

Profile-Driven Compilation

by

Alan Dain Samples

Abstract

As the size and complexity of software continues to grow, it will be necessary for software construction systems to collect, maintain, and utilize much more information about programs than systems do now. This dissertation explores compiler utilization of profile data.

Several widely held assumptions about collecting profile data are not true. It is not true that the optimal instrumentation problem has been solved, and it is not true that counting traversals of the arcs of a program flow graph is more expensive and complex than counting executions of basic blocks. There are simple program flow graphs for which finding optimal instrumentations is possibly exponential. An algorithm is presented that computes instrumentations of a program to count arc traversals (and therefore basic block counts also). Such instrumentations impose 10% to 20% overhead on the execution of a program, often less than the overhead required for collecting basic block execution counts.

An algorithm called Greedy Sewing improves the behavior of programs on machines with instruction caches. By moving basic blocks physically closer together if they are executed close together in time, miss rates in instruction caches can be reduced up to 50%. Arc-count profile data not only allows the compiler to know which basic blocks to move closer together, it also allows those situations that will have little or no effect on the final performance of the reorganized program to be ignored. Such a low-level compiler optimization would be difficult to do without arc-count profile data.

The primary contribution of this work is the development of `TYPESETTER`, a programming system that utilizes profile data to select implementations of program abstractions. The system integrates the development, evaluation, and selection of alternative implementations of programming abstractions into a package that is transparent to the programmer. Unlike previous systems, `TYPESETTER` does not require programmers to know details of the compiler implementation. Experience indicates that the `TYPESETTER` approach to system synthesis has considerable benefit, and will continue to be a promising avenue of research.

Report Documentation Page

Form Approved
OMB No. 0704-0188

Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.

1. REPORT DATE

APR 1991

2. REPORT TYPE

3. DATES COVERED

00-00-1991 to 00-00-1991

4. TITLE AND SUBTITLE

Profile-Driven Compilation

5a. CONTRACT NUMBER

5b. GRANT NUMBER

5c. PROGRAM ELEMENT NUMBER

6. AUTHOR(S)

5d. PROJECT NUMBER

5e. TASK NUMBER

5f. WORK UNIT NUMBER

7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)

University of California at Berkeley, Department of Electrical Engineering and Computer Sciences, Berkeley, CA, 94720

8. PERFORMING ORGANIZATION REPORT NUMBER

9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)

10. SPONSOR/MONITOR'S ACRONYM(S)

11. SPONSOR/MONITOR'S REPORT NUMBER(S)

12. DISTRIBUTION/AVAILABILITY STATEMENT

Approved for public release; distribution unlimited

13. SUPPLEMENTARY NOTES

14. ABSTRACT

As the size and complexity of software continues to grow, it will be necessary for software construction systems to collect, maintain, and utilize much more information about programs than systems do now. This dissertation explores compiler utilization of profile data. Several widely held assumptions about collecting profile data are not true. It is not true that the optimal instrumentation problem has been solved, and it is not true that counting traversals of the arcs of a program flow graph is more expensive and complex than counting executions of basic blocks. There are simple program flow graphs for which finding optimal instrumentations is possibly exponential. An algorithm is presented that computes instrumentations of a program to count arc traversals (and therefore basic block counts also). Such instrumentations impose 10% to 20% overhead on the execution of a program, often less than the overhead required for collecting basic block execution counts. An algorithm called Greedy Sewing improves the behavior of programs on machines with instruction caches. By moving basic blocks physically closer together if they are executed close together in time, miss rates in instruction caches can be reduced up to 50%. Arc-count profile data not only allows the compiler to know which basic blocks to move closer together, it also allows those situations that will have little or no effect on the final performance of the reorganized program to be ignored. Such a low-level compiler optimization would be difficult to do without arc-count profile data. The primary contribution of this work is the development of TYPESETTER, a programming system that utilizes profile data to select implementations of program abstractions. The system integrates the development, evaluation, and selection of alternative implementations of programming abstractions into a package that is transparent to the programmer. Unlike previous systems, TYPESETTER does not require programmers to know details of the compiler implementation. Experience indicates that the TYPESETTER approach to system synthesis has considerable benefit, and will continue to be a promising avenue of research.

15. SUBJECT TERMS

16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT Same as Report (SAR)	18. NUMBER OF PAGES 184	19a. NAME OF RESPONSIBLE PERSON
a. REPORT unclassified	b. ABSTRACT unclassified	c. THIS PAGE unclassified			

Standard Form 298 (Rev. 8-98)
Prescribed by ANSI Std Z39-18

Acknowledgements

My heartfelt appreciation goes first and foremost to my wife Pat for being patient (almost) every time I reported that I needed just six more months, nine at the most. That she could continue to be encouraging and excited about my work is almost more than I can understand. This dissertation is dedicated to her.

Whatever I have accomplished, it is because I have stood on the shoulders of two of the most loving giants I shall ever know. My dreams could not even have been dreamt without the many accomplishments achieved by my parents. I simply cannot find the words that sufficiently express my thanks and admiration.

I am grateful to Bruce MacLennan for more than he can possibly know. I thank him most for his friendship, his intellect, and encouragement. A lot of him, Gail, and Kimmie are in these pages. Thank you!

I may never have started on this task without the inspiration and downright prodding (goading?) of Dick Hamming and his perspective on the true meaning of a Ph.D.

“[W]e greatly appreciated the valuable suggestions by ... Sue Graham, ... and, especially, Paul Hilfinger — the reviewers who provided most of the comments that kept us busy for so long producing the final draft.” [13] I am particularly pleased that I never had to make use of Paul’s dimes. I also thank Stuart Dreyfus for his work on my thesis committee.

If I acknowledged everyone the way I want to, these acknowledgements could easily become longer than the dissertation itself. Therefore, I hope a simple thanks suffices for everyone who has helped me achieve this goal. Special thanks to Dr. Rodney Farrow, the unofficial fourth member of my committee, and to Dr. Wendy Sinclair-Brown for helping me see early on what was really important; thanks to Michelle and Allen for their valuable friendship; to Dave Ungar, Richard Probst, and other members of the Stiletto Club; to Eduardo and Vicki; to Michael and Karen; to Charlie and Kendall; to Doug; and especially to Toni, Mara, and everyone in 508-20 for putting up with me and helping me during the Final Days.

This dissertation was supported in part by an AT&T Bell Laboratories Scholarship, and by the Defense Advanced Research Projects Agency (DoD), monitored by Space and Naval Warfare Systems Command under Contract N00039-88-C-0292. Their support is very gratefully acknowledged.

Contents

List of Figures	v
List of Tables	vii
1 Introduction	1
1.1 Past uses of profile data	2
1.2 Research Contributions	5
2 Profiling Techniques	6
2.1 Profiling techniques	6
2.1.1 Monitoring	7
2.1.2 Tracing	8
2.1.3 Counting	9
2.2 Efficient Counting	9
2.2.1 Basic block counts	12
2.2.2 Arc counts	12
2.2.3 The ‘optimal’ algorithms are not optimal	17
2.3 Empirical data	19
2.4 Counter-example	21
2.5 Problems with counting	24
2.6 Conclusions	25
3 A Low-level Use: Code Reorganization for Instruction Cache Performance	27
3.1 Instruction cache utilization	27
3.2 Greedy Sewing	30
3.3 Results	34
3.3.1 The Programs and Traces	34
3.3.2 Miss Rates	38
3.3.3 Performance improvement	41
3.4 Limitations	42
3.5 Conclusions	47

4	A High-level use: Implementation selection of abstract data types	48
4.1	The Problem	48
4.2	Previous work	53
4.3	TYPESETTER: The System	58
4.3.1	Formalities	58
4.3.2	The ideal system	59
4.3.3	ADTs	60
4.3.4	Iterators	66
4.3.5	Optional parameters	66
4.3.6	Alternative implementations	68
4.3.7	Feasibility functions	71
4.3.8	Evaluation functions	71
4.4	TYPESETTER: The Implementation	73
4.4.1	The Implementation Selection Algorithm	75
4.4.2	Code sharing	83
4.4.3	Refinements	84
4.5	Examples	85
4.5.1	Small example	86
4.5.2	MINOPT	93
4.5.3	Implementing THERBLIG	94
5	Conclusion	102
5.1	Problems and future work	103
5.2	Summary	106
A	The Knuth-Stevenson Algorithm	112
B	Therblig	123

List of Figures

2.1	Block counts are insufficient. The sub-graph on the right is the same as the sub-graph on the left with the addition of an arc from block C to block A. Knowing the execution counts of the blocks does not allow the derivation of the arc counts.	10
2.2	Another graph with cheaper instrumentation.	15
2.3	A program flow-graph for which MINARC is not optimal	16
2.4	The problem configuration	17
2.5	ISPLIT vs. ISPLITCHEAP	18
2.6	The problem configuration.	22
2.7	Case 1: Estimating that the most frequent arc is split expensively. . .	22
2.8	Cases 2 and 3: Estimating that the least frequent arc is split expensively, or that both arcs are split expensively.	23
2.9	Case 4: Estimating that neither arc is split expensively.	23
3.1	Removing cache contention by reorganizing	28
3.2	Improving cache utilization by reorganizing	29
3.3	Two threads from if-then-else	32
3.4	Pseudo-inlining	34
3.5	Scrunch miss rates and percent improvement	38
3.6	Troff miss rates and percent improvement	39
3.7	Cc1 miss rates and percent improvement	40
3.8	Scrunch performance for $K = 2, 3, 4, 5, 6, 7, 10, 15, 20,$ and 25	43
3.9	Troff performance for $K = 2, 3, 4, 5, 6, 7, 10, 15, 20,$ and 25	44
3.10	Cc1 performance for $K = 2, 3, 4, 5, 6, 7, 10, 15, 20,$ and 25	45
4.1	Low's algorithm	54
4.2	Possible implementations of sets (* in prototype)	61
4.3	Specification of List ADT	62
4.4	Possible implementations of lists (* in prototype)	63
4.5	Specification of Map ADT	64
4.6	Possible implementations of maps (* in prototype)	65
4.7	TYPESETTER optionals.	69
4.8	Profiling implementation of <i>add</i>	70

4.9	An alternative implementation of <i>add</i>	71
4.10	Feasibility function for implementation <i>Set_bm</i>	72
4.11	Evaluation function for <i>add</i>	72
4.12	Set union using linked lists	74
4.13	Steps to process a User program	76
4.14	Steps to process implementations and build THERBLIG	77
4.15	Main routine for choosing representations	79
4.16	Routine <i>assignable</i> for choosing the implementation of an abstract function	80
4.17	Mapping parameters onto implemented functions' signature	80
4.18	Finding compatible function implementations	81
4.19	Miscellaneous functions	82
4.20	Mapping parameters onto implemented functions' signature with equiv- alence classes	82
4.21	Declarations of generic Set functions	83
4.22	Declarations of generic Set functions	84
4.23	Declarations of generic Set functions	84
4.24	Small example	87
4.25	The call sites sorted by profiling estimates of importance	88
4.26	The actual profiling implementation for the <i>add</i> function for Sets . . .	89
4.27	Estimates of the cost of the intersection function	89
4.28	Estimates of the cost of the add function on line 54	90
4.29	TYPESETTER's assignment of types to the program	91
4.30	Results from the example with LOOPCOUNT= 10	92
4.31	The sorted call sites of Set functions from one THERBLIG run	98

List of Tables

2.1	Profiling overhead	20
2.2	Comparison of profile data size requirements	21
3.1	Summary of traces: number of instruction words fetched	37
3.2	Summary of traces: number of basic blocks	37
4.1	Small example running times with various implementation assignments for <i>cSet</i>	92
4.2	Running times for the K-S algorithm.	94
4.3	Variable assignments based on profile of three runs of THERBLIG with $p = 0$	99
4.4	Variable assignments based on the profile of three runs of THERBLIG with $p = 1$	100
4.5	Variable assignments based on the profile of three runs of THERBLIG with $p = .9$	100
4.6	THERBLIG running times with various implementation assignments	101

Chapter 1

Introduction

The ‘ideal system of the future’ will keep profiles associated with source programs, using the frequency counts in virtually all phases of a program’s life. . . . [I]f it is to be a frequently used program the high counts in its profile often suggest basic improvements that can be made. An optimizing compiler can also make very effective use of the profile, since it often suffices to do time-consuming optimization on only one-tenth or one-twentieth of a program.

[28]

In spite of the fact that Knuth made this pronouncement twenty years ago, and in spite of the fact that programmers routinely ‘optimize’ programs by hand based on profile data, Knuth’s Dictum (as we will call it) still has not been fully implemented in an automated profiling system nor shown to be undesirable.

This dissertation examines the issues surrounding the utilization of profile data in the compilation of source code. This is a larger subject than that of generating or collecting profile data. It requires asking, at a minimum, the following questions: *Given that profile data exists for a program, how might a compiler make use of that data to produce better executable code? What kinds of profile data can be generated/collected? What kinds of profile data are useful? How expensive is this profiling?*

We start with two hypotheses:

1. Collecting profile data need not be prohibitively expensive.
2. Compilers can profitably use profile data at all levels of the compilation process.

Compilers that emit profiling code have been at least partially implemented by most modern systems. But no systems of which I am aware utilize a program’s profile data throughout the compilation process. Furthermore, while these hypotheses might be accepted in a general way, there are still some misconceptions about the cost of profiling, and room for improvement in the profiling algorithms themselves.

1.1 Past uses of profile data

The collection and use of profile data has a long history, beginning with Knuth's 1971 paper [28]. In this paper, the term 'profile' was first used, and defined to be the collection of execution frequency counts taken during executions of a program. Since that time, the term 'profile data' has come to mean any quantitative information gathered about the run-time behavior of a program, including execution counts of the program and its sub-parts, reference counts of the program's data objects, and real-time measures of algorithm executions.

Knuth's examination of execution profiles of running user programs uncovered two facts: (1) Most programmers do not know where their programs spend most of the time; and (2) even when programmers analyze their programs, they still don't know where their programs spend most of the time due to the fact that programmers almost never have access to sufficient information about system and library functions to deduce the runtime resources they consume. For example, a major culprit in the FORTRAN environment in which Knuth did his study was the formatting routines in the I/O statements.

Another result of Knuth's study was the rule-of-thumb that said that 90% of a program's time is spent in 10% of the code, variously called the 90-10, or 80-20, rule. Knuth never used either of these numbers but reported that in his studies 50% of the time was spent in 4% of the code. In fact, the actual numbers are different for each program. The 90-10 rule, or whatever you want to call it, is one of the guiding principles on which manual program optimization is based: find that section of your program that takes most of the runtime resources, and either modify the algorithm itself (e.g., change a bubble-sort into a quick-sort) or make the existing algorithm more efficient at the low level (e.g., hoist common expressions out of loops, do strength reduction on the index variables, turn repeated array indexing operations into pointer operations, etc.).

After Knuth's 1971 paper, Dan Ingalls published two papers describing descendants of the FORDAP profile tool used by Knuth. FORDAP was a basic-block counting profiler. The first technical report [23] gives details of how the FORTRAN Execution Time Estimator (FETE) adds execution time estimates to the frequency count displays of FORDAP. This enhancement was prompted by the obvious fact that not all statements are created equal. For example, the FORTRAN statement

$$A = B(I)$$

will execute at vastly different speeds depending on whether B is an array or a function. FETE used a value based on 'weights' assigned to expression operators and statement classes to give a rough estimate of the execution time. FETE was not sophisticated enough to handle calls on user functions. If, in the example above, B is not an array but a call on a function, FETE will count it as an array reference unless B is a standard FORTRAN function.

Ingalls' second paper [22] describes FORTUNE, which is simply a renamed, product version of FETE. FORTUNE and FETE modified the source program so that it contained FORTRAN statements incrementing elements of an array of counters. These counting statements were placed essentially at the beginnings of basic blocks. The analyzer reported statement execution counts and estimates of execution time for each statement.

Prof [5], a profile collector for C, Pascal, and FORTRAN programs on the UNIX system is an example of profiling tools in current use. *Prof* samples the program counter via timer interrupts to estimate the amount of time spent between the symbols of the program. *Prof*'s usability has been enhanced by *Gprof*, a program developed at Berkeley by Graham, Kessler and McKusick [16]. *Gprof* explicitly concentrates on procedure calls, providing both frequencies for calls through counting and estimates of the time spent in each procedure. Timing estimates are derived by sampling the program counter, as in *prof*. After the user's program has run, a separately invoked post-pass analyzer distributes timing estimates to the program's procedures based on the static and dynamic call graphs.

There has been a small amount of research published about the best way to profile a program. In the first volume of his *Art of Programming* series, Knuth gives the algorithm for determining a minimal instrumentation of a program for collecting the execution counts of arcs. Knuth and Stevenson [30] (about which more will be said later) published the definitive algorithm for finding a minimal instrumentation of a program that counts the execution of its basic blocks. Cheung [7] concurrently developed algorithms for finding minimal instrumentations that count the frequency of execution paths through a program. A paper by Sarkar [42], without referencing this body of work, developed an algorithm for instrumenting a program based on its dependence graph.

There has been some research into the potential uses of profile data. Gilbert Hansen's research [17] is an early investigation into behavior-driven optimization. He hypothesized that, for certain classes of software, the optimization of a program could be done at *run*-time more economically than at compile-time. Instead of the usual compile-a-file paradigm that most compiler systems utilize, Hansen's "adaptive" compiler consisted of two "phases". The first phase generated an interpretable form of a FORTRAN program in a fast, one-pass compilation (it produced 'quads' as the interpretable form). The second part consisted of an interpreter and optimizer loaded with the compiled program. When the interpreter detected that a basic block was being executed sufficiently often, interpretation was suspended while the optimizer was invoked to compile the basic block to a lower level. If a basic block were executed often enough, it would eventually be compiled down to machine language.

The one-pass compiler annotated each basic block with information about its size and complexity so the interpreter could predict profitable optimizations. The system was designed to expend effort only on optimizations that had a high probability of paying for themselves through improved execution of the program. There

were four levels of optimizations, three of which were performed on the interpretable form (constant folding, common subexpression elimination, and moving invariants out of loops), and the final of which compiled the quads resulting from the first three optimizations into machine language, an ‘optimization’ he called ‘fusion’. Therefore, the optimizer would be invoked up to four times on a basic block, each time optimizing it further. When the basic block had been optimized as much as possible in its interpreted form, the final optimization would compile it to machine language. At this time, several machine-dependent optimizations would also be applied and execution would become ‘threaded’; i.e. a mixture of interpretation and direct execution in which the program has been ‘fused’ into the interpreter.

It is not surprising that Hansen’s one-pass compiler executes more quickly than optimizing, multi-pass compilers. What is surprising is that those initial savings were almost never depleted. That is, the time required to do the one-pass compilation to ‘quads’, plus the time for interpretation and intermittent optimization was almost always less than the time required to do a full optimization of the original program plus the execution time of the optimized program!

Hansen’s system is appropriate for compiling and running throw-away, run-once programs; e.g. in a student programming environment, the overall CPU utilization is decreased. It is not an appropriate system for constructing software designed to be run many times (e.g. editors, the operating system, the adaptive compiler itself). The success of this system depended on the ratio of the number of compilations to the number of executions being very close to one. If a program is executed many times, it then becomes profitable to compile and optimize the whole program once.

Hansen’s research lends credence to our contention that profile-driven optimization is a useful adjunct to ‘traditional’ compilation. If his adaptive compiler could use heuristics to predict the future behavior of a program successfully enough, a static compiler using those same heuristics with complete profile data should do better. The major point is that, whereas the adaptive compiler system is forced to make the assumption that the performance of the program in the immediate past is predictive of its performance in the immediate future, a static compiler can accumulate profile data, smooth out anomalous behavior over several runs of the program, and make the more accurate assumption that the average past performance of the program is a good predictor of its average future behavior.

There has been a large amount of research using profile data to improve virtual memory performance. Most of this work has depended on profile data in the form of an address trace and has improved program performance by reorganizing the modules of programs to minimize page faults (Ferrari [12] gives a summary of this early work). Nearly all of this research has concentrated on post-compilation module-level reorganization of a program. There have been some techniques developed to reorganize programs dynamically based on their behavior. For example, K. D. Ryder [39] and J.-L. Baer and G. Sager [2] used dynamically collected profile data to

allocate physical memory for programs running on virtual systems.

In spite of all the use of profile data to improve program performance in various ways, modern programmers wishing to improve the performance of their programs using profile data have limited options. In almost every system that provides any profiling capability at all, it is still up to the users to analyze the data and manually reorganize or rewrite their programs based on that data.

1.2 Research Contributions

This dissertation presents several results related to the use of profile data, ranging from exactly how profile data should be collected, to actual uses of profile data in languages and their compilers.

Since at least the mid-seventies the problem of efficient profile collection via code instrumentation has been considered a solved problem. This research discovered problems with the solutions, and results are presented here that show that old assumptions about profiling are incorrect. Specifically, it is not the case that counting the execution frequencies of transfers of control in a program is expensive: it is often cheaper than simply counting execution frequencies of the basic blocks in a program. Also, I show that the algorithms that have heretofore been considered ‘optimal’ are not only *not* optimal, but that optimality is difficult to achieve. In Chapter 2, I present an algorithm that finds better instrumentations of programs.

A difficult problem from past research has been improving the performance of a program in a paging environment. While paging is no longer the issue that it once was, caching in a memory hierarchy has taken its place. In Chapter 3 I demonstrate that a dramatic percentage of the improvement in the instruction-cache behavior of a program can be obtained by reorganizing a small percentage of the actual code.

Users should be able to declare a variable to be of some abstract type without worrying about the implementation. Unfortunately, it has turned out to be very difficult to design an efficient system with this capability. In Chapter 4 I present the design of a system which uses profile data to assign implementations to variables. With appropriate language extensions that allow the writer of alternative implementations to specify what kinds of profile data are needed and how it is to be evaluated, the user can declare variables to be of a generic type (e.g. `Set(int)`) and let the system decide, based on the profile data, which implementation to use for the variable.

Chapter 2

Profiling Techniques

How do I count thee?
Let me love the ways ...
(*apologies to Ms. Browning*)

Programmers' intuitions about the runtime behavior of their programs are notoriously bad. Profiling counteracts this deficiency by providing objective measures.

After a brief discussion of various profiling techniques, we will focus on the insertion of counting code in programs as the technique of choice. We will see that the traditional solutions to the minimal instrumentation problem are not optimal. An algorithm that is more nearly optimal is presented, with evidence to show that true optimality may be quite difficult to achieve.

2.1 Profiling techniques

When designing a profile data collection system, two questions must be answered: how, and how much? This work explores the use of software instrumentation to collect profile data. This is in contrast to, say, using specialized hardware to monitor a system from 'outside'. Such specialized hardware has been built, but computer manufacturers that provide it do not provide high-level language programmers with useful tools that can take advantage of such equipment. Tools such as in-circuit emulators or digital oscilloscopes are primarily useful to the electrical engineer or, at best, to the writer of peripheral device drivers. There is no evidence to suggest that such specialized hardware can provide better data than software instrumentation can provide for high-level language applications.

There are three ways profile data can be collected with software instrumentation: monitoring, tracing, and counting. All three methods are discussed in detail below. Monitoring requires some hardware support, usually in the form of a countdown-timer interrupt. Tracing refers to recording in memory or on some external media the sequence of relevant operations of the program or system. Counting is

implemented with code inserted into the program to increment (an array of) counters to record the execution frequencies of the program.

The answer to the ‘how much’ question depends on the granularity of the profile data for any particular program. There are basically four granularities in use: procedure, basic block, statement, and instruction level granularities. There are many variations on these themes, but this gives us enough framework to discuss several profiling techniques.

2.1.1 Monitoring

A histogram of program execution is generated by observing the value of the PC (program counter) at frequent intervals (hence the alternative name of ‘PC sampling’). There are three parameters for programmers to specify for monitoring: the area (i.e. address range) of the program that is to be monitored, the number of data points to be generated for that address range (the granularity), and the sampling frequency.

For example, a programmer may specify that addresses in the range 0x2000 and 0x100000 are to be monitored, and that 65,536 data points are to be generated. So $0x100000 - 0x2000 = 1,040,384$ addresses are divided into 65,536 regions, meaning that one data point will represent 15 addresses: we say that the measurement granularity is 15. On most machines, more than one instruction can fit in a range of 15 addresses, particularly when the machine is byte-addressable. Therefore, information about multiple basic blocks totally contained in a single 15-address range, or about basic blocks that straddle the boundaries of multiple blocks, will be fuzzy, at best. The goal is to choose a granularity that does not generate too much information¹ and yet captures sufficient information about the program to make reasonable deductions about its performance.

Choosing a good sample rate is also fraught with tradeoffs. If the interrupts occur too infrequently, too much of the program’s behavior will occur between the interrupts. If interrupts occur too frequently, the program’s execution will be totally swamped by the interrupt overhead. Finding the proper value is hit-or-miss, and there are no published statistical studies showing what range of interrupt frequencies sufficiently capture the behavior of a program.

One problem shared by almost all profiling techniques is that of measuring system overhead. There is no way that a program in a multiprogramming environment can reasonably measure the behavior of the operating system activity due to the executing program. The best that can be hoped for is some measure of the frequency of the system function calls during the program’s execution.

All monitoring implementations depend on the sequentiality of instruction execution to extrapolate the statistical samples into information about the program’s

¹Make the granularity one, and the profile data tables will be at least as large as the program and possibly four times larger than the program.

execution. Since there is no single register that can be sampled to derive a profile of a program's data reference behavior, it is extremely difficult to derive measures of a program's use of data with monitoring. Even continuous monitoring of, say, a data bus (which would certainly demand hardware support) cannot provide very interesting information. For example, tracing the data references in the area of memory devoted to the execution stack provides little information. Since the stack varies dynamically as the execution of the program proceeds, functions and their associated data are not guaranteed to map into the same addresses, nor is there any guarantee that a given sequence of addresses will always be used by a given function. To understand the memory reference behavior of a program's stack memory requires some knowledge of the program's dynamic call tree, something a simple memory monitor cannot provide.

2.1.2 Tracing

Very often, simply knowing how many times a function was called, or having an estimate of how much time each function consumed, is not sufficient. Questions such as "how often was event X followed by event Y" turn out to be important in some contexts. In studies of a program's performance in a paging system, the question takes the form "What is the sequence of page references?", and similarly in cache performance studies the question is "What is the sequence of cache line references?". Both kinds of studies have traditionally used address traces (both instruction and data) of actual program executions to answer these questions.

There are three ways to gather traces (without special hardware assistance):

1. *Instrument the code.* This is a messy and laborious technique, particularly if the instrumentation must not interfere with the generated addresses. The trace can be considered legitimate for most purposes only if the recorded addresses are the same as they would be if there were no instrumentation code.
2. *Simulation.* Sometimes it is simpler to write a simulator for the machine in question. Each instruction simulated can then produce a trace of the instruction's address, all data references generated by the instruction, and perhaps other information such as timing estimates.
3. *Single-stepping.* Some microprocessor architectures have a single-step feature that interrupts a process after each instruction executed. A separate process (in a separate address space) handles the interrupt and generates the trace of address references.

Tracing a program usually is quite expensive, causing instrumented programs to run anywhere from 2 (if you're lucky) to 10 times slower, depending on the actual method of tracing and the number of execution features being traced. Furthermore, tracing can produce tremendous amounts of data. With memories and

programs getting larger, it may take many millions of instructions of trace data to capture interesting effects. Borg et al. report that some interesting characteristics of traces were not apparent until several billion instructions had executed [6]. Compaction techniques can alleviate some of the problem as I have demonstrated [41], but the management and assimilation of huge amounts of data is always difficult.

Again, even if these problems are solved, the trace of a program on a virtual memory system is not sufficient for capturing effects introduced by the operating system. Tracing the operating system alone may not provide the necessary information either. Tracing a complete system across system calls and interrupts is a gargantuan task, and produces a prodigious amount of information in a small amount of time. Agarwal, Sites, and Horowitz instrumented a processor's micro-code to collect such system-level traces [1], but it is difficult for programmers to utilize this technique to gain an understanding of how their programs and the system interact.

2.1.3 Counting

Many, if not most, compilers now have the ability to insert counting code into users' programs and provide very precise information about the number of times lines, basic blocks, or functions are called. The actual instrumentation is quite simple, even for separately compiled units of a program. Since this is my preferred method of collecting profile data, I discuss it in more detail in the next section.

Counting executions of functions is probably the least useful granularity. Programmers learn which are the most frequently called functions, but that may bear very little relation to the location of the most frequently executed inner loop. Counting lines is not always sufficient since "lines of code" are artifacts of particular languages and the styles of individual programmers. For example, some programs written in C can have many basic blocks concealed on a single line of code due to some programmers' proclivity for using C's conditional expression construct in deeply nested macros. The net result is insufficient information when more than one basic block is on a line, and superfluous counts when basic blocks span lines.

Any modern compiler should implement some form of counting, and at a minimum it should include basic block counting. In the next section, however, I argue that counting execution frequencies of transfers of control (*arc* counting) is best.

2.2 Efficient Counting

Most compilers that implement profiling via the insertion of code count lines or, at best, basic blocks. However, there are some applications where basic block counts are insufficient. Examples include code reorganization to improve the performance of multi-level memory hierarchies [21,35], jump optimization, and code

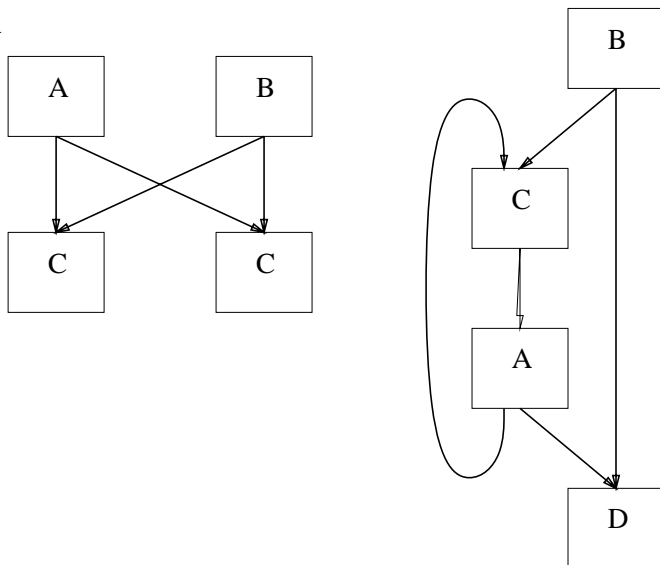


Figure 2.1: Block counts are insufficient. The sub-graph on the right is the same as the sub-graph on the left with the addition of an arc from block C to block A. Knowing the execution counts of the blocks does not allow the derivation of the arc counts.

generation and register assignment [25]. Such optimizations and code transformations depend on knowing branch probabilities (i.e. arc frequencies) and can make only imprecise use of block counts. Arc frequencies frequently cannot be deduced from block frequencies; an example program graph is in Figure 2.1, where knowing that each block is executed n times does not provide enough information to determine the number of times, say, execution of block A is followed by execution of block D . This is not a contrived example. Also shown in Figure 2.1 is an actual program flow-graph (PFG) that contains the problematic graph. This flow graph was generated frequently by a Pascal compiler for **while** loops.

Therefore, arc frequencies, not just block frequencies, are desired, although historically block frequencies have been considered more desirable. To elaborate further on the problem of minimal instrumentation for arc counts, we will need some definitions. We assume that the execution cost of inserted instrumentation code is constant and non-zero; call this cost K_I . In Section 2.5, we will discuss some subtleties in instrumentation code, but for the moment we will assume that each instance of instrumentation code is exactly the same (e.g., a memory-to-memory increment operation, or an equivalent sequence of operations). When instrumentation is inserted in code, it may be necessary to insert a jump instruction to maintain the semantics of the code. We will assume that the execution cost of this jump instruction is also constant, $K_J > 0$.

A program flowgraph is

- a set V of basic blocks, the vertices of the graph;
- a set E of directed edges, which are pairs of vertices, $e = (src(e), snk(e))$. We say that the edge e leaves $src(e)$ and enters $snk(e)$, the *source* and *sink* of the edge, respectively. The edge is an exit arc of $src(e)$ and an entrance arc of $snk(e)$.
- a distinguished edge of the graph, e_0 ; we assume, without loss of generality, that $snk(e_0)$, the *entrance block* of the graph, has no other predecessors, and $src(e_0)$, the *exit block* of the graph has no other successors; each flowgraph has exactly one entrance block and exactly one exit block.

For each edge e of the graph, $F(e)$ is the frequency count of e ; $J(e)$ is a boolean function that is *true* if e is an out-of-line *jump* arc, and *false* if e is a *fall-through* arc; C is a function that maps edges into instrumentation costs. $C(e)$ is the cost required to instrument edge e , and depends on $F(e)$, K_I , and K_J ; specifically, $C(e) = C_I(e) = F(e)K_I$ if it is not necessary to insert a jump instruction, or $C(e) = C_J(e) = F(e)(K_I + K_J)$ if a jump instruction is required.

Each block v is the sink of *at most* one fall-through arc, and the source of *at most* one fall-through arc. There is no limit on the number of arcs for which a block is a sink or a source. We define $in(v)$ to be the set of predecessors of v , and $out(v)$ to be the set of successors of v . We say that an arc e is *crowded at its sink* if $|in(snk(e))| > 1$. Likewise, it is *crowded at its source* if $|out(src(e))| > 1$. If an arc is crowded both at its source and at its sink, we simply say that it is *crowded*. We define the predicates $e.crowdedSnk$, $e.crowdedSrc$, and $e.crowded$ on the edge e for these conditions.

The remainder of this discussion assumes that $K_J > 0$ and $K_I > 0$; all of the examples we will display assume $K_J = K_I = 1$.

All of the following algorithms take advantage of the fact that, given a program flow graph, all execution frequencies of the arcs and nodes can be derived if we know an appropriate $|E| - |N| + 2$ arc frequencies. Such a subset of arcs can be selected by finding those arcs that form the complement of a (non-directed) spanning tree in the flow graph. If the arcs have associated costs, then a minimal (maximal) cost subset of arcs can be found by taking those arcs that form the complement of a maximal (minimal) spanning tree in the program flow graph; a proof can be found in Knuth, Vol 1., page 368 [29]. I refer to the algorithm for finding a spanning tree as the SPAN algorithm, and to the algorithm for finding the minimal (maximal) spanning tree as MINSPAN (MAXSPAN). All of this, with the exception of the arc characterization function J and slight differences in notation, is consistent with previous work.

2.2.1 Basic block counts

The most commonly implemented instrumentation technique counts the number of executions of every basic block, which we'll refer to as FULLNODE. Techniques that count every line or every programming language statement are even more inefficient variants of this technique. FULLNODE has the advantage of being the easiest to implement, but the disadvantage of being inefficient: a program can easily be slowed down by as much as 50% to 100%, depending on the execution cost of the instrumentation and the average size of a program's basic blocks.

However, it is not necessary to instrument each and every basic block to get complete block counts. Knuth and Stevenson [30] and Cheung [7] present algorithms that compute a minimal subset of basic blocks that when instrumented provide sufficient data to recover the execution frequencies of all other basic blocks. (Cheung also contains a much more detailed study of minimal instrumentation, including minimal instrumentation for determining path coverage.) A TYPESETTER version of Knuth and Stevenson's algorithm is included in Appendix A. I will refer to this algorithm as the K-S instrumentation algorithm.

Conceptually, the K-S algorithm is a graph transformation followed by an application of a spanning tree algorithm. Given a program-flow graph with basic blocks V and edges E , the relation \equiv between basic blocks is defined to be the smallest equivalence relation such that $a \equiv b$ if there exists vertex c and arcs $c \rightarrow a$ and $c \rightarrow b$. A *reduced* graph is produced whose vertices V_r are the equivalence classes of the original graph, and whose edges E_r correspond one-to-one with the basic blocks of the original graph; for each basic block $b \in V$, there exists an edge in V_r from the equivalence class containing b to the class containing the successors of b (by construction, they are all in the same class).

SPAN is applied to the reduced graph to find a spanning tree, and those edges of E_r not in the spanning tree specify the nodes in V that need to be instrumented. The K-S algorithm also computes the expressions that will later allow the frequencies of all nodes to be computed in one pass.

If we have some notion of the execution behavior of the program, then each node v of a program flow graph can be assigned an instrumentation cost, $K_I F(v)$. A minimal cost instrumentation of the program flow-graph can be found by following the steps for the K-S algorithm, but using MAXSPAN to find the spanning tree rather than SPAN. We'll call the minimum cost algorithm MINNODE.

2.2.2 Arc counts

A naive implementation of arc counting in a program flow graph would be to instrument each and every arc of the original graph. This would provide the necessary counts, but rather expensively. The cost comes from two facts: (1) $|E| > |V|$, meaning more space would be required by the instrumentation code; and, (2) to instrument some arcs requires the creation of a new basic block, and the

addition of a jump instruction. For such arcs, the instrumentation cost would be $K_I + K_J$, whereas the instrumentation cost for all nodes is simply K_I . These facts have contributed to a commonly held belief that arc counting is too expensive and complicated for practicality. As we shall see, this simply is not the case.

As in the case of basic block counts, it is not necessary to measure each and every arc to derive the number of times each was traversed. A minimum set of arcs to be measured in a graph consists of those arcs not in the spanning tree constructed by SPAN, and the minimal *cost* instrumentation is the set of arcs that form the complement of the spanning tree found by MAXSPAN. We will call the resulting algorithm MINARC.

As it stands, MINARC is too expensive. As noted above, instrumenting an arc requires the creation of a new basic block that is inserted in the arc between its source and sink nodes. For jump arcs, this new basic block must also contain a jump to the original target of the arc; this jump adds to the cost of the instrumentation. This cost is often ameliorated by using transformations to turn arc measurements into node measurements wherever possible. For instance, if an edge that is to be instrumented represents the fall-through of one basic block into another, and the edge is the only edge leaving the source block, then the instrumentation can be inserted in the edge's source block; similarly, if the edge is the only edge entering the target block, then the instrumentation can be inserted in the edge's target block.

Procedure 1 implements a heuristic algorithm using these transformations to reduce the cost of profiling arc traversals, where

ISINK(e)=Add instrumentation code to the front of the block $snk(e)$.

ISOURCE(e)=Add instrumentation code to the front of $src(e)$.

and

ISPLIT(e)=Replace edge e with a new basic block v_e and edges $e_1 = (src(e), v_e)$ and $e_2 = (v_e, snk(e))$. Furthermore, $T(e_1) = T(e_2) = T(e)$, and $F(e_1) = F(e_2) = F(e)$. (As we will see, determining $C(e_1)$ and $C(e_2)$ is problematic.)

These transformations allow us to keep the problem relatively simple. If we added transformations that allowed us, say, to move basic blocks to remove jump instructions, or to make the most frequent arc out of a block the fall-through arc from that block, the problem would be much more complicated. To keep the problem simple, we assume that the linear order of basic blocks in memory will remain the same throughout the process of instrumenting the program. Our only options are to determine *where* to add instrumentation code.

Before the MAXSPAN algorithm can be applied to the program flow graph to find which arcs are to be instrumented, the cost of instrumenting each arc must be estimated. If at all possible, we would like to insert instrumentation code without introducing extra control flow logic or extra basic blocks, so if we can identify where the above transformations can be applied, we can more accurately estimate the cost of instrumenting an arc. This leads to the following heuristic cost estimation algorithm.

Procedure 1 *Instrumenting an arc:**Input:* An arc that is to be instrumented.*Result:* Instrumentation code is added to the PFG to count executions of the arc.*Method:*

```

Instrument(Arc e)
{
  if not e.crowdedSrc then ISOURCE(e);
  elsif not e.crowdedSnk then ISINK(e);
  else ISPLIT(e);
  endif;
}

```

□

Function 2 *Estimate cost of instrumenting an arc:**Input:* An arc in the PFG.*Output:* An estimate of the cost of instrumenting the arc should it be chosen for instrumentation.*Method:*

```

InstrumentationCost(Arc e)
{
  if not e.crowded then
    return  $K_I \times F(e)$ ;
  elsif not  $J(e)$  then
    return  $K_I \times F(e)$ ;
  else return  $(K_I + K_J) \times F(e)$ ;
  endif;
}

```

□

If the instrumentation code can be appended to a basic block, then the cost of instrumenting the arc e is the cost of the instrumentation code itself, K_I , times the frequency of execution $F(e)$. If the arc is a fall-through arc then the cost is still $K_I F(e)$ even if the arc is crowded: the instrumentation can be inserted between the two blocks. In all other cases, the source block of the arc will be jumping to the instrumentation code, which will itself have to jump to the sink block; therefore, the cost of instrumenting the arc is $(K_I + K_J)F(e)$, where K_J is the cost of the

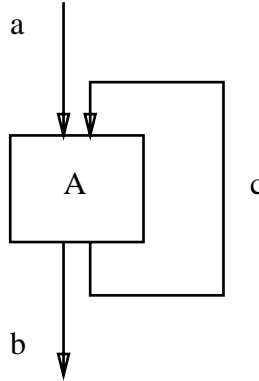


Figure 2.2: Another graph with cheaper instrumentation.

jump instruction. We will call MINARC augmented with the heuristic placement algorithm MINARC'.

This approach looks good, and I will present results below to show that it is effective, but there is nothing to suggest that it is complete or produces minimal instrumentations. Rather, it is a set of *ad hoc* rules for utilizing the results of the MINARC algorithm. That it is not complete can be seen from Figure 2.2². If arc c is to be instrumented, and if the execution frequencies of a or b can be derived independently of c , then the instrumentation on c can be moved into the node A . The frequency of arc c is then $F(c) = F(A) - F(a)$, since $F(a) = F(b)$.

This suggests another way of asking the question. We have an algorithm that will find the minimum set of nodes for computing node counts, and an algorithm for finding a minimum set of arcs for computing arc (and therefore node) counts: is there an algorithm that will find a minimum set of nodes *and* arcs for computing execution frequencies for a program flow graph?

An extension to the above algorithms produces a candidate algorithm, which I'll (presumptuously and inaccurately) call OPT. Given a flow graph (V, E) , we construct a new one (V', E') as follows. For each node in V we create a corresponding node $v' \in V'$. For each arc $e = v_1 \rightarrow v_2 \in E$, we construct a basic block $v'_e \in V'$ and two arcs $v'_1 \rightarrow v'_e$ and $v'_e \rightarrow v'_2$ in E' . So $|V'| = |V| + |E|$ and $|E'| = 2|E|$.

The K-S algorithm applied to the new graph (V', E') will yield a minimum set of nodes in V' required to compute all frequencies of all nodes in V' . Since each node selected for instrumentation in V' corresponds to either a node or an arc of the original graph, we also have a minimum set of nodes and arcs of (V, E) from whose measurement we can compute all other execution frequencies of nodes and arcs in F . We extend the algorithm to find the minimal cost set of nodes and arcs by assigning the same costs to the nodes in (V', E') as are estimated for the nodes and arcs in (V, E) from which they are constructed. Therefore, if a node $v' \in V'$ corresponds to

²Thanks to Jim Wilson of Cygnus Corporation for pointing this example out and for taking the time to convince me it was worth a second look.

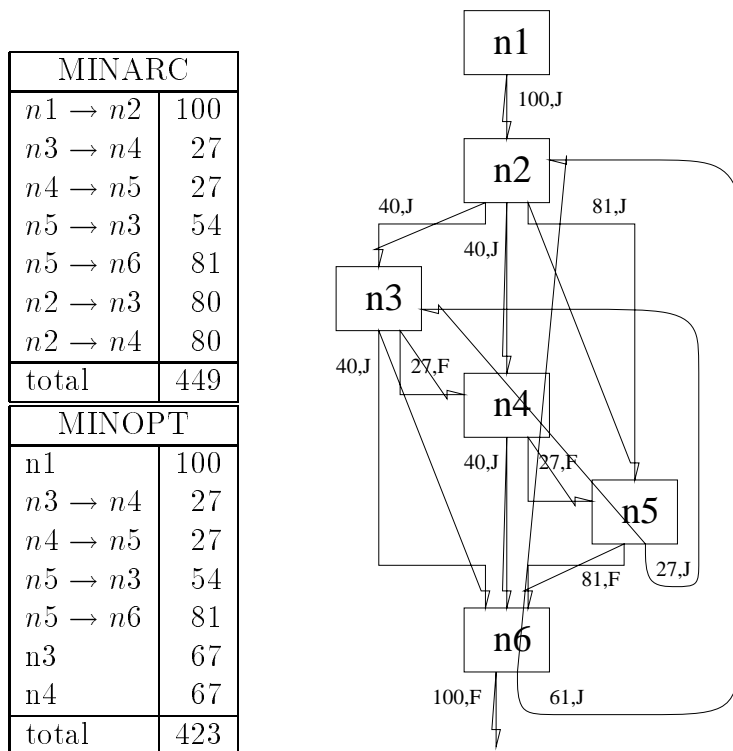


Figure 2.3: A program flow-graph for which MINARC is not optimal

node $v \in V$, then $C(v') = C(v)$ (which is always K_I). If a node $v' \in V'$ corresponds to the arc $e \in E$, then $C(v') = C(e)$. By assigning these costs and invoking a maximal spanning tree algorithm, we find the minimal cost instrumentation using nodes and arcs.

That this algorithm, which we'll call MINOPT, is not equivalent to MINARC', and can sometimes improve on the measurement costs of a program graph is easily proved. In Figure 2.2 MINOPT always instruments node A instead of arc c , unless of course c is executed a lot less than once per execution of A .

Finding PFG instances on which MINOPT produces *different* instrumentations than MINARC' is a bit more difficult. Figure 2.3 is the second simplest subgraph I have found that demonstrates this difference (the simplest is Figure 2.2). In Figure 2.3 each arc is labeled with an execution frequency, and with whether it is a jump arc or a fall-through arc. The two tables show the results of the MINARC' and MINOPT algorithms. The first column contains the objects chosen to be measured, and the second column contains the cost of measuring that object (assuming $K_I = K_J = 1$). The major difference in the results of the two algorithms is that MINOPT has chosen two nodes to measure, while MINARC' can choose only arcs, and then look for transformations to decrease the cost. The instrumentation transformations described previously do not help MINARC' in this example. The only one that applies is ISOURCE($n1 \rightarrow n2$).

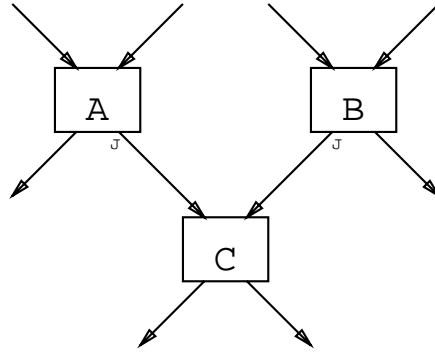


Figure 2.4: The problem configuration

While MINARC' does almost as well as MINOPT , it is heuristic and not minimal. MINOPT is provably minimal (with respect to a set of instrumentation cost estimates), and can find instrumentations that would be difficult to characterize easily as post-transformations for MINARC . Given that the complexity of both MINARC and MINOPT is $O(|E| \log |E|)$ MINOPT is to be preferred for its simplicity and minimality.

2.2.3 The ‘optimal’ algorithms are not optimal

The assignment of minimal instrumentation costs to the edges of a program flowgraph has been glossed over in the literature. In fact, such an assignment cannot always be done unambiguously so as to guarantee a minimal solution. That is to say, all optimal solutions shown in the literature are optimal with respect to a specific instrumentation cost assignment on the arcs. But until this time no one has examined the question of how those costs are assigned, or even if they can be assigned, and whether such an assignment still permits an efficient optimal solution. For instance, Cheung mentions that instrumenting some arcs requires extra flow-control instructions, [7, pp. 38-39], but his algorithms assume that these costs can be assigned in linear time, and that the cost of instrumenting one arc does not affect the cost of instrumenting other arcs.

That these assumptions do not hold for even very simple cases can be seen in Figure 2.4 where the arcs to node C are both jump arcs and are both crowded. The cost assignment algorithm described above would assign instrumentation costs of $F(A \rightarrow C)(K_I + K_J)$ and $F(B \rightarrow C)(K_I + K_J)$ to the arcs. However, if either or both arcs are chosen for instrumentation, it is obvious that one of them does not have to jump to block C , but rather can create a new basic block that simply falls through to C . In Figure 2.5 we assume that both arcs going to basic block C are to be instrumented. The arc $B \rightarrow C$ is instrumented by placing its instrumentation code in a separate basic block which falls through to the basic block C . The instrumentation for the arc $A \rightarrow C$ cannot be so configured and hence must use the more

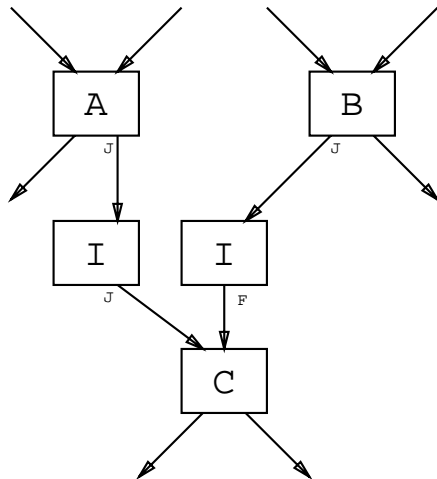


Figure 2.5: ISPLIT vs. ISPLITCHEAP

expensive method for splitting the arc. We can add this transformation to our list of transformation heuristics on page 13:

ISPLITCHEAP(e): Replace edge e with a new basic block v_e and edges $e_1 = (src(e), v_e)$ and $e_2 = (v_e, snk(e))$; $J(e_1) = J(e)$ and $J(e_2) = false$.

If in order to instrument edge e we must use ISPLIT(e), then the instrumentation cost assigned to edge e must be $C_J(e) = F(e)(K_I + K_J)$. If we use ISPLITCHEAP(e), then the cost is $C_I(e) = F(e)K_I$.

Given the above example, it is easy to see that assigning an accurate instrumentation cost to an arc when all that is known is that arc's frequency is not possible: we do not know until the completion of the algorithm which set of arcs must be instrumented, and therefore we don't know whether an arc will need to be ISPLIT or whether it can be ISPLITCHEAP. Surprisingly, it is not possible to assign correct instrumentation costs to the arcs in the above situation even when the frequencies of all arcs entering a node are taken into account. A proof by counter-example is given below in Section 2.4.

The two-step process of assigning instrumentation costs to arcs and then applying a maximal spanning tree algorithm does *not* always produce an optimal instrumentation. There may be a polynomial-time algorithm for finding an optimal instrumentation, but as of this writing, I do not know what it is. If there are p instances of nodes like C in Figure 2.4, then finding the optimal instrumentation could require examining 2^{n^p} spanning trees, trying all n cost assignments at each of the p problem nodes. I conjecture that the problem is NP-complete. A fruitful line of search for a proof might begin with Szymanski's NP-completeness proof for the variable-span branching problem [47], which has some of the same characteristics as the instrumentation-cost problem.

Even so, finding the *optimal* minimal-cost solution in practice does not result in sufficiently more efficient instrumentations to warrant an intense search for a general solution. MINOPT works quite adequately in practice. The non-optimality of the solution matters little because

- we show in Section 2.3 that profiling requires less than 20% overhead, anyway;
- the problem configuration arises only under unusual circumstances (it occurred in only 2.4% of the 20457 basic blocks in our experiments); and,
- of those instances in which it does occur, the performance degradation is estimated to be less than 1%.

Throughout, whenever I refer to ‘optimal instrumentation’ algorithms, I mean optimal with respect to the instrumentation estimates.

2.3 Empirical data

I created a system that inserts optimal arc-counting in programs and used it to instrument several C programs by some of the algorithms mentioned in the previous section. The programs were first instrumented without the benefit of profile data; the frequencies of all arcs and nodes was one, resulting in a more-or-less random selection of instrumentation points by the algorithms (‘random’ in the sense that the selected instrumentation points were the result of vagaries of selecting a spanning tree from quick-sorting equi-valued elements). Next, a heuristic was used to assign relative frequencies to arcs: back-arcs and their target nodes were given higher frequencies than the rest on the assumption that a back-arc indicates a loop. Finally, profile data was used to compute a minimal-cost instrumentation.

The results are presented in Table 2.1. Four programs were compiled with `gcc -O` and instrumented: *intmm*, an integer matrix multiply; *compress*, the UNIX compression utility; *troff*, the UNIX typesetting program; and *cc1* of the *gcc* compiler. The first column for each program shows the running time of the program in CPU seconds, while the second column shows the running times of the instrumented versions of the programs expressed as a percentage of the original running time (more precisely, if N is the running time of the program without any instrumentation, and P is the running time of the program with instrumentation inserted, then the second column = $100 \times ((N/P) - 1)$). All programs were run on a Sun 3/140 and were compiled with the GNU C compiler, version 1.37.1. All running times are the average of 10 runs to smooth out system-dependent fluctuations.

The *intmm* program had no input data, and the contents of the matrices were initialized the same for each run (whether they were or not would not have made any difference to the running of the algorithm). The running times shown reflect the best that the various instrumentations could do for that program based

means		intmm		compress		troff		cc1	
not instrumented		29.06		17.03		96.91		28.91	
<i>prof</i>		29.93	3.0%	18.37	7.9%	127.91	32.0%	36.08	24.8%
<i>gprof</i>		31.45	8.2%	21.22	24.6%	164.81	70.1%	46.38	60.4%
FULLNODE		32.08	10.4%	26.78	57.3%	175.18	80.8%	42.40	46.7%
MINNODE	random	31.92	9.8%	20.58	20.8%	138.83	43.3%	36.40	25.9%
	heuristic	31.83	9.5%	21.32	25.2%	139.97	44.4%	37.55	29.9%
	profile	31.73	9.2%	19.90	16.9%	130.44	34.6%	35.72	23.6%
MINARC'	random	33.36	14.8%	20.83	22.3%	131.40	35.6%	36.46	26.1%
	heuristic	33.38	14.9%	20.26	19.0%	131.68	35.9%	37.01	28.0%
	profile	33.32	14.7%	19.94	17.1%	117.40	21.1%	35.21	21.8%
MINOPT	random	33.31	14.6%	20.39	19.7%	131.24	35.4%	36.35	25.7%
	heuristic	32.93	13.3%	20.15	18.3%	141.56	46.1%	37.85	30.9%
	profile	31.89	9.7%	19.89	16.8%	118.99	22.8%	34.90	20.7%

Table 2.1: Profiling overhead

on the profile data: the profile would be exactly the same each run. MINOPT's 9.7% overhead *vs.* MINARC's 14.7% reflects the fact that the inner loop in *intmm* mirrors exactly the situation in Figure 2.2. MINOPT was able to find an instrumentation that did not require the expensive arc-splitting that MINARC' was required to do.

For each of the other three programs, different input was used to create the profile data than was used create the numbers in the table. For *compress*, the profile data was generated by compressing *compress.c*, the source file for the utility. The numbers in the table are from compressing */usr/dict/words*, a 200Kb file containing a sorted list of 25,144 words. *Troff*'s profile data was created by typesetting a 48Kb language reference summary, and the numbers in the table are from typesetting a 190Kb technical report on a bibliographic database browser [49].

The *cc1* profile data was created by compiling *gcc.c*, the 23Kb source file for the process-dispatching front-end of the *gcc* compiler, and *combine.c*, a 46Kb source file for compile-time constant expression evaluation for the same compiler. The measured run compiled *cccp.p*, the 73Kb source file for the Gnu C-preprocessor. The sizes of the source files are *after* all pre-processing commands and all comments were stripped.

Using different input for profiling than for timing runs is necessary to convince us that we simply aren't 'training' the profile algorithms to a specific program. However, it does introduce some anomalies in the numbers in Table 2.1. For instance, MINARC', using profile data to instrument *troff*, resulted in the instrumented program taking 21.1% longer to run than did the uninstrumented version. This contrasts with MINOPT using profile data on the same program: the instrumented version of *troff* required 22.8% longer (the difference is statistically significant and not due to variations in measurement). Obviously, the instrumentation selected by MINOPT

	intmm	compress	troff	cc1
prof	36813	44189	66309	324357
gprof	20992	25304	42128	204160
FULLNODE	88	1492	10332	68636
MINNODE	56	796	6336	48216
MINARC'	56	808	6408	49040
MINOPT	56	808	6408	49040

Table 2.2: Comparison of profile data size requirements

does not do as well on the set of *troff* input as did the instrumentation selected by MINARC'. Such variation due to variations in the input are to be expected.

The numbers also indicate that the heuristic I used to try to guess which would be the heavily used arcs and nodes in the PFG is quite inadequate: it is better simply to accept a random assignment (i.e., if there is no profile data available). Perhaps a closer analysis of average aggregate program behavior could produce a heuristic that would do better than just random chance.

From the data presented in Table 2.1 we can see that MINOPT is definitely competitive with MINNODE, compares favorably with *prof* and *gprof*, and is definitely better than FULLNODE, the traditional technique for instrumenting programs.

Another benefit of MINOPT is demonstrated in Table 2.2. It shows the number of bytes required for the generated profile data. In all fairness, comparing the output of *prof* and *gprof* with the others is comparing apples and oranges. You can't get execution counts of basic blocks from *prof* or *gprof*, but then you can't get timing estimates from the others. However, the difference between FULLNODE and the MIN algorithms is significant.

2.4 Counter-example

In general, it is not possible to assign instrumentation costs to arcs in such a way that an optimal instrumentation can be found using MAXSPAN. To prove this, it suffices to show that there exists one program flow graph for which this is true. To this end, we construct a subgraph and show that for any instrumentation cost assignment algorithm, the subgraph can be embedded in a larger graph that causes MAXSPAN to select a non-optimal instrumentation.

Figure 2.6 shows a portion of a program flow graph that satisfies the criteria. Basic block *C* has two entering arcs that are both crowded jump arcs. We more-or-less arbitrarily assign execution frequencies to the arcs. The arc $A \rightarrow C$ is executed 200 times, the arc $B \rightarrow C$ 90 times. We assume that the cost of instrumentation $K_I = 1$ and the cost of a jump instruction $K_J = 1$. The values for these constants could be any non-zero value and we would still be able to find our counter-example,

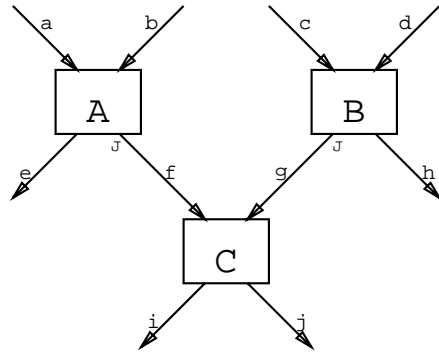


Figure 2.6: The problem configuration.

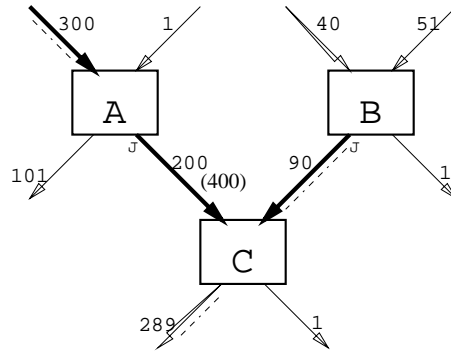


Figure 2.7: Case 1: Estimating that the most frequent arc is split expensively.

but that is not necessary to prove here: we need to show only that there exists one graph for which the algorithm is non-optimal.

The graph surrounding the sub-graph in Figure 2.6 is not shown, but is constructed such that the MAXSPAN algorithm will put nodes A , B , and C in the spanning tree last. This is easily done by creating the surrounding PFG such that each node has at least one entrance or exit arc with a frequency count higher than any arcs in the sub-graph. The nodes in the sub-graph are added to the spanning tree by selecting the arc in the sub-graph with the largest frequency. The arc is chosen from among all of the arcs shown in Figure 2.6, including the entrance and exit arcs of all three nodes. Exactly three of these arcs will be selected by the algorithm to complete the spanning tree.

Before invoking MAXSPAN, instrumentation costs must be assigned to each arc. The question is: does there exist an algorithm for assigning costs to the arcs $A \rightarrow C$ and $B \rightarrow C$ that has as its inputs only the frequencies of those arcs?

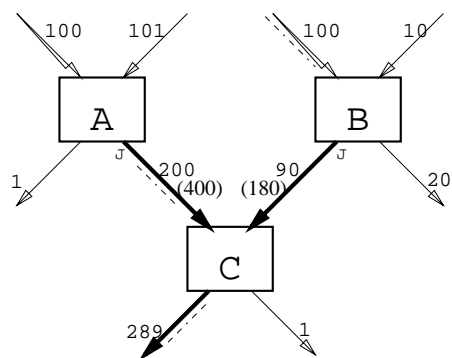


Figure 2.8: Cases 2 and 3: Estimating that the least frequent arc is split expensively, or that both arcs are split expensively.

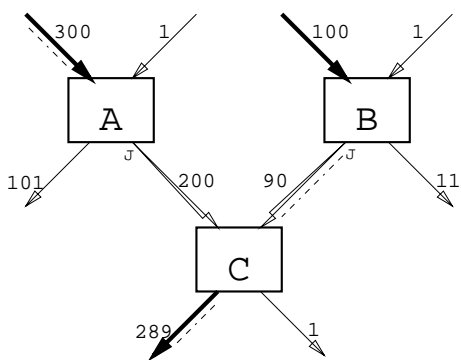


Figure 2.9: Case 4: Estimating that neither arc is split expensively.

If such an algorithm existed it would produce one of four results: it would assign the expensive-split cost to A , assign it to B , assign it to both arcs, or assign it to neither. Figures 2.7 through 2.9 show that no matter which assignment is made, there exists a consistent set of arc frequencies which will cause the MAXSPAN algorithm to misfire and select the wrong arcs for instrumentation – wrong in the sense of being non-optimal. In each figure, each arc is labeled with its frequency, and with its cost estimate in parenthesis if different from the frequency times K_I . For example, in Figure 2.7 the arc $A \rightarrow C$ has frequency 200, but its assigned cost is $K_I + K_J = 2$ times that, or 400. Given these instrumentation cost assignments, the arcs selected by the MAXSPAN algorithm are shown in bold. So in Figure 2.7, the cost of measuring the sub-graph is $1 + 40 + 51 + 101 + 1 + 289 + 1 = 492$. The dashed lines show a better spanning tree resulting in a cheaper instrumentation: in Figure 2.7 that cheaper instrumentation would cost $1 + 40 + 51 + 101 + 200 + 1 + 1 = 395$.

Figure 2.8 shows a set of frequencies for which guessing that the least frequent arc ($B \rightarrow C$) is expensively split, or that both arcs are expensively split, will also fail. Again, the bold arcs are the ones chosen by the MAXSPAN algorithm for inclusion in the spanning tree, while the dashed lines show a better selection, one resulting in a cheaper instrumentation.

Figure 2.9 shows that assuming both arcs will be ISPLITCHEAP (an impossibility in actuality) does not work either. The arcs selected by MAXSPAN result in an instrumentation that costs $1 + 1 + 101 + 200 + 180 + 11 + 1 = 495$. If neither $A \rightarrow C$ nor $B \rightarrow C$ is put in the spanning tree, then both will be measured. After the MAXSPAN algorithm chooses them for instrumentation, it is easy to see that of the two it is better to ISPLITCHEAP the most heavily used arc; hence the cost of 180 for instrumenting the lesser executed arc $B \rightarrow C$. It would have been better not to measure $B \rightarrow C$, as shown by the dashed lines. This better instrumentation would cost only $1 + 100 + 1 + 101 + 200 + 11 + 1 = 415$.

Therefore, an algorithm for assigning instrumentation estimates to arcs does not exist that has as its only inputs the frequencies of all arcs entering a node and that depends on the MAXSPAN algorithm to find the minimum instrumentation.

2.5 Problems with counting

There are several complications that must be considered when implementing a profiling system. The first is determining the point during compilation when the profiling code is to be inserted. The system I used operated on the assembly language output of the GNU C compiler. However, if the profiling code is inserted earlier by the compiler, then it is easy to finesse some of the task of counting. For instance simple loops (i.e., reducible sub-graphs with no mid-loop exits) that are executed a compile-time constant number of times do not need to increment a counter on each execution of the loop. In general, if the execution frequency of an arc in a program flow graph is known *a priori* that arc can be removed from the program-flow

graph prior to computing the set of instrumentation points. If the removal of the arc results in dis-connected sub-graphs, each sub-graph is treated separately. Even when the loop is not executed a constant number of times, but where strength-reduction can hoist the count increment out of the loop, the loop need not be instrumented. Rather the number of executions can be counted outside the loop and the counter incremented only once. Sarkar [42] shows one method for doing this using dependency graphs. However, it is impossible to tell from his paper exactly how much counting overhead is actually reduced by his technique. Further study and better numbers are needed here.

All of the the results in Section 2.3 computed each function’s instrumentation assuming that the number of times each function was called had to be counted. That is, the functions were instrumented one at a time without knowledge of how or from where the function was called³. However, if we have profiled the entire program, the number of times a function is executed is simply the sum of the execution frequencies of all call sites that call this function: there is no reason to recompute that number in the instrumentation of the function. Not all programs execute synchronously, as we have implicitly assumed throughout this discussion, nor do all call sites call only one function. If functions are called indirectly, for example by interrupt handling facilities, then it is mandatory that each function be instrumented separately to compute the number of times it was called.

The point to be made here is that these optimizations are possible only if the instrumentation code is inserted prior to the optimizing passes of the compiler.

Another problem for post-pass instrumentation is related to the machine architecture. If the processor’s instruction set makes use of condition flags, and if the life of condition flag values extends across basic blocks, then the profile code must preserve those values. This can be done either by saving and restoring the values of the condition flags around the instrumentation code (the method used in our experiments), using an instrumentation sequence that does not change the setting of the condition flags, or by scanning the basic block and inserting the instrumentation just before the first instruction that kills but does not use the condition flags [50]. This situation is much more naturally handled in the compiler proper than in a post-processor.

2.6 Conclusions

I have demonstrated that counting arcs is as cheap as counting nodes, and cheaper than counting every node. The MINOPT algorithm produces instrumentations that take 50-70% of the space required by instrumentations produced by the FULLNODE algorithm, execute in 70-80% of the time and provide significantly more

³This method was encouraged by *troff*, the only widely-used non-interactive program I know of that uses inter-procedural *gotos* as a major form of control flow.

information. MINOPT's instrumentations require approximately the same amount of time as MINNODE's and require slightly less space, but, again, they provide significantly more information. (In the next section I present an optimization that uses arc counts and would not work with only node counts.) MINOPT should be the instrumentation algorithm of choice for compilers/systems.

Chapter 3

A Low-level Use: Code Reorganization for Instruction Cache Performance

3.1 Instruction cache utilization

If a compiler has profile data available there are simple optimizations that can take advantage of the data. This chapter explores a simple optimization that requires profile data. Specifically, I will show how to utilize information about the runtime behavior of a program to enhance the performance of that program on architectures with an instruction cache.

There have been many investigations into improving computer performance by reorganizing programs' address spaces on virtual memory machines. In this chapter, I address the question of whether reorganization can be beneficial for machines with caches and examine the costs required to achieve improved performance. If an inexpensive way can be found to reorganize the address space of a program such that a small cache with code reorganization can have the performance of a larger cache without reorganization, the smaller inexpensive caches would be a more competitive choice.

Instructions that are executed close together in time are *temporally* local. Instructions that are close together in the address space are *physically* local. A cache turns temporal locality into physical locality by holding the most recently executed instructions in faster memory. Exactly how a cache should be implemented in hardware and which strategies for replacing the data in the cache are topics that have received a great deal of study. For an overview of cache designs and organizations, see Smith's survey article [45].

If we let C represent a cache, where each line i of the cache has an address $C(i).addr$ and contents $C(i).instr$ (the contents is an instruction for instruction caches), then when address a is referenced, the cache is examined to see if a is

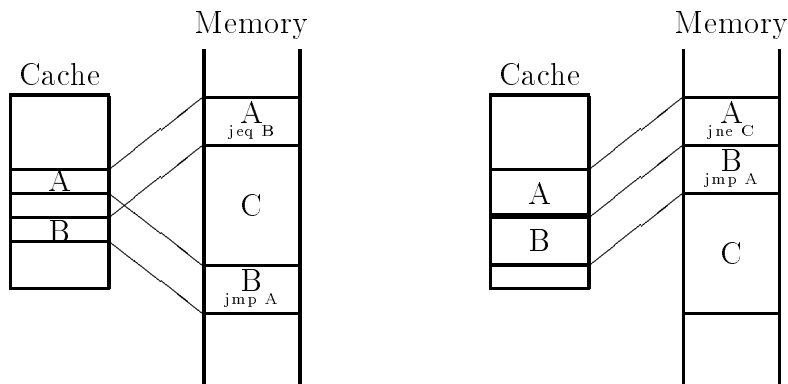


Figure 3.1: Removing cache contention by reorganizing

already in the faster memory. If $a = C(i).addr$ for some i , then the contents of that line is returned as the contents of the referenced memory address.

A fully-associative instruction cache is one that searches in parallel each line of the cache for the referenced address; i.e., if $C(i).addr = a$ for some i , then return $C(i).instr$. In a direct-mapped instruction cache, on the other hand, the address and its contents can be in only one line of the cache. Which line is (usually) determined by the low-order bits of the address; i.e., if $C(\text{lower bits of } a).addr = a$ then return $C(\text{lower bits of } a).instr$. Where in a fully associative cache an address and its contents can be placed in any line of the cache, in a direct-mapped cache, an address and its contents can be put in only one place; hence the name direct-mapped. The hardware required to do the parallel search is expensive to build, while a direct-mapped cache is much simpler and less expensive.

In either case, we are interested in several statistics as indicators of how well a cache performs. A critical statistic is the *miss ratio*: the number of times an address was referenced and it was not found in the cache. The dual of the miss ratio is the *hit ratio*: the number of times an address was referenced and found in the cache. By definition, then, $\text{miss ratio} = 1 - \text{hit ratio}$.

There are only two ways to improve the performance of a program in a cache: (1) decrease the probability that frequently-executed sections of the program compete for cache resources (Figure 3.1); and (2) increase the amount of useful information in the cache (Figure 3.2).

In Figure 3.1, assume the code at block A and the code at block B are the active portions of a loop. Assume further that, due to the size of the infrequently executed block C, blocks A and B are mapped to the same locations in the direct-mapped cache, as shown on the left. The loop containing these blocks can be made more efficient by moving A and B with respect to one another so that they do not

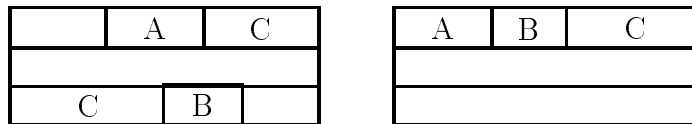


Figure 3.2: Improving cache utilization by reorganizing

conflict in the cache, as shown on the right.

In Figure 3.2 assume the blocks A , B , and C are of such a size that A and B could fit in a cache line. The cache lines on the left represent one way they might map into a direct-mapped cache, with the side-effect of loading infrequently executed code from block C into both lines. In the cache on the right the initial cache miss that loads the code from A also loads B , saving at least one cache miss in the execution of the loop; also, infrequently executed code from C takes up much less space.

For a fully associative cache there may be ways of reorganizing a program to improve its performance with respect to (2); little can be done as far as (1) is concerned. With direct-mapped caches, however, both (1) and (2) suggest easy ways to gain performance improvement. In a direct-mapped cache, contention is a function of the addresses of the competing program segments, which is easily controlled by a loader and/or compiler.

Mark Hill argues that direct-mapped caches are not only cheaper and easier to build [20], they also can give equivalent performance as more complex cache arrangements for the same silicon acreage invested. I have developed an algorithm called *Greedy Sewing* that uses arc counts to reorganize code for improved performance in direct-mapped caches, and that is independent of the parameters of the target cache.

While there have been published results for organizing *data* in memory to improve cache performance, there has been little published regarding rearranging the instruction space. For instance Janet Fabri [11] and K. O. Thabit [48] both discuss methods for improving the cache behavior of a program's data accesses, but say little about the behavior of the code itself.

Some recent work on code reorganization include Scott McFarling's work at Stanford [34], work done at Hewlett-Packard by Pettis and Hansen [36], and Hwu and Chang's work on combining reorganization with the in-lining of procedures [21]. McFarling's work differed from mine by concentrating on positioning basic blocks based on their frequency counts and by utilizing knowledge of the target cache. Hwu and Chang extended my work by doing actual in-lining (as opposed to my pseudo-

inlining). Pettis and Hansen pretty much duplicated my work, with the exception that their algorithm reorganizes the entire program instead of just those areas where the vast majority of the improvement is gained.

3.2 Greedy Sewing

A program control-flow digraph (or program flow graph or PFG for short) is a set of basic blocks V and directed arcs E . If $e = v \rightarrow w$ for nodes $v, w \in V$ then we define $src(e) = v$ and $snk(e) = w$. Associated with each arc e (basic block v) in the graph is a positive integer $F(e)$ ($F(v)$) representing the number of times this arc (basic block) was executed during the execution of the program (so $S = \sum_{e \in E} F(e) = \sum_{v \in V} F(v)$). Associated with each basic block v are functions $onThread(v)$ that returns the thread that basic block v is on, $onHead(v)$ that returns true if the block v is at the head of its thread, and $onTail(v)$ that returns true if v is at the tail of its thread. Given threads $t, s \in \mathcal{T}$, we define the procedure $append(t, s_1, s_2, \dots)$ to concatenate the threads s_i in order onto thread t , and delete threads s_i from \mathcal{T} . The functions $first(t)$ and $last(t)$ return the first and last blocks, respectively, on the thread t .

The basic idea is to sew threads together such that the order of the basic blocks in a thread tends to improve the correspondence between the static spatial locality of basic blocks with their dynamic temporal locality. We define the function $canStitch(u, v)$ to return true if the nodes u and v can be concatenated onto the same thread; this is true only if $u \neq v$ and $onHead(u)$ and $onTail(v)$ are both true.

We define the procedure $Stitch(e)$ to ‘sew’ two threads together:

Procedure 3 *Stitching Basic Blocks:*

Input: An edge e in a PFG; a set of threads of basic blocks \mathcal{T} .

Result: If the source and sink blocks of e can be concatenated, the set \mathcal{T} is modified such that it contains one less thread due to the concatenation of the two member threads.

Method:

```

Stitch( $e$ : Arc)
begin
  if canStitch( $e$ ) then
    append( onThread( $src(e)$ ) , onThread( $snk(e)$ ));
  end

```

□

Using these functions, we are now ready to lay out a preliminary version of the Greedy Sewing algorithm.

Algorithm 4 *Greedy Sewing Algorithm(1):*

Input: A PFG (V, E) , and a parameter p such that $0 \leq p \leq 1$; and a set \mathcal{T} of threads that are initialized such that each basic block is on its own thread; after initialization $onHead(v)$ and $onTail(v)$ are true for all $v \in V$.

Result: The set of threads \mathcal{T} is modified to indicate the relative ordering of the basic blocks in V . A thread $t \in \mathcal{T}$ specifies the order in which the basic blocks are to be placed contiguously in memory. There is no implied ordering of basic blocks on different threads.

Method: The parameter p is used to specify what portion of the arcs will be examined. That is, a set of arcs \mathcal{A} will be processed where for all $a \in \mathcal{A}$, $F(a) \geq F(b)$ for all $b \notin \mathcal{A}$, and $\sum_{a \in \mathcal{A}} F(a) < p * \sum_{e \in E} F(e)$. That is, setting $p = .90$ would cause the main loop of the algorithm to be repeated until sufficient arcs had been processed to account for 90% of all arc traversals. This means that usually 5-10% of the arcs, and even fewer basic blocks, need be reorganized.

```

Greedy( $p$ : real;  $\mathcal{A}$ : Set(Arc))
begin
  assert( $0 \leq p \leq 1$ );
   $S \leftarrow p \times \sum_{e \in \mathcal{A}} F(e)$ 
   $s \leftarrow 0$ 
  while ( $s < S$ ) do
    Select  $e \in \mathcal{A}$  such that  $F(e)$  is maximum.
     $E \leftarrow E - \{e\}$ 
     $s \leftarrow s + F(e)$ 
    if canStitch( $src(e)$ ,  $snk(e)$ ) then
      Stitch( $src(e)$ ,  $snk(e)$ )
    endwhile
end

```

□

In preliminary tests of the greedy algorithm, several situations were observed that this simple algorithm did not handle adequately. While the majority of program improvement comes from the simple version of Greedy Sewing, it does not take very many special cases to eat into those savings. Based on observations in these preliminary runs, the simple algorithm was enhanced with some checks for special cases. For instance, consider the flow graph in Figure 3.3. If the path $A \rightarrow B \rightarrow D$ is the more frequently executed path, then the thread ABD will be formed (1). Since C cannot be sewn to either A or D now, it remains a singleton thread (2). If it is very infrequently executed, then it makes little difference where C is placed relative to the thread ABD . However, if the path $A \rightarrow C \rightarrow D$ is only slightly less frequently executed than the path $A \rightarrow B \rightarrow D$, and if the sum of the sizes of A , B , C , and D

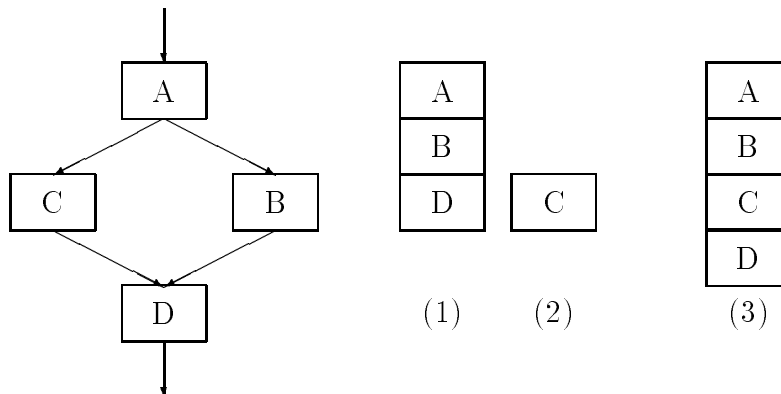


Figure 3.3: Two threads from if-then-else

are small enough that they might all fit in a cache, then the single thread $ABCD$ in Figure 3.3(3) is preferable over the two threads (1) and (2).

Since the Greedy Sewing Algorithm is general and does not depend on any particular cache configuration or size, it cannot know whether any set of basic blocks will fit in a cache, and so uses a heuristic to attempt to capture instances of this configuration of basic blocks. The function *isSmallEquiCondl* checks for basic blocks matching exactly this configuration—i.e. the basic block ends with a conditional jump instruction, the two arms of the conditional have at most one basic block in them and are very nearly equi-probable—and when found the procedure *StitchCondl* creates the longer thread. The actual mechanics of putting a small if-then-else on a thread is straightforward in procedure *StitchCondl* (see the next page).

A second common, but more complicated, situation is pictured in Figure 3.4 where a procedure P is called from a basic block A . We want to concatenate block P_r with block B because the same considerations that applied to the previous if-then-else example apply here: if the frequent path through the procedure is small enough such that A , the frequently executed portions of P , and B could fit in the cache, then we would like to construct the thread shown on the right of Figure 3.4. The procedure *StitchCall* effectively constructs the arc $P_r \rightarrow B$ such that P_r and B are eventually made contiguous. When the bottom of A is sewn to the top of P_0 , we say that procedure P has been *pseudo-inlined*. A basic block containing a single jump instruction is inserted between the call and the target to maintain the semantics of the original code.

That is the simple view of pseudo-inlining. It is complicated by the fact that at the time we invoke *StitchCall* on the arc e (using the notation in the example in Figure 3.4, it will be one of arc $A \rightarrow B$ or arc $A \rightarrow P_0$, depending on the vicissitudes of the sorting algorithm) we have not yet encountered the return block P_r , and may

Function 5 *isSmallEquiCondl*:

Input: An arc $e \in E$.

Output: Returns true if $src(e)$ is the head of a small if-then-else with approximately equi-probable true and false arms. A global constant δ defines what is meant by equi-probable.

Method:

```

isSmallEquiCondl( $e$ : Arc) return boolean
begin
  if  $src(e)$  ends with a conditional branch (and therefore has
    two exit arcs  $e$  and  $w$ )
    and  $snk(e)$  has one exit arc  $e_x$ 
    and  $snk(w)$  has one exit arc  $w_x$ 
    and  $snk(e_x) == snk(w_x)$  --they go to the same block
    and  $snk(e_x) != src(e)$  --they do not make a loop
    and  $snk(e_x) != src(e_x)$ 
    and  $(|F(e) - F(w)| / (F(e) + F(w))) < \delta$ 
  then
    return true
  else
    return false
end

```

□

Procedure 6 *StitchCondl*:

Input: An arc $e \in E$ that is one arm of a basic block that ends in a conditional branch instruction. Assumes that *isSmallEquiCondl*(e) is true.

Result: A new thread is added to \mathcal{T} that contains the four basic blocks of the if-then-else.

Method: Let w and e_x be defined as in the function *isSmallEquiCondl*. Then a new thread is created by concatenating

```

append(onThread( $src(e)$ ),
      onThread( $snk(e)$ ),
      onThread( $snk(w)$ ),
      onThread( $snk(e_x)$ )).

```

□

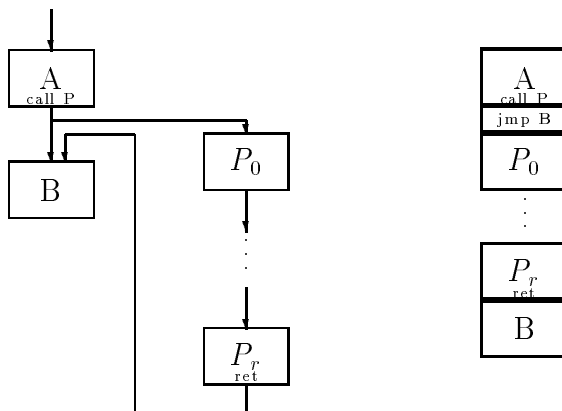


Figure 3.4: Pseudo-inlining

not encounter it if it is not one of the hot spots of the program. During the normal operation of Greedy Sewing all blocks ending with a return instruction will end up at the tail of a thread: there are no exit arcs from a return block. At the same time, we do not want basic block B , the one following the call instruction, to be threaded with a less frequently occurring basic block.

Let the function $target(v)$ return the basic block that is the target of the call instruction that ends block v (undefined otherwise), and let $returnsTo(v)$ be the basic block to which the called procedure returns. Then, whenever an arc e is selected for which $isCall(src(e))$ is true (i.e. the basic block $src(e)$ ends with a call instruction), the procedure *StitchCall* (Procedure 7) modifies block $B = returnsTo(src(e))$ so that $onHead(B)$ is false (even though B is still (on) a singleton thread) and adds $src(e)$ to a set of remembered basic blocks \mathcal{R} . At the end of the Greedy Sewing Algorithm, the procedure *append_R_blocks* (Procedure 8) is invoked to append all of the blocks $r \in \mathcal{R}$ to the appropriate threads.

The entire Greedy Sewing Algorithm used in our experiments is given in Algorithm 9.

3.3 Results

I used the profile-collection techniques of the previous chapter to collect arc frequencies of several programs. After profile data was collected, the program *reorgBBs* then read the original assembly language files and reorganizes them based on that profile.

3.3.1 The Programs and Traces

There were three programs chosen for experimentation and each program had four versions created: the normal, unreorganized version produced by the Gnu

Procedure 7 *StitchCall*:

Input: An arc $e \in E$ such that $isCall(src(e))$ is true.

Result: The target of the call in $src(e)$ is pseudo-inlined into the thread containing $src(e)$. The set of threads T and the set of remembered basic blocks R are modified so that the pseudo-inlining can be completed later.

Method:

```

    StitchCall( $e$ : Arc)
    begin
    assert( $isCall(src(e))$ )
     $t \leftarrow target(src(e))$ 
     $r \leftarrow returnsTo(src(e))$ 
     $append(onThread(src(e), t))$ 
     $onHead(r) \leftarrow \mathbf{false}$ ;
     $add\ src(e)\ to\ \mathcal{R}$ 
    end

```

□

Procedure 8 *append_R_blocks*:

Input: The set R of return blocks; the set of threads T .

Result: The threads containing the return blocks are appended to the threads containing the corresponding call blocks. All return blocks are at the head of a thread, even though *StitchCall* modified them to appear otherwise.

Method:

```

    append_R_blocks( $\mathcal{R}$ : Set(Node))
    begin
    for each  $e \in \mathcal{R}$  do
         $append(onThread(e), onThread(target(e)))$ 
    endfor
    end

```

□

Algorithm 9 *Greedy Sewing Algorithm(2):*

Input: A PFG (V, E) ; a parameter p such that $0 \leq p \leq 1$; and a set of threads \mathcal{T} that are initialized such that each basic block is on its own thread; after initialization of \mathcal{T} $onHead(v)$ and $onTail(v)$ are true for all $v \in V$.

Result: The set of threads \mathcal{T} is modified to indicate the relative ordering of the basic blocks in V .

Method: The parameter p is used to indirectly specify what portion of the arcs will be examined.

```

Greedy( $p$ : real;  $E$ : Set(Arc))
begin
 $S \leftarrow p \times \sum_{e \in E} F(e)$ 
 $s \leftarrow 0$ 
 $\mathcal{R} \leftarrow \emptyset$ 
while ( $s < S$ ) 'do
  Select  $e \in E$  such that  $F(e)$  is maximum.
   $E \leftarrow E - \{e\}$ 
   $s \leftarrow s + F(e)$ 
  if isSmallEquiCondl( $e$ ) then
    StitchCondl( $e$ )
  else if isCall( $e$ ) then
    StitchCall( $e$ )
  else if canStitch(src( $e$ ), snk( $e$ )) then
    Stitch(src( $e$ ), snk( $e$ ))
  endwhile
append_R_blocks( $\mathcal{R}$ );
end

```

□

name	trace length unreorganized	trace length reorganized		
		$p = .80$	$p = .90$	$p = .95$
SCRUNCH	9,405,156	9,656,437	9,715,946	9,693,277
TROFF	8,059,174	8,343,440	8,325,470	8,338,350
CC1	8,263,593	8,268,313	8,293,376	8,301,226

Table 3.1: Summary of traces: number of instruction words fetched

name	number of basic blocks	number of blocks reorganized		
		$p = .80$	$p = .90$	$p = .95$
SCRUNCH	1,233	22 (1.8%)	27 (2.2%)	32 (2.6%)
TROFF	4,000	149 (3.7%)	207 (5.2%)	318 (8.0%)
CC1	26,407	727 (2.8%)	1,317 (5.0%)	1,939 (7.3%)

Table 3.2: Summary of traces: number of basic blocks

C compiler, and three reorganized by the Greedy Sewing Algorithm with p set to .80, .90, and .95. A summary of the programs, the basic block counts, and trace sizes is in Tables 3.1 and 3.2. I collected a trace of each of these twelve programs which were then used as input to Mark Hill’s DineroIII cache simulation program [19]. Each trace was simulated on eleven different cache configurations: 256 bytes with 4 and 8 byte blocks; 1024 byte cache with 4, 8, and 16 byte blocks; and 4096 byte cache with 4, 8, and 16 bytes, each using single associativity (direct-mapped) and two-way set associativity.

The first program was *scrunch*, a Huffman encoding algorithm. The profile was generated by *scrunching* a 200K spelling dictionary. The trace was created by *scrunching* *scrunch.c*, a 42Kb C source file.

A second program, *troff*, was chosen because of its wide use in UNIX environments. The profile was generated by *troffing* three separate technical documents, chosen to represent a broad and typical use of the program. The first document consisted of 103K bytes after being preprocessed by *tbl*, *eqn*, and *grn*. This included 1933 lines (32K bytes) of *troff* commands, the remainder being plain text. The other two documents totaled 228K bytes and contained 4004 lines (73K bytes) of preprocessed *troff* commands. The trace was created by *troffing* a reduced version of the first document of length 7705 bytes, of which 273 lines (2728 bytes) were *troff* commands.

A third program was the Gnu C compiler itself. The profile was collected of the compiler compiling three Gnu C source files: *topev.c*, *loop.c*, and *recog.c*. They totaled 79K bytes, with 20, 12, and 15 C function definitions, respectively. The trace was collected while compiling *genemit.c*, a 6Kb file containing nine function definitions.

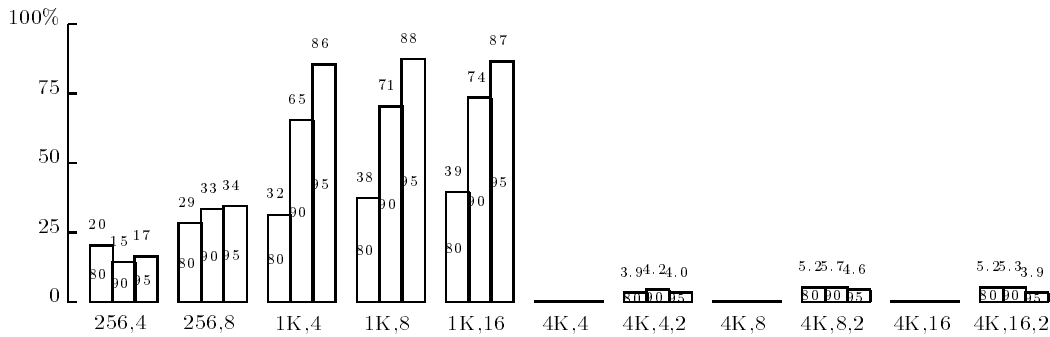
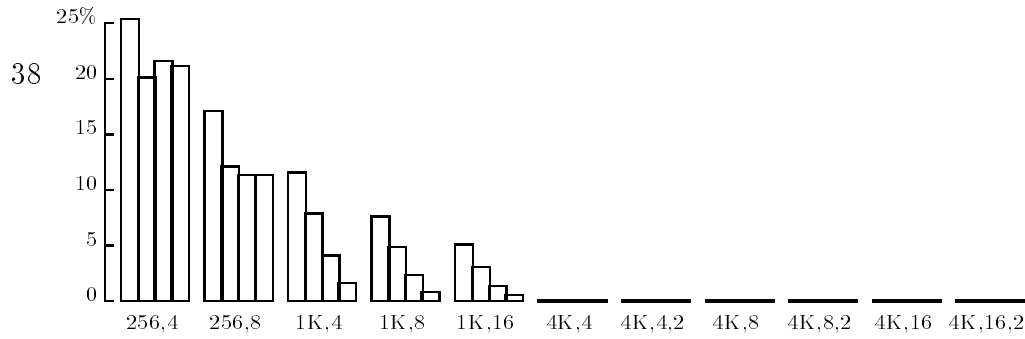


Figure 3.5: Scrunch miss rates and percent improvement

3.3.2 Miss Rates

Figures 3.5, 3.6, and 3.7 are histograms of the results of running the traces of the programs through the cache simulator. (Since the reorganization algorithm didn't take any cache parameters into account, the same traces are used in all eleven simulation runs for each program.) They show the miss rates for the four versions of each program on each of the eleven cache configurations. The leftmost bar of each group of four is the average miss rate of the unreorganized program. The other three of the group, from left to right, are the average miss rates for the reorganized versions for $p = 80\%$, 90% , and 95% respectively. The cache configuration is noted beneath each group. Figures 3.5, 3.6, and 3.7 show the improvement in miss rates of the reorganized versions over the miss rates of the unreorganized versions (i.e. $1 - (M_r/R_r)/(M_u/R_u)$).

There are instances where reorganization can buy the (miss rate) equivalent of a larger cache. For example, looking at Figures 3.5, 3.6, and 3.7 we see that the reorganized programs using a 256 byte cache with 8-byte blocks consistently had as good as or better miss rates than its unreorganized version running in a 1K cache with 4-byte blocks.

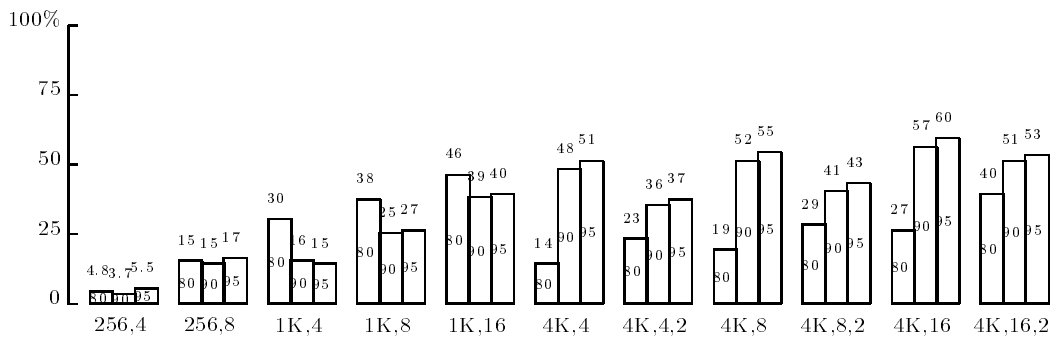
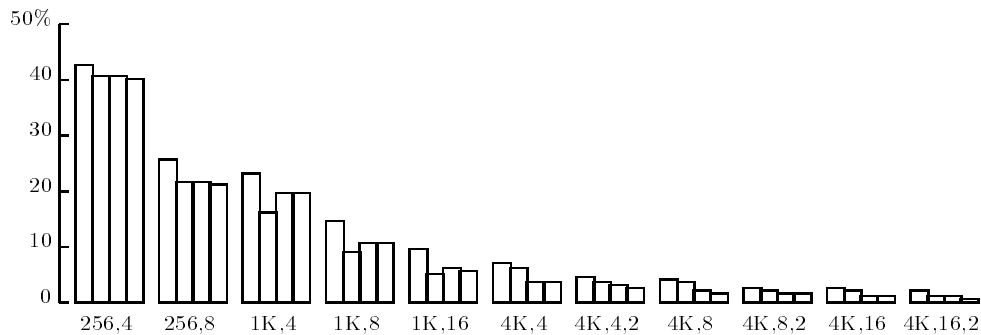


Figure 3.6: Troff miss rates and percent improvement

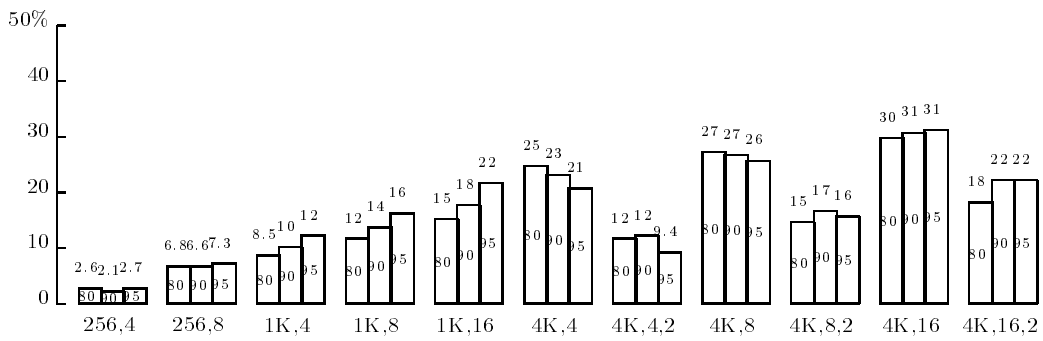
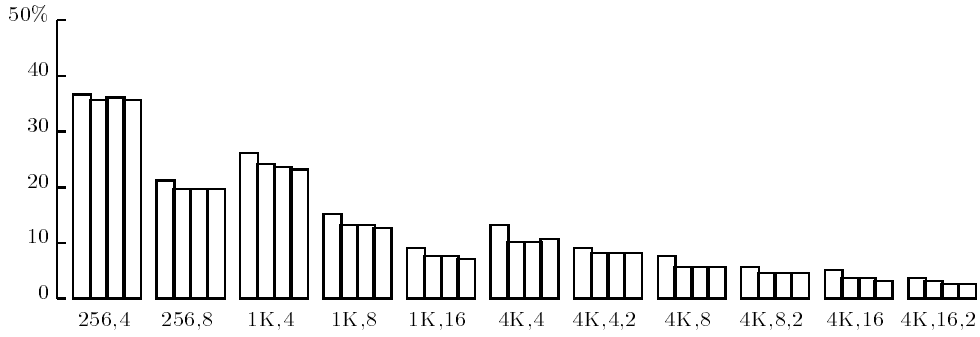


Figure 3.7: Cc1 miss rates and percent improvement

3.3.3 Performance improvement

Let R_u be the number of instruction fetches on the original, unreorganized program, and let R_r be the corresponding number for a reorganized version. For the sake of these estimates, we will assume $R = R_u = R_r$. We see from Table 3.2 that this is not strictly true, but they are sufficiently close for our purposes here. Let M_u be the number of cache misses and H_u the number of hits when the original, unreorganized program is run, and M_r and H_r be the corresponding values for the reorganized program. Define the miss rates $m_u = M_u/R$ and $m_r = M_r/R$, and let $m_\Delta = m_u - m_r$ be the difference in miss rates. Let t_h be the time required to handle a cache hit, and let t_m be the time required to handle a cache miss. The running time of the original program is then $T_u = t_h h_u + t_m m_u$ and the improved running time is $T_r = t_h h_r + t_m m_r$. Finally, define $f = t_m m_\Delta / T_u$, the fraction of the original program's time taken up by cache misses that are turned into cache hits, and $K = t_m / t_h$ the ratio of the cost to handle a miss to the cost to handle a hit. Then by Amdahl's Law:

$$T_r/T_u = (1 - f) + f/K. \quad (3.1)$$

We can now estimate the improvements in performance from reorganization taking our example from the specification of the SPUR memory architecture [18]. In general, the cost of a miss is very high on multi-processor, shared-bus systems due to bus contention or the length of the cache line. SPUR has a 512-byte on-chip cache and 128Kb off-chip cache. According to Mark Hill, a miss in the on-chip cache costs three times as much as a hit, assuming the instruction to be in the off-chip memory cache [20]. Let us assume the on-chip cache shows a normal miss rate of about 20%, and that we can improve that to 15% by reorganizing. Then $f = 3 * .05 / (3 * .20 + .80) = .107$. Plugging this into (1) above, we get $T_r/T_u = .929$, i.e. the program executes in only 92.9% of the time of the original, a 7.1% improvement. The maximum possible improvement is 29.6% assuming the unattainable miss rate of 0%.

For the SPUR architecture, an off-chip cache miss will cost 12 to 20 times that for handling a cache hit. SPUR therefore has a very large mixed cache to combat this penalty. If we assume that reorganization can reduce SPUR's miss rate by an absolute 0.25% (e.g. from 1% to 0.75%), then, assuming $K = 17$ (a number lifted from Katz and Eggers [26]), $f = 17 * .0025 / (17 * .01 + .99) = .0366$. Plugging this into (1) above, we get $T_r/T_u = .966$, a 3.4% improvement in performance. This is in addition to the performance improvement for the on-chip cache noted above. With these assumptions, we predict reorganization can improve SPUR's performance by about 10%.

This prediction is consistent with other numbers recently published for similar systems. Pettis and Hansen [36] at Hewlett-Packard Laboratories report 10%-26% improvement on a machine with a 16Kb unified cache (the HP-UX 825). When they increased the cache size to 128Kb, the improvements decreased to less

than 10% on the HP-UX 835; all five benchmarks' improvements averaged 5%, but with a very wide range of differences (0.8% to 9.3%). Their study concentrated on placing code to improve cache performance utilizing information about the size of the cache. They also used arc counts, rather than node counts, but their algorithm operated over the whole program, not just the hot spots. They also counted each and every arc in the flow graph, which may account for the compiler going twice as slowly when instrumenting the program. Their algorithm for positioning the code slowed the compiler down by about 15-20%.

Figures 3.8, 3.9, and 3.10 show the theoretical improvement in execution performance of the reorganized versions over the original unreorganized versions when the cost factor $K = 2, 3, 4, 5, 6, 7, 10, 15, 20,$ and 25. Each graph has a column for each of the eleven cache configurations. A line within a column plots the expected performance of the indicated program on that cache when reorganized with the three values $p = .80, .90, .95$ moving from left to right. $K = 2$ is the very top line in each column, and $K = 25$ is the bottom-most line in each column, yielding a range in which I expect reorganization to improve the performance of the programs. So we see that for a 1K cache with 4-byte blocks and $K = 2$, a version of *scrunch* reorganized with $p = 80\%$ would take about 97% as long as the unreorganized version (the left end of the topmost line in the third column from the left). It would take only about 77% as long if $K = 25$ (the left end of the bottom line in that column), and only about 38% as long when $K = 25$ and $p = 95\%$ (the rightmost end of the bottom line in that column).

3.4 Limitations

The algorithm for *StitchCall* is not quite correct, since it does not handle nested procedure calls correctly. The net effect on the numbers is not at all clear, but it should be slight in whichever direction it goes. There simply were not that many sequences of nested procedure calls that could fit in a cache in the code I used as test cases. This would not be true in languages that encouraged the use of many small procedures, e.g., C++.

Furthermore, it is not clear whether pseudo-inlining a procedure in only one location is sufficient. Further tests should be performed to determine whether it is worthwhile to copy the main thread of execution of a frequently executed basic block into multiple locations. If I had to guess as to which would have the most effect on the results—correcting the nested procedure call problem or making multiple copies of threads of procedure execution—I would say that thread copying would probably have more effect.

There is a potential problem in *isSmallEquiCondl* in that it will not handle correctly a small if-then-else where one of the arms of the conditional is empty, and the common exit block loops on itself. This situation never arose in my experiments, and so the potential problem was not detected until this writing.

Figure 3.8: Scrunch performance for $K = 2, 3, 4, 5, 6, 7, 10, 15, 20,$ and 25

Figure 3.9: Troff performance for $K = 2, 3, 4, 5, 6, 7, 10, 15, 20,$ and 25

Figure 3.10: Cc1 performance for $K = 2, 3, 4, 5, 6, 7, 10, 15, 20,$ and 25

I have not considered many of the parameters for designing a cache that may be relevant. For example, I have not considered the increase in bus contention caused by going to a larger block size. Nor have I considered alternative cache management strategies such as sub-block placement. My goal was simply to explore improving the performance of the cheapest of these design alternatives with simple compiler enhancements.

I have not considered cache effects such as cold start misses or cache flushes due to system interrupts or context switches. I did not consider any parameters of the target caches when reorganizing code. It was not clear at the beginning of this research how much the parameterization of the algorithms by the cache characteristics would benefit the program's performance; hence I went for the simple solution first.

I haven't solved the problem of *case* statements satisfactorily. Currently, it is possible for the reorganizer to move the code around to such an extent that the jump table can end up quite a distance away from one or more of its targets. On the 68020, this presents a practical problem since jump tables with half-word *pc*-relative entries are much faster than full word entries. If an item of a case is a "hot spot", it is very difficult to relocate it and still satisfy the distance constraint of jump tables. The only program that gave me real problems was the Gnu C compiler, for which I generated full-word jump table entries. This does not change any of the *miss* rate results significantly, but in real life, it would be an unacceptably slow implementation due to the slower execution of the table jump.

Finally, there are architectures that present difficulties for the Greedy Sewing algorithm. An example is the MIPS-X instruction set [8], which has asymmetric conditional branch instructions. Due to the nature of the MIPS-X pipeline, each conditional branch instruction is followed by two instructions that are fetched before the CPU has determined whether the branch will be taken. Each conditional branch instruction also has a *squash* bit that, if 'on', prevents the execution of these two instructions if the branch is not taken. Both of these delay slot instructions are always executed whenever the branch is taken: there is no way to *squash* their execution when the branch is taken. This makes sense if all programs follow the pattern of code generated by most compilers, where the conditional test is at the 'bottom' of the loop and the conditional branch is, therefore, almost always taken.

However, the results of Greedy Sewing's reorganization described here turns those statistics upside down. The algorithm almost guarantees that the head and tail of a frequently-executed loop will be made contiguous, and that execution will almost always fall through the conditional test to the head of the loop: the conditional branch is almost never taken. This leaves three options: (1) fill the delay slots as is currently done by MIPS-X compilers, followed by a jump instruction to the infrequent target (normally the loop-exit); (2) reverse the sense of the conditional and try to fill the delay slots with instructions that don't have to be squashed when the branch is *not* taken (because there is no squash bit for this direction); or, (3) punt and put no-ops in the delay slots.

Option (1) puts infrequently executed instructions right in the middle of high-frequency basic blocks, working against one of the aims of reorganization (better cache utilization in high frequency code). Option (2) sounds plausible, but it is difficult to find instructions that can always be executed no matter which way the branch goes. It may be possible to generate instructions to un-do the effects of the delay-slot instructions when the branch is finally taken, but this begins to get complicated and presents the possibility of really slowing down a frequently executed inner loop. Option (3) is an obvious loss. McFarling [34] implemented option (1), and reports that the size of repositioned code increases about 14%.

Due to the fact that the available MIPS-X compilers filled the delay slots before emitting assembly language code requiring my software to have a MIPS-X assembly language parser, and given the fact that the only profile data I had was the basic block counts generated by their compiler, applying Greedy Sewing to MIPS-X code was too far outside the reach of this research.

3.5 Conclusions

Profile driven code reorganization definitely improves the performance of programs. In envisioned programming environments where profile data is a permanent part of the information manipulated by both programmer and compiler, these improvements would come simply and cheaply. My experiments have shown improvements in miss rates on the order of 30% to 50%, and sometimes as high as 50% to 80%. These figures were obtained by relocating only 3% to 8% of the basic blocks of typical programs.

Chapter 4

A High-level use: Implementation selection of abstract data types

In the previous chapter, I demonstrated one way in which profile data could be used profitably by an optimizing compiler at a low level. Code reorganization is attractive since it improves performance at a small cost, that cost being the time it takes to decide the order in which to emit a few of the basic blocks of a program.

In this chapter, we will look at what a compiler might do with profile data at a very high level. Specifically, I developed `TYPESETTER`, a system for selecting implementations for abstract data type representations and functions. By assuming that profile data exists for a program, we have seen that low-level program transformations can use very simple algorithms to achieve improvements comparable to much more complicated algorithms. With `TYPESETTER` I tested whether comparable simplifications could be made in selecting implementations of high-level abstract data types.

4.1 The Problem

Before stating the problem, it is useful to differentiate between two classes of programmers that would make use of `TYPESETTER`. The *User* of `TYPESETTER` would write a program using only the available abstract data types, and would not be concerned with how those abstractions were eventually implemented (as long as the implementations were relatively inexpensive, of course). The *Implementor* is the programmer that adds implementations to the `TYPESETTER` system. I do not envision that `TYPESETTER` can be (or should be) an extensible language system at the User level. If new implementations are to be added, it is an enhancement to the system, and not simply the shipment of a new library. This difference between User and Implementor allows us to discuss efficiently the difference between using an abstraction and implementing it by referring to the specific programmer.

The general problem can be stated simply: what is the most efficient imple-

mentation of a User's program? This includes selecting the most efficient instruction sequences as well as the best implementations for the program's data structures. Most compilers completely side-step this latter problem by giving the programmer a specific set of data types with unique implementations and letting the programmer construct the necessary data structures with the pre-defined data types provided by the system.

For example, there are many ways sets of objects might be implemented: linked lists of various kinds, bit maps, arrays, trees organized by various techniques, etc. Which of these implementation should be used for a particular program depends on the algorithms used in the program and the characteristics of the data. For the most part, letting the compiler select which implementation to use based only on static declarations has proven viable, but difficult and expensive. The problem is difficult even when attention is focused on a small set of abstractions, as the SETL language effort has shown [46,51]. Barstow's PECOS system [3,4] is an attempt to collect a database of rules and heuristics that allow a programmer's specification of a program to be given an implementation. Elaine Kant [24] extended the system to consider rules and heuristics regarding the efficiency of various implementations.

My approach is to assume that programming environments of the future will be collecting and maintaining much more information about a program than the programmer's static declarations. In particular, the collection and utilization of profile data will become a matter of course. In this chapter, I address the following questions:

- Can profile data reduce the complexity of the representation selection problem to a level that compilers can make such choices effectively?
- What information needs to be collected by the profile mechanism so the selector can make effective choices?
- How much control over the collection of profile data can be put in the hands of the designers and implementors of abstract data types?
- Is there a general algorithm for selecting representations that works for a wide variety of abstract data types? That is, can we limit the overall task of the Implementor to implementation of the ADT, specification of what profile data to collect, and specification of the runtime resources used by the implementation?

Obviously all these questions are interrelated: for example, the selection algorithm will influence the kinds of profile data that will be necessary, and possibly the detail to which the Implementor must go to describe the behavior of an implementation.

We are interested in improvements only from data representation selection, in contrast to algorithm transformations, traditional examples of which include such optimizations as code motion, finite differencing, strength reduction, etc. (Low [31] calls these *representation dependent* optimizations.)

As an example, consider the implementation of:

$$S \leftarrow S \cup \{a, b, c\}$$

which may be more efficiently implemented as

```

put a in S
put b in S
put c in S

```

depending on the representation selected for the set S . Consider also the boolean expression:

$$X \in (S1 \cup S2)$$

which is usually much more efficiently computed as

$$(X \in S1) \vee (X \in S2)$$

For the purposes of this research, we will assume that such optimizations are discovered by a high-level optimization module. This ignores problems introduced by the interaction of this ‘higher level’ optimization and our representation selection mechanism, but allows us to investigate the use of profile data in the selection process. Assuming the latter is possible, the interaction problem can be investigated later.

Unlike Low, we do not want to limit expressions to be homogeneous in representations. That is, in the expression

$$S1 \leftarrow S1 \cup S2;$$

if $S1$ and $S2$ are two sets, it may be the case that they have two different representations. We are interested in seeing if there are situations in which it is profitable to handle the overhead imposed. In this example, either

1. One of $S1$ or $S2$ must be converted to the same representation as the other, or
2. $S1$ and $S2$ are both converted to a third representation, or
3. a routine to take the union of objects of type $itype(S1)$ and $itype(S2)$ must be generated by the compiler, or
4. there must already exist a procedure that can explicitly handle the union of these two representations.

Larry Rowe [38] explored the problem of generating implementations and we will not pursue it in this study. The `TYPESETTER` prototype requires that the explicit function must exist (corresponding to the fourth option above).

`TYPESETTER` has been designed to allow experimentation with various kinds of profile data in addition to execution counts. For instance, knowing that the

attempts to add an element to a set usually have no effect because the element is already in the set may affect which implementation is chosen for that set. It is easy in `TYPESETTER` to count how many times the add-an-element function was called to add an object already a member of the target set. It would be exceptionally difficult to determine analytically when such a situation holds.

The traditional metric (or objective function) used for selecting one representation over another has been the space-time product, but there are complications in dealing with parameters of the system in which the software is to be run. For example, it may be that representation X is better as long as the amount of memory used does not exceed physical memory limits, otherwise using the more compact representation Y will cause less thrashing of the virtual memory system. The inefficiency due to a denser encoding would be offset by the improved performance of the system as a whole. It is problematic how to specify these kinds of limits, particularly those that depend on system parameters that are difficult to profile or can vary dynamically and orthogonally to the actions of the program (e.g. dynamic paging rates). Therefore, this study has not attempted to determine the ‘best’ way of characterizing program behavior. The system is designed such that each interface function has exactly one evaluation function, which returns a real number that represents the relative behavior of the interface function at a particular call site.

Also, this study limits itself to implementations that do not require automatic changes to User-defined data structures; all abstract data types implemented in the system will be ‘exomorphic’—only pointing to objects the user is manipulating. For example, an exomorphic implementation of a list would manipulate only references to the objects in the list. In contrast, an *endomorph*ic implementation might include the links of the list as part of the User object, one (set of) link(s) for each list to which the object might belong. An endomorphic strategy might be particularly attractive when, for example, it is known that each object can be on only one list at a time.

The problem of generating or modifying structures to take advantage of such implementations is orthogonal to the problem of using profile data to determine which implementation is best. Once the use of profile data is shown to be viable, then the same techniques can be applied to endomorphic implementations.

Programs that exhibit *phase behavior* are problematic. A user’s program could exhibit phase behavior by manipulating data one way early in the execution of the program (say, during the initial construction of an aggregate variable) and utilizing that data quite differently in later stages of the program (say, during access and modification of that aggregate variable). An interesting problem is the detection of the behavior and the optimal points for changing the implementation of the ADT from one representation suitable for the first phase into another representation more suitable for the later phase. A general solution to this problem is beyond the scope of our work. In our model, each static instance of a variable will have exactly one implementation. We can approximate some of the advantages of phase behavior

detection by limiting conversion from one representation to another at the assignment of a variable. Consider the following:

```

Set(int) A;
:
: (section one)
:
-- end of phase one of program
:
: (section two)
:

```

It may very well be the case that the behavior of the program in section one demands that the variable *A* be implemented as a singly-linked list, whereas the behavior of section two would be more efficient if *A* were implemented as a doubly-linked list. In our current model, *A* will be assigned only one implementation that will minimize the cost of running the program. If the program looked like the following:

```

Set(int) A;
Set(int) B;
:
: (section one uses A)
:
-- end of phase one of program
B = A;
:
: (section two uses B)
:

```

then we can look for the possibility of converting representations when *B* is assigned, and releasing resources used by *A*. This would, of course, require live-dead analysis, something that is currently beyond the prototype.

Another complication is introduced by user-defined functions. Consider the following:

```

main() {
  Set(int) A;
  Set(int) B;
  :
  user_fcn(A,B); -- call site 1
  :
  user_fcn(B,A); -- call site 2
  :
}

void user_fcn(Set(int) X, Set(int) Y) {
  :
}

```

Under our model, A and B must have the same implementation since X can have only one implementation (the same is true for Y). It may very well be the case that the program would perform better if there were two user functions, each with a different type signature. However, creation of multiple copies of the user function are beyond the capabilities of the prototype. At any rate, the problem is again orthogonal to the problem of using profile data, so we apply the general rule stated above for variables to the signature of functions: each statically declared object in the User's program will have exactly one representation associated with it during the running of the program.

Real-time applications can impose severe constraints on the behavior of a program, and are not considered further in this work. The work of Kenny and Lin [27] provides an important parallel to our work in the area of real-time control. In particular, we discuss later how their evaluation function generation techniques could be used in our system to improve the precision and portability of evaluation functions (see page 104).

4.2 Previous work

TYPESETTER is the first system to use a general technique for collecting ADT-specific profile data, and using that data to choose implementations. Almost all previous systems (with the exception of Low's) attempt to synthesize data representations: TYPESETTER chooses the implementation of a function to use, and thereby indirectly selects, but does not synthesize, the representations of the program's variables.

Low did the original work on implementation selection using profile data [31,32,33]. His system attempted to provide implementations for three abstractions: sets, lists, and a ternary relation which is unique to the base language SAIL. Each

```

partition the variables and expressions into equivalence classes (eq.c.);
determine which operations are performed on which eq.c.;
for each eq.c do
  S ← all representations ;
  remove from S all representations which are not feasible;
    -- may not have sufficient information at compile time
    -- may require an operation not implemented in a rep.
  predict time and space requirements for all s ∈ S
  for all s1, s2 ∈ S
    if s1 requires both more time and space than s2 then
      remove s1 from S
    endifor
  rank remaining representations in S by likelihood of being
  the best representation; -- uses a cost fcn
  use a hill-climbing heuristic to finalize implementations

```

Figure 4.1: Low's algorithm

ADT had several implementations that could be used to implement the User's variables. Each of the functions making up the interface of an ADT was written in assembler, and had associated with it an evaluation function that, given a frequency of execution and an aggregate size (e.g., the number of elements in a set), would return an estimate of the cost of using the interface function. His evaluation functions returned an estimate of the number of machine cycles and bytes required on any one invocation of a function.

The system required four passes over a program, with human interaction as one of the passes. The first step ran the subject program (using default implementations for the abstractions) with software monitoring inserted to collect a profile of the performance of the program in terms of statement counts. The system then prompted the user for information too difficult or impossible to derive analytically (e.g. "What is the average size of set *foo*?"). A penultimate static analysis pass computed the possible contents of variables in terms of other variables. This had the side-effect of partitioning the variables of program into equivalence classes; each equivalence class identified the variables that had to have the same implementation as all other variables in the class.

Low's algorithm for selecting representations (Figure 4.1) uses call sites solely for feasibility testing. Once a set of feasible assignments have been established, then an initial set of implementations are assigned. The final heuristic (Rowe called it incremental search) continually 'perturbs' the assignment of implementations by making changes that seem likely to improve the overall performance of the program

and keeping only those changes that do.

Low says that his system could not take into account certain features of set insertion (e.g., elements are always inserted in a specific order) and that he thinks it would be hard to include. `TYPESETTER` allows collecting this kind of information (i.e. ‘the data is always added in increasing/decreasing/non-increasing/non-decreasing/extremal order) as well as other information (e.g., sizes of sets, average length of lists, etc.).

Low’s system did not allow operators to work on multiple representations: a union operator’s two operands had to have the same representation. In our approach, a particular implementation function can be assigned to a call site if the actual parameters at the call site can be assigned the types of the implementation function’s formal parameters. It is up to the Implementor(s) which of these mixed-representation functions to implement.

`TYPESETTER` explores several aspects of Low’s general technique. Low’s ADT interface functions were written in assembler so he could make the evaluation of an invocation of one of those functions as precise as possible. He did not want to tackle the problem of writing evaluation functions for compiler-generated code. Given that precision is lost in any estimation of future performance of a real program, and that the performance of a function depends on more than just its frequency of execution and the size of the aggregate-type object, `TYPESETTER`’s evaluation functions accept inexactness as inevitable, and assume that programs satisfying the 90-10 rule are skewed enough to make the loss of precision irrelevant to the final decisions.

Also, `TYPESETTER`’s interface *and* evaluation functions are all in a higher-level language (C++), and the evaluation functions are in terms of this language’s constructs. That is, whereas Low’s system required the Implementor to count cycles in instructions in order to write an evaluation function, `TYPESETTER`’s evaluation functions are written in terms of the high-level language’s constructs. `TYPESETTER` has not solved the problem of providing an evaluation of compiler-generated code, rather it finesses the whole problem by admitting up front that evaluation is inexact. Precision is not possible *a priori* with compiled high-level language code (e.g., a different compiler’s optimizer will produce different code), but in exchange we get portable, understandable, easily tuned code, both for the implementations of the interface functions, but also for the evaluation functions.

And finally, Low’s system concentrated on finding implementations for a program’s variables, using program structure solely to determine the feasibility of the various implementations. `TYPESETTER` turns that around, and concentrates on finding implementations for the interface functions, and lets that specify what the implementation of the variables must be. Section 4.4.1 below explains this method of implementation selection in detail.

The other major work relevant to `TYPESETTER` is Hansen’s work on adaptive compilation, which we already discussed in some detail in Chapter 1.1. Hansen’s

work is the primary justification for this research in compiler utilization of profile data, though Hansen's system was targeted toward 'one-shot' compilation.

Other work relevant to what we are doing has concentrated primarily on how to assign implementations analytically. That is, given some program specification, most work has concentrated on finding means for determining implementations solely on the basis of that specification.

Barstow's PECOS system [3,24] processes a program specification via a knowledge database of program transformation rules. The program is declaratively, as opposed to imperatively, specified. Kant's LIBRA system [24] extended PECOS and attempted to apply the same rule-based, synthesizing approach to performance prediction. The knowledge database was enhanced with rules about estimating potential performance of partially constructed programs. Although LIBRA could allow the use of profile data, it was neither integral nor essential to the approach.

TYPESSETTER does not attempt to synthesize programs analytically, nor does it attempt to work with program synthesis at as high a level as does PECOS. Rather, my goal was to explore the possibility of providing implementation selection in the context of modern day compilers. Rather than seek a Copernican revolution and invent a totally new language in which to specify programs, I sought a more evolutionary approach to give existing languages and systems as much capability as possible.

Ramirez [37] used zero-one integer programming to assign implementations. His approach required condensing the behavior of a program into two matrices $s(i, j)$ and $t(i, j)$ where s is the estimated storage space consumed by implementation j when used to implement variable (substructure, he calls it) i , and t is the corresponding time estimate. His claim that the behavior of an implementation of an ADT can be summarized by two numbers $s(i, j)$ and $t(i, j)$ is highly suspect. It ignores, for example, how the behavior of a function or operator may change when provided with arguments of differing implementations. That is, he assumes that if implementation j is assigned to variable i , then $t(i, j)$, the amount of time required by that assignment, is independent of any other assignments. This is almost never the case, particularly when operators can accept operands with differing implementations (e.g., a *union* of a set implemented as a bitmap with a set implemented as a linked list). TYPESSETTER moves the focus of evaluation functions from the variable to the implementation of the interface functions. This allows the interacting costs of assignments to be taken into account at the expense of losing the ability to use zero-one integer programming to achieve an optimal solution.

Work within the SETL project [9,43] derives representations from declarations in the language and from analysis; e.g., frequencies are estimated by an analysis of the program text. I know of no work using profile data in the synthesis of SETL programs.

The SETL optimizing compiler attempts to determine a good implementation of for the set and mapping abstractions in the language (there is only one

representation for tuples). The default representation for sets and maps uses hash tables. If the analysis can determine *bases* for the elements of the sets, or if the programmer declares elements to belong to specific bases, then other more efficient implementations are possible for subsets of the bases. A subset can be represented as a bit in the structures for the elements of the bases (if the bit is one, then the element belongs to that subset, if zero, then not). If all elements of a base set are assigned unique integers, then a subset can be implemented as a bit-map. Or a subset might be represented with a separate hash table of pointers into the base set.

Straub's *Taliere* system improves on the optimization phase of the SETL compiler by considering estimates of performance, including symbolic analysis of execution frequencies. However, since he does not utilize profile data, the User must answer questions¹ of the form *What is the average size of s*t in line 215?*; or even *What is the expected number of iterations in an average execution of the loop starting at line 1235?* Even worse examples of the kinds of dialogue the system forces on the User are questions about probabilities: *What is the probability of the CASE statement of line 1113 taking the alternative of line 1126?* It seems extremely doubtful to me that a User would know this information with any precision or confidence without profile data.

Weiss [51] worked on finding types of recursive SETL variables, and presented methods for implementing such structures. However, he does not worry about selection of 'best' implementations by numeric criteria.

Sherman's dissertation [44] presents a very comprehensive approach to the problem through language design. The primary contribution of his programming language Paragon is the idea that implementations are subclasses, or refinements, in the ADT hierarchy (with multiple inheritance). That is to say, an implementation is just a refinement of an abstract data type and is specified using the same notation as that used to specify the abstraction. Paragon is an ambitious system that attempts to solve many problems at once, including selection of a refinement of an ADT based solely on the program text. Presumably, profile data could be used, but he does not discuss this in any depth. In the Paragon model, the User (our terminology) is responsible for writing the complete evaluation function (Sherman calls it the *policy* procedure) that selects the implementations of the variables of the program. This puts the onus of selection on the wrong member of our programming duo. We have attempted to design a system that puts the onus of implementation evaluation on the Implementor, and selection of implementations for functions and variables on the system, not on the User.

Rowe's system [38] approached the problem from the direction of selecting an implementation from a description of the desired data relations and functionality. His modeling domain language is implementation independent and is used to search for implementations that satisfy the described relations and operations. In those cases where there does not exist an implementation satisfying the description, Rowe

¹The questions are taken from his dissertation.

investigated ways of generating an implementation. However, he did not investigate the use of profile data in his work.

4.3 TYPESETTER: The System

In the discussion of the TYPESETTER system, it is useful to distinguish between different modes of programming which I identify by reference to two different programmer groups: the Users and the Implementors. Since Users need not be as facile with the theory behind TYPESETTER as Implementors need to be, the distinction between them is not one simply of convenience of notation. Users write programs that utilize the abstractions provided by the system; the system selects from among the implementations installed by the Implementors.

Examples in the following sections are based on the existing prototype, described in greater detail starting in Section 4.4. Language enhancements required to support TYPESETTER are not extensive and the notations should be relatively transparent to anyone who has used an object-oriented programming language like C++.

4.3.1 Formalities

An abstract data type (ADT) is a set of *function signatures* indexed by function names. For each ADT, there is a set of representation types; we'll write $\mathcal{A} \mapsto \mathcal{R}$ to mean that ADT \mathcal{A} can be represented by representation \mathcal{R} ; $impl(\mathcal{A})$ is the set of possible representations of \mathcal{A} . A function signature has the form T_0, T_1, \dots, T_n , for $n \geq 0$ and types T_i . By convention, T_0 is the type returned by the function. For all signatures in ADT \mathcal{A} , the T_i are themselves ADTs.

For each abstract function in an ADT, there is a set of *implementation functions*. Like their abstract function counterparts, implementation functions consist of a name and a signature. But where the abstract functions' parameter types T_i are ADTs, the implementation functions' parameter types T'_i are representation types. Furthermore, $T_i \mapsto T'_i$ for all i .

Our task is to assign a representation type to each variable in a user's program, and hence an implementation function to each call site in the program.

Define a program to be a set of variables \mathcal{V} and function call sites \mathcal{C} . Each variable $v \in \mathcal{V}$ has been declared to have one of the ADTs in \mathcal{T} , $atype(v)$. Each function call site $c \in \mathcal{C}$ consists of the name of an abstract function, $absfcn(c)$, and a list of actual argument variables $actuals(c)$. When an implementation function is assigned as the implementation of the abstract function at a call site, then the implementation type of v_j , the j^{th} actual, is assigned to be T'_j , the j^{th} type in the signature of the implementation function; i.e., $itype(v_j) = T'_j$.

4.3.2 The ideal system

The prototype system is a subset of a larger vision. I have argued previously that profile data should be collected all the time, and that the overhead for doing so is quite minimal. Then the programming system would allow the user to develop a debugged and efficient program in the following steps.

During initial implementation and debugging, the User rarely needs to decide the implementation of many, if not most, of the abstractions used in the program. The TYPESETTER system will utilize the simple program execution counts described in chapter 2 to pick a reasonable implementation for the abstractions.

As the User's system is implemented, and its structure is tested with more complex, more complete, and perhaps larger sets of input data, the actual implementation of the abstractions becomes more important (if for no other reason than that debugging a slow program can be particularly aggravating). By this time, the User will have gained enough experience with the program that he can conjecture which implementations of the abstractions may be reasonably efficient for the program. This conjecturing is important only to the extent that it allows the User to determine what additional information might be helpful to supply about the abstractions: TYPESETTER will determine which implementations are actually best. This additional information is supplied as part of the declarations of the variables, and is discussed in more detail in section 4.3.5.

It is also true that simple execution counts are insufficient for selecting an implementation. Alternative implementations of abstractions are created by programmers to take advantage of the interaction of properties of the abstractions with specific properties of sets of input data. Therefore, to determine that a bit-mapped implementation of a set is preferred over a linked-list implementation requires knowing not only how the program makes use of the data (e.g., number of insertions *vs.* number of deletions; the mix of element access and destructive operations; etc.) but also requires some information about the input data itself (e.g., is it read in increasing/decreasing order; is it 'sparse'; is it locally dense; etc.). This information can only be gathered directly and intentionally, and cannot be inferred from execution counts except at great expense, if at all.

At this point in the development process, the abstractions in the program are assigned *profiling implementations*, and a couple of runs of the instrumented program (over whatever data it can handle at this stage of development) will provide further data upon which TYPESETTER can assign more efficient implementations. Once implementations are determined based on this data, development and debugging of the program can proceed using this more appropriate implementation of its abstractions.

The TYPESETTER prototype described here has concentrated on implementing only the intermediate step: using profiling implementations of abstractions to collect abstraction-specific profile data that TYPESETTER can use to select from among a set of implementations. Lacking any profile information, the prototype al-

ways links in the profiling implementations for all abstractions; running the program then generates profile data for selecting more efficient implementations on a future run of the system.

4.3.3 ADTs

Among the ADTs available to programmers in `TYPESETTER` there are the usual built-in data types (e.g. integer, real, structures, arrays, pointers, etc.), each of which has a fixed implementation in User's programs, and a library of more complex abstract data types (ADTs) with implementations of functions that make up their interface tailored to specific representations. Users, however, make no reference to specific implementations of the functions or variables; they simply make use of the publicly declared interface (or protocol) of the ADT. The compiler will then choose implementations for the variables and functions to minimize a cost function based on data collected by a profiling version of the ADT. `TYPESETTER` supplies three ADTs: sets, lists, and maps. (These correspond closely to the SETL data types of finite sets, tuples, and maps [43].) Their definitions below are in `TYPESETTER`'s C++ dialect; in particular the first argument to a function is understood to be a pointer to the object by which the function is invoked.

Sets The Set abstract data type is generic in the type of the contained objects, which is denoted as `Any`. Figure ?? contains the definition of the Set ADT, and Figure 4.2 lists some possible implementations of exomorphic sets. Those with asterisks are currently implemented in the prototype. There are many possible representations of Sets, a few of which are briefly described in Figure 4.2. (The reader may wish to compare this list with the implementations provided by the SETL compiler; see pg. 57.)

Lists The List abstraction is generic in the type of the contained objects. Figure 4.3 lists the functions comprising the interface to the ADT. There are many possible representations of Lists, a few of which are described briefly in Figure 4.4.

Maps Maps, or finite functions, are generic in the type of the domain element and the type of the range element. Figure 4.5 lists the functions forming the interface to the ADT, and Figure 4.6 lists possible implementations. Figure 4.6 lists a few of the many possible implementations of Maps.

***Linked list** .

Doubly-linked list .

***Sorted linked list**: A linked list on which the elements are maintained in sorted order.

***Bit vector**: Elements must be, or map to, integers. Requires knowing max and min integers. Three variations: noElements kept as part of the set, fast array lookup element count for bytes, ditto for words.

Hash table: Useful when the key is not an integer, but an arbitrary collection of bits. Information about the range of the hash function and the density of the resulting hash values would help select good parameters for the hash table. For sets of arbitrary objects, the programmer must supply a hash function.

Sorted array: Keeps a sorted list of the actual elements of a set. Requires knowing max and min elements; knowing the maximum size of a set, and the average size of sets may help select better parameters for the implementation.

Sorted variable length array: Ditto. Requires extra overhead for the dope vector.

Linked array: Requires knowing max and min elements, as well as the fact that the elements tend to cluster. Optimizes space at the expense of time. To be used in environments where reallocating sets due to growth or memory compaction may be more expensive than just chasing pointers.

Figure 4.2: Possible implementations of sets (* in prototype)

`List()`: The constructor of List.

`boolean empty()`: Returns `true` if the list is empty.

`void makeEmpty()`: The list is emptied.

`boolean in(Any elt)`: Returns `true` if `elt` is in the list.

`int cardinality()`: Returns the size of the list.

`void rest(List L)`: Removes the first element of the list.

`Any first(List L)`: Returns the first element on the list.

`append(Any e)`: The element `e` is appended to the list.

`prepend(Any e)`: The element is pushed onto the front of the list.

`delete(Any e)`: All instances of the element `e` are removed from the list.

`sort(CmpFcn f(Any,Any))`: The list is sorted in place using the comparison function `f`.

`iterInit(Iterator i)`: Initialize an iterator over the list.

`iterate(Iterator i, Any &elt)`: Assign `elt` the next element of the list in the iteration and return `true`, else return `false`.

`iterDone(Iterator i)`: Return `true` if an invocation of `iterate` would return `false`.

`iterCleanup(Iterator i)`: Return resources allocated to the iterator.

`iterCopy(Iterator i, Iterator &j)`: The iterator `i` is copied into a new iterator `j`.

Figure 4.3: Specification of List ADT

***Linked list.**

Doubly-linked list.

Fixed-length array: Each list has an array of maximum size allocated for it.

Linked array: The list is kept in a list of arrays, each sub-array allocated/deallocated as the list is manipulated. Requires knowing max and min elements, as well as the fact that the elements tend to cluster. Optimizes space at the expense of time. To be used in environments where reallocating lists due to growth or memory compaction may be more expensive than just chasing pointers.

Figure 4.4: Possible implementations of lists (* in prototype)

`Map()`: The constructor of `Map`.
`boolean empty()`: Returns `true` if the map is empty.
`void makeEmpty()`: The map is emptied.
`boolean inDomain(Any elt)`: Returns `true` if e is in the domain of the map.
`boolean inRange(Any elt)`: Returns `true` if e is in the range of the map.
`int cardinality()`: Returns the size of the map the number of elements defined in the range).
`define(Any d, Any r)`: The element d is added to the domain of the map so that it returns the element r .
`delete(Any e)`: All instances of the element e in the domain are removed from the map.
`sort(CmpFcn f(Any,Any))`: The map is sorted in place using the comparison function f .
`iterInit(Iterator i)`: Initialize an iterator over the map.
`iterate(Iterator i, Any &d, Any &r)`: Assign d the next element in the domain of the map, r the corresponding element in the range, and return `true`; else return `false`.
`iterDone(Iterator i)`: Return `true` if an invocation of `iterate` would return `false`.
`iterCleanup(Iterator i)`: Return resources allocated to the iterator.
`iterCopy(Iterator i, Iterator &j)`: The iterator i is copied into a new iterator j .

Figure 4.5: Specification of Map ADT

***Linked map:** The map may be singly or doubly linked, and consists of pairs of elements as related by calls to *define*.

Fixed-length array: Each map has an array of maximum size allocated for it, if a maximum size is known. The array is two dimensional, one each for the domain and range.

Linked array: The map is kept in a map of arrays, each sub-array allocated/deallocated as the map is manipulated.

Hash table: Hashed by domain elements for faster lookup.

Binary tree: So seeks on domain elements are $O(\log n)$, for n the number of elements in the domain.

Figure 4.6: Possible implementations of maps (* in prototype)

4.3.4 Iterators

In a complete language definition, much of the functional interface for iteration would be hidden from the Users with some syntactic sugar. I envision something close to the Alghard paradigm for iterators, and the User would write something like the following to have the compiler invoke the appropriate iteration functions.

```

Set(SomeType) S;
...
for i in S do
  ... -- use of element i
endfor;

```

The above code would be translated into something like the following using the TYPESETTER iterator paradigm:

```

Set(SomeType) S;
...
Iterator S_iter;
S.iterInit(S_iter);
while S.iterate(S_iter,i) do
  ...
  endwhile;
S.iterCleanup(S_iter);

```

The iterators are associated with the object being iterated. Their exact form is never available at the User level and, therefore, the ADT implementation is free to create the iterator object necessary to successfully traverse the aggregate type. In other words, every ADT exports the `Iterator` type, but not the internals of the `Iterator` type.

4.3.5 Optional parameters

Users should be able to write simple declarations of their program variables and have the system select an appropriate implementation of those variables based on that declaration and on knowledge of the behavior of the program containing those declarations.

```

Set(Bar) foo;

```

declares that the variable *foo* contains a set of objects of type *Bar*. TYPESETTER recognizes several objects in this declaration. The first is, of course, the use of *foo*, which is the name of the variable being declared, and the second is the name of the ADT. The remaining parameters supply information to the ADT itself, and are of two types: those *required* for even a minimal implementation of the ADT, and

optional parameters which supply further information that may enable (or exclude) implementations for consideration. Required parameters are positional, and optional parameters are named. In our example, *Bar* is the name of a User type and is required by all implementations of Sets. A declaration of a variable of type Map has two required parameters, the type of the *key*, and the type of the *data*. For example,

```
Map( T1, T2) vmap;
```

declares *vmap* to be a mapping from key objects of type *T1* to data objects of type *T2*.

Use of the variables follows standard object-oriented form. For example,

```
iset.add(3);
```

adds the element '3' to the set named *iset*.

Optional parameters (or just *optionals*) are not required for an implementation to be assigned to a variable: there is always at least one implementation of an ADT that can be assigned to any variable of that type. Optionals supply information that allow the system to consider other, possibly more efficient, implementations for a variable. For instance, consider:

```
Set(int) iset;
```

As declared, the variable *iset* could be implemented with one of a variety of bit-map implementations, but with only those implementations that can handle bitmaps of unknown and possibly varying size, and perhaps even negative values as elements. This implies a relatively complex implementation of bitmaps. If the User were aware that the only correct integer values for this set were positive, that information could be provided with the optional parameter *lowerb*:

```
Set(int, lowerb=0) iset;
```

The optional parameters are named parameters. The additional information they provide could enable the consideration of other bit-mapped implementations that need it. Obviously, the more information provided about the user's objects, the more implementations that can be considered.

Consider the following declaration:

```
Set(Utype) uset;
```

The User has declared a set of objects of type *Utype*, a user-declared type. In this situation, bit-mapped implementations are not at all feasible since the compiler has no way of mapping objects of type *Utype* onto integers. If such a mapping is possible, the User may declare the mapping function and its inverse:

```
Set(Utype, objToInt=f, intToObj=g) uset;
```

The system is now free to consider bit-mapped implementations as before.

In an ideal system, the User might be given hints as to which optionals might provide performance improvements. How such hints might be automatically generated is a topic for further research. In the `TYPESETTER` system, the Implementor is responsible for providing the documentation for the optionals supported by an implementation. That documentation will include a general description of the possible effects on `TYPESETTER`'s ultimate choices of implementations. That is, we depend on documentation to help Users decide which optionals might be beneficial for their programs.

The actual number of optionals required for the prototype has been few. Table 4.7 describes optionals envisioned as useful. Asterisks mark the optionals actually implemented in the prototype. The optionals in the table fall into one of two categories: those that provide information that would be difficult to derive from the program even with profile data, and those that are a convenience for the current implementation but could be eliminated with appropriate programming by the Implementor. For instance, the *upperb* optional for set implementations cannot in general be derived from program source and profile data.

The declaration optional *addedDecreasing* could be detected at run time and encoded in the profile data. If different input data were fed to an implementation that attempted a more efficient representation by assuming the elements were added in order, more than likely the User's program would run slower, but would not fail. Of course, it is possible to design an implementation whose correct operation depended on the assumption, in which case detection by a profiling implementation would not be sufficient: a contract with the user in the form of a declaration would be required. The 'order' of the objects in this example is an internal ordering; if the implementation depends on a User-defined ordering, then another optional declaring the order function would need to be defined by the Implementor and declared by the User. Specifically, the implementation *Set_slistord* takes advantage of the fact that sometimes a program creates all the objects that are in a set, and adds them in the order they are created. This often results in the heap allocator allocating the objects such that their memory addresses correlate with their time of creation. *Set_slistord* attempts to cut down on lookup time of elements in a set by keeping the elements on a list in increasing order of the addresses of the objects.

4.3.6 Alternative implementations

The User sees a system complete with a set of possible implementations of abstract data types. These implementations were provided with the system and are, conceptually at least, part of the system. It must be possible for Implementors to specify under what conditions their implementations can be selected, what profile data needs to be collected, and how that data is to be evaluated.

One of the first responsibilities of the first Implementor of an abstraction is

Optional	value	description
Sets		
* IntToObj	function	the function that converts integers into objects of the correct type for this set
* ObjToInt	function	converts objects into integers
* lowerb	integer	the lower-bound value of the elements of a set of integers
* upperb	integer	the upper-bound value of the elements of a set of integers
nofElts	integer	the number of base elements of the set; must equal upperb−lowerb+1, if they are specified
* ObjsAreInts	(none)	declares that the base objects of this set are in fact integral, and the compiler will perform the correct coercions; this is a convenience optional for the prototype
compareFcn	function	accepts pointers to two objects and returns −1, 0, or 1 depending on whether the first is less than, equal to, or greater than the second.
addedDecreasing	(none)	elements are added in decreasing order
addedIncreasing	(none)	elements are added in increasing order
Lists		
maxLength	integer	User contracts that list will never be longer than this
Maps		
IntToObj	function	the function that converts integers into objects of the correct type for this set
ObjToInt	function	converts objects into integers
lowerb	integer	the lower-bound value of the elements of a set of integers
upperb	integer	the upper-bound value of the elements of a set of integers
hashFcn	function	returns a 32-bit integer that can be used in a hash table implementation of a map
compareFcn	function	accepts pointers to two objects and returns −1, 0, or 1 depending on whether the first is less than, equal to, or greater than the second.

Figure 4.7: TYPESETTER optionals.

```

function Set_P::add(Set this, any e)
{
  profiler
    pcnt, psizeSum, pwasIn;
  Link lp;
  pcnt++;
  psizeSum += length;
  lp = this→first;
  while (lp != nil && e != lp→data) {
    lp = lp→next;
  }
  if (lp == nil) {
    // e not in the set
    Link newp = new Link;
    newp→data = e;
    newp→next = first;
    first = newp;
  }
  else {
    pwasIn++;
  }
}

```

Figure 4.8: Profiling implementation of *add*

to define the functionality of the abstract data type and provide its first implementation. This first implementation must also be the profiling implementation for this ADT, and it must be the most general: it is a requirement in `TYPESETTER` that no ADT is defined with functions in the interface that cannot be implemented in the profiling implementation; otherwise, there would be no way to collect profile data about that function.

Figure 4.8 shows the code for the profiling implementation of the add-an-element function in the interface for sets. This implementation, called *Set_P*, uses a very general structure (in this case a linked list) to ensure that any function in the interface can somehow be implemented and profiled. (The converse is not true: an alternative implementation does not have to implement every function in the interface. Any program using that function could not have that implementation assigned to the involved variables, however.)

Profile variables (declared as **profilers** in Figure 4.8) are allocated per call site in the User's program. That is, if the User's program calls *add* from three distinct

```

function Set_bm::add(Set_bm this, any e)
{
    int i = (*this→objToInt)(e);
    int w = (i / (sizeof(integer)*sizeof(byte)));
    int b = (i mod (sizeof(integer)*sizeof(byte)));
    this→setbits[w] |= (1 << b);
}

```

Figure 4.9: An alternative implementation of *add*

sites, then a total of three instances each of *pcnt*, *psizeSum*, and *pwasIn* are allocated. On each call of the *add* function, the invocation counter *pcnt* is incremented, and the *psizeSum* profile variable is incremented by the current length of the set. From this information, evaluation functions can compute the average size of the set per call per call site. Finally, it may be useful to an implementation to know how many times *add* was invoked to add an element that was already a member: the profiling variable *pwasIn* allows us to compute that statistic. As other implementations are added to the collection of implementations for sets, more information may need to be collected by the profiling implementation. The Implementor of an implementation that requires new profile data will modify the profiling implementation to collect it.

4.3.7 Feasibility functions

Continuing with the example of the add-an-element function, Figure 4.9 shows its implementation when sets are implemented as a bit map. Before a User’s variable can be assigned *Set_bm* as its implementation, TYPESETTER must first check that this is a feasible assignment. Therefore, the Implementor of *Set_bm* must provide a *feasibility function* that TYPESETTER can call to check feasibility. The function returns either **true** or **false**; the feasibility function for a bit-map implementation having 32 elements or less is shown in Figure 4.10.

4.3.8 Evaluation functions

Traditional profiling techniques cannot capture the wealth of detail required for intelligent selection of implementations. For instance, from knowing the number of times allocation and deallocation routines are executed, it is extremely difficult to deduce the average size of sets, say, at any particular call site in a program. With our profiling schema, it is particularly easy. Furthermore, rather complex information can be acquired such as “Is this a sparse set?”, “Does this list ever have items deleted from it?”, “Are the elements of this set entered in any particular order that yields advantage to any implementation?”, etc.

FEASIBILITY

```

{
  if (!u→IntToObj→def && !u→ObjToInt→def    &&
      u→ObjsAreInts→def &&
      u→upperb→def    && u→lowerb→def    &&
      u→upperb→ival > 0 &&
      u→lowerb→ival <= 0 &&
      u→lowerb→ival < u→upperb→ival &&
      (u→upperb→ival - u→lowerb→ival) < 32
  ) return true;
  else return false;
}

```

Figure 4.10: Feasibility function for implementation *Set_bm***Eval** *Set_bm::add*(*CallSite c*)

```

{
  return c.pcnt *
    (idividePwr2_op + modPwr2_op + orAssign_op + array_op + shift_op);
}

```

Figure 4.11: Evaluation function for *add*

TYPESETTER determines which implementation of an ADT is best based on the estimates returned by the evaluation functions supplied with each implementation. For each function in the interface of an ADT, the Implementor must supply an **Eval** function. For example, the evaluation function for *Set_bm::add* is in Figure 4.11. When called with a call site as a parameter, these functions return an estimate of the runtime resources required by this implementation.

The variable *pcnt*, declared in the profiling implementation as a profile variable, is used here to estimate how much time this implementation of *add* would take at a particular call site. The other profiler variables are not used for evaluating this implementation of *add*. The variables in Figure 4.11 ending in ‘*_op*’ are constants that estimate the relative execution times of each of the indicated high-level language operations. These times will in general be approximate, if only because the execution time of any construct is context-dependent. However, the purpose of these constants is merely to provide an estimate of the time required by this function *relative* to other functions implementing the same functionality for different representations. It is my assertion that such an evaluation technique is ‘close enough’ to allow reasonably

‘correct’ assignment of implementations. It pays for itself by allowing the library of implementations to be ported to a new system by simply changing the values of the indicated constants, if necessary.

Therefore, all of the examples here, and all of the implementations in the prototype, estimate the amount of time used in a direct fashion (assuming infinite physical memory, etc.), and do not attempt to handle the complexities of more precise estimates of performance.

Not all functions are as straightforward as the *add* example above. Consider the function for computing the union of two sets. The `TYPESETTER` code for the `Set_slist` implementation of the two-operand union function (i.e. the union of the two sets is assigned to one of the sets) is shown in Figure 4.12. This implementation of *union* determines for each element in the set *sB* if it is in the current set (the *this* set, in C++ terminology). If not, it is added to the current set. The evaluation function for this function must therefore have access to the evaluation function for the *in* function, passing to it the information it needs to create an estimate of its behavior at the call site in the *union* function. The profiling implementation may not have invoked the *in* function (and in the prototype’s profiling implementation for sets, it doesn’t) so there is in general no profile data specific to call sites contained in the alternative (i.e., non-profiling) implementations of an abstraction’s interface functions. Therefore, in this case, the *union* evaluation function must provide estimates for the profiling data, which it passes as parameters to the evaluation function for *in*.

4.4 TYPESETTER: The Implementation

I use the name `TYPESETTER` to refer to the whole system and to the language that results from the enhancements made to C++. The actual implementation of `TYPESETTER` consists of several parts, some of which are written in `TYPESETTER`, the language. The extensions to C++ have been discussed already, and are straightforward. There is the addition of evaluation and feasibility functions to the declaration of class member functions (see pg. 71ff), the use of profiling variables (see pg. 70ff), and the User’s ability to optionally declare extra information in a variable declaration (see pg. 67ff).

All of these language features currently take the form of macro invocations. A macro processor first scans all source files and extracts the relevant information via the macro invocations. This information is then fed to the analysis program which makes the actual implementation decisions. The macros are written in `m5` [40], a powerful macro language designed for the manipulation of name-scoped text, such as is found in programming language text. In addition to the alternative implementations, the program that makes the actual implementation decisions is also written in `TYPESETTER`; I call this program `THERBLIG` after Frank Gilbreth’s qualitative unit of work-motion [14]. Figure 4.13 shows the steps necessary to compile

```

void
FUNCTION(union1)(Set_slist sB)
{
    Set_slist_Link *lp = sB.first;
    while (lp != NULL) {
        if (!Set__in(lp->data)) prepend(lp->data);
        lp = lp->next;
    }
}
EVALSUB(Pcnt,Psza,Pszb,Povrlp)
{
    @@ executed Pcnt times;
    @@ there were Povrlp elements of s2 already in s1;
    @@ each time loop exec'd avgBsz times;
    @@ and prepend exec'd avgNotIn times;
    if (Pcnt == 0) return 0;
    double avgBsz = Pszb / Pcnt;
    double avgIn = Povrlp / Pcnt;
    double avgNotIn = avgBsz - avgIn;
    return Pcnt *
        (assign_op // startup
         + (avgBsz * (cmpZero_op + deref_op + assign_op + not_op))
         + EVALSUBfor(in)(avgBsz, Psza, avgNotIn)
         + (avgNotIn * EVALSUBfor(prepend)()));
}
EVALUATE
{
    return EVALSUB(p_cnt,p_sza,p_szb,p_ovrlp);
}
END_FUNCTION(union1)

```

Figure 4.12: Set union using linked lists

a program written in `TYPESETTER` into C++. The first invocation of `m5` collects information about the declaration of variables and the call sites of abstract functions, and emits a description of the User's program. `THERBLIG` analyzes this description, along with profile data about the program (if it exists), and emits an assignment of implementations for all variables and call sites. These assignments, along with a library of ADT implementation sources and the original source of the User program, are processed by another run of `m5` to transform the User program into C++.

`THERBLIG` contains all of the functions necessary to evaluate the use of ADT implementations in the User's program, including the evaluation and feasibility functions provided by the Implementors for their contributions to the library of implementations. As can be seen in Figure 4.14, these are compiled as part of `THERBLIG` after processing of the implementations declared for each of the ADTs that are part of the system.

These steps are discussed in more detail in the following sections. First, I will discuss the algorithm used in `THERBLIG` to assign implementations to variables. Next, I will present the mechanism used in `TYPESETTER` for implementing code sharing among the implementations assigned to the User's variables.

4.4.1 The Implementation Selection Algorithm

In the concluding chapters of his dissertation, Low [31] observed that his hill-climbing heuristic seemed to display the property that implementation decisions were made early and were rarely re-made. Straub [46] notes that even though his representation selection algorithm was run many times with widely varying expected values for the program's variables, "the choice of data structures made by the system tended to be independent of the responses to the queries made to the user." He concluded that this indicates that the selection depends much more on the operations performed than on the expected values of the program variables. (It is curious that he does not even consider collecting profile data as a better source of information than interactive querying of the user. Nor does he consider profile data as a source for finding those important operations.)

Apparently, the representation selection process is being made much more complicated than it really is: good selections depend on a small part of the information that has been utilized in previous research. My hypothesis is that this is due to the 90-10 nature of most programs. That is, if more weight is given to the more frequently executed sections of code (as they are in Low's technique) then the implementation costs of these sections will dominate the overall execution costs of the program. It also means that implementation decisions that are good for these 'hot spots' will be good for the whole program, and, conversely, implementation decisions for very infrequently executed code have little effect on the performance of the program.

This observation, bolstered by the observations of previous researchers, plus

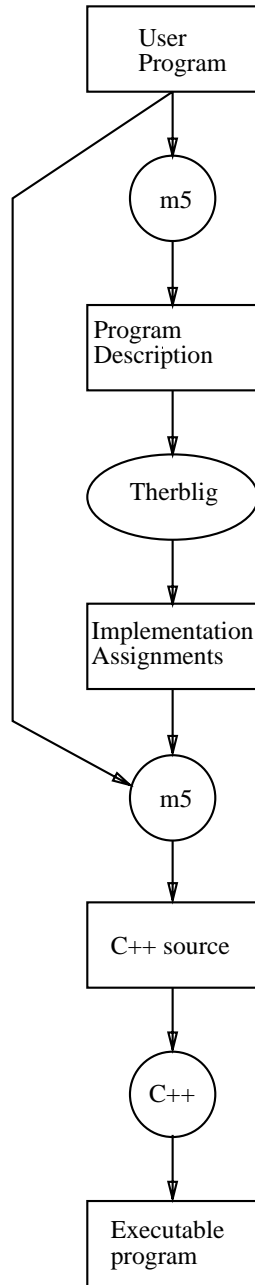


Figure 4.13: Steps to process a User program

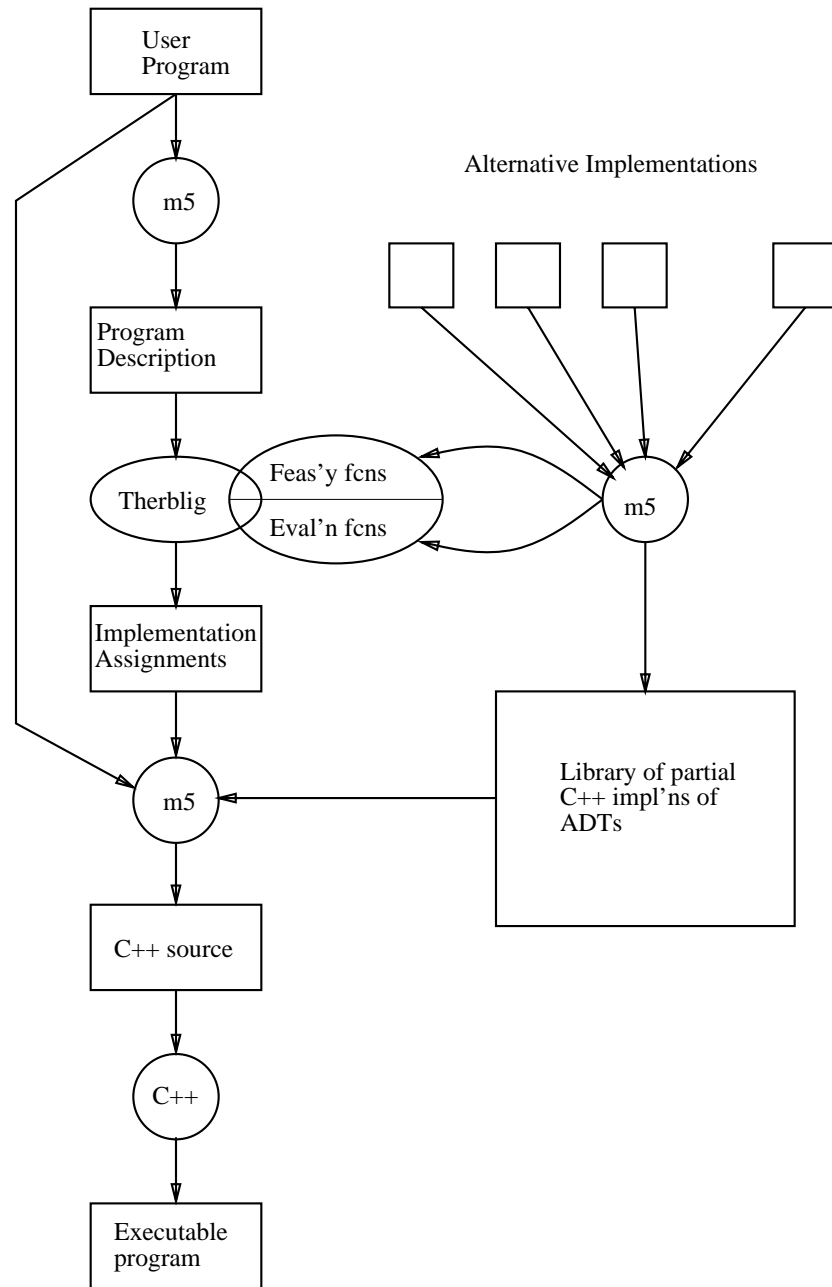


Figure 4.14: Steps to process implementations and build THERBLIG

the success of the Greedy Sewing algorithm for code reorganization (covered in Chapter 3), plus Hansen's success with focused optimization [17], plus the time-honored success of hand-tailored optimizations based on profile data all suggest that in general a greedy algorithm using profile data will work for assigning implementations to abstract data types. This results in a two stage process: stage one decides which of the existing implementations are feasible, and stage two chooses from among the feasible implementations the one that minimizes the cost of running the program. Furthermore, the algorithm does *not* require analysis of the program flow graph, as does Low's, Rowe's, and Straub's techniques: it reduces to a problem of matching function implementations with function call sites.

Three questions must be answered to find an implementation for a variable.

1. Does enough information exist to take advantage of a representation? That is, does the information exist that would allow variable v to be implemented with implementation R ? For example, if a list is implemented as a fixed size array, we have to know the maximum size of the list. If the maximum size of a list is unknown, or there is no maximum size, then the fixed-size array implementation of a list is not feasible.
2. Does enough of a representation exist to implement a variable? That is, if variable v is given implementation R , for each call site c which has v in its actual parameter list, does there exist an implementation f of F_c consistent with that and all previous assignments of implementations?
3. Which implementation of variable v will provide the best performance for the program?

The first two questions come under the realm of *feasibility*: is it possible to select an implementation for the program? The last question seeks to find an implementation that minimizes the cost of executing the program, using the Implementor-provided evaluation functions as objective functions for that minimization.

The assignment algorithm which I have implemented is in Figures 4.15 through 4.19. The pseudo-code in is meant to be descriptive rather than a rigorous program in a well-defined language. A description of some functions that are not otherwise defined can be found in Figure 4.19. The code is not specific as to how the assignment of variables and functions are recorded. We assume a global data structure *assignment* contains a consistent assignment when *chooseImplementation* returns, or it is empty if no assignments were possible.

The basic procedure is summarized as: sort all call sites in decreasing order by some preliminary metric and assign implementations to variables based on the cheapest implementation of the functions called at each site. While I have emphasized the evaluation functions for alternative implementations, the profiling implementation for each ADT also has evaluation functions that are used to estimate

```

proc chooseImplementation(Set(callSite)  $\mathcal{C}$ )
  List(callSite)  $\mathcal{S}$ ;
   $\mathcal{S} = \text{sortByImportance}(\mathcal{C})$ ;
  if not assignable( $\mathcal{S}$ ) then error fi
endproc

```

Figure 4.15: Main routine for choosing representations

the potential impact of a call site. Consider a program that creates a set of n objects, sorts the set into a list, and then accesses the objects in the list. The call site (assume there is only one) that adds elements to the set is executed n times. The list is accessed kn times, for some integral k . But the sort routine is called exactly once. Nevertheless, that sort routine has complexity $O(n \log n)$ to $O(n^2)$, meaning it has the potential of swamping the significance of the add-an-element function. The evaluation functions of the profiling implementation return values that reflect this potential impact of each function.

The recursive selection algorithm is encoded primarily in the function *assignable* (Figure 4.16) which is passed a list of call sites; the first call site c on the list will be assigned an implementation, if possible. That is, at call site c where the abstract function F_c is invoked, *assignable* will pick the (next) cheapest implementation of F_c , which thereby determines the types of the variables that are passed as parameters to F_c . The function *findCompatible* (Figure 4.18) finds all implementations of F_c that are compatible with this call site. This set of compatible implementations is then sorted in increasing order of the estimated costs provided by the implementations' evaluation functions.

For each function in the sorted sequence of implementations \mathcal{S} , if the actual parameters to F can be assigned the implementation types required by the implementation function f , then an assignment is attempted on the next call site on list \mathcal{L} . Otherwise, we back out of any implementation assignments made in this invocation of *assignable*, and another implementation function f on \mathcal{S} is tried as the implementation function for this call site. If every implementation function has been tried and no consistent assignment of implementations to variables and call sites has been made, then *assignable* returns false. The function *parmsImplementable* (Figure 4.17) determines if the implementation function f can be used to implement the abstract function F by checking that the actual parameters to F can be assigned the implementation types required by the formal parameters of f .


```

function assignable(List(callSite)  $\mathcal{L}$ ) returns boolean:
  if isEmpty( $\mathcal{L}$ ) then return true;
   $F \leftarrow \text{absFcn}(\text{first}(\mathcal{L}))$ ;
   $S \leftarrow \text{sortByCost}(\text{findCompatible}(\text{first}(\mathcal{L})))$ ;
  foreach  $f$  on  $S$  do
    if not parmsImplementable( $f, F$ ) or
      not assignable(rest( $\mathcal{L}$ )) then --backtrack
        undoImplementations;
    else return true;
  fi
  endfor
return false;
endfcn;

```

Figure 4.16: Routine *assignable* for choosing the implementation of an abstract function

```

function parmsImplementable(impFunc  $f$ , absFunc  $F$ ) returns boolean:
  foreach  $v \in \text{signature}(F)$ , each  $t \in \text{signature}(f)$  do --parallel
    if not alreadyImpltd( $v$ ) then
      if implementable( $v, t$ ) then
        implement( $v, t$ );
      else --conflict; backtrack
        return false;
      fi
    fi
  endfor
return true;
endfcn

function implementable(variable  $v$ , implType  $t$ ) returns boolean:
  if not atype( $v$ )  $\mapsto t$  then return false;
  if not t.feasible( $v$ ) then return false;
  return true;
endfcn

```

Figure 4.17: Mapping parameters onto implemented functions' signature

```

function findCompatible(callSite c) returns Set(implType):
  I ← implementations(absFcn(c));
  Ic ← emptySet;
  foreach f ∈ I do
    if isCompatible(f, c) then
      Ic ← Ic ∪ {f};
    enddo
  return Ic;
endfcn

function isCompatible(absFunc f, callSite c) returns boolean:
  return if signature(f) ↦ actuals(c) then true else false;
endfcn;

```

Figure 4.18: Finding compatible function implementations

The function *implementable* checks that v can be assigned the implementation type t . The function call

$$t.\text{feasible}(v)$$

calls the Implementor-supplied feasibility function for the implementation t to verify that t is a feasible representation for v .

Classes of variables due to aliasing in user functions

User-declared functions require some twists on the algorithm as we have presented it so far. User-declared functions can result in equivalence classes of variables caused by aliasing of variables with formal parameters in the signature of these functions. In our example on page 50 if representation R is assigned to $S1$, then it also has to be assigned to $S2$. These equivalence classes are reminiscent of Low's equivalence classes; however, his classes were imposed by the design decision to not allow mixed representation functions in the ADT interfaces and were much more restrictive; the equivalence classes of aliased variable names are much less so.

Figure 4.20 contains the modifications to the algorithm to handle this complication. The prime difference between Figure 4.17 and Figure 4.20 is that the former looks only at the variable, while the latter looks at all variables that are equivalent to the variable because of parameter passing to calls on user functions. This acts as an implementation optimization to make sure implementation selections are compatible without having to backtrack.

<i>absFcn</i> (<i>callSite</i> <i>c</i>)	→ --abstract function invoked at call site <i>c</i> ;
<i>signature</i> (<i>absFunc</i> <i>f</i>)	→ --signature of abstract function;
<i>actuals</i> (<i>callSite</i> <i>c</i>)	→ --list of actual parameters at call site <i>c</i> ;
<i>sortByImportance</i> (Set (<i>callSite</i>) <i>C</i>)	→ --list of call sites sorted in decreasing --order of the results of the evaluation functions --in the profiling implementation;
<i>sortByCost</i> (Set (<i>impFunc</i>) <i>F</i>)	→ --list of implementation functions sorted in --increasing order of estimated execution cost;
<i>first</i> (List (<i>any</i>) <i>T</i>)	→ --first element of list;
<i>implement</i> (<i>variable</i> <i>v</i> , <i>implType</i> <i>t</i>)	→ --update implementation list contained --in global variable <i>assignment</i>
<i>assignImplType</i> (<i>variable</i> <i>v</i> , <i>implType</i> <i>t</i> , <i>signature</i> <i>s</i>)	→ --assign all occurrences of variable <i>v</i> --the <i>implType</i> <i>t</i> in the signature <i>s</i>
<i>assignImplType</i> (<i>variable</i> <i>v</i> , <i>implType</i> <i>t</i> , <i>signature</i> <i>s</i>)	→ --return the signature <i>s</i> with all --instances of <i>v</i> assigned implementation type <i>t</i>

Figure 4.19: Miscellaneous functions

```

proc parmsImplementable(impFunc f, absFunc F) returns boolean:
  foreach v ∈ signature(F), each t ∈ signature(f) do
    if not alreadyImpltd(v) then
      if implementable(w,t) ∀w ∈ class(v) then
        implement(w,t) ∀w ∈ class(v);
      else --conflict; backtrack
        return false;
      fi
    fi
  endfor
return true;
endproc

```

Figure 4.20: Mapping parameters onto implemented functions' signature with equivalence classes

```

class Set {
    void in(Any e);
    void add(Any e);
    void subset(Set s);
};

```

Figure 4.21: Declarations of generic Set functions

4.4.2 Code sharing

For each user variable declared to be, say, a set of some user-type *UType*, a naive implementation of the generic abstraction of sets would create a new copy of the implementation code for sets with all instances of the generic parameter replaced with *UType*. In general, this is quite unnecessary, and particularly so in the context of `TYPESETTER` where all ADTs are restricted to be exomorphic and *UType* is restricted to be a pointer to a user type. In this case, code can be written once to handle sets of pointers to objects. In the case of sets of various sized integers, only one copy need be created for each size of integer. If pointers are the same size as one of the sizes of integers, the same version of Set code can be used for both.

However, it is important not to give up strong type checking to gain this savings in code space. Users should still be notified when their programs violate the declarations they themselves have made.

For instance, in Figure 4.21 are the generic declarations of some of the interface functions for sets. The abstraction of sets is generic in one type, that of the type of the elements of the sets. If the user declares types `Token` and `String`, for instance, and then declares variables to be `Set(Token)` and others to be `Set(String)`, it would be wasteful to have two implementations of the set functions, one for `Tokens` and one for `Strings`, since both are actually implemented as a set of pointers to tokens, and set of pointers to strings. All that is needed as an implementation of sets that can handle pointers.

All that is needed to maintain strong type-checking is the declaration of *coercion types* for a set of `Tokens` and a set of `Strings`, each of which calls the appropriate function for handling sets of pointers. Therefore, `TYPESETTER` generates one implementation of the set functions capable of handling pointers (and incidentally four-byte integers on many systems) and then generates a coercion class for maintaining strong type checking. The declarations for an implementation of ‘set of pointers’ are in Figure 4.22; Figure 4.23 contains the coercion implementations. This model of typing and implementation of generic functions builds on the *template* idea first proposed by Stroustrup [10] for C++ but is much more powerful in that it allows the Implementor more control over how much new source code is generated. Such control cannot be created easily in C++ without extending the language fur-

```

class Set_ptr {
    void    in(void* e);
    void    add(void* e);
    boolean subset(Set_ptr& s);
};

```

Figure 4.22: Declarations of generic Set functions

```

class Set-Token : public Set_ptr {
    void    in(Token* e) { Set_ptr::in((void*) e); }
    void    add(Token* e){ Set_ptr::add((void*) e); }
    boolean subset(Set-Token& s) { return Set_ptr::subset((Set_ptr&)s); }
};
class Set-String : public Set_ptr {
    void    in(String* e) { Set_ptr::in((void*) e); }
    void    add(String* e){ Set_ptr::add((void*) e); }
    boolean subset(Set-String& s) { return Set_ptr::subset((Set_ptr&)s); }
};

```

Figure 4.23: Declarations of generic Set functions

ther than is suggested in the latest C++ language reference by Ellis and Stroustrup [10].

4.4.3 Refinements

In order to simplify the above discussion of THERBLIG, I have not included the enhancements added to the code that allows THERBLIG to search the solution space in a controlled manner. In the actual implementation (see Appendix B), using a mechanism very similar to that used in the Greedy Sewing Algorithm, THERBLIG can be invoked with a parameter p that specifies indirectly the portion of call sites that are ‘optimally’ assigned implementations; i.e., the search for their implementations is exhaustive, with every possible combination of implementations examined.

The set of call sites is sorted in decreasing order of the values returned by the profiling implementation’s evaluation functions. Let $S = \sum_i^n C(i)$, where $C(i)$ is the cost estimate returned for the function at the i^{th} location in the list. The sum of these values, S , is multiplied by the parameter p , a number between 0 and 1, to determine a cutoff point in the list of sorted call sites. For a call site at location i in the list of sorted call sites, it is above the cutoff point if $\sum_{j<i} C(j) < p * S$, and

it is below the cutoff point if $\sum_{j < i} C(j) \geq p * S$. At each point in the assignment algorithm, if the call site being considered is below the cutoff point, only the first consistent implementation is returned, and all others are ignored. If the call site is above the cutoff point, then each consistent implementation is examined to see if it improves the program's implementation.

Invoking THERBLIG with $p = .9$, all possible implementations are examined for the call sites that account for 90% of the estimated runtime resources. By setting $p = 1$, all possible implementations are examined. The results and methods of the algorithms as I have described them above are achieved by setting $p = 0$; i.e., THERBLIG returns the first consistent implementation for the program.

4.5 Examples

I have demonstrated that TYPESETTER code is not difficult to write, either for the User or for the Implementor. No specialized knowledge of compilers or profiling technology is required by either the User or the Implementor. The Implementor specifies the information required to make a reasonable implementation decision with normal-looking programming language statements; the only difference is that profiling variables are allocated per call site rather than per function. In an ideal implementation of TYPESETTER, the User would need only (1) to re-compile the system as directed by the system (although this cycling could certainly be automated), and (2) to be aware of the different kinds of optional information that may be specified for a data type (e.g., upper and lower bounds on elements of sets).

Given profile data and User declarations, TYPESETTER gives the User program a 'reasonable' implementation. I cannot claim that TYPESETTER constructs 'optimal' implementations: the whole process of software construction is too heuristic to allow such a claim. Future work can concentrate on determining exactly how 'optimal' an implementation of a User's program is possible. For this exploratory work, I have concentrated on demonstrating that the implementations chosen are not 'wrong', that is, that TYPESETTER chooses an implementation for an abstract data type that a human programmer would agree is a reasonable candidate.

To convince the reader, I will present some results using three examples to demonstrate TYPESETTER's flexibility. The first example is a small program (approximately 60 lines) that is useless except to the extent it displays some of the capabilities of TYPESETTER. The second example is an implementation of the MINOPT algorithm presented in section 2.2.2. Finally, we will look at TYPESETTER itself, and examine how it chooses its own implementation. The TYPESETTER prototype has nine implementations spread among the three abstractions Set, List, and Map. Set has five implementations, and the other two have two apiece. Since Sets have more possibilities than the other two ADTs, we will concentrate on showing how TYPESETTER performs on variables declared to be sets of User-defined objects.

There are two distinct questions that the prototype was designed to answer.

The first is to test our hypothesis that a greedy assignment algorithm works well. To recap, the implementation assignment algorithm used in `TYPESETTER` sorts the call sites in decreasing order of importance (where importance is estimated by evaluation functions provided by the Implementor), assigns the most efficient implementation to the first call site, and then, in decreasing order of importance, assigns to all other call sites the most efficient implementation that is consistent with previous assignments. We want to know how quickly an initial assignment of implementations is made, and how close that assignment is to the ‘optimal’ solution, assuming that the performance estimates returned by the evaluation functions is accurate.

The second question is: how accurate are the estimates returned by the Implementor’s evaluation functions? Or, in other words, how closely does the final performance of the implemented program correlate with the predictions made by the Implementor’s evaluation functions?

4.5.1 Small example

Figure 4.24 contains the `TYPESETTER` code for an example program that constructs three sets of integers. Given that the declarations of the variables contain some optional declarations that tell `TYPESETTER` the sets really are sets of integers, it is not too surprising that `TYPESETTER` selects a bit-mapped implementation for them over a linked-list implementation. The primary point of this example, however, is `TYPESETTER`’s ability to share code between instantiations of abstract types. Even though the set *cSet* has more elements than do sets *aSet* and *bSet*, they all three have few enough base elements that they can fit in a 32-bit word, hence they will all use the same implementation of one-word bitmaps. However, because *aSet* and *bSet* were declared to be sets of *short* integers, they will share a coercion class that is different from *cSet*’s.

Figure 4.25 shows the order of priority given to the call sites of our example program based on estimates provided by the profiling implementation’s evaluation routines. Each line shows the name of the function being invoked, the line in the file where the invocation occurs, the index in the profiling array for this file (there is one for each source file making up a program), and the values of the profiling variables. Finally, the “profiling costs” (actually an estimate of the cost) is given. The call sites are sorted in decreasing order based on those estimates.

We’ll look closely at the information printed for the *Set_add* function on line 35 of our test program (the line numbers are not contiguous due to some irrelevant material not included). It has seven profiling variables corresponding to the numbers given in parentheses: *p_cnt*, *p_szA*, *p_appended*, *p_prepende*d, *p_wasIn*, *p_inserte*d, and *p_looke*dAt. Respectively, they count the number of times this call site was executed (*p_cnt*=8), the sum of the size of the set at each invocation (*p_szA*=28), the number of times the element could be appended to a list in which the elements of the set were sorted by their memory address (*p_appended*=8), or prepended

```

23 #define LOOPSIZE 1
24
25 DECLARE(aSet, Set, short, ObjsAreInts, upperb=15, lowerb=0);
26 DECLARE(bSet, Set, short, ObjsAreInts, upperb=15, lowerb=0);
27 DECLARE(cSet, Set, int, ObjsAreInts, upperb=31, lowerb=0);
29
30 main()
31 {
32     int i;
33     for (i = 0; i < 16; i++) {
34         if ((i & 2) != 0) {
35             Set_add(aSet,i);
36         }
37     }
38     Set_add(aSet,1);
39     Set_add(aSet,10);
43     for (i = 0; i < 16; i++) {
44         if ((i & 4) != 0) {
45             Set_add(bSet,i);
46         }
47     }
51     //
52     for (i = 0; i < 32; i++) {
53         if ((i & 15) == 15) {
54             Set_add(cSet,i);
55         }
56     }
60     Set_intersect1(bSet, aSet);
64     //
65     // for every integer in the set c
66     //
67     for (int j=0; j < LOOPSIZE; j++) {
68         forAll(i, cSet,
69             if (j == 0) {
70                 cout << "Found " << i << "\n";
71             }
72         );
73     }
74 }

```

Figure 4.24: Small example


```

Sorted call sites:
Set__intersect1 (line 60 file impltest p[40]=(1,8,9,13,4,4)
                                     profiling costs 72)
Set__add (line 35 file impltest p[5]=(8,28,8,0,0,0,28) profiling costs 36)
Set__add (line 45 file impltest p[26]=(8,28,8,0,0,0,28)
                                     profiling costs 36)
Set__add (line 38 file impltest p[12]=(1,8,0,1,0,0,0) profiling costs 9)
Set__add (line 39 file impltest p[19]=(1,9,0,0,1,0,5) profiling costs 9)
Set__add (line 54 file impltest p[33]=(2,1,2,0,0,0,1) profiling costs 3)
Set__iterate (line 72 file impltest p[47]=(3,6) profiling costs 3)
Set (line 25 file impltest p[0]=(1) profiling costs 1)
Set (line 26 file impltest p[1]=(1) profiling costs 1)
Set (line 27 file impltest p[2]=(1) profiling costs 1)
Set__iterInit (line 72 file impltest p[46]=(1) profiling costs 1)
Set__iterCleanup (line 72 file impltest p[49]=(1) profiling costs 1)

```

Figure 4.25: The call sites sorted by profiling estimates of importance

(*p_prepended*=0), the number of times the element being added was already in the set (*p_wasIn*=0), the number of times the element being added had to be inserted into the interior of a list sorted by address (*p_inserted*=0), and the sum of the number of elements that had to be examined over all calls to this functions (*p_lookedAt*=28). The actual code for the profiling implementation for sets is given in Figure 4.26, from which we can see that the profiling implementation also keeps the elements of the set on a list sorted by their memory address.

From Figure 4.25, we can see that the profiling evaluation routines consider the intersection operation on line 60 to be the dominating factor in this program, giving it a weight (72) twice the nearest competitor (the two *adds*, weight 36). Since the intersection function has *aSet* and *bSet* as parameters, assigning an implementation to the intersection function on line 60 will also assign implementations to those two variables. The evaluation functions for the four possible implementations of sets produced the estimates in Figure 4.27 for the *intersect* function. The bitmapped-word implementation is the cheapest, while the most expensive implementation is the one that keeps the elements on a sorted list (in this case, the list is sorted by the values of the integers): apparently, the number of elements in the two sets is not sufficient to pay for the extra overhead of keeping the lists sorted. Therefore, the intersection function on line 60 was assigned *Set_bm wrd_intersect1*, the single word bit-map implementation for sets whose base size is less than or equal to 32. Once this implementation for the function is decided upon, then the arguments to *intersect1* (*aSet* and *bSet*) will be assigned types corresponding to the formal

```

void
FUNCTION(add)( Any e)
{
  Set_P_Link *lp = firstp;
  Set_P_Link **bp = &firstp;
  p_cnt++;
  p_szA += len;
  if (lp != nil && e < lp->data) {
    p_prepended++;
  }
  else {
    while (lp != nil && e > lp->data) {
      bp = &lp->next;
      lp = lp->next;
      p_lookedAt++;
    }
    if (lp == nil) p_appended++;
    else if (e == lp->data) p_wasIn++;
    else p_inserted++;
  }
  if (lp == nil || e < lp->data) {
    Set_P_Link *tp = new Set_P_Link;
    tp->data = e;
    tp->next = lp;
    *bp = tp;
    len++;
  }
  assert(p_cnt == p_prepended + p_appended + p_wasIn + p_inserted);
}

```

Figure 4.26: The actual profiling implementation for the *add* function for Sets

```

Callsite(40): Set_bmwrdr::intersect1__=1.5
Callsite(40): Set_slist::intersect1__=68.1125
Callsite(40): Set_bmarr::intersect1__=5.9
Callsite(40): Set_slistord::intersect1__=91.4

```

Figure 4.27: Estimates of the cost of the intersection function

```

Callsite(40): Set_bmwrdd::add__=12
Callsite(40): Set_slist::add__=23.4
Callsite(40): Set_bmarr::add__=16.8
Callsite(40): Set_slistord::add__=27.6

```

Figure 4.28: Estimates of the cost of the add function on line 54

types of *Set_bmwrdd_intersect1*, which in this case is also *Set_bmwrdd*. Our prototype has only one implementation of the intersect function with two parameters defined. `TYPESETTER` is designed to allow as many implementations as Implementors may deem usable in various situations. This is easily incorporated into `TYPESETTER` because we concentrate on assigning implementations to *functions*; in other words, implementations of variables occurs as a side effect of assigning implementations to functions.

After assigning an implementation to the intersection function, the only variable implementation remaining to be decided is that of *cSet*. Since it does not interact with either *aSet* or *bSet* in a function call, its implementation is independent of theirs. The call on *Set__add* on line 54 of the program is the most important function, according to the profiling implementation's estimates. Figure 4.28 gives the implementations' estimates of the cost of calling their respective versions of the the *Set__add* function. Therefore, *cSet* is also assigned the word bitmap implementation. Figure 4.29 shows `TYPESETTER`'s output, specifying the types of the variables of our small program based on these considerations. The specifications are interpreted as follows:

`INSTANTIATE(I,N,P...)` Create source code for implementation I (name the class N) with parameters P. In our example in Figure 4.29, the instantiations require two parameters: the functions for taking an object to an integer and the inverse.

`COERCE(I,N,C,P...)` Create a coercion class C which converts calls on the functional interface of implementation I into the instantiation class N, using the parameters P.

`DECLARE_M(V,C)` Declare variable V to be of type (coercion class) C.

In the sample program is a constant that determines the number of times the loop containing the iteration over the elements of *cSet* is executed. If that constant is set to ten, instead of one, then the profile data induces `TYPESETTER` to make a different implementation assignment to *cSet*. Figure 4.30 contains a summary of the output from `TYPESETTER`, emphasizing the differences with the previous run of the program. The intersection function is still the most important, but the iterator functions have moved up in importance. Again, because of the independence

```

INSTANTIATE(Set_bmwrđ, Set_bmwrđ_int_int, int, int)@;
COERCE(Set_bmwrđ, Set_bmwrđ_int_int, Set_bmwrđ_int_int_of_int, int)@;
DECLARE_M(cSet, Set_bmwrđ_int_int_of_int, 32)@;
@;
INSTANTIATE(Set_bmwrđ, Set_bmwrđ_int_int, int, int)@;
COERCE(Set_bmwrđ, Set_bmwrđ_int_int, Set_bmwrđ_int_int_of_short, short)@;
DECLARE_M(bSet, Set_bmwrđ_int_int_of_short, 16)@;
@;
INSTANTIATE(Set_bmwrđ, Set_bmwrđ_int_int, int, int)@;
COERCE(Set_bmwrđ, Set_bmwrđ_int_int, Set_bmwrđ_int_int_of_short, short)@;
DECLARE_M(aSet, Set_bmwrđ_int_int_of_short, 16)@;

```

Figure 4.29: TYPESETTER’s assignment of types to the program

of *cSet* from the other variables in the program, there is no effect except on *cSet*’s implementation. Now its most important call site is the call on *Set__iterate*, and the various implementations’ estimates of cost are shown in Figure 4.30. *Set_slist* and *Set_slistord* evaluate the same since they are both linked-list implementations, differing only in the order in which the elements of the set are returned. Selecting between them more or less at random results in assigning the *Set_slist* implementation to *cSet*.

The above implementations were assigned by THERBLIG with $p = 0$; that is, the implementations chosen were the first consistent set of implementations. While a reasonable argument can be made for the implementations that were selected, two questions still remain. Are they the best implementations possible, given the results of the evaluation functions? And does the performance of the program improve?

To be able to answer the second question, we have to have a program that requires a non-trivial amount of time. To that end, we modify our `LOOPSIZE` macro to 100,000. To answer the first question, we run THERBLIG on the program with $p = 1$; i.e. all possible assignments of implementations are evaluated. On this small example, it made no difference: an exhaustive search across all possible implementations of the program still assigned *Set_slist* to *cSet*, and *Set_bmwrđ* to *aSet* and *bSet*.

Table 4.1 shows the various running times of our small example program when *cSet* is implemented with each of the possible implementations for Set. From it, we can see that THERBLIG correctly chose *Set_slist* as one of the best implementations possible for *cSet*. Running THERBLIG with $p = .9$ produced exactly the same result as $p = 1$, corroborating my hypothesis.

```

Sorted call sites:
Set__intersect1 (line 60 file impltest p[40]=(1,8,9,13,4,4)
                                     profiling costs 72)
Set__add (line 35 file impltest p[5]=(8,28,8,0,0,0,28) profiling costs 36)
Set__add (line 45 file impltest p[26]=(8,28,8,0,0,0,28)
                                     profiling costs 36)
Set__iterate (line 72 file impltest p[47]=(30,60) profiling costs 30)
Set__iterInit (line 72 file impltest p[46]=(10) profiling costs 10)
Set__iterCleanup (line 72 file impltest p[49]=(10) profiling costs 10)
Set__add (line 38 file impltest p[12]=(1,8,0,1,0,0,0) profiling costs 9)
Set__add (line 39 file impltest p[19]=(1,9,0,0,1,0,5) profiling costs 9)
Set__add (line 54 file impltest p[33]=(2,1,2,0,0,0,1) profiling costs 3)
Set (line 25 file impltest p[0]=(1) profiling costs 1)
Set (line 26 file impltest p[1]=(1) profiling costs 1)
Set (line 27 file impltest p[2]=(1) profiling costs 1)

Callsite(58): Set_bmwrdd::iterate__=377.830
Callsite(58): Set_bmarrr::iterate__=367.916
Callsite(58): Set_slist::iterate__=105
Callsite(58): Set_slistord::iterate__=105

INSTANTIATE(Set_slist, Set_slist)@;
COERCE(Set_slist, Set_slist, Set_slist_of_int, int)@;
DECLARE_M(cSet, Set_slist_of_int)@;
@;
INSTANTIATE(Set_bmwrdd, Set_bmwrdd_int_int, int, int)@;
COERCE(Set_bmwrdd, Set_bmwrdd_int_int, Set_bmwrdd_int_int_of_short, short)@;
DECLARE_M(bSet, Set_bmwrdd_int_int_of_short, 16)@;
@;
INSTANTIATE(Set_bmwrdd, Set_bmwrdd_int_int, int, int)@;
COERCE(Set_bmwrdd, Set_bmwrdd_int_int, Set_bmwrdd_int_int_of_short, short)@;
DECLARE_M(aSet, Set_bmwrdd_int_int_of_short, 16)@;

```

Figure 4.30: Results from the example with LOOPCOUNT= 10

Set_slist	2.30s
Set_slistord	2.30s
Set_bmwrdd	6.32s
Set_bmarrr	7.63s

Table 4.1: Small example running times with various implementation assignments for *cSet*

4.5.2 MINOPT

Appendix A contains `TYPESETTER` code for an implementation of the MINOPT algorithm discussed in section 2.2.2. This program, we'll call it *minopt*, has three set variables, two list variables, and one map. The map is a dictionary mapping tokens onto node and arc names. In order to compute execution frequencies of a graph object (arc or node) that is not instrumented, each non-instrumented object maintains a list of other graph objects from which its execution count is com-

Graph	gozintas	gozoutas	time
Set_bmarr	Set_slist	Set_slist	2.50s
Set_slist	Set_slist	Set_slist	2.92s
Set_slistord	Set_slistord	Set_slistord	3.37s
Set_bmarr	Set_bmarr	Set_bmarr	3.79s

Table 4.2: Running times for the K-S algorithm.

puted; this is class member variable *flowlist* in class *GraphObj_obj*. The other list is a sorted list of all graph objects (*sortedObjList*).

The remaining three set variables are: *Graph*, the set of all objects (nodes and arcs) that comprise the current graph; *Node_obj::gozintas*, the set of all arcs that enter a node; and *Node_obj::gozoutas*, the set of all arcs that exit a node.

Profile data was generated by running two PFGs through *minopt*, one is a small five-node graph that Knuth and Stevenson used as an example in their paper [30], and the other is the graph in Figure 2.3. Based on that profile data, THERBLIG, with $p = 0$, selected *Set_bmarr* for the variable *Graph*, and *Set_slist* for the two arc lists, *gozintas* and *gozoutas*. This seems a reasonable assignment of implementations, since *Graph* is added to and iterated over, but nothing else. Since it is a completely full set, there are no penalties to pay in a bitmap implementation for having to check bits in that map that aren't set. This is not the case for the *gozintas* and *gozoutas* variables: the number of arcs coming into or leaving an arc is never more than three in our example graphs; a linked list would do much better for these two variables.

With $p = 1$, THERBLIG makes exactly the same choices, again in support of the hypothesis that implementation decisions made early are close to the 'optimal'. Table 4.2 are *minopt*'s running times when the variables are assigned as shown. The input data is a 364-node graph made by replicating and concatenating the graph in Figure 2.3. The first entry in the table uses the implementations chosen by THERBLIG, and the remainder show that it was indeed a reasonable set of implementations.

4.5.3 Implementing THERBLIG

THERBLIG is the analysis software for the TYPESETTER system. From the descriptions of the available abstractions and their implementations, and the description of the User's program, it selects implementations for the variables declared, and functions invoked, in the User's program. THERBLIG is the most complex software written in TYPESETTER and, therefore, will be our next example.

THERBLIG consists of over 8500 lines of TYPESETTER code and comments. This includes almost 2500 lines of TYPESETTER code for the analysis portion of the software, with the other 6000 lines taken up by the nine implementations of the three

abstractions of Sets, Lists, and Maps. Sets has five implementations, including the profiling implementation, while Lists and Maps have two apiece.

In the body of `THERBLIG`, there are 23 variables utilizing these abstractions: four are Lists, seven are Maps, and eleven are Sets. Given that Sets have a more

complete set of implementations than do the other abstractions we will look at how they are implemented by `TYPESETTER`. The complete listing of `THERBLIG` is given in Appendix B for reference.

The eleven Set variables in `THERBLIG` are:

`ADTcalls`: The set of all call sites in the User's program.

`Ic`: A formal parameter to the function *findCompatibleImplementations* (see page 81).

`adt_afcns`: A member of the class *ADType_obj* representing abstract data types; it is the set of all functions that define the interface to the abstraction.

`adtaf_impl_fcns`: A member of the class *ADTabsFcn_obj* representing the abstract functions, each of which will have a set of functions that are its implementations; this is that set of implementation functions.

`callSites`: A variable, local to the function *implementable*, containing all call sites that have a specific variable in their argument lists.

`callSitesp`: A formal parameter which contains all call sites with a specific variable in their argument lists.

`changed`: A local variable (in function *assignable*) which keeps track of variables which have been given tentative assignments.

`changedp`: A formal parameter which keeps track of variables which have been given tentative assignments.

`implSet`: A local variable to function *assignable* that keeps track of all function implementations that are compatible with the current state of assignments and a particular call site.

`ivars`: A formal parameter to *undoImplementations* that is a set of variables whose tentative assignments are to be undone.

`vd_inSigsOf`: A member of the class *VarDecl_obj* that contains the set of all call sites which have this variable as an actual argument.

`as_set`: Each variable used as a parameter to a User-defined function will be aliased to other variables; this is the set of aliases for a variable. It is a member of the class *AliasSet_obj*.

I have not attempted to solve the problem of detecting input-dependent behavior: that is still up to Users to realize about their own programs, and to take appropriate actions, specifically to run the programs a sufficient number of times with typical input data to cover all the important behaviors of their programs. The problem of determining when sufficient 'typical' input data has been utilized is also

beyond the scope of this dissertation. THERBLIG is interesting because it exhibits input-dependent behavior: if there is no profile data for the User's program, it selects only profiling implementations for the program. Because this is a straightforward procedure that does not require much computation, multiple profiling runs of THERBLIG are required to insure that the profiling data is representative of the average behavior of the program. Once profile data exists for THERBLIG to analyze, then the more extensive analysis described in section 4.4.1 is executed to determine an implementation for the program under consideration. Looking at how THERBLIG modifies its idea of a good selection will give us some insight into how it works.

The first step is to build a version of THERBLIG with all variables implemented by the profiling versions of their underlying abstraction. The second step is to run THERBLIG feeding it its own source code. Since there is not yet any profile data, this run simply reads the program description, and writes a specification file giving all of its variables a profiling implementation: this is more-or-less a do-nothing run of THERBLIG.

The third step is to run THERBLIG again, but this time there is profile data generated from its first run. Figure 4.31 shows the call sites of all Set interface functions sorted in the order indicated by the Set profiling implementation's evaluation functions using the profile data generated on the first run. There were a total of 208 call sites in the THERBLIG sources, 174 of which are in the main code. Of these 174, only 53 involve calls on Set interface functions. (Only the non-zero cost call sites that invoke functions in the Set abstraction are shown in the figures.)

The Set profiling implementation's evaluation functions estimate that the call site on line 761, which constructs Set objects, is possibly the greatest bottleneck in the program, with an invocation of *add* trailing a very close second. In general, the profiling implementations' evaluation functions attempt to estimate the *potential* impact of the function at a particular call site—it is a worst case evaluation. In the case of the first *add* on the list, if the implementation chosen were a linked-list implementation, then adding an element could mean having to examine all of the current members looking for duplication; hence the large estimate. Given the number of times the *add* function on line 1078 was called ($p_cnt= 208$), and the sum of the set sizes across those calls ($p_szA= 21528$), we can estimate the average size of the set for each call to be 103.5 (and we note that $208/2 = 104$). Of these 208 calls, exactly $p_appended= 208$ elements were appended to the address-sorted list, and $p_prepending= 0$ were prepended. There were no elements already in the set ($p_wasIn= 0$), and no elements had to be inserted in the list ($p_inserted= 0$). Because all of the elements were appended to the list, then $p_lookedAt= 21528$.

Based on this data, THERBLIG assigned the implementations to the variables as shown in Table 4.3, first column. (Only the Set variables are shown.)

The next step is to run THERBLIG a third time. The last run added to the existing profile data, and did so while executing code that it did not execute during the first run. The question naturally arises as to whether this would change the as-

```

Sorted call sites:
Set (line 761 file main p[88]=(685) profiling costs 21920)
Set__add (line 1078 file main p[176]=(208,21528,208,0,0,0,21528)
          profiling costs 21736)
Set__add (line 1073 file main p[168]=(213,2407,211,0,2,0,2405)
          profiling costs 2618)
Set (line 318 file main p[15]=(276) profiling costs 2208)
Set (line 302 file main p[9]=(239) profiling costs 1912)
Set (line 330 file main p[21]=(119) profiling costs 952)
Set__add (line 763 file main p[89]=(685,0,685,0,0,0,0)
          profiling costs 685)
Set__add (line 827 file main p[103]=(56,519,56,0,0,0,519)
          profiling costs 575)
Set (line 310 file main p[12]=(51) profiling costs 408)
Set (line 610 file main p[85]=(56) profiling costs 392)
Set__add (line 882 file main p[121]=(114,126,114,0,0,0,126)
          profiling costs 240)
Set (line 487 file main p[63]=(4) profiling costs 28)
Set__union1 (line 769 file main p[96]=(4,5,4,5,0) profiling costs 20)
Set (line 130 file main p[3]=(1) profiling costs 8)
Set__iterate (line 389 file main p[33]=(8,8) profiling costs 8)
Set__iterInit (line 389 file main p[32]=(4) profiling costs 4)
Set__iterCleanup (line 389 file main p[35]=(4) profiling costs 4)
:

```

Figure 4.31: The sorted call sites of Set functions from one THERBLIG run

Variable	Assignment 1	Assignment 2	Assignment 3
ADTcalls	Set_bmarr	Set_bmarr	Set_bmarr
Ic	Set_bmarr	* Set_slist	Set_slist
adt_afcns	Set_bmarr	Set_bmarr	Set_bmarr
adtaf_impl_fcns	Set_bmarr	* Set_slist	Set_slist
as_set	Set_slist	Set_slist	Set_slist
callSites	Set_bmarr	Set_bmarr	Set_bmarr
callSitesp	Set_bmarr	Set_bmarr	Set_bmarr
changed	Set_bmarr	* Set_slist	Set_slist
changedp	Set_bmarr	* Set_slist	Set_slist
implSet	Set_bmarr	* Set_slist	Set_slist
ivars	Set_bmarr	* Set_slist	Set_slist
vd_inSigsOf	Set_bmarr	Set_bmarr	Set_bmarr

Table 4.3: Variable assignments based on profile of three runs of THERBLIG with $p = 0$

signments of implementations to variables. And indeed it does, as the second column of Table 4.3 shows. The changed implementations are marked with an asterisk.

The third column of Table 4.3 shows the assignments when THERBLIG is run for a fourth time. But by now the statistics have stabilized, and the assignments do not change.

All of the runs of THERBLIG above were with $p = 0$. Table 4.4 shows the results of running Therblig with $p = 1$. The asterisk beside the entry for the first assignment means that the first choice with $p = 1$ differed from the first choice when $p = 0$. The asterisks in later columns means, as before, that THERBLIG changed the implementation based on more profile data. Again, the assignments have stabilized by the third run.

There are only three differences between the selections made when $p = 0$ and $p = 1$: the variables *callSites*, *callSitesp*, and *vd_inSigsOf* were formerly *Set_bmarr*, a bit mapped array. Looking at all possible combinations of assignment resulted in those implementations being changed to a *Set_slistord*, a simple linked list that keeps its member in the order of their memory addresses.

And finally, Table 4.5 shows the implementations selected when running with $p = .9$. The important point to note here is that running with $p = .9$ means that only about thirty out of 208 call sites (about 15%) are exhaustively analyzed: the remaining 170-some-odd call sites are assigned the first consistent implementation found. The asterisks in the first column indicate that the first assignment differs from the first assignment in Table 4.4. Asterisks in later columns highlight differences from the preceding column. It is interesting to note that the assignments had not stabilized by the third assignment. I did not determine how many iterations $p = .9$ would have required to stabilize.

Variable	Assignment 1	Assignment 2	Assignment 3
ADTcalls	Set_bmarr	Set_bmarr	Set_bmarr
Ic	Set_bmarr	* Set_slist	Set_slist
adt_afcns	Set_bmarr	Set_bmarr	Set_bmarr
adtaf_impl_fcns	* Set_bmarr	* Set_slist	Set_slist
as_set	Set_slist	Set_slist	Set_slist
callSites	* Set_slistord	Set_slistord	Set_slistord
callSitesp	* Set_slistord	Set_slistord	Set_slistord
changed	Set_bmarr	* Set_slist	Set_slist
changedp	Set_bmarr	* Set_slist	Set_slist
implSet	Set_bmarr	* Set_slist	Set_slist
ivars	Set_bmarr	* Set_slist	Set_slist
vd_inSigsOf	* Set_slistord	Set_slistord	Set_slistord

Table 4.4: Variable assignments based on the profile of three runs of THERBLIG with $p = 1$

Variable	Assignment 1	Assignment 2	Assignment 3
ADTcalls	Set_bmarr	Set_bmarr	Set_bmarr
Ic	Set_bmarr	* Set_slist	Set_slist
adt_afcns	Set_bmarr	Set_bmarr	Set_bmarr
adtaf_impl_fcns	* Set_slist	Set_slist	Set_slist
as_set	Set_slist	Set_slist	Set_slist
callSites	* Set_bmarr	Set_bmarr	* Set_slist
callSitesp	* Set_bmarr	Set_bmarr	* Set_slist
changed	Set_bmarr	* Set_slist	Set_slist
changedp	Set_bmarr	* Set_slist	Set_slist
implSet	Set_bmarr	* Set_slist	Set_slist
ivars	Set_bmarr	* Set_slist	Set_slist
vd_inSigsOf	* Set_bmarr	Set_bmarr	* Set_slist

Table 4.5: Variable assignments based on the profile of three runs of THERBLIG with $p = .9$

1	$p = 0$	34.92s
2	$p = .9$	32.60s
3	$p = 1$	35.98s
4	Set_slist	36.21s
5	Set_slistord	36.36s
6	Set_bmarr	152.17s
7	profiling	44.73s

Table 4.6: THERBLIG running times with various implementation assignments

To measure the effectiveness of the assignments, I examined the output from THERBLIG to see if I could have done better with the existing implementations of sets. In essence, I manually performed Low’s search heuristic: perturbing an existing assignment of implementations to see if another would be better. For THERBLIG, there are only three implementations of sets that are feasible:

Set_slist: A simple list, with a single link to successive elements, and a single pointer to the first element of the list.

Set_slistord: A singly-linked list as for Set_slist, with the addition of a pointer to the last element of the list, and the elements are kept on the list in the order of their memory addresses.

Set_bmarr: An array of bits; requires functions to map objects to integers and integers to objects.

The *Set_bmarr* implementation is not feasible since all sets in THERBLIG have more than 32 elements. I ran seven versions of THERBLIG: one that uses only the profiling implementations; one with implementations assigned by THERBLIG running with $p = 0$ (assignment 3 from Table 4.3); one with implementations assigned with $p = 1$ (assignment 3 from Table 4.4); one with implementations assigned with $p = .9$ (assignment 3 from Table 4.5); and three others with all sets assigned *Set_bmarr*, *Set_slist*, and *Set_slistord*. The timing runs had $p = 1$ to exercise THERBLIG as fully as possible. The results are in Table 4.6.

Finally, we look at exactly how much it costs us to profile THERBLIG. A sense of the cost can be had by comparing the running times of the profiling implementation *vs.* the *Set_slist* implementation in Table 4.6. The *slist* implementation of sets was initially derived by removing all profiling code from the profiling implementation. From this, I estimate that profiling using counters in a special implementation slows down the program by 10-20%; that is, it runs 10-20% slower than it would if all the counting code were removed. However, in some cases, this will often be insignificant, particularly where the default implementation is ill-suited to the program being profiled. In that case, the major slow down of the program will be due to algorithmic unsuitability, and not to counting.

Chapter 5

Conclusion

For the last twenty years, no one has agreed with Knuth's Dictum (see page 1) enough to implement the idea, nor has anyone proven the assertion false. There are several assumptions in the Dictum, two of which have formed the central hypotheses of this work. They are:

- that profiling can be done efficiently enough so as not to be perceived as onerous by the programmer; and,
- that compilers and other tools can automatically extract useful information from profile data.

In the process of investigating the first of these hypotheses, I determined that an implicit assumption held by many programmers is false. Most programmers (myself included) have believed that counting executions of basic blocks is sufficient and more efficient than getting the more complete information about arc traversals. I demonstrated in Chapter 2 that this simply is not so. I presented an algorithm MINOPT which finds the 'optimal' instrumentation of a program by automatically placing the instrumentation code in the nodes or on the arcs. Previous algorithms have found optimal solutions for nodes, or for arcs. MINOPT is the first provably minimal algorithm for both nodes and arcs. I also pointed out that the 'optimal' algorithms aren't, that all have assumed the ability to compute instrumentation costs in linear time. I do not know whether there exists such an algorithm or not, but I have shown that if it does exist, it cannot be 'local'. That is, when estimating the instrumentation costs of a node's incoming and outgoing arcs, more information is required than just the execution frequencies of that node and its arcs.

My measurements showed that profiling in the form of in-line execution counts imposes anywhere from 10% to 20% overhead. This can be predicted solely from the observation that most programs' basic blocks average from four to ten instructions in size, and from the not too unrealistic assumption that incrementing a counter in memory requires about the average number of cycles for the execution of an instruction on a machine. Therefore, putting instrumentation in the most

frequently executed basic block will produce a slowdown of 10-20% for the program as a whole.

Programmers would complain if a programs were slowed down 10-20% for no reason. That is, if no one (or thing) was making any use of the profile data, then programmers would turn off profile collection. (That is why all compilers today require the programmer to specify when to collect profile data.) However, the second of the hypotheses above would alleviate the problem considerably. If some part of the programming system were able to utilize the profile data to produce superior programs, then the profile collecting overhead is not onerous. It is comparable to the overhead of non-optimized bounds-checking code. While there has been some research in improving the overhead of profile collection (in particular, see Sarkar's paper on using dependency graphs to optimize profile counting [42]), there has yet to be a definitive exploration of the optimization of profile counting.

For there to be such research, it has to be shown that continual collection of profile data is a win. Therefore, I concentrated in Chapters 3 and 4 in exploring ways a compiler might make use of profile data. In Chapter 3 I presented an algorithm I call Greedy Sewing for improving the behavior of programs on machines with instruction caches. By physically moving basic blocks closer together that are executed close together in time, miss rates in instruction caches can be reduced up to 50%. Profile data not only allows the compiler to know which basic blocks to move closer together, it also allows it to ignore those situations where it will not matter to the final performance of the program.

The primary contribution of this work is the development of a programming system that utilizes profile data to select implementations of program abstractions. The TYPESETTER system integrates the development, evaluation, and selection of alternative implementations of programming abstractions into a package that is transparent to the User. Unlike previous systems, TYPESETTER does not require specialized compiler knowledge of the User or the Implementor. From the data collected so far, the TYPESETTER approach to system synthesis appears to be a promising avenue of research.

5.1 Problems and future work

I have only scratched the surface of the body of engineering problems that need to be solved before TYPESETTER can be considered a complete system. Some of these are related to problems inherent in using profile data to predict the future performance of a program, but others are related to the specific approach taken by TYPESETTER.

Execution counts: During this work, I fell into an assumption that I think is widely shared, but which can cause problems. I had assumed that summing profile counts across multiple runs of a program was a reasonable approach to understanding

the behavior of a program. But consider a program that has a function that is called once for each element on a list. For 99% of the elements, the function requires $O(1)$ time to execute. But for 1 out of 100 elements, it requires much more time. For example, let us assume that the occurrence of a certain kind of element requires that it be put in a separate list, and that sorting this list n -element list requires $O(n^2)$ time (it uses an inefficient sorting algorithm) where $n = 1$. If the program is run M times, and the profile counts used as measures of the complexity of this function are summed, then there comes a point where the one-in-a-hundred event dominates the analysis. If we assume a list that is 100 elements long, and one of the elements causes a re-sort, then running the program 100 times could make the list look like it was 10,000 elements long, with 100 re-sorts, implying that the sorting of the special elements requires as much time as the processing of the non-special elements, when in fact it never sorts a list longer than one element.

In general, this problem will rear its head when evaluation functions are non-linear in the values of the profile variable. For profilers like *prof* and *gprof*, this may not cause any particular problem, even though their output does not indicate how many runs of the program produced the data on which they base their analysis. THERBLIG was modified to count the number of executions of a program in addition to the counters specified in the profiling implementations. During analysis, all counters were divided by the number of program runs to try to avoid problems similar to the ones described in the previous paragraph. However, I am not satisfied that this avoids *all* problems of analysis from execution counts derived from multiple runs. This needs to be examined further.

Evaluation functions: The most difficult functions to write in TYPESETTER are the evaluation functions. While some of the difficulty is due to the fact that I've never had to write functions that evaluate the potential performance of other functions in such numbers before, they bring their own set of problems. For one thing, they are hardly ever 'wrong', at least not in the sense that inaccuracies produce obviously aberrant behavior on the part of the program. I have serendipitously discovered several instances where evaluation expressions I have written do not accurately reflect the performance of the actual function; even ignoring the fact that these are all estimates anyway, the results returned were misleading. Debugging these routines to a reasonable level of accuracy is difficult.

Kenny and Lin [27] report a technique for capturing the behavior of functions that might be usable in a THERBLIG-like environment. The Implementor would specify an expression with free variables that he suspects would adequately capture the behavior of the function in question; for example, $A * x + B * y^2 + C$, where x and y are parameters such as the length of a list, or size of set. By executing the function many times on many inputs, an average behavior for the function based on x and y can be found by determining appropriate values for A , B and C with a curve fitting algorithm. While this may be an approach for rigorously and more

automatically producing evaluation functions, it will not reduce the amount of work required by the Implementor and may impede the Implementor from taking advantage of *logical* information contained in the profile data. For example, a curve-fitting approach may not be able to handle knowledge about the density of bit vectors, order of presentation of elements to a function, etc.

In general, evaluation functions need to be easier to write and debug.

Evaluating ADT-invoked User functions: There are several optionals that require the names of User-defined functions. The ones I have identified are `ObjToInt`, `IntToObj`, and `compareFcn`. They present problems when used because `TYPESETTER` has no way of estimating the runtime resources of the indicated functions. Presumably, future systems will have the User give some indication of the cost of executing these functions so that the evaluation functions can give better estimates of the cost of using implementations that require them. I would like to avoid forcing the User to write evaluation functions: that is mixing the roles of User and Implementor too much. Exactly how to achieve the same result without User-written evaluation routines is yet to be determined.

The prototype finesses the problem entirely. Currently, the `ObjToInt` function must always be a reference to an integer field of the object, and `IntToObj` must be an array reference. This has not been terribly restrictive up to this point, but since the maintenance of the array of objects must be done by the User, it imposes some overhead that should be eliminated. Ideally, a map from integers to objects, and its inverse, should not be in the final implementation of a program unless it is needed. Currently, it will always be there, whether `THERBLIG` selects implementations that use them or not.

Second-order effects: Another problem arises when there are dependencies in the User program that are not part of the information available to a `THERBLIG`-like analyzer. Consider a program that keeps objects sorted on a list, but has its own sorted-list code rather than using a library routine. The list is created from a set of these objects, the implementation of said set assigned by the system. It could turn out that the implementation of the set causes the elements to be returned in an order that interferes with the efficient execution of the User's code: i.e. one implementation of set returns the elements in the order of their memory address which corresponds to the order in which they were constructed which in turn corresponds to the order data was read from a file. It is easy to see that there could be interference between the User's implementation and any implementation chosen by the system for the set, and no amount of analysis of the User's use of the *set* would uncover it.

This is outside the scope of a `THERBLIG` style system. One of its major premisses is that looking at the use of the ADTs alone is sufficient to make a reasonable assignment, and extra-ADT information is simply not made available to it. I have not encountered this kind of second-order effect in any of the programs I have

run through THERBLIG, but theoretically it is possible.

Implementation containment: When a bit vector implementation of a set of size N is instantiated, then any declaration of a smaller set could share the code for the larger set. This would decrease the program's memory size further, at the expense of making the space allocated for some sets larger. Discovering and taking advantage of these tradeoffs would require the evaluation functions to consider space as well as time in their analysis. Since Low, for one, has already considered the more complex space-and-time integral objective function for minimization, I felt that duplicating this was not necessary to my objectives and I have concentrated on the simpler time-analysis.

Even if THERBLIG were capable of handling the space analysis, there is nothing in its analysis framework that would allow the kinds of implementation containment described above. In other words, there is no way for the evaluation functions written by the Implementor to conclude "Use implementation X unless condition Y holds, in which case use implementation Z." Again, future work will have to show, first, that this is an optimization that needs to be available and, second, how to obtain it.

Design of implementation libraries: I have barely begun to explore the possibilities in a library of implementations. As mentioned before, it may be desirable to have several profiling implementations, each capable of collecting certain kinds of information that is otherwise difficult to obtain. For example, once a bit vector implementation of a set is determined to be desirable, another bit-vector oriented profiling implementation could be used to determine which of the many bit vector implementations would be best for this program.

In the interests of simplicity, I have also avoided making use of the more complex language features available in my base language, C++. For instance, the implementations List_slist, Set_slist, and Map_slist all use the same implementations of a linked list as their underlying representation. Currently, they each have their own copies of this code, primarily because the *kinds* of profiling information collected differs between the implementations. It is possible that they could all be derived from a linked-list class, increasing even further the possibilities for code sharing. Future work is needed to look at integrating the class hierarchy and attendant inheritance into the library of implementations.

5.2 Summary

I have explored in some detail the proposition that compilers and language systems can make use of profile data in the generation of code for programs, and in the synthesis of large software systems. I have improved the existing 'optimal'

instrumentation algorithms, and shown how arc counts can be used to improve the execution time of programs on machines with instruction caches. I have presented the design of a language and attendant system that can select for a User the implementations of variables declared to be of an abstract data type. I have also demonstrated that such a system can make reasonable choices for those implementations based on the profile data collected by abstraction-specific profiling implementations.

References

1. AGARWAL, A., SITES, R. AND HOROWITZ, M. ATUM: A New Technique for Capturing Address Traces Using Microcode. *Proceedings 13th Annual Symposium on Computer Architecture* (June 1986).
2. BAER, J-L. AND SAGER, G. R. Dynamic Improvement of Locality in Virtual Memory Systems. *IEEE Transactions on Software Engineering*, SE-2, 1 (March 1976), 54-62.
3. BARSTOW, D. *Automatic Construction of Algorithms and Data Structures using a Knowledge Base of Programming Rules*. PhD Dissertation, Stanford University, November 1977.
4. BARSTOW, D. On Convergence Toward a Database of Program Transformations. *ACM Transactions on Programming Languages and Systems*, 7, 1 (January 1985), 1-9.
5. BELL LABORATORIES. *UNIX Programmer's Manual*. Murray Hill, NJ, 1984.
6. BORG, A., KESSLER, R. E., LAZANA, G. AND WALL, D. W. Long Address Traces for RISC Machines: Generation and Analysis. DEC Western Research Lab, WRL Research Report 89/14, September 1989.
7. CHEUNG, R. C-YEE A Structural Theory for Improving Software Reliability. University of California, Berkeley, PhD Dissertation, December 1974.
8. CHOW, P. MIPS-X Instruction Set and Programmer's Manual. CSLSU, CSL-86-289, May 1986.
9. DEWAR, R. B. K., GRAND, A., LIU, S-C. AND SCHWARTZ, J. T. Programming by Refinement, as Exemplified by the SETL Representation Sublanguage. *ACM Transactions on Programming Languages and Systems*, 1, 1 (January 1979), 27-49.
10. ELLIS, M. AND STROUSTRUP, B. *The Annotated C++ Reference Manual*. Addison Wesley, Reading, MA, 1990.
11. FABRI, J. *Automatic Storage Optimization*. PhD Dissertation, Courant Institute of Mathematical Sciences, NYU, 1979.
12. FERRARI, D. Improving Program Locality by Strategy-Oriented Restructuring. *Information Processing, Proceedings IFIP Congress 74*, New York, NY, Amsterdam (1974).
13. FISCHER, C. N. AND LEBLANC, R. J. *Crafting a Compiler*. Benjamin/Cummings, Menlo Park, CA, 1988.
14. GILBRETH, JR., F. B. AND CAREY, E. G. *Cheaper by the Dozen*. T.Y. Crowell, New York, 1948.

15. GOLDSTINE, H. H. AND VON NEUMANN, J. *Planning and Coding of Problems for an Electronic Computing Instrument*. Institute for Advanced Study, Princeton, NJ, 1947, (3 vols.).
16. GRAHAM, S. L., KESSLER, P. B. AND MCKUSICK, M. K. An Execution Profiler for Modular Programs. *Software-Practice Experience*, 13 (August 1983), 671-685.
17. HANSEN, G. J. *Adaptive Systems for the Dynamic Run-time Optimization of Programs*. PhD Dissertation, Carnegie-Mellon University, Pittsburgh, PA, 1974.
18. HILL, M. D., EGGERS, S. J., LARUS, J. R., TAYLOR, G. S., ADAMS, G., BOSE, B. K., GIBSON, G. A., HANSEN, P. M., KELLER, J., KONG, S. I., LEE, C. G., LEE, D., PENDLETON, J. M., RITCHIE, S. A., WOOD, D. A., ZORN, B. G., HILFINGER, P. N., HODGES, D., KATZ, R. H., OUSTERHOUT, J. AND PATTERSON, D. A. Design Decisions in SPUR. *IEEE Computer*, 19, 11 (November 1986).
19. HILL, M. D. DineroIII Cache Simulator. University of California, Berkeley, UNIX Programmer's Manual, August 1985.
20. HILL, M. D. Aspects of Cache Memory and Instruction Buffer Performance. Computer Science Division, EECS, University of California, Berkeley, Technical Report UCB/CSD 87/381, PhD Dissertation, November 1987.
21. HWU, W-MEI W. AND CHANG, P. P. Achieving High Instruction Cache Performance with an Optimizing Compiler. *Proceedings 16th Annual Symposium on Computer Architecture* (June 1989).
22. INGALLS, D. H. H. *The Execution Time Profile as a Programming Tool*. In *Design and Optimization of Compilers*, R. Rustin, Ed. Prentice Hall, Englewood Cliffs, NJ, 1972, 107-128.
23. INGALLS, D. H. H. FETE - A FORTRAN Execution Time Estimator. SUCSD, STAN-CS-71-204, February 1971.
24. KANT, E. Efficiency Considerations in Program Synthesis. Stanford University, PhD Dissertation, 1981.
25. KARR, M. Code Generation by Coagulation. *Proceedings of the ACM-SIGPLAN 1984 Symposium on Compiler Construction, SIGPLAN Notices*, 19, 6 (June 1984).
26. KATZ, R. H., EGGERS, S. J., GIBSON, G. A., HANSEN, P. M., HILL, M. D., PENDLETON, J. M., RITCHIE, S. A., TAYLOR, G. S., WOOD, D. A. AND PATTERSON, D. A. Memory Hierarchy Aspects of a Multiprocessor RISC: Cache and Bus Analyses. Computer Science Division, EECS, University of California, Berkeley, UCB/CSD 85/221, January 1985.
27. KENNY, K. B. AND LIN, K-J. Performance Polymorphism: Integrating Performance Constraints in a Class Hierarchy. University of Illinois at Urbana-Champaign, Report No. UIUCDCS-R-89-1523, Urbana, Illinois, May 1989.

28. KNUTH, D. E. An Empirical Study of FORTRAN Programs. *Software-Practice Experience*, 1 (1971), 105-133.
29. KNUTH, D. E. *The Art of Computer Programming: Vol. 1, Fundamental Algorithms*. Addison Wesley, Reading, MA, 1973, 2nd ed..
30. KNUTH, D. E. AND STEVENSON, F. R. Optimal measurement points for program frequency counts. *BIT*, 13 (1973), 313-322.
31. LOW, J. R. Automatic Coding: Choice of Data Structures. Computer Science Department, Stanford University, PhD Dissertation, Technical Report CS-452, August 1974.
32. LOW, J. R. Automatic Data Structure Selection: An Example and Overview. *Communications of the ACM*, 21, 5 (May 1978), 376-385.
33. LOW, J. R. AND ROVNER, P. Techniques for the Automatic Selection of Data Structures. *Conference Record of the Third ACM Symposium on Principles of Programming Languages* (January 1976).
34. MCFARLING, S. Program Optimization for Instruction Caches. *Symposium on Architectural Support for Programming Languages and Operating Systems*, Boston, MA (April 3-6, 1989).
35. PARIS, J-F. Application of Restructuring Techniques to the Optimization of Program Behavior in Virtual Memory Systems. PhD Dissertation, University of California, Berkeley, PROGRES report 81.5, May 1981.
36. PETTIS, K. AND HANSEN, R. C. Profile Guided Code Positioning. *Proceedings of the ACM-SIGPLAN 1990 Conference on Programming Language Design and Implementation*, 25, 6 (June 20-22, 1990), 16-27.
37. RAMIREZ, R. J. Efficient Algorithms for Selecting Efficient Data Storage Structures. Faculty of Mathematics, University of Waterloo, PhD Dissertation, Technical Report CS-80-18, March 1980.
38. ROWE, L. A. A Formilization for Modelling Structures and the Generation of Efficient Implementation Structures. University of California, PhD Dissertation, 1976.
39. RYDER, K. D. Optimizing Program Placement in Virtual Systems. *IBM Systems Journal*, 13 (1974), 292.
40. SAMPLES, A. D. User's Guide to the M5 Macro Language. Technical Report UCB/CSD 91/621, March 1991.
41. SAMPLES, A. D. Mache: No-Loss Trace Compaction. *Proceedings of the ACM SIGMETRICS Conference on Measurement and Modeling of Computer Systems*, Berkeley, CA (May 1989).
42. SARKAR, V. Determining Average Program Execution Times and their Variance. *Proceedings of the ACM-SIGPLAN 1989 Conference on Programming Language Design and Implementation*, 24, 7 (June 21-23, 1989), 298-312.

43. SCHWARTZ, J. T., DEWAR, R. B. K., DUBINSKY, E. AND SCHONBERG, E. *Programming with Sets: An Introduction to SETL*. Springer Verlag, Berlin, 1986.
44. SHERMAN, M. S. *Paragon: A Language Using Type Hierarchies for the Specification, Implementation and Selection of Abstract Data Types*. Lecture Notes in Computer Science, 189, Springer Verlag, Berlin, 1985.
45. SMITH, A. J. Cache Memories. *ACM Computing Surveys*, 14, 3 (September 1982), 473–530.
46. STRAUB, R. M. *Taliere: An Interactive System for Data Structuring SETL Programs*. PhD Dissertation, Courant Institute of Mathematical Sciences, NYU, May 1988.
47. SZYMANSKI, T. B. Assembling Code for Machines with Span Dependent Instructions. *Communications of the ACM*, 21, 4 (April 1978), 300–308.
48. THABIT, K. O. *Cache management by the compiler*. PhD Dissertation, Rice University, November 1981.
49. VAN DE VANTER, M. L. Designing BiblioText: An Experiment in User Interface Design. Computer Science Division, EECS, University of California, Berkeley, 88/454, October 24, 1988.
50. WEINBERGER, P. J. Cheap Dynamic Instruction Counting. *Bell System Technical Journal*, 63(1984), 1815.
51. WEISS, G. Recursive Data Types in SETL: Automatic Determination, Data Language Description, and Efficient Implementation. Courant Institute of Mathematical Sciences, NYU, Department of Computer Science, Technical Report 102, March 1986.

Appendix A

The Knuth-Stevenson Algorithm

```

1 //
2 // A program that implements the MINOPT algorithm.  MINOPT finds the minimal
3 // set of arcs and nodes required for implementation in a program that allows
4 // execution counts of all arcs and nodes to be computed.
5 //   Syntax:
6 //       ks < graphDescriptionFile
7 //
8 // This program also implements Knuth and Stevenson's (K-S) algorithm for
9 // finding a minimal set of nodes to instrument in a program.
10 //   Syntax:
11 //       ks -o < graphDescriptionFile
12 //
13 // The differences between the two algorithms are controlled by the boolean
14 // flag MINOPTing; these areas are highlighted with a right comment:
15 //                               f                               // MINOPT
16 //
17
18 #include <stream.h>
19 #include <stdio.h>
20 #include <assert.h>
21 #include "util.H"
22 #include "Tokens.H"
23 #include "userTypes.H"
24 #include "ks_ADTs.H"
25 #include "IMPLSRCS.H"
26
27 boolean MINOPTing = true;                               // MINOPT
28
29 void
30 error(char *msg)
31 {

```

```

32  cerr <<"Error: " << msg << "\n";
33 }
34
35 void
36 fatal(char *msg)
37 {
38  cerr << "Fatal error: " << msg << "\n";
39  abort();
40 }
41
42 Define('maxNofGraphObjs',1000)@;
43 Define('maxNofGraphObjsm1','Eval(maxNofGraphObjs - 1)')@;
44 Define('go_bminfo','lowerb=0, upperb=maxNofGraphObjsm1,
45         IntToObj=GraphObjNo, ObjToInt=NofGraphObj')@;
46
47 //
48 // class definitions
49 //
50
51 CLASS(Registration_obj)
52 {
53  friend main();
54  Any*  reg;
55  int   size;
56  int   last;
57 public:
58  Registration_obj(int maxsz) { last = -1; size = maxsz;
59                               reg = new Any[size]; }
59  Any&  operator[](int i) {
60      if (i < 0 || i >= size) fatal("illegal registration number");
61      return reg[i];
62  }
63  int   next(Any obj, char* str);
64 }SSALC
65
66 CLASS(Flowitem_obj)
67 {
68 public:
69  Flowitem_obj(GraphObj o, boolean p) { obj = o; plus = p; }
70  GraphObj      obj;
71  boolean       plus; // false if value is to be subtracted
72 }SSALC
73
74
75 CLASS(GraphObj_obj)

```

```

76 {
77  friend class Registration_obj;
78  friend int NofGraphObj(GraphObj go);
79  friend GraphObj GraphObjNo(int i);
80  boolean      IsArc; // false if node
81  int          rnum;
82 public:
83  static Registration_obj go_R;
84  // data
85  int          freq;
86  double       cost; // the cost of measuring this object; defaults to 1
87  Token        name;
88  boolean      instrument; // true if obj is to be instrumented
89  DECLARE(flowlist,List,Flowitem);
90  // structors
91  GraphObj_obj(Token t, boolean b );
92  // functions
93  boolean      isArc() { return IsArc; }
94  boolean      isNode() { return !IsArc; }
95  int          sum(); // add up the flowitems for this object
96  //
97  // data and functions specific to the Knuth-Stevenson algorithm
98  //   I have maintained the variable names from the Knuth and
99  //   Stevenson paper in BIT 13 (1973), pp 323-337.
100 void          arcto(GraphObj); // to make a node for each node and arc
101                                     // in the original graph
102 GraphObj      equivTo; // A graph object that belongs to the
103                                     // same equivalence class as THIS object.
104                                     // == nil if THIS object is the
105                                     // representative of its equivalence class.
106 //
107 GraphObj      super(); // Returns the unique graph object
108                                     // that is the representative object
109                                     // of the equivalence class to which
110                                     // THIS graph object belongs (a
111                                     // function that chases the
112                                     // equivTo chain).
113 //
114 GraphObj      follow; // creates an super-arc in the reduction graph
115                                     // from THIS object to the FOLLOW object.
116 //
117 boolean      d;
118 GraphObj      u;
119 GraphObj      compfather;
120 // These three attributes are used only for vertices in the reduced

```

```

121 // graph. We build the spanning tree of connected components in the
122 // reduced graph as we go. There is one arc in the reduced graph
123 // for each vertex and arc in the original graph. If COMPFATHER ==
124 // nil then this vertex represents an (super) arc in the reduced
125 // graph, viz an arc from u.super() to u.follow().super(). If d,
126 // then that arc goes from THIS to COMPFATHER, else the arc goes
127 // from COMPFATHER to THIS.
128 //
129 GraphObj      comp();
130 // The representative of this supervertex's component in the reduced
131 // program flow graph constructed so far; analogous to super(), above.
132 void          makeComponent();
133 void          createRedArc();
134 }SSALC
135
136
137 CLASS(Node_obj) : public GraphObj_obj
138 {
139 public:
140 // data
141 DECLARE(gozoutas, Set, Arc, go_bminfo);
142 DECLARE(gozintas, Set, Arc, go_bminfo);
143 // functions
144 void      gozinta(Arc);
145 void      gozouta(Arc);
146 int      sumGozintas();
147 int      sumGozoutas();
148 // structors
149 Node_obj(Token t);
150 }SSALC
151
152 CLASS(Arc_obj) : public GraphObj_obj
153 {
154 public:
155 // data
156 Node      from;
157 Node      to;
158 // structors
159 Arc_obj(Token t);
160 // functions
161 void      goes(Node F, Node T)
162 {
163 (from=F)->gozouta(this);
164 (to=T)->gozinta(this);
165 }

```

116

```
166 }SSALC
167
168 @@ ===== Registration =====
169
170 Registration_obj::next(Any obj, char* str)
171 {
172     if (last == size-1) {
173         int newsize = (9*size)/8;
174         Any* newreg = new Any[newsize];
175         int i;
176         for (i = 0; i < size; i++) newreg[i] = reg[i];
177         delete [size]reg;
178         size = newsize;
179         reg = newreg;
180         cerr << "Warning: registration for " << str << " increased to " << size
181             << "\n";
182     }
183     reg[++last] = obj;
184     return last;
185 }
186
187 // ===== Graph objects =====
188
189 Registration_obj GraphObj_obj::go_R(maxNofGraphObjs);
190 int NofGraphObj(GraphObj go)
191 {
192     return (((GraphObj)go)->rnum);
193 }
194 GraphObj GraphObjNo(int i)
195 {
196     return (GraphObj_obj::go_R[i]);
197 }
198
199 CONSTRUCTOR(GraphObj_obj, Token t, boolean b)
200 {
201     name = t;
202     IsArc = b;
203     freq = -99999999;
204     cost = 1;
205     instrument = false;
206     equivTo = nil;
207     follow = nil;
208     u = nil;
209     compfather = nil;
210     rnum = go_R.next(this,"GraphObj");
```

```
211 }
212
213 void
214 GraphObj_obj::arcto(GraphObj tgt)
215 {
216     if (follow == nil) {
217         follow = tgt;
218     }
219     else {
220         // make tgt and follow equivalent
221         tgt = tgt->super();
222         if (tgt != follow->super()) {
223             tgt->equivTo = follow->super();
224         }
225     }
226 }
227
228 GraphObj
229 GraphObj_obj::super()
230 {
231     return (equivTo == nil? this : equivTo->super());
232 }
233
234 GraphObj
235 GraphObj_obj::comp()
236 {
237     return (compfather == nil? this : compfather->comp());
238 }
239
240 void
241 GraphObj_obj::makeComponent()
242 {
243     // transform the d, u, and compfather attributes of the
244     // supervertices so that this supervertex is the representative of
245     // its reduction graph spanning tree component.
246     if (compfather != nil) {
247         compfather->makeComponent();
248         compfather->compfather = this;
249         compfather->d = !d;
250         compfather->u = u;
251         compfather = nil;
252     }
253 }
254
255 void
```

```

256 GraphObj_obj::createRedArc()
257 {
258     // create the arc in the reduction graph that will correspond to
259     // THIS graph object
260     GraphObj v, w;
261     v = super();
262     v->makeComponent();
263     w = follow->super();
264     // if v and w not in the same tree of the forest of intermediate
265     // spanning trees, ...
266     if (v != w->comp()) {
267         // put in spanning tree
268         v->compfather = w;
269         v->d = true;
270         v->u = this;
271         instrument = false;
272     }
273     else {
274         // update the flow items
275         instrument = true;
276         while( w != v ) {
277             Flowitem fi = new Flowitem_obj(this, w->d);
278             List_append1(w->u->flowlist, fi);
279             w = w->compfather;
280         }
281     }
282 }
283
284 // ===== Nodes =====
285
286 CONSTRUCTOR(Node_obj, Token t) CC GraphObj_obj(t,false)
287 {
288 }
289
290 void
291 Node_obj::gozinta(Arc a)
292 {
293     Set_add(gozintas, a);
294 }
295
296 void
297 Node_obj::gozouta(Arc a)
298 {
299     Set_add(gozoutas, a);
300 }

```

```

301
302 int
303 Node_obj::sumGozintas()
304 {
305     Arc a;
306     int sum = 0;
307     forAll(a, gozintas,
308           sum += a->freq;
309           );
310     return sum;
311 }
312
313 int
314 Node_obj::sumGozoutas()
315 {
316     Arc a;
317     int sum = 0;
318     forAll(a, gozoutas,
319           sum += a->freq;
320           );
321     return sum;
322 }
323
324 // ===== Arcs =====
325
326 CONSTRUCTOR(Arc_obj, Token t) CC GraphObj_obj(t, true)
327 {
328     from = nil;
329     to = nil;
330 }
331
332 //
333 // ===== the program =====
334 //
335
336 DECLARE(Graph,Set,GraphObj, go_bminfo);
337 DECLARE(dict,Map,Token,GraphObj);
338
339 Node
340 readNode()
341 {
342     Node n;
343     Token t;
344     cin >> t;
345     if (Map_in(dict,t)) {

```


120

```
346     Map_value(dict, t, n);
347     }
348 else {
349     n = new Node_obj(t);
350     Map_define(dict, t, n);
351     }
352 return n;
353 }
354
355 boolean
356 readArc()
357 {
358     // input is a set of lines of the form:
359     // arcname freq cost nodename nodename ,
360     //
361     Token t;
362     double k;
363     cin >> t;
364     if (cin.eof()) return;
365     assert(!Map_in(dict,t));
366     Arc a = new Arc_obj(t);
367     cin >> a->freq;
368     cin >> a->cost;
369     a->cost *= a->freq;
370     Map_define(dict, t, a);
371     // now read the from and to nodes
372     Node F = readNode();
373     Node T = readNode();
374     Set_add(Graph,a);
375     Set_add(Graph,F);
376     Set_add(Graph,T);
377     a->goes(F,T);
378     if (MINOPTing) {
379         // now do the equivalent of K-S arcto function, except we do it for
380         // the nodes and the arcs
381         F->arcto(a);
382         a->arcto(T);
383     }
384     else {
385         F->arcto(T);
386     }
387     // done
388     cin >> t;
389     if (t == commaToken) return true;
390     return false;
```

```
391 }
392
393 void
394 readGraph()
395 {
396     while (readArc());
397 }
398
399 void
400 propagateCounts()
401 {
402     GraphObj o;
403     forAll(o, Graph,
404         if (o->isArc()) {
405             assert(o->freq >= 0);
406         }
407         else {
408             Node_obj& n = *((Node)o);
409             int t = n.sumGozintas();
410             assert(n.freq < 0 || n.freq == t);
411             n.freq = t;
412             n.cost = (double)n.freq;
413             assert(n.freq == n.sumGozoutas());
414         }
415     );
416 }
417
418 int
419 gobjCmp(GraphObj f1, GraphObj f2)
420 {
421     if (f1->cost > f2->cost) return -1;
422     if (f1->cost < f2->cost) return 1;
423     return 0;
424 }
425
426 main(int argc, char **argv)
427 {
428     GraphObj go;
429     double sum = 0;
430     DECLARE(sortedObjList, List, GraphObj);
431     if (argc > 1) {
432         if (strcmp(argv[1], "-o") == 0) MINOPTing = false;
433         else {
434             fatal("Unrecognized option");
435         }

```

122

```
436     }
437     readGraph();
438     propagateCounts();
439     Set_sort2(Graph, sortedObjList, gobjCmp);
440     forAll(go, sortedObjList,
441           if (MINOPTing || go->isNode()) {           // MINOPT
442               go->createRedArc();
443           }
444     );
445     forAll(go, sortedObjList,
446           if (!MINOPTing && go->isArc()) continue;   // MINOPT
447           cout << go->name;
448           if (go->instrument) {
449               cout << ": instrument cost=" << go->cost << "\n";
450               sum += go->cost;
451           }
452           else {
453               // print equation that computes count for this object
454               cout << "= ";
455               Flowitem item;
456               forAll(item, go->flowlist,
457                     if (item->plus) cout << "+";
458                     else cout << "-";
459                     cout << item->obj->name;
460                 );
461               cout << "\n";
462           }
463     );
464     cout << "Instrumentation cost = " << sum << "\n";
465 }
```

Appendix B

Therblig

```

1 @ FILE: types.t
2 @ terminology:
3 @   The class that is a specific implementation of an ADT is called
4 @   a "representation class" of that ADT, since an ADT is not just a
5 @   set of functions but also a representation of the data.
6 @
7 @   The ADT is defined as a set of functions operating on objects of that
8 @   type. A representation class will have zero, one, or more
9 @   "implementations" of those interface functions.
10 @
11 #include <stream.h>
12 #include <stdio.h>
13 #include "userTypes.H"
14 #include "Tokens.H"
15 #include "IMPLSRCS.H"
16 #include "option_types.H"
17
18 Define('maxNofVarDecls',1000)@;
19 Define('maxNofVarDeclsm1','Eval(maxNofVarDecls - 1)')@;
20 Define('vd_bminfo','lowerb=0, upperb=maxNofVarDeclsm1,
21         IntToObj=VarDeclNo, ObjToInt=NofVarDecl')@;
22
23 Define('maxNofADTabsFcn',200)@;
24 Define('maxNofADTabsFcnm1','Eval(maxNofADTabsFcn-1)')@;
25 Define('adtaf_bminfo','lowerb=0, upperb=maxNofADTabsFcnm1,
26         IntToObj=ADTabsFcnNo, ObjToInt=NofADTabsFcn')@;
27
28 Define('maxNofADTimpFcn',200)@;
29 Define('maxNofADTimpFcnm1','Eval(maxNofADTimpFcn-1)')@;
30 Define('afd_bminfo','lowerb=0, upperb=maxNofADTimpFcnm1,
31         IntToObj=ADTimpFcnNo, ObjToInt=NofADTimpFcn')@;

```

```

32
33 Define('maxNofADTcallSite',225)@;
34 Define('maxNofADTcallSitem1','Eval(maxNofADTcallSite-1)')@;
35 Define('acs_bminfo','lowerb=0, upperb=maxNofADTcallSitem1,
36         IntToObj=ADTcallSiteNo, ObjToInt=NofADTcallSite')@;
37
38 CLASS(Registration_obj)
39 {
40     friend main();
41     Any*   reg;
42     int    size;
43     int    last;
44 public:
45     Registration_obj(int maxsz) { last = -1; size = maxsz;
                                     reg = new Any[size]; }
46     Any&   operator[](int i) {
47         if (i < 0 || i >= size) fatal("illegal registration number");
48         return reg[i];
49     }
50     int    next(Any obj, char* str);
51 }SSALC
52
53
54 CLASS(Profarray_obj)
55 {
56     Token    filename; // @N the filename for this array
57     long     size;     // @N the size of this array
58     long     *iarray; // @N the array itself
59     double   *array;  // @N avg'd over runs
60     boolean  Valid;   // @N did the read work?
61 public:
62     Profarray_obj(FILE *pf, Token fn, int sz);
63     ~Profarray_obj() { delete [size+1]array; }
64     double& operator[](int i) { return array[i]; }
65     boolean valid() { return Valid; }
66     Token   file() { return filename; }
67 }SSALC
68
69 #define maxNofOptsPerADT 10
70
71 CLASS(Optional_obj)
72 { @ optionals as they are read in from the input file
73   @ this class is actually overloaded.  When the ADT is being defined,
74   @ a list of these is created for each optional that is possible on a
75   @ variable declaration.  The ival is then the index of that optional for

```

```

76  @ that ADT.
77  @ When the optionals are read on an actual variable declaration,
78  @ then the fields are filled in as documented below.
79  public:
80  int    idx;    @ the index of the optional
81  optionType *table; @ the table it is an index into
82  Token sval;   @ the value represented as a token
83  int    ival;   @ the value as an integer, if ivalid
84  bool   ivalid; @ is the value a valid integer?
85  @ fcns
86  Optional_obj(Token val, optType intp, int idx, optionType *tbl);
87  Optional_obj(int value);
88  Optional_obj();
89  ostream& print(ostream&);
90  ostream& printForm(ostream&);
91 }SSALC
92
93 CLASS(VarDecl_obj)
94 {
95     void                init_VarDecl(void);
96 public:
97     static Registration_obj vd_R;
98     Token                name;           @ my name;
99     ADType               vd_ADT;        @ my abstract type;
100    boolean               implemented;    @ means that this variable has been
101                                @ assigned a representation (i.e.
102                                @ vd_repr contains a valid impln type)
103     ADTRepr              vd_repr;       @ my representation type; used by assign
104     ADTRepr              vd_bestRepr;
105     DECLARE(vd_adtParms, List, Token);  @@ parms in my DECLARE
106                                @ (including optionals?)
107     DECLARE(vd_inSigsOf, Set, ADTcallSite, acs_bminfo);
108     AliasSet             vd_as;
109     int                  rnum;
110     @@
111     @@ OPTIONALS
112     @@
113     boolean              optsParsed;
114     Optional              vd_opts[maxNofOptsPerADT]; // not nil if present
115     String                instance_name, instance_parm;
116     String                coercion_name, coercion_parm;
117     String                constructor_parms;
118     @ fcns
119     VarDecl_obj(Token, ADType, ADTRepr); @ #1
120     VarDecl_obj(Token, ADType, Token);   @ #1

```

```

121 VarDecl_obj(Token, ADType);           @ #2
122 VarDecl_obj(Token, Token); @ check that the 2nd is an ADType name
123 void aliasOf(VarDecl);
124 void betterRepr(){assert(implemented); vd_bestRepr = vd_repr; }
125 @ Note that the difference between #1 and #2 is that #1 assigns a
126 @ representation, while #2 does not.
127 boolean operator==(VarDecl_obj &that);
128 boolean operator!=(VarDecl_obj &that) { return !(*this == that); }
129 ostream& print(ostream&);
130 ostream& printForm(ostream&);
131 }SSALC
132
133 CLASS(Signature_obj)
134 {
135 public:
136 DECLARE(sig_sig, List, VarDecl);
137 @ fcns
138 Signature_obj();
139 void add(VarDecl V);
140 void add(Token, Token, Token); @create the VarDecl yourself
141 void add(VarDecl, Token, Token);
142 void add(Token, ADType, Token); @ ditto
143 void add(Token, ADType, ADTRepr);
144 void add_dontCare();
145 int len(void);
146 ostream& print(ostream&);
147 ostream& printForm(ostream&);
148 }SSALC
149
150 CLASS(ADType_obj)
151 { @ An AbstractDataType; defined by a set of AbstractFunctionDefinitions.
152 public:
153 bool adt_initied;
154 int adt_number;
155 Token name;
156 ADTRepr adt_profileImpl;           @ my profiling representation
157 DECLARE(adt_reprs, Map, Token, ADTRepr); @ my representations
158 DECLARE(adt_afcns, Set, ADTabsFcn, adtaf_bminfo);
159                                     @ abstract fcns defining my interface
159 @@DECLARE(adt_optionals, Map, Token, Optional); @
160 @ fcns
161 ADType_obj(Token);
162 ADType_obj(char *);
163 ostream& print(ostream&);
164 ostream& dump(ostream&);

```

```

165   boolean operator==(ADType_obj r) { return (name == r.name); }
166   boolean operator!=(ADType_obj r) { return (name == r.name); }
167 }SSALC
168
169 CLASS(ADTRepr_obj)
170 { @ the representation of an AbstractDataType;
171   @ an ADT implementation consists of a set of sets of implementations of the
172   @ ADTs functions.
173   @ Note: name == adtr_of->name + '_' + adtr_suffix
174 public:
175     bool          adtr_initied;
176     Token         name;           @ My name (Set_1, List_3, Map_2, etc.)
177     ADType       adtr_of;       @ the ADT I'm a representation of
178     Token         adtr_suffix;
179     int          adtr_number;   @ for accessing tables
180     @ fcns
181     ADTRepr_obj(Token);
182     ADTRepr_obj(Token, ADType);
183     ostream& print(ostream&);
184     void printName(ostream&);
185     boolean operator==(ADTRepr_obj r) { return (name == r.name); }
186 }SSALC
187
188 CLASS(ADTabsFcn_obj)
189 { @ abstract function (interface function)
190 public:
191     static Registration_obj adtaf_R;
192     int          adtaf_uid;     @ for registration
193     Token        name;         @ the name by which the user invokes me
194     ADType       adtaf_for;    @ the ADT I'm in the interface of
195     Signature_obj adtaf_sig;   @ my parameters
196     ADTimpFcn   evalFcn;     @ the function representing the profiling
197                                     @ implementation evaluation function.
198     DECLARE(adtaf_impl_fcns, Set, ADTimpFcn, afd_bminfo); @ my implementations
199     int          nofProfVars;  @ number of profiling variables for this fcn
200     @ functions
201     ADTabsFcn_obj(Token, ADType);
202     ostream& print(ostream&);
203     void printName(ostream&);
204     boolean operator==(ADTabsFcn_obj r) { return (name == r.name); }
205 }SSALC(adtaf_sig())
206
207 CLASS(ADTimpFcn_obj)
208 { @ An implementation of an abstract (interface) function
209 public:

```



```

210 static Registration_obj afd_R;
211 int          afd_uid;          @ index into evalFcns
212 Token        name;            @ the name by which I am ref'd in the library
213 ADTabsFcn    afd_impl_of;     @ the abstract function I implement
214 ADTRepr      repr;           @ the representation I'm a member of
215 Signature_obj afd_sig;        @ my parameters;
216 @ only types, not names, are important in my signature; should
217 @ probably be a derived class of Signature_obj
218 @ functions
219 ADTimpFcn_obj(Token, ADTabsFcn, ADTRepr);
220 ostream& print(ostream&);
221 void printName(ostream&);
222 void typeName(ostream& o) { o << repr->name << "::" << name; }
223 boolean operator==(ADTimpFcn_obj r)
224     { return (name == r.name); }
225 }SSALC(afd_sig())
226
227
228 CLASS(ADTcallSite_obj)
229 { @ A call site where the user's program has invoked one of the interface
230   @ functions for an ADT.
231   boolean implemented;
232   ADTimpFcn    implementation;
233   ADTimpFcn    BetterImpl;
234 public:
235   static Registration_obj acs_R;
236   int          acs_ruid; @ for registration
237   int          acs_upid; @ each call site has a unique id which is its
238   @ base index in the profile eval-function arrays.
239   Profarray    acs_parr; @ the profile array for this call site
240   ADTabsFcn    acs_afcn; @ the abstract function being invoked
241   Signature_obj acs_sig; @ my actual parameters; names given, types computed
242   int          acs_line; @ the line number of the file I'm in
243   double       acs_rank;
244   @ constructors
245   ADTcallSite_obj(int, Profarray, int, ADTabsFcn);
246   @ functions
247   double eval(ADTimpFcn);
248   double eval() { assert(implemented); return eval(implementation); }
249   ostream& print(ostream&);
250   ostream& printForm(ostream&);
251   void printName(ostream& fout) { fout << acs_line; }
252   void implement(ADTimpFcn f) { implemented = true; implementation = f; }
253   void unimplement() { implemented = false; implementation = nil; }
254   void betterImpl() { assert(implemented && implementation!=nil);

```

```

255             BetterImpl=implementation; }
256 }SSALC(acs_sig())
257
258
259 CLASS(UserFcnDecl_obj)
260 {
261     @ We have to know the structure of some user functions' signatures;
262     @ user fcns can cause either conversion functions to be invoked
263     @ or force certain bindings (e.g. if this global is assigned this
264     @ representation, then this formal parm must also have it). Depends
265     @ on whether the analyzer has been implemented with conversion-on-calls
266     @ implemented.
267 public:
268     int             ufd_upid; @ each user function has a unique id
269     Token           name; @ the name of the user function
270     Signature_obj  ufd_sig; @ my parameters; abstract given, repr computed.
271     @ fcns
272     UserFcnDecl_obj(int, Token );
273     ostream& print(ostream&);
274 }SSALC(ufd_sig())
275
276
277 CLASS(UserFcnCall_obj)
278 {
279 public:
280     int             ufc_upid; @ each call site of a user function has a unique id
281     UserFcnDecl    ufc_decl; @ the function being called.
282     Signature_obj  ufc_sig; @ my actual parameters
283     @ fcns
284     UserFcnCall_obj(int, UserFcnDecl);
285     ostream& print(ostream&);
286 }SSALC(ufc_sig())
287
288
289 CLASS(AliasSet_obj)
290 {
291 friend class VarDecl_obj;
292 public:
293     DECLARE(as_set, Set, VarDecl, vd_bminfo);
294     @ fcns
295     AliasSet_obj(VarDecl);
296     void merge(AliasSet);
297 }SSALC()
298
299

```

130

300

```
1 @@ FILE: main.t
2 #include <stream.h>
3 #include <stdio.h>
4 #include <assert.h>
5 #include "util.H"
6 #include <math.h>
7
8 #ifdef DBG_MALLOC
9 extern void malloc_verify();
10 #define MALLOCK do { if (DebugMalloc) malloc_verify(); } while (0)
11 #else
12 #define MALLOCK
13 #endif
14
15 @@
16 @@ read the input file that has all the information we need:
17 @@ (we might note here that EVERY name in this file MUST be unique)
18 @@
19 @@ (1) a list of all ADTs available for analysis, and names of
20 @@                                     implementations:
21 @@         (1.1) a list of the interface functions and abstract parameter
22 @@                                     types;
23 @@         (1.2) a list of the implementations of the interface functions and
24 @@         the parameters' implementation types.
25 @@ (The Therblig system puts these in file <adt>.th in the <adt> directory
26 @@ impls/<adt>. The user has a conventional way of creating the
27 @@ appropriate declaration file, named, ADTs.th, for his program.)
28 @@ Syntax:
29 @@ The ADT and its implementations on one line, followed by several
30 @@ lines of declarations of the abstract functions, followed by the
31 @@ lines describing the implementations of the abstract functions. I.e:
32 @@
33 @@ <ADT> <ADT_i> ... ;
34 @@ = <AbsFcnName> <signature> ,
35 @@ :
36 @@ <AbsFcnName> <ImpFcnName> <signature> ,
37 @@ :
38 @@ ;
39 @@ : another block like the above
40 @@ ;
41 @@ .
42 @@
43 @@ I'll use the equal sign as a flag that this is an abstract function
44 @@ declaration, and so I won't have to worry about the order of the
```

```

43 @@ declarations if I decide to change it later.  The abstract defns
44 @@ will come from the ADT_P.H file, and the implementation defn's will
45 @@ come from the ADT_i.H files.  They end up in their respective ADT_i.th
46 @@ files, which are included into one file ADTs.th in the user's
47 @@ directory by the user's makefile.
48 @@ E.g.:
49 @@
50 @@ Set Set_1 Set_2 ... ,
51 @@ = union1 Set Set ,
52 @@ = union2 Set Set Set ,
53 @@ = add Set ? ,
54 @@ :
55 @@ union1 union1_1_1 Set_1 Set_1 ,
56 @@ union1 union1_1_2 Set_1 Set_2 ,
57 @@ :
58 @@ union2 union2_1 Set_1 Set_1 Set_1 ,
59 @@ union2 union2_2 Set_2 Set_2 Set_2 ,
60 @@ :
61 @@ ;
62 @@ List List_1 List_2 ... ,
63 @@ :
64 @@ ;
65 @@ .
66 @@
67 @@
68 @@ (2) all variable declarations (Therblig puts them in ADT_vars.th).
69 @@ Syntax: var-name ADT-name p1 p2 p3 ... ;
70 @@ A Set 10 int ... ;
71 @@ :
72 @@ .
73 @@
74 @@ (3) all user function declarations of interest
75 @@ (Therblig puts them in ADT_ufcns.th)
76 @@ Syntax: fid user-fcn-name p1-name p1-type p2-name p2-type ... ;
77 @@ 145 userFcn A Set B ? C List ... ;
78 @@ :
79 @@ .
80 @@
81 @@ (4) all ADT function call sites (Therblig puts them in ADT_csites)
82 @@ Syntax: upid ADT-fcn-name var1 var2 ... ;
83 @@ 1234 union1 A B ;
84 @@ :
85 @@ .
86 @@
87 @@ (5) all user function call sites (Therblig puts them in ADT_ucsites)

```

```

88 @@ Syntax: upid user-fcn-name var1 var2 ... '
89 @@ 134 userFcn A ? C ... ;
90 @@ :
91 @@ .
92 @@
93 @@
94 @@ (1) must be written by the compiler when doing the implementations.
95 @@
96 @@ (2)-(5) must be written by the compiler when doing the user's program.
97 @@
98 @@
99
100
101 #include "userTypes.H"
102 #include "main_ADTs.H"
103 Include(types.t)
104 @@ also defines ADTinfoTable
105 #include "OptArrays.H"
106
107 #include "EvalFcns.H"
108
109 #define openFile(fn,io,mode,filename,die)          \
110     name2(io,stream) fn(filename,mode);          \
111     if (die && fn.fail()) fatal("Could not open file");
112
113 #define fopenFile(fn,mode,filename,die)           \
114     FILE *fn = fopen(filename, mode);           \
115     if (die && fn == 0) fatal("Could not fopen file");
116
117
118 @@COERCN_CLASSES
119
120
121 @@ ===== GLOBALS =====
122 @@Debug(Std)
123 @@DebugStack
124 @@DebugPools
125 DECLARE(ADTs,      Map, Token, ADType); // abstract data types
126 @@Debug(Off)
127 DECLARE(ADTReprs,  Map, Token, ADTRepr); // their representations
128 DECLARE(ADTafcns,  Map, Token, ADTabsFcn); // abstract functions
129 @@ DECLARE(ADTifcns,  Map, Token, ADTimpFcn); // their implementations
130 DECLARE(ADTcalls,  Set, ADTcallSite, acs_bminfo); // their call sites
131 DECLARE(Vars,      Map, Token, VarDecl); // variables in the program
132 DECLARE(UserFcns,  Map, Token, UserFcnDecl); // user declared functions

```

```

133 @@ DECLARE(UserCalls, List, UserFcnCall);      // user call sites
134 DECLARE(ProfArrays,Map, Token, Profarray);    // the profile data
135
136 bool    profileDataValid;
137 double  curAssignCost;
138 int     cutOffIndex;
139 int     cutOffPercent;
140
141 ADType  dontCareADT;
142 ADRepr  dontCareRepr;
143
144 // The following defines are due to therblig's minimal parsing ability
145 #define VarDecl_obj_print      (void *)VarDecl_obj::print
146 #define VarDecl_obj_printForm  (void *)VarDecl_obj::printForm
147 #define Token_obj_print       (void *)Token_obj::print
148 #define ADRepr_obj_printName  (void *)ADRepr_obj::printName
149 #define ADTabsFcn_obj_print   (void *)ADTabsFcn_obj::print
150 #define ADTabsFcn_obj_printName (void *)ADTabsFcn_obj::printName
151 #define ADTimpFcn_obj_print    (void *)ADTimpFcn_obj::print
152 #define ADTimpFcn_obj_printName (void *)ADTimpFcn_obj::printName
153 #define c_print                (void *)c->print
154 #define ADTcallSite_obj_print  (void *)ADTcallSite_obj::print
155 #define ADTcallSite_obj_printName (void *)ADTcallSite_obj::printName
156
157
158 @@ THE CLASS ROUTINES FOR THERBLIG
159
160 @@ ===== print routines =====
161
162 #define DefPrinter(type) \
163     ostream& operator<<(ostream& fout, type v) { return v->print(fout); }
164 #define DefPrintForm(type) \
165     ostream& operator<<(ostream& fout, type v) { return v->printForm(fout); }
166 DefPrinter(Optional)
167 DefPrinter(VarDecl)
168 ostream& operator<<(ostream& fout, Signature_obj &v)
169     { return v.printForm(fout); }
170 DefPrinter(ADType)
171 DefPrinter(ADRepr)
172 DefPrinter(ADTabsFcn)
173 DefPrinter(ADTimpFcn)
174 DefPrintForm(ADTcallSite)
175 DefPrinter(UserFcnDecl)
176 DefPrinter(UserFcnCall)
177

```

```

178 @@ ===== Registration =====
179
180 Registration_obj::next(Any obj, char* str)
181 {
182     if (last == size-1) {
183         int newsize = (9*size)/8;
184         Any* newreg = new Any[newsize];
185         int i;
186         for (i = 0; i < size; i++) newreg[i] = reg[i];
187         delete [size]reg;
188         size = newsize;
189         reg = newreg;
190         cerr << "Warning: registration for " << str <<
191             " increased to " << size << "\n";
192     }
193     reg[++last] = obj;
194     return last;
195 }
196
197
198 @@ ===== Profarray_obj =====
199
200 CONSTRUCTOR(Profarray_obj, FILE *pfile, Token fn, int sz)
201 {
202     @@ New requirement: the last entry of each profile array is the number
203     @@ of times that profile array was written to.  When we read the array in
204     @@ it is converted from long to double, with each entry divided by the
205     @@ execution count.
206     Valid = false;
207     filename = fn;
208     size = sz;
209     iarray = new long[sz+1];
210     array = new double[sz+1];
211     if (fread(iarray, sizeof(long), size+1, pfile) != size+1) {
212         cerr << "Profile data array " << fn << " corruption? Not valid.\n";
213         cerr << "    size+1=" << size+1 << "?\n";
214         fclose(pfile);
215         return;
216     }
217     double nofExecutions = (double)iarray[sz];
218     assert(nofExecutions > 0);
219     for (int i = 0; i < sz; i++) {
220         array[i] = iarray[i] / nofExecutions;
221     }
222     array[sz] = iarray[sz];

```

```

223 delete [sz+1]iarray;
224 Valid = true;
225 }
226
227 @@ ===== Optional_obj =====
228
229 CONSTRUCTOR(Optional_obj,Token val,optType type,int index,optionType *tbl)
230 {
231     sval = val;
232     idx = index;
233     table = tbl;
234     if (type == int_opt_type) {
235         assert(val != nil);
236         if (val->type == num_tkn) {
237             ivalid = true;
238             ival = val->val;
239         }
240     } else {
241         cerr << "Warning: " << val << " is not an integer\n";
242         ivalid = false;
243     }
244 }
245 else if (type == tmpint_opt_type) {
246     cerr << "Warning: " << val << " is not a valid optional\n";
247 }
248 }
249
250 CONSTRUCTOR(Optional_obj) // used by feasibility/eval routines
251 {
252     ivalid = false;
253     sval = nil;
254     idx = 0x4FFFFFFF; // something to cause a problem if used
255     table = nil;
256 }
257
258 CONSTRUCTOR(Optional_obj, int value) // used by feasibility/eval routines
259 {
260     ivalid = true;
261     ival = value;
262     sval = nil;
263     idx = 0x4FFFFFFF; // something to cause a problem if used
264     table = nil;
265 }
266
267 ostream &

```



```

268 Optional_obj::print(ostream& fout)
269 {
270     return fout << table[idx].name;
271 }
272
273 ostream &
274 Optional_obj::printForm(ostream& fout)
275 {
276     return fout << table[idx].name;
277 }
278
279 @@ ===== VarDecl_obj =====
280
281 Registration_obj VarDecl_obj::vd_R(maxNofVarDecls);
282 int NofVarDecl(VarDecl vd)
283 {
284     return (((VarDecl)vd)->rnum);
285 }
286 VarDecl VarDeclNo(int i)
287 {
288     return (VarDecl_obj::vd_R[i]);
289 }
290
291 void
292 VarDecl_obj::init_VarDecl()
293 {
294     for (int i = 0; i < maxNofOptsPerADT; i++) {
295         vd_opts[i] = nil;
296     }
297     vd_bestRepr = nil;
298     rnum = vd_R.next(this,"VarDecl");
299     vd_as = new AliasSet_obj(this);
300 }
301
302 CONSTRUCTOR(VarDecl_obj, Token t, ADType a, ADTRepr r)
303 {
304     name = t;
305     vd_ADT = a;
306     vd_repr = r; implemented = true;
307     init_VarDecl();
308 }
309
310 CONSTRUCTOR(VarDecl_obj, Token t, ADType a)
311 {
312     name = t;

```

```

313  vd_ADT = a;
314  vd_repr = nil; implemented = false;
315  init_VarDecl();
316 }
317
318 CONSTRUCTOR(VarDecl_obj, Token t, ADType a, Token tr)
319 {
320  ADTRepr r = new ADTRepr_obj(tr,a);
321  if (tr != dontCareToken && !Map_in(a->adt_reprs, r->name)) {
322      fatal("VarDecl passed Token not in ADType's repr List");
323  }
324  name = t;
325  vd_ADT = a;
326  vd_repr = r; implemented = true;
327  init_VarDecl();
328 }
329
330 CONSTRUCTOR(VarDecl_obj, Token t, Token ta)
331 {
332  ADType a;
333  name = t;
334  Map_value(ADTs, ta, a); // ADTs->value(ta,a);
335  if (a == nil) {
336      a = new ADType_obj(ta);
337      Map_define(ADTs, ta, a);
338  }
339  vd_ADT = a;
340  vd_repr = nil; implemented = false;
341  init_VarDecl();
342 }
343
344 boolean
345 VarDecl_obj::operator==(VarDecl_obj &that)
346 {
347  if (name != that.name) return false;
348  if (vd_ADT != that.vd_ADT) return false;
349  if (!List_equal(vd_adtParms, that.vd_adtParms)) return false;
350  return true;
351 }
352
353 ostream&
354 VarDecl_obj::print(ostream& fout)
355 {
356  fout << name << "(" << vd_ADT->name;
357  if (implemented) {

```

```

358     if (vd_repr != nil) fout << "/" << vd_repr->name;
359     else fout << "<nil>";
360     }
361   if (vd_bestRepr != nil) fout << "/" << vd_bestRepr->name;
362   return fout << " ";
363 }
364
365 ostream&
366 VarDecl_obj::printForm(ostream& fout)
367 {
368   if (implemented) {
369     if (vd_repr != nil) fout << vd_repr->name;
370     else fout << "<nil>";
371     }
372   if (vd_bestRepr != nil) fout << "/" << vd_bestRepr->name;
373   else fout << "(" << vd_ADT->name << " ";
374   return fout << " " << name;
375 }
376
377 void
378 VarDecl_obj::aliasOf(VarDecl v)
379 {
380   AliasSet AS = v->vd_as;
381   VarDecl var;
382   if (AS == vd_as) {
383     // then this is a redundant call: they are both pointing to the same
384     // alias set.
385     return;
386     }
387   forAll(var, v->vd_as->as_set,
388         if (var != v) var->vd_as = vd_as; // because iterCheck gets upset
389         );
390   v->vd_as = vd_as;
391   vd_as->merge(AS); // also frees up set AS
392 }
393
394 @@ ===== Signature_obj =====
395
396 CONSTRUCTOR(Signature_obj)
397 {
398   // do nothing special
399 }
400
401 int
402 Signature_obj::len()

```

```
403 {
404     return List_length(sig_sig);
405 }
406
407
408 void
409 Signature_obj::add(VarDecl v)
410 {
411     List_append1(sig_sig, v);
412 }
413
414 void
415 Signature_obj::add_dontCare()
416 {
417     VarDecl v = new VarDecl_obj(dontCareToken, dontCareADT, dontCareToken);
418     List_append1(sig_sig, v);
419 }
420
421 void
422 Signature_obj::add(Token vn, Token an, Token rn)
423 {
424     ADType a;
425     Map_value(ADTs, an, a);
426     if (a == nil) {
427         fatal("Signature passed non_adt Token");
428     }
429     VarDecl v = new VarDecl_obj(vn, a, rn);
430     List_append1(sig_sig, v);
431 }
432
433 void
434 Signature_obj::add(VarDecl v, Token an, Token rn)
435 {
436     ADType a;
437     Map_value(ADTs, an, a);
438     if (a == nil) {
439         fatal("Signature passed non_adt Token");
440     }
441     List_append1(sig_sig, v);
442 }
443
444 void
445 Signature_obj::add(Token vn, ADType a, Token rn)
446 {
447     VarDecl v = new VarDecl_obj(vn, a, rn);
```

```

448 List_append1(sig_sig, v);
449 }
450
451 void
452 Signature_obj::add(Token vn, ADType a, ADTRepr r)
453 {
454     VarDecl v = new VarDecl_obj(vn, a, r);
455     List_append1(sig_sig, v);
456 }
457
458 ostream&
459 Signature_obj::print(ostream& fout)
460 {
461     fout << "{Sig: ";
462     List_print(sig_sig, fout, VarDecl_obj_printForm);
463     return fout << "}";
464 }
465
466 ostream&
467 Signature_obj::printForm(ostream& fout)
468 {
469     fout << "(";
470     List_print(sig_sig, fout, VarDecl_obj_printForm);
471     return fout << ")";
472 }
473
474
475 @@ ===== ADType_obj =====
476
477 int
478 ADTinfoTableLookup(Token t)
479 {
480     for (int i=0; i < nofADTs; i++) {
481         if (*t == ADTinfoTable[i].name) return i;
482     }
483     cerr << t << ": ";
484     fatal("Unknown ADT in ADTinfoTableLookup");
485 }
486
487 CONSTRUCTOR(ADType_obj, Token t)
488 {
489     name = t;
490     adt_initied = false;
491     adt_number = ADTinfoTableLookup(name);
492 }

```

```

493
494 CONSTRUCTOR(ADType_obj, char *sp)
495 {
496     name = new Token_obj(sp, id_tkn, 0);
497     adt_initied = false;
498     adt_number = ADTinfoTableLookup(name);
499 }
500
501 ADType
502 isADType(Token t)
503 {
504     ADType adt;
505     Map_value(ADTs, t, adt);
506     if (adt == nil) {
507         adt = new ADType_obj(t);
508         Map_define(ADTs, t, adt);
509     }
510     return adt;
511 }
512
513
514 ostream&
515 ADType_obj::dump(ostream& fout)
516 {
517     fout << "{ADType: " << BOOL(adt_initied)
518         << " " << name << "\nreprs = ";
519     Map_print(adt_reprs, fout,
520              Token_obj_print,
521              ADTRepr_obj_printName);
522     fout << "\nabs. fcns. = ";
523     Set_print(adt_afcns, fout, ADTabsFcn_obj_printName);
524     return fout << "}";
525 }
526
527 ostream&
528 ADType_obj::print(ostream& fout)
529 {
530     return fout << name;
531     return fout << "}";
532 }
533
534 @@ ===== ADTRepr_obj =====
535
536 CONSTRUCTOR(ADTRepr_obj, Token t, ADType a)
537 {

```

```
538 name = t;
539 adtr_of = a;
540 adtr_initied = false;
541 adtr_suffix = t->suffix();
542 }
543
544 CONSTRUCTOR(ADTRepr_obj, Token t)
545 {
546     name = t;
547     adtr_of = dontCareADT;
548     adtr_initied = false;
549     adtr_suffix = t->suffix();
550 }
551
552 // overload isADTRepr;
553
554 ADTRepr
555 isADTRepr(Token t)
556 {
557     ADTRepr adtr;
558     Map_value(ADTReprs, t, adtr);
559     if (adtr == nil) {
560         adtr = new ADTRepr_obj(t);
561         Map_define(ADTReprs, t, adtr);
562     }
563     return adtr;
564 }
565
566 ADTRepr
567 isADTRepr(Token t, ADType a)
568 {
569     ADTRepr adtr;
570     Map_value(ADTReprs, t, adtr);
571     if (adtr == nil) {
572         adtr = new ADTRepr_obj(t,a);
573         Map_define(ADTReprs, t, adtr);
574     }
575     else {
576         if (*adtr->adtr_of != *a) {
577             fatal("error in isADTRepr(t,a)");
578         }
579     }
580     return adtr;
581 }
582
```

```

583
584 void
585 ADTRepr_obj::printName(ostream& fout)
586 {
587     fout << name;
588 }
589
590 ostream&
591 ADTRepr_obj::print(ostream& fout)
592 {
593     fout << "{ADTRepr_" << adtr_suffix << ": "
594         << BOOL(adtr_initied) << " " << name << "(" << adtr_of << ")" " ;
595     return fout << "}";
596 }
597
598 @@ ===== ADTabsFcn_obj
599
600 Registration_obj ADTabsFcn_obj::adtaf_R(maxNofADTabsFcn);
601 int NofADTabsFcn(ADTabsFcn af)
602 {
603     return (((ADTabsFcn)af)->adtaf_uid);
604 }
605 ADTabsFcn ADTabsFcnNo(int i)
606 {
607     return (ADTabsFcn_obj::adtaf_R[i]);
608 }
609
610 CONSTRUCTOR(ADTabsFcn_obj, Token t, ADType a)
611 {
612     name = t;
613     adtaf_for = a;
614     adtaf_uid = adtaf_R.next(this,"ADTabsFcn");
615 }
616
617 ostream&
618 ADTabsFcn_obj::print(ostream& fout)
619 {
620     fout << "{ADTabsFcn(" << name << "(" << adtaf_for << ")" << adtaf_sig
621         << ")\nimpls:";
622     Set_print(adtaf_impl_fcns, fout, ADTimpFcn_obj_printName);
623     return fout << "}";
624 }
625
626 void
627 ADTabsFcn_obj::printName(ostream& fout)

```



```

628 {
629     fout << name;
630 }
631
632 @@ ===== ADTimpFcn_obj =====
633
634 Registration_obj ADTimpFcn_obj::afd_R(maxNofADTimpFcn);
635 int NofADTimpFcn(ADTimpFcn aif)
636 {
637     return (((ADTimpFcn)aif)->afd_uid);
638 }
639 ADTimpFcn ADTimpFcnNo(int i)
640 {
641     return (ADTimpFcn_obj::afd_R[i]);
642 }
643
644 CONSTRUCTOR(ADTimpFcn_obj, Token t, ADTabsFcn af, ADTrepr r)
645 {
646     name = t;
647     afd_impl_of = af;
648     repr = r;
649     afd_uid = afd_R.next(this,"ADTimpFcn");
650 }
651
652 void
653 ADTimpFcn_obj::printName(ostream& fout)
654 {
655     fout << name;
656 }
657
658 ostream&
659 ADTimpFcn_obj::print(ostream& fout)
660 {
661     fout << "{ADTimpFcn:" << name << "(" << afd_impl_of << ")"
662         << afd_sig;
663     return fout << "}";
664 }
665
666 @@ ===== ADTcallSite_obj =====
667
668 @@ a presumably invalid value
669 #define IllegalRank -999999.0
670
671 Registration_obj ADTcallSite_obj::acs_R(maxNofADTcallSite);
672 int NofADTcallSite(ADTcallSite acs)

```

```

673 {
674     return (((ADTcallSite)acs)->acs_ruid);
675 }
676 ADTcallSite ADTcallSiteNo(int i)
677 {
678     return (ADTcallSite_obj::acs_R[i]);
679 }
680
681 CONSTRUCTOR(ADTcallSite_obj,int lno,Profarray array,int upid,ADTabsFcn af)
682 {
683     acs_upid = upid;
684     acs_afcn = af;
685     acs_line = lno;
686     acs_parr = array;
687     implemented = false;
688     acs_ruid = acs_R.next(this,"ADTcallSite");
689     acs_rank = IllegalRank;
690 }
691
692 double
693 ADTcallSite_obj::eval(ADTimpFcn f)
694 {
695     return (*evalFcns[f->afd_uid])(this);
696 }
697
698 ostream&
699 ADTcallSite_obj::print(ostream& fout)
700 {
701     fout << "{ADTcallSite(" << acs_upid << ")" << acs_afcn
702         << acs_sig;
703     return fout << "}";
704 }
705
706 ostream&
707 ADTcallSite_obj::printForm(ostream& fout)
708 {
709     fout << acs_afcn->name << " (line " << acs_line << " file "
710         << acs_parr->file() << " p[" << acs_upid << "]=(";
711     for (int i = 0; i < acs_afcn->nofProfVars; i++) {
712         if (i > 0) fout << ",";
713         fout << (*acs_parr)[acs_upid + i];
714     }
715     fout << ") ";
716     if (implemented) {
717         fout<<implementation->repr->name<<" cost "<<eval(implementation);

```

```

718     }
719   else if (BetterImpl != nil) {
720     fout << BetterImpl->repr->name << " cost " << eval(BetterImpl);
721     }
722   else {
723     fout << "profiling costs " << eval(acs_afcn->evalFcn);
724     }
725   fout << ")";
726   return fout;
727 }
728
729 @@ ===== UserFcnDecl_obj =====
730
731 CONSTRUCTOR(UserFcnDecl_obj, int i, Token t)
732 {
733   ufd_upid = i;
734   name = t;
735 }
736
737 ostream&
738 UserFcnDecl_obj::print(ostream& fout)
739 {
740   fout << "{UserFcnDecl(" << ufd_upid << ") " << name << ufd_sig;
741   return fout << "}";
742 }
743
744 @@ ===== UserFcnCall_obj =====
745
746 CONSTRUCTOR(UserFcnCall_obj, int i, UserFcnDecl uf)
747 {
748   ufc_upid = i;
749   ufc_decl = uf;
750 }
751
752 ostream&
753 UserFcnCall_obj::print(ostream& fout)
754 {
755   fout << "{UserFcnCall(" << ufc_upid << ") " << ufc_decl << ufc_sig;
756   return fout << "}";
757 }
758
759 @@ ===== AliasSet_obj =====
760
761 CONSTRUCTOR(AliasSet_obj, VarDecl v)
762 {

```

```

763 Set_add(as_set, v);
764 }
765
766 void
767 AliasSet_obj::merge(AliasSet a)
768 {
769     Set_union1(as_set,a->as_set);
770     delete a;
771 }
772
773 @@ ===== here is the beginning of the input routines for therblig =====
774
775 void
776 readADTs(char *filename)
777 {
778     Token t;
779     int adtrnumber = 0;
780     istream fin(filename,"r");
781     if (fin.fail()) fatal("Could not open file in readADTs");
782     @@ read the names of the abstract data types (ADTs)
783     fin >> t;
784     while (t != dotToken) {
785         @@ - <adt> <adt_i> ... ,
786         assert(t == minusToken);
787         fin >> t;
788         ADType adt = isADType(t);
789         if (adt->adt_initied) {
790             fatal("Duplicate ADT types declared");
791         }
792         @@ read the names of the implementations of this adt
793         @@ the profiling implemenation is the ADT name appended with '_P'
794         @@ and is created automatically.
795         @@ It is not added to the dictionary for this type, since it never
796         @@ enters into the assignment computation. But it must exist for
797         @@ printAssignments procedure to work generally.
798         adt->adt_profileImpl = isADTRepr(t->append("_P"));
799         adt->adt_profileImpl->adtr_initied = true;
800         adt->adt_profileImpl->adtr_number = adtrnumber++;
801         fin >> t;
802         while (t!=commaToken) {
803             ADTRepr adti = isADTRepr(t);
804             if (adti->adtr_initied) {
805                 fatal("Duplicate ADT Reprs declared");
806             }
807             adti->adtr_number = adtrnumber++;

```

```

808         Map_define(adti->adt_reprs, t, adti);
809         adti->adtr_inited = true;
810         fin >> t;
811     }
812     assert(t==commaToken);
813     @@ read the names of the (abstract) functions in the interface
814     @@ now we have one of two kinds of lines:
815     @@     = <absfcname> <Signature> ,
816     @@     or
817     @@     <absfcname> <impfcname> <Signature> ,
818     @@     down to the first semi-colon.
819     @@
820     fin >> t;
821     while (t!=semiToken) {
822         if (t==eqToken) {
823             int nofProfVars;
824             @@ = <absfcname> <nofProfVars> <parmtyp1> <parmtyp2> ... ,
825             fin >> t;
826             ADTabsFcn adaf = new ADTabsFcn_obj(t, adt);
827             Set_add(adti->adt_afcns, adaf);
828             if (Map_in(ADTafcns, t))
829                 fatal("duplicate abstract function names");
830             Map_define(ADTafcns, t, adaf);
831             @@ read the number of profiling variables for this function
832             fin >> adaf->nofProfVars;
833             @@ now read the parameters to the abstract function
834             fin >> t;
835             while (t!=commaToken) {
836                 VarDecl v = new VarDecl_obj(dontCareToken, t);
837                 adaf->adtaf_sig.add(v);
838                 fin >> t;
839             }
840             @@ the profiling evaluation functions must be accessible, but
841             @@ they are declared implicitly.
842             adaf->evalFcn=new ADTimpFcn_obj(t,adaf,adt->adt_profileImpl);
843         }
844         else if (t == plusToken) {
845             @@ + optionalName optionalName ... ,
846             do {
847                 fin >> t;
848             } while (t != commaToken);
849         }
850         else {
851             @@ <ADT> <reprname> <absfcname> <impfcname> <type> <type>
852             @@ (1) (2) (3) (4)

```

```

853      @@ Each function is associated with several names:
854      @@ (1) the name of the ADT it is an impl'n function for;
855      @@ (2) the name of the representation it is an an impl'n
856      @@      function for;
857      @@ (3) the name of the abstract function it implements;
858      @@      These must be unique across all ADTs.
859      @@ (4) the name of the member function by which it is invoked;
860      @@      (this is finessed right now: all member functions
861      @@      are invoked by the same name as their abstract
862      @@      function)
863
864      ADType adt2;
865      Map_value(ADTs, t, adt2);
866      if (adt2 == nil)
867          fatal("Unknown ADT in imp fcn decln");
868      @@ read the repr name (Set_1, Map_2, etc.)
869      Token rn;
870      fin >> rn;
871      ADTRepr repr = isADTRepr(rn);
872      @@ read the abstract fcn name;
873      ADTabsFcn adaf2;
874      Token aftok;
875      fin >> aftok;
876      Map_value(ADTafcns, aftok, adaf2);
877      if (adaf2 == nil)
878          fatal("Unknown abstract function name in imp fcn decln");
879      @@ read an implementation name of an abstract function
880      fin >> t;
881      ADTimpFcn fd = new ADTimpFcn_obj(t, adaf2, repr);
882      Set_add(adaf2->adtaf_impl_fcns, fd);
883      @@ read the parameters of the implementation function
884      fin >> t;
885      while (t != commaToken) {
886          if (t == questToken) {
887              fd->afd_sig.add_dontCare();
888          }
889          else {
890              ADTRepr adtr = isADTRepr(t);
891              assert(t==dontCareToken || adtr->name!=dontCareToken);
892              fd->afd_sig.add(dontCareToken, adtr->adtr_of, adtr);
893          }
894          fin >> t;
895      }
896      }
897      @@ this declaration line processed

```

```

898         fin >> t;
899     }
900     @@ an adt group processed;
901     adt->adt_inited = true;
902     fin >> t;
903     }
904     @@ all ADTs are read
905     @@ need a check that all ADTs have been inited
906     @@ and that all ADTReprs have been inited;
907     ADType adtelt; Token name;
908     forAll('name, adtelt', ADTs,
909         if (!adtelt->adt_inited) {
910             error(name->str);
911             fatal("An uninitialized ADT");
912         }
913     );
914     ADTRepr reprelt;
915     forAll('name, reprelt', ADTReprs,
916         if (!reprelt->adtr_inited) {
917             error(name->str);
918             fatal("An uninitialized ADT representation");
919         }
920     );
921     }
922
923 void
924 readVarDecls(Token ftok, char *filename)
925 {
926     Token var;
927     openFile(fin,i,"r",filename,true);
928     fin >> var;
929     while (var != dotToken) {
930         Token typ,param;
931         int cnt = 0;
932         fin >> typ;
933         @@ assert(isanadt(typ));
934         @@ read all the variable names declared in this program
935         ADType adt3;
936         Map_value(ADTs, typ, adt3);
937         VarDecl vd = new VarDecl_obj(var, adt3);
938         fin >> param;
939         while (param != commaToken &&
940             cnt < ADTinfoTable[adt3->adt_number].nofReqd) {
941             List_append1(vd->vd_adtParms, param);
942             fin >> param;

```

```

943         cnt++;
944     }
945     assert(cnt == ADTinfoTable[adt3->adt_number].nofReqd);
946     @@ just to be on the safe side, we'll check that if the var name
947     @@ is already defined, its declaration matches exactly what we
948     @@ have already defined.
949     if (Map_in(Vars, var)) {
950         VarDecl vddefd;
951         Map_value(Vars, var, vddefd);
952         if (*vd != *vddefd) {
953             cerr << vd << "\n" << vddefd << "\n";
954             fatal("Duplicate declarations of variable do not match");
955         }
956         delete vd;
957     }
958     else {
959         Map_define(Vars, var, vd);
960     }
961     // is also the place to read the optionals
962     // parse_optionals(vd, ADTinfoTable[adt3->adt_number].tbl);
963     // the optionals are of the form dd or dd=token
964     optionType *tbl = ADTinfoTable[adt3->adt_number].tbl;
965     while (parm != commaToken) {
966         // the parm better be an integer
967         int idx = parm->integer();
968         fin >> parm;
969         if (parm == eqToken) {
970             fin >> parm;
971             vd->vd_opts[idx] =
972                 new Optional_obj(parm, tbl[idx].type, idx, tbl);
973             fin >> parm;
974         }
975         else {
976             vd->vd_opts[idx] =
977                 new Optional_obj(nil, tbl[idx].type, idx, tbl);
978         }
979     }
980     fin >> var;
981 }
982 }
983
984 void
985 readUserFcnDecls(Token ftok, char *filename)
986 {
987     Token t;

```



```

988  openFile(fin,i,"r",filename,true);
989  fin >> t;
990  while (t!=dotToken) {
991      int num = t->integer();
992      Token fcnnamet, varnamet;
993      fin >> fcnnamet;
994      UserFcnDecl ufd = new UserFcnDecl_obj(num, fcnnamet);
995      assert(!Map_in(UserFcns,fcnnamet));
996      Map_define(UserFcns,fcnnamet,ufd);
997      @@ read the parm types
998      fin >> varnamet;
999      while (varnamet!=commaToken) {
1000          Token vartypet;
1001          VarDecl var;
1002          fin >> vartypet;
1003          assert(vartypet != commaToken);
1004          if (varnamet == dontCareToken) {
1005              assert(vartypet == dontCareToken);
1006              ufd->ufd_sig.add_dontCare();
1007          }
1008          else {
1009              Map_value(Vars, varnamet, var);
1010              if (var == nil) {
1011                  cerr << "Undeclared variable read in readUserFcnDecls: "
1012                      << varnamet << "\n";
1013                  exit(1);
1014              }
1015              ufd->ufd_sig.add(var, vartypet, dontCareToken);
1016          }
1017          fin >> varnamet;
1018      }
1019      fin >> t;
1020  }
1021 }
1022
1023 void
1024 readADTcallSites(Token ftok, char *filename)
1025 { @@ call sites of calls on functions in the interface of an adt
1026   Token t;
1027   Profarray profileArray;
1028   openFile(fin,i,"r",filename,true);
1029   fin >> t;
1030   Map_value(ProfArrays, ftok, profileArray);
1031   while (t!=dotToken) {
1032       Iterator apli;

```

```

1033     Token nt;
1034     ADTabsFcn aaf1;
1035     int profnum = t->integer();
1036     fin >> t;
1037     int linenum = t->integer();
1038     fin >> nt;
1039     Map_value(ADTafcns, nt, aaf1);
1040     assert(aaf1 != nil);
1041     ADTcallSite afc =
1042         new ADTcallSite_obj(linenum, profileArray, profnum, aaf1);
1043     @@ read the actual parms
1044     fin >> t;
1045     List_iterInit(aaf1->adtaf_sig.sig_sig, apli);
1046     while (t != commaToken) {
1047         VarDecl vd3;
1048         VarDecl abvd3;
1049         assert(!List_iterDone(aaf1->adtaf_sig.sig_sig, apli));
1050         List_iterate(aaf1->adtaf_sig.sig_sig, apli, abvd3);
1051         Map_value(Vars,t,vd3);
1052         @@ if the variable is not found, it is a dontCare;
1053         @@ That is, there is not a dontCareVar;
1054         @@ this can be confirmed by seeing if the corresponding parameter
1055         @@ of the abstract function is a dont care.  Otherwise, error;
1056         if (vd3 == nil) {
1057             if (abvd3 != nil) {
1058                 if (abvd3->vd_ADT == dontCareADT) {
1059                     afc->acs_sig.add_dontCare();
1060                 }
1061             } else {
1062                 fatal("Unrecognized var in call site parm list");
1063             }
1064         }
1065         else {
1066             fatal("List_iterate(aaf1->adtaf_sig.sig_sig,apli,abvd3)
1067                 == nil");
1068         }
1069     }
1070     else {
1071         assert(abvd3->vd_ADT != dontCareADT);
1072         assert(vd3->name != dontCareToken);
1073         afc->acs_sig.add(vd3); @N type to be computed;
1074         Set_add(vd3->vd_inSigsOf, afc);
1075     }
1076     fin >> t;
1077 }

```

```

1077     List_iterCleanup(aaf1->adtaf_sig.sig_sig, apli);
1078     Set_add(ADTcalls,afc);
1079     fin >> t;
1080     }
1081 }
1082
1083 void
1084 readUserFcnCalls(Token ftok, char *filename)
1085 { @@ call sites of calls on (interesting) user functions
1086   Token t;
1087   VarDecl fcnParm;
1088   openFile(fin,i,"r",filename,true);
1089   fin >> t;
1090   while (t!=dotToken) {
1091     int num = t->integer();
1092     Token nt;
1093     fin >> nt;
1094     UserFcnDecl ufd3;
1095     Map_value(UserFcns,nt,ufd3); @@ must exist;
1096     assert(ufd3 != nil);
1097     UserFcnCall ufc = new UserFcnCall_obj(num, ufd3);
1098     Iterator sigIter;
1099     Signature_obj& formalParms = ufc->ufc_decl->ufd_sig;
1100     List_iterInit(formalParms.sig_sig, sigIter);
1101     @@ read the actual parms
1102     fin >> t;
1103     while (t!=commaToken) {
1104       @@ assert((t == dontCare) || isavar(t));
1105       VarDecl vd4;
1106       Map_value(Vars, t, vd4);
1107       ufc->ufc_sig.add(vd4);
1108       // alias this variable with the formal
1109       List_iterate(formalParms.sig_sig, sigIter, fcnParm);
1110       if (vd4 == nil) {
1111         assert(fcnParm != nil && fcnParm->name == dontCareToken);
1112       }
1113       else {
1114         vd4->aliasOf(fcnParm);
1115         if (DebugDetails) {
1116           cerr << ">>>Aliases for "
1117             << hex((int)vd4) << " " << vd4->name << ":";
1118           Set_print(vd4->vd_as->as_set, cerr, VarDecl_obj_printForm);
1119           cerr << "\n";
1120           cerr << ">>>Aliases for "
1121             << hex((int)fcnParm) << " " << fcnParm->name << ":";

```

```

1122             Set_print(fcnParm->vd_as->as_set, cerr,
1123                        VarDecl_obj_printForm);
1124             cerr << "\n";
1125         }
1126     fin >> t;
1127 }
1128 fin >> t;
1129 List_iterCleanup(formalParms.sig_sig, sigIter);
1130 }
1131 }
1132
1133 void
1134 readProfileData(Token ftok, char *filename)
1135 {
1136     if (!profileDataValid) return;
1137     @@ Form of profile data input:
1138     @@ It is an array of 32-bit integers. Each slot in the array
1139     @@ corresponds to a profile variable declared in the ADT profiling
1140     @@ implementations, or to the static call site of an adt function or
1141     @@ user function. The number read with each function in
1142     @@ readprogramdesc is the beginning location of the profile
1143     @@ variables for all functions. E.g. if f1 collects 5 profile
1144     @@ variables, and f1s number is 15, then locations 15 through 19
1145     @@ are the locations in the profile array of f1s profile variables.
1146     @@ these variables are not accessed by anything here except that the
1147     @@ first location of each function read is its execution frequency count.
1148     @@ (this is a required convention and is not currently enforced by
1149     @@ any software checks. this is easy to fix by always declaring
1150     @@ p_cnt for ADT interface functions.)
1151     @@ Note that we are talking about the unique profile id (upid), not the
1152     @@ function id (fid).
1153     @@ If the profile data file does not exist, then our choices are simple:
1154     @@ all implementations are profiling (*_P) implementations.
1155     @@ New requirement: the last entry of each profile array is the number
1156     @@ of times that profile array was written to. When we read the array in
1157     @@ it is converted from long to double, with each entry divided by the
1158     @@ execution count.
1159     long size;
1160     FILE *pfile = fopen(filename, "r");
1161     if (pfile == 0) {
1162         profileDataValid = false;
1163         if (DebugFiles) {
1164             cerr << "    Profile file " << filename << " not there.\n";
1165         }

```

156

```
1166     return;
1167     }
1168     fread((char *)&size, sizeof(long), 1, pfile);
1169     Profarray pa = new Profarray_obj(pfile, ftok, size);
1170     if (!pa->valid()) {
1171         profileDataValid = false;
1172         return;
1173     }
1174     Map_define(ProfArrays, ftok, pa);
1175     fclose(pfile);
1176     return;
1177 }
1178
1179 @@ =====
1180 @@ end of input section of program
1181 @@ =====
1182
1183 @@ main section
1184
1185 @@ global variables
1186 DECLARE(sortedCallSites, List, ADTcallSite);
1187
1188 @@ little functions
1189
1190 boolean
1191 mapsto(ADTRepr r, ADType t)
1192 {
1193     @@ can t be implemented/represented by r
1194     if (DebugDetails)
1195         cerr << ">>mapsto(ADTRepr " << r << "::ADType " << t << ")\n";
1196     /* if (r == dontCareRepr) {
1197         if (t == dontCareADT) return true;
1198         else return false;
1199     }
1200     else if (t == dontCareADT) return false;
1201     -- the above is done by have a pseudo-ADT called dontCareADT, with
1202     -- one don't care ADTRepr dontCareRepr
1203     */
1204     return Map_in(t->adt_reprs,r->name);
1205 }
1206
1207 boolean
1208 mapsto(VarDecl actual, VarDecl formal)
1209 {
1210     assert(formal->implemented);
```

```

1211   if (DebugDetails)
1212       cerr << ">>>mapsto(VarDecl " <<actual<<"::VarDecl "<<formal<<")\n";
1213   if (actual->implemented) {
1214       return (*actual->vd_repr == *formal->vd_repr);
1215   }
1216   else return mapsto(formal->vd_repr, actual->vd_ADT);
1217 }
1218
1219 boolean
1220 mapsto(Signature_obj &actual, Signature_obj &formal)
1221 {
1222   @@ do the types of the variables map to the abstract types in the first
1223   @@ parameter List. they map if each parameter maps.
1224   if (actual.len() != formal.len()) return false;
1225   Iterator nexta, nextf;
1226   VarDecl av, fv;
1227   if (DebugDetails)
1228       cerr << ">>>mapsto(Sig " << actual << "::Sig " << formal << ") \n";
1229   List_iterInit(actual.sig_sig, nexta);
1230   List_iterInit(formal.sig_sig, nextf);
1231   while (List_iterate(actual.sig_sig, nexta, av) &&
1232          List_iterate(formal.sig_sig, nextf, fv))
1233   {
1234       boolean r_mapsto, r_feas;
1235       if (!(r_mapsto = mapsto(av, fv))
1236           || (r_feas = (fv->vd_repr != dontCareRepr
1237                        && !(*feasibilityFcns[fv->vd_repr->adtr_number])(av)))) {
1238           List_iterCleanup(actual.sig_sig, nexta);
1239           List_iterCleanup(formal.sig_sig, nextf);
1240           if (DebugDetails) {
1241               cerr << " >>>Doesn't map because ";
1242               if (!r_mapsto) {
1243                   cerr << "the actual does not map to the formal\n";
1244               }
1245               else {
1246                   assert(r_feas);
1247                   cerr << av->name << " cannot be implemented by "
1248                        << fv->vd_repr->name << "\n";
1249               }
1250           }
1251           return false;
1252       }
1253   }
1254   List_iterCleanup(actual.sig_sig, nexta);

```

158

```
1256 List_iterCleanup(formal.sig_sig, nextf);
1257 return true;
1258 }
1259
1260 // overload isCompatible;
1261
1262 boolean
1263 isCompatible(ADTimpFcn fi, ADTcallSite c)
1264 {
1265     if (DebugAssign)
1266         cerr << ">>>isCompatible(" << fi << ", " << c << ") == ";
1267     if (mapsto(c->acs_sig, fi->afd_sig)) {
1268         if (DebugAssign) cerr << "true\n";
1269         return true;
1270     }
1271     else {
1272         if (DebugAssign) cerr << "false\n";
1273         return false;
1274     }
1275 }
1276
1277 DeclareUserFcn(findCompatibleImplementations, ?, ?, Ic, Set)@;
1278
1279 void
1280 findCompatibleImplementations(ADTcallSite &c,
1281                               DeclareParm(Ic, Set, ADTimpFcn, afd_bminfo))
1282 { @@ find all implementations of c->acs_afcn compatible with the
1283   @@ parameters in call site c;
1284   ADTimpFcn fi;
1285   Set_makeEmpty(Ic);
1286   forAll(fi, c->acs_afcn->adtaf_impl_fcns,
1287         if (isCompatible(fi, c)) {
1288             Set_add(Ic, fi);
1289         }
1290   );
1291   return;
1292 }
1293
1294 DeclareUserFcn(callSitesContaining, ?, ?, callSitesp, Set)@;
1295
1296 void
1297 callSitesContaining(VarDecl v, DeclareParm(callSitesp, Set, ADTcallSite,
1298                                             acs_bminfo))
1299 {
1300     @@ union is overkill: callSitesp is always empty. An interesting
```

```

1301  @@ data point for therblig analysis, though.
1302  Set_union1(callSitesp, v->vd_inSigsOf);
1303 }
1304
1305 void
1306 assignImplType(VarDecl v, ADTRepr t)
1307 {
1308     assert(!v->implemented);
1309     if (DebugDetails)
1310         cerr << ">>>assignImplType: " << v << " <- " << t << "\n";
1311     v->vd_repr = t;
1312     v->implemented = true;
1313 }
1314
1315
1316 void
1317 unassignImplType(VarDecl v)
1318 {
1319     assert(v->implemented);
1320     v->implemented = false;
1321 }
1322
1323 boolean
1324 implementable(VarDecl callSiteVar, ADTRepr r)
1325 {
1326     @@ The inherent feasibility of assigning callSiteVar the
1327     @@ representation r was checked in isCompatible.
1328     assert((*feasibilityFcns[r->adtr_number])(callSiteVar));
1329     @@
1330     @@ if callSiteVar is assigned the impl. type impFcnFormalVar->vd_repr,
1331     @@ then for every call site c that has callSiteVar in its actual parameter
1332     @@ List, check that there still exists AT LEAST ONE impl'n function
1333     @@ that can be used to implement the function called at c.
1334     @@ This check is not absolutely necessary, but I suspect it may cut
1335     @@ down on the amount of backtracking.
1336     @@ Side effect: callSiteVar is assigned impln
1337     @@         impFcnFormalVar->vd_repr, if feasible.
1338     @@
1339     DECLARE(callSites, Set, ADTcallSite, acs_bminfo);
1340     ADTcallSite c;
1341     @@ by defn of findCompatibleImplementations:
1342     assert(mapsto(r, callSiteVar->vd_ADT));
1343     assignImplType(callSiteVar, r);
1344     if (DebugAssign) {
1345         cerr << ">>>trying " << r << " for " << callSiteVar->name << "\n";

```



```

1346     }
1347     CallUserFcn(callSitesContaining, callSiteVar, callSites)@;
1348     callSitesContaining(callSiteVar, callSites);
1349     if (DebugAssign) {
1350         cerr << ">>>callSitesContaining: ";
1351         Set_print(callSites, cerr, ADTcallSite_obj_printName);
1352         cerr << "\n";
1353     }
1354     if (Set_empty(callSites)) return true;
1355     @@ for each call site
1356     Iterate(next, c, callSites,
1357         ADTimpFcn fi;
1358         Iterate(nextfi, fi, c->acs_afcn->adtaf_impl_fcns,
1359             if (mapsto(c->acs_sig, fi->afd_sig)) {
1360                 Set_iterCleanup(c->acs_afcn->adtaf_impl_fcns, nextfi);
1361                 goto SUCCESS;
1362             }
1363         );
1364     unassignImplType(callSiteVar);
1365     if (DebugAssign) {
1366         cout << "No implementations for " << c->acs_afcn << "\n";
1367     }
1368     Set_iterCleanup(callSites, next);
1369     return false;
1370     SUCCESS: ;
1371     );
1372     return true;
1373 }
1374
1375 DeclareUserFcn(parmsImplementable, ?, ?, ?, ?, changedp, Set)@;
1376
1377 boolean
1378 parmsImplementable(ADTcallSite c, ADTimpFcn f,
1379     DeclareParm(changedp, Set, VarDecl, vd_bminfo))
1380 {
1381     @@ Check that the impl. fcn f can be used to implement the fcn
1382     @@ called at call site c by checking that the variables in the actual
1383     @@ parameter List to c can be assigned the
1384     @@ impl. types required by the formal parms of f. Note that the
1385     @@ parallel iteration over the formal parms of f and the
1386     @@ actual parms of c works by the
1387     @@ definition of the function findCompatibleImplementations.
1388     @@ Refinement: this check has to be performed for the actual var and every
1389     @@ variable to which it is aliased in every function call in which they
1390     @@ occur.

```

```

1391  assert(Set_empty(changedp));
1392  Iterator nextf;      @@ points to the implementation formal
1393  VarDecl impforml;
1394  Iterator nextv;      @@ points to the call site actual
1395  VarDecl actual;
1396  if (DebugAssign)
1397      cerr << ">>>parmsImplementable:" << c->acs_sig << f->afd_sig << "\n";
1398  List_iterInit(f->afd_sig.sig_sig, nextf);
1399  List_iterInit(c->acs_sig.sig_sig, nextv);
1400  while (List_iterate(f->afd_sig.sig_sig, nextf, impforml) &&
1401         List_iterate(c->acs_sig.sig_sig, nextv, actual))
1402      {
1403          assert(impforml->implemented); @@ it is a formal impln parm
1404          if (!actual->implemented) {
1405              //      for each variable aliased to actual,
1406              //      for each call site using that variable,
1407              //      see if it is implementable using the type of impformal
1408              VarDecl aliasv;
1409              MALLOC;
1410              if (DebugAssign) {
1411                  cerr << ">>>Aliases for " << actual->name << ":";
1412                  Set_print(actual->vd_as->as_set, cerr, VarDecl_obj_printForm);
1413                  cerr << "\n";
1414              }
1415              forAll(aliasv, actual->vd_as->as_set,
1416                    if (DebugAssign)
1417                        cerr << ">>>Alias loop: " << aliasv << "\n";
1418                    if (!aliasv->implemented) {
1419                        if (implementable(aliasv, impforml->vd_repr)) {
1420                            // implementable assigns implementation to actual
1421                            if (DebugAssign) {
1422                                cout << "Implementing " << aliasv->name <<
1423                                     " as " << impforml->vd_repr->name << "\n";
1424                            }
1425                            Set_add(changedp, aliasv);
1426                        }
1427                    } else {
1428                        List_iterCleanup(f->afd_sig.sig_sig, nextf);
1429                        List_iterCleanup(c->acs_sig.sig_sig, nextv);
1430                        if (DebugAssign) {
1431                            cout << "Could not implement " << aliasv->name
1432                                 << " as " << impforml->vd_repr << "\n";
1433                        }
1434                        return false;
1435                    }

```

162

```
1436         }
1437     );
1438 }
1439 }
1440 List_iterCleanup(f->afd_sig.sig_sig, nextf);
1441 List_iterCleanup(c->acs_sig.sig_sig, nextv);
1442 return true;
1443 }
1444
1445 DeclareUserFcn(undoImplementations, ivars, Set, ?, ?)@;
1446
1447 void
1448 undoImplementations(DeclareParm(ivars, Set, VarDecl, vd_bminfo), int index)
1449 {
1450     @@ undo the implementations of the variables in ivars
1451     VarDecl v;
1452     forAll(v, ivars,
1453         if (DebugAssign){
1454             cout << "    Undoing " << v->name << " " << index << "\n";
1455         }
1456         unassignImplType(v);
1457     );
1458 }
1459
1460 @@ the call site of interest in the current invocation of assignable
1461 @@ an unfortunately necessary global
1462
1463 ADTcallSite curCallSite;
1464
1465 int
1466 compareCosts(ADTimpFcn f1, ADTimpFcn f2)
1467 {
1468     @@ compare the resource costs of the two functions
1469     // these computations should be cached in the call site (e.g. a list
1470     // of costs for each possible afd_uid).
1471     double f1r, f2r;
1472     f1r = curCallSite->eval(f1);
1473     f2r = curCallSite->eval(f2);
1474     if (DebugCosts) {
1475         cerr << "Callsite(" << curCallSite->acs_upid << "): ";
1476         f1->typedName(cerr);
1477         cerr << "=" << f1r << " and ";
1478         f2->typedName(cerr);
1479         cerr << "=" << f2r << "\n";
1480     }
```

```

1481   if (f1r < f2r) return -1;
1482   if (f1r > f2r) return 1;
1483   return 0;
1484 }
1485
1486 void
1487 printSortedCallSites(boolean better)
1488 {
1489   ADTcallSite cs;
1490   if (DebugSortCallSites) {
1491     cout << "\nSorted call sites:\n";
1492     forAll(cs, sortedCallSites,
1493           cout << cs << "\n";
1494           if (better)
1495             cs->betterImpl();
1496           );
1497     cout << "\n";
1498   }
1499 }
1500
1501 void
1502 RecordCurAssignments()
1503 {
1504   VarDecl var;
1505   Token t;
1506   printSortedCallSites(true);
1507   forAll('t,var', Vars,
1508         if (t != dontCareToken) {
1509           cout << "Implemented " << var->name
1510                << " as " << var->vd_repr->name;
1511           if (var->vd_bestRepr!=nil && var->vd_bestRepr!=var->vd_repr) {
1512             cout << " (vs. " << var->vd_bestRepr->name << ")";
1513           }
1514           cout << "\n";
1515           var->betterRepr();
1516         }
1517         );
1518 }
1519
1520 boolean
1521 assignable( Iterator iter, int listIndex, double cost )
1522 {
1523   @@ take the next call site c, and assign it the cheapest
1524   @@ implementation you can. whether c can be assigned the cheapest
1525   @@ implementation is determined by parmsImplementable.

```

```

1526 ADTcallSite c;
1527 ADTimpFcn fi;          @@ a candidate implementation of this abs. fcn
1528 DECLARE(implSet, Set, ADTimpFcn, afd_bminfo); @@ Set of impl'n fcns
1529                                     @@ compatible with ADTcallSite c
1530 DECLARE(implList, List, ADTimpFcn);@@ implSet sorted;
1531 boolean worked;@@ true if parmsImplementable succeeded;
1532 assert(cost >= 0);
1533 @@ prune this search branch if we are already too costly;
1534 if (curAssignCost != -1.0 && cost >= curAssignCost) {
1535     cout << ".Pruned at " << listIndex << ".\n";
1536     return false;
1537 }
1538 cout << "." << listIndex;
1539 if (List_iterDone(sortedCallSites, iter)) {
1540     // then we are at the bottom of the file.
1541     if (curAssignCost == -1.0 || cost < curAssignCost) {
1542         RecordCurAssignments();
1543         if (curAssignCost == -1.0) cout << "First ";
1544         else cout << "Better ";
1545         cout << "implementation: " << form("%10.2f",cost);
1546         if (curAssignCost != -1.0) {
1547             double delta = (curAssignCost - cost ) / curAssignCost;
1548             cout << " delta=" << form("%10.8f",delta);
1549         }
1550         cout << "\n";
1551         curAssignCost = cost;
1552     }
1553     else {
1554         cout << "Not better implementation: " << cost << "\n";
1555     }
1556     return true;
1557 }
1558 if (!List_iterate(sortedCallSites,iter,c)) {
1559     assert(false);
1560 }
1561 if (DebugAssign) {
1562     cerr << ">>>LOOP: assignable call site: " << c << "\n";
1563 }
1564 MALLOC;
1565 CallUserFcn(findCompatibleImplementations, c, implSet)@;
1566 findCompatibleImplementations(c, implSet);
1567 if (DebugAssign) {
1568     cerr << ">>>Compatible implementations: ";
1569     Set_print(implSet, cerr, ADTimpFcn_obj_print);
1570     cerr << "\n";

```

```

1571     }
1572     curCallSite = c;
1573     Set_sort2(implSet, implList, &compareCosts);
1574     @@ for each implementation fi for c, assign the types implied by
1575     @@ fi's Signature to the variables for c.
1576     @@ if this assignment is feasible, then recurse and try the next
1577     @@ call site on the List.  if the recursion returns true, then
1578     @@ return true, else the assignment is not feasible.
1579     @@ if the assignment is not feasible try the next implementation fcn.
1580     @@ if no impl fcns are feasible then return false.
1581     Iterate(next, fi, implList,
1582             @@ the Set of all variables assigned
1583             DECLARE(changed, Set, VarDecl, vd_bminfo);
1584             @@ on a call to parmsImplementable.
1585             if (DebugAssign) {
1586                 cerr << ">>>assignable loop on " << fi << " index " <<
1587                 listIndex << "\n";
1588             }
1589             CallUserFcn(parmsImplementable, c, fi, changed);
1590             worked = parmsImplementable(c, fi, changed);
1591             if (DebugAssign) {
1592                 cerr << ">>>parmsImplementable: " << BOOL(worked) << "\n";
1593             }
1594             if (worked) {
1595                 Iterator iter2;
1596                 List_iterCopy(sortedCallSites, iter, iter2);
1597                 c->implement(fi);
1598                 double fcnCost = c->eval();
1599                 if (isnan(fcnCost) || isinf(fcnCost) || fcnCost < 0.0) {
1600                     cerr << "Evaluation function problem:\n";
1601                     cerr << "The fcn for "<<fi->repr->name<<fi->name<<"\n";
1602                     cerr << "fubarred.  It returned " << fcnCost << "\n";
1603                     exit(1);
1604                 }
1605                 if (assignable(iter2, listIndex+1, cost + c->eval())) {
1606                     // we're walking back up the list toward the more
1607                     // important call sites;
1608                     if (listIndex < cutOffIndex) {
1609                         // then we want to try alternatives;
1610                         List_iterCleanup(sortedCallSites, iter2);
1611                         CallUserFcn(undoImplementations,changed,listIndex);
1612                         undoImplementations(changed, listIndex);
1613                         c->unimplement();
1614                         continue; @@ very important: continues the forAll!!!
1615                     }

```

166

```
1616         else {
1617             // then we won't try alternatives;
1618             List_iterCleanup(sortedCallSites, iter2);
1619             List_iterCleanup(implList, next);
1620             if (DebugAssign || DebugCosts) {
1621                 cerr<<">>assigned "<< c <<" "<< fi->repr->name
1622                 << "; cost: " << c->eval(fi) << "\n";
1623             }
1624             return true;
1625         }
1626     }
1627     else {
1628         CallUserFcn(undoImplementations, changed, listIndex);
1629         undoImplementations(changed, listIndex);
1630         c->unimplement();
1631     }
1632     List_iterCleanup(sortedCallSites, iter2);
1633 }
1634 else {
1635     if (DebugAssign) {
1636         cout << "Did not implement call site " << c << "\n";
1637     }
1638     CallUserFcn(undoImplementations, changed, listIndex);
1639     undoImplementations(changed, listIndex);
1640     c->unimplement();
1641 }
1642 );
1643 return false;
1644 }
1645
1646 //=====
1647
1648 void
1649 findVariableAliases()
1650 {
1651     // for each call site
1652     //     alias the formal and the actual
1653     // until done
1654     //     for each variable
1655     //         merge alias sets.
1656 }
1657
1658 void
1659 assignEverythingToProfile()
1660 {
```

```

1661 Token t;
1662 VarDecl vd;
1663 forAll('t, vd', Vars,
1664     vd->vd_repr = vd->vd_bestRepr = vd->vd_ADT->adt_profileImpl;
1665     vd->implemented = true;
1666     );
1667 }
1668
1669 int
1670 compareCallSitesImportance(ADTcallSite c1, ADTcallSite c2)
1671 {
1672     if (c1->acs_rank == IllegalRank)
1673         c1->acs_rank = (*evalFcns[c1->acs_afcn->evalFcn->afd_uid])(c1);
1674     if (c2->acs_rank == IllegalRank)
1675         c2->acs_rank = (*evalFcns[c2->acs_afcn->evalFcn->afd_uid])(c2);
1676     if (c1->acs_rank < c2->acs_rank) return 1;
1677     if (c1->acs_rank > c2->acs_rank) return -1;
1678     return 0;
1679 }
1680
1681 void
1682 sortByImportance(void)
1683 {
1684     @@ sorts the Set ADTcalls (the Set of callSites) into the List
1685     @@ sortedCallSites.  the key is the frequency of the call sites.
1686     Set_toList(ADTcalls, sortedCallSites);
1687     List_sort1(sortedCallSites, &compareCallSitesImportance);
1688     @@ better:
1689     @@ Set_sort(ADTcalls, sortedCallSites, &compareCallSitesImportance);
1690     @@ now determine how many of these items we're going to iterate over;
1691     int size = List_length(sortedCallSites);
1692     if (cutOffPercent == 100) {
1693         cutOffIndex = size;
1694     }
1695     else if (cutOffPercent == 0) {
1696         cutOffIndex = 0;
1697     }
1698     else {
1699         double rankSum = 0;
1700         ADTcallSite cs;
1701         forAll(cs, sortedCallSites,
1702             assert(cs->acs_rank != IllegalRank);
1703             rankSum += cs->acs_rank;
1704             );
1705         double rankCutOff = rankSum * cutOffPercent / 100;

```



```

1706     rankSum = 0;
1707     cutOffIndex = 0;
1708     forAll(cs, sortedCallSites,
1709         if (rankCutOff <= rankSum) break;
1710         rankSum += cs->acs_rank;
1711         cutOffIndex++;
1712     );
1713 }
1714 printSortedCallSites(false);
1715 }
1716
1717 void
1718 printAssignments(char *foname)
1719 {
1720     // print out the assignments
1721     openFile(afile,o,"w",foname,true);
1722     Token t;
1723     VarDecl vd;
1724     afile << m5Comment << "This is automatically created by therblig\n";
1725     afile << m5Comment << "\n";
1726     afile << m5PushPool_VarDecl_pool << "VARDECLS_HDR\n";
1727     afile << m5Comment2;
1728     forAll('t, vd', Vars,
1729         if (vd->vd_bestRepr == nil) {
1730             afile << "!!!Not implemented: " << vd->name << "\n";
1731             cerr << "!!!Not implemented: " << vd->name << "\n";
1732         }
1733         else {
1734             @@ first, call the instantiation routine for the ADT
1735             @@ on this variable. It will define the strings necessary
1736             @@ to declare the variable, make sure the sources of the code
1737             @@ exist, and that the appropriate coercion class exists.;
1738             (*instantiationFcns[vd->vd_bestRepr->adtr_number])(vd);
1739             afile << m5Comment2 << m5Comment << "    " << vd->name << "\n"
1740                 << m5Comment2;
1741             @@;
1742             @@ first, put out the INSTANTIATE_ADT_i macro;
1743             @@;
1744             afile << "INSTANTIATE(" << vd->vd_bestRepr->name << ","
1745                 << vd->instance_name;
1746             if (!vd->instance_parm->empty())
1747                 afile << "," << vd->instance_parm->string();
1748             afile << ")" << m5Comment2;
1749             @@;
1750             @@ second put out the COERCE_ADT_i macro;

```

```

1751         @@;
1752         afile << "COERCE(" << vd->vd_bestRepr->name << ","
1753             << vd->instance_name
1754             << "," << vd->coercion_name->string();
1755         if (!vd->coercion_parm->empty())
1756             afile << "," << vd->coercion_parm->string();
1757         afile << ")" << m5Comment2;
1758         @@;
1759         @@ third put out the DECLARE_M macro for the variable itself.;
1760         @@;
1761         afile << m5DECLARE_M << vd->name << ","
1762             << vd->coercion_name->string();
1763         if (!vd->constructor_parms->empty())
1764             afile << "," << vd->constructor_parms->string();
1765         afile << "''" << m5Comment2;
1766         afile << m5Comment2;
1767         @@ Token t;
1768         @@ forAll(t, vd->vd_adtParms, api, afile << "," << t; );
1769         }
1770     );
1771     afile << "VARDECLS_TLR\n";
1772 }
1773
1774
1775 void
1776 initialize()
1777 {
1778     dontCareADT = isADType(dontCareToken);
1779     dontCareRepr = isADTRepr(dontCareToken);
1780     Map_define(dontCareADT->adt_reprs, dontCareToken, dontCareRepr);
1781
1782     dontCareADT->adt_initied = true;
1783     dontCareADT->adt_number = 0x7FFFFFFF;
1784
1785     dontCareRepr->adtr_of = dontCareADT;
1786     dontCareRepr->adtr_initied = true;
1787     dontCareRepr->adtr_number = 0x7FFFFFFF;
1788 }
1789
1790 #define ReadThisForAllFiles(varname, fcn)
1791 fnp = argv; do {
1792     Token t_f;
1793     strcpy(fbuf, *fnp);
1794     strcat(fbuf, "_");
1795     strcat(fbuf, varname);

```

170

```
1796         t_f = new Token_obj(*fnp, id_tkn, 0);           \
1797         if (DebugFiles) {                               \
1798             cerr<<"Doing function "<<#fcn<<" on file "<<fbuf<<"\n"; \
1799             }                                           \
1800         fcn(t_f,fbuf); } while (++fnp != 0)
1801
1802 main(int argc, char **argv)
1803 {
1804     Iterator scs_iter;
1805     char **fnp;
1806     char fbuf[64];
1807     char *vardeclsName = "vardecls.m5";
1808     curAssignCost = -1.0;
1809     argv++; argc--;
1810     cutOffPercent = 0; @@ this should yield the same results as the original
1811     @@ version;
1812     while (argv[0] != 0 && argv[0][0] == ch_minus) {
1813         char *cp = &argv[0][1];
1814         if (*cp == ch_D) {
1815             while (++cp != 0) {
1816                 if (*cp == ch_a) { DebugAssign = true; }
1817                 else if (*cp == ch_c) { DebugCosts = true; }
1818                 else if (*cp == ch_f) { DebugFiles = true; }
1819                 else if (*cp == ch_i) { DebugInput = true; }
1820                 else if (*cp == ch_s) { DebugSortCallSites = true; }
1821                 else if (*cp == ch_d) { DebugDetails = true; }
1822 #ifdef DBG_MALLOC
1823                 else if (*cp == ch_m) { DebugMalloc = true; malloc_debug(1); }
1824                 else if (*cp == ch_M) { DebugMalloc = true; malloc_debug(2); }
1825 #endif
1826                 else fatal("Unknown debugging flag");
1827             }
1828         }
1829         else if (*cp == ch_o) { vardeclsName = argv[1]; argv++; }
1830         else if (*cp == ch_P) {
1831             if (*(cp+1) != ch_null) {
1832                 cutOffPercent = atoi(cp+1);
1833             }
1834             else {
1835                 cutOffPercent = atoi(argv[1]); argv++;
1836             }
1837             assert( 0 <= cutOffPercent && cutOffPercent <= 100 );
1838             if (cutOffPercent < 0) cutOffPercent = 0;
1839             if (cutOffPercent > 100) cutOffPercent = 100;
1840         }

```

```

1841     else fatal("Unknown argument");
1842     argv++;
1843     argc--;
1844     }
1845 if (argc < 1) fatal("Must specify file to work on");
1846 initialize();
1847 readADTs("ADTs.th"); @@ this should really be compiled in and not read.
1848 profileDataValid = true;
1849 ReadThisForAllFiles("profData.dat",readProfileData);
1850 ReadThisForAllFiles("'ADT_vars'.th",readVarDecls);
1851 ReadThisForAllFiles("'ADT_ufcns'.th",readUserFcnDecls);
1852 ReadThisForAllFiles("'ADT_csites'.th",readADTcallSites);
1853 ReadThisForAllFiles("'ADT_ufcalls'.th",readUserFcnCalls);
1854 findVariableAliases();
1855 if (profileDataValid) {
1856     cout << "Attempting assignment ... \n";
1857     sortByImportance();
1858     cout << "Cutoff: " << cutOffPercent << " results in "
1859         << cutOffIndex << " of " << List_length(sortedCallSites)
1860         << " being recursed over. \n";
1861     List_iterInit(sortedCallSites, scs_iter);
1862     assignable(scs_iter, 0, 0.0); @@ always returns false
1863     if (curAssignCost >= 0.0) {
1864         cout << "Writing " << vardeclsName << ". \n";
1865         printAssignments(vardeclsName);
1866     }
1867     else {
1868         cout << "\nCannot assign for some reason:
1869             assigning default profiling implementations\n";
1870         assignEverythingToProfile();
1871         cout << "Writing " << vardeclsName << ". \n";
1872         printAssignments(vardeclsName);
1873         exit(1);
1874     }
1875     List_iterCleanup(sortedCallSites, scs_iter);
1876 }
1877 else {
1878     cout <<
1879         "No profile data: assigning default profiling implementations\n";
1880     assignEverythingToProfile();
1881     printAssignments(vardeclsName);
1882 }
1883 /*                                     */ MALLOC;
1884 if (DebugAssign || DebugCosts) {
1885     cerr<<"Nof VarDecls      = "<< VarDecl_obj::vd_R.last + 1 <<"\n";

```

```

1884     cerr<<"Nof ADTabsFcns   = "<< ADTabsFcn_obj::adtaf_R.last + 1 <<"\n";
1885     cerr<<"Nof ADTimpFcns   = "<< ADTimpFcn_obj::afd_R.last + 1 <<"\n";
1886     cerr<<"Nof ADTcallSites = "<< ADTcallSite_obj::acs_R.last + 1 <<"\n";
1887     }
1888     exit(0);
1889 }

```

```

1  @@ FILE: Tokens.t
2  #define TOKENS_MAIN
3  #include <stream.h>
4  #include <ctype.h>
5  #include "util.H"
6  #include "Tokens.H"
7  #include "charClasses.h"
8  #include "Tokens_ADTs.H"
9  #include "userTypes.H"
10
11 String_obj::~String_obj()
12 {
13     delete [maxlen]str;
14 }
15
16 void
17 String_obj::realloc(int add)
18 {
19     // assumes that len is already set to new length
20     // if add == 0, then this copies the string.
21     char *cp = new char[len];
22     strcpy(cp, str);
23     delete [len-add]str;
24     str = cp;
25     maxlen = len;
26 }
27
28 String_obj&
29 String_obj::operator<<(String_obj& t)
30 {
31     if (maxlen <= (len += t.len)) {
32         realloc(t.len);
33     }
34     strcat(str, t.str);
35     return *this;
36 }
37
38 String_obj&
39 String_obj::operator<<(char *cp)

```

```

40 {
41     int l = strlen(cp);
42     if (maxlen <= (len += 1)) {
43         realloc(1);
44     }
45     strcat(str, cp);
46     return *this;
47 }
48
49 String_obj&
50 String_obj::operator<<(int i)
51 {
52     char buf[32]; // should be big enough for an int
53     sprintf(buf,"%d",i);
54     int l = strlen(buf);
55     if ((len += 1) >= maxlen) {
56         realloc(1);
57     }
58     strcat(str,buf);
59     return *this;
60 }
61
62 String_obj&
63 String_obj::operator<<(Token_obj& T)
64 {
65     if ((len += T.len) >= maxlen) {
66         realloc(T.len);
67     }
68     strcat(str,T.str);
69     return *this;
70 }
71
72 #define HASH_SZ MAX_NOF_TOKENS
73
74 #define NOF_SECHASH 16
75 static Token hash[HASH_SZ];
76 static int hashprime[NOF_SECHASH] = {
77     13,31,41,71,131,139,149,157,163,173,227,373,389,457,461,499 };
78
79 @@ Tokens consist of a type (id_tkn, punct_tkn, num_tkn, ...) and a
80 @@ string.
81
82 Token_obj::Token_obj(char *sp, typetype t, int v)
83 {
84     int hashv;

```

```

85  int i;
86  int l = strlen(sp);
87  int first, phash, shash;
88  char *cp = sp;
89  Token T;
90  bool done;
91  phash = 0;
92  for (i=0; i < l; i++) {
93      phash = (phash << 1) + *cp++;
94      }
95  phash += (int)t + (phash >> 8);
96  first = phash & (HASH_SZ-1);
97  shash = hashprime[(phash >> 10) & (NOF_SECHASH - 1)];
98  done = false;
99  do {
100     if ((T=hash[first]) == 0) {
101         @@ not in table;
102         if (this == 0) {
103             T = hash[first] = (Token) new char[sizeof(Token_obj)];
104             }
105         else T = hash[first] = this;
106         T->str = new char[l+1];
107         strcpy(T->str, sp);
108         T->hashv = phash;
109         T->val = v;
110         T->len = l;
111         T->type = t;
112         T->uid = Token_nextuid++;
113         done = true;
114         }
115     else {
116         if (T->hashv == phash && T->len == l && strcmp(T->str, sp) == 0)
117             @@ found it in the table
118             done = true;
119         else {
120             first = (first + shash) & (HASH_SZ-1);
121             if (first == (phash & (HASH_SZ-1))) {
122                 cerr << "Token table full\n";
123                 exit(1);
124             }
125         }
126     }
127 } while (!done);
128 this = T;
129 }

```

```

130
131 Token Token_obj::suffix()
132 {
133     // return the suffix of the string of the form: xxxx_ss
134     char *cp = strrchr(str, ch__);
135     if (cp == nil) return new Token_obj("?", id_tkn, 0);
136     return new Token_obj(cp+1, id_tkn, 0);
137 }
138
139 Token Token_obj::append(char *sfx)
140 {
141     char buffer[256];
142     strcpy(buffer, str);
143     strcat(buffer, sfx);
144     return new Token_obj(buffer, type, 0);
145 }
146
147 ostream&
148 operator<<(ostream &s, String str)
149 {
150     if (str == 0) { return s << "<string not allocated!>"; }
151     if (str->str == 0) { return s<<"<char array not allocated for string!>"; }
152     else return s << str->str;
153 }
154
155 ostream&
156 operator<<(ostream &s, Token gt)
157 {
158     gt->print(s);
159     return s;
160 }
161
162 static char buf[256];
163
164 #define casech(ch,tkn) case ch:{*cp++ = c; s.get(c); toktype = tkn; break;}
165
166 istream&
167 operator>>(istream& s, Token& gt)
168 {
169     char c;
170     char *cp;
171     typetype toktype;
172     int tokval;
173     int toklen;
174     cp = buf;

```


176

```
175 LOOP:
176   if (!s.good()) {
177       if (s.eof())
178           cerr << "Tokens>> should never reach EOF in any input files\n";
179       if (s.fail()) cerr << "Token >> => _fail?\n";
180       if (s.rdstate() == _bad) cerr << "Token >> => _bad?\n";
181       cerr << "Token >> not good!\n";
182       exit(1);
183   }
184   s.get(c);
185   switch (c) {
186
187       case_is_C_firsttid:
188       {
189           *cp++ = c;
190           s.get(c);
191           while (s.rdstate() != _eof && (isalnum(c) || c == ch_)) {
192               *cp++ = c;
193               s.get(c);
194           }
195           toktype = id_tkn;
196           break;
197       }
198
199       case_ch_twiddle:
200       c = ' ';
201       /* NOTE fall through! */
202
203       case_isspace:
204       goto LOOP;
205
206       case_isdigit:
207       {
208           tokval = c - ch_0;
209           *cp++ = c;
210           s.get(c);
211           if (c == ch_x || c == ch_X) { @@ hex numbers
212               *cp++ = ch_x;
213               s.get(c);
214               while (s.rdstate() != _eof && (isxdigit(c))) {
215                   int h = c - ch_0;
216                   if (h > 9) h = c - ch_A + 10;
217                   if (h > 15) h = c - ch_a + 10;
218                   tokval = tokval*16 + h;
219                   *cp++ = c;
```

```
220             s.get(c);
221         }
222     }
223     else if (isdigit(c)) {
224         tokval = tokval * 10 + c - ch_0;
225         *cp++ = c;
226         s.get(c);
227         while (s.rdstate() != _eof && isdigit(c)) {
228             tokval = tokval * 10 + c - ch_0;
229             *cp++ = c;
230             s.get(c);
231         }
232     }
233     toktype = num_tkn;
234     break;
235 }
236 casech(ch_dot, punct_tkn);
237 casech(ch_colon, punct_tkn);
238 casech(ch_semi, punct_tkn);
239 casech(ch_quest, punct_tkn);
240 casech(ch_comma, punct_tkn);
241 casech(ch_eq, punct_tkn);
242 casech(ch_minus, punct_tkn);
243 casech(ch_plus, punct_tkn);
244
245 default:
246     *cp++ = c;
247     s.get(c);
248     toktype = unknown_tkn;
249     error("Unknown token");
250     *cp = ch_null;
251     error(buf);
252 }
253
254 if (s.rdstate() != _eof) s.putback(c);
255 *cp = ch_null;
256 toklen = cp - buf;
257 gt = new Token_obj(buf, toktype, tokval);
258 if (DebugInput) {
259     cerr << buf << " ";
260     if (gt == commaToken || gt == semiToken) cerr << "\n";
261 }
262 return s;
263 }
```