

A Tree-Based Alternative to Java Byte-Codes

Thomas Kistler and Michael Franz

Department of Information and Computer Science
University of California at Irvine
Irvine, CA 92697-3425

Abstract. Despite the apparent success of the Java Virtual Machine, its lackluster performance makes it ill-suited for many speed-critical applications. Although the latest just-in-time compilers and dedicated Java processors try to remedy this situation, optimized code compiled directly from a C program source is still considerably faster than software transported via Java byte-codes. This is true even if the Java byte-codes are subsequently further translated into native code.

In this paper, we claim that these performance penalties are not a necessary consequence of machine-independence, but related to Java's particular intermediate representation and runtime architecture. We have constructed a prototype and are further developing a software transportability scheme founded on a tree-based alternative to Java byte-codes. This tree-based intermediate representation is not only twice as compact as Java byte-codes, but also contains more high-level information, some of which is critical for advanced code optimizations.

Our architecture not only provides on-the-fly code generation from this intermediate representation, but also continual re-optimization of the existing code-base by a low-priority background process. The re-optimization process is guided by up-to-the-minute profiling data, leading to superior runtime performance.

1 Introduction

In recent months, the Java Virtual Machine [LYJ96] has rapidly become a standard platform for building portable Internet “applets” and applications. For these applications, portability is achieved by compiling Java source files into Java byte-codes (instruction sequences for the Java Virtual Machine) that are completely independent of the eventual target architecture. These byte-codes can easily be distributed over the Internet and interpreted on any given machine.

For small Internet applets, electronics, and household appliances, interpreting Java byte-codes yields adequate performance in most cases. For most other application areas, however, the performance penalty associated with interpreting byte-codes makes such an approach unsuitable—higher performance is required.

To remedy this situation, major software distributors have introduced *just-in-time* compilers. Just-in-time compilers translate Java byte-codes into a sequence of native machine instructions at the time of method-activation (presupposed that the method's activation count exceeds a certain threshold). The compiled version is then cached for subsequent activations. According to Sun Microsystems

[Sun95], the quality of the generated code is “reasonably good” and “almost indistinguishable from native C or C++”.

Just-in-time compilers improve the situation relative to interpreted execution, but they still cannot compete with true optimizing compilers. Since code generation is carried out during a method’s activation, a just-in-time compiler needs to be fast so that the execution is disrupted only minimally. Unfortunately, most advanced optimization techniques are profoundly time-consuming.

The problem is worsened by the fact that most high-level information, although present in the source program, is lost in the transition to Java byte-codes. Not only is the reconstruction extraordinarily difficult but it is also time-intensive.

Hence, just-in-time compilers that take Java byte-codes as their input are not easily capable of performing aggressive optimizations such as global inlining, global trace scheduling, or code-parallelizations. This makes them intrinsically inferior to optimizing C or C++ compilers.

Since just-in-time compilers cannot always meet the required performance goals, many designers of time-critical Web-applications have to rule out the use of Java for their purposes and have to resort to natively compiled code. The current proliferation of native plug-ins (software programs that extend the capabilities of Web-browsers) rather than Java applets for high-performance applications clearly illustrates this point (e.g. Shockwave, PDFViewer, Live3D).

In this paper, we demonstrate that portability and high-performance are two goals that need not necessarily be irreconcilable. In the first part we describe an alternative intermediate representation that is based on high-level abstract syntax-trees rather than on low-level byte-codes. Abstract syntax-trees provide the necessary foundation for advanced code optimization and impose no artificial barriers to it since they do not require extensive information-reconstruction.

In the second part, we introduce the concepts of *dynamic runtime optimization* and *adaptive profiling*. In our system, a dynamic runtime optimizer performs code optimizations *continually* in the background based on runtime profile data. Performing optimizations in the background instead of at the time of method invocation eliminates the time constraint for code optimization faced by today’s just-in-time compilers. Basing compilation on an adaptive profiler further allows the code optimizer to make superior optimization decisions, improving even the quality of already optimized code on subsequent re-optimization iterations.

2 A Tree-Based Intermediate Representation

Rather than compiling source files into a sequence of Java byte-codes or into a register transfer language [Wal86], source files in our implementation are translated into an intermediate representation called *Slim Binaries* [Fra94, FK97]. The Slim Binary representation is based on abstract syntax-trees. It describes the actions of the original program similar to a parse tree in which every node is strongly typed by a reference to the symbol table. In contrast, abstract machine representations such as Java byte-codes are linear and type information as well

as the block structure of the program are only implicitly present and not directly accessible.

The Slim Binary representation, as its name suggests, is exceptionally dense, more so than compressed source code or compressed object code, accelerating the transfer of executable content over a network. It utilizes a variation of adaptive compression schemes, such as the popular LZW algorithm [Wel84], tailored towards syntax trees. It is based on the observation that different parts of programs often look very similar. As an example, expressions like `j++` or subexpressions like `...*pi/360` might be used several times within the same scope. The same holds for procedure calls. Procedures might be called repeatedly with similar parameter sets (e.g. `formatfloat(..., 10, 2)`). These similarities can be exploited by the use of a predictive algorithm that encodes recurring expressions and subexpressions efficiently both in terms of space and time.

In our implementation, the abstract syntax tree is reconstructed at load-time and a first version of native code is generated on-the-fly at this point. Because code-generation is performed at load-time, and because generating code takes more time than merely linking programs, we have built a code generating loader with the explicit design goal of fast loading times. In this context, the importance of Slim Binaries being compact becomes even more significant. The time saved by the faster downloading of object files compensates for the on-the-fly compilation phase. Measurements show that the resulting loading times are well within the range of what users are willing to tolerate—even for large applications. Surprisingly the goal of fast load-times does not even go at the expense of code quality. The code generated by our loader is comparable in quality to commercial C and Java just-in-time compilers. In contrast to Java interpreters and just-in-time compilers, however, the full native speed of applications is brought into action from the very beginning of executing an application.

At first sight, the disadvantage of Slim Binaries is that they cannot be easily interpreted at runtime. Their structure is less suited in the area for which Java was originally invented—embedded systems, and advanced consumer electronics. This area distinguishes itself by limited memory capacity and modest computing resources. However, this argument is less significant considering the recent increase of computing power and the recent reduction of memory-prices. For personal computers, this argument does not hold at all.

However, Slim Binaries have several advantages over Java byte-codes. First, Slim Binaries are much denser than Java byte-codes considerably reducing download-time over a network. Second, and more important, the information available in our intermediate representation is better suited for advanced and aggressive code optimizations. In contrast to Java, essential data about control-flow and data-flow is preserved in the abstract syntax-tree and does not have to be reconstructed prior to program optimization.

3 Advanced Code Optimizations

In a runtime environment that is based on byte-codes and just-in-time compilation such as Java's, two categories of optimizations are possible. The first category encompasses optimizations that are completely *independent* of the target architecture. They can entirely be performed at *compile-time* and on the level of the source language. For these optimizations, the time used for code-generation is not significant, because optimizations are carried out at compile-time. Examples are constant folding, dead-code elimination, loop-invariant code motion, and to some extent, even procedure inlining.

The second category comprises optimizations that depend on *target-specific* information (e.g. the number of available registers). Because this information is not available until *load-time*, these optimizations must operate on byte-code sequences at the time of method-invocation. Therefore, efficiency of compilation for these optimizations is crucial in order not to disrupt the running program while generating code. Optimizations are limited to elementary techniques such as rearranging instructions to achieve a better instruction mix, or eliminating unnecessary and expensive register-spills by smart register allocation algorithms. Peephole optimizations can also be classed with this category of optimizations.

The current state of the art in optimization techniques, however, is much further advanced than that. Highly beneficial optimizations are known to operate on the level of the source language and also to depend on processor specific information that is only available at load-time. These optimizations cannot currently be performed on the Java virtual machine architecture because reconstructing data-flow and control-flow information from byte-codes and performing these optimizations at the time of method-activation is too time-intensive. Cache blocking [WL91] and loop-unrolling are two examples of these techniques. Analyzing and recognizing access patterns, as well as having precise information about important cache parameters (e.g. cache size, line size) is a prerequisite for these optimizations. While the former can be accomplished at compile-time, the latter cannot be accomplished in practice. Value numbering [CS70] poses a similar problem. If done at all at compile-time, byte-code instructions that cannot be mapped to the underlying architecture on a one-to-one basis, but have to be translated into a sequence of native instructions (e.g. `invokevirtual`, `invokestatic`, `invokeinterface`) cannot reasonably be taken into consideration. Delaying value numbering until load-time is also impractical. As a further example, parallelizing instruction streams requires analyzing properties of data-structures that can only be realized at compile-time. However, important information about underlying hardware parameters (e.g. number of processors) is not available until load-time. Not being able to perform any of these optimizations is an immense disadvantage, which will play an ever greater role in the near future.

Moreover, Java byte-codes have additional disadvantages. Directly mapping byte-codes onto the underlying architecture is much more difficult than generating machine instructions from an abstract syntax-tree. Code generators that are based on a high-level representation do not have to deal with unfavorable peculiarities of Java byte-codes but can tailor their output towards advanced and

specific processor features, such as special purpose instructions, size of register sets, and cache architectures. This is especially true for today's RISC processors that are less well suited for modeling the stack operations that are used in byte-codes. Whether dedicated Java processors, such as Sun Microsystems recently announced picoJava architecture, will overcome this disadvantage is still an open question.

In contrast to Java byte-codes, Slim Binaries are equally well suited for *all categories* of code optimizations and do not have to deal with any of the disadvantages of byte-codes. At the time of loading, the abstract syntax tree, which can be efficiently decoded, contains the same amount of information as is available at compile-time. Not only is the control-flow and data-flow of programs preserved, but also the structure and property of data-structures and data-types. This information is essential for aggressive code optimizations and does not have to be recomputed from a lower-level representation in a time-intensive process.

4 Runtime Optimization and Adaptive Profiling

Performing optimizations at load-time or at the time of method-activation is usually not feasible because of the mentioned lack of time to perform these optimizations. In many cases optimization takes at least 5 times as long as just compiling the program [Bra95]. This might be feasible for small applications, or large numerical applications in which the time saved by the optimizations is much more substantial than the additional time required to optimize the program. For all other applications, however, a different solution is necessary.

Therefore, in our design (Fig. 1), program optimization is performed in the background at runtime, taking advantage of idle cycles (in a typical GUI-based environment, we measured processor idle times of more than 90%). At load-time, the fast code-generating loader transforms the intermediate representation into a first unoptimized code-image. The optimizer then continually generates faster versions of the program in the background, replacing older code-images. This step is repeated until a fixed-point is reached and further optimizations do not increase the overall system performance anymore.

Performing optimizations at runtime also facilitates a completely new set of *intermodular* optimizations. Since the configuration of the system (i.e. which components are active, and how they interact) is known at runtime, optimizations are not restricted to local algorithms. As previous studies have shown, the impact of intermodular optimizations on runtime performance can be dramatic in some cases [Höl94]. Examples of intermodular optimizations are intermodular inlining, intermodular register allocation, and global cache optimizations [Kis96].

Runtime optimization is only one aspect of our architecture. Equally important is the adaptive profiler that continually collects information about the system's runtime behavior. The profiler's primary goal is to pinpoint the program parts that account for most of the execution time. That way, optimizations can be concentrated on high payoff areas rather than being applied uniformly to each section of the program. Less frequently executed sections are optimized

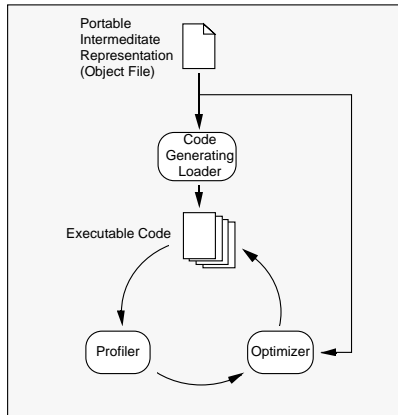


Fig. 1. Architecture

sparsely, and no optimization is performed on rarely executed sections or sections in which optimizations would not yield profitable results.

Further, with the availability of accurate profiling-data at the time of optimization, the optimizer never has to resort to inexact heuristics. This often leads to superior results. Many of today’s aggressive optimization algorithms are based on heuristics, in order to achieve good results. However, this can be a double-edged sword. On the one hand, if the system’s runtime behavior is properly predicted, considerable performance increases may be expected. On the other hand, if predictions do not come true, these optimizations will actually lead to *performance penalties*. As an example, in trace scheduling, traces (also called execution-paths) are selected and scheduled in decreasing order of their execution frequency. The most frequently executed path is scheduled first, as if it were one single basic block. However, in order to preserve semantic correctness, “compensation code” has to be introduced in off-trace paths. If, at runtime, the trace which was assumed to be executed most often is indeed executed most of the time, this optimization yields superior results. If that is not the case, and off-trace paths are executed more often, then this optimization will deteriorate the overall performance. Other examples that highly depend on heuristics are loop-unrolling which depends on loop-frequency estimates and cache parameters, or inlining and partial evaluation which depend on call-frequency estimates.

In order to make our profiler as unobtrusive as possible, it uses a combination of dynamic instrumentation of the object code and statistical profiling techniques. It also varies the granularity at which it monitors the system’s execution, and is only applied when it can contribute to the overall system performance, pushing the profiling overhead below 2%. Previous studies have reported profiling overheads of 2%–91% [ABD+97, BL94].

5 Results

In the last few months, we have implemented an experimental system that is based on our proposed architecture. The system, named “Juice”, enables the seamless integration of Slim Binary encoded executables into HTML-pages. It is based on a family of Netscape plug-ins that contain an on-the-fly code generator. Juice is currently publicly available for Intel based computers running Windows 95, and for PowerPC based Macintosh computers.

Besides being reliable and simple to use, Juice is also efficient. Table 1 shows time-measurements for basic operations, such as assignments, additions, and method calls. The benchmark was executed on an Intel Pentium processor clocked at 166Mhz (Dell OptiPlex GXM 5166). Since neither the optimizer nor the profiler have yet been fully implemented and integrated into the Juice architecture, they have not been taken into account for all of the benchmarks (this special configuration that only applies on-the-fly compilation but no optimizations is subsequently called Juice Level I). Juice does very well in comparison to just-in-time compilers. The runtime-differences are only minimal. Both runtime systems achieve an average speed-up factor of 12 to 18 in contrast to byte-code interpretation.

	Internet Explorer 3.0 (Interpreted)	Netscape Navigator 3.0 (Just-In-Time)	Internet Explorer (Just-In-Time)	Juice Level I (No Opt.)
Local Var. Assignment	0.220	0.011	0.006	0.015
Instance Var. Assignment	0.440	0.010	0.007	0.046
Array Elem. Assignment	0.590	0.050	0.051	0.045
Byte Addition	0.680	0.044	0.030	0.021
Short Addition	0.660	0.044	0.030	0.047
Int Addition	0.570	0.015	0.013	0.017
Float Addition	0.570	0.046	0.045	0.054
Double Addition	0.500	0.140	0.044	0.110
Method Call	1.500	0.092	0.091	0.120
Average	0.637	0.050	0.035	0.053

Table 1. Suite 1—Basic Operations. All numbers are given in microseconds per operation.

Yet, speed-up factors in this range are not realistic in most cases. This has to be attributed to the fact that larger applications often call Java library routines that are distributed as native binaries, already optimized for speed. The more native libraries are called, the less just-in-time compilers boost performance. The results of the second benchmark, which comprises of several long-running, computational intensive tasks, emphasizes this statement. It compares the execution times of Juice, just-in-time compilers, and optimized C++ to the execution time of byte-code interpretation (Table 2). The speed-ups are remarkably smaller than the ones measured in the first test suite. Performance comprehensibly degrades with the number of library calls down to disappointing ratios of

2:1–4:1. Examples are the “String Sort” benchmark that frequently invokes the native “System.arraycopy” method and the “Fourier Analysis” benchmark that frequently calls the math library (Math.sin, Math.cos, Math.exp).

	Netscape Navigator 3.0 (Just-In-Time)	Internet Explorer 3.0 (Just-In-Time)	Juice Level I (No Opt.)	C++ Optimized
Numeric Sort	11.17	13.21	9.63	79.69
String Sort	3.50	4.72	1.55	6.70
Bitfield Operations	17.61	15.92	20.82	64.94
Fourier Analysis	0.87	2.76	2.45	4.27
IDEA Encryption	4.54	3.24	6.69	16.30
Huffman Compression	11.87	16.14	20.68	35.32
LU Decomposition	7.63	7.18	6.18	36.69
Average	8.43	9.19	10.05	35.33

Table 2. Suite 2—Computational Intensive Operations. All numbers are given in multiples of the performance of interpreted byte-codes using Internet Explorer.

This benchmark also unequivocally demonstrates that native code, compiled from an abstract syntax tree, is at least equivalent in quality to code generated by just-in-time compilers (Fig. 2). In some cases, the results even surpass the fastest available Java runtime systems—notably without applying any optimizations or profiling. However, the benchmarks also clearly demonstrate the current deficiencies of just-in-time compilers—they cannot yet compete with true optimizing compilers. Optimized C++ code is still much faster. In order to narrow this gap, aggressive and advanced optimizations are a necessity. The proposed tree-based intermediate representation coupled with background optimization fulfills all the requirements for achieving this goal.

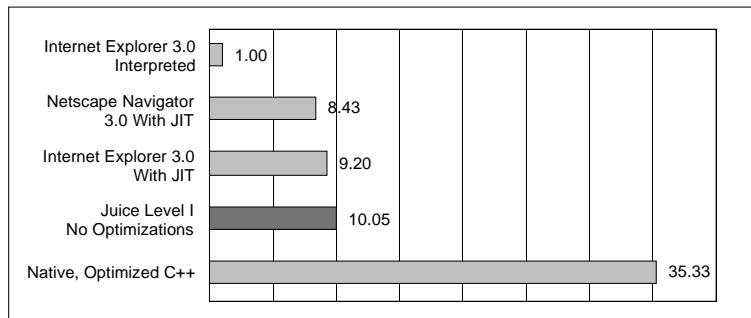


Fig. 2. Average Speed-Up Results

One of the initial claims of this paper was that a tree-based intermediate

representation not only provides the basis for closing the efficiency gap between Java byte-codes and optimized C++ code, but also reduces the overhead for transferring files over a network. Not only are Slim Binary object files more than twice as dense as Java class files, and a factor of 3 to 4 smaller than traditional native object files, Slim Binaries are even smaller than compressed native code. Figure 3 summarizes the results for the above test suite (consisting of 12 source files, approximately 130kBytes in size).

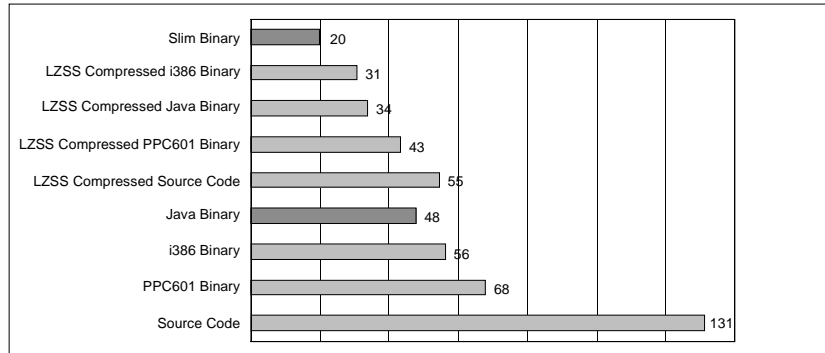


Fig. 3. Size Comparison Between Different Distribution Formats. All numbers are given in kBytes.

Although Slim Binaries reduce network traffic by a factor of 2 in terms of network packets, the differences between downloading Java class files and Juice Slim Binaries are almost indistinguishable in terms of download-time. This stems from the fact that for small execution units, it is mostly the time to set up the connection to a server that accounts for most of the waiting-time. However, with the introduction of *component packaging* concepts (compact archive formats for packaging the components of a Java or Juice application) the size of executable units will become more significant.

Finally, we have also measured the time that is required to compile the abstract syntax tree into native code. As mentioned earlier, the compilation time is hardly noticeable to the user. On a PowerPC based computer (Power Macintosh 8500/120) it takes approximately 470ms to compile the 12 benchmark-files. In comparison to the time required to download these files using a fast connection (4s), or using a slow connection (40s), the overhead of on-the-fly code generation can be neglected.

6 Related Work

The most notable related project in the area of portable intermediate representations is TDF [OSF91, DRA93], the result of the *architecture neutral distribution*

format (ANDF) initiative by the *Open Software Foundation* (OSF). Like Slim Binaries, TDF is also based on a tree-structured intermediate language in conjunction with an embedded symbol table. In contrast to Slim Binaries however, TDF files are generally around twice the size of the native binary code and individual modules that comprise an application program are linked together statically and are then translated into native code in an *off-line* process. Conversely, Slim Binaries have been designed to support *dynamic module loading*, which is a prerequisite for user-serviceable plug-in software components.

In the area of program reoptimization, pioneering research was done by Hansen [Han74]. He was the first person to build a fully automated optimization system. For speed considerations, his system was based on a non-portable intermediate representation that could be interpreted directly. Only when exceeding a certain runtime threshold was the representation translated just-in-time to native code and optimized for speed. Several similar systems have since then been described in literature. Among them are work by Wall [Wal92], work on the SELF system [Höl94], DIGITAL's FX!32 technology [HH97], and the Morph system from Harvard [ZWG96+]. Although some of these systems perform dynamic adaptive optimizations at run-time and use profiling data to improve optimization results, none of the aforementioned systems was developed in the context of portable executable formats.

7 Conclusion

In the last few months, Java and the Java Virtual Machine have become a standard environment for building portable Internet “applets” and applications. However, despite their success, their lack of performance makes them ill-suited for many performance critical applications. Although just-in-time compilers try to remedy this situation and achieve speed-ups of 10–15 compared to interpreted byte-codes, they still cannot compete with true optimizing compilers.

In this paper, we have shown that portability and high performance need not necessarily be irreconcilable. We have proposed an alternative intermediate representation that is based on abstract syntax trees. Our intermediate representation is twice as dense as Java byte-codes and does not require expensive information reconstruction prior to code-optimization. We have also shown that by delaying code-optimization, the timing-constraints witnessed by today's just-in-time compilers can be eliminated. This facilitates advanced and aggressive code-optimizations.

We have shown that our current version of the proposed architecture can compete with today's fastest just-in-time compilers, although no optimizations have yet been implemented. With the availability and integration of the runtime optimizer and the adaptive profiler, we will be able to level the performance with optimized C++ and dismantle the current performance deficiency of portable transportability schemes.

Additional information about the Juice project and related research topics at the University of California, Irvine can be found on the World Wide Web at <http://www.ics.uci.edu/~juice>.

8 Acknowledgement

Part of this work is being funded by the National Science Foundation under grant CCR-97014000.

References

- [ABD+97] J. Anderson, L. Berc, J. Dean, S. Ghemawat, M. Henzinger, S-T. Leung, R. Sites, M. Vandevoorde, C. Waldspurger, and W. Wehl. Continuous Profiling: Where Have All the Cycles Gone? In *Proceedings of the 16th ACM Symposium on Operating Systems Principles*, St. Malo, France, October 1997. Also published as Technical Note 1997-016.
- [Bra95] M. M. Brandis. *Optimizing Compilers for Structured Programming Languages*. (Doctoral Dissertation) Eidgenössische Technische Hochschule Zürich, 1995.
- [BL94] T. Ball, J. R. Larus. Optimally Profiling and Tracing Programs. In *ACM Transactions on Programming Languages and Systems*, 16(4), pp 1319–1360, July 1994.
- [CS70] J. Cocke, J. Schwartz. *Programming Languages and Their Compilers: Preliminary Notes*. Courant Institute of Mathematical Sciences, New York University, April 1970.
- [DRA93] United Kingdom Defence Resesearch Agency. *TDF Specification, Issue 2.1*. June 1993.
- [OSF91] Open Software Foundation. *OSF Architecture-Neutral Distribution Format Rationale*. 1991.
- [FK97] M. Franz and T. Kistler; Slim Binaries. In *Communications of the ACM*, 40(12), pp 87-94, December 1997. Also published as Technical Report No. 96-24, Department of Information and Computer Science, University of California, Irvine, June 1996.
- [Fra94] M. Franz. *Code-Generation On-the-Fly: A Key to Portable Software*. (Doctoral Dissertation) Verlag der Fachvereine, Zürich, 1994.
- [Han74] G. Hansen. *Adaptive Systems for the Dynamic Run-Time Optimization of Programs*. Ph.D. Dissertation, Department of Computer Science, Carnegie-Mellon University, 1974.
- [HH97] R. Hookway and M. Herdeg. *DIGITAL FX!32: Combining Emulation and Binary Translation*. In *Digital Technical Journal* 9(1):3–12, 1997.

- [Höl94] U. Hölzle. *Adaptive Optimization for SELF: Reconciling High Performance with Exploratory Programming*. (Ph.D. Dissertation) Department of Computer Science, Stanford University, 1994.
- [Kis96] T. Kistler. Dynamic Runtime Optimization. In *Proceedings of the Joint Modular Languages Conference, JMLC'97*, pp 53-66. Published as Springer Lecture Notes in Computer Science No. 1204, March 1997. Also published as Technical Report No. 96-54, Department of Information and Computer Science, University of California, Irvine, November 1996.
- [LYJ96] T. Lindholm, F. Yellin, B. Joy, K. Walrath. *The Java Virtual Machine Specification*. Addison-Wesley, 1996.
- [Mot93] Motorola, Inc. *PowerPC 601: RISC Microprocessor User's Manual*. 1993.
- [Sun95] Sun Microsystems. *The Java Language: An Overview*. <http://java.sun.com/doc/Overviews/java/java-overview-1.html>, 1995.
- [Wal86] D. W. Wall. Global Register Allocation at Link Time. In *Proceedings of SIGPLAN '86 Symposium on Compiler Construction*, pp 264–275, July 1986.
- [Wal92] D. Wall. *Systems for Late Code Modification*. WRL Research Report 92/3, Digital Equipment Corporation, Western Research Laboratory, Palo Alto CA, May 1992.
- [Wel84] T. A. Welch. A Technique for High-Performance Data Compression. *IEEE Computer*, 17(6), pp 8–19, June 1984.
- [Wir88] N. Wirth. The Programming Language Oberon. In *Software-Practice and Experience* 18(7), pp 671–690, July 1988.
- [WL91] M. Wolf, M. Lam. A Data Locality Optimization Algorithm. In *Proceedings of the SIGPLAN '91 Conference on Programming Language Design and Implementation*, pp 30–44, Published as *SIGPLAN Notices* 26(6), June 1991.
- [ZWG+97] X. Zhang, Z. Wang, N. Gloy, J. Chen, and M. Smith. System Support for Automatic Profiling and Optimization. In *Proceedings of the 16th ACM Symposium of Operating Systems Principles*, September 1997.