*, volume 30, pages 296–305. ACM.
- Malewicz, G., Austern, M. H., Bik, A. J., Dehnert, J. C., Horn, I., Leiser, N., and Czajkowski, G. (2010). Pregel: A system for large-scale graph processing. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of Data*, SIGMOD ’10, pages 135–146, New York, NY, USA. ACM.
- Nguyen, H. A., Dyer, R., Nguyen, T. N., and Rajan, H. (2014). Mining preconditions of apis in large-scale code corpus. In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*, pages 166–177. ACM.
- Nielson, F., Nielson, H. R., and Hankin, C. (2010). *Principles of Program Analysis*. Springer Publishing Company, Incorporated.
- Ramanathan, M. K., Grama, A., and Jagannathan, S. (2007). Path-sensitive inference of function precedence protocols. In *Proceedings of the 29th International Conference on Software Engineering*, ICSE ’07, pages 240–250, Washington, DC, USA. IEEE Computer Society.
- Siek, J. and Taha, W. (2007). Gradual typing for objects. In *Proceedings of the 21st European Conference on Object-Oriented Programming*, ECOOP ’07, pages 2–27, Berlin, Heidelberg. Springer-Verlag.
- Thummalapenta, S. and Xie, T. (2009). Alattin: Mining alternative patterns for detecting neglected conditions. In *Proceedings of the 2009 IEEE/ACM International Conference on Automated Software Engineering*, pages 283–294. IEEE Computer Society.
- Upadhyaya, G. and Rajan, H. (2017). On accelerating ultra-large-scale mining. In *Proceedings of the 39th International Conference on Software Engineering: New Ideas and Emerging Results Track*, pages 39–42. IEEE Press.

- Vallée-Rai, R., Co, P., Gagnon, E., Hendren, L., Lam, P., and Sundaresan, V. (1999). Soot : Java bytecode optimization framework. In *Proceedings of the 1999 conference of the Centre for Advanced Studies on Collaborative research*, page 13. IBM Press.
- Weimer, W. and Necula, G. C. (2005). Mining temporal specifications for error detection. In *Proceedings of the 11th International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, TACAS'05, pages 461–476, Berlin, Heidelberg. Springer-Verlag.
- Yang, J., Evans, D., Bhardwaj, D., Bhat, T., and Das, M. (2006). Perracotta: Mining temporal api rules from imperfect traces. In *Proceedings of the 28th International Conference on Software Engineering*, ICSE '06, pages 282–291, New York, NY, USA. ACM.
- Zhang, S., Saff, D., Bu, Y., and Ernst, M. D. (2011). Combined static and dynamic automated test generation. In *Proceedings of the 2011 International Symposium on Software Testing and Analysis*, ISSTA '11, pages 353–363, New York, NY, USA. ACM.
- Zhong, H., Xie, T., Zhang, L., Pei, J., and Mei, H. (2009). Mapo: Mining and recommending api usage patterns. *ECOOP 2009—Object-Oriented Programming*, pages 318–343.

CHAPTER 16. Appendix

Other source code analyses is present in <https://github.com/boalang/compiler/tree/program-analysis/test/program-analysis>.

Listing 16.1: API Precondition Mining Analysis

```

1 p: Project = input;
2 results: output collection of string;
3 verylarge: output sum of int;
4
5 apis: array of string;
6 apis = {"next","equals","get","substring","append","put","add","indexOf",
    "length","charAt","startsWith","getString","split","arraycopy","nextElement",
    "equalsIgnoreCase","getInt","setInt","setString","elementAt","getProperty",
    "max","endsWith","getClass","lastIndexOf","parseInt","size","replace",
    "trim","intValue","replaceAll","toString","exit","remove","min","nextToken",
    "toArray","valueOf","format","prepareStatement","contains","sleep","group",
    "getAbsolutePath","println","pow","forName","iterator","addElement",
    "getName","abs","toLowerCase","setLong","info","compareTo","compile",
    "containsKey","getLong","sqrt","delete","matcher","setProperty","invoke",
    "hashCode","setTimestamp","mkdirs","matches","print","getTimestamp",
    "fillRect","parseDouble","fine","set","write","keySet","parseLong","addAll",
    "toUpperCase","drawLine","getBytes","executeQuery","nextInt",
    "setAutoCommit","round","sin","setLength","getPath","cos","getX","getY",
    "isDirectory","listFiles","log","drawString","close","mkdir",
    "entrySet","createTempFile","newInstance","getMethod","getWidth",
    "drawImage","getHeight","pop"," setDate","warning","getDate",
    "booleanValue","getParentFile","read",
    "getRow","asList","last","getTime","getDouble","concat","push",
    "getBoolean","replaceFirst","subList","fill","values","encode","clear",
    "setColor","exists","toHexString","getObject","floor","firePropertyChange",
    "getSuperclass","getTimeZone","getBigDecimal","hasNext","toCharArray",
    "insert","getInstance","createNewFile","getBundle","parseFloat",
    "sort","finest","getNewValue","setDaemon","getLogger","doubleValue",
    "end","setObject","isWhitespace","start","parse","store","readLine",
    "ceil","getPoint","setAccessible","isNaN","list","removeElementAt",
    "setBoolean","drawRect","setStroke"," setTime","isAssignableFrom",
    "setBigDecimal","position","isDigit","exp","getAnnotation",
    "decode","parseBoolean","removeFirst","interrupt","executeUpdate",
    "setPaint","getChars","getComponent","getParameterTypes",
    "severe","finer","renameTo","getConstructor","copyInto",
    "getKey","peek","getFloat","setMaximumFractionDigits",
    "stringWidth","getMessage","lastModified","execute",
    "getenv","getCanonicalPath","insertElementAt","isFile",
    "setScale","compareToIgnoreCase","hasMoreTokens",
    "longValue","getValue","getColumnName","getParent",
    "setNull","deleteCharAt","previous","allocate",
    "getConnection","getDeclaredMethod","cast","getSimpleName"};
7 types: array of string;
8 types = {"Iterator","String","List","StringBuffer","Map",
    "IllegalArgumentException","StringBuilder","ArrayList","ResultSet",
    "Dimension","HashMap","System","Enumeration","PreparedStatement","Vector",
    "Math","Object","Integer","Insets","File","RuntimeException",
    "IllegalStateException","Hashtable","Set","StringTokenizer","Exception",
    "Connection","NullPointerException","ResourceBundle","Thread",
    "SimpleDateFormat","Matcher","Properties","Font","SortedMap","Color",
    "Rectangle","PrintWriter","Class","BigDecimal","IOException",
    "GridLayout","Logger","Error","Collection","Pattern","ListIterator",
    "HashSet","Long","Method","Boolean","GridBagConstraints",
    "Graphics","Double","Calendar","Writer","BufferedWriter",
    "SQLException"};
9 fulltypes: array of string;
10 fulltypes = {"java.util.Iterator","java.lang.String",
    "java.util.List","java.lang.StringBuffer",
    "java.util.Map","java.lang.IllegalArgumentException",
    "java.lang.StringBuilder","java.util.ArrayList",
    "java.sql.ResultSet","java.awt.Dimension",
    "java.util.HashMap","java.lang.System",
    "java.util.Enumeration","java.sql.PreparedStatement",
    "java.util.Vector","java.lang.Math",
    "java.lang.Object","java.lang.Integer",
    "java.awt.Insets","java.io.File",
    "java.lang.RuntimeException","java.lang.IllegalStateException",
    "java.util.Hashtable","java.util.Set",
    "java.util.StringTokenizer","java.lang.Exception",
    "java.sql.Connection","java.lang.NullPointerException",
    "java.util.ResourceBundle",

```

Listing 16.2: API Precondition Mining Analysis

```

11 "java.lang.Thread","java.text.SimpleDateFormat","java.util.regex.Matcher","java
   .util.Properties","java.awt.Font","java.util.SortedMap","java.awt.Color",
   "java.awt.Rectangle","java.io.PrintWriter","java.lang.Class","java.math.
   BigDecimal","java.io.IOException","java.awt.GridLayout","java.util.logging.
   Logger","java.lang.Error","java.util.Collection","java.util.regex.Pattern
   ","java.util.ListIterator","java.util.HashSet","java.lang.Long","java.lang.
   reflect.Method","java.lang.Boolean"};
12 consts: array of string;
13 consts = {"IllegalArgumentException","Dimension","StringBuffer","Integer",
   "Insets","File","RuntimeException","IllegalStateException","Exception",
   "ArrayList","StringTokenizer","NullPointerException","SimpleDateFormat",
   "StringBuilder","Font","Color","Rectangle","BigDecimal","IOException",
   "GridLayout","Error","String","Long","Boolean","GridBagConstraints","Double
   ","SQLException","Point","DecimalFormat","UnsupportedOperationException",
   "BorderLayout","Vector","InputStreamReader","HashMap","URL",
   "OutputStreamWriter","RandomAccessFile","Thread","FlowLayout","BasicStroke
   ","BufferedImage","PrintWriter","FileInputStream","FileOutputStream",
   "FileWriter","Date","HashSet","Float","ActionEvent","StringReader",
   "FileReader","ParseException","GradientPaint","PropertyChangeEvent",
   "ByteArrayOutputStream","NoSuchElementException","AssertionError","Hashtable
   ","SecurityException","Character","URI","ArithmeticException","Random",
   "Locale","BufferedReader","BufferedOutputStream","Button",
   "IndexOutOfBoundsException","NumberFormatException","BufferedWriter",
   "BufferedInputStream","ParsePosition","BigInteger","PropertyDescriptor",
   "ZipEntry","MenuItem","InternalError"};
14 path := "";
15 decls: stack of string;
16 locals: map[string] of string;
17 fields: map[string] of string;
18 parent_fields: map[string] of string;
19
20 getTypeOfLocals := function(method: Method): map[string] of string {
21 l: map[string] of string;
22 foreach (i:int; def(method.arguments[i])) {
23 arg := method.arguments[i].name;
24 argType := method.arguments[i].variable_type.name;
25 l[arg] = argType;
26 }
27 visit(method, visitor {
28 before expr: Expression -> {
29 if (expr.kind == ExpressionKind.VARDECL) {
30 foreach (j: int; def(expr.variable_decls[j])) {
31 var := expr.variable_decls[j].name;
32 varType := expr.variable_decls[j].variable_type.name;
33 l[var] = varType;
34 }
35 }
36 }
37 });
38 return l;
39 };
40
41 getTypeOfFields := function(decl: Declaration): map[string] of string {
42 l: map[string] of string;
43 foreach(i:int; def(decl.fields[i])) {
44 arg := decl.fields[i].name;
45 argType := decl.fields[i].variable_type.name;
46 l[arg] = argType;
47 }
48 return l;
49 };
50
51 hasJDKType := function(t: string): bool {

```

Listing 16.3: API Precondition Mining Analysis

```

52 if (t == "String" || t == "java.lang.String")
53 return true;
54 return false;
55 #     foreach (i: int; def(types[i])) {
56 #         if (t == types[i])
57 #             return true;
58 #     }
59
60 #     foreach (i: int; def(fulltypes[i])) {
61 #         if (t == fulltypes[i])
62 #             return true;
63 #     }
64 #
65 #     return false;
66 };
67
68 getVariableType := function(variable: string): string {
69 lKeys := keys(locals);
70 t := "";
71 foreach (i: int; def(lKeys[i])) {
72 if (variable == lKeys[i]) {
73 t = locals[lKeys[i]];
74 if (hasJDKType(t))
75 return t;
76 }
77 }
78 lKeys = keys(fields);
79 foreach (i: int; def(lKeys[i])) {
80 if (variable == lKeys[i]) {
81 t = fields[lKeys[i]];
82 if (hasJDKType(t))
83 return t;
84 }
85 }
86 lKeys = keys(parent_fields);
87 foreach (i: int; def(lKeys[i])) {
88 if (variable == lKeys[i]) {
89 t = parent_fields[lKeys[i]];
90 if (hasJDKType(t))
91 return t;
92 }
93 }
94 return "";
95 };
96
97 type T= {dom: set of string, dummy : int};
98 cfg: CFG;
99 cfgnode_ids:set of string;
100
101 allnode_ids := traversal(node: CFGNode): string {
102 add(cfgnode_ids, string(node.id));
103 return string(node.id);
104 };
105
106 cfg_dom := traversal(node: CFGNode): T {
107 cur_value : T;
108 if (node.id == 0) {
109 self_dom: set of string;
110 cur_value = {self_dom, 0};
111 }
112 else {

```

Listing 16.4: API Precondition Mining Analysis

```

112 else {
113   cur_value = {setClone(cfgnode_ids), 0};
114 }
115 if (def(getValue(node))) {
116   cur_val1 := getValue(node);
117   cur_value = clone(cur_val1);
118 }
119 preds := node.predecessors;
120 foreach (i: int; def(preds[i])) {
121   pred_value := getValue(preds[i]);
122   if (def(pred_value)) {
123     cur_value.dom = intersection(cur_value.dom, pred_value.dom);
124   }
125 }
126 gen_kill := getValue(node, allnode_ids);
127 add(cur_value.dom, gen_kill);
128 return cur_value;
129 };
130
131 fixp1 := fixp(curr, prev: T) : bool {
132   if (difference(curr.dom, prev.dom) == 0)
133     return true;
134   return false;
135 };
136
137 apiCallNodes: set of string;
138 predicateExprsAtNodes: map[string] of string;
139
140 hasAPICall := function(expression: Expression): bool {
141   ret := false;
142   visit(expression, visitor {
143     before expr: Expression -> {
144       m := expr.method;
145       if (expr.kind == ExpressionKind.METHODCALL) {
146         foreach (i: int; def(apis[i])) {
147           if (m == apis[i]) {
148             # check the type of the receiver
149             if (len(expr.expressions) > 0) {
150               recv := expr.expressions[0];
151               rv := "";
152               visit(recv, visitor {
153                 before e: Expression -> {
154                   if (e.kind == ExpressionKind.VARACCESS)
155                     rv = e.variable;
156                 }
157               });
158               rvType := getVariableType(rv);
159               if (rv != "" && rvType != "") {
160                 ret = true;
161                 break;
162               }
163             }
164           }
165         }
166       } else if (expr.kind == ExpressionKind.NEW) {
167         foreach (i: int; def(consts[i])) {
168           if (expr.new_type.name == consts[i]) {
169             ret = true;
170             break;
171           }
172         }
173       }

```

Listing 16.5: API Precondition Mining Analysis

```

173 }
174 }
175 });
176 return ret;
177 };
178
179 mineAPICallsAndPredicates := traversal(node: CFGNode): bool {
180 ret := false;
181 if (def(node.expr) && hasAPICall(node.expr)) {
182 add(apiCallNodes, string(node.id));
183 } else if (node.name == "IF") {
184 predicateExpr := format("%s", node.expr);
185 predicateExprsAtNodes[string(node.id)] = predicateExpr;
186 }
187 return ret;
188 };
189
190 apiPredicates: set of string;
191 dom_result := traversal(node: CFGNode): bool {
192 ret := false;
193 cur_value := getValue(node, cfg_dom);
194 nodeId := string(node.id);
195 predicates: set of string;
196
197 if (def(cur_value)) {
198 if (contains(apiCallNodes, nodeId)) {
199 foreach(dom: string=cur_value.dom) {
200 if (haskey(predicateExprsAtNodes, dom)) {
201 predicate := predicateExprsAtNodes[dom];
202 add(predicates, predicate);
203 }
204 }
205 }
206 }
207 if (len(predicates) > 0) {
208 if (node.expr.method=="charAt" ||
209 node.expr.method=="codePointAt" ||
210 node.expr.method=="codePointBefore" ||
211 node.expr.method=="codePointCount" ||
212 node.expr.method=="compareTo" ||
213 node.expr.method=="compareToIgnoreCase" ||
214 node.expr.method=="concat" ||
215 node.expr.method=="contains" ||
216 node.expr.method=="contentEquals" ||
217 node.expr.method=="copyValueOf" ||
218 node.expr.method=="endsWith" ||
219 node.expr.method=="equals" ||
220 node.expr.method=="equalsIgnoreCase" ||
221 node.expr.method=="format" ||
222 node.expr.method=="getBytes" ||
223 node.expr.method=="getChars" ||
224 node.expr.method=="hashCode" ||
225 node.expr.method=="indexOf" ||
226 node.expr.method=="intern" ||
227 node.expr.method=="isEmpty" ||
228 node.expr.method=="lastIndexOf" ||
229 node.expr.method=="length" ||
230 node.expr.method=="matches" ||
231 node.expr.method=="offsetByCodePoints" ||
232 node.expr.method=="regionMatches" ||
233 node.expr.method=="replace" ||
234 node.expr.method=="replaceAll" ||

```

Listing 16.6: API Precondition Mining Analysis

```

234 node.expr.method=="replaceAll" ||
235 node.expr.method=="replaceFirst" ||
236 node.expr.method=="split" ||
237 node.expr.method=="startsWith" ||
238 node.expr.method=="subSequence" ||
239 node.expr.method=="substring" ||
240 node.expr.method=="toCharArray" ||
241 node.expr.method=="toLowerCase" ||
242 node.expr.method=="toString" ||
243 node.expr.method=="toUpperCase" ||
244 node.expr.method=="trim" ||
245 node.expr.method=="valueOf") {
246 apiCall := format("%s(%d)", node.expr.method, len(node.expr.method_args));
247 preds := format("%s", predicates);
248 add(apiPredicates, format("%s --- %s", apiCall, preds));
249 }
250 }
251 return ret;
252 };
253
254 q_all := visitor {
255 before node: CodeRepository -> {
256 snapshot := getsnapshot(node, "SOURCE_JAVA_JLS");
257 foreach (i: int; def(snapshot[i]))
258 visit(snapshot[i]);
259 stop;
260 }
261 before chfl: ChangedFile -> path = chfl.name;
262 after chfl: ChangedFile -> path = "";
263 before decl: Declaration -> {
264 push(decls, decl.name);
265 if (len(keys(fields)) > 0) {
266 mKeys := keys(fields);
267 foreach (i: int; def(mKeys[i])) {
268 k := mKeys[i];
269 v := fields[k];
270 parent_fields[k] = v;
271 }
272 }
273 clear(fields);
274 fields = getTypeOfFields(decl);
275 }
276 after decl: Declaration -> {
277 pop(decls);
278 clear(fields);
279 if (len(decl.nested_declarations) > 0)
280 clear(parent_fields);
281 }
282
283 before method: Method -> {
284 locals = getTypeOfLocals(method);
285
286 clear(cfgnode_ids); clear(apiCallNodes); clear(predicateExprsAtNodes);
287 clear(mineAPICallsAndPredicates); clear(allnode_ids); clear(cfg_dom); clear(
287     dom_result);
288 cfg = getcfg(method);
289 if (len(cfg.nodes) > 10000) {
290 verylarge << 1;
291 stop;
292 }
293 if(len(cfg.nodes) != 0) {
294 traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID,
294     mineAPICallsAndPredicates);

```

Listing 16.7: API Precondition Mining Analysis

```
295 traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID, allnode_ids);
296 traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID, cfg_dom, fixp1)
297 ;
298 }
299
300 }
301 after method: Method -> {
302 clear(locals);
303 }
304 };
305
306 visit(p, q_all);
307
308 foreach(apiPred: string = apiPredicates) {
309 results << apiPred;
310 }
```

Listing 16.8: Dominator analysis

```

1 m: output collection[string][int] of string;
2 mt: output collection[int][string] of string;
3 p: Project = input;
4
5 type T= {dom: set of string, dummy : int};
6
7 cfg: CFG;
8 cfgnode_ids:set of string;
9
10 allnode_ids := traversal(node: CFGNode) : string{
11     add(cfgnode_ids, string(node.id));
12     return string(node.id);
13 };
14
15 cfg_dom := traversal(node: CFGNode): T {
16     cur_value : T;
17     if(node.id==0) {
18         self_dom:set of string;
19         cur_value = {self_dom, 0};
20     }
21     else
22         cur_value = {setClone(cfgnode_ids), 0};
23     if(def(getValue(node))) {
24         cur_val1 := getValue(node);
25         cur_value = clone(cur_val1);
26     }
27     preds:=node.predecessors;
28     foreach(i:int;def(preds[i])) {
29         pred_value := getValue(preds[i]);
30         if(def(pred_value)) {
31             cur_value.dom = intersection(cur_value.dom,pred_value.dom);
32         }
33     }
34     gen_kill := getValue(node, allnode_ids);
35     add(cur_value.dom, gen_kill);
36     return cur_value;
37 };
38
39 fixp1 := fixp(curr, prev: T) : bool {
40     if (difference(curr.dom, prev.dom) == 0)
41         return true;
42     return false;
43 };
44
45 q_all := visitor {
46     before node: CodeRepository -> {
47         snapshot := getsnapshot(node, "SOURCE_JAVA_JLS");
48         foreach (i: int; def(snapshot[i]))
49             visit(snapshot[i]);
50         stop;
51     }
52     before node: Method -> {
53         clear(allnode_ids);clear(cfgnode_ids);clear(cfg_dom);
54
55         cfg = getcfg(node);
56         traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID,
57                  allnode_ids);
58         traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID, cfg_dom
59                  , fixp1);
60     }
61 };
62 visit(p, q_all);

```

Listing 16.9: Post Dominator analysis

```

1 m: output collection[string][int] of string;
2 mt: output collection[int] of string;
3 p: Project = input;
4
5 type T= {dom: set of string, dummy : int};
6
7 cfg: CFG;
8 cfgnode_ids:set of string;
9 exitId : int;
10
11 allnode_ids := traversal(node: CFGNode) : string {
12     add(cfgnode_ids, string(node.id));
13     return string(node.id);
14 };
15
16 cfg_dom := traversal(node: CFGNode): T {
17     cur_value : T;
18     if(node.id==exitId) {
19         self_dom:set of string;
20         cur_value = {self_dom, 0};
21     }
22     else
23         cur_value = {cfgnode_ids, 0};
24     if(def(getValue(node))) {
25         cur_val1 := getValue(node);
26         cur_value = clone(cur_val1);
27     }
28     start_dom:set of string;
29     add(start_dom, string(node.id));
30     preds:=node.successors;
31     foreach(i:int;def(preds[i])) {
32         pred_value := getValue(preds[i]);
33         if(def(pred_value)) {
34             cur_value.dom = intersection(cur_value.dom,pred_value.dom);
35         }
36     }
37     gen_kill := getValue(node, allnode_ids);
38     if(def(gen_kill)) {
39         add(cur_value.dom, gen_kill);
40     }
41     return cur_value;
42 };
43
44 fixp1 := fixp(curr, prev: T) : bool {
45     if (difference(curr.dom, prev.dom) == 0)
46         return true;
47     return false;
48 };
49
50 q_all := visitor {
51     before node: Method -> {
52         clear(allnode_ids);clear(cfgnode_ids);clear(cfg_dom);
53
54         cfg = getcfg(node);
55         exitId = len(cfg.nodes) - 1;
56         traverse(cfg, TraversalDirection.BACKWARD, TraversalKind.HYBRID,
57                 allnode_ids);
58         traverse(cfg, TraversalDirection.BACKWARD, TraversalKind.HYBRID,
59                 cfg_dom, fixp1);
60     };
61 };

```

Listing 16.10: Live variable analysis

```

1 m: output collection[string][string][int] of int;
2 p: Project = input;
3
4 type T_gen_kill= {gen: set of string, kill: string, dummy : int};
5 type T_inout= {in: set of string, out: set of string};
6
7 cfg: CFG;
8 cur_cfg_node: CFGNode;
9
10 genset : set of string;
11 killset : set of string;
12
13 init := traversal(node: CFGNode): T_gen_kill {
14     cur_value : T_gen_kill;
15     cur_value = {node.useVariables, node.defVariables, 0};
16     return cur_value;
17 };
18
19 live := traversal(node: CFGNode): T_inout {
20     succs := node.successors;
21     in_set : set of string;
22     out_set : set of string;
23     cur_val : T_inout = {in_set, out_set};
24     if(def(getValue(node))) {
25         cur_val1 := getValue(node);
26         cur_val = clone(cur_val1);
27     }
28     foreach(succ_node:CFGNode=succs) {
29         succ := getValue(succ_node);
30         if(def(succ)) {
31             cur_val.out = union(cur_val.out,succ.in);
32         }
33     }
34     gen_kill := getValue(node, init);
35     if(def(gen_kill)) {
36         remove(cur_val.out, gen_kill.kill);
37         cur_val.in = union(gen_kill.gen, cur_val.out);
38     }
39     return cur_val;
40 };
41
42 fixp1 := fixp(curr, prev: T_inout) : bool {
43     if (difference(curr.in, prev.in) == 0)
44         return true;
45     return false;
46 };
47
48 q_all := visitor {
49     before node: Method -> {
50         clear(init);clear(live);
51
52         cfg = getcfg(node);
53         traverse(cfg, TraversalDirection.BACKWARD, TraversalKind.HYBRID, init);
54         traverse(cfg, TraversalDirection.BACKWARD, TraversalKind.HYBRID, live,
55                 fixp1);
56     }
57 };
58 visit(p, q_all);

```

Listing 16.11: Reaching definition analysis

```

1 p: Project = input;
2 m: output collection[string][int] of string;
3
4 type T= {in: set of string, out: set of string};
5 type T1= {gen: string, kill: string};
6
7 cfg: CFG;
8 cur_node_id: int;
9 cur_cfg_node: CFGNode;
10
11 gensex: set of string;
12 vardef: string;
13 killsex: set of string;
14
15 cfg_def := traversal(node: CFGNode) : T1 {
16     cur_val : T1 = {"", ""};
17     if(node.defVariables!="") {
18         cur_val.gen = node.defVariables+"@"+string(node.id);
19         cur_val.kill = node.defVariables;
20     }
21     return cur_val;
22 };
23
24 cfg_reach_def := traversal(n: CFGNode): T {
25     preds := n.predecessors;
26     in_set : set of string;
27     out_set : set of string;
28     cur_val : T = {in_set, out_set};
29     if(def(getValue(n))) {
30         cur_val1 := getValue(n);
31         cur_val = clone(cur_val1);
32     }
33     foreach(pred_node:CFGNode=preds) {
34         pred := getValue(pred_node);
35         if(def(pred))
36             cur_val.in = union(cur_val.in, pred.out);
37     }
38     cur_val.out = setClone(cur_val.in);
39     genkill := getValue(n, cfg_def);
40     if(genkill.kill!="") {
41         tmp_out:=setClone(cur_val.out);
42         foreach(tmp:string=tmp_out) {
43             tmp1:=stringClone(tmp);
44             str_array:=splitall(tmp1,"@");
45             if(str_array[0] == genkill.kill) {
46                 remove(cur_val.out, tmp1);
47             }
48         }
49         add(cur_val.out, genkill.gen);
50     }
51     return cur_val;
52 };
53
54 fixp1 := fixp(curr, prev: T) : bool {
55     if (difference(curr.out, prev.out) == 0)
56         return true;
57     return false;
58 };

```

Listing 16.12: Reaching definition analysis

```
59
60  reach_def := visitor {
61      before node: Method -> {
62          clear(cfg_def); clear(cfg_reach_def);
63          cfg = getcfg(node);
64          if(cfg.isValid==0) {
65              traverse(cfg, TraversalDirection.FORWARD, TraversalKind.ITERATIVE,
66                      cfg_def);
67              traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID,
68                      cfg_reach_def, fixp1);
69          }
70      };
71
72  visit(p, reach_def);
```

Listing 16.13: Available expression analysis

```

1 m: output collection[string][int] of string;
2 p: Project = input;
3
4 type T1 = {in: set of set of string, out: set of set of string};
5 allExprset: set of set of string;
6 cfg: CFG;
7
8 str: set of string;
9 str1: set of set of string;
10
11 node_expr := visitor {
12     before node: Expression -> {
13         switch (node.kind) {
14             case ExpressionKind.LT, ExpressionKind.OP_ADD, ExpressionKind.OP_SUB
15                 , ExpressionKind.OP_INC, ExpressionKind.OP_MULT, ExpressionKind
16                     .OP_DIV, ExpressionKind.OP_MOD, ExpressionKind.OP_DEC,
17                     ExpressionKind.GT, ExpressionKind.EQ, ExpressionKind.NEQ,
18                     ExpressionKind.LTEQ, ExpressionKind.GTEQ, ExpressionKind.
19                     LOGICAL_NOT, ExpressionKind.LOGICAL_AND, ExpressionKind.
20                     LOGICAL_OR, ExpressionKind.BIT_AND, ExpressionKind.BIT_OR,
21                     ExpressionKind.BIT_NOT, ExpressionKind.BIT_XOR, ExpressionKind.
22                     BIT_LSHIFT, ExpressionKind.BIT_RSHIFT, ExpressionKind.
23                     BIT_UNSIGNEDRSHIFT:
24                 add(str, string(node.kind));
25                 foreach(j:int;def(node.expressions[j])) {
26                     visit(node.expressions[j]);
27                 }
28                 break;
29             case ExpressionKind.ASSIGN:
30                 foreach(j:int;def(node.expressions[j])) {
31                     if(j!=0) {
32                         visit(node.expressions[j]);
33                     }
34                 }
35                 break;
36             case ExpressionKind.VARACCESS:
37                 add(str, node.variable);
38                 break;
39             case ExpressionKind.VARDECL:
40                 visit(node.variable_decls[0].initializer);
41                 break;
42             case ExpressionKind.LITERAL:
43                 add(str, node.literal);
44                 break;
45             case ExpressionKind.METHODCALL:
46                 clear(str);
47                 stop;
48             default:break;
49         }
50         stop;
51     }
52 };
53
54 allExprTraversal := traversal(node: CFGNode): set of string {
55     init_set : set of string;
56     str = init_set;
57     if(def(node.expr)) {
58         visit(node.expr, node_expr);
59         if(len(str)!=0) {
60             add(allExprset, setClone(str));
61         }
62     }
63 }
64 return setClone(str);

```

Listing 16.14: Available expression analysis

```

59     str1 = init_set1;
60     if(def(node.defVariables)) {
61         foreach(aa:set of string=allExprset) {
62             if(contains(aa, node.defVariables)) {
63                 add(str1, aa);
64             }
65         }
66     }
67     return setClone(str1);
68 };
69
70 avail_expr := traversal(node: CFGNode): T1 {
71     in_set : set of set of string;
72     out_set : set of set of string;
73     cur_value : T1;
74     if(node.id==0) {
75         cur_value = {in_set, out_set};
76     }
77     else
78         cur_value = {setClone(allExprset), out_set};
79     if(def(getValue(node))) {
80         cur_val1 := getValue(node);
81         cur_value = clone(cur_val1);
82     }
83     preds := node.predecessors;
84     foreach(pred_node:CFGNode=preds) {
85         pred := getValue(pred_node);
86         if(def(pred))
87             cur_value.in = intersection1(cur_value.in, pred.out);
88     }
89     genkill := getValue(node, allExprTraversal);
90     killset := getValue(node, killTraversal);
91     cur_value.out = setClone(cur_value.in);
92     removeAll(cur_value.out, killset);
93     if(len(genkill)!=0)
94         add(cur_value.out, genkill);
95     return cur_value;
96 };
97
98
99 fixp1 := fixp(curr, prev: T1) : bool {
100     if (difference1(curr.out, prev.out) == 0)
101         return true;
102     return false;
103 };
104
105 def := visitor {
106     before node: Method -> {
107         clear(allExprTraversal); clear(avail_expr); clear(killTraversal); clear(
108             allExprset);
109         cfg = getcfg(node);
110         traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID,
111             allExprTraversal);
112         traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID,
113             killTraversal);
114         traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID,
115             avail_expr, fixp1);
116     }
117 };
118
119 visit(p, def);

```

Listing 16.15: Used defined variable

```

1 m: output collection[string][string][int] of int;
2 p: Project = input;
3
4 type T_gen_kill= {gen: set of string, kill: string};
5
6 cfg: CFG;
7 cur_cfg_node: CFGNode;
8
9 genset : set of string;
10 killset : string;
11
12 node_def := visitor {
13     before expr:Expression -> {
14         switch(expr.kind) {
15             case ExpressionKind.VARDECL:
16                 var_decls := cur_cfg_node.expr.variable_decls;
17                 if(len(var_decls)!=0) {
18                     killset = var_decls[0].name;
19                 }
20                 break;
21             case ExpressionKind.ASSIGN:
22                 exprs := cur_cfg_node.expr.expressions;
23                 if(len(exprs)!=0) {
24                     killset = exprs[0].variable;
25                 }
26                 break;
27             default:
28                 break;
29         }
30     }
31 };
32
33 node_use := visitor {
34     before expr:Expression -> {
35         if(def(expr.variable)) {
36             add(genset, expr.variable);
37         }
38     }
39 };
40
41 init := traversal(node: CFGNode): T_gen_kill {
42     cur_value : T_gen_kill;
43     killset = "";
44     init_set1 : set of string;
45     genset = init_set1;
46     cur_cfg_node = node;
47     if(def(node.expr)) {
48         visit(node.expr, node_use);
49         visit(node.expr, node_def);
50     }
51     cur_value = {genset, killset};
52     return cur_value;
53 };
54
55 q_all := visitor {
56     before node: Method -> {
57         clear(init);
58         cfg = getcfg(node);
59         traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID, init);
60     }
61 };
62
63 visit(p, q_all);

```

Listing 16.16: Nullness Analysis

Listing 16.17: Nullness Analysis

```

59                     remove(cur_val.out, node.defVariables);
60                     break;
61                 }
62             }
63         }
64         if(flag == false) {
65             add(cur_val.out, node.defVariables);
66         }
67     }
68 }
69 }
70 if(node.expr.kind == ExpressionKind.ASSIGN) {
71     if(contains(local, node.defVariables)) {
72         if(node.expr.expressions[1].kind == ExpressionKind.LITERAL) {
73             if(node.expr.expressions[1].literal == "null") {
74                 remove(cur_val.out, node.defVariables);
75             }
76             else {
77                 add(cur_val.out, node.defVariables);
78             }
79         }
80     }
81     else {
82         flag1 := false;
83         foreach(use:string=node.useVariables) {
84             if(contains(local, use)) {
85                 if(!contains(cur_val.out, use)) {
86                     flag1 = true;
87                     remove(cur_val.out, node.defVariables);
88                     break;
89                 }
90             }
91             if(flag1 == false) {
92                 add(cur_val.out, node.defVariables);
93             }
94         }
95     }
96 }
97 }
98 return cur_val;
99 };
100 fixp1 := fixp(curr, prev: T) : bool {
101     if (difference(curr.out, prev.out) == 0)
102         return true;
103     return false;
104 };
105
106
107
108 reach_def := visitor {
109     before node: Method -> {
110         clear(allVarTraversal); clear(nullness); clear(local);
111
112         cfg = getcfg(node);
113         traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID,
114                  allVarTraversal);
115         traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID,
116                  nullness, fixp1);
117     };
118
119 visit(p, reach_def);

```

Listing 16.18: Local may alias

```

1 m: output collection[string][int] of string;
2 p: Project = input;
3 type T= {gen : string, kill : string};
4 type T1= {in : set of set of string, out : set of set of string};
5 cfg : CFG;
6
7 initPhase := traversal(node: CFGNode) : T {
8     cur_val : T = {"", ""};
9     if(node.defVariables!="")
10        cur_val.kill = node.defVariables;
11    if(def(node.expr)) {
12        if(node.expr.kind == ExpressionKind.VARDECL || node.expr.kind ==
13            ExpressionKind.ASSIGN) {
14            if(def(node.rhs)) {
15                if(node.rhs.kind == ExpressionKind.VARACCESS) {
16                    cur_val.gen = node.rhs.variable;
17                }
18            }
19        }
20    return cur_val;
21 };
22
23 analysisPhase := traversal(node: CFGNode) : T1 {
24     inset : set of set of string;
25     outset : set of set of string;
26     cur_val : T1 = {inset, outset};
27     if(def(getValue(node))) {
28         cur_val = getValue(node);
29     }
30     preds := node.predecessors;
31     foreach(pred_node:CFGNode=preds) {
32         pred := getValue(pred_node);
33         if(def(pred))
34             cur_val.in = union1(cur_val.in, pred.out);
35     }
36     cur_val.out = setClone(cur_val.in);
37     genkill := getValue(node, initPhase);
38     if(genkill.kill != "") {
39         foreach(tmp:set of string=cur_val.in) {
40             if(contains(tmp, genkill.kill)) {
41                 tmp2:=setClone(tmp);
42                 remove(cur_val.out, tmp2);
43                 remove(tmp2, genkill.kill);
44                 if(len(tmp2)!=0) {
45                     add(cur_val.out, tmp2);
46                 }
47             }
48         }
49         if(genkill.gen != "") {
50             flag := false;
51             tmpSet := setClone(cur_val.out);
52             foreach(tmp:set of string=tmpSet) {
53                 if(contains(tmp, genkill.gen)) {
54                     tmp3:=setClone(tmp);
55                     flag = true;
56                     remove(cur_val.out, tmp3);
57                     add(tmp3, genkill.kill);
58                     add(cur_val.out, tmp3);
59             }
}

```

Listing 16.19: Local may alias

```

59         }
60     }
61     if(flag == false) {
62         tmp1 : set of string;
63         add(tmp1, genkill.gen);add(tmp1, genkill.kill);
64         add(cur_val.out, setClone(tmp1));
65     }
66 }
67 else {
68     tmp4 : set of string;
69     add(tmp4, genkill.kill);
70     add(cur_val.out, setClone(tmp4));
71 }
72 }
73 return cur_val;
74 };
75
76 fixp1 := fixp(curr, prev: T1) : bool {
77     if (difference1(curr.out, prev.out) == 0)
78         return true;
79     return false;
80 };
81
82 q_all := visitor {
83     before node: Method -> {
84         clear(initPhase);clear(analysisPhase);
85
86         cfg = getcfg(node);
87         traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID,
88                 initPhase);
89         traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID,
90                 analysisPhase, fixp1);
91     }
92 };
93 visit(p, q_all);

```

Listing 16.20: Resource status

```

1 m: output collection[string][int] of string;
2 p: Project = input;
3
4 type T= {in: set of string, out : set of string};
5 type T1= {gen: string, kill : string};
6
7 cfg: CFG;
8 cfgnode_ids:set of string;
9
10 genStr : string;
11 killStr : string;
12
13 evisitor := visitor {
14     before node:Expression -> {
15         switch(node.kind) {
16             case ExpressionKind.METHODCALL:
17                 if(node.method == "write" || node.method == "read" || node.
18                     method == "open") {
19                     if(len(node.expressions)>0) {
20                         genStr = node.expressions[0].variable;
21                     }
22                     else if(node.method == "close") {
23                         if(len(node.expressions)>0) {
24                             killStr = node.expressions[0].variable;
25                         }
26                     }
27                     break;
28                 default: break;
29             }
30         }
31     };
32
33 track := traversal(node: CFGNode) : T1 {
34     cur_val : T1 = {"", ""};
35     genStr = "";
36     killStr = "";
37     if (def(node.expr)) {
38         visit(node.expr, evisitor);
39         cur_val = {genStr, killStr};
40     }
41     return cur_val;
42 };
43
44 cfg_def := traversal(node: CFGNode) : T {
45     inset : set of string;
46     outset : set of string;
47     cur_val : T = {inset, outset};
48     cur_val = getValue(node);
49     cur_val = {inset, outset};
50     preds := node.predecessors;
51     foreach(pred_node:CFGNode=preds) {
52         pred := getValue(pred_node);
53         if(def(pred))
54             cur_val.in = union(cur_val.in, pred.out);
55     }
56     track_val := getValue(node, track);
57     addAll(cur_val.out, cur_val.in);
58     if(track_val.gen != "")
59         add(cur_val.out, track_val.gen);

```

Listing 16.21: Resource status

```
59         add(cur_val.out, track_val.gen);
60     if(track_val.kill != "")
61         remove(cur_val.out, track_val.kill);
62     return cur_val;
63 };
64
65 fixp1 := fixp(curr, prev: T) : bool {
66     if (difference(curr.out, prev.out) == 0)
67         return true;
68     return false;
69 };
70
71 q_all := visitor {
72     before node: Method -> {
73         clear(track); clear(cfg_def);
74
75         cfg = getcfg(node);
76         traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID, track);
77         traverse(cfg, TraversalDirection.FORWARD, TraversalKind.HYBRID, cfg_def
78                 , fixp1);
79     }
80 };
81 visit(p, q_all);
```