



No Assembly Required: Compiling Standard ML to C

DAVID TARDITI, PETER LEE,
and
ANURAG ACHARYA
Carnegie Mellon University

C has been used as a portable target language for implementing languages like Standard ML and Scheme. Previous efforts at compiling these languages to C have produced efficient code, but have compromised on portability and proper tail recursion. We show how to compile Standard ML to C without making such compromises. The compilation technique is based on converting Standard ML to a continuation-passing style λ -calculus intermediate language and then compiling this language to C. The code generated by this compiler achieves an execution speed that is about a factor of two slower than that generated by a native code compiler. The compiler generates highly portable code, yet still supports advanced features like garbage collection and first-class continuations. We analyze the performance and determine the aspects of the compilation method that lead to the observed slowdown. We also suggest changes to C compilers that would better support such compilation methods.

Categories and Subject Descriptors: D.3.2 [**Programming Languages**]: Language Classifications; D.3.4 [**Programming Languages**]: Processors—*compilers*

General Terms: Languages

Additional Key Words and Phrases: Compilation to C, continuation-passing style, Scheme, Standard ML

1. INTRODUCTION

With the profusion of new computer architectures and programming languages, compiler writers are increasingly faced with making difficult compromises between the efficiency of compiled programs and retargetability of the compiler. One approach that has been used in the past is to compile programs into the C programming language [15], in essence, using C as a universal intermediate language. C makes for an efficient intermediate language

This research was partially supported by the National Science Foundation under PYI grant CCR-9057567. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the National Science Foundation or the U.S. Government.

Authors' address: School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

© 1992 ACM 1057-4514/92/0600-0161 \$01.50

ACM Letters on Programming Languages and Systems, Vol. 1, No. 2, June 1992, Pages 161-177.

because it allows relatively low-level access to the machine, yet in a number of important respects is mostly machine independent. Furthermore, it is usually safe to assume that a C compiler will be available on almost any general-purpose machine.

Bartlett has shown that it is possible to compile Scheme programs into efficient C programs [8]. In his approach, constructs in Scheme are mapped directly to similar constructs in C. Although this is a conceptually simple approach, it also leads to certain compromises. First, Bartlett's Scheme implementation, Scheme \rightarrow C, fails to reflect the pragmatics of the language. Specifically, "proper" tail recursion [21] is lost, despite a compile-time analysis of procedures for tail-recursive behavior. Second, features such as first-class continuations and garbage collection are implemented with the use of assembly code. The main point of compiling to C, though, is to avoid such machine dependencies.

We are led, then, to consider several questions: How can languages such as Scheme and Standard ML (SML), which support first-class procedures and other powerful control constructs such as exceptions and first-class continuations, be compiled without any use of assembly language? Can languages that depend on proper tail recursion be faithfully compiled to C? Can this be done without an unacceptable loss of efficiency, and if not, what extensions to C would be necessary to make this possible?

In this paper we address these questions by describing our experience with compiling SML [19] to C. A compiler is presented that compiles the entire SML language into efficient, portable C code. In doing so, the use of assembly language is completely avoided, and only two assumptions about the architecture are made: that both pointers and integers are represented in 32 bits and that integer arithmetic is two's-complement.

Our compiler, *sml2c*, is based entirely on the Standard ML of New Jersey (SML/NJ) compiler [6]. We have replaced the code generator and slightly modified the run-time system.¹ All other modules of the SML/NJ compiler have been used "as is." As a result, our compiler is source-code compatible with the SML/NJ compiler and handles recently proposed extensions to SML such as first-class continuations (*call/cc*) [10], asynchronous signal handling [22], lightweight concurrent processes [9], and multiprocessing [20]. Our compiler also handles separate compilation. The implementation has been tested on a number of nontrivial programs, including bootstrapping of the compiler itself, and the generated C code has been tested, without any modification, on Sun-3s, Sparcstations, Decstations, IBM RTs, a Motorola 88000-based machine, a Sequent Symmetry (where it was used for multiprocessing ML programs) [14], and an 80486-based PC.

Our primary considerations in the design of the compiler were, in decreasing order of importance, ability to handle the full language (with extensions), faithfulness to the pragmatics of the language, and efficiency of the compiled

¹ The present run-time system is written in C and requires an operating system compatible with UNIX 4.3 BSD.

code. The first two goals were met completely. The price we pay in run-time overhead is typically about 70 to 100 percent in running time, as compared to the native machine code generated by the SML/NJ compiler. In this paper our basic approach is described, and then some limitations and the optimizations we used to improve significantly the performance of the generated C code are discussed. This is followed by a performance analysis and finally some conclusions. The system is available via anonymous ftp from **dravido.soar.cs.cmu.edu: /usr / nemo / sml2c**.

2. RELATED WORK

Besides Scheme \rightarrow C, several other compilers use C as their target language. Examples include **f2c** [11], which compiles Fortran-77 to C; **mtc** [18], which compiles Modula-II to C; **p2c** [12], which compiles Pascal to C; **CParaOPS5** [1], which compiles OPS5 to C; and Portable Cedar [7], which compiles Cedar to C.

The effort most closely related to sml2c is Portable Cedar. Cedar, like SML, has exceptions and garbage collection. Compiling Cedar to C, however, is easier than compiling SML to C, since Cedar does not have proper tail-recursion or first-class functions. Furthermore, the Portable Cedar compiler compromises machine-independence by generating code that traverses the C call stack to implement exceptions; this is dependent on the machine and C compiler being used. The sml2c compiler generates fully machine-independent code that is portable across a large class of machines.

3. THE DESIGN OF THE COMPILER

SML is a lexically scoped, call-by-value language with higher-order functions. SML has polymorphic types that are automatically inferred and checked by the compiler. Support for developing large programs is provided by a sophisticated modules system that provides for static type-checking of the interfaces between modules, as in Ada and Modula-II. A type-safe, dynamically scoped exception mechanism allows programs to handle unusual or deviant conditions. Garbage collection is provided to automate the management of heap storage.

Clearly, there is considerable semantic distance between SML and C. Features such as dynamically scoped exceptions, garbage collection, and higher-order functions pose significant problems for an efficient and portable implementation of SML in C. Recently proposed extensions to SML for first-class continuations and concurrent processes further increase this semantic distance.

3.1 Intermediate Representation

The key element in the design of sml2c was the decision to replace the code generator of the SML/NJ compiler with one that produced C code. The SML/NJ compiler is a publicly available, freely redistributable, and modifiable compiler for SML developed at AT & T and Princeton University. Input

to the SML/NJ code generator is represented in a continuation-passing, closure passing style λ -calculus [4, 5].

There were several reasons for this decision. First, by focusing on the translation of continuation-passing, closure-passing style programs to C, techniques developed would be applicable to many languages, in particular to Scheme. The second reason was the remarkable similarity between programs represented in a continuation-passing, closure-passing style and C programs. It seemed likely that these programs could be translated to highly portable C code, although it was unclear whether the C code could be made reasonably efficient. Another reason was pragmatic. The compiler had to be built by one person working over a summer. This time constraint made it necessary to modify an existing compiler.

The continuation-passing, closure-passing style λ -calculus used by the code generator is a refinement of a continuation-passing style (CPS) λ -calculus. A program is transformed into CPS by adding a continuation argument to all user functions. The continuation argument is a function that represents the remaining execution of the program. When a function has finished computing its result, instead of returning, the function calls its continuation with the result as the argument. Programs that have been transformed into the particular CPS used by the SML/NJ compiler always satisfy the syntactic invariants that function applications are never nested and are always tail recursive. Since function calls are always tail recursive, the first call to return is essentially the only call to return. Thus, a function call can be considered as a **goto** with arguments. Conversion of Scheme or ML programs into CPS can be performed by a simple $O(n)$ transformation [4].

There are several advantages to using a CPS intermediate representation [16, 25]. First, all intermediate values are explicitly named. Second, control-flow is explicitly represented by continuations. This makes it easy to implement exceptions and first-class continuations. Third, since all function calls turn into jumps with arguments, tail-recursion elimination is achieved automatically. In addition to the SML/NJ compiler, CPS has been used successfully in the RABBIT [25] and ORBIT [17] compilers.

A continuation-passing, closure-passing style λ -calculus is a refinement of CPS in which environment operations are made explicit. All functions are flattened out to one lexical level, and explicit record operations are used to represent operations on closures. A description of the conversion of CPS programs into continuation-passing, closure-passing style can be found in [4].

The target machine for a continuation-passing, closure-passing style program is conceptually simple, consisting of a heap, a set of general-purpose registers, and a set of reserved registers for the heap pointer, heap limit pointer, current exception handler, and arithmetic temporary storage. The instructions required are those typically found on most conventional machines, including loads, stores, logical operations, arithmetic operations with and without overflow checking, floating-point operations, register-to-register moves, conditional branches, and jumps.

The internal representation of the CPS used by the SML/NJ compiler is shown in Figure 1. Continuation expressions, or values of type **cexp**, repre-

```

datatype cexp =
  RECORD of (value * accesspath) list * lvar * cexp
| SELECT of int * value * lvar * cexp
| OFFSET of int * value * lvar * cexp
| APP of value * value list
| FIX of (lvar * lvar list * cexp) list * cexp
| SWITCH of value * cexp list
| PRIMOP of primop * value list * lvar list * cexp list

datatype value =
  VAR of lvar
| LABEL of lvar
| INT of int
| REAL of string
| STRING of string

```

Fig. 1. Internal representation of CPS used by the Standard ML of New Jersey compiler.

sent CPS programs. Most variants of the **cexp** type bind zero or more variables in the scope of another continuation expression. Variable names are represented by values of type **lvar**. **RECORD**(a, b, c) creates a tuple of values from the values in a , binds the tuple to b , and continues execution in c . **SELECT**(a, b, c, d) selects the a th field of b , binds it to c , and continues execution in d . **OFFSET**(a, b, c, d) takes the tuple bound to b , binds a new tuple starting at the a th field of b to c , and continues execution in d . **APP**(a, b) applies a to the arguments given in b . **FIX**(a, b) binds a list of potentially mutually recursive functions in a and continues execution into b . Each element of a is a tuple ($f, args, body$), where f is the name of the function to be bound, and the function is $\lambda args. body$. **SWITCH**(a, b) is a case statement where execution continues with the a th element of b . **PRIMOP**(a, b, c, d) applies a primitive operator a to a tuple of values b , binds some variables c , and continues execution into one of the continuation expressions given in d . Primitive operators include integer arithmetic operators, floating-point arithmetic, integer relational operators, floating-point relational operators, and operators on arrays.

The intermediate representations for a simple SML program are given in Figures 2 to 5. These examples are similar to the ones in [4, p. 104–105]. By convention, the intermediate representations are presented in a simplified pseudocode style. In particular, infix operators have not been expanded to their true CPS form. The notation $c.x$ denotes the tuple field selection; it selects the x th field from the tuple bound to c . Figure 2 shows the SML source program. Figure 3 shows its CPS intermediate representation. Continuation arguments, named by variables of the form k_n , where n is an integer, have been added to each user function. Figure 4 shows the closure-passing, continuation-passing representation. Closure arguments, named by variables of the form c_n have been added to each function. A closure argument is a record containing a pointer to the code for the function and the values of

Fig. 2. Source program.

```

fun add x =
  let fun add' y = x+y
    in add'
  end

```

```

fun add (x,k1) =
  let fun add' (y,k2) = k2 (x+y)
    in k1 add'
  end

```

Fig. 3. CPS intermediate representation.

Fig. 4. Closure-converted, continuation-passing style intermediate representation.

```

fun add (c1,x,k1) =
  let fun add' (c2,y,k2) =
    let val x' = c2.2
      val k2' = k2.1
    in k2' (k2,x'+y)
    end
    val add'' = (add',x)
    val k1' = k1.1
  in k1' (k1,add')
  end

```

```

fun add(c1,x,k1) =
  let val add'' = (add',x)
    val k1' = k1.1
  in k1' (k1,add'')
  end
and add' (c2,y,k2) =
  let val x' = c2.2
    val k2' = k2.1
  in k2' (k2,x'+y)
  end

```

Fig. 5. After flattening.

(formerly) free variables of the function. The closure for the function **add'** is created by the body of the **add** function and bound to the variable **add''**. After conversion to closure-passing, continuation-passing style, functions have no free variables other than other function names. All functions can, therefore, be flattened out to one lexical level. Figure 5 shows the closure-passing, continuation-passing representation after this has been done. Note that the flattened closure-converted, CPS code is remarkably similar to C code, except for the fact that all function calls are tail recursive.

Another example is given in Figures 6 to 9. In the closure-converted, CPS representation for this program, note that **bar** and **baz** do not have closure arguments. This is because they are used only in the function position of applications. Closures do not have to be built, and free variables (of which there are none) can instead be added as extra arguments. Again, note the similarity of the flattened, closure-converted, CPS code to C code.

```

fun foo x =
  let fun bar 0 = "bar"
      | bar y = baz(y-1)
      and baz 0 = "baz"
      | baz z = bar(z-1)
  in bar x
  end

```

Fig. 6. Source program.

```

fun foo (x,k1) =
  let fun bar (0,k2) = k2 "bar"
      | bar (y,k2) = baz(y-1,k2)
      and baz (0,k3) = k3 "baz"
      | baz (z,k3) = bar(z-1,k3)
  in bar(x,k1)
  end

```

Fig. 7. CPS intermediate representation.

```

fun foo (c,x,k1) =
  let fun bar (0,k2) =
      let val k2' = k2.1 in k2'(k2,"bar") end
      | bar (y,k2) = baz(y-1,k2)
      and baz (0,k3) =
      let val k3' = k3.1 in k3'(k3,"baz") end
      | baz (z,k3) = bar(z-1,k3)
  in bar(x,k1)
  end

```

Fig. 8. Closure-converted, continuation-passing style intermediate representation.

```

fun foo (c,x,k1) = bar(x,k1)
and bar (0,k2) =
  let val k2' = k2.1 in k2'(k2,"bar") end
  | bar (y,k2) = baz(y-1,k2)
and baz (0,k3) =
  let val k3' = k3.1 in k3'(k3,"baz") end
  | baz (z,k3) = bar(z-1,k3)

```

Fig. 9. After flattening.

3.2 The Compilation Model

One may note the lack of a stack in the target machine. In fact, the issue of whether to allow a stack in the target machine is a fundamental design decision. In principle, the presence of higher-order functions means that activation records for procedures must be allocated on the heap since the environment structure will have the structure of a tree. In practice, heap allocation of activation records is common in compilers that use a CPS representation, for example, the RABBIT and SML/NJ compilers.

With a suitable amount of care, a stack can be used to store most activation records [17]. This leads to the possibility of flattening out all functions in an SML program to a single lexical level and then mapping SML functions directly to corresponding C functions, thereby making use of the C call stack. Closures would then be represented by records containing a C function pointer and the environment for the function. (Bartlett [8] uses a similar approach in Scheme \rightarrow C.) Unfortunately, this makes it difficult to implement garbage collection in a portable manner. This is due to the fact that finding the roots for garbage collection requires a nonportable traversal of the C stack. It also makes it difficult to preserve proper tail recursion. Some C compilers implement some tail-recursion optimization for C functions that are immediately tail recursive (i.e., they call themselves). They generally fail on mutually recursive tail-calls and tail-calls to functions that are not known until run time, which are both quite frequent in a programming style using higher-order functions. Finally, the use of the stack makes it impossible to obtain a portable implementation of the call/cc extension to SML (though efficient stack-based approaches have been proposed [13]). Not even the *longjmp* and *setjmp* features provided by some C implementations can provide this functionality.

For these reasons, we adhere to the compilation model used by the SML/NJ compiler, which uses heap allocation of activation records and makes no use of a stack.

4. CODE GENERATION

It is straightforward to implement the target machine resources in C. The registers and heap are implemented by integer arrays. (Use of actual machine registers is discussed in the next section.) Most of the target machine instructions are also straightforward. Only the jump and arithmetic operations (with overflow checking) complicate matters. Recall that function calls are tail recursive after CPS conversion and are, thus, compiled into jumps with arguments passed in registers. The implementation of a jump is problematic because labels are not first-class values in C. In particular, we cannot store a label in memory or jump to it from an arbitrary point in a C program. The only way to obtain the address of a block of C code is to encapsulate it in a C function. But, if C function calls are used naively, the stack would quickly overflow.

Instead, we borrow a mechanism from the RABBIT compiler, called the *UUO handler*. We refer to it as the *dispatch loop*. The dispatch loop emulates the apply operation. Code execution begins with the dispatch loop calling the first function to be executed. When this function wants to call another function, it returns to the dispatch loop with the address of the next function to call. In general, when a function f needs to call another function g , it returns the address of g to the dispatch loop, which then calls g . Arguments to functions are passed in memory in the simulated machine registers. In this way, the depth of the stack of activation records never grows to more than two.

The implementation of integer arithmetic is complicated by the fact that overflow checking is not normally provided in C. The definition of SML, however, mandates overflow checking. To implement overflow checking, bit-level checks on the operands and the result are used.

Register allocation is based on the spilling transformation used in the SML/NJ compiler [4]. This guarantees that, at every point in the program, a sufficient number of registers are available for all variables. Our code generator then performs register assignment on the fly via register tracking. A variety of heuristics are used in order to minimize shuffling of registers at function calls.

Figure 10 contains a simplified version of the C code that would be generated for the intermediate code shown in Figure 9. The figure also shows the dispatch loop stripped of its initialization code.

5. RUN-TIME INTERACTION

The implementation of the CPS target machine as a C program is supported by a modified version of the SML/NJ run-time system [3], which provides operating-system services, garbage collection, and asynchronous signal handling. The generated code and the SML/NJ run-time system work together as coroutines. For the purposes of multiprocessing, the per-processor state is threaded throughout the run-time system and the generated code. As far as the generated code is concerned, a pointer to the per-processor state is maintained in a dedicated element of the register array.

Execution begins in the run-time system, which then transfers control to the generated code. The generated code must transfer control back to the run-time system for successful program completion, requests for operating-system services, and machine faults such as floating-point overflow, garbage collection, and asynchronous signal handling.

The code transfers control to the run-time system by doing a *longjmp* back into the function containing the dispatch loop, and returning. The target of the *longjmp* is set by doing a *setjmp* at the top of the function containing the dispatch loop. Of course, the jump buffer is part of the per-processor state. An advantage of this approach is that it allows the heart of the dispatch loop procedure to be extremely simple. In fact, the main loop of the procedure in actual use is merely an unrolled version of the code shown in Figure 10.

Garbage collection is initiated by having the generated code call the garbage collector as a subroutine when the heap is exhausted. Like the SML/NJ compiler, our compiler avoids heap checks on every allocation by inserting a single heap check at the beginning of every function. The heap check is accomplished by adding the heap pointer and an offset that is the maximum amount of allocation that could be done along any path through the function, and checking that this value is less than the heap limit pointer. If this is not true, the garbage collector is called. Code generation ensures that at such points all roots are in the register array.

The modified run-time system contains no assembly language code. The original run-time system uses assembly language code for certain language

```

int foo(),bar(),baz();

int dispatch(start,r)
int (*start)();
int *r;
{ while (1)
    start = (int (*)()) (*start)(r);
}

int foo (r)
int *r;
{ return((int) bar); }

int bar (r)
int *r;
{ if (r[1]==0)
    { r[1] = "bar"; return r[2]; }
  else
    { r[1] = r[1] - 1; return ((int) baz); }
}

int baz (r)
int *r;
{ if (r[1]==0)
    { r[1] = "baz"; return r[2]); }
  else
    { r[1] = r[1] - 1; return ((int) bar); }
}

```

Fig. 10. C code and dispatch loop (simplified).

primitives and to implement the transfer from the C run-time system to the generated ML code. The language primitives have been recoded in portable C, and the dispatch loop has replaced the transfer code.

6. OPTIMIZATIONS

There are a number of problems with the implementation described in Section 4. The generated C code fails to make effective use of actual machine registers, function calls are expensive, and arithmetic involves a high overhead for portable overflow checking. The following subsections describe optimizations that address these problems.

6.1 Register Caching

Recall that the target machine registers are implemented by an integer array. Global variables could be used, but most C compilers will not move global variables into real machine registers [23]. Furthermore, this approach is not viable for a shared-memory multiprocessor implementation. To make effective use of real registers, we cache target machine registers in local C

variables for the duration of function calls. We then depend on the C compiler to place the local variables in registers when possible.

The register caching optimization uses a simple static count of the number of uses of a register along each possible path through a function body to decide whether to cache a target machine register for the duration of the function call. A consequence of the continuation-passing, closure-passing style is that function bodies are extended basic blocks. In other words, they have only forward branches and no loops. If the number of uses along any path is greater than two, we cache the register. Care must be taken when caching certain registers that must be valid whenever a machine fault might occur. These registers are the heap pointer register and the exception continuation register. We never cache the exception continuation register, and the heap pointer register is spilled back to the corresponding global variable before executing any instruction that may cause a machine fault. Fortunately, the only instructions that may cause machine faults in our implementation are division by zero and floating-point operations, so the need to spill is fairly rare.

6.2 Optimizing Function Calls

Functions calls are expensive since every function call involves a return to the dispatch loop and then a call from the dispatch loop to the target function. This involves at least two jumps. The expense is usually more than this, since many C compilers implement the return by a jump to a piece of “epilogue” code at the end of a function that cleans things up and then does the actual jump back to the dispatch loop. In addition, values stored in actual machine registers might be saved to memory upon entry to a function and then restored upon exit. This occurs, in practice, on callee-save systems even though the dispatch loop has just one variable live across function calls. Thus, the true cost of a function call is usually three jumps and some number of loads and store.

In order to eliminate some of this overhead, we optimize calls to known functions. A *known* function is one whose possible call sites are all known. If all call paths to a known function f pass through a single function g , we can move f into the body of g and turn all calls to f into **gotos**. We also perform direct tail-recursion elimination, turning all calls to g within its body to **gotos**. Known function-call optimization is useful for avoiding passes through the dispatch loop. In particular, it is useful for optimizing tight tail-recursive loops.

Computing whether all call paths to a known function pass through another function can be cast as the classical problem of computing dominators in a call graph with multiple start nodes, where each unknown function is a start node. An algorithm for computing dominator information can be found in [2]. Define a maximal dominator to be a function that is only dominated by itself. This captures the intuitive concept of a highest-level dominator. First, the set of dominators is computed for each known function. For each known function f , it is an easily proved fact that there is a unique

maximal dominator for f . If the maximal dominator of f is not itself, then f is moved into its maximal dominator.

Consider the intermediate code of Figure 9. The functions **bar** and **baz** are known, and **foo** is the maximal dominator for them. Thus, **bar** and **baz** can be integrated into **foo**. Figure 9 shows the pseudo-CPS code that represents these functions as they are presented to the code generator. Figure 11 shows a simplified version of C code generated for these functions with function integration and register caching. Note that the mutually tail-recursive functions have been compiled into a tight loop with the loop counter placed in a register. (The redundant **gotos** appearing in the code are optimized away by most C compilers.)

6.3 Overflow Checking

Finally, integer arithmetic is expensive because most C compilers provide no portable way to use the integer overflow signal. The implementation of addition and subtraction with portable overflow checking requires complicated explicit bit-wise tests, multiplication requires a function call, and division requires some comparisons and branches. The overflow checks can be simplified by constant-folding them in the cases where one operand is known at compile time. For example, the overflow check for addition can be reduced from a complicated Boolean expression involving five operators to a single comparison and branch, and the function call for multiplication can be replaced with two compares and a branch.

7. BENCHMARKS

To measure the performance of sml2c, we compared sml2c with and without the optimizations turned on, with the SML/NJ compiler (version 0.75) on two different architectures.

The benchmark suite is a subset of the suite used in [4] for benchmarking the SML/NJ compiler. It consists of the following programs: *Life*, a program for the Game of Life; *Lex*, a lexical analyzer generator; *Yacc*, a LALR(1) parser generator; *Knuth-B*, an implementation of the Knuth-Bendix completion algorithm; and *VLIW*, an instruction scheduler that performs software pipelining on register-transfer machines. These programs range in size from 117 to 5800 lines of SML code. *Life* was run on 50 generations of a glider gun. *Lex* was run on a lexical description of SML. *Yacc* was run on a grammar for SML.

The benchmarks were run on a Decstation 5000/200 with 48 Mbytes of memory running Mach 2.5, and a Sun 4/330 with 24 Mbytes running Mach 2.5. We used built-in timing functions available in the SML/NJ system that are based on the UNIXTM *getrusage* utility. We factored out garbage-collection costs, so our numbers represent only code speed. Garbage-collection costs are irrelevant for the sake of comparison, since all versions of code generated

TM UNIX is a registered trademark of AT & T Bell Laboratories.

```

int foo (r)
int *r;
{ register int r1;
  r1=r[1];
  goto bar;
bar:
  if (r1==0) { r1 = "bar"; r[1]=r1; return r[2]; }
  else { r1 = r1 - 1; goto baz; }
baz:
  if (r1==0) { r1 = "baz"; r[1]=r1; return r[2]; }
  else { r1 = r1 - 1; goto bar; }
}

```

Fig. 11. Optimized C code.

by `sml2c` and the SML/NJ compiler generate the same amount of garbage and use the same garbage collector. Ignoring garbage-collection times does not significantly affect the basic performance ratios, since we found garbage-collection times to be between 1–20 percent of the overall execution times for the native code programs. For both `sml2c` and the native code compiler, we used the same high optimization settings provided by the SML/NJ compiler. The C code was compiled on both machines using version 1.40 of `gcc` [24] with the `-O` and the `-fomit-frame-pointer` flags.

Tables I and II show the execution times for the various benchmark programs. The column labeled “UnOpt” gives the times for the code generated by `sml2c` for the benchmark programs with the optimizations turned off, whereas the column labeled “Opt” gives the times with the optimizations turned on. The code generated by `sml2c` for the benchmark programs with the optimizations turned on generally runs at about half the speed of the native code generated by the SML/NJ compiler.

We also measured compilation times for our benchmarks to show that the times were still reasonable. The compilation times include the time spent by `sml2c` and the time spent by `gcc` compiling the generated C code. The compilations were done on a Decstation 5000/200 using `sml2c` with optimization and version 0.75 of the SML/NJ compiler, both run with the settings described previously. The times are shown in Table III. The compilation times for the long programs, *Lex*, *Yacc*, and *VLIW*, are less than a factor of two slower than the SML/NJ times.

Several aspects of our compilation method lead to the observed slowdown in the benchmark programs. One possible source of overhead comes from the loads and stores to the register array. Another possible source are operations that are more expensive to implement in C than in assembly language, namely, the jumps (which are implemented by transfers through the dispatch loop) and the heap check (which in the native code typically requires a single instruction, but in C requires a conditional test and branch). There is also added overhead for integer arithmetic with overflow checking, but a previous benchmark study showed this to have a negligible effect [26].

Table I. Benchmarks for Decstation 500/200

Program ()	SML/NJ (α) (s)	Opt (β) (s)	β/α ()	Unopt (δ) (s)	δ/α ()
<i>Life</i>	25.3	50.2	1.98	90.0	3.56
<i>Lex</i>	17.3	30.4	1.76	57.3	3.31
<i>Yacc</i>	5.4	10.8	2.00	19.2	3.56
<i>Knu-B</i>	14.4	32.1	2.23	57.5	3.99
<i>VLIW</i>	30.0	77.7	2.60	142.4	4.75

Table II. Benchmarks for Sun 4/330

Program ()	SML/NJ (α) (s)	Opt (β) (s)	β/α ()	Unopt (δ) (s)	δ/α ()
<i>Life</i>	56.3	105.9	1.88	161.0	2.85
<i>Lex</i>	39.3	64.0	1.62	109.3	2.78
<i>Yacc</i>	14.8	24.2	1.63	35.2	2.37
<i>Knu-B</i>	43.4	78.2	1.80	103.0	2.37
<i>VLIW</i>	80.4	133.2	1.65	207.6	2.58

Table III. Compilation Times on Decstation 5000/200

Program ()	SML/NJ (α) (s)	Opt (β) (s)	β/α ()
<i>Life</i>	11.2	28.2	2.5
<i>Lex</i>	74.2	134.6	1.8
<i>Yacc</i>	395.8	602.5	1.5
<i>Knu-B</i>	26.3	62.1	2.4
<i>VLIW</i>	367.0	646.8	1.8

In order to better understand the nature of these sources of overhead, we instrumented the generated C code. The instrumentation counts the following: s , the number of C statements executed, excluding assignment statements that move values between the register array and local variables; w , the number of assignments to the register array; r , the number of reads from the register array; c , the number of function calls; and ifc , the number of calls to integrated functions. Calls to integrated functions are included in c . The results of instrumenting the benchmark programs are given in Table IV.

It is clear from the data that function calls occur with enormous frequency and, furthermore, that most function calls transfer through the dispatch loop. Table V gives the average number of statements executed in a function (including the last jump), and the percentage of function calls that are calls to integrated functions, that is, turned into **gotos**. It also seems clear that there is a tremendous amount of reading and writing of values in the register array, relative to the total number of statements executed. Table VI gives the number of reads and writes as a fraction of the statement count. The reading

Table IV. Results from Instrumented Programs

Program ()	<i>s</i> (‘000s)	<i>w</i> (‘000s)	<i>r</i> (‘000s)	<i>c</i> (‘000s)	<i>ifc</i> (‘000s)
<i>Life</i>	256,421	71,502	109,710	37,763	10,348
<i>Lex</i>	172,358	31,255	54,376	15,070	5,976
<i>Yacc</i>	55,948	17,685	10,875	4,248	832
<i>Knu-B</i>	158,533	43,410	59,928	12,630	284
<i>VLIW</i>	170,849	52,781	72,624	18,267	2,993

Table V. Frequency of Function Calls and Integrated Function Calls

Program	<i>s/c</i>	<i>ifc/c</i> (%)
<i>Life</i>	6.8	27
<i>Lex</i>	11.4	39
<i>Yacc</i>	13.2	20
<i>Knu-B</i>	12.6	2.2
<i>VLIW</i>	9.4	16.4

Table VI. Overhead of Register Array Usage

Program	$(r + w)/s$
<i>Life</i>	0.70
<i>Lex</i>	0.50
<i>Yacc</i>	0.51
<i>Knu-B</i>	0.65
<i>VLIW</i>	0.73

and writing are to be expected with small function bodies and the passing of parameters in memory via the register array.

This overhead could be reduced greatly by extending C to support global register variables and by modifying C compilers to support proper tail recursion. The introduction of global register variables would allow one to eliminate the register array. Thus, the generated C code would no longer make reads and writes to the register array. Extending C compilers to support proper tail recursion or some form of interprocedural **goto** would allow one to eliminate the use of the dispatch loop in the generated code. Instead, the generated code could use C function calls as its function-call mechanism.

8. CONCLUSIONS

Our experience shows that it is possible to compile languages such as Scheme and SML without using any assembly language. It is possible to develop portable implementations of such languages without sacrificing either proper tail recursion or an inordinate amount of efficiency.

ACKNOWLEDGMENTS

We would like to thank Joel Bartlett and Henry Baker for their comments on an earlier version of this paper, and Mark Leone for comments on the present paper. Andrew Appel provided valuable feedback when we first began this work. We would also like to thank the Venari group at Carnegie-Mellon University for giving us computing time on their workstations.

REFERENCES

1. ACHARYA, A., AND KALP, D. Release notes for CParaOPS5 5.3 and ParaOPS5 4.4. Available with the CParaOPS5 release from School of Computer Science, Carnegie-Mellon Univ. Pittsburgh, Pa., May 1990.
2. AHO, A. V., SETHI, R., AND ULLMAN, J. D. *Compilers: Principles, Techniques and Tools*. Addison-Wesley, Reading, Mass., 1986.
3. APPEL, A. W. A runtime system. *Lisp Symbolic Comput.* 3, 4 (Nov. 1990), 343–380.
4. APPEL, A. W. *Compiling with Continuations*. Cambridge University Press, New York, 1992.
5. APPEL, A. W., AND JIM, T. Y. Continuation-passing, closure passing style. In *Proceedings of the 16th Annual ACM Symposium on Principles of Programming Languages* (Austin, Tex., Jan. 11–13), ACM, New York, 1989, pp. 293–302.
6. APPEL, A. W., AND MACQUEEN, D. B. A standard ML compiler. In *Functional Programming Languages and Computer Architecture*, vol. 274, G. Kahn, Ed. Springer-Verlag, New York, 1987, pp. 301–324.
7. ATKINSON, R., DEMERS, A., HAUSER, C., JACOBI, C., KESSLER, P., AND WEISER, M. Experience creating a Portable Cedar. In *Proceedings of the SIGPLAN 89 Conference on Programming Language Design and Implementation* (Portland, Oreg., June 21–23). ACM, New York, 1989, pp. 322–329.
8. BARTLETT, J. F. SCHEME → C: A portable Scheme-to-C compiler. Tech. Rep., DEC Western Research Laboratory, Palo Alto, Calif., Jan. 1989.
9. COOPER, E. C., AND MORRISETT, J. G. Adding threads to Standard ML. Tech. Rep. CMU-CS-90-186, School of Computer Science, Carnegie-Mellon Univ., Pittsburgh, Pa., Dec. 1990.
10. DUBA, B. F., HARPER, R., AND MACQUEEN, D. Typing first-class continuations in ML. In *Proceedings of the 18th Annual ACM Symposium on Principles of Programming Languages* (Orlando, Fla., Jan. 21–23), ACM, New York, 1991, pp. 163–173.
11. FELDMAN, S., GAY, D., MAIMONE, M. W., AND SCHYER, N. A Fortran-to-C converter. Comput. Sci. Tech. Rep. 149, AT & T Bell Laboratories, Murray Hill, N.J., May 1990.
12. GILLESPIE, D. The p2c translator. Available by anonymous ftp from csvax.cs.caltech.edu under the GNUcplyleft, 1989.
13. HIEB, R., DYBVIK, R. K., AND BRUGGEMAN, C. Representing control in the presence of first-class continuations. In *Proceedings of the SIGPLAN 90 Conference on Programming Language Design and Implementation* (White Plains, N.Y., June 20–22). ACM, New York, 1990, pp. 66–77.
14. HUELSBERGEN, L., AND LARUS, J. Dynamic program parallelization. In *Proceedings of the 1992 ACM Conference on LISP and Functional Programming* (San Francisco, Calif., June 22–24). ACM, New York, 1992, pp. 311–323.
15. KERNIGHAN, B. W., AND RITCHIE, D. M. *The C Programming Language*. 2nd Ed. Prentice-Hall, Englewood Cliffs, N.J., 1988.
16. KRANZ, D. ORBIT: An optimizing compiler for Scheme. Ph.D. thesis, Yale Univ., Dept. of Computer Science, New Haven, Conn., Feb. 1988.
17. KRANZ, D., KELSEY, R., REES, J., HUDAK, P., PHILBIN, J., AND ADAMS, N. ORBIT: An optimizing compiler for Scheme. In *Proceedings of the SIGPLAN 86 Conference on Programming Language Design and Implementation* (Palo Alto, Calif., June 25–27). ACM, New York, 1986, pp. 219–233.
18. MARTIN, M. Entwurf und Implementierung eines bersetzers von Modula-2 nach C. Master's thesis, Fakultät für Informatik, Univ. of Karlsruhe, Germany, Feb. 1990.

19. MILNER, R., TOFTE, M., AND HARPER, R. *The Definition of Standard ML*. MIT Press, Cambridge, Mass., 1990.
20. MORRISETT, J. G., AND TOLMACH, A. A portable multiprocessor interface for Standard ML of New Jersey. Tech. Rep. CMU-CS-92-155, School of Computer Science, Carnegie-Mellon Univ., Pittsburgh, Pa., June 1992.
21. REES, J., AND CLINGER, W. Revised report on the algorithmic language Scheme. *SIGPLAN Not. (ACM)* 21, 12 (Dec. 1986), 37–79.
22. REPPY, J. H. Asynchronous signals in Standard ML. Tech. Rep. 90-1144, Dept. of Computer Science, Cornell Univ., Ithaca, N.Y., Aug. 1990.
23. SANTHANAM, V., AND ODNERT, D. Register allocation across procedure and module boundaries. In *Proceedings of the SIGPLAN 90 Conference on Programming Language Design and Implementation* (White Plains, N.Y., June 20–22), ACM, New York, 1990, pp. 28–39.
24. STALLMAN, R. M. Using and porting GNU CC. GNU CC is a widely available C compiler developed by the Free Software Foundation, Cambridge, Mass., Sept. 1989.
25. STEEL, G. L. JR. RABBIT: A compiler for Scheme (a study in compiler optimization). Master's thesis, Tech. Rep. AI-TR-474, AI Lab., MIT, Cambridge, Mass. May 1978.
26. TARDITI, D., ACHARYA, A., AND LEE, P. No assembly required: Compiling Standard ML to C. Tech. Rep. CMU-CS-90-187, School of Computer Science, Carnegie-Mellon Univ., Pittsburgh, Pa., Nov. 1990.

Received May 1992; revised July 1992; accepted August 1992