

A Brief History of Just-In-Time

JOHN AYCOCK

University of Calgary

Software systems have been using “just-in-time” compilation (JIT) techniques since the 1960s. Broadly, JIT compilation includes any translation performed dynamically, after a program has started execution. We examine the motivation behind JIT compilation and constraints imposed on JIT compilation systems, and present a classification scheme for such systems. This classification emerges as we survey forty years of JIT work, from 1960–2000.

Categories and Subject Descriptors: D.3.4 [Programming Languages]: Processors; K.2 [History of Computing]: Software

General Terms: Languages, Performance

Additional Key Words and Phrases: Just-in-time compilation, dynamic compilation

1. INTRODUCTION

Those who cannot remember the past are condemned to repeat it.

George Santayana, 1863–1952 [Bartlett 1992]

This oft-quoted line is all too applicable in computer science. Ideas are generated, explored, set aside—only to be reinvented years later. Such is the case with what is now called “just-in-time” (JIT) or dynamic compilation, which refers to translation that occurs after a program begins execution.

Strictly speaking, JIT compilation systems (“JIT systems” for short) are completely unnecessary. They are only a means to improve the time and space efficiency of programs. After all, the central problem JIT systems address is a solved one: translating programming languages

into a form that is executable on a target platform.

What is translated? The scope and nature of programming languages that require translation into executable form covers a wide spectrum. Traditional programming languages like Ada, C, and Java are included, as well as little languages [Bentley 1988] such as regular expressions.

Traditionally, there are two approaches to translation: compilation and interpretation. Compilation translates one language into another—C to assembly language, for example—with the implication that the translated form will be more amenable to later execution, possibly after further compilation stages. Interpretation eliminates these intermediate steps, performing the same analyses as compilation, but performing execution immediately.

This work was supported in part by a grant from the National Science and Engineering Research Council of Canada.

Author’s address: Department of Computer Science, University of Calgary, 2500 University Dr. N. W., Calgary, Alta., Canada T2N 1N4; email: aycock@cpsc.ucalgary.ca.

Permission to make digital/hard copy of part or all of this work for personal or classroom use is granted without fee provided that the copies are not made or distributed for profit or commercial advantage, the copyright notice, the title of the publication, and its date appear, and notice is given that copying is by permission of ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists requires prior specific permission and/or a fee.

©2003 ACM 0360-0300/03/0600-0097 \$5.00

JIT compilation is used to gain the benefits of both (static) compilation and interpretation. These benefits will be brought out in later sections, so we only summarize them here:

- Compiled programs run faster, especially if they are compiled into a form that is directly executable on the underlying hardware. Static compilation can also devote an arbitrary amount of time to program analysis and optimization. This brings us to the primary constraint on JIT systems: speed. A JIT system must not cause untoward pauses in normal program execution as a result of its operation.
- Interpreted programs are typically smaller, if only because the representation chosen is at a higher level than machine code, and can carry much more semantic information implicitly.
- Interpreted programs tend to be more portable. Assuming a machine-independent representation, such as high-level source code or virtual machine code, only the interpreter need be supplied to run the program on a different machine. (Of course, the program still may be doing nonportable operations, but that's a different matter.)
- Interpreters have access to run-time information, such as input parameters, control flow, and target machine specifics. This information may change from run to run or be unobtainable prior to run-time. Additionally, gathering some types of information about a program before it runs may involve algorithms which are undecidable using static analysis.

To narrow our focus somewhat, we only examine software-based JIT systems that have a nontrivial translation aspect. Keppel et al. [1991] eloquently built an argument for the more general case of run-time code generation, where this latter restriction is removed.

Note that we use the term *execution* in a broad sense—we call a program representation executable if it can be executed by the JIT system in any manner, either

directly as in machine code, or indirectly using an interpreter.

2. JIT COMPILATION TECHNIQUES

Work on JIT compilation techniques often focuses around implementation of a particular programming language. We have followed this same division in this section, ordering from earliest to latest where possible.

2.1. Genesis

Self-modifying code has existed since the earliest days of computing, but we exclude that from consideration because there is typically no compilation or translation aspect involved.

Instead, we suspect that the earliest published work on JIT compilation was McCarthy's [1960] LISP paper. He mentioned compilation of functions into machine language, a process fast enough that the compiler's output needn't be saved. This can be seen as an inevitable result of having programs and data share the same notation [McCarthy 1981].

Another early published reference to JIT compilation dates back to 1966. The University of Michigan Executive System for the IBM 7090 explicitly notes that the assembler [University of Michigan 1966b, p. 1] and loader [University of Michigan 1966a, p. 6] can be used to translate and load during execution. (The manual's preface says that most sections were written before August 1965, so this likely dates back further.)

Thompson's [1968] paper, published in *Communications of the ACM*, is frequently cited as "early work" in modern publications. He compiled regular expressions into IBM 7094 code in an ad hoc fashion, code which was then executed to perform matching.

2.2. LC²

The Language for Conversational Computing, or LC², was designed for interactive programming [Mitchell et al. 1968]. Although used briefly at Carnegie-Mellon University for teaching, LC² was

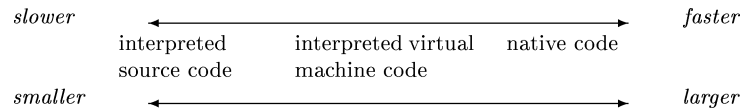


Fig. 1. The time-space tradeoff.

primarily an experimental language [Mitchell 2000]. It might otherwise be consigned to the dustbin of history, if not for the techniques used by Mitchell in its implementation [Mitchell 1970], techniques that later influenced JIT systems for Smalltalk and Self.

Mitchell observed that compiled code can be derived from an interpreter at run-time, simply by storing the actions performed during interpretation. This only works for code that has been executed, however—he gave the example of an if-then-else statement, where only the else-part is executed. To handle such cases, code is generated for the unexecuted part which reinvokes the interpreter should it ever be executed (the then-part, in the example above).

2.3. APL

The seminal work on efficient APL implementation is Abrams’ dissertation [Abrams 1970]. Abrams concocted two key APL optimization strategies, which he described using the connotative terms *drag-along* and *beating*. Drag-along defers expression evaluation as long as possible, gathering context information in the hopes that a more efficient evaluation method might become apparent; this might now be called *lazy evaluation*. Beating is the transformation of code to reduce the amount of data manipulation involved during expression evaluation.

Drag-along and beating relate to JIT compilation because APL is a very dynamic language; types and attributes of data objects are not, in general, known until run-time. To fully realize these optimizations’ potential, their application must be delayed until run-time information is available.

Abrams’ “APL Machine” employed two separate JIT compilers. The first trans-

lated APL programs into postfix code for a D-machine,¹ which maintained a buffer of deferred instructions. The D-machine acted as an “algebraically simplifying compiler” [Abrams 1970, p. 84] which would perform drag-along and beating at run-time, invoking an E-machine to execute the buffered instructions when necessary.

Abrams’ work was directed toward an architecture for efficient support of APL, hardware support for high-level languages being a popular pursuit of the time. Abrams never built the machine, however; an implementation was attempted a few years later [Schroeder and Vaughn 1973].² The techniques were later expanded upon by others [Miller 1977], although the basic JIT nature never changed, and were used for the software implementation of Hewlett-Packard’s APL\3000 [Johnston 1977; van Dyke 1977].

2.4. Mixed Code, Throw-Away Code, and BASIC

The tradeoff between execution time and space often underlies the argument for JIT compilation. This tradeoff is summarized in Figure 1. The other consideration is that most programs spend the majority of time executing a minority of code, based on data from empirical studies [Knuth 1971]. Two ways to reconcile these observations have appeared: mixed code and throw-away compiling.

Mixed code refers to the implementation of a program as a mixture of native code and interpreted code, proposed independently by Dakin and Poole [1973] and Dawson [1973]. The frequently executed parts of the program would be

¹ Presumably *D* stood for *Deferral* or *Drag-Along*.

² In the end, Litton Industries (Schroeder and Vaughn’s employer) never built the machine [Mauriello 2000].

in native code, the infrequently executed parts interpreted, hopefully yielding a smaller memory footprint with little or no impact on speed. A fine-grained mixture is implied: implementing the program with interpreted code and the libraries with native code would *not* constitute mixed code.

A further twist to the mixed code approach involved customizing the interpreter [Pittman 1987]. Instead of mixing native code into the program, the native code manifests itself as special virtual machine instructions; the program is then compiled entirely into virtual machine code.

The basic idea of mixed code, switching between different types of executable code, is still applicable to JIT systems, although few researchers at the time advocated generating the machine code at run-time. Keeping both a compiler and an interpreter in memory at run-time may have been considered too costly on the machines of the day, negating any program size tradeoff.

The case against mixed code comes from software engineering [Brown 1976]. Even assuming that the majority of code will be shared between the interpreter and compiler, there are still two disparate pieces of code (the interpreter proper and the compiler's code generator) which must be maintained and exhibit identical behavior.

(Proponents of partial evaluation, or program specialization, will note that this is a specious argument in some sense, because a compiler can be thought of as a specialized interpreter [Jones et al. 1993]. However, the use of partial evaluation techniques is not currently widespread.)

This brings us to the second manner of reconciliation: throw-away compiling [Brown 1976]. This was presented purely as a space optimization: instead of static compilation, parts of a program could be compiled dynamically on an as-needed basis. Upon exhausting memory, some or all of the compiled code could be thrown away; the code would be regenerated later if necessary.

BASIC was the testbed for throw-away compilation. Brown [1976] essentially characterized the technique as a

good way to address the time-space trade-off; Hammond [1977] was somewhat more adamant, claiming throw-away compilation to be superior except when memory is tight.

A good discussion of mixed code and throw-away compiling may be found in Brown [1990].

2.5. FORTRAN

Some of the first work on JIT systems where programs automatically optimize their “hot spots” at run-time was due to Hansen [1974].³ He addressed three important questions:

- (1) What code should be optimized? Hansen chose a simple, low-cost frequency model, maintaining a frequency-of-execution counter for each block of code (we use the generic term *block* to describe a unit of code; the exact nature of a block is immaterial for our purposes).
- (2) When should the code be optimized? The frequency counters served a second rôle: crossing a threshold value made the associated block of code a candidate for the next “level” of optimization, as described below. “Supervisor” code was invoked between blocks, which would assess the counters, perform optimization if necessary, and transfer control to the next block of code. The latter operation could be a direct call, or interpreter invocation—mixed code was supported by Hansen’s design.
- (3) How should the code be optimized? A set of conventional machine-independent and machine-dependent optimizations were chosen and ordered, so a block might first be optimized by constant folding, by common subexpression elimination the second

³ Dawson [1973] mentioned a 1967 report by Barbieri and Morrissey where a program begins execution in interpreted form, and frequently executed parts “can be converted to machine code.” However, it is not clear if the conversion to machine code occurred at run-time. Unfortunately, we have not been able to obtain the cited work as of this writing.

time optimization occurs, by code motion the third time, and so on. Hansen [1974] observed that this scheme limits the amount of time taken at any given optimization point (especially important if the frequency model proves to be incorrect), as well as allowing optimizations to be incrementally added to the compiler.

Programs using the resulting Adaptive FORTRAN system reportedly were not always faster than their statically compiled-and-optimized counterparts, but performed better overall.

Returning again to mixed code, Ng and Cantoni [1976] implemented a variant of FORTRAN using this technique. Their system could compile functions at run-time into “pseudo-instructions,” probably a tokenized form of the source code rather than a lower-level virtual machine code. The pseudo-instructions would then be interpreted. They claimed that run-time compilation was useful for some applications and avoided a slow compile-link process. They did not produce mixed code at run-time; their use of the term referred to the ability to have statically compiled FORTRAN programs call their pseudo-instruction interpreter automatically when needed via linker trickery.

2.6. Smalltalk

Smalltalk source code is compiled into virtual machine code when new methods are added to a class [Goldberg and Robson 1985]. The performance of naive Smalltalk implementations left something to be desired, however.

Rather than attack the performance problem with hardware, Deutsch and Schiffman [1984] made key optimizations in software. The observation behind this was that they could pick the most efficient representation for information, so long as conversion between representations happened automatically and transparently to the user.

JIT conversion of virtual machine code to native code was one of the optimization techniques they used, a process they

likened to macro-expansion. Procedures were compiled to native code lazily, when execution entered the procedure; the native code was cached for later use. Their system was linked to memory management in that native code would never be paged out, just thrown away and regenerated later if necessary.

In turn, Deutsch and Schiffman [1984] credited the dynamic translation idea to Rau [1978]. Rau was concerned with “universal host machines” which would execute a variety of high-level languages well (compared to, say, a specialized APL machine). He proposed dynamic translation to microcode at the granularity of single virtual machine instructions. A hardware cache, the dynamic translation buffer, would store completed translations; a cache miss would signify a missing translation, and fault to a dynamic translation routine.

2.7. Self

The Self programming language [Ungar and Smith 1987; Smith and Ungar 1995], in contrast to many of the other languages mentioned in this section, is primarily a research vehicle. Self is in many ways influenced by Smalltalk, in that both are pure object-oriented languages—everything is an object. But Self eschews classes in favor of prototypes, and otherwise attempts to unify a number of concepts. Every action is dynamic and changeable, and even basic operations, like local variable access, require invocation of a method. To further complicate matters, Self is a dynamically-typed language, meaning that the types of identifiers are not known until run-time.

Self’s unusual design makes efficient implementation difficult. This resulted in the development of the most aggressive, ambitious JIT compilation and optimization up to that time. The Self group noted three distinct generations of compiler [Hölzle 1994], an organization we follow below; in all cases, the compiler was invoked dynamically upon a method’s invocation, as in Deutsch and Schiffman’s [1984] Smalltalk system.

2.7.1. First Generation. Almost all the optimization techniques employed by Self compilers dealt with type information, and transforming a program in such a way that some certainty could be had about the types of identifiers. Only a few techniques had a direct relationship with JIT compilation, however.

Chief among these, in the first-generation Self compiler, was customization [Chambers et al. 1989; Chambers and Ungar 1989; Chambers 1992]. Instead of dynamically compiling a method into native code that would work for any invocation of the method, the compiler produced a version of the method that was customized to that particular context. Much more type information was available to the JIT compiler compared to static compilation, and by exploiting this fact the resulting code was much more efficient. While method calls from similar contexts could share customized code, “overcustomization” could still consume a lot of memory at run-time; ways to combat this problem were later studied [Dieckmann and Hölzle 1997].

2.7.2. Second Generation. The second-generation Self compiler extended one of the program transformation techniques used by its predecessor, and computed much better type information for loops [Chambers and Ungar 1990; Chambers 1992].

This Self compiler’s output was indeed faster than that of the first generation, but it came at a price. The compiler ran 15 to 35 times more slowly on benchmarks [Chambers and Ungar 1990, 1991], to the point where many users refused to use the new compiler [Hölzle 1994]!

Modifications were made to the responsible algorithms to speed up compilation [Chambers and Ungar 1991]. One such modification was called *deferred compilation of uncommon cases*.⁴ The compiler

⁴ In Chambers’ thesis, this is referred to as “lazy compilation of uncommon branches,” an idea he attributes to a suggestion by John Maloney in 1989 [Chambers 1992, p. 123]. However, this is the same technique used in Mitchell [1970], albeit for different reasons.

is informed that certain events, such as arithmetic overflow, are unlikely to occur. That being the case, no code is generated for these uncommon cases; a stub is left in the code instead, which will invoke the compiler again if necessary. The practical result of this is that the code for uncommon cases need not be analyzed upon initial compilation, saving a substantial amount of time.⁵

Ungar et al. [1992] gave a good presentation of optimization techniques used in Self and the resulting performance in the first- and second-generation compilers.

2.7.3. Third Generation. The third-generation Self compiler attacked the issue of slow compilation at a much more fundamental level. The Self compiler was part of an interactive, graphical programming environment; executing the compiler on-the-fly resulted in a noticeable pause in execution. Hölzle argued that measuring pauses in execution for JIT compilation by timing the amount of time the compiler took to run was deceptive, and not representative of the user’s experience [Hölzle 1994; Hölzle and Ungar 1994b]. Two invocations of the compiler could be separated by a brief spurt of program execution, but would be perceived as one long pause by the user. Hölzle compensated by considering temporally related groups of pauses, or “pause clusters,” rather than individual compilation pauses.

As for the compiler itself, compilation time was reduced—or at least spread out—by using adaptive optimization, similar to Hansen’s [1974] FORTRAN work. Initial method compilation was performed by a fast, nonoptimizing compiler; frequency-of-invocation counters were kept for each method to determine when recompilation should occur [Hölzle 1994; Hölzle and Ungar 1994a, 1994b]. Hölzle makes an interesting comment on this mechanism:

... in the course of our experiments we discovered that the trigger mechanism (“when”) is

⁵ This technique can be applied to dynamic compilation of exception handling code [Lee et al. 2000].

much less important for good recompilation results than the selection mechanism (“what”). [Hölzle 1994, p. 38]⁶

This may come from the slightly counterintuitive notion that the best candidate for recompilation is *not* necessarily the method whose counter triggered the recompilation. Object-oriented programming style tends to encourage short methods; a better choice may be to (re)optimize the method’s caller and incorporate the frequently invoked method inline [Hölzle and Ungar 1994b].

Adaptive optimization adds the complication that a modified method may already be executing, and have information (such as an activation record on the stack) that depends on the previous version of the modified method [Hölzle 1994]; this must be taken into consideration.⁷

The Self compiler’s JIT optimization was assisted by the introduction of “type feedback” [Hölzle 1994; Hölzle and Ungar 1994a]. As a program executed, type information was gathered by the run-time system, a straightforward process. This type information would then be available if and when recompilation occurred, permitting more aggressive optimization. Information gleaned using type feedback was later shown to be comparable with, and perhaps complementary to, information from static type inference [Agesen and Hölzle 1995; Agesen 1996].

2.8. Slim Binaries and Oberon

One problem with software distribution and maintenance is the heterogeneous computing environment in which software runs: different computer architectures require different binary executables. Even within a single line of backward-compatible processors, many variations in capability can exist; a program statically

compiled for the least-common denominator of processor may not take full advantage of the processor on which it eventually executes.

In his doctoral work, Franz addressed these problems using “slim binaries” [Franz 1994; Franz and Kistler 1997]. A slim binary contains a high-level, machine-independent representation⁸ of a program module. When a module is loaded, executable code is generated for it on-the-fly, which can presumably tailor itself to the run-time environment. Franz, and later Kistler, claimed that generating code for an entire module at once was often superior to the method-at-a-time strategy used by Smalltalk and Self, in terms of the resulting code performance [Franz 1994; Kistler 1999].

Fast code generation was critical to the slim binary approach. Data structures were delicately arranged to facilitate this; generated code that could be reused was noted and copied if needed later, rather than being regenerated [Franz 1994].

Franz implemented slim binaries for the Oberon system, which allows dynamic loading of modules [Wirth and Gutknecht 1989]. Loading and generating code for a slim binary was not faster than loading a traditional binary [Franz 1994; Franz and Kistler 1997], but Franz argued that this would eventually be the case as the speed discrepancy between processors and input/output (I/O) devices increased [Franz 1994].

Using slim binaries as a starting point, Kistler’s [1999] work investigated “continuous” run-time optimization, where parts of an executing program can be optimized *ad infinitum*. He contrasted this to the adaptive optimization used in Self, where optimization of methods would eventually cease.

Of course, reoptimization is only useful if a new, better, solution can be obtained; this implies that continuous optimization is best suited to optimizations whose input varies over time with the program’s

⁶ The same comment, with slightly different wording, also appears in Hölzle and Ungar [1994a, p. 328].

⁷ Hansen’s work in 1974 could ignore this possibility; the FORTRAN of the time did not allow recursion, and so activation records and a stack were unnecessary [Sebesta 1999].

⁸ This representation is an abstract syntax tree, to be precise.

execution.⁹ Accordingly, Kistler looked at cache optimizations—rearranging fields in a structure dynamically to optimize a program’s data-access patterns [Kistler 1999; Kistler and Franz 1999]—and a dynamic version of trace scheduling, which optimizes based on information about a program’s control flow during execution [Kistler 1999].

The continuous optimizer itself executes in the background, as a separate low-priority thread which executes only during a program’s idle time [Kistler 1997, 1999]. Kistler used a more sophisticated metric than straightforward counters to determine when to optimize, and observed that deciding *what* to optimize is highly optimization-specific [Kistler 1999].

An idea similar to continuous optimization has been implemented for Scheme. Burger [1997] dynamically reordered code blocks using profile information, to improve code locality and hardware branch prediction. His scheme relied on the (copying) garbage collector to locate pointers to old versions of a function, and update them to point to the newer version. This dynamic recompilation process could be repeated any number of times [Burger 1997, page 70].

2.9. Templates, ML, and C

ML and C make strange bedfellows, but the same approach has been taken to dynamic compilation in both. This approach is called *staged compilation*, where compilation of a single program is divided into two stages: static and dynamic compilation. Prior to run-time, a static compiler compiles “templates,” essentially building blocks which are pieced together at run-time by the dynamic compiler, which may also place run-time values into holes left in the templates. Typically these templates are specified by user annotations, although some work has been done on deriving them automatically [Mock et al. 1999].

⁹ Although, making the general case for run-time optimization, he discussed intermodule optimizations where this is not the case [Kistler 1997].

As just described, template-based systems arguably do not fit our description of JIT compilers, since there would appear to be no nontrivial translation aspect. However, templates may be encoded in a form which requires run-time translation before execution, or the dynamic compiler may perform run-time optimizations after connecting the templates.

Templates have been applied to (subsets of) ML [Leone and Lee 1994; Lee and Leone 1996; Wickline et al. 1998]. They have also been used for run-time specialization of C [Consel and Noël 1996; Marlet et al. 1999], as well as dynamic extensions of C [Auslander et al. 1996; Engler et al. 1996; Poletto et al. 1997]. One system, Dynamo,¹⁰ proposed to perform staged compilation and dynamic optimization for Scheme and Java, as well as for ML [Leone and Dybvig 1997].

Templates aside, ML may be dynamically compiled anyway. In Cardelli’s description of his ML compiler, he noted:

[Compilation] is repeated for every definition or expression typed by the user... or fetched from an external file. Because of the interactive use of the compiler, the compilation of small phrases must be virtually instantaneous. [Cardelli 1984, p. 209]

2.10. Erlang

Erlang is a functional language, designed for use in large, soft real-time systems such as telecommunications equipment [Armstrong 1997]. Johansson et al. [2000] described the implementation of a JIT compiler for Erlang, HiPE, designed to address performance problems.

As a recently designed system without historical baggage, HiPE stands out in that the user must explicitly invoke the JIT compiler. The rationale for this is that it gives the user a fine degree of control over the performance/code space tradeoff that mixed code offers [Johansson et al. 2000].

HiPE exercises considerable care when performing “mode-switches” back and

¹⁰ A name collision: Leone and Dybvig’s “Dynamo” is different from the “Dynamo” of Bala et al. [1999].

forth between native and interpreted code. Mode-switches may be needed at the obvious locations—calls and returns—as well as for thrown exceptions. Their calls use the mode of the *caller* rather than the mode of the called code; this is in contrast to techniques used for mixed code in Lisp (Gabriel and Masinter [1985] discussed mixed code calls in Lisp and their performance implications).

2.11. Specialization and O’Caml

O’Caml is another functional language, and can be considered a dialect of ML [Rémy et al. 1999]. The O’Caml interpreter has been the focus of run-time specialization work.

Piumarta and Riccardi [1998] specialized the interpreter’s instructions to the program being run, in a limited way.¹¹ They first dynamically translated interpreted bytecodes into direct threaded code [Bell 1973], then dynamically combined blocks of instructions together into new “macro opcodes,” modifying the code to use the new instructions. This reduced the overhead of instruction dispatch, and yielded opportunities for optimization in macro opcodes which would not have been possible if the instructions had been separate (although they did not perform such optimizations). As presented, their technique did not take dynamic execution paths into account, and they noted that it is best suited to low-level instruction sets, where dispatch time is a relatively large factor in performance.

A more general approach to run-time specialization was taken by Thibault et al. [2000]. They applied their program specializer, Tempo [Consel et al. 1998], to the Java virtual machine and the O’Caml interpreter at run-time. They noted:

While the speedup obtained by specialization is significant, it does not compete with results obtained with hand-written off-line or run-time compilers. [Thibault et al. 2000, p. 170]

¹¹ Thibault et al. [2000] provided an alternative view on Piumarta and Riccardi’s work with respect to specialization.

But later in the paper they stated that

... program specialization is entering relative maturity. [Thibault et al. 2000, p. 175]

This may be taken to imply that, at least for the time being, program specialization may not be as fruitful as other approaches to dynamic compilation and optimization.

2.12. Prolog

Prolog systems dynamically compile, too, although the execution model of Prolog necessitates use of specialized techniques. Van Roy [1994] gave an outstanding, detailed survey of the area. One of SICStus Prolog’s native code compilers, which could be invoked and have its output loaded dynamically, was described in Haygood [1994].

2.13. Simulation, Binary Translation, and Machine Code

Simulation is the process of running native executable machine code for one architecture on another architecture.¹² How does this relate to JIT compilation? One of the techniques for simulation is binary translation; in particular, we focus on dynamic binary translation that involves translating from one machine code to another at run-time. Typically, binary translators are highly specialized with respect to source and target; research on retargetable and “resourceable” binary translators is still in its infancy [Ung and Cifuentes 2000]. Altman et al. [2000b] have a good discussion of the challenges involved in binary translation, and Cmelik and Keppel [1994] compared pre-1995 simulation systems in detail. Rather than duplicating their work, we will take a higher-level view.

May [1987] proposed that simulators could be categorized by their implementation technique into three generations. To

¹² We use the term *simulate* in preference to *emulate* as the latter has the connotation that hardware is heavily involved in the process. However, some literature uses the words interchangeably.

this, we add a fourth generation to characterize more recent work.

- (1) First-generation simulators were interpreters, which would simply interpret each source instruction as needed. As might be expected, these tended to exhibit poor performance due to interpretation overhead.
- (2) Second-generation simulators dynamically translated source instructions into target instruction one at a time, caching the translations for later use.
- (3) Third-generation simulators improved upon the performance of second-generation simulators by dynamically translating entire blocks of source instructions at a time. This introduces new questions as to what should be translated. Most such systems translated either basic blocks of code or extended basic blocks [Cmelik and Keppel 1994], reflecting the static control flow of the source program. Other static translation units are possible: one anomalous system, DAISY, performed page-at-a-time translations from PowerPC to VLIW instructions [Ebcioglu and Altman 1996, 1997].
- (4) What we call fourth-generation simulators expand upon the third-generation by dynamically translating paths, or traces. A path reflects the control flow exhibited by the source program at run-time, a dynamic instead of a static unit of translation. The most recent work on binary translation is concentrated on this type of system.

Fourth-generation simulators are predominant in recent literature [Bala et al. 1999; Chen et al. 2000; Deaver et al. 1999; Gschwind et al. 2000; Klaiber 2000; Zheng and Thompson 2000]. The structure of these is fairly similar:

- (1) *Profiled execution.* The simulator's effort should be concentrated on "hot" areas of code that are frequently executed. For example, initialization code that is executed only once should not be translated or optimized. To deter-

mine which execution paths are hot, the source program is executed in some manner and profile information is gathered. Time invested in doing this is assumed to be recouped eventually.

When source and target architectures are dissimilar, or the source architecture is uncomplicated (such as a reduced instruction set computer (RISC) processor) then interpretation of the source program is typically employed to execute the source program [Bala et al. 1999; Gschwind et al. 2000; Transmeta Corporation 2001; Zheng and Thompson 2000]. The alternative approach, direct execution, is best summed up by Rosenblum et al. [1995, p. 36]:

By far the fastest simulator of the CPU, MMU, and memory system of an SGI multiprocessor is an SGI multiprocessor.

In other words, when the source and target architectures are the same, as in the case where the goal is dynamic optimization of a source program, the source program can be executed directly by the central processing unit (CPU). The simulator regains control periodically as a result of appropriately modifying the source program [Chen et al. 2000] or by less direct means such as interrupts [Gorton 2001].

- (2) *Hot path detection.* In lieu of hardware support, hot paths may be detected by keeping counters to record frequency of execution [Zheng and Thompson 2000], or by watching for code that is structurally likely to be hot, like the target of a backward branch [Bala et al. 1999]. With hardware support, the program's program counter can be sampled at intervals to detect hot spots [Deaver et al. 1999].

Some other considerations are that paths may be strategically *excluded* if they are too expensive or difficult to translate [Zheng and Thompson 2000], and choosing good stopping points for paths can be as important as choosing good starting points in terms

of keeping a manageable number of traces [Gschwind et al. 2000].

- (3) *Code generation and optimization.* Once a hot path has been noted, the simulator will translate it into code for the target architecture, or perhaps optimize the code. The correctness of the translation is always at issue, and some empirical verification techniques are discussed in [Zheng and Thompson 2000].
- (4) *“Bail-out” mechanism.* In the case of dynamic optimization systems (where the source and target architectures are the same), there is the potential for a negative impact on the source program’s performance. A bail-out mechanism [Bala et al. 1999] heuristically tries to detect such a problem and revert back to the source program’s direct execution; this can be spotted, for example, by monitoring the stability of the working set of paths. Such a mechanism can also be used to avoid handling complicated cases.

Another recurring theme in recent binary translation work is the issue of hardware support for binary translation, especially for translating code for legacy architectures into VLIW code. This has attracted interest because VLIW architectures promise legacy architecture implementations which have higher performance, greater instruction-level parallelism [Ebcioglu and Altman 1996, 1997], higher clock rates [Altman et al. 2000a; Gschwind et al. 2000], and lower power requirements [Klaiber 2000]. Binary translation work in these processors is still done by software at run-time, and is thus still dynamic binary translation, although occasionally packaged under more fanciful names to enrapture venture capitalists [Geppert and Perry 2000]. The key idea in these systems is that, for efficiency, the target VLIW should provide a superset of the source architecture [Ebcioglu and Altman 1997]; these extra resources, unseen by the source program, can be used by the binary translator for aggressive optimizations or to simulate troublesome aspects of the source architecture.

2.14. Java

Java is implemented by static compilation to bytecode instructions for the Java virtual machine, or JVM. Early JVMs were only interpreters, resulting in less-than-stellar performance:

Interpreting bytecodes is slow. [Cramer et al. 1997, p. 37]

Java isn’t just slow, it’s *really* slow, *surprisingly* slow. [Tyma 1998, p. 41]

Regardless of how vitriolic the expression, the message was that Java programs had to run faster, and the primary means looked to for accomplishing this was JIT compilation of Java bytecodes. Indeed, Java brought the term *just-in-time* into common use in computing literature.¹³ Unquestionably, the pressure for fast Java implementations spurred a renaissance in JIT research; at no other time in history has such concentrated time and money been invested in it.

An early view of Java JIT compilation was given by Cramer et al. [1997], who were engineers at Sun Microsystems, the progenitor of Java. They made the observation that there is an upper bound on the speedup achievable by JIT compilation, noting that interpretation proper only accounted for 68% of execution time in a profile they ran. They also advocated the direct use of JVM bytecodes, a stack-based instruction set, as an intermediate representation for JIT compilation and optimization. In retrospect, this is a minority viewpoint; most later work, including Sun’s own [Sun Microsystems 2001], invariably began by converting JVM code into a register-based intermediate representation.

The interesting trend in Java JIT work [Adl-Tabatabai et al. 1998; Bik et al. 1999; Burke et al. 1999; Cierniak and Li 1997; Ishizaki et al. 1999; Krall and Graf 1997; Krall 1998; Yang et al. 1999] is the implicit assumption that mere

¹³ Gosling [2001] pointed out that the term *just-in-time* was borrowed from manufacturing terminology, and traced his own use of the term back to about 1993.

translation from bytecode to native code is not enough: code optimization is necessary too. At the same time, this work recognizes that traditional optimization techniques are expensive, and looks for modifications to optimization algorithms that strike a balance between speed of algorithm execution and speed of the resulting code.

There have also been approaches to Java JIT compilation besides the usual interpret-first-optimize-later. A compile-only strategy, with no interpreter whatsoever, was adopted by Burke et al. [1999], who also implemented their system in Java; improvements to their JIT directly benefited their system. Agesen [1997] translated JVM bytecodes into Self code, to leverage optimizations already existing in the Self compiler. Annotations were tried by Azevedo et al. [1999] to shift the effort of code optimization prior to runtime: information needed for efficient JIT optimization was precomputed and tagged on to bytecode as annotations, which were then used by the JIT system to assist its work. Finally, Plezbert and Cytron [1997] proposed and evaluated the idea of “continuous compilation” for Java in which an interpreter and compiler would execute concurrently, preferably on separate processors.¹⁴

3. CLASSIFICATION OF JIT SYSTEMS

In the course of surveying JIT work, some common attributes emerged. We propose that JIT systems can be classified according to three properties:

- (1) *Invocation*. A JIT compiler is explicitly invoked if the user must take some action to cause compilation at runtime. An implicitly invoked JIT compiler is transparent to the user.
- (2) *Executability*. JIT systems typically involve two languages: a source language to translate from, and a target language to translate to (although

these languages can be the same, if the JIT system is only performing optimization on-the-fly). We call a JIT system *monoexecutable* if it can only execute one of these languages, and *polyexecutable* if it can execute more than one. Polyexecutable JIT systems have the luxury of deciding when compiler invocation is warranted, since either program representation can be used.

- (3) *Concurrency*. This property characterizes how the JIT compiler executes, relative to the program itself. If program execution pauses under its own volition to permit compilation, it is not concurrent; the JIT compiler in this case may be invoked via subroutine call, message transmission, or transfer of control to a coroutine. In contrast, a concurrent JIT compiler can operate as the program executes concurrently: in a separate thread or process, even on a different processor.

JIT systems that function in hard real time may constitute a fourth classifying property, but there seems to be little research in the area at present; it is unclear if hard real-time constraints pose any unique problems to JIT systems.

Some trends are apparent. For instance, implicitly invoked JIT compilers are definitely predominant in recent work. Executability varies from system to system, but this is more an issue of design than an issue of JIT technology. Work on concurrent JIT compilers is currently only beginning, and will likely increase in importance as processor technology evolves.

4. TOOLS FOR JIT COMPILATION

General, portable tools for JIT compilation that help with the dynamic generation of binary code did not appear until relatively recently. To varying degrees, these toolkits address three issues:

- (1) *Binary code generation*. As argued in Ramsey and Fernández [1995], emitting binary code such as machine language is a situation rife with opportunities for error. There are associated

¹⁴ As opposed to the ongoing optimization of Kistler’s [2001] “continuous optimization,” only compilation occurred concurrently using “continuous compilation,” and only happened once.

Table 1. Comparison of JIT Toolkits

Source	Binary code generation	Cache coherence	Execution	Abstract interface	Input
Engler [1996]	•	•	•	•	ad hoc
Engler and Proebsting [1994]	•	•	•	•	tree
Fraser and Proebsting [1999]	•	•	•	•	postfix
Keppel [1991]		•	•	•	n/a
Ramsey and Fernández [1995]	•				ad hoc

Note: n/a = not applicable.

bookkeeping tasks too: information may not yet be available upon initial code generation, like the location of forward branch targets. Once discovered, the information must be backpatched into the appropriate locations.

- (2) *Cache coherence.* CPU speed advances have far outstripped memory speed advances in recent years [Hennessy and Patterson 1996]. To compensate, modern CPUs incorporate a small, fast cache memory, the contents of which may get temporarily out of sync with main memory. When dynamically generating code, care must be taken to ensure that the cache contents reflect code written to main memory before execution is attempted. The situation is even more complicated when several CPUs share a single memory. Keppel [1991] gave a detailed discussion.
- (3) *Execution.* The hardware or operating system may impose restrictions which limit where executable code may reside. For example, memory earmarked for data may not allow execution (i.e., instruction fetches) by default, meaning that code could be generated into the data memory, but not executed without platform-specific wrangling. Again, refer to Keppel [1991].

Only the first issue is relevant for JIT compilation to interpreted virtual machine code—interpreters don’t directly execute the code they interpret—but there is no reason why JIT compilation tools cannot be useful for generation of nonnative code as well.

Table I gives a comparison of the toolkits. In addition to indicating how well the toolkits support the three areas above, we have added two extra categories. First, an *abstract interface* is one that is architecture-independent. Use of a toolkit’s abstract interface implies that very little, if any, of the user’s code needs modification in order to use a new platform. The drawbacks are that architecture-dependent operations like register allocation may be difficult, and the mapping from abstract to actual machine may be suboptimal, such as a mapping from RISC abstraction to complex instruction set computer (CISC) machinery.

Second, *input* refers to the structure, if any, of the input expected by the toolkit. With respect to JIT compilation, more complicated input structures take more time and space for the user to produce and the toolkit to consume [Engler 1996].

Using a tool may solve some problems but introduce others. Tools for binary code generation help avoid many errors compared to manually emitting binary code. These tools, however, require detailed knowledge of binary instruction formats whose specification may itself be prone to error. Engler and Hsieh [2000] presented a “metatool” that can automatically derive these instruction encodings by repeatedly querying the existing system assembler with varying inputs.

5. CONCLUSION

Dynamic, or just-in-time, compilation is an old implementation technique with a fragmented history. By collecting this historical information together, we hope to shorten the voyage of rediscovery.

ACKNOWLEDGMENTS

Thanks to Nigel Horspool, Shannon Jaeger, and Mike Zastre, who proofread and commented on drafts of this paper. Comments from the anonymous referees helped improve the presentation as well. Also, thanks to Rick Gorton, James Gosling, Thomas Kistler, Ralph Mauriello, and Jim Mitchell for supplying historical information and clarifications. Evelyn Duesterwald's PLDI 2000 tutorial notes were helpful in preparing Section 2.9.

REFERENCES

- ABRAMS, P. S. 1970. An APL machine. Ph.D. dissertation. Stanford University, Stanford, CA. Also, Stanford Linear Accelerator Center (SLAC) Rep. 114.
- ADL-TABATABAI, A.-R., CIERNIAK, M., LUEH, G.-Y., PARIKH, V. M., AND STICHNOY, J. M. 1998. Fast, effective code generation in a just-in-time Java compiler. In *PLDI '98*. 280–290.
- AGESEN, O. 1996. Concrete type inference: Delivering object-oriented applications. Ph.D. dissertation. Stanford University, Stanford, CA. Also Tech. Rep. SMLI TR-96-52, Sun Microsystems, Santa Clara, CA (Jan. 1996).
- AGESEN, O. 1997. Design and implementation of Pep, a Java just-in-time translator. *Theor. Prac. Obj. Syst.* 3, 2, 127–155.
- AGESEN, O. AND HÖLZLE, U. 1995. Type feedback vs. concrete type inference: A comparison of optimization techniques for object-oriented languages. In *Proceedings of OOPSLA '95*. 91–107.
- ALTMAN, E., GSCHWIND, M., SATHAYE, S., KOSONOCKY, S., BRIGHT, A., FRITTS, J., LEDAK, P., APPENZELLER, D., AGRICOLA, C., AND FILAN, Z. 2000a. BOA: The architecture of a binary translation processor. Tech. Rep. RC 21665, IBM Research Division, Yorktown Heights, NY.
- ALTMAN, E. R., KAELI, D., AND SHEFFER, Y. 2000b. Welcome to the opportunities of binary translation. *IEEE Comput.* 33, 3 (March), 40–45.
- ARMSTRONG, J. 1997. The development of Erlang. In *Proceedings of ICFP '97* (1997). 196–203.
- AUSLANDER, J., PHILIPPOSE, M., CHAMBERS, C., EGGERS, S. J., AND BERSHAD, B. N. 1996. Fast, effective dynamic compilation. In *Proceedings of PLDI '96*. 149–159.
- AZEVEDO, A., NICOLAU, A., AND HUMMEL, J. 1999. Java annotation-aware just-in-time (AJIT) compilation system. In *Proceedings of JAVA '99*. 142–151.
- BALA, V., DUESTERWALD, E., AND BANERJIA, S. 1999. Transparent dynamic optimization. Tech. Rep. HPL-1999-77, Hewlett-Packard, Palo Alto, CA.
- BARTLETT, J. 1992. *Familiar Quotations* (16th ed.). J. Kaplan, Ed. Little, Brown and Company, Boston, MA.
- BELL, J. R. 1973. Threaded code. *Commun. ACM* 16, 6 (June), 370–372.
- BENTLEY, J. 1988. Little languages. In *More Programming Pearls*. Addison-Wesley, Reading, MA, 83–100.
- BIK, A. J. C., GIRKAR, M., AND HAGHIGHAT, M. R. 1999. Experiences with Java JIT optimization. In *Innovative Architecture for Future Generation High-Performance Processors and Systems*. IEEE Computer Society Press, Los Alamitos, CA, 87–94.
- BROWN, P. J. 1976. Throw-away compiling. *Softw.—Pract. Exp.* 6, 423–434.
- BROWN, P. J. 1990. *Writing Interactive Compilers and Interpreters*. Wiley, New York, NY.
- BURGER, R. G. 1997. Efficient compilation and profile-driven dynamic recompilation in scheme. Ph.D. dissertation, Indiana University, Bloomington, IN.
- BURKE, M. G., CHOI, J.-D., FINK, S., GROVE, D., HIND, M., SARKAR, V., SERRANO, M. J., SREEDHAR, V. C., AND SRINIVASAN, H. 1999. The Jalapeño dynamic optimizing compiler for Java. In *Proceedings of JAVA '99*. 129–141.
- CARDELLI, L. 1984. Compiling a functional language. In *1984 Symposium on Lisp and Functional Programming*. 208–217.
- CHAMBERS, C. 1992. The design and implementation of the self compiler, an optimizing compiler for object-oriented programming languages. Ph.D. dissertation. Stanford University, Stanford, CA.
- CHAMBERS, C. AND UNGAR, D. 1989. Customization: optimizing compiler technology for Self, a dynamically-typed object-oriented programming language. In *Proceedings of PLDI '89*. 146–160.
- CHAMBERS, C. AND UNGAR, D. 1990. Iterative type analysis and extended message splitting: Optimizing dynamically-typed object-oriented programs. In *Proceedings of PLDI '90*. 150–164.
- CHAMBERS, C. AND UNGAR, D. 1991. Making pure object-oriented languages practical. In *Proceedings of OOPSLA '91*. 1–15.
- CHAMBERS, C., UNGAR, D., AND LEE, E. 1989. An efficient implementation of Self, a dynamically-typed object-oriented programming language based on prototypes. In *Proceedings of OOPSLA '89*. 49–70.
- CHEN, W.-K., LERNER, S., CHAIKEN, R., AND GILLIES, D. M. 2000. Mojo: a dynamic optimization system. In *Proceedings of the Third ACM Workshop on Feedback-Directed and Dynamic Optimization* (FDDO-3, Dec. 2000).
- CIERNIAK, M. AND LI, W. 1997. Briki: an optimizing Java compiler. In *Proceedings of IEEE COMPCON '97*. 179–184.
- CMELIK, B. AND KEPPEL, D. 1994. Shade: A fast instruction-set simulator for execution profiling. In *Proceedings of the 1994 Conference on Measurement and Modeling of Computer Systems*. 128–137.

- CONSEL, C., HORNOF, L., MARLET, R., MULLER, G., THIBAUT, S., VOLANSCHI, E.-N., LAWALL, J., AND NOYÉ, J. 1998. Tempo: Specializing systems applications and beyond. *ACM Comput. Surv.* 30, 3 (Sept.), 5pp.
- CONSEL, C. AND NOËL, F. 1996. A general approach for run-time specialization and its application to C. In *Proceedings of POPL '96*. 145–156.
- CRAMER, T., FRIEDMAN, R., MILLER, T., SEBERGER, D., WILSON, R., AND WOLCZKO, M. 1997. Compiling Java just in time. *IEEE Micro* 17, 3 (May/June), 36–43.
- DAKIN, R. J. AND POOLE, P. C. 1973. A mixed code approach. *The Comput. J.* 16, 3, 219–222.
- DAWSON, J. L. 1973. Combining interpretive code with machine code. *The Comput. J.* 16, 3, 216–219.
- DEAVER, D., GORTON, R., AND RUBIN, N. 1999. Wiggins/Redstone: An on-line program specializer. In *Proceedings of the IEEE Hot Chips XI Conference* (Aug. 1999). IEEE Computer Society Press, Los Alamitos, CA.
- DEUTSCH, L. P. AND SCHIFFMAN, A. M. 1984. Efficient implementation of the Smalltalk-80 system. In *Proceedings of POPL '84*. 297–302.
- DIECKMANN, S. AND HÖLZLE, U. 1997. The space overhead of customization. Tech. Rep. TRCS 97-21. University of California, Santa Barbara, Santa Barbara, CA.
- EBCIOĞLU, K. AND ALTMAN, E. R. 1996. DAISY: Dynamic compilation for 100% architectural compatibility. Tech. Rep. RC 20538. IBM Research Division, Yorktown Heights, NY.
- EBCIOĞLU, K. AND ALTMAN, E. R. 1997. Daisy: Dynamic compilation for 100% architectural compatibility. In *Proceedings of ISCA '97*. 26–37.
- ENGLER, D. R. 1996. VCODE: a retargetable, extensible, very fast dynamic code generation system. In *Proceedings of PLDI '96*. 160–170.
- ENGLER, D. R. AND HSIEH, W. C. 2000. DERIVE: A tool that automatically reverse-engineers instruction encodings. In *Proceedings of the ACM SIGPLAN Workshop on Dynamic and Adaptive Compilation and Optimization* (Dynamo '00). 12–22.
- ENGLER, D. R., HSIEH, W. C., AND KAASHOEK, M. F. 1996. C: A language for high-level, efficient, and machine-independent dynamic code generation. In *Proceedings of POPL '96*. 131–144.
- ENGLER, D. R. AND PROEBSTING, T. A. 1994. DCG: An efficient, retargetable dynamic code generation system. In *Proceedings of ASPLOS VI*. 263–272.
- FRANZ, M. 1994. *Code-generation on-the-fly: A key to portable software*. Ph.D. dissertation. ETH Zurich, Zurich, Switzerland.
- FRANZ, M. AND KISTLER, T. 1997. Slim binaries. *Commun. ACM* 40, 12 (Dec.), 87–94.
- FRASER, C. W. AND PROEBSTING, T. A. 1999. Finite-state code generation. In *Proceedings of PLDI '99*. 270–280.
- GABRIEL, R. P. AND MASINTER, L. M. 1985. *Performance and Evaluation of Lisp Systems*. MIT Press, Cambridge, MA.
- GEPPERT, L. AND PERRY, T. S. 2000. Transmeta's magic show. *IEEE Spectr.* 37, 5 (May), 26–33.
- GOLDBERG, A. AND ROBSON, D. 1985. *Smalltalk-80: The Language and its Implementation*. Addison-Wesley, Reading, MA.
- GORTON, R. 2001. Private communication.
- GOSLING, J. 2001. Private communication.
- GSCHWIND, M., ALTMAN, E. R., SATHAYE, S., LEDAK, P., AND APPENZELLER, D. 2000. Dynamic and transparent binary translation. *IEEE Comput.* 33, 3, 54–59.
- HAMMOND, J. 1977. BASIC—an evaluation of processing methods and a study of some programs. *Softw.—Pract. Exp.* 7, 697–711.
- HANSEN, G. J. 1974. Adaptive systems for the dynamic run-time optimization of programs. Ph.D. dissertation. Carnegie-Mellon University, Pittsburgh, PA.
- HAYGOOD, R. C. 1994. Native code compilation in SICStus Prolog. In *Proceedings of the Eleventh International Conference on Logic Programming*. 190–204.
- HENNESSY, J. L. AND PATTERSON, D. A. 1996. *Computer Architecture: A Quantitative Approach*, 2nd ed. Morgan Kaufmann, San Francisco, CA.
- HÖLZLE, U. 1994. *Adaptive optimization for Self: Reconciling high performance with exploratory programming*. Ph.D. dissertation. Carnegie-Mellon University, Pittsburgh, PA.
- HÖLZLE, U. AND UNGAR, D. 1994a. Optimizing dynamically-dispatched calls with run-time type feedback. In *Proceedings of PLDI '94*. 326–336.
- HÖLZLE, U. AND UNGAR, D. 1994b. A third-generation Self implementation: Reconciling responsiveness with performance. In *Proceedings of OOPSLA '94*. 229–243.
- ISHIZAKI, K., KAWAHITO, M., YASUE, T., TAKEUCHI, M., OGASAWARA, T., SUGANUMA, T., ONODERA, T., KOMATSU, H., AND NAKATANI, T. 1999. Design, implementation, and evaluation of optimizations in a just-in-time compiler. In *Proceedings of JAVA '99*. 119–128.
- JOHANSSON, E., PETERSSON, M., AND SAGONAS, K. 2000. A high performance Erlang system. In *Proceedings of PPDP '00*. 32–43.
- JOHNSTON, R. L. 1977. The dynamic incremental compiler of APL\3000. In *APL '79 Conference Proceedings*. Published in *APL Quote Quad* 9, 4 (June), Pt. 1, 82–87.
- JONES, N. D., GOMARD, C. K., AND SESTOFT, P. 1993. *Partial Evaluation and Automatic Program Generation*. Prentice Hall, Englewood Cliffs, NJ.
- KEPPEL, D. 1991. A portable interface for on-the-fly instruction space modification. In *Proceedings of ASPLOS IV*. 86–95.
- KEPPEL, D., EGGERS, S. J., AND HENRY, R. R. 1991. A case for runtime code generation. Tech. Rep.

- 91-11-04. Department of Computer Science and Engineering, University of Washington, Seattle, WA.
- KISTLER, T. 1997. Dynamic runtime optimization. In *Proceedings of the Joint Modular Languages Conference (JMLC '97)*. 53–66.
- KISTLER, T. 1999. *Continuous program optimization*. Ph.D. dissertation. University of California, Irvine, Irvine, CA.
- KISTLER, T. 2001. Private communication.
- KISTLER, T. AND FRANZ, M. 1999. The case for dynamic optimization: Improving memory-hierarchy performance by continuously adapting the internal storage layout of heap objects at run-time. Tech. Rep. 99-21 (May). University of California, Irvine, Irvine, CA. Revised September, 1999.
- KLAIBER, A. 2000. The technology behind Crusoe processors. Tech. Rep. (Jan.), Transmeta Corporation, Santa Clara, CA.
- KNUTH, D. E. 1971. An empirical study of Fortran programs. *Softw.—Pract. Exp.* 1, 105–133.
- KRALL, A. 1998. Efficient JavaVM just-in-time compilation. In *Proceedings of the 1998 International Conference on Parallel Architectures and Compilation Techniques (PACT '98)*. 205–212.
- KRALL, A. AND GRAFL, R. 1997. A Java just-in-time compiler that transcends JavaVM's 32 bit barrier. In *Proceedings of PPOPP '97 Workshop on Java for Science and Engineering*.
- LEE, P. AND LEONE, M. 1996. Optimizing ML with run-time code generation. In *Proceedings of PLDI '96*. 137–148.
- LEE, S., YANG, B.-S., KIM, S., PARK, S., MOON, S.-M., EBICIOĞLU, K., AND ALTMAN, E. 2000. Efficient Java exception handling in just-in-time compilation. In *Proceedings of Java 2000*. 1–8.
- LEONE, M. AND DYBVIK, R. K. 1997. DynamO: A staged compiler architecture for dynamic program optimization. Tech. Rep. 490. Computer Science Department, Indiana University, Bloomington, IN.
- LEONE, M. AND LEE, P. 1994. Lightweight run-time code generation. In *Proceedings of the ACM SIGPLAN Workshop on Partial Evaluation and Semantics-Based Program Manipulation*. 97–106.
- MARLET, R., CONSEL, C., AND BOINOT, P. 1999. Efficient incremental run-time specialization for free. In *PLDI '99*. 281–292.
- MAURIELLO, R. 2000. Private communication.
- MAY, C. 1987. Mimic: A fast System/370 simulator. In *Proceedings of the SIGPLAN '87 Symposium on Interpreters and Interpretive Techniques* (June). ACM Press, New York, NY, 1–13.
- MCCARTHY, J. 1960. Recursive functions of symbolic expressions and their computation by machine, part I. *Commun. ACM* 3, 4, 184–195.
- MCCARTHY, J. 1981. History of LISP. In *History of Programming Languages*, R. L. Wexelblat, Ed. Academic Press, New York, NY, 173–185.
- MILLER, T. C. 1977. Tentative compilation: A design for an APL compiler. In *APL '79 Conference Proceedings*. Volume 9 Published in *APL Quote Quad* 9, 4 (June), Pt. 1, 88–95.
- MITCHELL, J. G. 1970. The design and construction of flexible and efficient interactive programming systems. Ph.D. dissertation. Carnegie-Mellon University, Pittsburgh, PA.
- MITCHELL, J. G. 2000. Private communication.
- MITCHELL, J. G., PERLIS, A. J., AND VAN ZOEREN, H. R. 1968. LC²: A language for conversational computing. In *Interactive Systems for Experimental Applied Mathematics*, M. Klerer and J. Reinfelds, Eds. Academic Press, New York, NY. (Proceedings of 1967 ACM Symposium.)
- MOCK, M., BERRYMAN, M., CHAMBERS, C., AND EGGERS, S. J. 1999. Calpa: A tool for automating dynamic compilation. In *Proceedings of the Second ACM Workshop on Feedback-Directed and Dynamic Optimization*. 100–109.
- NG, T. S. AND CANTONI, A. 1976. Run time interaction with FORTRAN using mixed code. *The Comput. J.* 19, 1, 91–92.
- PITTMAN, T. 1987. Two-level hybrid interpreter/native code execution for combined space-time program efficiency. In *Proceedings of the SIGPLAN Symposium on Interpreters and Interpretive Techniques*. ACM Press, New York, NY, 150–152.
- PIUMARTA, I. AND RICCARDI, F. 1998. Optimizing direct threaded code by selective inlining. In *Proceedings of PLDI '98*. 291–300.
- PLEZBERT, M. P. AND CYTRON, R. K. 1997. Does “just in time” = “better late than never”? In *Proceedings of POPL '97*. 120–131.
- POLETTI, M., ENGLER, D. R., AND KAASHOEK, M. F. 1997. tcc: A system for fast, flexible, and high-level dynamic code generation. In *Proceedings of PLDI '97*. 109–121.
- RAMSEY, N. AND FERNÁNDEZ, M. 1995. The New Jersey machine-code toolkit. In *Proceedings of the 1995 USENIX Technical Conference*. 289–302.
- RAU, B. R. 1978. Levels of representation of programs and the architecture of universal host machines. In *Proceedings of the 11th Annual Microprogramming Workshop (MICRO-11)*. 67–79.
- RÉMY, D., LEROY, X., AND WEIS, P. 1999. Objective Caml—a general purpose high-level programming language. *ERCIM News* 36, 29–30.
- ROSENBLUM, M., HERROD, S. A., WITCHEL, E., AND GUPTA, A. 1995. Complete computer system simulation: The SimOS approach. *IEEE Parall. Distrib. Tech.* 3, 4 (Winter), 34–43.
- SCHROEDER, S. C. AND VAUGHN, L. E. 1973. A high order language optimal execution processor: Fast Intent Recognition System (FIRST). In *Proceedings of a Symposium on High-Level-Language*

- Computer Architecture*. Published in *SIGPLAN* 8, 11 (Nov.), 109–116.
- SEBESTA, R. W. 1999. *Concepts of Programming Languages* (4th ed.). Addison-Wesley, Reading, MA.
- SMITH, R. B. AND UNGAR, D. 1995. Programming as an experience: The inspiration for Self. In *Proceedings of ECOOP '95*.
- SUN MICROSYSTEMS. 2001. The Java HotSpot virtual machine. White paper. Sun Microsystems, Santa Clara, CA.
- THIBAUT, S., CONSEL, C., LAWALL, J. L., MARLET, R., AND MULLER, G. 2000. Static and dynamic program compilation by interpreter specialization. *Higher-Order Symbol. Computat.* 13, 161–178.
- THOMPSON, K. 1968. Regular expression search algorithm. *Commun. ACM* 11, 6 (June), 419–422.
- TRANSMETA CORPORATION. 2001. Code morphing software. Available online at http://www.transmeta.com/technology/architecture/code_morphing.html. Transmeta Corporation, Santa Clara, CA.
- TYMA, P. 1998. Why are we using Java again? *Commun. ACM* 41, 6, 38–42.
- UNG, D. AND CIFUENTES, C. 2000. Machine-adaptable dynamic binary translation. In *Proceedings of Dynamo '00*. 41–51.
- UNGAR, D. AND SMITH, R. B. 1987. Self: The power of simplicity. In *Proceedings of OOPSLA '87*. 227–242.
- UNGAR, D., SMITH, R. B., CHAMBERS, C., AND HÖLZLE, U. 1992. Object, message, and performance: How they coexist in Self. *IEEE Comput.* 25, 10 (Oct.), 53–64.
- UNIVERSITY OF MICHIGAN. 1966a. The System Loader. In *University of Michigan Executive System for the IBM 7090 Computer*, Vol. 1. University of Michigan, Ann Arbor, MI.
- UNIVERSITY OF MICHIGAN. 1966b. The “University of Michigan Assembly Program” (“UMAP”). In *University of Michigan Executive System for the IBM 7090 Computer*, Vol. 2. University of Michigan, Ann Arbor, MI.
- VAN DYKE, E. J. 1977. A dynamic incremental compiler for an interpretive language. *Hewlett-Packard J.* 28, 11 (July), 17–24.
- VAN ROY, P. 1994. The wonder years of sequential Prolog implementation. *J. Logic Program.* 19–20, 385–441.
- WICKLINE, P., LEE, P., AND PFENNING, F. 1998. Runtime code generation and Modal-ML. In *Proceedings of PLDI '98*. 224–235.
- WIRTH, N. AND GUTKNECHT, J. 1989. The Oberon system. *Softw.—Pract. Exp.* 19, 9 (Sep.), 857–893.
- YANG, B.-S., MOON, S.-M., PARK, S., LEE, J., LEE, S., PARK, J., CHUNG, Y. C., KIM, S., EBCIOĞLU, K., AND ALTMAN, E. 1999. LaTTe: A Java VM just-in-time compiler with fast and efficient register allocation. In *Proceedings of the International Conference on Parallel Architectures and Compilation Techniques*. 128–138. IEEE Computer Society Press, Los Alamitos, CA.
- ZHENG, C. AND THOMPSON, C. 2000. PA-RISC to IA-64: Transparent execution, no recompilation. *IEEE Comput.* 33, 3 (March), 47–52.

Received July 2002; revised March 2003; accepted February 2003