

Call Graph Correction Using Control Flow Constraints

by

Kevin Resnick

Thesis

Presented to the Faculty of
The University of Texas at Austin

in Partial Fulfillment

of the Requirements

for the Degree of

Bachelor of Science

The University of Texas at Austin

May 2006

Acknowledgments

This thesis represents joint work with Byeongcheol Lee, Mike Bond, and Prof. Kathryn McKinley. Byeongcheol’s contributions are the design and implementation of the call graph correction algorithms and the FDOM (frequency dominator) computation algorithms for general control flow graph. I would like to thank my supervisor Kathryn McKinley for advising this research. This accomplishment would be not be possible without the valuable advice from Mike.

I would like to thank Prof. Calvin Lin and Prof. Emmett Witchel for reviewing this thesis. I thank Xianglong Huang, Robin Garner, David Grove, and Matthew Arnold for help with Jikes RVM and the benchmarks. I thank Jennifer Sartor and Curt Reese for their helpful suggestions for improving the paper.

Call Graph Correction Using Control Flow Constraints

Kevin Resnick, B.S.

The University of Texas at Austin, 2006

Supervisor: Prof. Kathryn S. McKinley

Dynamic optimizers for object-oriented languages collect a variety of profile data to drive optimization decisions. In particular, the *dynamic call graph* (DCG) informs key structural optimizations such as which methods to optimize and how to optimize them. Unfortunately, current low-overhead call-stack hardware and software sampling methods are subject to sampling bias, which loses accuracy of 40 to 50% when compared with a perfect call graph.

This paper introduces *DCG correction*, a novel approach that uses static and dynamic *control-flow graphs* (CFGs) to improve DCG accuracy. We introduce the static *frequency dominator* (FDOM) relation, which extends the dominator relation on the CFG to capture relative execution frequencies and expose static constraints on DCG edges, which we use to correct DCG edge frequencies. Using conservation of flow principles, we further show how to use dynamic CFG basic block profiles to correct DCG edge frequencies intraprocedurally and interprocedurally.

We implement and evaluate DCG correction in Jikes RVM on the SPEC JVM98 and DaCapo benchmarks. Default DCG sampling attains an average accuracy of 52–59% compared with perfect, whereas FDOM correction improves average accuracy to 64–68%, while adding 0.2% average overhead. The dynamic correction raises accuracy to 85% on average, while adding 1.2% average overhead. We then provide dynamically corrected DCGs to the inliner with mixed results—1% average degradations and improvements across a variety of configurations. However, prior work shows that increased DCG accuracy in production VMs has benefits. We believe that high-accuracy DCGs will become more important in the future as the complexity and modularity of object-oriented programs increases.

Contents

Acknowledgments	ii
Abstract	iii
Chapter 1 Introduction	1
Chapter 2 Background and Related Work	5
2.1 Using and Collecting Dynamic Call Graphs	5
2.2 The Dominator Relation and Strong Regions	7
2.3 Constructing the DCG using Control-Flow Information	8
Chapter 3 Call Graph Correction Algorithms	9
3.1 Terminology	9
3.2 The Frequency Dominator (FDOM) Relation	11
3.3 Static FDOM Correction	12
3.4 Dynamic Basic Block Profile Constraints	13
3.5 Dynamic Basic Block Profile Correction	16
3.6 Implementing Online DCG Correction	19
Chapter 4 Computing FDOM	21
4.1 General Control Flow Graphs	21
4.2 Reducible Control Flow Graphs	24

Chapter 5 Methodology	33
Chapter 6 Results	36
6.0.1 Accuracy	36
6.0.2 Overhead	38
6.0.3 Performance	38
Chapter 7 Conclusion	42
Bibliography	43

Chapter 1

Introduction

Dynamic optimizers for object-oriented languages ameliorate the overhead of online compilation by detecting and optimizing the most frequently executed *hot* methods. To detect hot methods and call edges, dynamic optimizers use hardware or software call-stack sampling to build a *dynamic call graph* (DCG) and selectively target optimizations. A key optimization for these systems is method inlining because well designed object-oriented programs achieve reusability, reliability, and maintainability by decomposing functionality into many small methods, and virtual dispatch obscures the hot target. Inlining exposes opportunities for other optimizations and decreases call overhead, but inlining must be applied judiciously because it also increases code size and compilation time. Good decisions depend on the accuracy of the DCG.

Dynamic optimizers trade accuracy for low overhead by using sampling to collect the DCG. Hardware-based sampling lowers overhead by examining hardware performance counters [16] instead of the call stack but gives up portability. Software-based sampling periodically examines the call stack and updates the DCG [3, 13, 19, 22]. All DCG sampling approaches suffer from sampling error, and approaches that use a timer to take samples suffer from timing bias [3]. Arnold and Grove [3] were

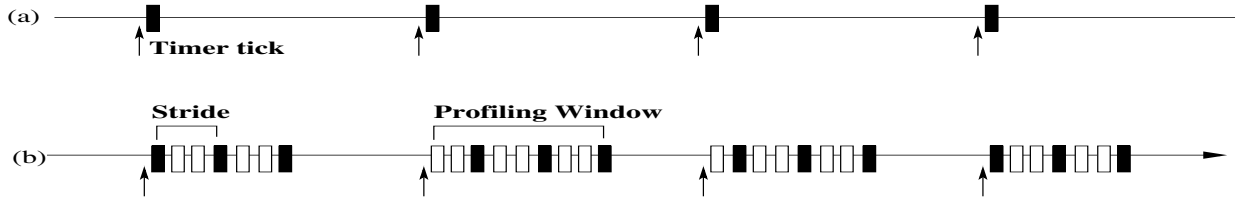


Figure 1.1: Timer-based sampling. Filled boxes are taken samples; unfilled boxes are skipped samples. (a) One sample per timer tick. (b) Counter-based Sampling (CBS) takes multiple samples per timer tick and strides between samples.

the first to measure and note that the DCG is not very accurate. They introduce counter-based sampling (CBS), which takes multiple samples per timer tick and strides over some samples to reduce sampling error and timing bias. The overhead of CBS is directly proportional to the number of samples, and they recommend configurations that achieve an average of 58% accuracy of a perfect call graph with 0.1% additional overhead in Jikes RVM and J9.

This paper presents new *DCG correction* algorithms to improve DCG accuracy with extremely low overhead (1% on average). The key insight is that a program’s static and dynamic *control-flow graph* (CFG) constrains possible DCG frequency values. For example, two calls must execute the same number of times if their basic blocks execute the same number of times. To leverage this insight, we introduce the static *frequency dominator* (FDOM) relation, which extends the dominator and strong region relations as follows: given statements x and y , x FDOM y if and only if x executes at least as many times as y . DCG correction applies FDOM constraints to the sampled DCG to improve its accuracy. For example, given two call sites cs_1 and cs_2 , if cs_1 FDOM cs_2 and in the sampled DCG, $|cs_1| < |cs_2|$, then DCG correction sets $|cs_1|$ to $|cs_2|$. Because these relationships are static, the correction algorithm computes them once and then periodically combines them with the DCG.

DCG correction uses dynamic *basic block profiles* to improve DCG accuracy.

Most dynamic optimizers collect high-accuracy control-flow profiles such as basic block and edge profiles to make better optimization decisions [1, 2, 4, 6, 9, 13, 18, 19]. We show how to combine these constraints and method counter profiling to further improve the accuracy of the DCG. This correction requires a single pass over the basic block profile, which we perform periodically.

We evaluate DCG correction in the Jikes RVM on the SPEC JVM98 and DaCapo benchmarks. We compare our approach to the CBS sampling configurations recommended by Arnold and Grove [3] that are now standard in Jikes RVM. For our benchmarks, CBS attains 52 to 60% average accuracy compared to a perfect call graph. FDOM constraints improve DCG accuracy to 64% to 68% on average, while adding less than 0.3% average overhead. Static and dynamic control-flow information together raise DCG accuracy to 85% on average, while adding just 1% overhead.

We evaluate the performance impact of corrected DCGs using inlining in Jikes RVM with its adaptive hotspot compiler that periodically recompiles and inlines hot methods. We add to this system DCG correction immediately before the system recompiles. We measure (1) application only, and (2) application plus compile times and find that DCG correction does change optimization decisions and performance but that average performance does not change. However, we find modest 3% average improvements when we provide an initial profile with a perfect call graph, and prior work shows that an accurate DCG can improve performance in other virtual machines [3]. These and other results [10] suggest that there are inlining policies that could be tuned to benefit from DCG correction. The main contributions of this paper are:

- Call graph correction algorithms using control-flow graph consistency constraints
- The frequency dominator relation and its computation

- Implementation and evaluation of call graph correction schemes
- A technique that yields highly accurate call graphs with very little overhead that can easily be added to many virtual machines

Chapter 2

Background and Related Work

This section presents background material and compares dynamic call graph (DCG) correction to previous work. We first discuss how dynamic optimizers collect a high-accuracy DCG that represents the relative frequencies of direct and virtual method calls. We then compare the frequency dominator relation to previous work. Finally we compare DCG correction to previous work that uses control-flow information to generate the DCG.

2.1 Using and Collecting Dynamic Call Graphs

Dynamic optimizers must balance increases in compilation time and code space costs with application improvement. They use call graph and control-flow information to help select optimization candidates, tailor optimizations to a specific run, and balance compile and application time. For example, while *static inlining* makes inlining decisions based on trivial criteria in dynamic optimizers (e.g., always inline a direct method call that is smaller than the calling sequence), *dynamic inlining* decisions are based on the DCG and method sizes.

Dynamic optimizers can profile every call in order to collect a *perfect DCG*,

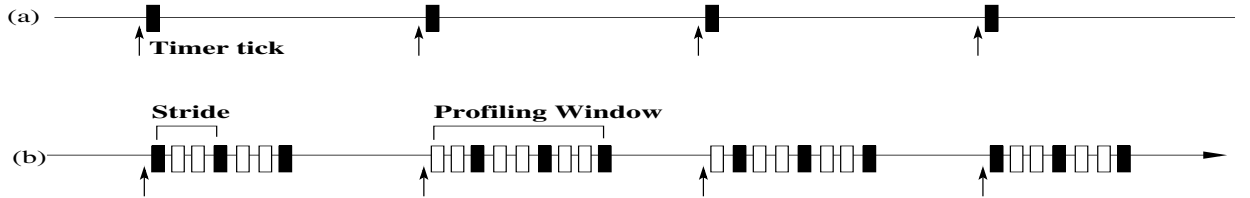


Figure 2.1: Timer-based sampling. Filled boxes are taken samples; unfilled boxes are skipped samples. (a) One sample per timer tick. (b) Counter-based Sampling (CBS) takes multiple samples per timer tick and strides between samples.

but this overhead is too high [3, 14]. Some dynamic optimizers profile calls fully for some period of time and then turn off profiling to keep overhead down. For example, HotSpot adds call graph instrumentation only for unoptimized code [19]. Suganama et al. use *code patching* to insert call instrumentation, collect call samples for a period of time, and then remove the instrumentation [22]. These *one-time profiling* approaches keep overhead down but may lose accuracy if behavior changes. Jikes RVM uses a similar approach for computing the control-flow edge profiles only for unoptimized code. DCG correction can improve the accuracy of an out-of-date DCG using up-to-date control-flow information, although this paper does not specifically evaluate DCG correction for this purpose.

Many current dynamic optimizers use software sampling to profile calls and identify hot methods while keeping overhead low [3, 5, 13, 25]. Software-based approaches examine the call stack periodically and update the DCG with the call(s) on the top of the stack. For example, Jikes RVM and J9 use a periodic timer that sets a flag that causes the next *yield-point* in the code to examine the call stack and update the DCG [5, 13]. Figure 2.1(a) illustrates basic timer-based sampling. Arnold and Grove show that this approach suffers from insufficient samples and *timing bias*: some yield-points are more likely to be sampled than others, which skews the results and possibly leads to poor inlining decisions. They present *counter-based sampling* (CBS), which takes multiple samples per timer tick and *strides* to skip some samples

in the profiling window and thus reduces timing bias. Figure 1.1(b) shows CBS configured to take 3 samples for each timer tick and to stride by 3 (take every third sample). By widening the profiling windows, CBS improves DCG accuracy, but the large profiling windows also increase profiling overhead. For example, DCG correction achieves 85% accuracy with 1% overhead whereas CBS’s overhead is 6 to 20% in J9 to achieve the same accuracy. In Jikes RVM, achieving 75% accuracy comes at an overhead of 14% or more (85% accuracy hits some pathological case, costing 1000% overhead).

Other dynamic optimizers periodically examine hardware performance counters such as those in the Itanium [16] to update the DCG. All sampling approaches suffer from sampling error, and timer-based sampling approaches suffer from timing bias as well. DCG correction can improve the accuracy without introducing significant overhead of any DCG collected by sampling.

2.2 The Dominator Relation and Strong Regions

This paper introduces the *frequency dominator* (FDOM) relation, which extends the *dominator* and *strong regions* relations [7]. Prosser first introduced dominators, which have a rich history [11, 21, 23, 24]. The set of dominators and post-dominators of x is the set of y that will execute at least once if x does. The set which frequency dominates x , on the other hand, is the subset which executes at least as many times as x . While strong regions find vertices x and y that must execute the same number of times, FDOM goes further and also finds vertices x and y where y must execute at least as many times as x .

2.3 Constructing the DCG using Control-Flow Information

This section compares DCG correction to previous work that uses control-flow information to construct the DCG. Hashemi et al. use static heuristics to construct an estimated call frequency profile [15]. Wu and Larus construct an estimated edge profile, which they use to construct an estimated call frequency profile [26]. These approaches rely solely on control-flow information to estimate call frequencies, whereas DCG correction starts with an inaccurate DCG and applies control-flow constraints to improve its accuracy. Also, these approaches use static *heuristics* to estimate frequencies, while DCG correction uses static *constraints* and dynamic profile information. Hashemi et al. and Wu and Larus report high accuracy but use an accuracy metric that considers only the relative rank of call sites, while our overlap accuracy metric uses call edge frequencies. They construct profiles for C programs, while we target Java, which has richer DCGs and multiple call targets for a call site because of virtual dispatch.

Chapter 3

Call Graph Correction

Algorithms

This section describes DCG correction algorithms. We first present formal definitions for a control flow graph (CFG) and the dynamic call graph (DCG). We introduce the *frequency dominator* (FDOM) relation that captures the relative frequencies of program statements based on control-flow graph constraints. We show how to apply these static constraints to improve the accuracy of the DCG, and how to combine them with dynamic CFG frequencies to further improve the DCG. In particular, we show how DCG correction updates the relative frequencies of call sites (which we refer to as the call sites' *OUTFLOW*) in the DCG to reflect the relative execution frequencies of basic blocks from the CFG.

3.1 Terminology

A *control-flow graph* (CFG) represents the *static* intraprocedural control flow in a method, and consists of basic blocks and edges between them. Figure 3.1 shows an example control-flow graph CFG_p that consists of basic blocks *ENTRY*, *a*, *b*, *c*, *d*, *e*,

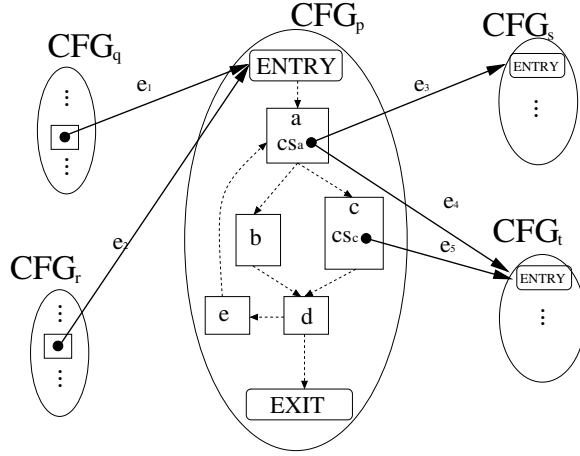


Figure 3.1: Example dynamic call graph (DCG) and control flow graphs (CFGs).

and *EXIT*, as well as edges between them. A *basic block profile* gives the *dynamic* execution frequency of each basic block in the program, from some execution.

A *call edge* represents a method call, and consists of a *call site* and a *callee*. An example call edge in Figure 3.1 is e_5 , the call from cs_c to CFG_t . The DCG of a program includes the *dynamic* frequency of each call edge, from some execution. For a call site cs , $OutEdges(cs)$ is the set of call edges that start at call site cs . $OutEdges(cs_a) = \{e_3, e_4\}$ in Figure 3.1. For a method m , $InEdges(m)$ is the set of call edges that end at m . $InEdges(CFG_t) = \{e_4, e_5\}$ in Figure 3.1.

We define the *INFLOW* of a method m as the total *flow* (execution frequency) coming into m :

$$INFLOW(m) \equiv \sum_{e \in InEdges(m)} f(e)$$

where $f(e)$ is the execution frequency of call edge e . $INFLOW(m)$ in a perfectly accurate DCG is thus equal to the number of times m executes.

We define the *OUTFLOW* of a call site cs as the total flow going out of cs :

$$OUTFLOW(cs) \equiv \sum_{e \in OutEdges(cs)} f(e)$$

$OUTFLOW(cs)$ in a perfectly accurate DCG is thus equal to the number of times cs executes.

Because the collected DCG is not accurate due to sampling error and timing bias, the DCG reports inaccurate *OUTFLOW* and *INFLOW* values. DCG correction’s goal is to correct the *OUTFLOW* of call sites using constraints provided by static and dynamic control-flow information (doing so indirectly corrects method *INFLOW* as well).

DCG correction does not correct the relative execution frequencies of multiple call edges coming out of the same call site. Instead, it uses the existing relative call edge frequencies from the uncorrected DCG, and scales them to update call sites’ *OUTFLOW* in the corrected DFG. So for example in Figure 3.1, if sampling data yields flow such that $f(e_5) > f(e_3) + f(e_4)$, DCG correction increases the total flow of cs_a , and scales it by their current relative frequencies.

3.2 The Frequency Dominator (FDOM) Relation

This section introduces the *frequency dominator* (FDOM) relation, a major contribution of this paper. FDOM is a static property of CFGs that represents constraints on program statements’ relative execution frequencies. We show two constraints (theorems) it provides from the CFG on the DFG, and we prove these relations in the appendix.

DEFINITION 1. *Frequency Dominator (FDOM).* Given statements s and t in the

same method, s *FDOM* t if and only if for every possible path through the method, s must execute at least as many times as t . We also define $FDOM(t) \equiv \{s \mid s \text{ FDOM } t\}$.

Like the dominator relation, *FDOM* is reflexive and transitive.

THEOREM 1. *FDOM OUTFLOW Constraint:* *Given method m and two call sites cs_1 and cs_2 in m , if cs_1 FDOM cs_2 ,*

$$OUTFLOW(cs_1) \geq OUTFLOW(cs_2) \tag{3.1}$$

Intuitively, the *OUTFLOW* constraint tells us that flow on two call edges is related if they are related by *FDOM*. For example, in Figure 3.1, cs_a *FDOM* cs_c and thus each time cs_c executes, cs_a must have executed.

THEOREM 2. *FDOM INFLOW Constraint:* *Given method m , if cs FDOM ENTRY (m 's entry basic block),*

$$INFLOW(m) \leq OUTFLOW(cs) \tag{3.2}$$

Intuitively, the *INFLOW* constraint specifies that a call site must execute at least as many times as a method that always executes the call site.

3.3 Static *FDOM* Correction

This section shows how to use the *INFLOW* and *OUTFLOW* constraints to correct DCG frequencies. (The next section shows how to use dynamic CFG basic block profiles to further correct *INFLOW* and *OUTFLOW*.)

Figure 3.2 applies the *FDOM OUTFLOW* constraint to a sampled DCG. The algorithm *FDOMOutflowConstraint* compares the sampled *OUTFLOW* of pairs of call sites that satisfy the *FDOM* relation. If their *OUTFLOW*s violate the *FDOM OUTFLOW* constraint, *FDOMOutflowConstraint* sets both *OUTFLOW*s to the maximum of their two *OUTFLOW*s. After processing a method, *FDOMOutflowConstraint* scales the *OUTFLOW*s of all the method’s call sites so the sum remains the same as before.

Similarly, Figure 3.3 applies the *FDOM INFLOW* constraint to a sampled DCG. For every call site *cs* that must execute each time the method executes (*cs FDOM ENTRY*), *FDOMInflowConstraint* sets *OUTFLOW(cs)* to *INFLOW(m)* if *OUTFLOW(cs) < INFLOW(m)*.

3.4 Dynamic Basic Block Profile Constraints

This section describes constraints on DCG frequencies provided by basic block profiles, and the following section shows how to correct the DCG with them. The appendix proves these relations.

The *Dynamic OUTFLOW* constraint says that the ratio between the execution frequencies of two call sites specified by the basic block profile can be applied to the *OUTFLOW* of these two call sites:

THEOREM 3. *Dynamic OUTFLOW Constraint* *Given two call sites cs_1 and cs_2 , and execution frequencies $f_{\text{bprof}}(cs_1)$ and $f_{\text{bprof}}(cs_2)$ provided by a basic block profile,*

$$\frac{\text{OUTFLOW}(cs_1)}{\text{OUTFLOW}(cs_2)} = \frac{f_{\text{bprof}}(cs_1)}{f_{\text{bprof}}(cs_2)} \tag{3.3}$$

procedure FDOMOutflowCorrection**input:***CALLSITES*: a set of call sites*FDOM*: a frequency dominator relationship between call sites $f_{sample}(e)$: a function that returns the frequency of call edge e from sampling $f_{sample}(cs)$: a function that returns the frequency sum of call edges in $OutEdge(cs)$ from sampling**output:** $f_{corrected}(e)$: corrected frequency of call edge e

```

1: {STEP1: Initialize outflow for each call site and its sum.}
2:  $sum_{old} \leftarrow 0$ 
3: for all  $cs \in CALLSITES$  do
4:    $Outflow(cs) \leftarrow f_{sample}(cs)$ 
5:    $sum_{old} \leftarrow sum_{old} + Outflow(cs)$ 
6: end for
7: {STEP2: satisfy FDOM Outflow constraint.}
8:  $sum_{new} \leftarrow sum_{old}$ 
9: for all  $cs_y \in CALLSITES$  do
10:  for all  $cs_x \in FDOM(cs_y)$  do
11:    {Force  $OUTFLOW(cs_x) \geq OUTFLOW(cs_y)$ }
12:     $outflow_{old} \leftarrow Outflow(cs_x)$ 
13:     $Outflow(cs_x) \leftarrow Max(Outflow(cs_x), Outflow(cs_y))$ 
14:     $diff = Outflow(cs_x) - outflow_{old}$ 
15:     $sum_{new} \leftarrow sum_{new} + diff$ 
16:  end for
17: end for
18: {STEP3: Use new outflow to derive the corrected frequency.}
19:  $scale \leftarrow sum_{old} / sum_{new}$ 
20: for all  $cs \in CALLSITES$  do
21:  for all  $e \in OutEdges(cs)$  do
22:     $fraction = f_{sample}(e) / f_{sample}(cs)$ 
23:    {Preserve the call target fraction and the frequency sum.}
24:     $f_{corrected}(e) \leftarrow Outflow(cs) \times scale \times fraction$ 
25:  end for
26: end for

```

Figure 3.2: DCG Correction with FDOM OUTFLOW Constraints

procedure FDOMInflowCorrection

input:

p : a procedure to apply FDOM inflow constraint

$FDOM(ENTRY)$: a set of call sites that frequency-dominates the entry of the procedure P

$f_{sample}(e)$: a function that returns the frequency of call edge e from sampling

$f_{sample}(cs)$: a function that returns the frequency sum of call edges in $OutEdge(cs)$ from sampling

output:

$f_{corrected}(e)$: corrected frequency of call edge e

```
1: {STEP1: Initialize outflow for each call site and its sum.}
2:  $inflow \leftarrow 0$ 
3: for all  $e \in InEdges(p)$  do
4:    $inflow \leftarrow inflow + f_{sample}(e)$ 
5: end for
6: {STEP2: satisfy FDOM inflow constraint.}
7: for all  $cs_x \in FDOM(ENTRY)$  do
8:   {Force  $OUTFLOW(cs_x) \geq INFLOW(p)$ }
9:    $outflow_{old} \leftarrow Outflow(cs_x)$ 
10:   $Outflow(cs_x) \leftarrow Max(Outflow(cs_x), inflow)$ 
11: end for
12: {STEP3: Use new outflow to derive the corrected frequency.}
13: for all  $cs \in FDOM(ENTRY)$  do
14:   for all  $e \in OutEdges(cs)$  do
15:      $fraction = f_{sample}(e) / f_{sample}(cs)$ 
16:     {Preserve the call target fraction}
17:      $f_{corrected}(e) \leftarrow Outflow(cs) \times fraction$ 
18:   end for
19: end for
```

Figure 3.3: DCG Correction with FDOM INFLOW Constraints

The *Dynamic OUTFLOW* constraint can be applied to two call sites in different methods if basic block frequencies from different methods are accurate relative to each other (i.e., if the basic block profiles have *interprocedural accuracy*). In our implementation, basic block profiles do *not* have interprocedural accuracy. Thus, we do not apply the *Dynamic OUTFLOW* constraint to call sites in different methods. If basic block profiles do not have interprocedural accuracy, then the *Dynamic OUTFLOW* constraint provides no help for correcting the *OUTFLOW* of call sites in methods with a single basic block. We experiment with using low-overhead method invocation counters to give basic block profiles interprocedural accuracy, and in this case we apply *Dynamic OUTFLOW* to call sites in different methods (Section 3.6).

The *Dynamic INFLOW* constraint says that the call edge flow (frequency) coming into a method with a single basic block constrains the flow leaving any call site in the method:

THEOREM 4. *Dynamic INFLOW Constraint:* *Given a method m with a single basic block and a call site cs in m ,*

$$\text{INFLOW}(m) = \text{OUTFLOW}(cs) \tag{3.4}$$

The *Dynamic INFLOW* constraint is useful for methods with a single basic block because the *Dynamic OUTFLOW* constraint cannot constrain the *OUTFLOW* of call sites in the single basic block (when basic block profiles do not have interprocedural accuracy). The *Dynamic INFLOW* constraint uses the total flow (frequency) coming into the method to constrain call sites' *OUTFLOW*.

3.5 Dynamic Basic Block Profile Correction

procedure DynamicOutflowCorrection

input:

CALLSITES: a set of call sites

$f_{bprof}(cs)$: a function that returns the frequency of the call site cs from basic block profiles

$f_{sample}(e)$: a function that returns the frequency of call edge e from sampling

output:

$f_{corrected}(e)$: a function that returns the corrected frequency for the call edge e

- 1: **{STEP1:}** Iterate call sites to find scale factor from basic block profile count to the sample count.
- 2: $sum_{sample} \leftarrow \sum_{cs \in CALLSITES} f_{sample}(cs)$
- 3: $sum_{prof} \leftarrow \sum_{cs \in CALLSITES} f_{prof}(cs)$
- 4: $scale \leftarrow sum_{sample} / sum_{prof}$
- 5: **{STEP2:}** assign corrected call edge frequency
- 6: **for all** $cs \in CALLSITES$ **do**
- 7: **for all** $e \in OutEdges(cs)$ **do**
- 8: $fraction = f_{sample}(e) / f_{sample}(cs)$
- 9: {Preserve the call target fraction and the frequency sum.}
- 10: $f_{corrected} \leftarrow f_{bprof}(cs) \times scale \times fraction$
- 11: **end for**
- 12: **end for**

Figure 3.4: DCG Correction with Dynamic OUTFLOW Constraints

procedure DynamicInflowCorrection

input:

p : a single basic block procedure

$f_{sample}(e)$: a function that returns the frequency of call edge e from sampling

$f_{sample}(cs)$: a function that returns the frequency sum of call edges in $OutEdge(cs)$ from sampling

output:

$f_{corrected}(e)$: corrected frequency of call edge e

- 1: $CALLSITES \leftarrow getCallSiteInsideProcedure(p)$
- 2: {**STEP1**: Compute maxflow for the procedure p .}
- 3: $inflow \leftarrow \sum_{e \in InEdges(p)} f_{sample}(e)$
- 4: $maxoutflow \leftarrow \max_{cs \in CALLSITES} f_{sample}(cs)$
- 5: $maxflow \leftarrow \max(inflow, maxoutflow)$
- 6: {**STEP2**: assign corrected frequency}
- 7: **for all** $cs \in CALLSITES$ **do**
- 8: **for all** $e \in OutEdges(cs)$ **do**
- 9: $fraction = f_{sample}(e) / f_{sample}(cs)$
- 10: {constraint: INFLOW(p) = OUTFLOW(cs).}
- 11: $f_{corrected}(e) \leftarrow maxflow \times fraction$
- 12: **end for**
- 13: **end for**

Figure 3.5: DCG Correction with Dynamic INFLOW Constraints

Correction algorithm	Correction range	Algorithms
Static FDOM CF Correction	Call sites within a method to be optimized	Figures 3.2 and 3.4
Dynamic Interprocedural CF Correction	Call sites within a method to be optimized	Figures 3.4 and 3.5
Dynamic Intraprocedural CF Correction	All call sites in the DCG	Figures 3.4 and 3.5

Table 3.1: Call Graph Correction Implementations

Figure 3.4 presents the algorithm for applying the *Dynamic OUTFLOW* constraint. *DynamicOutflowCorrection* sets the *OUTFLOW* of each call site cs to $f_{bprof}(cs)$, its frequency from the basic block profile. The algorithm then scales all the *OUTFLOW* values so that the method’s total *OUTFLOW* is the same as before. This scaling helps to maintain the frequencies due to sampling across disparate parts of the DCG.

Figure 3.5 presents the algorithm for applying the *Dynamic INFLOW* constraint to the DCG. For each method with a single basic block, *DynamicInflowCorrection* sets the *OUTFLOW* of each call site in the method to the *INFLOW* of the method.

3.6 Implementing Online DCG Correction

Dynamic compilation systems perform profiling while they execute and optimize the application, and therefore DCG correction needs to be done at the same time with minimal overhead.

We minimize the call graph correction overhead by limiting the frequency and scope of DCG correction. We limit DCG correction’s frequency by delaying DCG correction until the optimizing compiler requests DCG information. The correction overhead is thus proportional to the number of methods optimized during an execution. Correction overhead is thus naturally minimized as a result of the dynamic optimizers’ selective choices about when and which methods to recompile.

We limit the scope of DCG correction by localizing the range of correction.

When the compiler optimizes a method m , it does not require the entire DCG, but instead considers a localized portion of the DCG relative to m . Because we preserve the call edge frequency sum in the *OUTFLOW* correction algorithm, we can correct m and all the methods it invokes without compromising the correctness of the other portions of the DCG. Because we preserve the DCG frequency sum, the normalized frequency of a call site in a method remains the same, independent of whether call edge frequencies in other methods are corrected or not.

Table 3.1 summarizes the correction algorithms and the correction range. All the correction algorithms take as input the set of call sites to be corrected. Clearly, for FDOM correction, the basic unit of correction is the call sites within a procedure boundary.

For dynamic basic block profile correction, there are two options. The first one limits the call site set to be within a procedural boundary, and the second one corrects all the reachable methods. Since many dynamic compilation systems support only high precision intraprocedural basic profiles, the first configuration indicates how much DCG correction would benefit these systems.

Because our system does not collect interprocedural basic block profiles, we implement interprocedural correction by instrumenting each method with a method counter. We multiply the counter value by the normalized intraprocedural basic block frequency. We find this mechanism is a good approximation to interprocedural basic block profiles.

Chapter 4

Computing FDOM

This section presents our algorithm for computing the *frequency dominator* (FDOM) relation for statements in a method. We first show a correlation between *simple cycles* and FDOM that applies to both *irreducible* and *reducible* graphs. We then present our algorithm for computing FDOM, which applies to reducible graphs only. Throughout this section, we use s in place of *ENTRY* and e in place of *EXIT* for brevity. Both algorithms assume the existence of a back edge from e to s . This edge simplifies the analysis but does not affect the FDOM relation, since following this edge is equivalent to executing the method again.

4.1 General Control Flow Graphs

In this section, we show that the FDOM relation relates to *simple cycles* for general CFGs (both irreducible and reducible). This relation could be used to compute FDOM for irreducible CFGs. However, we only present an algorithm that computes FDOM for reducible CFGs, which we present in the next section. Lemma 1 shows that FDOM can be computed by considering *cycles* in the CFG.

LEMMA 1. *Given vertices x and y in V ,*

$x \text{ FDOM } y$
if and only if
 x is contained in every cycle containing y

Proof. Show both forward and backward directions.

(\Rightarrow) (by contradiction) Suppose there is a cycle that contains y but not x . Let this cycle be $c_y = \langle a, \dots, y, \dots, a \rangle$. From the definition of CFG, there is one path from the CFG entry, s , to a and another path from a to the CFG exit, e . By concatenating these three paths, we can build a path from s to e , $c'_y = \langle s, \dots, a, \dots, y, \dots, a, \dots, e \rangle$. If x is not in c'_y then y is already executed more times than x , a contradiction. So let x in some execution path that includes c'_y be executed n times, then we can repeat c_y $n + 1$ times, resulting in y being executed more than x , a contradiction.

(\Leftarrow) (by contrapositive) Suppose there is a path p from s to e where the number of executions of y exceeds that of x . For the number of y to exceed x , p must contain at least one y . If the number of y in the path is n , then path p should be of the form, $p = \langle s, \dots, y_1, \dots, y_2, \dots, y_n, \dots, e \rangle$. The path, p , can be divided into $(n + 1)$ subpaths: $\langle s, \dots, y_1 \rangle, \langle y_1, \dots, y_2 \rangle, \dots, \langle y_n, \dots, e \rangle$. Because the number of x is less than the number of y in p , the number of x in all of the subpaths is less than or equal to $(n - 1)$. From the pigeonhole principle, there exist at least two x -free subpaths. If one of the two x -free subpaths is $\langle y_i, \dots, y_{i+1} \rangle$, this x -free subpath contradicts our assumption that there exists no cycle containing y and not x . If two subpaths are $\langle s, \dots, y_1 \rangle$ and $\langle y_n, \dots, e \rangle$, another x -free cycle $\langle y_n, \dots, e, s, \dots, y_1 \rangle$ can be constructed since we have a backedge from e to s , again a contradiction. \square

Lemma 1 states that FDOM can be computed by considering every cycle in the CFG. However, it is impractical to check if x is contained in every cycle that contains y because the number of cycles may be unbounded. The number of simple cycles in a method is bounded, and Lemma 2 shows that *simple cycles* are sufficient.

LEMMA 2. *Given vertices x and y in V ,*

*x is in every cycle containing y
if and only if
 x is in every simple cycle containing y*

Proof. Forward direction is trivial because every simple cycle is a cycle. To show the backward direction:

(by contradiction) Suppose that x is in every simple cycle containing y , but let there exist x -free cycle c that contains y , $c = \langle a, \dots, y, \dots, a \rangle$ or $c = \langle y, \dots, y \rangle$. By our assumption, c cannot be simple, so there must exist some element in c other than the beginning or end. There exists a cycle, therefore, in c that is from $\langle w, \dots, w \rangle$ that is simple, which we will call c' . Since by our assumption, x is in every simple cycle containing y , if there exists a y in c' , then there also exists an x , and since c' is simple, there can exist at most one of each. Therefore if c is a valid cycle including c' , then another valid cycle is c with c' replaced with the single element w . Further, this new c still has the property that there exists an x -free cycle that contains y . Since c is of finite length, this process can be continued until there are no simple cycles left in c . Note that c still contains an x -free path that contains y , but now c is simple too, a contradiction. This proof holds for the entire program because there exists a back edge, by assumption, from e to s , and hence it is a cycle.

□

Theorem 5 provides a method for computing FDOM, and it is easily proved from the lemmas.

THEOREM 5. *Given vertices x and y in V ,*

x FDOM y

if only if
x is in every simple cycle containing y.

Proof. From Lemma 1 and Lemma 2. □

4.2 Reducible Control Flow Graphs

This section presents our near-linear-time algorithm for computing the FDOM relation for reducible CFGs. The algorithm uses CFG properties from previous work to compute the FDOM relation. We first present background information on these properties. We then prove that FDOM can be defined in terms of these properties. Finally, we present an algorithm that first computes these properties and then uses the resulting data structures to compute the FDOM relation.

We first describe CFG properties from previous work. For each loop header $h \in V$, the following set is defined by Ball [7]:

$$\begin{aligned}
 \text{backsrcs}^*(h) &= \{v \mid v \rightarrow h \text{ is a backedge.}\} \\
 \text{natloop}^*(h) &= \{h \mid \text{there is an } h\text{-free path from} \\
 &\quad v \text{ to a basic block in } \text{backsrcs}(h)\} \\
 \text{exits}^*(h) &= \{v \mid \exists v \rightarrow w \text{ such that } v \in \text{natloop}(h) \\
 &\quad \text{and } w \notin \text{natloop}(h)\}
 \end{aligned}$$

We extend these three definitions to include any basic block $y \in V$.

$$\begin{aligned}
loophead(y) &= \text{loop entry of innermost loop that contains } y \\
backsrcs(y) &= backsrcs^*(loophead(y)) \\
natloop(y) &= natloop^*(loophead(y)) \\
exits(y) &= exits^*(loophead(y))
\end{aligned}$$

Ball defines “ w pd v with respect to a set of vertices V ” to capture post dominance within a natural loop [7]. For our algorithm, we only use two sets of vertices, back edges and loop exits for a given natural loop. Thus we use $PDBE$ y to mean $backsrcs(y)$ and $PDLE$ y to mean $exits(y)$. We combine these terms and make a new term, PDL $y \equiv PDBE$ $y \cup PDLE$ y .

Theorem 6 defines FDOM in terms of other CFG properties. Ball [7] and Johnson et al. [17] use sufficient and necessary conditions for FDOM in Lemma 1 to characterize *control regions* [12]. The algorithm in this section is motivated by Ball’s paper [7].

THEOREM 6. *Given two vertices x and y in V ,*

$$\begin{aligned}
&x \text{ FDOM } y \\
&\textit{if and only if} \\
&x \in \text{natloop}(y) \textit{ and } (x \text{ DOM } y \textit{ or } x \text{ PDL } y)
\end{aligned}$$

Proof. (\Rightarrow) If x is not in $\text{natloop}(y)$, then there exists an x -free cycle that contains y , a contradiction due to Lemma 1. So suppose that $x \text{ DOM } y$ or $x \text{ PDL } y$ is not true. If $x \text{ DOM } y$ is not true, then there exists a path from $loophead(y)$ to y that does not include x . If $x \text{ PDL } y$ is not true, then there exists a path from y to either a back edge or an exit of $\text{natloop}(y)$. Hence there exists either a path $\langle loophead(y), \dots, y, \dots, exit(y) \rangle$, which means there exists a path from s to e that

includes y and not x , or there exists a cycle $\langle loophead(y), \dots, y, \dots, backsrcs(y) \rangle$, which is a cycle that includes y and not x , also a contradiction. Hence $x \text{ DOM } y$ or $x \text{ PDL } y$ must also be true.

(\Leftarrow) Show that every cycle that starts and ends at y , contains x . Let h_y be $loophead(y)$, and every path from y to y is one of this form, $c_y = \langle y, \dots, z, \dots, h_y, \dots, y \rangle$, where z is in $backsrcs(y) \cup exits(y)$. Suppose that $x \in natloop(y)$, and $x \text{ DOM } y$, then every path from h_y to y must contain x . If $x \text{ PDL } y$, then this implies that $\langle y, \dots, z \rangle \subset c_y$ always contain x . Therefore c_y contains x if either of the two conditions is satisfied. \square

$FDOM(y)$ is the set of all x s.t. $x \text{ FDOM } y$, and from Theorem 6, the following equations hold:

$$\begin{aligned} FDOM(y) &= natloop(y) \cap (DOM(y) \cup PDL(y)) \\ &= FDOM_D(y) \cup FDOM_P(y) \\ FDOM_D(y) &= natloop(y) \cap DOM(y) \\ FDOM_P(y) &= natloop(y) \cap PDL(y) \end{aligned}$$

We create an algorithm that can compute $FDOM_D(y)$ and $FDOM_P(y)$ in near-linear time.

Pingali and Bilardi give the paradox of linear computation time for super-linear-sized relations such as *control dependence* and *post-dominance* [20]. This paradox is resolved by the fact that these relations are transitive, and these relations can be factored into the transitive reduction form. The preprocessing time is used to construct some data structure that describes this reduction form, and a query is performed on this data structure. For instance, the post-dominance relation is transitive, and its transitive reduction form is a *post-dominator tree*, which can be


```

procedure setupFDOM (G:CFG)
1: Compute LoopStructureTree(G)
2: Compute DominatorTree(G)
3: Transform(G)
4: Compute PostdominatorTree(G)
procedure Transform(G:CFG)
1: for all  $L \in \text{natloops}(G)$  do
2:   Create  $\text{tempnode}_L$ 
3:   {Redirect each exit edge to go through the temp node.}
4:   for all edges  $b \rightarrow k$  such that  $b \in L, k \notin L, b$  and  $k \neq \text{tempnode}_L$  do
5:     remove  $b \rightarrow k$ 
6:     add  $b \rightarrow \text{tempnode}_L$ 
7:     add  $\text{tempnode}_L \rightarrow k$ 
8:   end for
9:   {Redirect each back edge to go through the temp node.}
10:  for all edges  $b \rightarrow k$  such that  $b \in L, k = \text{loophead}(L), b \neq \text{tempnode}_L$  do
11:    remove  $b \rightarrow k$ 
12:    add  $b \rightarrow \text{tempnode}_L$ 
13:    add  $\text{tempnode}_L \rightarrow k$ 
14:  end for
15: end for

```

Figure 4.1: FDOM setup algorithm

computed in $O(|E|)$ time.¹ We describe an algorithm to compute FDOM in $O(|E|)$ if a dominator tree and a loop structure tree is precomputed. We face this same issue, and will construct a data structure containing super-linear information in super-linear time.

There are a few observations from the transformation in Figure 4.1. First, for each loop l there exists a unique tempnode_l . Second, due to the *forall* statement, the source vertex of every loop exit or back edge in l now points to tempnode_l . Further, tempnode_l is the only source vertex of a back edge or loop exit in l . Finally, because of it being the only source vertex for an exit, tempnode_l post-dominates every node in $\text{natloop}(l)$ in the transformed graph.

¹Cooper et al. describe the history of algorithms for dominance relation in detail [11].

```

procedure retrieveFDOM (y:Vertex):Set
1:  $s \leftarrow \{y\}$ 
2: {Compute  $FDOM_D(y)$ }
3:  $current \leftarrow \text{dominatorParent}(y)$ 
4: while  $current \neq \text{null}$  AND  $current \neq \text{loophead}(y)$  do
5:    $s \leftarrow s \cup \{current\}$ 
6:    $current \leftarrow \text{dominatorParent}(current)$ 
7: end while
8: {Compute  $FDOM_P(y)$ }
9:  $current \leftarrow \text{postdominatorParent}(y)$ 
10: while  $current \neq \text{null}$  AND  $current \neq \text{tempnode}(y)$  do
11:    $s \leftarrow s \cup \{current\}$ 
12:    $current \leftarrow \text{postdominatorParent}(current)$ 
13: end while
14: return  $s$ 

```

Figure 4.2: FDOM retrieval algorithm

Theorem 7 shows that we can compute $PDL(y)$ inside a loop by applying the post-dominator algorithm on the transformed graph G' .

THEOREM 7. *Given two vertices x and y in a CFG G , $l = \text{natloop}(y)$, $x \in l$, and $G' = T(G)$, T as defined above,*

$$\begin{aligned}
 &x \text{ PDL } y \text{ in } G \\
 &\quad \text{if and only if} \\
 &x \text{ PDOM } y \text{ in } G'
 \end{aligned}$$

Proof. (\Rightarrow) Assume $x \text{ PDL } y$ in G is true.

\Rightarrow In G , there is no x -free path from y to any z in $\text{backsrcs}(y) \cup \text{exit}(y)$.

\Rightarrow In G' , there is no x -free path from y to tempnode_l .

\Rightarrow In G' , since tempnode_l postdominates every node in l , there is no x -free path from y to $EXIT$

$\Rightarrow x \text{ PDOM } y$ in G'

(\Leftarrow)(by contrapositive) Assume $x \text{ PDL } y$ in G is false.

x not *PDL* y

\Rightarrow In G , there is an x -free path from y to some z in $backsrcs(y) \cup exits(y)$.

\Rightarrow In G' , there is an x -free path from y to $tempnode_l$.

\Rightarrow Since $tempnode_l$ post-dominates every node in l , in G' , there is an x -free path from y to *EXIT*.

\Rightarrow In G' , x *PDOM* y is not true. □

Figure 4.1 shows a fast algorithm that computes $FDOM(y)$ in near-linear time based on these proofs. For a high-level description, we compute the dominator tree and the loop structure tree for the CFG, then we apply a transformation to the CFG and compute the post-dominator tree on the new graph. The transformation creates a single node per loop that the loop exits and back edges all go through. This is all the computation that is necessary to compute *FDOM*, and hence is near-linear time. This computation is bounded by the time complexity of constructing the loop structure tree, or $O(|V| + |E|\alpha(|E|, |V|))$ [7]. If the loop structure tree has already been computed, as it is in many compilers, then this algorithm is $O(|E|)$.

To retrieve the *FDOM* information for a given node y , we simply iterate up the dominator tree from y until we hit y 's loop header. We then iterate up the post-dominator tree on the transformed graph until we hit y 's tempnode. We are guaranteed that from y to $loophead(y)$ in the dominator tree and from y to $tempnode(y)$ in the postdominator tree, we will not leave $natloop(y)$. Further, once we go beyond these nodes, we will no longer be in $natloop(y)$. Hence we are only considering the set x s.t. $x \in natloop(y)$ (see Figure 4.2).

FDOM Algorithm History

We investigated several *FDOM* algorithms before arriving at the one above. This section discusses these algorithms and the insights we developed. The one presented earlier is our final and simplest version of the algorithm.

The original algorithm still computed the dominator and post-dominator portion of the control flow graph separately. But the way it did so was more complicated. Instead of simply computing the dominator tree and then walking up it, the original algorithm had the information flow down the dominator tree. Also, instead of transforming the graph as described above and computing the post-dominator over the entire graph at once, it condensed each loop and computed the post-dominator information on each individual loop.

To have the dominator information flow down the dominator tree, we employed a persistent stack. A persistent stack is one in which every version of the stack is available at any time. This data structure can be implemented by modifying a singly linked list with its pointers reversed, so it is pointing towards the "head". To push, simply create a new node and point it to the top of the stack. To pop, move the given "top" pointer to top's next pointer. If there exists multiple "top" pointers, and the data in the stack does not change, then there essentially exists multiple stacks all sharing data.

By using persistent stacks, we can share information and hence prevent producing the $O(n^2)$ relation. Each element on the stack is another singly linked list. At a given node in the dominator tree Y , each level of the stack represents the lowest loop nest depth which must be traveled through on a path from X to Y . Hence the top-most list on the stack represents all the nodes which $FDOM_D Y$, as at any point Y the stack has as many elements as the loop nest depth of Y .

Now the problem is how to create the above data structure at every point. Going from the root of the dominator tree down, each child will have its own "top" pointer. In addition, each node will have a length of the top list of the stack. At each node Y , there are three possibilities for the stack from Y 's immediate dominator (Y 's parent). The stack size is less than, equal to, or greater than the loop nest depth of Y . If it is less than, it can only be less than by one (because the control flow graph

is reducible, we know Y must be a loop header), so we will create a new linked list and push it onto the stack. If it is equal to the loop nest depth, then we will use the linked list on top of the stack. If it is greater than, we need to combine the topmost lists until the stack is the right size. By pointing the tail of the topmost list to the head of the second topmost list, and repeating, this can easily be accomplished. Since each node has a list length, the chaining together will not change the list it has. We then add the node itself to the head of the front of the linked list.

There is one additional complication to the above algorithm. If Y is a loop header, then the stack size needs to be one less than its loop nest depth. This requirement is due to there being some point immediately prior to Y which went one level below Y 's loop nest depth. The above algorithm works with this small change.

Now at every node there exists a stack of linked lists, the topmost of which is the set that $FDOM_D Y$. To retrieve this set, simply iterate over the linked list, keeping in mind the length stored at the node.

As to the $FDOM_P$ aspect of the algorithm, we used an idea from Ball [7]. For computing Strong Regions, Ball reduced each loop subgraph to a single point from the point of view of other loops. So at the end of the transformation, there is a set of loops which contain no inner loops. On each element of this set we apply the transformation to turn it into a DAG as in our current implementation, by having all loop exits and backedges point to a temporary node. We then compute the post-dominator. In a given DAG there may be a node which represents a more deeply nested loop. If one of these nodes, Y , post-dominates X in a DAG, then all nodes which post-dominate the loop header of Y in Y 's DAG must also post-dominate X .

In time we realized a few key aspects to help simplify the algorithm. First, there is no need to reduce the control flow graph into several independent DAGs. Instead, we apply the transformation to get a single point of exit for each loop in

our current control flow graph. Further, we found that instead of traversing down the dominator tree, traversing up was more natural, and already in the proper form. Hence we were able to reduce this initial algorithm to simply computing the dominator tree, transforming the graph, and then computing post-dominator.

Although as of now unproven, we believe there exists a simple fast algorithm to compute FDOM on a non-reducible control flow graph which is similar to the one described earlier in this paper. In our previous transformation make a single point of exit of a loop and then compute post-dominator. Let there be another transformation which makes a single point of entry to every loop and compute the dominator. It is unclear whether both transformations can be applied to the same graph, or if each must work on a copy of the original.

Chapter 5

Methodology

This section describes our benchmarks, platform, implementation, and VM compiler configurations. We describe our methodologies for accuracy measurements against the perfect dynamic call graph (DCG), overhead measurements, and performance measurements.

We implemented and evaluated DCG correction algorithms in Jikes RVM (CVS HEAD 2005/08/13 21:10:06 UTC) a Java-in-Java VM in its production configuration [2]. This configuration pre-compiles the VM methods (e.g., compiler and garbage collector) and any libraries it calls into the boot image. Jikes RVM contains two compilers: the *baseline compiler* and *optimizing compiler* with three optimization levels. (There is no interpreter in this system.) When a method is first executed, the baseline compiler generates assembly code (x86 in our experiments). A call-stack sampling mechanism identifies frequently executed (*hot*) methods. Based on these method sample counts, the *adaptive compilation system* then recompiles methods at progressively higher levels of optimization. Because it is sample based, the adaptive compiler is non-deterministic.

We use the SPEC JVM98 benchmarks and the DaCapo benchmarks [8]. We omit the DaCapo benchmarks *xalan*, and *hsqldb* because we could not get *xalan* to

run correctly with Jikes RVM, and *hsqldb* showed very large run-to-run variation in its execution due to Jikes RVM’s thread scheduling implementation. We omit the DaCapo benchmark *antlr* from the inlining performance experiments because we could not get it to run correctly for 25 iterations.

We perform our experiments on a 3.2 GHz Pentium 4 with hyper-threading enabled. It has a 64-byte L1 and L2 cache line size, an 8KB 4-way set associative L1 data cache, a 12K μ ops L1 instruction trace cache, a 512KB unified 8-way set associative L2 on-chip cache, and 1GB main memory, and runs Linux 2.6.0.

Accuracy Methodology. To measure the accuracy of our technique against the perfect DCG for each application, we first generate a perfect DCG by modifying Jikes RVM call graph sampling to sample every method call (instead of skipping). We also turn off the optimizing compiler to eliminate non-determinism due to sampling and since call graph accuracy is not influenced by code quality. By the end of the application’s execution, the sampler will have collected the perfect call graph. We restrict the call graph to the application methods by excluding all call edges with both the source and target in the boot image, and calls from the boot image to the application. We include calls edges into the boot image, since these represent calls to libraries that the compiler may want to inline into the application.

To measure and compare call graph accuracy, we compare the final perfect DCG to the final corrected DCG generated by our approach. Because DCG clients use incomplete graphs to make optimization decisions, we could have compared the accuracy of the instantaneous perfect and corrected DCGs as a function of time. We follow prior work in comparing the final graphs [3] rather than a time series, and believe these results are representative of the instantaneous DCGs.

We measure accuracy against two (final) perfect DCGs: (1) *without inlining* and (2) *with trivial inlining*. The first configuration has the largest DCG, but the second is more representative of the base call graph presented to the optimizing

compiler. Trivial inlining in Jikes RVM inlines a call site if the size of the callee is smaller than the calling sequence. The inliner will therefore never need to use the frequency information for these call sites.

Overhead Methodology. To measure the overhead of DCG correction without including its influence on optimization decisions, we configure the adaptive compiler to perform only *trivial inlining*. This configuration does not use the DCG to make inlining decisions, so correction does not affect inlining decisions.

Performance Methodology. We use the following configuration to measure the performance of using corrected DCGs to drive inlining. We correct the DCG on the fly by performance correction as the VM optimizes the application, providing a realistic measure of DCG correction’s ability to affect inlining decisions. We measure steady-state performance by executing 25 iterations of the application in one run of the VM, and we take the median time of iterations 13 through 25, which we find decreases variation due to the non-determinism of the adaptive compiler.

Chapter 6

Results

This section evaluates the accuracy, overhead, and performance effects of the DCG correction algorithms.

6.0.1 Accuracy

We use the overlap accuracy metric from prior work to compare the accuracy of DCGs [3].

$$\begin{aligned} \text{overlap}(DCG_1, DCG_2) = \\ \sum_{e \in \text{CallEdges}} \min(\text{weight}(e, DCG_1), \text{weight}(e, DCG_2)) \end{aligned}$$

where *Call Edges* is the intersection of the two call edge sets in DCG_1 and DCG_2 respectively, and $\text{weight}(e, DCG_x)$ is the normalized frequency for a call edge e in DCG_x . We use this function to compare the perfect DCG to the other DCGs.

We compute an upper bound for our correction schemes based on the call edges in the sampled DCG . For this accuracy upper bound, we consider DCG_{sample} , a sampled DCG, and $DCG_{perfect}$. Given $DCG_{perfect}$ and DCG_{sample} , we define DCG_{bound} by taking the call nodes and edges from the DCG_{sample} and by taking the edge frequencies from $DCG_{perfect}$. Then, the frequencies of call edges in DCG_{bound}

are completely unbiased, but DCG_{bound} only contains the methods and calls edges in the sampled graph. Therefore, DCG_{bound} is the best that we can hope for correction to achieve.

Figures 6.1 and 6.2 compare DCG accuracy for the no-inlining and trivially-inlining configurations. For the individual benchmarks, we report average profile accuracy over 25 trials (shown as dots). The perfect DCG is 100% (not shown). Since the DCGs are similar between the two figures, the accuracies are as expected also similar. The graphs compare the perfect DCG to the base system without correction (*Base*), *Static FDOM CF Correction* (only), *Dynamic Intraprocedural CF Correction* (only), *Dynamic Interprocedural CF Correction*, and the *Correction Upper Bound* (DCG_{bound}).

The DCG with no correction (*Base*) uses the recommended sampling and striding settings from Arnold and Grove [3] that take more samples to increase DCG accuracy, but keep average overhead under a few percent. This configuration has an average accuracy of 52% (no inlining) and 60% (static inlining), and drops as low as 20%.

Static and dynamic DCG correction together significantly improve accuracy to an average of 85%. The second bar, *Static FDOM CF Correction*, teases apart the contribution just from static FDOM, which improves over sampling by 7 to 10% on average. Intraprocedural correction performs about the same, but interprocedural correction applies global dynamic constraints and shows much better accuracy than the local methods. In addition, interprocedural correction comes within 10% of the bound, which loses only on average 5% to missing call edges. Interprocedural correction is 10% less than the upper bound because (1) DCG correction does not correct the relative frequencies of call edges coming out of the same call site, and (2) the basic block profiles are slightly inaccurate because they measure only initial execution behavior [9].

6.0.2 Overhead

Figure 6.3 presents the execution times for various DCG correction configurations using the adaptive compiler in which the corrected DCGs are never used, but are computed each time the optimizing compiler recompiles a method. Correction could occur on every sample, but this approach aggregates the work and eliminates repeatedly correcting the same edges. We take the median out of 25 trials (shown as dots) to eliminate high variability of the adaptive runs. Static FDOM Correction and Dynamic Intraprocedural CF Correction add no detectable overhead. The overhead of the interprocedural correction is on average 1% and at most 6% on mpegaudio and jython. Further examination of these benchmarks show that this overhead stems from method counter instrumentation (Section 3.6). We plan to reduce overhead by eliminating method counter instrumentation from methods compiled at the highest optimization level since method counts are no longer useful for such methods.

6.0.3 Performance

We evaluate the costs and benefits of using DCG correction to drive optimization in Jikes RVM. Figure 6.4 shows steady-state performance (median of iterations 13 through 25) with several DCG correction configurations. We repeat each experiment 10 times (shown as dots) and take the median. The graphs are normalized to the execution time without correction. *Static FDOM CF Correction* shows the improvement from static FDOM correction, which is about 0.5% on average. *Dynamic Intraprocedural CF Correction* improves performance by almost 1% on average. *Dynamic Interprocedural CF Correction* uses method invocation counters (Section 3.6) to get basic block profiles with interprocedural accuracy. While the method counters provide higher accuracy, they hurt performance considerably (by 4% on average) in steady state, particularly for *compress* and *fop*. In the future, we plan to not add method counters to methods optimized at the highest optimization level, when they

are not needed anymore but hurt performance the most.

We also evaluate the performance of *Perfect DCG*, which feeds a perfect DCG to the inliner at the beginning of execution, rather than computing and correcting it on the fly as in the other configurations. The perfect DCG improves performance by only about 0.5% on average, suggesting that Jikes RVM's inliner cannot benefit significantly from high-accuracy DCGs. Previous work confirms that higher accuracy does not help Jikes RVM's inliner much but that other VMs can benefit by up to 9% from higher accuracy [3].

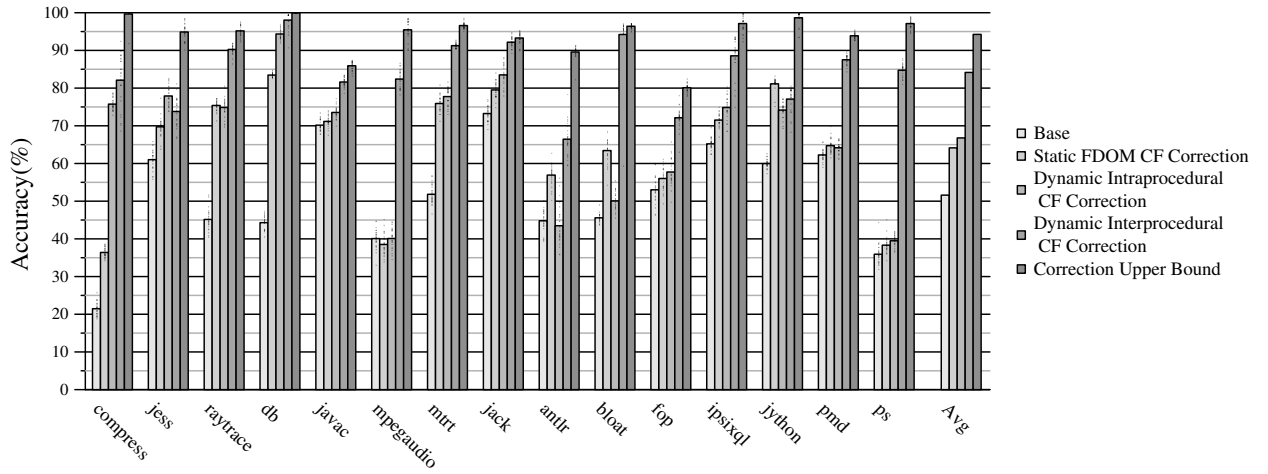


Figure 6.1: Accuracy of DCG correction on the complete DCG.

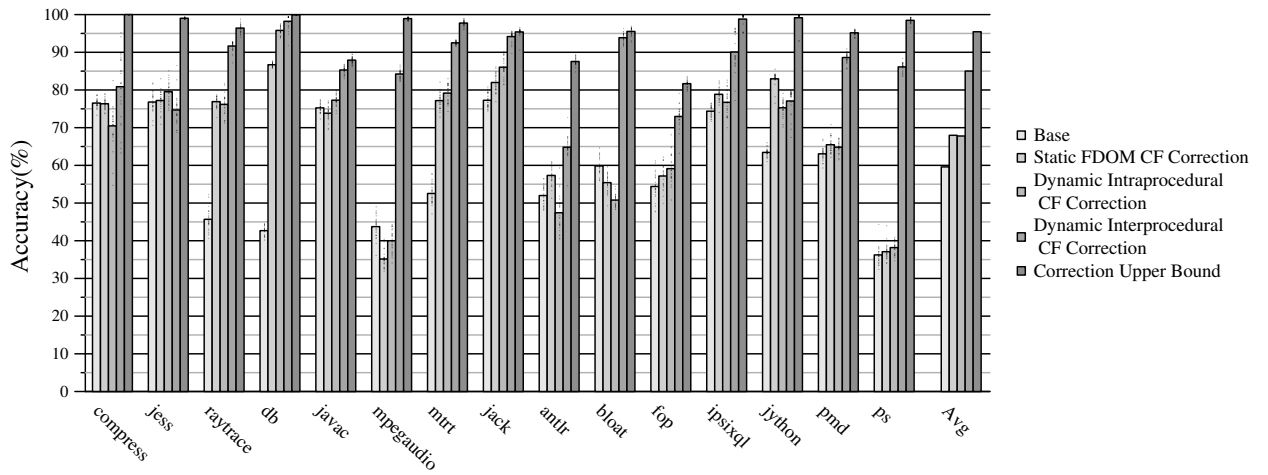


Figure 6.2: Accuracy of DCG correction on the DCG with trivially inlined call sites.

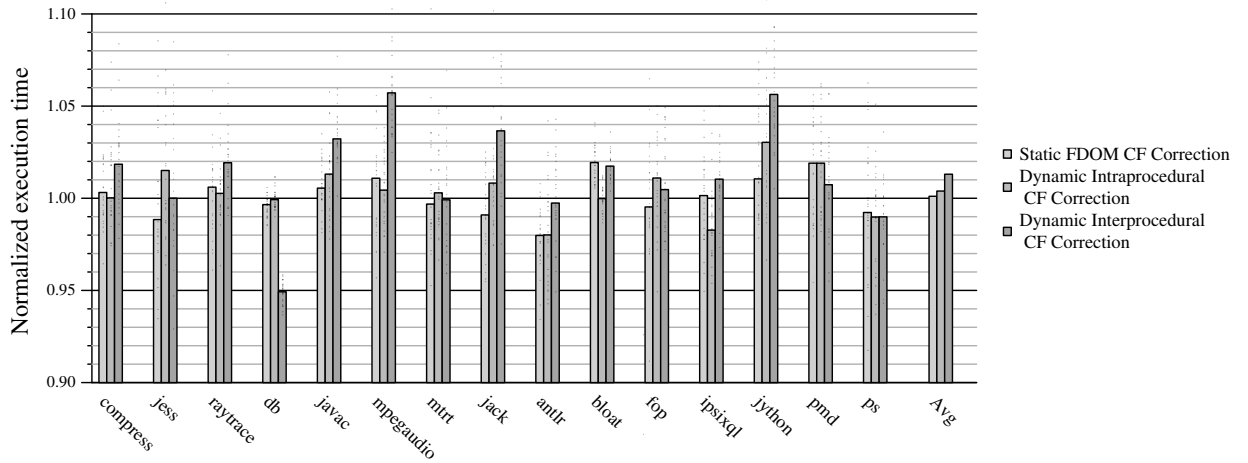


Figure 6.3: The runtime overhead of call graph correction.

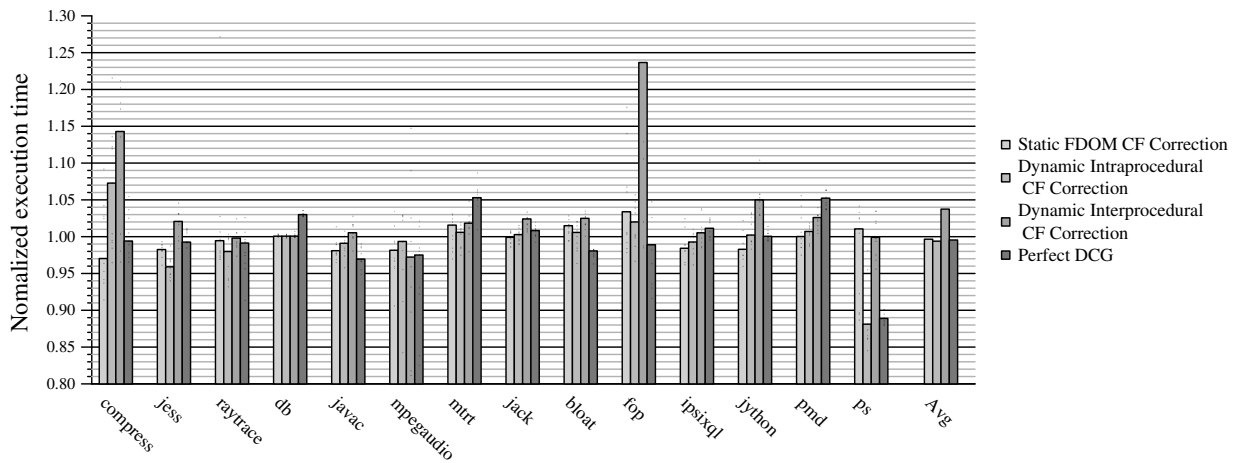


Figure 6.4: The steady state performance of correcting inlining decision using adaptive methodology.

Chapter 7

Conclusion

This paper introduces *dynamic call graph (DCG) correction*, a novel approach for increasing DCG accuracy with existing static and dynamic control-flow information. We introduce the *frequency dominator (FDOM)* relation to constrain and correct DCG frequencies, and also use intraprocedural and interprocedural basic block profiles to correct the DCG. By adding just 1% overhead on average, DCG correction increases average DCG accuracy from 52-60% to 85%. Although we obtain only a modest performance boost from using corrected DCGs to drive inlining, prior work shows other VMs benefit from higher accuracy DCGs, and we believe DCG correction will be increasingly useful in the future as object-oriented programs become more complex and more modular.

Bibliography

- [1] J. M. Anderson, L. M. Berc, J. Dean, S. Ghemawat, M. R. Henzinger, S.-T. A. Leung, R. L. Sites, M. T. Vandevoorde, C. A. Waldspurger, and W. E. Weihl. Continuous Profiling: Where Have All the Cycles Gone? In *Symposium on Operating Systems Principles*, pages 1–14, 1997.
- [2] M. Arnold, S. J. Fink, D. Grove, M. Hind, and P. Sweeney. Adaptive optimization in the Jalapeño JVM. In *ACM Conference on Object Oriented Programming, Systems, Languages, and Applications*, pages 47–65, Minneapolis, MN, October 2000.
- [3] M. Arnold and D. Grove. Collecting and exploiting high-accuracy call graph profiles in virtual machines. In *Symposium on Code Generation and Optimization*, pages 51–62, Mar. 2005.
- [4] M. Arnold and B. G. Ryder. A Framework for Reducing the Cost of Instrumented Code. In *Conference on Programming Language Design and Implementation*, pages 168–179, Snowbird, UT, 2001.
- [5] M. Arnold and P. F. Sweeney. Approximating the calling context tree via sampling. Technical Report RC 21789, IBM T.J. Watson Research Center, July 2000.
- [6] V. Bala, E. Duesterwald, and . Banerjia. Dynamo: A Transparent Dynamic Op-

- timization System. In *ACM Workshop on Dynamic and Adaptive Compilation and Optimization*, pages 1–12, Boston, MA, July 2000.
- [7] T. Ball. What’s in a region?: or computing control dependence regions in near-linear time for reducible control flow. *ACM Letters on Programming Languages and Systems*, 2(1-4):1–16, 1993.
- [8] S. M. Blackburn, K. S. McKinley, J. E. B. Moss, S. Augart, P. Cheng, A. Diwan, S. Guyer, M. Hirzel, C. Hoffman, A. Hosking, X. H. M. Jump, A. Khan, P. McGachey, D. Stefanović, and B. Wiedermann. The dacapo benchmarks. Technical report, 2006. <http://ali-www.cs.umass.edu/DaCapo/Benchmarks>.
- [9] M. D. Bond and K. S. McKinley. Continuous path and edge profiling. In *IEEE/ACM International Symposium on Microarchitecture*, pages 130–140, Barcelona, Spain, 2005.
- [10] J. Cavazos and M. F. P. O’Boyle. Automatic tuning of inlining heuristics. In *ACM/IEEE Conference on Supercomputing*, page 14, Washington, DC, 2005.
- [11] K. D. Cooper, T. J. Harvey, and K. Kennedy. A simple, fast dominance algorithm. *Software Practice & Experience*, 4:1–10, 2001.
- [12] J. F., K. J. Ottenstein, and J. D. Warren. The program dependence graph and its use in optimization. *ACM Transactions on Programming Languages and Systems*, 9(3):319–349, 1987.
- [13] N. Grcevski, A. Kielstra, K. Stoodley, M. G. Stoodley, and V. Sundaresan. Java just-in-time compiler and virtual machine improvements for server and middle-ware applications. In *Virtual Machine Research and Technology Symposium*, pages 151–162, 2004.
- [14] D. Grove, J. Dean, C. Garrett, and C. Chambers. Profile-guided receiver class

- prediction. In *ACM Conference on Object Oriented Programming, Systems, Languages, and Applications*, pages 108–123, New York, NY, USA, 1995.
- [15] A. Hashemi, D. Kaeli, and B. Calder. Procedure mapping using static call graph estimation. In *Workshop on Interaction between Compiler and Computer Architecture*, san Antonio, TX, 1997.
- [16] Intel Corporation. Intel itanium 2 processor. <http://www.intel.com/-/products//processor//itanium2/index.htm>.
- [17] R. Johnson, D. Pearson, and K. Pingali. The program structure tree: computing control regions in linear time. In *ACM Conference on Programming Language Design and Implementation*, pages 171–185, 1994.
- [18] M. C. Merten, A. R. Trick, E. M. Nystrom, R. D. Barnes, and W. mei W. Hmu. A Hardware Mechanism for Dynamic Extraction and Relayout of Program Hot Spots. In *International Symposium on Computer Architecture*, pages 59–70, 2000.
- [19] M. Paleczny, C. Vick, and C. Click. The java hotspot server compiler. In *Usenix Java Virtual Machine Research and Technology Symposium (JVM'01)*, pages 1–12, April 2001.
- [20] K. Pingali and G. Bilardi. Optimal control dependence computation and the roman chariots problem. *ACM Transactions on Programming Languages and Systems*, 19(3):462–491, 1997.
- [21] R. Prosser. Applications of boolean matrices to the analysis of flow diagrams. In *Eastern Joint Computer Conference*, pages 133–138, NY, 1959. Spartan Books.
- [22] T. Suganuma, T. Yasue, M. Kawahito, H. Komatsu, and T. Nakatani. A Dynamic Optimization Framework for a Java Just-in-Time Compiler. In *ACM*

- Conference on Object Oriented Programming, Systems, Languages, and Applications*, pages 180–195, 2001.
- [23] R. E. Tarjan. Finding dominators in directed graphs. *SIAM Journal of Computing*, 3(1):62–89, 1974.
- [24] R. E. Tarjan. Fast algorithms for solving path problems. *Journal of the ACM*, 28(3):594–614, 1981.
- [25] J. Whaley. A portable sampling-based profiler for Java virtual machines. In *ACM Conference on Java Grande*, pages 78–87, 2000.
- [26] Y. Wu and J. R. Larus. Static branch frequency and program profile analysis. In *ACM/IEEE International Symposium on Microarchitecture*, pages 1–11, 1994.