# A Verified Compiler for VLISP PreScheme

DINO P. OLIVA[*]                                             (*oliva@cse.ogi.edu*)

*College of Computer Science*
*161 Cullinane Hall*
*Northeastern University*
*Boston, MA 02115*

JOHN D. RAMSDELL[†]                                          (*ramsdell@mitre.org*)

*The MITRE Corporation*
*202 Burlington Road*
*Bedford, MA 01730-1420*

MITCHELL WAND[‡]                                             (*wand@ccs.neu.edu*)

*College of Computer Science*
*161 Cullinane Hall*
*Northeastern University*
*Boston, MA 02115*

**Keywords:** verified, compiler

## Contents

**Abstract.** This paper describes a verified compiler for PreScheme, the implementation language for the VLISP run-time system. The compiler and proof were divided into three parts: A transformational front end that translates source text into a core language, a syntax-directed compiler that translates the core language into combinator-based tree-manipulation language, and a linearizer that translates combinator code into code for an abstract stored-program machine with linear memory for both data and code. This factorization enabled different proof techniques to be used for the different phases of the compiler, and also allowed the generation of good code. Finally, the whole process was made possible by carefully defining the semantics of VLISP PreScheme rather than just adopting Scheme's. We believe that the architecture of the compiler and its correctness proof can easily be applied to compilers for languages other than PreScheme.

## 1. Introduction

As part of the VLISP project, we have developed an architecture for the specification and implementation of verified optimizing compilers. We have used this architecture to develop a verified compiler for PreScheme.

This architecture divides the compiler into three components. In the case of PreScheme, these components are:

1. A transformational front end that translates source text into a core language called Pure PreScheme.

2. A syntax-directed compiler that translates Pure PreScheme into a

combinator-based tree-manipulation language. We call this language *combinator code*.

3. An assembler that translates combinator code into code for an abstract stored-program machine with linear memory for both data and code.

We believe that this is a good architecture for the specification and implementation of verified optimizing compilers. A smaller example of this proof architecture (without the front end) is given in [28]; another example, with a more elaborate assembler, is given in [17].

The correctness of each transformation of the program is justified relative to an appropriate semantics:

- The source language and its core subset (PreScheme and Pure PreScheme in our example) share a denotational semantics.

- The language of combinators has both a denotational semantics, expressed using the same domains as the source language, and an operational semantics. The latter may be derived from the former by general theorems of the $\lambda$-calculus.

- The stored-program machine is specified using an operational semantics.

This separation of semantics allows each transformation to be proved correct using a suitable proof technique:

1. The transformational front end is shown to preserve the denotational semantics of the source language, by induction on the number of transformation steps.

2. The syntax-directed compiler is shown correct by denotational reasoning, using structural induction on the phrases of the core language. The proof establishes a relation between the denotational semantics of the core language and the denotational semantics of the combinator language. The proof is made feasible by the fact that the core and combinator semantics are expressed using the same domains, and hence the same $\lambda$-theory.

3. The assembler is proved correct by operational reasoning, using the operational semantics of the combinator language and the operational semantics of the stored-program machine. We define a relation, called

| Level | Language | Acronym | Section |
|---|---|---|---|
| Source Language | VLISP PreScheme | VLPS | 2 |
| Parsed Source Language | Macro-Free PreScheme | MFPS | 2 |
| Core Language | Pure PreScheme | PPS | 2 |
| Tree Structured Language | Combinator Code | CC | 5 |
| Linear Data Language | Linear-Data Code | LDC | 6 |
| Linear Program Language | Stored-Program Code | SPC | 7 |
| Target Language | Assembly Code | AC | 9 |

Figure 1: Levels in the PreScheme compilation process

> a *storage-layout relation*, that determines when the combinator machine and the stored-program machine are in corresponding states, and we show that operation of the machines preserves this relation.

The verification of the PreScheme compiler is at the same level as that of the VLISP system itself: the algorithms are verified with respect to a formal denotational semantics of PreScheme and a formal model of an abstract target machine. The implementation of the algorithms was carried out in Scheme, and code for the abstract target machine was translated into assembly language for two real target machines (Motorola 68000 and SPARC). While neither of these implementation steps was verified, both were straightforward; for a discussion of the trustworthiness of this approach, see [9].

In the case of the PreScheme compiler, the assembler proof is performed in two steps: we first relate the combinator machine to a machine, called the *linear-data machine*, in which the data store is linear, but the program is still a tree. We then relate the linear-data machine to the *stored-program machine*, in which both the data and code are kept in linear memories.

This gives a total of five language or machine levels and four major proofs. In general, we will present level $n$, followed by the proof or proofs relating level $n$ to level $n - 1$. Figure 1 tabulates the various levels.

The rest of the paper is organized as follows:

## 1.1. Introduction

Following [3], we will give a synopsis of each section in the paper. Section $n$ will be summarized in section 1.$n$.

## 1.2.   VLISP PreScheme

PreScheme is a dialect of Scheme intended for systems programming. PreScheme was carefully designed so that it syntactically looks like Scheme and has similar semantics. With a little care, PreScheme programs can be run and debugged as if they were ordinary Scheme programs.

However, PreScheme is particularly suitable for systems programming because every valid PreScheme program can be executed using only a C-like run-time system. A compiler for this language will reject any program that requires dynamic type checking or creation of closures at run-time. Furthermore, it is intended that PreScheme programs be runnable without the use of automatic storage reclamation.

There are a number of dialects of PreScheme that played a part in this project. VLISP PreScheme is the source language that our compiler translates. It is the least restrictive PreScheme dialect we consider. Macro-Free PreScheme programs result from VLISP PreScheme by expanding all derived syntax and performing a few simple transformations. Pure PreScheme programs are syntactically restricted, strongly typed Macro-Free PreScheme programs. The syntactic restrictions used to define Pure PreScheme imply that these programs will meet all of the run-time conditions for a valid VLISP PreScheme. Pure PreScheme programs are strongly typed, so no operator will be applied to data of the wrong type. `lambda` expressions in Pure PreScheme programs may occur only as initializers in top-level `letrec` bindings, or in the operator position of a procedure call. As a result, there is no need to represent closures at run-time. Finally, Pure PreScheme's syntax forbids non-tail-recursive procedure calls.

We present in more detail the various dialects of PreScheme used in the VLISP project and their design rationales. We also discuss the process by which the denotational semantics of PreScheme was derived.

## 1.3.   The transformational front end and its denotational correctness

The transformational front end transforms VLISP PreScheme into Pure PreScheme. We give an overview of the transformation process and describe some of the transformations, including simplifiers, $\beta$-reduction, inlining, and lambda-lifting. We give an example to illustrate the operation of the compiler. Finally, we give an example of the correctness proof for a particular transform (letrec-lifting).

## 1.4.   Compiler-Oriented Semantics

As discussed in Section 2, Pure PreScheme's semantics is inherited from
the semantics of Macro-Free PreScheme. However, in order to justify the
compiler presented in Section 5, we reframed the semantics to distinguish
some special cases syntactically and split both the environment and the
continuation into a compile-time and run-time components. We call this
the *compiler semantics*. The compiler reorganizes these quantities so that
all the static arguments are handled before all the dynamic ones. Thus
the semantics can be thought of as taking all the compile-time data and
producing a code-like object $\pi$ that can be applied to the run-time data,
including the local data stack (*cf* the discussion in [23, p. 253]). The
code-like objects are built from combinators, which will play the role of
instructions in the later operational semantics.

We present some pieces of the compiler-oriented semantics. For each
valuation in the semantics, we present a formal specification showing the
relation that must hold between the compiler-oriented valuation and the
corresponding valuation in the original semantics. These specifications are
based on the induction hypotheses of [6]. The proof that the compiler ori-
ented semantics is the same as the original then becomes a straightforward
structural induction. We present some sample cases in Section 4.6.

## 1.5.   The Combinator Code Machine and the Compiler

The combinator machine manipulates tuples of rational trees built from
the combinators. A machine state is the 4-tuple $\langle q, u, z, h \rangle$, where $q$ is the
code which operates on a runtime environment, $u$, a stack of stackable val-
ues, $z$, and a heap $h$. These rational trees are given two semantics: one
denotational and one operational. The denotational semantics gives trees
meanings in the same semantic domains as those used for the compiler-
oriented semantics. Given the denotational semantics, the operational se-
mantics is derived by performing $\beta$-conversion on the terms used in the
compiler oriented valuations.

The Pure PreScheme compiler is similar to the compiler oriented seman-
tics except the compiler produces syntax rather than a denotation. Given a
program, the compiler is required to produce a code tree whose denotation is
the same as that given by the compiler-oriented semantics. The compiler is
produced by taking the definition of the valuations in the compiler-oriented
semantics, and replacing domain transformations by operations that pro-
duce trees. An adequacy theorem relates the operational semantics of the
combinator machine to the denotational semantics of the original program.

We conclude by discussing some of the choices made in the design of the

compiler-oriented semantics and the combinator machine.

## 1.6.  The Linear-Data Machine

The linear-data machine uses the same machine language as the combinator machine, but instead of manipulating rational trees, its programs manipulate representations of those trees in a linear data store. It also has a heap which corresponds to the heap of the combinator machine, except that it contains untagged quantities. The state of the linear-data machine is of the form $\langle q, up, sp, s, h \rangle$, where $sp$ and $up$ are pointers into the store representing the locations of the stack and the environment.

The representation of the stack and environment by these pointers is formalized. To aid in the intuition behind these formalisms, we first give an informal sketch of the representation. We then present the formal definitions. This is done by defining a *storage-layout relation* between states of the combinator-code machine and the linear-data machine. A pair of states is in this relation iff the linear-data-machine state represents the combinator-machine state. Some of the basic properties of this relation are sketched.

We define the operational semantics of the linear-data machine so that the linear-data machine simulates the behavior of the combinator machine. That is, if $L_1$ corresponds to $C_1$ and $C_1$ rewrites to $C_2$, then we want the linear-data machine to send $L_1$ to some state $L_2$ such that $L_2$ corresponds to $C_2$.

We prove that the linear-data machine satisfies the desired simulation theorem. The simulation property will need some refining, since the combinator machine keeps tag information and the linear-data machine does not. The proof relies on the global invariant that every validly compiled Pure PreScheme program runs in bounded control space.

## 1.7.  The Stored-Program Machine

The stored-program machine is very similar to the linear-data machine, except that the code is represented in cells of a linear instruction store, much like the cells of the data store. The correspondence between states of the stored-program machine and states of the linear-data machine is defined, and a simulation theorem is proved.

## 1.8.  The Linearizer and its Correctness

The assembler takes code for the combinator or linear-data machine and produces a state of the linear instruction store which corresponds to the original according to the definition in the preceding section. We therefore

call it the *linearizer*. The major technical difficulty is the treatment of branches and joins, so that code is not duplicated. The proof is by induction on the structure of the combinator code.

### 1.9. Implementation

A compiler for PreScheme was developed in parallel with this specification. It generates code for the Motorola 68000. Representing the abstract machines with the 68000 requires a mapping similar to the ones just detailed. This section informally describes that mapping. Some performance results are presented.

### 1.10. Conclusions

We present some of the conclusions we have drawn from this effort. Some comparisons with other work are made. Some alternative design decisions are sketched. Directions for future work are suggested.

## 2. VLISP PreScheme

PreScheme is a dialect of Scheme intended for systems programming. PreScheme was carefully designed so that it syntactically looks like Scheme and has similar semantics. With a little care, PreScheme programs can be run and debugged as if they were ordinary Scheme programs. This section describes VLISP PreScheme and various related dialects used in the VLISP project.

VLISP PreScheme was inspired by Scheme48 PreScheme, but differs from it in that VLISP PreScheme has no user-defined syntax, macros, or compiler directives, and in that it provides a different set of standard procedures.

Because VLISP PreScheme is intended for systems programming, programs in the language are restricted so that they make as few assumptions as possible about the facilities available at run time. In particular, the computation model underlying VLISP PreScheme has the following properties:

- VLISP PreScheme programs manipulate data objects that fit in machine words. The type of each data object is an integer, a character, a boolean, a string, a port, a pointer to an integer, or a procedure (represented as a pointer). A PreScheme data object may be a full word, without room for run-time tags. Therefore no type predicates, like those of Scheme, are possible in the language. It becomes the compiler's responsibility to ensure statically that operators are never applied to data of the wrong type.

- A valid VLISP PreScheme program can run without creating closures at run time. Thus, if `p` is a pointer to a procedure, the free variables of that procedure must be allocatable at compile time. However, a VLISP PreScheme program may textually contain procedures with free variables that are lambda-bound. The compiler must transform these programs so that they meet the run-time restriction.

- Implementations of VLISP PreScheme are required to be tail-recursive, which means that iterative processes can be expressed by means of procedure calls. When the last action taken by a procedure is a call, a tail-recursive implementation is required to eliminate the control information of the calling procedure so that the order of space growth of iterative processes is constant. Such a call is said to be tail-recursive. The requirement that implementations be tail-recursive is inherited from Scheme [1].

- All procedure calls in a running VLISP PreScheme program must be tail-recursive. Thus, the VLISP PreScheme compiler used in the VLISP project can only be used to specify iterative processes. This restriction was made for historical reasons, as described in Section 2.5. One of the lessons learned was how to eliminate this restriction, and new versions of the compiler allow non-tail-recursive calls [**?**].

A PreScheme program is said to be *legal* only if the compiler can verify statically that these properties can be guaranteed at run time. Different implementations of PreScheme may accept different sets of programs, depending on how clever the compiler is in reasoning about these properties. The strategies used by the VLISP PreScheme compiler are discussed in Section 3.

## 2.1.  VLISP PreScheme Syntax

VLISP PreScheme was carefully designed so that it syntactically looks like Scheme and has similar semantics. With a little care, VLISP PreScheme programs can be run and debugged as if they were ordinary Scheme programs.

The syntax of the VLISP PreScheme language is identical to the syntax of the language defined in the Scheme standard [1, Chapter 7] with the following exceptions:

- Every defined procedure takes a fixed number of arguments.

- The only variables that can be modified are those introduced at top level using the syntax

$$(\texttt{define } \langle \text{variable} \rangle \ \langle \text{expression} \rangle),$$

and whose name begins and ends with an asterisk and is at least three characters long. Variables so defined are called mutable variables. Note that a variable introduced at other than top level may have a name which begins and ends with an asterisk, but this practice is discouraged.

- If ⟨expression⟩ is `lambda` expression, variables can also be defined using the syntax

    (`define-integrable` ⟨variable⟩ ⟨expression⟩).

  When ⟨variable⟩ occurs in the operator position of a combination, compilers must replace it with ⟨expression⟩.

- No variable may be defined more than once.

- `letrec` is not a derived expression. The initializer for each variable bound by a `letrec` expression must be a `lambda` expression.

- Constants are restricted to integers, characters, booleans, and strings.

- Finally, a different set of primitive procedures is specified. These are listed in Figure 2

## 2.2.  Macro-Free PreScheme

Macro-Free PreScheme programs result from VLISP PreScheme programs by expanding all derived syntax except the `case` expression, identifying which variables refer to standard procedures, and replacing single armed conditionals (`if` E E) with (`if` E E (`if` #f #f)) (while both these expressions give an unspecified result in the event that the test expression is false, the latter form allows later transformation rules to be applied more uniformly; see Section 3.1). A grammar for the resulting language is shown in Figure 3.

Not every string generated by the grammar in Figure 3 is a legal Macro-Free PreScheme program. A string is a legal program only if it is generated by the grammar and the compiler can verify that the program can be transformed to meet the runtime restrictions set forth at the beginning of this section. Macro-Free PreScheme programs must also satisfy the restriction that $n$-ary standard procedures must be used at one fixed arity.

VLISP PreScheme's complete formal denotational semantics is given in [21, Section 2.3]. Following the usual conventions, it assigns meanings to constructs of an abstract syntax. The abstract syntax is very similar to the

```
not, zero?, positive?, negative?, <, <=, =, >=, >, abs, +, -,
*, quotient, remainder, ashl, ashr, low-bits,
integer->char, char->integer, char=?, char<?, char<=, char>,
char>=,
make-vector, vector-ref, vector-set!, vector-byte-ref,
vector-byte-set!,
addr<, addr=, addr+, addr-, addr<=, addr>, addr>=,
addr->integer, integer->addr, addr->string,
port->integer, integer->port,
read-char, peek-char, eof-object?, write-char, write-int,
write, newline,
force-output, null-port?,
open-input-file, close-input-port, open-output-file,
close-output-port, current-input-port, current-output-port,
read-image, write-image,
bytes-per-word, useful-bits-per-word,
exit, err
```

Figure 2: Primitive Operators in VLISP PreScheme

$$
\begin{array}{ll}
\text{K} \in \text{Con} & \text{constants} \\
\text{I} \in \text{Ide} & \text{variables} \\
\text{O} \in \text{Op} & \text{primitive operators} \\
\text{E} \in \text{Exp} & \text{expressions} \\
\text{B} \in \text{Bnd} & \text{bindings} \\
\text{P} \in \text{Pgm} & \text{programs}
\end{array}
$$

$$
\begin{array}{ll}
\text{P} & ::= (\texttt{define } \text{I})^* \text{ E} \\
\text{B} & ::= (\text{I } (\texttt{lambda } (\text{I}^*) \text{ E}))^* \\
\text{E} & ::= \text{K} \mid \text{I} \mid (\text{E } \text{E}^*) \mid (\texttt{lambda } (\text{I}^*) \text{ E}) \\
& \quad \mid (\texttt{begin } \text{E}^* \text{ E}) \mid (\texttt{letrec } (\text{B}) \text{ E}) \\
& \quad \mid (\texttt{if } \text{E E E}) \mid (\texttt{if \#f \#f}) \mid (\texttt{set! } \text{I E}) \\
& \quad \mid (\text{O } \text{E}^*) \mid (\texttt{case } \text{E } ((\text{K}) \text{ E})^* \text{ } ((\text{K}) \text{ E}))
\end{array}
$$

Figure 3: Macro-Free PreScheme Syntax

syntax of Macro-Free PreScheme, so in the interest of space, it will be used as VLISP PreScheme's abstract syntax in what follows.

The semantics is presented in a form that very closely resembles Scheme's semantics [1, Appendix A]. It uses the same mathematical conventions, and many of the standard's definitions. The VLISP PreScheme formal semantics differs from Scheme's in a small number of ways:

- All variables must be defined before they are referenced or assigned and no variable may be defined more than once.

- Lambda bound variables are immutable, so a location need not be allocated for each actual parameter of an invoked procedure. Consequently, variables may be bound to either locations or expressed values.

- VLISP PreScheme procedure values do not have a location associated with them because there is no comparison operator for procedures.

- VLISP PreScheme `letrec` is no longer a derived expression, because the immutability of lambda bound variables makes Scheme's definition of `letrec` inapplicable. Instead, `letrec` is defined by an explicit fixed point.

- Procedures always return exactly one value, so expression continuations map a single expressed value to a command continuation.

- The domain of answers, unspecified in Scheme, is specified to be the coalesced sum of a one-point domain (representing a run-time error) and the flat domain of the integers. Thus there are exactly three outcomes of a VLISP PreScheme computation: non-termination, termination in an error, or an integer. The initial continuation is specified to be $(\lambda \epsilon \sigma.(\epsilon \mid N) \text{ in } A)$.

- Finally, memory is assumed to be infinite, so the storage allocator *new* always returns a location.

The following fragments of the VLISP PreScheme semantics illustrate some of these distinctions. For example, the semantics of variable reference follows. The first two lines contains the domain equations for denoted

values and environments.

$$\delta \in D = E + L$$
$$\rho \in U = \text{Ide} \to D$$

$$\mathcal{E} : \text{Exp} \to U \to K \to C$$
$$\mathcal{E}[\![\text{I}]\!] = \lambda\rho\kappa.\,hold(\rho\text{I})(\lambda\epsilon.\,\epsilon = undefined \to wrong \text{ ``undefined variable''}, \kappa\epsilon)$$

$$hold : D \to K \to C$$
$$hold = \lambda\delta\kappa\sigma.\,\delta \in E \to send(\delta \mid E)\kappa\sigma, send(\sigma((\delta \mid L) \downarrow 1))\kappa\sigma$$

These semantics differ from Scheme's in that some variables are bound to expressed values instead of locations.

The semantics of a `lambda` expression provide another example.

$$\mathcal{L} : \text{Exp} \to U \to E$$
$$\mathcal{L}[\![(\texttt{lambda } (\text{I}^*) \text{ E})]\!] =$$
$$\quad \lambda\rho.\,(\lambda\epsilon^*\kappa.\,\#\epsilon^* = \#\text{I}^* \to \mathcal{E}[\![\text{E}]\!](extends\ \rho\text{I}^*\epsilon^*)\kappa,$$
$$\qquad\qquad wrong \text{ ``wrong number of arguments''})$$
$$\qquad in\ E$$

Notice that the meaning of a `lambda` expression depends only on its environment. In contrast, the meaning of a `lambda` expression in Scheme depends on the store as well as its environment. This difference makes it easier to justify changing the time at which a PreScheme `lambda` expression is evaluated.

The semantics of a `letrec` expression are quite different from Scheme's.

$$\mathcal{B} : \text{Bnd} \to \text{Ide}^* \to U \to E^* \to E^*$$
$$\mathcal{E} : \text{Exp} \to U \to K \to C$$
$$\mathcal{B}[\![\,]\!] = \lambda\text{I}^*\rho\epsilon^*.\,\langle\rangle$$
$$\mathcal{B}[\![(\text{I } (\texttt{lambda } (\text{I}^*) \text{ E})) \text{ B}]\!] =$$
$$\quad \lambda\text{I}_0^*\rho\epsilon^*.\,\langle\mathcal{L}[\![(\texttt{lambda } (\text{I}^*) \text{ E})]\!](extends\ \rho\text{I}_0^*\epsilon^*)\rangle \,\S\, \mathcal{B}[\![\text{B}]\!]\text{I}_0^*\rho\epsilon^*$$
$$\mathcal{E}[\![(\texttt{letrec } (\text{B}) \text{ E})]\!] = \lambda\rho\kappa.\,\mathcal{E}[\![\text{E}]\!](extends\ \rho(\mathcal{I}[\![\text{B}]\!])(fix\,(\mathcal{B}[\![\text{B}]\!](\mathcal{I}[\![\text{B}]\!])\rho)))\kappa$$

The meaning of the body of a `letrec` expression is given with an environment extended by values derived from the `letrec` bindings. As with `lambda` expressions, the values depend only on the environment and not the store.

The specification of the meaning of a procedure call is textually unchanged from the Scheme semantics. However, the semantics is somewhat different, because the domains are different. Scheme's expression continuations map a sequence of expressed values into a command continuation,

but VLISP PreScheme maps a single expressed value into a command continuation.

$$\mathcal{E}[\![(\mathrm{E}_0 \ \mathrm{E}^*)]\!] =$$
$$\lambda \rho \kappa. \, \mathcal{E}^* \, (permute \, (\langle \mathrm{E}_0 \rangle \ \S \ \mathrm{E}^*))$$
$$\rho$$
$$(\lambda \epsilon^*. \, ((\lambda \epsilon^*. \, applicate \, (\epsilon^* \downarrow 1)(\epsilon^* \dagger 1)\kappa)$$
$$(unpermute \ \epsilon^*)))$$

The motivation for all these changes was efficiency for the Scheme48 version of PreScheme. However, these changes profoundly impact the ability to prove interesting properties of programs based on their semantics. In particular, these changes greatly facilitate proofs that certain program transformations preserve the meanings of programs. The utility of the changed semantics will be demonstrated later in the paper.

### 2.3. Static Semantics

Macro-Free PreScheme programs may be strongly typed. As in Standard ML [15], types are inferred, not declared, but unlike Standard ML, there are no polymorphic variables. All expressions are monomorphic except (if #f #f) and (set! I E). The static semantics is given as a conventional set of inference rules, except that when an expression is unconstrained by the rules, the expression is assigned the integer type [21, Section 3.3].

The standard results from the theory of type inference imply that the operators in a strongly typed Macro-Free PreScheme program will never be applied to data of the wrong type.

### 2.4. Pure PreScheme

Pure PreScheme programs are syntactically restricted, strongly typed Macro-Free PreScheme programs. The syntax is given in Figure 4. The syntactic restrictions used to define Pure PreScheme imply that these programs will meet all of the run-time conditions of a VLISP PreScheme program presented at the beginning of this section. Pure PreScheme programs are strongly typed, so no operator will be applied to data of the wrong type. lambda expressions in Pure PreScheme programs may occur only as initializers in top-level letrec bindings, or in the operator position of a procedure call. As a result, there is no need to represent closures at run-time. Finally, Pure PreScheme's syntax forbids non-tail-recursive procedure calls.

There is a further syntactic restriction. The first selection criteria of a case clause must be zero and the selection criteria for other clauses must be the successor of the previous clause's selection criterion. This restriction is

$$
\begin{array}{lll}
\text{K} & \in \text{Con} & \text{constants} \\
\text{I} & \in \text{Ide} & \text{variables} \\
\text{O} & \in \text{Op} & \text{primitive operators} \\
\text{S} & \in \text{Smpl} & \text{simple expressions} \\
\text{T} & \in \text{Tail} & \text{tail recursive expressions} \\
\text{E} & \in \text{Exp} & \text{top level expressions} \\
\text{B} & \in \text{Bnd} & \text{bindings} \\
\text{P} & \in \text{Pgm} & \text{programs}
\end{array}
$$

P ::= (define I)$^*$ E
E ::= (letrec (B) T)
B ::= (I (lambda (I$^*$) T))$^*$
T ::= S | (S S$^*$) | (if S T T) | (case S ((K) T)$^*$)
       | (begin S$^*$ T) | ((lambda (I$^*$) T) S$^*$)
S ::= K | I | (O S$^*$) | (if S S S) | (if #f #f) | (set! I S)
       | (begin S$^*$ S) | ((lambda (I$^*$) S) S$^*$) | (case S ((K) S)$^*$)

Figure 4: Pure PreScheme Syntax

imposed to allow an `case` expression to be compiled into a computed goto.

Pure PreScheme's semantics are inherited from Macro-Free PreScheme's semantics, but they assume particular values for *permute* and *unpermute*, that is, the arguments of a call are always evaluated left-to-right, and then the operator is evaluated.

## 2.5.   Challenges in the Design of VLISP PreScheme

The presentation of the PreScheme dialects used in the VLISP project has been made as logical and tidy as possible. This presentation, however, deprives the reader of an appreciation of the complexities of the design process. At the beginning of the VLISP project, we knew that the architecture of VLISP would be modeled after Scheme48, but we did not know how to build a verified compiler for a language like PreScheme. We also knew that the VLISP byte-code interpreter would mostly be an iterative process.

We began by concentrating on the run-time conditions for VLISP Pre-Scheme programs presented at the beginning of this section. Pure Pre-Scheme was the first language defined, and its syntax and semantics directly reflected the run-time conditions. When it was defined, Pure PreScheme's syntax and semantics were given independently; only later did we work out

the simple embedding of that semantics in the more general semantics of VLISP PreScheme discussed here.

Unfortunately, the VLISP byte-code interpreter is not conveniently described using only tail-recursive calls. While the byte-code interpreter is a finite state machine, it would have been more natural to allow non-tail-recursive calls to the interpreter's garbage collector. The implementation of a purely iterative byte-code interpreter is described in [8, Section 1.2.2]. In retrospect, we suffered from attempting to produce a language too finely tuned to our specific task. When we found our analysis of the task was wrong, the language was not general enough to conveniently handle the task when properly specified.

## 3. The Transformational Front End

The Front End translates VLISP PreScheme into Pure PreScheme. The translation occurs in three phases.

**Parse:** Expands usages of derived syntax by rules consistent with those presented in the Scheme standard. In addition, the program's bound variables are changed so that no variable occurs both bound and free, and no variable is bound more than once. This implies that a variable bound in the program is changed if it is also bound to a standard procedure. Other syntactic checks are made. The result is a Macro-Free PreScheme program.

**Apply transformation rules:** Translates a Macro-free PreScheme program into the syntax of Pure PreScheme using meaning-refining transformations. More will be said about this phase later.

**Type check:** Ensures that a syntactically restricted Macro-Free PreScheme program is strongly typed, and is therefore a Pure PreScheme program. This phase implements Algorithm W. That algorithm's correctness was demonstrated by Robin Milner [14]. The unifier is based on a published program by Laurence Paulson [19, p. 381].

A compiler for for VLISP PreScheme cannot accept all syntactically correct programs, because many Macro-Free PreScheme programs cannot be translated into Pure PreScheme. Furthermore, the set of Macro-Free programs that can be translated by a particular compiler depends on the transformation rules used by that compiler.

While there is no precise characterization of what constitutes a VLISP PreScheme program, knowledgeable programmers know one when they see

one. Their intuition is based on an understanding of the run-time conditions and the knowledge that existing compilers perform $\beta$-conversion, procedure inlining, and closure hoisting.

For example, the following VLISP PreScheme program cannot be translated into Pure PreScheme. The procedure `loop` must be unwound if there are to be only tail-recursive procedure calls, but the compiler cannot predict how many times it should be unwound because `*x*` is a mutable variable.

```
(define *x* 2)
(define-integrable (loop x)
  (if (positive? x) (loop (- x 1)) x))
(some-hairy-command-which-modifies-*x*!)
(+ (loop *x*) 3)
```

In contrast, the following loop will be unwound at compile-time.

```
(define bytes-per-word 4)
(define-integrable (floor-log2 x a)
  (if (<= x 1)
      a
      (floor-log2 (quotient x 2) (+ 1 a))))
(define log-bytes-per-word
  (floor-log2 bytes-per-word 0))
log-bytes-per-word
```

The most complex and error-prone phase of the VLISP PreScheme Front End translates Macro-free PreScheme programs into a syntactically restricted form. The program is transformed by applying a set of rules. Each rule can be shown to be *meaning-refining*, in a sense to be defined below.

The selection and application of rules is performed by an intricate set of procedures. However, since each rule is meaning-refining, the output of the transformation process will be a program that is a refinement of the original, regardless of the choice of transformation rules. Therefore, the verification effort focused solely on the transformation rules. The consequences of this strategy are discussed in Section 3.3.

## 3.1.   Transformation Rules

Each rule is a conditional rewrite rule. It has a pattern, a predicate, and a replacement. An expression matches a pattern if there is an assignment of pattern variables which makes the two expressions equal. The rewrite is performed if the matching expression satisfies the predicate.

To avoid name conflicts, the matching system works only on expressions that are $\alpha$-converted. An expression is *$\alpha$-converted* if no variable occurs

both bound and free, and no variable is bound more than once. The output of the parsing phase is already $\alpha$-converted, and the matching system $\alpha$-converts the output of each transformation step.

A rule with pattern $E_0$ and replacement $E_1$ is written $E_0 \Longrightarrow E_1$, and its predicate is given in the text. If no predicate is given, the predicate is assumed to be true of any $\alpha$-converted expression.

The Front End implements a variety of rules [21, Section 5.1.2]. There are roughly five categories of rules. One category contains rules specific to primitive operators. These rules help implement constant folding. For example, the rules associated with the plus operator follow.

| | |
|---|---|
| $(+ \ K_0 \ K_1) \Longrightarrow K_0 + K_1$ | constant folding |
| $(+ \ 0 \ E) \Longrightarrow E$ | idempotency |
| $(+ \ E \ K) \Longrightarrow (+ \ K \ E)$ | commutativity |
| $(+ \ E_0 \ (+ \ E_1 \ E_2)) \Longrightarrow (+ \ (+ \ E_0 \ E_1) \ E_2)$ | associativity |

The next category of rules focus on conditional expressions. As with the previous category of rules, the form of the rules are quite conventional, as the following example shows.

**Definition 1** *The* `if`*-in-an-*`if`*'s-Consequence Rule.*

$$(\text{if } E_0 \ (\text{if } E_0 \ E_1 \ E_2) \ E_3) \Longrightarrow (\text{if } E_0 \ E_1 \ E_3)$$
$$(\text{if } E_0 \ E_1 \ (\text{if } E_0 \ E_2 \ E_3)) \Longrightarrow (\text{if } E_0 \ E_1 \ E_3)$$

*when* $E_0$ *is side-effect-free.*

The predicate for this rule requires that a certain expression is side-effect-free. Side-effect-free expressions are defined so as to allow them to be recognized by a simple recursive function on syntax.

**Definition 2** *The* side-effect-free *expressions are defined inductively by the following rules:*

- K *is side-effect-free.*

- I *is side-effect-free.*

- *If* O *is side-effect-free, and* $E_1 \ldots E_n$ *are side-effect-free, then so is* $(O \ E_1 \ldots E_n)$.

- $(\text{lambda } (I^*) \ E)$ *is side-effect-free.*

- *If* $E_0, \ldots, E_n$ *are side-effect-free, then so is* $(\text{begin } E_0 \ldots E_n)$.

- *If* E *is side-effect-free, then so is* (`letrec` (B) E).

- *If* E$_0$, E$_1$, *and* E$_2$ *are side-effect-free, then so is* (`if` E$_0$ E$_1$ E$_2$).

- *If* E$_0$, . . . , E$_n$ *are side-effect-free, then so is* (`case` E$_0$ ((K) E$_1$) . . .).

Intuitively, if an expression is side-effect-free, it returns a value without modifying the store, as will be shown in Theorem 1. This expression can be eliminated when its value is ignored.

There are several other recursively defined predicates on expression syntax which are used by a variety of rules. There is a predicate which determines if a variable is free in an expression. Another predicate is true of expressions which are side-effect-free and have values which do not depend on modifiable values. Yet another one recognizes expressions that do not modify the store until their last action.

The predicates have been carefully chosen so as to work with the rule application mechanism. For example, the definition of side-effect-free expressions could have included the following clause.

- If E$_0$, . . . , E$_n$ are side-effect-free, then ((`lambda` (I$_1$ . . . I$_n$) E$_0$) E$_1$ . . . E$_n$) is side-effect-free.

This clause was eliminated not because its presence adversely affected proofs, but because it caused looping of the rule application mechanism.

The remaining categories of rules implement $\beta$-conversion, closure hoisting, and defined constant substitution. Many of the rules in these categories look unusual. The forms of these rules was intended to facilitate correctness proofs. For example, two distinct rules produce the effect of $\beta$-conversion. Closures are hoisted by the use of several other rules, including the rule below.

**Definition 3** *The* `letrec` *Expression Merging Rule.*

$$(\text{\texttt{letrec}} \ (B_0) \ (\text{\texttt{letrec}} \ (B_1) \ E)) \Longrightarrow (\text{\texttt{letrec}} \ (B_0 \ B_1) \ E)$$

The rules result in program transformations similar to those produced by other compilers [11, 12]. Consider the even-odd program given in Figure 5. It specifies a simple iterative process. Figures 6–11 show how this program might be translated into Pure PreScheme using transformations implemented in the Front End. Notice that unlike [11, 12], this compiler does not convert the program into continuation-passing style [2, 22].

```
(define number 77)                      ; This program
(define (dec x) (- x 1))                ; produces zero
(define (even x)                        ; if NUMBER is even
  (if (zero? x)                         ; and non-negative,
      0                                 ; otherwise it
      (odd (dec x))))                   ; produces one.
(define (odd x)
  (if (zero? x)
      1
      (even (dec x))))
(if (negative? number)
    1
    (even number))
```

Figure 5: Example VLISP PreScheme Program

```
(define number)
(define dec)
(define even)
(define odd)
(begin
  (set! number 77)
  (set! dec (lambda (x) (- x 1)))
  (set! even (lambda (y) (if (= 0 y) 0 (odd (dec y)))))
  (set! odd (lambda (z) (if (= 0 z) 1 (even (dec z)))))
  (if (> 0 number) 1 (even number)))
```

Figure 6: Expand into Macro-Free PreScheme

```
(define number) (define dec) (define even) (define odd)
(begin
  (set! number 77)
  (set! dec (letrec ((dec-0 (lambda (x) (- x 1))) dec-0)))
  (set! even (letrec ((even-0 (lambda (y)
                                (if (= 0 y)
                                    0
                                    (odd (dec y)))))
                      even-0)))
  (set! odd (letrec ((odd-0 (lambda (z)
                              (if (= 0 z)
                                  1
                                  (even (dec z)))))
                     odd-0)))
  (if (> 0 number) 1 (even number)))
```

Figure 7: Name Anonymous Lambda Expressions

```
(define number) (define dec) (define even) (define odd)
(letrec
    ((dec-0 (lambda (x) (- x 1)))
     (even-0 (lambda (y) (if (= 0 y) 0 (odd (dec y)))))
     (odd-0 (lambda (z) (if (= 0 z) 1 (even (dec z))))))
  (begin
    (set! number 77)
    (set! dec dec-0)
    (set! even even-0)
    (set! odd odd-0)
    (if (> 0 number) 1 (even number))))
```

Figure 8: Closure Hoisting

```
(define number) (define dec) (define even) (define odd)
(letrec
    ((dec-0 (lambda (x) (- x 1)))
     (even-0 (lambda (y) (if (= 0 y) 0 (odd-0 (dec-0 y)))))
     (odd-0 (lambda (z) (if (= 0 z) 1 (even-0 (dec-0 z))))))
  (begin
    (set! number 77) (set! dec dec-0)
    (set! even even-0) (set! odd odd-0)
    (if (> 0 77)                           ; Which simplifies
        1                                  ; to (even-0 77).
        (even-0 77))))
```

Figure 9: Constant Folding

```
(define number) (define dec) (define even) (define odd)
(letrec
    ((dec-0 (lambda (x) (- x 1)))
     (even-0 (lambda (y) (if (= 0 y) 0 (odd-0 (- y 1)))))
     (odd-0 (lambda (z) (if (= 0 z) 1 (even-0 (- z 1))))))
  (begin
    (set! number 77) (set! dec dec-0)
    (set! even even-0) (set! odd odd-0)
    (even-0 77)))
```

Figure 10: Inline Non-Tail-Recursive Procedure Calls

```
(define number) (define dec) (define even) (define odd)
(letrec
    ((even-0 (lambda (y) (if (= 0 y) 0 (odd-0 (- y 1)))))
     (odd-0 (lambda (z) (if (= 0 z) 1 (even-0 (- z 1))))))
  (begin
    (set! number 0) (set! dec 0)
    (set! even 0) (set! odd 0)
    (even-0 77)))
```

Figure 11: Eliminate Unused Lambda Expressions and Initializers

## 3.2.  Justification of the Rules

The application of a rule is justified if it transforms a Macro-free Pre-Scheme program into another, and both programs have the same meaning as given by the formal semantics. Some of the rules have another interesting property—they can transform a program which has bottom denotation into a program which produces a non-bottom answer. For example, the program

```
(define two (+ 1 one))
(define one 1)
two
```

is transformed into a program which produces the answer 2!

This odd behavior is tolerated so as to allow constant propagation without performing a dependency analysis. In the above example, 1 is substituted for the occurrence of the immutable variable `one` even though *undefined* should have been substituted.

**Definition 4** *A rule is* meaning-refining *if its application does not affect non-bottom computational results.*

The proof that a rule is meaning-refining usually has the following form. Suppose $E_1$ is the expression that results from applying the rule to expression $E_0$. The proof shows that $\mathcal{E}[\![E_0]\!]\rho\kappa\sigma \sqsubseteq \mathcal{E}[\![E_1]\!]\rho\kappa\sigma$.

The reason the proof shows that a rule is meaning-refining follows. Let $P_1$ be a program which results from the application of the rule to an expression in program $P_0$, and let $\mathcal{P}$ be the semantic function for programs. The proof implies that $\mathcal{P}[\![P_0]\!] \sqsubseteq \mathcal{P}[\![P_1]\!]$, because the semantic functions are compositional. Hence, if $\mathcal{P}[\![P_0]\!]$ is not bottom, we have $\bot \sqsubset \mathcal{P}[\![P_0]\!] \sqsubseteq \mathcal{P}[\![P_1]\!]$. But since the answer domain is flat, $\mathcal{P}[\![P_0]\!] = \mathcal{P}[\![P_1]\!]$.

The proof of each rule usually requires the consideration of many details. Many rules are proved using structural induction on the syntax of Macro-Free PreScheme. Some rules have predicates which oblige certain expressions not modify the store, others require certain expressions not modify the store or reference mutable data. Some of the proofs in [21] are incomplete as only the cases that appeared to us to be "interesting" were selected for detailed consideration.

As a result, there is no proof which is representative of all of the others, However, the proofs of the `if`-*in-an*-`if`'s-Consequence Rule and the `letrec` Expression Merging Rule give some flavor of the other proofs.

### 3.2.1.   The `if`-*in-an*-`if`'s-Consequence Rule

This rule has a predicate which requires one of its expressions to be side-effect-free. Side-effect-free expressions have the following property $\Pi$.

**Definition 5** E *has property* $\Pi$ *iff* $(\forall \rho, \sigma)(\exists \epsilon)(\forall \kappa)$

$$\mathcal{E}[\![E]\!]\rho\kappa\sigma = (\epsilon = undefined \to \bot_A, \kappa\epsilon\sigma)$$

To understand this definition, ignore the case in which $\epsilon = undefined$. Notice the store given to the command continuation $\mathcal{E}[\![E]\!]\rho\kappa$ on the left-hand-side is the same as the one given to the command continuation $\kappa\epsilon$ on the right-hand side. Thus, the meaning of E is either erroneous, or is the same as passing some value $\epsilon$ to $\kappa$, without changing the store. This characterizes the the meaning of side-effect-free expressions in a form useful for proofs.

**Theorem 1** *If* E *is side-effect-free, then* E *has property* $\Pi$.

> **Proof:** This is proved by structural induction on side-effect-free expressions. The cases of E being I and (`if` $E_0$ $E_1$ $E_2$) are shown.
>
> Case E = I: Let $\epsilon = \rho I \in E \to \rho I \mid E, \sigma(\rho I \mid L) \downarrow 1$. Expanding definitions gives
>
> $$\mathcal{E}[\![I]\!]\rho\kappa\sigma = (\epsilon = undefined \to \bot_A, \kappa\epsilon\sigma).$$
>
> Notice $\epsilon$ is independent of $\kappa$ so
>
> $$(\forall\kappa)\mathcal{E}[\![I]\!]\rho\kappa\sigma = (\epsilon = undefined \to \bot_A, \kappa\epsilon\sigma).$$
>
> Case E = (`if` $E_0$ $E_1$ $E_2$): By the induction hypothesis, there is at least one $\epsilon_0$ such that
>
> $$(\forall\kappa)\mathcal{E}[\![E_0]\!]\rho\kappa\sigma = (\epsilon_0 = undefined \to \bot_A, \kappa\epsilon_0\sigma).$$
>
> If $\epsilon_0 = undefined$ the result is immediate. Otherwise,
>
> $$\begin{aligned} &\mathcal{E}[\![(\texttt{if } E_0 \texttt{ } E_1 \texttt{ } E_2)]\!]\rho\kappa\sigma \\ &= \mathcal{E}[\![E_0]\!]\rho(\lambda\epsilon.\ truish\ \epsilon \to \mathcal{E}[\![E_1]\!]\rho\kappa, \mathcal{E}[\![E_2]\!]\rho\kappa)\sigma \\ &= truish\ \epsilon_0 \to \mathcal{E}[\![E_1]\!]\rho\kappa\sigma, \mathcal{E}[\![E_2]\!]\rho\kappa\sigma. \end{aligned}$$
>
> When $\epsilon_0 = false$,
>
> $$\mathcal{E}[\![(\texttt{if } E_0 \texttt{ } E_1 \texttt{ } E_2)]\!]\rho\kappa\sigma = \mathcal{E}[\![E_2]\!]\rho\kappa\sigma,$$
>
> otherwise
>
> $$\mathcal{E}[\![(\texttt{if } E_0 \texttt{ } E_1 \texttt{ } E_2)]\!]\rho\kappa\sigma = \mathcal{E}[\![E_1]\!]\rho\kappa\sigma.$$
>
> Use of the induction hypothesis verifies both alternatives. ∎

**Theorem 2** *If* $E_0$ *is side-effect-free,*

$$\mathcal{E}[\![(\texttt{if } E_0 \texttt{ (if } E_0 \texttt{ } E_1 \texttt{ } E_2\texttt{) } E_3)]\!]\rho\kappa\sigma = \mathcal{E}[\![(\texttt{if } E_0 \texttt{ } E_1 \texttt{ } E_3)]\!]\rho\kappa\sigma.$$

*Hence the* `if`-*in-an-*`if`*'s-Consequence Rule is meaning-refining.*

**Proof:**   By Theorem 1, for any $\rho$, $\kappa$, and $\sigma$, there exists an $\epsilon_0$ such that

$$(\forall\kappa)\mathcal{E}[\![\text{E}_0]\!]\rho\kappa\sigma = (\epsilon_0 = \textit{undefined} \rightarrow \perp_A, \kappa\epsilon_0\sigma).$$

If $\epsilon_0 = \textit{undefined}$, the proof is immediate, so assume $\epsilon_0 \neq \textit{undefined}$.

$\mathcal{E}[\![(\text{if } \text{E}_0 \ (\text{if } \text{E}_0 \ \text{E}_1 \ \text{E}_2) \ \text{E}_3)]\!]\rho\kappa\sigma$
$\quad = \mathcal{E}[\![\text{E}_0]\!]\rho(\lambda\epsilon.\, \textit{truish}\,\epsilon \rightarrow \mathcal{E}[\![(\text{if } \text{E}_0 \ \text{E}_1 \ \text{E}_2)]\!]\rho\kappa, \mathcal{E}[\![\text{E}_3]\!]\rho\kappa)\sigma$
$\quad = \textit{truish}\,\epsilon_0 \rightarrow \mathcal{E}[\![(\text{if } \text{E}_0 \ \text{E}_1 \ \text{E}_2)]\!]\rho\kappa\sigma, \mathcal{E}[\![\text{E}_3]\!]\rho\kappa\sigma$
$\quad = \textit{truish}\,\epsilon_0 \rightarrow \mathcal{E}[\![\text{E}_0]\!]\rho(\lambda\epsilon.\, \textit{truish}\,\epsilon \rightarrow \mathcal{E}[\![\text{E}_1]\!]\rho\kappa, \mathcal{E}[\![\text{E}_2]\!]\rho\kappa)\sigma, \mathcal{E}[\![\text{E}_3]\!]\rho\kappa\sigma$
$\quad = \textit{truish}\,\epsilon_0 \rightarrow \textit{truish}\,\epsilon_0 \rightarrow \mathcal{E}[\![\text{E}_1]\!]\rho\kappa\sigma, \mathcal{E}[\![\text{E}_2]\!]\rho\kappa\sigma, \mathcal{E}[\![\text{E}_3]\!]\rho\kappa\sigma$
$\quad = \textit{truish}\,\epsilon_0 \rightarrow \mathcal{E}[\![\text{E}_1]\!]\rho\kappa\sigma, \mathcal{E}[\![\text{E}_3]\!]\rho\kappa\sigma$
$\quad = \mathcal{E}[\![\text{E}_0]\!]\rho(\lambda\epsilon.\, \textit{truish}\,\epsilon \rightarrow \mathcal{E}[\![\text{E}_1]\!]\rho\kappa, \mathcal{E}[\![\text{E}_3]\!]\rho\kappa)\sigma$
$\quad = \mathcal{E}[\![(\text{if } \text{E}_0 \ \text{E}_1 \ \text{E}_3)]\!]\rho\kappa\sigma$

The proof of the other form of the rule is similar.  ∎

### 3.2.2.   The `letrec` Expression Merging Rule

The correctness proof of this rule is interesting because it demonstrates the use of fixed points in a denotational semantics to prove something that would be difficult using an operational semantics.

The proof makes use of the following lemmas which are presented without proofs. The *extends* auxiliary function is similar to Scheme's except that it associates identifiers with expressed values instead of locations. Recall from [1, Appendix A] that § denotes concatenation of sequences.

**Lemma 1**  *If* $\text{I}_0^* \ \S \ \text{I}_1^*$ *is a sequence of distinct identifiers, then*

$$\textit{extends}\,(\textit{extends}\,\rho\text{I}_0^*\epsilon_0^*)\text{I}_1^*\epsilon_1^* = \textit{extends}\,\rho(\text{I}_0^* \ \S \ \text{I}_1^*)(\epsilon_0^* \ \S \ \epsilon_1^*).$$

**Lemma 2**

$\mathcal{B}[\![\text{B}_0\text{B}_1]\!](\mathcal{I}[\![\text{B}_0\text{B}_1]\!])\rho$
$\quad = \lambda\epsilon^*.\, \mathcal{B}[\![\text{B}_0]\!](\mathcal{I}[\![\text{B}_0]\!])$
$\qquad\qquad\qquad (\textit{extends}\,\rho(\mathcal{I}[\![\text{B}_1]\!])(\textit{dropfirst}\,\epsilon^*\#\text{B}_0))$
$\qquad\qquad\qquad (\textit{takefirst}\,\epsilon^*\#\text{B}_0)$
$\qquad\qquad \S\ \mathcal{B}[\![\text{B}_1]\!](\mathcal{I}[\![\text{B}_1]\!])$
$\qquad\qquad\qquad (\textit{extends}\,\rho(\mathcal{I}[\![\text{B}_0]\!])(\textit{takefirst}\,\epsilon^*\#\text{B}_0))$
$\qquad\qquad\qquad (\textit{dropfirst}\,\epsilon^*\#\text{B}_0)$

The following theorem shows that the `letrec` expression merging rule is meaning-refining:

**Theorem 3** *If $\mathcal{I}[\![B_0 B_1]\!]$ is a sequence of distinct identifiers, then*

$$\mathcal{E}[\![(\texttt{letrec } (B_0) \ (\texttt{letrec } (B_1) \ E))]\!]\rho\kappa\sigma$$
$$= \mathcal{E}[\![(\texttt{letrec } (B_0 \ B_1) \ E)]\!]\rho\kappa\sigma.$$

**Proof:**  Let

$$\text{Let } f_0 = \mathcal{B}[\![B_0]\!](\mathcal{I}[\![B_0]\!])\rho,$$
$$\rho' = extends \ \rho(\mathcal{I}[\![B_0]\!])(fix \ f_0),$$
$$f_1 = \mathcal{B}[\![B_1]\!](\mathcal{I}[\![B_1]\!])\rho'.$$

Then we have

$$\mathcal{E}[\![(\texttt{letrec } (B_0) \ (\texttt{letrec } (B_1) \ E))]\!]\rho\kappa\sigma$$
$$= \mathcal{E}[\![(\texttt{letrec } (B_1) \ E)]\!]\rho'\kappa\sigma$$
$$= \mathcal{E}[\![E]\!](extends \ \rho'(\mathcal{I}[\![B_1]\!])(fix \ f_1))\kappa\sigma$$

Because expressions are $\alpha$-converted,

$$extends \ \rho'(\mathcal{I}[\![B_1]\!])(fix \ f_1)$$
$$= extends \ \rho(\mathcal{I}[\![B_0 B_1]\!])(fix \ f_0 \ \S \ fix \ f_1).$$

Let $f_{01} = \mathcal{B}[\![B_0 B_1]\!](\mathcal{I}[\![B_0 B_1]\!])\rho.$

$$\mathcal{E}[\![(\texttt{letrec } (B_0 \ B_1) \ E)]\!]\rho\kappa\sigma$$
$$= \mathcal{E}[\![E]\!](extends \ \rho(\mathcal{I}[\![B_0 B_1]\!])(fix \ f_{01}))\kappa\sigma$$

The proof is completed by showing $fix \ f_{01} = fix \ f_0 \ \S \ fix \ f_1$.

$$
\begin{aligned}
f_{01} = \ &\lambda\epsilon^*.\, \mathcal{B}[\![B_0]\!](\mathcal{I}[\![B_0]\!]) \\
&\qquad (extends \ \rho(\mathcal{I}[\![B_1]\!])(dropfirst \ \epsilon^* \#B_0)) \\
&\qquad (takefirst \ \epsilon^* \#B_0) \\
&\quad \S \ \mathcal{B}[\![B_1]\!](\mathcal{I}[\![B_1]\!]) \\
&\qquad (extends \ \rho(\mathcal{I}[\![B_0]\!])(takefirst \ \epsilon^* \#B_0)) \\
&\qquad (dropfirst \ \epsilon^* \#B_0) \\
= \ &\lambda\epsilon^*.\, \mathcal{B}[\![B_0]\!](\mathcal{I}[\![B_0]\!])\rho(takefirst \ \epsilon^* \#B_0) \\
&\quad \S \ \mathcal{B}[\![B_1]\!](\mathcal{I}[\![B_1]\!]) \\
&\qquad (extends \ \rho(\mathcal{I}[\![B_0]\!])(takefirst \ \epsilon^* \#B_0)) \\
&\qquad (dropfirst \ \epsilon^* \#B_0) \\
= \ &\lambda\epsilon^*.\, f_0(takefirst \ \epsilon^* \#B_0) \\
&\quad \S \ \mathcal{B}[\![B_1]\!](\mathcal{I}[\![B_1]\!]) \\
&\qquad (extends \ \rho(\mathcal{I}[\![B_0]\!])(takefirst \ \epsilon^* \#B_0)) \\
&\qquad (dropfirst \ \epsilon^* \#B_0)
\end{aligned}
$$

because no binding in $B_0$ references a variable bound by $B_1$.

Let $g = \lambda\epsilon^*.\, f_0(takefirst \ \epsilon^* \#B_0) \ \S \ dropfirst \ \epsilon^* \#B_0$. Superscripts will denote function iteration: $f^0 = \lambda\epsilon^*.\, \epsilon^*$ and $f^{n+1} = f \circ f^n$. Observe that $f_{01}^n(fix \ g) = fix \ f_0 \ \S \ f_1^n \bot$, therefore, $\bigsqcup\{f_{01}^n(fix \ g)\} = fix \ f_0 \ \S \ fix \ f_1$.

$\mathit{fix}\ f_0\ \S\ \mathit{fix}\ f_1$ is a fixed point of $f_{01}$ because

$$f_{01}(\mathit{fix}\ f_0\ \S\ \mathit{fix}\ f_1)$$
$$= f_{01}(\bigsqcup\{f_{01}^n(\mathit{fix}\ g)\})$$
$$= \bigsqcup\{f_{01}^{n+1}(\mathit{fix}\ g)\} \qquad \text{by continuity}$$
$$= \bigsqcup\{f_{01}^n(\mathit{fix}\ g)\} \qquad \text{as } \mathit{fix}\ g \sqsubseteq f_{01}(\mathit{fix}\ g)$$
$$= \mathit{fix}\ f_0\ \S\ \mathit{fix}\ f_1.$$

$\mathit{fix}\ f_0\ \S\ \mathit{fix}\ f_1$ is the least fixed point of $f_{01}$ because, by construction, $g^m\bot \sqsubseteq f_{01}^m\bot$ so $f_{01}^n(g^m\bot) \sqsubseteq f_{01}^{m+n}\bot$.

$$f_{01}^n(\mathit{fix}\ g) = f_{01}^n(\bigsqcup\{g^m\bot\}) = \bigsqcup\{f_{01}^n(g^m\bot)\}$$
$$\sqsubseteq \bigsqcup\{f_{01}^n(f_{01}^m\bot)\} = f_{01}^n(\mathit{fix}\ f_{01}) = \mathit{fix}\ f_{01}$$

Therefore $f_{01}^n(\mathit{fix}\ g) \sqsubseteq \mathit{fix}\ f_{01}$ and $\mathit{fix}\ f_{01} = \mathit{fix}\ f_0\ \S\ \mathit{fix}\ f_1$. ∎

## 3.3.  Selection and Proof of Rules

The most significant contribution of the work on the Front End is the identification of a collection of transformation rules that can both be verified relative to the formal semantics of the source language, and can also form the basis of a practical optimizing compiler.

It is difficult to quantify the quality of the rule set implemented in the Front End. The Front End implements many of the rules used in successful compilers, and so its performance should be comparable. We studied the result of translating the VLISP Byte-Code Interpreter into Pure PreScheme and could identify no other plausible general transformation that would improve the code. The general nature of each implemented rule suggests the rule set should be useful for a wide range of VLISP PreScheme programs.

There were two distinct prototypes of the Front End before the final version was implemented. Early guesses of the form the transformation rules were quite naïve, because we tried to use rules designed for unverified compilers. The first prototype explored various rule sets. It also became apparent that the mechanism used in the prototype to control rule application would not scale so as to allow the compilation of a program as large as the VLISP Byte-Code Interpreter. The second prototype had the control mechanism nearly right and its rule set seemed to generate reasonable code. Attempts to prove the correctness of its rules failed miserably. The exercise led to changes in both the rules and VLISP PreScheme's semantics. The prototypes gave us extremely valuable information and made the identification of the final, finely tuned rule set possible. Our motto is prototype but verify!

There are several unverified aspects of the Front End. The algorithm used to translate VLISP PreScheme into Macro-Free PreScheme was not

studied rigorously. The lack of rigor is of little concern because the process is clearly specified in the Scheme Report section on derived expression types [1, Section 7.4].

The most unsatisfying part of the Front End is the fact that the algorithms describing the control procedures are unverified. All parts of the Front End were written in a simple and direct style. We studiously avoided tricky or obscure coding practices, partially so that someone reading the code which applies transformation rules will easily conclude that the control procedures only cause changes to the program which are consistent with the set of rules, but there is no formal verification of this fact. Further work is needed to design a formalism that when given a rule set produces an efficient and verified rule application mechanism.

## 4. Compiler-Oriented Semantics

As discussed in Section 2, Pure PreScheme's semantics is inherited from the semantics of Macro-Free PreScheme. However, in order to justify the compiler presented in Section 5, we reframed the semantics to distinguish some special cases syntactically and split both the environment and the continuation into a compile-time and run-time components. We call this the *compiler semantics*.

The architecture of this semantics is motivated by a binding-time analysis of the semantics:

- The program is always compile-time data (that is, it is available at compile time).

- The environment is split into a symbol table $\gamma$ (available at compile time), and a run-time environment or display $\mu$ (available at run-time).

- The continuation is also split into static (compile-time) and dynamic (run-time) components.

- The heap $h$ is always run-time data.

In addition, the compiler semantics introduces a data stack for the values of subexpressions. This stack is "partially static" in that its size is statically determinable, but its contents are dynamic.

The compiler semantics reorganizes these quantities so that all the static arguments are handled before all the dynamic ones. Thus the semantics can be thought of as taking all the compile-time data and producing a code-like

object $\pi$ that can be applied to the run-time data, including the local data stack (*cf* the discussion in [23, p. 253]).

In addition, the semantics distinguishes between simple expressions and tail-recursive expressions, following the grammar of Figure 4. An easy proof via structural induction shows that valuations of tail-recursive expressions always use the initial continuation $\kappa_0$, a fact exploited in the new semantics.

The compiler-oriented semantics also distinguishes commands (expressions executed only for effect) from expressions executed for their value. An easy proof via structural induction shows that valuations of commands always use a continuation which ignores its value. This separation allows the compiler to know that the value of a simple expression can be ignored. In principle, we could have preserved the semantic principle of "one non-terminal, one valuation" by creating a separate syntactic category for commands, but the programming principle of orthogonality ("a phrase should be usable anywhere it is meaningful") won out, since in Scheme in any (simple) expression can be evaluated for effect only.

As noted above, the new semantics relies on the usual division of environments into a compile-time and a run-time component. We introduce a new domain $D_c$ of compile-time denotations (lexical addresses), and represent an environment $\rho : Ide \to D$ by a symbol table $\gamma : Ide \to D_c$ and a runtime environment $\mu : D_c \to D$ such that $\rho = (\mu \circ \gamma)$. We need to ensure that the run-time and compile-time extension functions behave consistently. The behavior of these functions must satisfy the following for any $\mu$ and $\gamma$:

$$extends(\mu \circ \gamma)\, \mathrm{L}^* \epsilon^* = (extends_{rl}\ \mu\ \epsilon^*) \circ (extends_{cl}\ \gamma\ \mathrm{L}^*)$$

with similar conditions for global variables.[1] Because the scoping structure of Pure PreScheme programs is very simple, our semantics uses integers for lexical addresses.

The domain of code-like objects $\pi$ is

$$Q = U_r \to E^* \to S \to A$$

where $U_r$ is the domain of run-time environments (displays), $E$ is the domain of expressed values, $S$ is the domain of heaps, and $A$ is the domain of answers. These elements are the compile-time component of a continuation. In general, the elements of $Q$ will be built out of certain combinators, which we call *instructions*. The justification for this terminology will become evident in Section 5.

---

[1] For historical reasons, this semantics assumes that the identifiers have been partitioned into local and global identifiers. Global identifiers can only be used for defined variables, and local identifiers can only be used for `letrec` and `lambda` bound variables.

We next present some pieces used to construct the compiler oriented semantics, $\mathcal{CP}[\![P]\!]$. For each valuation we give the following data:

1. an informal specification of the valuation

2. the domain of the valuation

3. the induction hypothesis (formal specification) for the valuation

4. the rules for each production, and last

5. The semantics of typical instructions introduced in those rules

The proof that the compiler oriented semantics is the same as the original, i.e., $\mathcal{CP}[\![P]\!] = \mathcal{P}[\![P]\!]$, then becomes a straightforward structural induction. We will present some sample cases in Section 4.6.

## 4.1. Simple Expressions

Compile Simple Expressions. Generates code that evaluates the expression and pushes the value onto the runtime stack.[2]

**type:** $\qquad \mathcal{CSE} : \ \text{Smpl} \to U_c \to Q \to Q$
**specification:** $\quad (\mathcal{CSE}[\![\text{S}]\!]\gamma\pi)\mu\zeta = \mathcal{E}[\![\text{S}]\!](\mu \circ \gamma)(\lambda\epsilon.\,\pi\mu(\epsilon :: \zeta))$

The induction hypothesis states that running a simple expression compiled using symbol table $\gamma$ with code to follow $\pi$ in runtime environment $\mu$ with runtime stack $\zeta$ is the same as the semantics for that expression using composition of $\mu$ and $\gamma$ for an environment with a continuation executing $\pi$ in with the same runtime environment but with the result pushed onto the runtime stack.

Note that $\mathcal{CSE}$ takes all the static arguments before all the dynamic arguments, as discussed above.

**definition:**
$\mathcal{CSE}[\![\text{L}]\!] = \lambda\gamma\pi.\,fetch_l\,(\gamma\,\text{L})\,\pi$
$\mathcal{CSE}[\![(\texttt{if}\ \text{S}_0\ \text{S}_1\ \text{S}_2)]\!] = \lambda\gamma\pi.\,\mathcal{CSE}[\![\text{S}_0]\!]\gamma(brf\,(\mathcal{CSE}[\![\text{S}_1]\!]\gamma(label\ \pi))$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad (\mathcal{CSE}[\![\text{S}_2]\!]\gamma(goto\ \pi)))$
$\mathcal{CSE}[\![(\texttt{set!}\ \text{G}\ \text{S})]\!] = \lambda\gamma\pi.\,\mathcal{CSE}[\![\text{S}]\!]\gamma(update\text{-}store\,(\gamma\,\text{G})\,\pi)$
$\mathcal{CSE}[\![(\text{O}\ \text{S}^*)]\!] = \lambda\gamma\pi.\,\mathcal{CSE}^*[\![\text{S}^*]\!]\gamma(prim\text{-}apply\ \#\text{S}^*\,\mathcal{O}[\![\text{O}]\!]\,\pi)$

---

[2]The runtime stack in this semantics contains expressed values. The actual compiler allows environments as well [18]. The use of $\zeta$ instead of $\epsilon^*$ reminds us of the fact that the presentation has been simplified.

These clauses introduce six new instructions.

$fetch_l : \ D_c \to Q \to Q$
$fetch_l = \lambda \iota \pi. \ \lambda \mu \zeta. (\mu \ \iota) \ \text{E} \ E \to$
$\qquad\qquad\qquad \pi \mu (((\mu \ \iota) \mid E) :: \zeta),$
$\qquad\qquad\qquad (wrong \ \text{"Local variable given storage."})$

$brf : \ Q \to Q \to Q$
$brf = \lambda \pi_0 \pi_1 \lambda \mu \zeta. \ truish((\zeta \downarrow 1)) \to (\pi_0 \mu (\zeta \dagger 1)), \ (\pi_1 \mu (\zeta \dagger 1))$

$label : \ Q \to Q$
$label = \lambda \pi \mu \zeta. \pi \mu \zeta$

$goto : \ Q \to Q$
$goto = \lambda \pi \mu \zeta. \pi \mu \zeta$

The *label* and *goto* instructions are used as directives to the linearizer (Section 8).

$update\text{-}store : \ D_c \to Q \to Q$
$update\text{-}store = \lambda \iota \pi. \ \lambda \mu \zeta. (\mu \ \iota) \ \text{E} \ L \to$
$\qquad\qquad\qquad\qquad (assign \ ((\mu \ \iota) \mid L) \ (\zeta \downarrow 1) \ (\pi \mu \zeta)),$
$\qquad\qquad\qquad\qquad (wrong \ \text{"Can't assign to a local variable"})$

$prim\text{-}apply : \ N \to P \to Q \to Q$
$prim\text{-}apply = \lambda \nu \upsilon \pi. \ \lambda \mu \zeta. (apply\text{-}primitive \ \upsilon \ (take\text{-}first \ \nu \ \zeta)$
$\qquad\qquad\qquad\qquad (\lambda \epsilon. \ \pi \mu (\epsilon :: (pop\text{-}first \ \nu \ \zeta))))$

## 4.2.  Simple Expression Sequences

Compile Simple Expression Sequences. Generates code that evaluates the sequence of simple expressions and pushes the sequence of values onto the runtime stack.

**type:**        $\mathcal{CSE}^* : \ \text{Smpl}^* \to U_c \to Q \to Q$
**specification:**   $(\mathcal{CSE}^*[\![\text{S}^*]\!] \gamma \pi) \mu \zeta = \mathcal{E}^*[\![\text{S}^*]\!] (\mu \circ \gamma)(\lambda \epsilon^*. \ \pi \mu (\epsilon^* \ \S \ \zeta))$

The induction hypothesis states that running simple expression sequences compiled using symbol table $\gamma$ with code to follow $\pi$ in runtime environment $\mu$ with runtime stack $\zeta$ is the same as the semantics using composition of $\mu$ and $\gamma$ for an environment with a continuation executing $\pi$ in with the same runtime environment but with the resulting values pushed onto the runtime stack.

$\mathcal{CSE}^*[\![ \ ]\!] = \lambda \gamma \pi. \ \pi$
$\mathcal{CSE}^*[\![\text{S S}^*]\!] = \lambda \gamma \pi. \mathcal{CSE}[\![\text{S}]\!] \gamma (\mathcal{CSE}^*[\![\text{S}^*]\!] \gamma \pi)$

## 4.3.  Simple Commands

Compile Simple Commands.  Generates code that evaluates the expression and but does not push the value onto the runtime stack.

**type:** $\quad\quad\quad\mathcal{CSC} : \text{ Smpl} \to U_c \to Q \to Q$
**specification:** $\quad(\mathcal{CSC}[\![\text{S}]\!]\gamma\pi)\mu\zeta = \mathcal{E}[\![\text{S}]\!](\mu \circ \gamma)(\lambda\epsilon.\,\pi\mu\zeta)$

The induction hypothesis states that running simple commands compiled using symbol table $\gamma$ with code to follow $\pi$ in runtime environment $\mu$ with runtime stack $\zeta$ is the same as the semantics using composition of $\mu$ and $\gamma$ for an environment with a continuation executing $\pi$ in with the same runtime environment and stack (i.e. the value is ignored).

**definition:**
$\mathcal{CSC}[\![\text{L}]\!] = \lambda\gamma\pi.\,\pi$
$\mathcal{CSC}[\![(\text{if } \text{S}_0 \ \text{S}_1 \ \text{S}_2)]\!] = \lambda\gamma\pi.\,\mathcal{CSE}[\![\text{S}_0]\!]\gamma(brf(\mathcal{CSC}[\![\text{S}_1]\!]\gamma\pi)\ (\mathcal{CSC}[\![\text{S}_2]\!]\gamma\pi))$
$\mathcal{CSC}[\![(\text{set! } \text{G } \text{S})]\!] = \lambda\gamma\pi.\,\mathcal{CSE}[\![\text{S}]\!]\gamma(update\text{-}store_i\ (\gamma\ \text{G})\ \pi)$
$\mathcal{CSC}[\![(\text{O } \text{S}^*)]\!] = \lambda\gamma\pi.\,\mathcal{CSE}^*[\![\text{S}^*]\!]\gamma(prim\text{-}apply_i\ \#\text{S}^*\ \mathcal{O}[\![\text{O}]\!]\ \pi)$

This valuation introduces two new instructions, $update\text{-}store_i$ and $prim\text{-}apply_i$, which are like their expression counterparts but do not push a value onto the runtime stack.

## 4.4.  Tail-Recursive Expressions

Compile Tail-recursive Expressions.  Generates code that evaluates the tail-recursive expression and pushes the value onto the runtime stack.  As the expressions are tail-recursive, no static continuation is necessary for compilation.

**type:** $\quad\quad\quad\mathcal{CT} : \text{ Tre} \to U_c \to Q$
**specification:** $\quad(\mathcal{CT}[\![\text{T}]\!]\gamma)\mu\langle\rangle = \mathcal{E}[\![\text{T}]\!](\mu \circ \gamma)\kappa_0$

The induction hypothesis states that running tail recursive expressions compiled using symbol table $\gamma$ in runtime environment $\mu$ with an empty runtime stack is the same as the semantics using composition of $\mu$ and $\gamma$ for an environment.

$\mathcal{CT}[\![\text{S}]\!] = \lambda\gamma.\,\mathcal{CSE}[\![\text{S}]\!]\gamma(halt)$
$\mathcal{CT}[\![(\text{if } \text{S } \text{T}_0 \ \text{T}_1)]\!] = \lambda\gamma.\,\mathcal{CSE}[\![\text{S}]\!]\gamma(brf\,(\mathcal{CT}[\![\text{T}_0]\!]\gamma)\ (\mathcal{CT}[\![\text{T}_1]\!]\gamma))$

$\mathcal{CT}[\![(\text{begin } S^* \text{ T})]\!] = \lambda\gamma.\mathcal{CSC}[\![S^*]\!]\gamma(\mathcal{CT}[\![T]\!]\gamma)$
$\mathcal{CT}[\![((\text{lambda } (L^*) \text{ T}) S^*)]\!] = \lambda\gamma.\mathcal{CSE}^*[\![S^*]\!]\gamma$
$$(\textit{add-to-env} (\mathcal{CT}[\![T]\!](\textit{extends}_{c\,l} \ \gamma \ L^*)))$$
$\mathcal{CT}[\![(S \ S^*)]\!] = \lambda\gamma.\mathcal{CSE}^*[\![S^*]\!]\gamma(\mathcal{CSE}[\![S]\!]\gamma(\textit{tail-call}))$

These clauses introduce two more new instructions, *add-to-env* and *tail-call*:

$\textit{add-to-env} : \ Q \to Q$
$\textit{add-to-env} = \lambda\pi.\,\lambda\mu\zeta.\,\pi(\textit{extends}_{r\,l} \ \mu \ \zeta)\langle\rangle$

$\textit{tail-call} : \ Q$
$\textit{tail-call} = \lambda\mu\zeta.\,(\textit{tail-apply} \ (\zeta \downarrow 1) \ (\zeta \dagger 1))$

where *tail-apply* is defined by

$\textit{tail-apply} : \ E \to E^* \to C$
$\textit{tail-apply} = \lambda\epsilon\epsilon^*.\,\epsilon \text{ E } F \to (\epsilon \mid F)\epsilon^*\kappa_0, \ (\textit{wrong} \text{ "Non-function to apply"})$

## 4.5.  Recursive Procedures

Compile Recursive Declarations. Takes a set of local procedure declarations and generates code that extends the run-time environment by making the indicated mutually-recursive declarations.

**type:**              $\mathcal{CE} : \ \text{Exp} \to U_c \to Q$
**specification:**
$(\mathcal{CE}[\![(\text{letrec } (B) \text{ T})]\!]\gamma)\mu\langle\rangle = \mathcal{E}[\![(\text{letrec } (B) \text{ T})]\!](\mu \circ \gamma)\kappa_0$
**definition:**
$\mathcal{CE}[\![(\text{letrec } (B) \text{ T})]\!]\gamma = \textit{closerecs}\,(\mathcal{CB}[\![B]\!](\mathcal{GP}[\![B]\!])\gamma)(\mathcal{CT}[\![T]\!]\gamma)$

The induction hypothesis says that executing the code generated by the compiler is equivalent to executing the code for T in an appropriate recursively-extended environment.

The auxiliary valuation $\mathcal{CB}$ produces a sequence of code segments which will turn into the bodies of the procedures.

**type:**              $\mathcal{CB} : \ \text{Bnd} \to \text{LocIde}^* \to U_c \to Q^*$
**specification:**
$\textit{closerecs}\,(\mathcal{CB}[\![B]\!](\mathcal{GP}[\![B]\!])\gamma)\pi\mu = \pi(\textit{extends}_{rl}\,\mu(\textit{fix}(\mathcal{B}[\![B]\!](\mathcal{GP}[\![B]\!])(\mu \circ \gamma))))$
**definition:**
$\mathcal{CB}[\![ \ ]\!] = \lambda L^*\gamma.\langle\rangle$
$\mathcal{CB}[\![(L \ (\text{lambda}(L^*) \text{ T})) \text{ B}]\!] =$
  $\lambda L_0^*\gamma.\mathcal{CT}[\![T]\!](\textit{extends}_{cl}(\textit{extends}_{cl} \ \gamma L_0^*)L^*) :: (\mathcal{CB}[\![B]\!]L_0^*\gamma)$

This valuation introduces a single new instruction *closerecs*, defined by:

$$closerecs \, \pi^* \pi = \lambda\mu\zeta. \, \pi(\mathit{fix}(\lambda\mu'. \, extends_{rl} \, \mu(map(close_r \, \mu')\pi^*)))\zeta$$

where

$$close_r = \lambda\mu\pi. \, (\lambda\epsilon^*. \, \pi(extends_{rl} \, \mu\epsilon^*)\langle\rangle) \, \mathrm{in} \, E$$

## 4.6. Correctness of the Compiler-Oriented Semantics

The main statement of correctness of the compiler is:

**Theorem 4 (Correctness of the compiler oriented semantics)** *For any program* P, $\mathcal{P}[\![P]\!] = \mathcal{CP}[\![P]\!]$.

The proof is a tedious but straightforward structural induction. There are a total of 13 simultaneous induction hypotheses, one for each function in the compiler. There is one induction step for each clause in the definitions. Most of the induction steps are easy calculations: the only exceptions are the ones for recursive declarations.

Here we include a few sample cases to illustrate the calculations. We show a few simple cases and the cases for recursive declarations in more detail.

We begin with two cases for simple expressions. The induction hypothesis to be verified is

**specification:**
$$(\mathcal{CSE}[\![S]\!]\gamma\pi)\mu\zeta = \mathcal{E}[\![S]\!](\mu \circ \gamma)(\lambda\epsilon. \, \pi\mu(\epsilon :: \zeta))$$

The following three cases are typical.

$(\mathcal{CSE}[\![L]\!]\gamma\pi)\mu\zeta$
$=$ [by definition of $\mathcal{CSE}$]
$(\mathit{fetch}_l \, (\gamma \, L) \, \pi)\mu\zeta$
$=$ [by definition of $\mathit{fetch}_l$]
$((\mu \circ \gamma)L) \, \mathrm{E} \, E \to \pi\mu((((\mu \circ \gamma)L) \mid E) :: \zeta),$
$\qquad\qquad (\mathit{wrong} \text{ "Local variable given storage."})$
$=$ [by definition of $\mathcal{E}$]
$\mathcal{E}[\![L]\!](\mu \circ \gamma)(\lambda\epsilon. \, \pi\mu(\epsilon :: \zeta))$

$(\mathcal{CSE}[\![(\text{if } S_0 \, S_1 \, S_2)]\!]\gamma\pi)\mu\zeta$
$=$ [by definition of $\mathcal{CSE}$]

$(\mathcal{CSE}[\![S_0]\!]\gamma(\mathit{brf}\,(\mathcal{CSE}[\![S_1]\!]\gamma(\mathit{label}\,\pi))\,(\mathcal{CSE}[\![S_2]\!]\gamma(\mathit{goto}\ \pi))))\mu\zeta$

$=$     [by induction hypothesis at $S_0$]

$\mathcal{E}[\![S_0]\!](\mu\circ\gamma)(\lambda\epsilon.\,(\mathit{brf}\,(\mathcal{CSE}[\![S_1]\!]\gamma(\mathit{label}\ \pi))\ (\mathcal{CSE}[\![S_2]\!]\gamma(\mathit{goto}\ \pi)))\mu(\epsilon::\zeta))$

$=$     [by definitions of $\mathit{brf}$, $\mathit{label}$, $\mathit{goto}$]

$\mathcal{E}[\![S_0]\!](\mu\circ\gamma)(\lambda\epsilon.\,(\mathit{truish}(\epsilon)\rightarrow((\mathcal{CSE}[\![S_1]\!]\gamma\pi)\mu\zeta),$
$\qquad\qquad\qquad\qquad\qquad\qquad ((\mathcal{CSE}[\![S_2]\!]\gamma\pi)\mu\zeta)))$

$=$     [by induction hypothesis at $S_1$, $S_2$]

$\mathcal{E}[\![S_0]\!](\mu\circ\gamma)(\lambda\epsilon.\,(\mathit{truish}(\epsilon)\rightarrow\mathcal{E}[\![S_1]\!](\mu\circ\gamma)(\lambda\epsilon.\,\pi\mu(\epsilon::\gamma))$
$\qquad\qquad\qquad\qquad\qquad\quad \mathcal{E}[\![S_2]\!](\mu\circ\gamma)(\lambda\epsilon.\,\pi\mu(\epsilon::\zeta))))$

$=$     [by definition of $\mathcal{E}$]

$\mathcal{E}[\![(\texttt{if}\ S_0\ S_1\ S_2)]\!](\mu\circ\gamma)(\lambda\epsilon.\,\pi\mu(\epsilon::\zeta))$

<br>

$(\mathcal{CSE}[\![(\texttt{set!}\ G\ S)]\!]\gamma\pi)\mu\zeta$

$=$     [by definition of $\mathcal{CSE}$]

$(\mathcal{CSE}[\![S]\!]\gamma(\mathit{update\text{-}store}\ (\gamma\ G)\ \pi))\mu\zeta$

$=$     [by induction hypothesis at $S$]

$\mathcal{E}[\![S]\!](\mu\circ\gamma)(\lambda\epsilon.\,(\mathit{update\text{-}store}\ (\gamma\ G)\ \pi)\mu(\epsilon::\zeta))$

$=$     [by definition of $\mathit{update\text{-}store}$]

$\mathcal{E}[\![S]\!](\mu\circ\gamma)(\lambda\epsilon.\,((\mu\circ\gamma)G)\ \text{\sc e}\ L\rightarrow(\mathit{assign}\,(((\mu\circ\gamma)G)\mid L)\,\epsilon\,(\pi\mu(\epsilon::\zeta))),$
$\qquad\qquad\qquad\qquad\qquad\quad (\mathit{wrong}\ \text{``Can't assign to a local variable''}))$

$=$     [by definition of $\mathcal{E}$]

$\mathcal{E}[\![(\texttt{set!}\ G\ S)]\!](\mu\circ\gamma)(\lambda\epsilon.\,\pi\mu(\epsilon::\zeta))$

<br>

Several valuations recur on sequences of syntactic elements; $\mathcal{CSE}^*$ is typical. The induction hypothesis is

**specification:**

$(\mathcal{CSE}^*[\![S^*]\!]\gamma\pi)\mu\zeta=\mathcal{E}^*[\![S^*]\!](\mu\circ\gamma)(\lambda\epsilon^*.\,\pi\mu(\epsilon^*\,\S\,\zeta))$

There are two productions (clauses), which are proved as follows:

$(\mathcal{CSE}^*[\![\ ]\!]\gamma\pi)\mu\zeta$

$=$     [by definition of $\mathcal{CSE}^*$]

$\pi\mu\zeta$

$=$     [by definition of $\mathcal{E}^*$]

$\mathcal{E}^*[\![\ ]\!](\mu\circ\gamma)(\lambda\epsilon^*.\,\pi\mu(\epsilon^*\,\S\,\zeta))$

<br>

$(\mathcal{CSE}^*[\![S\ S^*]\!]\gamma\pi)\mu\zeta$

$=$     [by definition of $\mathcal{CSE}^*$]

$(\mathcal{CSE}[\![S]\!]\gamma(\mathcal{CSE}^*[\![S^*]\!]\gamma\pi))\mu\zeta$

$=$     [by induction hypothesis at $S$]

$\mathcal{E}[\![S]\!](\mu \circ \gamma)(\lambda \epsilon. (\mathcal{CSE}^*[\![S^*]\!]\gamma \pi)\mu(\epsilon :: \zeta))$
$= \quad$ [by induction hypothesis at $S^*$]
$\mathcal{E}[\![S]\!](\mu \circ \gamma)(\lambda \epsilon. \mathcal{E}^*[\![S^*]\!](\mu \circ \gamma)(\lambda \epsilon^*. \pi\mu(\epsilon^* \,\S\, (\epsilon :: \zeta))))$
$= \quad$ [by definition of $\mathcal{E}^*$]
$\mathcal{E}^*[\![S\ S^*]\!](\mu \circ \gamma)(\lambda \epsilon^*. \pi\mu(\epsilon^* \,\S\, \zeta))$

The only difficult proof is the one for recursive procedure declarations. The induction hypothesis is

**specification:**
$closerecs\,(\mathcal{CB}[\![B]\!](\mathcal{GP}[\![B]\!])\gamma)\pi\mu = \pi(extends_{rl}\,\mu(fix(\mathcal{B}[\![B]\!](\mathcal{GP}[\![B]\!])(\mu \circ \gamma))))$

This calculation is complicated by the necessity to reason both about sequences of values and about fixed points. To simplify this calculation, we write $\langle e \dots \rangle$ for a sequence of values whose typical element is given by $e$, and "X... " for a sequence of syntactic elements whose typical element is given by X.

To verify the induction hypothesis, let

$$\begin{aligned} B &= (L\ (\texttt{lambda}\ (L^*)\ T))\dots \\ L_0^* &= \mathcal{GP}[\![B]\!] \\ \gamma' &= extends_{cl}\,\gamma L_0^* \end{aligned}$$

Then, calculating from the left-hand side of the induction hypothesis, we get:

$closerecs\,(\mathcal{CB}[\![B]\!]L_0^*\gamma)\pi\mu$
$= \quad$ [by definition of $closerecs$ ]
$\pi(fix\,(\lambda\mu'.\,extends_{rl}\,\mu(map\ (close\ _r\ \mu')(\mathcal{CB}[\![B]\!]L_0^*\gamma))))$
$= \quad$ [by definition of $\mathcal{CB}$]
$\pi(fix\,(\lambda\mu'.\,extends_{rl}\,\mu(map\ (close\ _r\ \mu')\langle\mathcal{CT}[\![T]\!](extends_{cl}\,\gamma' L^*)\dots\rangle)))$
$= \quad$ [by definition of $map$ and $close\ _r$]
$\pi(fix\,(\lambda\mu'.\,extends_{rl}\,\mu\langle(\lambda\epsilon^*.\mathcal{CT}[\![T]\!](extends_{cl}\,\gamma' L^*) \circ (extends_{rl}\,\mu'\epsilon^*)\langle\rangle)\,\text{in}\,E$
$\quad\quad\quad\quad\quad\quad\quad\quad \dots\rangle))$
$= \quad$ [by induction hypothesis at T]
$\pi(fix\,(\lambda\mu'.\,extends_{rl}\,\mu\langle(\lambda\epsilon^*.\mathcal{E}[\![T]\!]((extends_{cl}\,\gamma' L^*) \circ (extends_{rl}\,\mu'\epsilon^*))\kappa_0)\,\text{in}\,E$
$\quad\quad\quad\quad\quad\quad\quad\quad \dots\rangle))$
$= \quad$ [by consistency of $extends_{rl}$, $extends_{cl}$]
$\pi(fix\,(\lambda\mu'.\,extends_{rl}\,\mu\langle(\lambda\epsilon^*.\mathcal{E}[\![T]\!](extends(\mu' \circ \gamma')L^*\epsilon^*)\kappa_0)\,\text{in}\,E\dots\rangle))$
$= \quad$ [by definition of $\mathcal{L}$]
$\pi(fix\,(\lambda\mu'.\,extends_{rl}\,\mu\langle\mathcal{L}[\![(\texttt{lambda}\ (L^*)\ T)]\!](\mu' \circ \gamma')\dots\rangle))$

Working from the right-hand side, we calculate:

$\pi(extends_{rl}\,\mu(fix(\mathcal{B}[\![\mathrm{B}]\!]\mathrm{L}_0^*(\mu \circ \gamma))))$
$=$    [by definition of $\mathcal{B}$]
$\pi(extends_{rl}\,\mu(fix(\lambda\epsilon^*.\,\langle\mathcal{L}[\![(\texttt{lambda } (\mathrm{L}^*)\ \mathrm{T})]\!](extends(\mu \circ \gamma)\mathrm{L}_0^*\epsilon^*)\ldots\rangle)))$

Let $f = \lambda\mu'.\,extends_{rl}\,\mu\langle\mathcal{L}[\![(\texttt{lambda } (\mathrm{L}^*)\ \mathrm{T})]\!](\mu' \circ \gamma')\ldots\rangle$
and $g = \lambda\epsilon^*.\,\langle\mathcal{L}[\![(\texttt{lambda } (\mathrm{L}^*)\ \mathrm{T})]\!](extends(\mu \circ \gamma)\mathrm{L}_0^*\epsilon^*)\ldots\rangle$

The theorem is proved once we show that:

$$fix\,f = extends_{rl}\,\mu(fix\,g)$$

The right-hand side is a fixed point of $f$ because:

$$
\begin{aligned}
f(extends_{rl}\,\mu\epsilon^*) &= extends_{rl}\,\mu(g\epsilon^*) & \text{[by substitution]}\\
f(extends_{rl}\,\mu(fix\,g)) &= extends_{rl}\,\mu(g(fix\,g)) & \text{[by use of previous line]}\\
&= extends_{rl}\,\mu(fix\,g)
\end{aligned}
$$

The right-hand side is a least fixed point of $f$ because:

$$extends_{rl}\,\mu(g^n\bot) \sqsubseteq f^{n+1}\bot$$

The proof is by induction on $n$.

$$
\begin{aligned}
extends_{rl}\,\mu(g^0\bot) &= extends_{rl}\,\mu\bot\\
&\sqsubseteq extends_{rl}\,\mu\langle\mathcal{L}[\![(\texttt{lambda } (\mathrm{L}^*)\ \mathrm{T})]\!](\bot \circ \gamma')\ldots\rangle\\
&= f^1\bot
\end{aligned}
$$

$$
\begin{aligned}
extends_{rl}\,\mu(g^{n+1}\bot) &= extends_{rl}\,\mu(g(g^n\bot))\\
&= f(extends_{rl}\,\mu(g^n\bot))\\
&\sqsubseteq f(f^{n+1}\bot) & \text{[by induction hypothesis]}\\
&= f^{n+2}\bot
\end{aligned}
$$

## 5. The Combinator Machine and the Compiler

### 5.1. The Combinator Machine

The combinator machine manipulates tuples of rational trees. Rational (or regular) trees [7] are finite or infinite trees that have only finitely many distinct subtrees. Such a tree can always be represented as a finite graph that unwinds to the possibly-infinite tree.

A machine state is the 4-tuple $\langle q, u, z, h\rangle$, where $q$ is the code which operates on a runtime environment, $u$, a stack of stackable values, $z$, and a heap $h$. A portion of the grammar for permissible states $\langle q, u, z, h\rangle$ is

$$
\begin{aligned}
q \quad ::= \quad & (\textit{halt}) \mid (\textit{goto } q) \mid (\textit{label } q) \mid (\textit{fetch}_l \ i \ q) \mid (\textit{brf } q_1 \ q_2) \\
\mid \quad & (\textit{save-env } q) \mid (\textit{restore-env } q) \mid (\textit{update-store } i \ q) \\
\mid \quad & (\textit{prim-apply } n \ p \ q) \mid (\textit{update-store}_i \ i \ q) \mid (\textit{prim-apply}_i \ n \ p \ q) \\
\mid \quad & (\textit{add-to-env } n \ q) \mid (\textit{add-to-env}_g \ q) \\
\mid \quad & (\textit{closerecs } (q_1 \ldots q_n) \ q) \mid (\textit{tail-call}) \\[1em]
u \quad ::= \quad & \textit{emptydisplay} \mid (\textit{extends}_{rl} \ u(v_1 \ldots v_n)) \mid (\textit{extends}_{rg} \ ua) \\[1em]
z \quad ::= \quad & v^* \\[1em]
v \quad ::= \quad & \langle \mathbf{proc}, \langle q, u \rangle \rangle \mid \langle \mathbf{env}, u \rangle \mid \langle \mathbf{int}, n \rangle \mid \langle \mathbf{bool}, b \rangle \\
\mid \quad & \langle \mathbf{char}, c \rangle \mid \langle \mathbf{string}, s \rangle \mid \langle \mathbf{hptr}, n \rangle \mid \langle \mathbf{quote}, d \rangle \\[1em]
a \quad ::= \quad & \langle \mathbf{hptr}, n \rangle
\end{aligned}
$$

Figure 12: Grammar for states of the combinator machine

shown in Figure 12. The combinator machine has a total of 20 instructions, of which 15 are shown. The permissible states of the machine are those configurations $\langle q, u, z, h \rangle$ where $q$, $u$, $z$, and $h$ are elements of the greatest fixed point defined by the operators of the grammar. In fact, the only infinite trees that arise are those generated for recursive environments.

The stack $z$ consists of a list of stackable values $v$. These values are procedures (closed in the representation $u_0$ of a global environment, to be described later), environments, and other values, which we call *immediate* values. These immediate values are integers, booleans, characters, strings, pointers into the heap (L-values), and quotations. L-values are tagged integers. Quotations represent data returned by primitives for use only with other primitives (primitives like *make-vector*, *vector-ref*, etc).

We have not included a grammar for the heap; it consists of three components: a map from heap pointers (as above) to immediate values, an integer (representing a free-location counter), and an unspecified third element. The first two components are used for ordinary mutable variables. The third component can be manipulated only by primitives (primitives like *make-vector*, *vector-ref*, etc); it plays a role analogous to that of the file system in a conventional language semantics. We use the notation $h.1$, $h.2$, and $h.3$ for the three components.

The operational semantics of each machine instruction is shown in Figure 13. The machine halts normally by executing a *halt* instruction, returning the value $\langle \mathbf{ok}, n \rangle$ for some integer $n$. The last rule in Figure 13

specifies that if the machine state does not match any of the preceding left-hand sides, then the machine goes into an error state and halts, returning the value $\langle \mathbf{error} \rangle$. $\mathtt{fix}$ builds an infinite rational tree with a loop in it. For example, if $t = \mathtt{fix}(\lambda u.\, extends_{rg}\, ua)$, then $t$ is a tree such that $t = extends_{rg}\, ta$. It is easy to see that there is exactly one infinite tree that has this property, namely the one generated by the finite graph in which the left subtree of the $extends_{rg}$ loops back to the top of the tree [7]. The auxiliary $\mathtt{lookup}$ is an operation on syntax trees that simulates environment lookup. $\mathtt{truish}$ is a test on rational trees, simulating the operator $truish$ of [1, Appendix A]. $\mathtt{apply\text{-}prim}$ handles the application of primitive procedures $p$.

Each set of terms manipulated has an analogous domain in the compiler oriented semantics for Pure PreScheme. The relation between syntactic and denotational variables is shown in Figure 14. Each set also has a natural valuation motivated by the form of the equations of the compiler oriented semantics. Some sample clauses of these valuations are:

$$\mathcal{U}[\![emptydisplay]\!] = \mu_0 \qquad \text{the initial runtime environment}$$

$$\mathcal{Q}[\![(brf\ q_0\ q_1)]\!] = \lambda\mu\zeta.\, truish((\zeta \downarrow 1)) \to \mathcal{Q}[\![q_0]\!]\mu(\zeta \dagger 1),\ \mathcal{Q}[\![q_1]\!]\mu(\zeta \dagger 1)$$

$$\mathcal{V}[\![\langle \mathbf{proc}, \langle q, u \rangle \rangle]\!] = \lambda\epsilon^*.\, \mathcal{Q}[\![q]\!](extends_{rl}(\mathcal{U}[\![u]\!])\epsilon^*) \text{ in } E$$

The semantics for machine states is given by

$$\mathcal{M}[\![\langle q, u, z, h \rangle]\!] = (\mathcal{Q}[\![q]\!])(\mathcal{U}[\![u]\!])(\mathcal{Z}[\![z]\!])(\mathcal{H}[\![h]\!])$$

We have given these valuations as if they were for finite trees. To accommodate infinite trees, we actually define a sequence of approximate valuations; for example, we write

$$\mathcal{V}_0[\![v]\!] = \bot$$
$$\mathcal{V}_{n+1}\langle \mathbf{proc}, \langle q, u \rangle \rangle = \lambda\epsilon^*.\, \mathcal{Q}_n[\![q]\!](extends_{rl}(\mathcal{U}_n[\![u]\!])\epsilon^*) \text{ in } E$$

Given the valuations, the transitions of the machine were derived by performing $\beta$-conversion on the terms used in the compiler oriented valuations. The $brf$ rule is derived by simple calculations:

$$\mathcal{M}[\![\langle (brf\ q_0\ q_1), u, (v :: z), h \rangle]\!]$$
$$=\quad [\text{by definition of } \mathcal{M}]$$
$$\mathcal{Q}[\![(brf\ q_0\ q_1)]\!](\mathcal{U}[\![u]\!])(\mathcal{Z}[\![(v :: z)]\!])(\mathcal{H}[\![h]\!])$$
$$=\quad [\text{by definition of } \mathcal{Z}]$$
$$\mathcal{Q}[\![(brf\ q_0\ q_1)]\!](\mathcal{U}[\![u]\!])(\mathcal{V}[\![v]\!] :: \mathcal{Z}[\![z]\!])(\mathcal{H}[\![h]\!])$$

Abstract machine:

$$\langle(halt), u, (v :: z), h\rangle \implies \langle\mathbf{ok}, v\rangle$$
$$\langle(goto\ q), u, z, h\rangle \implies \langle q, u, z, h\rangle$$
$$\langle(label\ q), u, z, h\rangle \implies \langle q, u, z, h\rangle$$
$$\langle(fetch_l\ i\ q), u, z, h\rangle \implies \langle q, u, ((\mathtt{lookup}\ u\ i) :: z), h\rangle$$
$$\langle(brf\ q_0\ q_1), u, (v :: z), h\rangle \implies \langle(\mathtt{truish}(v) \rightarrow q_0, q_1), u, z, h\rangle$$
$$\langle(save\text{-}env\ q), u, z, h\rangle \implies \langle q, u, (\langle\mathbf{env}, u\rangle :: z), h\rangle$$
$$\langle(restore\text{-}env\ q), u, z, h\rangle \implies \langle q, u', (v :: z'), h\rangle$$
$$\text{where } z = (v :: \langle\mathbf{env}, u'\rangle :: z')$$
$$\langle(update\text{-}store\ i\ q), u, (v :: z), h\rangle \implies \langle q, u, (v :: z), (\mathtt{update}\ (\mathtt{lookup}\ u\ i)\ v\ h')\rangle$$
$$\langle(prim\text{-}apply\ n\ p\ q), u, z, h\rangle \implies \langle q, u, (w_1 :: z'), h_1\rangle$$
$$\text{where } z = (v_1 :: \ldots :: v_n :: z')$$
$$\text{and } (w_1, h_1) = (\mathtt{apply\text{-}prim}\ p\ (v_1 \ldots v_n)h)$$
$$\langle(add\text{-}to\text{-}env\ n\ q), u, z, h\rangle \implies \langle q, (extends_{rl}\ u\ (v_1 \ldots v_n)), z', h\rangle$$
$$\text{where } z = (v_1 :: \ldots :: v_n :: z')$$
$$\langle(add\text{-}to\text{-}env_g\ q), u, (v :: z), h\rangle \implies$$
$$\langle q, (extends_{rg}\ u\ (\mathtt{new}\ h)), z, (\mathtt{update}\ (\mathtt{new}\ h)\ v\ h)\rangle$$
$$\langle(closerecs\ (q_1 \ldots q_n)\ q), u, z, h\rangle \implies$$
$$\langle q, (\mathtt{fix}\ (\lambda u'.\ (extends_{rl}\ u\ (\langle\mathbf{proc}, \langle q_1, u'\rangle\rangle \ldots \langle\mathbf{proc}, \langle q_n, u'\rangle\rangle)))), z, h\rangle$$
$$\langle(tail\text{-}call), u, z_1, h\rangle \implies \langle q, (extends_{rl}\ u_1\ (v_1 \ldots v_n)), \langle\rangle, h\rangle$$
$$\text{where } z_1 = (\langle\mathbf{proc}, \langle q, u_1\rangle\rangle :: v_1 :: \ldots :: v_n :: \langle\rangle)$$
$$\text{otherwise} \implies \langle\mathbf{error}\rangle$$

Auxiliaries:

$$\mathtt{lookup}\ (extends_{rl}\ u(v_0 \ldots v_{n-1}))\ (i) = \begin{cases} v_i & 0 \le i < n \\ \mathtt{lookup}\ u(i - n) & n \le i \end{cases}$$
$$\mathtt{lookup}\ (extends_{rg}\ u\ a)\ (0) = a$$
$$\mathtt{lookup}\ (extends_{rg}\ u\ a)\ (i + 1) = \mathtt{lookup}\ u\ i$$
$$\mathtt{lookup}\ emptydisplay\ i \text{ is an error}$$

Figure 13: Operational Semantics of the Combinator Machine

| Syntactic | Denotational | Description |
|-----------|--------------|-------------|
| $v$ | $\epsilon$ | expressed value |
| $q$ | $\pi$ | code sequence |
| $u$ | $\mu$ | runtime environment |
| $z$ | $\zeta$ | stack |
| $h$ | $\sigma$ | store, also called the heap |
| $i$ | $\iota$ | global variable reference |
| $a$ | $\alpha$ | location |
| $g$ | $\gamma$ | symbol table |

Figure 14: Analogous syntactic and denotational variables

$$= \quad [\text{by definition of } \mathcal{Q} \text{ and } brf]$$
$$truish(\mathcal{V}[\![v]\!]) \to \mathcal{Q}[\![q_0]\!](\mathcal{U}[\![u]\!])(\mathcal{Z}[\![z]\!])(\mathcal{H}[\![h]\!]), \mathcal{Q}[\![q_1]\!](\mathcal{U}[\![u]\!])(\mathcal{Z}[\![z]\!])(\mathcal{H}[\![h]\!])$$
$$= \quad [\text{by definition of } \texttt{truish}\,]$$
$$\texttt{truish}(v) \to \mathcal{Q}[\![q_0]\!](\mathcal{U}[\![u]\!])(\mathcal{Z}[\![z]\!])(\mathcal{H}[\![h]\!]), \mathcal{Q}[\![q_1]\!](\mathcal{U}[\![u]\!])(\mathcal{Z}[\![z]\!])(\mathcal{H}[\![h]\!])$$
$$= (\texttt{truish}(v) \to \mathcal{Q}[\![q_0]\!], \mathcal{Q}[\![q_0]\!])\,(\mathcal{U}[\![u]\!])\,(\mathcal{Z}[\![z]\!])\,(\mathcal{H}[\![h]\!])$$
$$= \mathcal{Q}[\![\texttt{truish}(v) \to q_0,\, q_1]\!](\mathcal{U}[\![u]\!])(\mathcal{Z}[\![z]\!])(\mathcal{H}[\![h]\!])$$
$$= \mathcal{M}[\![\langle(\texttt{truish}(v) \to q_0,\, q_1), u, z, h\rangle]\!]$$

Observe that the reduction semantics of the instructions in the compiler oriented semantics determine the transitions of the combinator machine. One could regard each instruction as a combinator, so that each transition corresponds to the reduction of one combinator. This interpretation of the machine gives it its name.

The soundness of the machine can be stated as follows:

**Theorem 5 (Soundness of the Combinator Machine)** *If* $\langle q, u, z, h\rangle \Longrightarrow$ $\langle q', u', z', h'\rangle$, *then* $\mathcal{M}[\![\langle q, u, z, h\rangle]\!] = \mathcal{M}[\![\langle q', u', z', h'\rangle]\!]$.

> **Proof:**  By calculation for each transition in Figure 13. In the case of a *closerecs*, we unwind the fixed-point one time to see that $u'$ is a legal runtime environment. The valuation $\mathcal{Q}$ is carefully written to check the length of the stack and the tags of the operands so that the error transition is sound. ∎

## 5.2.   The Combinator-Code Compiler

The Pure PreScheme compiler is similar to the compiler oriented semantics except the compiler produces syntax rather than a denotation. The compiler is patterned after the semantics in that there is function for each

valuation. For example, for the compilation of simple expressions, the analog of $\mathcal{CSE}$ is cse, which has the following specification:

$$\mathcal{Q}[\![\text{cse}\,[\![\text{S}]\!]gq]\!] = \mathcal{CSE}[\![\text{S}]\!](\mathcal{G}[\![g]\!])(\mathcal{Q}[\![q]\!])$$

where $\mathcal{G}$ is the valuation for symbol tables.

The function cse could be the lambda term constructed by replacing the denotational variables with their analogs in the definition of $\mathcal{CSE}$. Normal order reduction could be used to specify when cse is defined. This definition of cse obviously meets its specification. In practice, Scheme program text defined cse. For example, we could write

```
(define cse
  (lambda (exp symtab q)
    (record-case exp
      (local-var (l)
        (fetch-l (apply-symtab symtab l) q))
      ...)))
```

This Scheme code could be used to simulate $\mathcal{CSE}$ by defining, following Section 4,

```
(define fetch-l
  (lambda (lex-addr q)
    (lambda (mu zeta)
      (if (expressed-value? (mu lex-addr))
        (q mu (cons (restrict-to-E (mu lex-addr))
                    zeta))
        (wrong "Local variable given storage.")))))
```

or it could be used as part of the definition of cse by defining

```
(define fetch-l
  (lambda (lex-addr q)
    (list 'fetch-l lex-addr q)))
```

which emits a representation of a *fetch_l* instruction.

The syntactic structure of Pure PreScheme programs allows us to characterize the structure of a compiled program $\text{cp}[\![\text{P}]\!]$. Every compiled program consists of a prelude that allocates globals with a sequence of *add-to-env_g* instructions, followed by a *closerecs* instruction that creates all the procedures. This prelude is followed by the code for the body of the program. The effect of the prelude is to create an environment $u_0$ in which all the global procedures are closed.

We can combine this with the correctness of the compiler-oriented semantics (Theorem 4) to get the following:

**Theorem 6 (Adequacy)** *For any program* P *and answer* $a \neq \bot$, $\mathcal{P}[\![P]\!] = a$ *if and only if* $\langle \mathtt{cp}\,[\![P]\!], emptydisplay, \langle\rangle, h_0 \rangle \Longrightarrow a$

> **Proof:** The "if" follows immediately from the correctness of the compiler-oriented semantics, (Theorem 4), the soundness of the combinator machine, and the operational semantics of the *halt* instruction. Hence if the machine halts (either with a number or an error), the machine's answer and the denotational answer coincide.
>
> To prove the "only if", we must show that if the machine runs infinitely, then the denotation must be bottom. To do this, recall the definition of a reduction strategy from [3, Chapter 13]. A reduction strategy is a function $F$ on $\lambda$-terms such that for all $M$, $M \rightarrow^* F(M)$. Then the soundness theorem above asserts that the combinator machine induces a reduction strategy on all terms of the form $quzh$. Furthermore, this strategy is quasi-leftmost, since each machine-induced reduction begins with a head reduction. Since any quasi-leftmost strategy is normalizing (Theorem 13.2.6 of [3]), if the machine runs infinitely long from state $\langle q, u, z, h\rangle$, then $quzh$ (thought of as a $\lambda$-term) has no normal form. In fact, since the reduction sequence has infinitely many head reductions, this term has no head normal form, and therefore denotes $\bot$ in the semantic domains. The desired result then follows from Theorem 4. ∎

An important property of the combinator machine is that the size of the terms (except for the heap) is statically bounded. To express this, we define $dom(u) = \{i \mid \mathtt{lookup}\ u\ i\,\text{is not an error}\}$. Then we have:

**Theorem 7 (Bounded Space)** *If* P *is a valid Pure PreScheme program, then there exist integers* $N_0$ *and* $N$ *such that for any state* $\langle q, u, z, h\rangle$ *in the computation sequence of* $\langle \mathtt{cp}\,[\![P]\!], emptydisplay, \langle\rangle, h_0 \rangle$ , $|dom(u)| \leq N$. *Furthermore, after the execution of the single closerecs instruction in* $\mathtt{cp}\,[\![P]\!]$,

$$N_0 \leq |dom(u)| \leq N$$

> **Proof:** (Sketch) A Pure PreScheme program must consist of a set of global simple declarations followed by a set of local procedure declarations followed by a tail-recursive expression. Let $N_0$ be total number of global simple declarations and local procedure declarations, and let $N$ be $N_0$ plus the deepest lexical depth in the program. Since all procedures are closed in the global environment, which has $N_0$ elements, it must be that $N_0 \leq |dom(u)|$; the usual argument about static scoping establishes $|dom(u)| \leq N$. A completely formal proof would have to consider the details of the syntax of Pure PreScheme and the compiler. ∎

## 5.3. Design Choices in the Compiler

The translation from Pure PreScheme to combinator code was patterned after [6]. In this section, we discuss some of the innovations and design choices that we made in the process of developing our translator.

An important change from [6] is the explicit separation of the syntax and semantics of the combinator code. This separation was already present in [24], but was dropped in [25, 26] in favor of punning between $\lambda$-terms and their denotations. At the time, the pun was useful in deriving the compiler-oriented semantics, but many people found it confusing. Furthermore, a clear separation between syntax trees and domain elements allowed us to understand when continuity arguments were permissible, and when they were not. This was even more important once rational trees became part of the structure.

We slightly modified the architecture of the machine. Clinger used a value register $\epsilon$ and a value stack $\epsilon^*$ where we use the value stack $\zeta$. Clinger's value register serves to cache the top element of the stack. As a result, his induction hypotheses are asymmetrical because they need to keep track of when the cache is active: note that in Clinger's induction hypotheses A and B, $\epsilon$ occurs free on the left but not on the right, and in his induction hypothesis C, it appears free, but one has a precondition that says its value is ignored [6, page 360]. He also needs to generate *push* instructions to copy the value register onto the top of the stack. By replacing these with a single value stack, our induction hypotheses became more symmetrical and easier to work with; we did not need any pesky side conditions such as the one for Clinger's hypothesis C. We could rely on a later register allocation phase, such as that in [17], to deal with caching strategies.

Another change was in the separation of the different compiler modes. Clinger's compiler used two different compiling valuations, $t$ (for general expressions) and $e$ (for expressions in tail position). The choice of which valuation to use was embedded in the code of the compiler [6, Figure 5]. We chose instead to separate tail-recursive expressions from simple expressions in the syntax. This made it clearer which valuation should be used at any point in the compiler.

Another distinction in the compiler was between simple expressions evaluated for value or for effect. It is advantageous for the compiler to distinguish between these cases so that it need not worry about handling the value of an expression evaluated for effect. There were a number of different ways of organizing the compiler to treat this distinction. We could have made a separate syntactic category for commands. This would be the obvious approach for an ordinary procedural language. However, in our language, the set of commands would have been exactly the same as the

set of simple expressions, so introducing a new syntactic category would
have been redundant. We could have had commands leave their values on
the top of the stack, and introduced an instruction that would drop the top
value of the stack when it was to be ignored. This would have led to unac-
ceptably inefficient code. In the end, we simply used a separate compiler
valuation for expressions in command position, and we introduced a few
extra instructions, such as $update\text{-}store_i$, to perform primitive operations
without pushing a value on the stack. This is also consistent with modern
ideas about machine architecture, which condone and even encourage new
instructions to handle common special cases.

A more subtle change was in the role of the specifications. In [6], induc-
tion hypotheses were introduced only as part of the internals of the proof.
Indeed, Clinger writes, "... the algorithm was designed and a compiler
built before any thought was given to a formal correctness proof." [6, page
356]. Instead, in keeping with the ideas of [27], we viewed the induction
hypotheses as *specifications* for the compiler: they described what the code
generated by the compiler was supposed to do. This idea, along with our
simplified machine design, allowed us to write down the specifications for
the different valuations in a relatively systematic way. Given the specifi-
cations and a few general ideas about the form of the code, it was easy
to write the clauses for the various productions. Thus the design of the
specifications preceded the writing of the compiler. Most (though not all)
of the proofs were easy, as was shown in Section 4.6. We believe that the
induction hypotheses were a major work product, and should be considered
a major work product of future compiler verification projects.

## 6.    The Linear-Data Machine

The linear-data machine differs from the combinator machine in that it
uses a linearly-addressed data store to represent the environment and stack
of the combinator machine. It uses two pointers into the store, $sp$ and
$up$, to represent the stack and the environment. It also has a heap which
corresponds to the heap of the combinator machine. Thus the state of
the linear-data machine is of the form $\langle q, up, sp, s, h \rangle$, where $sp$ and $up$ are
integers, and $s$ is a map from integers to tree structured values.

In this section we will formalize the way in which the $sp$ and $up$ pointers
represent the stack and the environment. To aid in the intuition behind
these formalisms, we will first give an informal sketch of the representation.
We will then present the formal definitions.

The representation utilizes the fact that any PreScheme program uses
only a bounded, statically-determinable amount of space in the store. By
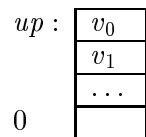analyzing the program, we may define two constants for that program.

The first, $N$, is the total amount of environment space required by the program. The second, $N_0$, is the amount of space initially required by the environment for globals and procedure declarations. These constants are easily calculated (see Theorem 7) and are used for setting the initial parameters of the linear-data machine.

## 6.1. Informal Presentation

We can take advantage of this by using locations 0 through $N$ of the linear-data machine's store for the environment $u$, and the locations above $N$ for the local stack $z$.

We can describe the representation pictorially as follows:

1. Stackable values (integers, booleans, characters, strings, procedures, environments). In keeping with the assumptions of the runtime model discussed in Section 2, runtime tags are not used. This restricts our correctness result but only for those programs that result in an error. All quantities are represented by a single machine word (i.e. one location in $s$). Strings and procedures are represented by pointers into a static space; since all procedures in a PreScheme program share the same environment, there is no need for a separate environment pointer in a procedure object. Environments need to be stacked, but they are not expressible; when stacked they are represented as pointers to the environment representation. Since all other data is either immediate or a pointer to static space, these environment pointers are the only real "pointers" in the system (that is, they are the only references that can potentially dangle).

2. Runtime environment $u$. As all values are represented by a single word in memory, this suggests that $u$ be represented by locations 0 through $up$ in $s$:

$$up: \begin{array}{|c|} \hline v_0 \\ \hline v_1 \\ \hline \dots \\ \hline \phantom{v_0} \\ \hline \end{array}$$

$0$

where $v_0 = \mathtt{lookup}\, u\, 0$, $v_1 = \mathtt{lookup}\, u\, 1$, etc. By Theorem 7, $0 \leq up \leq N$.

3. The local stack $z$. Since all values are represented by a single word in $s$ and that $up$ ranges between 0 and $N$, a local stack $z = (v_1 :: v_2 :: \dots :: v_n :: \langle\rangle)$ can be represented using $s$ as a stack growing upward from position $N + 1$:

$$sp: \begin{array}{|c|} \hline v_1 \\ \hline v_2 \\ \hline \dots \\ \hline v_n \\ \hline \\ \hline \end{array}$$
$$N$$

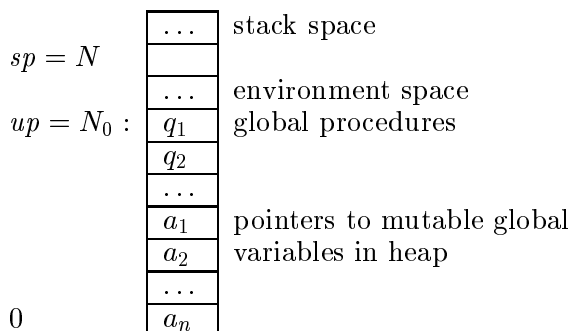In the initial startup of the machine, $up$ is started at $-1$ and the stack pointer $sp$ starts at $N$, so that both the environment and stack are empty. The prelude to the program allocates globals with a sequence of $add\text{-}to\text{-}env_g$ instructions, followed by a *closerecs* instruction. The effect of these instructions is to create an environment $u_0$ in which all the global procedures are closed. At this point $up$ has value $N_0$. This is the ground configuration which is used as a reference for all tail calls.

After the globals and procedures have been allocated, $up$ will have value $N_0$ and will not drop below this value for the remainder of program execution. This configuration is the configuration restored at tail call:

$$
\begin{array}{rl}
 & \begin{array}{|c|} \hline \dots \\ \hline \end{array} \quad \text{stack space} \\
sp = N & \begin{array}{|c|} \hline \\ \hline \end{array} \\
 & \begin{array}{|c|} \hline \dots \\ \hline \end{array} \quad \text{environment space} \\
up = N_0: & \begin{array}{|c|} \hline q_1 \\ \hline \end{array} \quad \text{global procedures} \\
 & \begin{array}{|c|} \hline q_2 \\ \hline \end{array} \\
 & \begin{array}{|c|} \hline \dots \\ \hline \end{array} \\
 & \begin{array}{|c|} \hline a_1 \\ \hline \end{array} \quad \text{pointers to mutable global} \\
 & \begin{array}{|c|} \hline a_2 \\ \hline \end{array} \quad \text{variables in heap} \\
 & \begin{array}{|c|} \hline \dots \\ \hline \end{array} \\
0 & \begin{array}{|c|} \hline a_n \\ \hline \end{array}
\end{array}
$$

## 6.2.   Formal Presentation

To formalize this representation, we will introduce the notion of *storage layout relation* [10] which is defined in [8, Section 3.3]. A storage layout relation is a predicate asserting that some concrete values (that is, values in the linear-data machine) represent some abstract term (that is, a term in the language of the combinator machine). A typical example is the relation for environments, which is a judgement of the form $(s, b) \models_U p \simeq u$. This predicate asserts that pointer $p$ into store $s$ represents run-time environment $u$. (Here the parameter $b$ marks the upper boundary of environment space). Such a relation is typically defined by a least fixed-point induction on a powerset $P(X \times Y)$ for suitable sets $X$ and $Y$.

We will have one storage layout relation for each kind of data manipulated by the combinator and linear-data machines. These are: immediate data,

stackable values, stacks, environments, heaps, and programs. These will be defined as the least fixed point of a simultaneous induction. There is no difficulty in defining these relations on infinite trees; certain infinite trees will be in the relation because one of the base cases in the induction (the one for recursive environments) will require it. The definition of the relations is a relatively straightforward transcription of the data in the diagrams above; the primary difficulty is in guarding against off-by-one errors.

### 6.2.1. Relating Immediate Data

If $i$ and $i'$ are immediate data, then $i \simeq_{imm} i'$ (read "$i$ represents $i'$") iff $i$ is the untagged version of $i'$. This is a degenerate storage-layout relation, since it does not involve the store.

### 6.2.2. Relating Stackable Values

$$
\begin{aligned}
(s, b) \models_V c \simeq v \iff & \\
\text{either} \quad & v = \langle \mathbf{proc}, \langle q, u_0 \rangle \rangle \text{ and } c = q \\
\text{or} \quad & v = \langle \mathbf{env}, u \rangle \\
& \text{and } (s, b) \models_U c \simeq u \\
& \text{and } c \leq b \\
\text{or} \quad & c \simeq_{imm} v
\end{aligned}
$$

Here $u_0$ denotes the environment tree at the end of the prelude. This rule is noteworthy for two reasons: First, in this representation scheme only procedures closed in this environment are representable. However, all procedures in a PreScheme program are closed in this environment, so this restriction is compatible with the semantics of PreScheme. Second, $u_0$ may in general be an infinite rational tree, so any solution to this equation, including its least fixed point, will have infinite trees in the relation. We could have made $u_0$ a parameter to the definition of $\models$, as we did with $s$ and $b$, but that would have made our notation even more cumbersome with little gain in clarity.

The parameter $b$ marks the upper boundary of environment space; the condition on environments ensures that no environment pointer is ever dangling. This parameter will almost always be *up*.

We use $c$ to range over concrete values (the values in the linear-data machine), and $v$ to range over rational trees as in Section 5.

### 6.2.3.   Relating Pointers and Stacks

$$(s, b) \models_Z \langle p, p' \rangle \simeq z \iff$$
$$\text{either} \quad z = \langle \rangle \text{ and } p = p'$$
$$\text{or} \qquad z = (v :: z') \text{ and } p > p' \text{ and } (s, b) \models_V s(p) \simeq v$$
$$\text{and } (s, b) \models_Z \langle p - 1, p' \rangle \simeq z'$$

Here the stack elements are in locations $p' + 1$ through $p$, so $p = p'$ marks an empty stack. In the proof, $p'$ will always be $N$.

### 6.2.4.   Relating Pointers and Environments

$$(s, b) \models_U p \simeq u \iff$$
$$\text{either} \quad u = emptydisplay \text{ and } p = -1$$
$$\text{or} \qquad u = (extends_{rl}(v_1 \ldots v_n)u') \text{ and } p \leq up$$
$$\text{and } (s, b) \models_V s(p) \simeq v_n \ldots (s, b) \models_V s(p - (n - 1)) \simeq v_1$$
$$\text{and } (s, b) \models_U (p - n) \simeq u'$$
$$\text{or} \qquad u = (extends_{rg} au') \text{ and } p \leq up$$
$$\text{and } s(p) \simeq a$$
$$\text{and } (s, b) \models_U (p - 1) \simeq u'$$

### 6.2.5.   Relating Programs

Programs in the linear-data and combinator machines differ only in the format of the literals embedded in them. Thus we define

$$q_l \simeq q_c \iff$$
$$\text{either} \quad q_c = (constant\ i'q_c') \text{ and } q_l = (constant\ iq_l')$$
$$\text{and } i \simeq i' \text{ and } q_l' \simeq q_c'$$
$$\text{or} \qquad q_c = (fetch_l\ iq_c') \text{ and } q_l = (fetch_l\ iq_l')$$
$$\text{and } q_l' \simeq q_c'$$
$$\text{or} \qquad q_c = (fetch_g\ iq_c') \text{ and } q_l = (fetch_g\ iq_l')$$
$$\text{and } q_l' \simeq q_c'$$
$$\text{or} \qquad q_c = (goto\ q_c') \text{ and } q_l = (goto\ q_l')$$
$$\text{and } q_l' \simeq q_c'$$
$$\text{etc.}$$

### 6.2.6.   Relating Heaps

The linear-data heap is like the combinator-machine heap except that its L-values are *untagged* integers, and its first component contains immediate values in their linear-data representation. Hence the first component maps untagged integers to immediate values in their linear-data representation, the second component is an integer-valued free-storage counter, and the

third component is the same as that for the combinator machine. Heaps correspond if their free-storage counters and third components agree, and if the allocated portions of their first components correspond:

$$h_l \simeq h_c \iff$$
$$h_c.2 = h_l.2$$
$$\text{and } (\forall i : 0 \leq i \leq h_l.2)(h_l.1(i) \simeq h_c.1(\langle \mathbf{hptr}, i \rangle))$$
$$\text{and } h_c.3 = h_l.3$$

## 6.3.   Relating the Combinator and Linear-Data Machines

We can now define the correspondence between combinator and linear-data machine states.

The effect of the prelude is to establish the invariants

$$N_0 \leq up \leq N \leq sp$$

and

$$(s, u_0) \models_U N_0 \simeq u_0$$

The first invariant expresses the disjointness of the environment and stack spaces, and the second establishes the correct representation of the globals.

**Definition 6** *We say a linear-data machine state $\langle q, up, sp, s, h \rangle$ corresponds to an combinator machine state $\langle q', u, z, h' \rangle$, written*

$$\langle q, up, sp, s, h \rangle \simeq \langle q', u, z, h' \rangle$$

*if and only if the following conditions are satisfied:*

*1. $q \simeq q'$,*

*2. $(s, up) \models_U up \simeq u$ ,*

*3. $(s, up) \models_Z (sp, N) \simeq z$ ,*

*4. $h \simeq h'$*

*5. $N_0 \leq up \leq N \leq sp$, and*

*6. $(s, u_0) \models_U N_0 \simeq u_0$*

### 6.4.   Basic Properties of the Representation

Before proceeding, we state several lemmas that express the basic properties of these relations. These are used constantly in the proof.

**Lemma 3 (Immediate Values)** *If $c \simeq_{imm} v$, then for all $s$ and $b$, $(s, b) \models_V c \simeq v$ .*

**Definition 7** *We say $s =_{[p,p']} s'$ if $p \leq p'$ and $(\forall x : p \leq x \leq p')(s(x) = s'(x))$.*

**Lemma 4 (Free-Storage Lemma)**    1. *If $(s, b) \models_V c \simeq v$  and $s =_{[0,b]} s'$ then $(s', b) \models_V c \simeq v$ .*

    2. *If $(s, b) \models_Z \langle p, p' \rangle \simeq z$  and $s =_{[0,b]} s'$ and $s =_{[p'+1,p]} s'$ then $(s', b) \models_Z \langle p, p' \rangle \simeq z$ .*

    3. *If $(s, b) \models_U p \simeq u$  and $s =_{[0,p]} s'$ then $(s', b) \models_U p \simeq u$ .*

    **Proof:**   By induction on the definitions of $\models_V, \models_Z,$ and  $\models_U$. ∎

**Lemma 5** *If $(s, b) \models_V c \simeq v$  and $b \leq b'$ then $(s, b') \models_V c \simeq v$ .*

    **Proof:**    The parameter $b$ is only used in the case $(s, b) \models_V c \simeq \langle \mathbf{env}, u \rangle$ . In that case, $c \leq b \leq b'$, so $(s, b') \models_V c \simeq \langle \mathbf{env}, u \rangle$ . The rest of the proof follows by tedious induction. ∎

**Lemma 6** *If $(s, b) \models_U p \simeq u$  then $(s, b) \models_V s(p - n) \simeq \mathtt{lookup}\ u\ n$ .*

### 6.5.   Definition of the Linear-Data Machine

We define the operational semantics of the linear-data machine by considering the representation of each possible combinator machine state, so that the linear-data machine simulates the behavior of the combinator machine. That is, if $L_1 \simeq C_1$ and $C_1 \Longrightarrow C_2$, then we want the linear-data machine to send $L_1$ to some state $L_2$ such that $L_2 \simeq C_2$.

In general, we will only be able to perform this simulation when $C_1$ is a state arising from reducing a properly compiled Pure PreScheme program. There are a number of things that might go wrong:

- The linear-data machine does not have tags to distinguish different data types. The combinator machine uses these tags to distinguish an

error state. Since the linear-data machine cannot mimic this behavior, we will require the linear-data machine to behave correctly only if $A_2$ is not an error state. If $A_2$ is an error state, the behavior of the linear-data machine is unspecified. For type-checked PreScheme programs, such as the Scheme 48 Virtual Machine, this restriction is moot, because the typing rules guarantee that the program will never reach an error state.

- Similarly, if $A_1$ is an combinator-machine state with some arbitrary program $q$, executing it may cause other errors, such as underflowing the stack or overflowing the environment. However, this will not happen if $A_1$ is a state that is reached when the abstract machine is started with a correctly compiled Pure PreScheme program. The simulation theorem will have to take this into account.

The operational semantics of the linear-data machine is shown in Figure 15. The various auxiliaries used by the machine are shown in Figure 16 `apply-prim`$'$ is unspecified, but it must satisfy the constraint that when `apply-prim` and `apply-prim`$'$ are applied to congruent arguments, they return congruent immediate values and congruent heaps.

## 6.6. Correctness of the Linear-Data Machine

The combinator and linear-data machines are related by a property like the following:

> Let $C_1$ and $C_2$ be states of the combinator machine and $L_1$ and $L_2$ be states of the linear-data machine. If $L_1 \simeq C_1$, $C_1 \Longrightarrow C_2$, and $L_1 \Longrightarrow L_2$, then $L_2 \simeq C_2$.

As suggested above, this property does not hold in general. However this result does hold where $C_1$ is a state that arises when the abstract machine is started with a correctly compiled Pure PreScheme program.

Let `conc` be the operation that converts an abstract machine program to a concrete machine program by removing the tags on the operands of all the *constant* instructions, so for all abstract-machine programs $q$, $\mathtt{conc}(q) \simeq q$. Now we can state the simulation theorem:

**Theorem 8 (Correctness of the Linear-Data Machine)** *If* P *is a Pure PreScheme program, and* $\mathcal{P}[\![\mathrm{P}]\!] = \langle \mathbf{ok}, \langle \mathbf{int}, n \rangle \rangle$*, then there are values of* $N_0$ *and* $N$ *such that*

$$\langle \mathtt{conc}(\mathtt{cp}\,[\![\mathrm{P}]\!]), -1, N, s_0, h_0 \rangle \stackrel{*}{\Longrightarrow} \langle \mathbf{ok}, n \rangle$$

$\langle (halt), up, sp, s, h \rangle \qquad\qquad \Longrightarrow \quad \langle \mathbf{ok}, s(sp) \rangle$

$\langle (goto\ q), up, sp, s, h \rangle \qquad\qquad \Longrightarrow \quad \langle q, up, sp, s, h \rangle$

$\langle (label\ q), up, sp, s, h \rangle \qquad\qquad \Longrightarrow \quad \langle q, up, sp, s, h \rangle$

$\langle (fetch_l\ i\ q), up, sp, s, h \rangle \qquad\quad \Longrightarrow \quad \langle q, up, sp+1, s', h \rangle$
   where $s' = s[(\mathtt{lookup}\ i\ up\ s)/sp+1]$

$\langle (brf\ q_1 q_2), up, sp, s, h \rangle \qquad\quad \Longrightarrow \quad \langle (\mathtt{truish}(s(sp)) \rightarrow q_1, q_2), up, sp-1, s, h \rangle$

$\langle (save\text{-}env\ q), up, sp, s, h \rangle \qquad\quad \Longrightarrow \quad \langle q, up, sp+1, s[up/sp+1], h \rangle$

$\langle (restore\text{-}env\ q), up, sp, s, h \rangle \qquad \Longrightarrow \quad \langle q, s(sp-1), sp-1, s_1, h \rangle$
   where $s_1 = s[s(sp)/sp-1]$

$\langle (update\text{-}store\ iq), up, sp, s, h \rangle \quad \Longrightarrow \quad \langle q, up, sp, s, h_1 \rangle$
   where $h_1 = (\mathtt{update}\ s(up-i)\ s(sp)\ h)$

$\langle (prim\text{-}apply\ n\ p\ q), up, sp, s, h \rangle \quad \Longrightarrow \quad \langle q, up, sp-(n-1), s_1, h_1 \rangle$
   where $(w_1, h_1) = (\mathtt{apply\text{-}prim}'\ p\ (\mathtt{collect}\ s\ sp\ n)\ h)$
   and $s_1 = s[w_1/\ sp-(n-1)]$

$\langle (add\text{-}to\text{-}env\ n\ q), up, sp, s, h \rangle \quad \Longrightarrow \quad \langle q, (up+n), (sp-n), s_1, h \rangle$
   where $s_1 = (\mathtt{copy}\ s\ n\ (up+1)\ sp)$

$\langle (add\text{-}to\text{-}env_g\ q), up, sp, s, h \rangle \qquad \Longrightarrow$
   $\langle q, up+1, sp-1, (s[(\mathtt{new}\ h)/(up+1)]), (\mathtt{update}\ (\mathtt{new}\ h)\ s(sp)\ h) \rangle$

$\langle (closerecs\ (q_1 \ldots q_n)\ q), up, sp, s, h \rangle$
$\qquad\qquad\qquad\qquad\qquad \Longrightarrow \quad \langle q, up+n, sp, s_1, h \rangle$
   where $s_1 = (\mathtt{spread}\ s(up+n)\ (q_1 \ldots q_n))$

$\langle (tail\text{-}call), up, sp, s, h \rangle \qquad\quad \Longrightarrow \quad \langle s(sp), (N_0 + n), (N-1), s_1, h \rangle$
   where $n = sp-N-1$
   and $s_1 = (\mathtt{copy}\ s\ n\ (N_0 + 1)\ (sp-1))$

otherwise $\qquad\qquad\qquad\qquad\qquad \Longrightarrow \quad \langle \mathbf{error} \rangle$

Figure 15: Operational Semantics of the Linear-Data Machine

$$(\texttt{new } h) = h.2 + 1$$

$$(\texttt{update } l \, v \, h) = \langle h.1[v/l], l > h.2 \rightarrow l, h.2, h.3 \rangle$$

$$(\texttt{lookup } i \, up \, s) = s(up - i)$$

$$(\texttt{copy } s \, n \, d \, s\,) = (n = 0) \rightarrow s, (\texttt{copy } s[s(s)/d] \, (n - 1) \, (d + 1) \, (s - 1))$$

$$(\texttt{spread } s \, n \, \langle \rangle) = s$$
$$(\texttt{spread } s \, n \, v :: (v_0 \ldots v_n)) = (\texttt{spread } s[v/n] \, (n - 1) \, (v_0 \ldots v_n))$$

$$(\texttt{collect } s \, i \, 0) = \langle \rangle$$
$$(\texttt{collect } s \, i \, n) = s(i) :: (\texttt{collect } s \, (i - 1) \, (n - 1))$$

Figure 16: Auxiliaries for the Linear-Data Machine

**Proof:** Choose $N_0$ and $N$ as in Theorem 7. Let $C_k$ and $L_k$ denote the states of the combinator and linear-data machines, respectively, after $k$ steps. We will show that for all sufficiently large $k$, either $L_k \simeq C_k$ or both machines halt with corresponding answers.

We observe that the programs for both the linear-data and combinator machines begin with a sequence of *add-to-env*$_g$ instructions, followed by a *closerecs* instruction. Let $k_0$ be the number of steps to execute these instructions.

The machines start in states

$$\langle q, -1, N, s_0, h_0 \rangle$$

for the linear-data machine and

$$\langle q', emptydisplay, \langle \rangle, h_0' \rangle$$

for the combinator machine. Comparing these states with the definition of $\simeq$ for states, we see that all of the conditions for correspondence hold, except for $N_0 \leq up$ and $(s, N) \models_U N_0 \simeq u_0$ Let $\simeq_0$ denote this weakened relation. Since there are no procedures represented during the preamble, this relation does not depend on $u_0$.

We then show the following:

> For all $k < k_0$, $L_k \simeq_0 C_k$.
> For all $k \geq k_0$, $L_k \simeq C_k$.

The programs $q$ and $q'$ begin with a sequence of *add-to-env*$_g$ instructions. A straightforward calculation, using the definitions of the stor-

age layout relations and Lemma 4, shows that these instructions preserve $\simeq_0$. This establishes the first proposition.

Next, we turn to the *closerecs* instruction that ends the prelude. Assume we have

$$\langle (closerecs\,(q_1 \ldots q_n)\,q), up, sp, s, h \rangle$$
$$\simeq_0 \langle (closerecs\,(q'_1 \ldots q'_n)\,q'), u, z, h' \rangle$$

We must show that

$$\langle q, up +n, sp, s', h \rangle \simeq \langle q', u_0, z, h' \rangle$$

where

$$s' = (\texttt{spread}\ s(up +n)\,(q_1 \ldots q_n))$$
$$u_0 = (\texttt{fix}(\lambda u'.\,(extends_{rl}\ u\,(\langle \mathbf{proc}, \langle q'_1, u' \rangle \rangle \ldots \langle \mathbf{proc}, \langle q'_n, u' \rangle \rangle))))$$

Note that here we need to establish $\simeq$, not $\simeq_0$.

At the end of this instruction, we have $up = N_0$, by the definition of $N_0$.

By the definition of $s'$, for $1 \le i \le n$ we have $s'(up +n - i + 1) = q_i$, so by the definition of $\models_V$ we have $(s', up) \models_V s'(up +n - i + 1) \simeq \langle \mathbf{proc}, \langle q'_i, u_0 \rangle \rangle$ , and $(s', up) \models_U up \simeq u$ , so by the definition of $\models_U$ we have $(s', up) \models_U up +n \simeq u_1$ , where

$$u_1 = (extends_{rl}\ u\,(\langle \mathbf{proc}, \langle q'_1, u_0 \rangle \rangle \ldots \langle \mathbf{proc}, \langle q'_n, u_0 \rangle \rangle))$$

But $u_1$ is just one unwinding of $u_0$, so $u_1 = u_0$.

All the other conditions are established trivially. This establishes $L_{k_0} \simeq C_{k_0}$.

We next turn to the case of $k_0 < k$ We must show that if $k_0 \le k$ and $L_k \simeq C_k$, then $L_{k+1} \simeq C_{k+1}$. We do this by a tedious analysis of each instruction, using the definitions of the storage layout relations and Lemma 4.

■

## 6.7.  Development of Storage-Layout Relations

Storage-layout relations, as we have described them, allow us to structure a proof by induction on computation steps. Many early compiler-correctness proofs used congruences between domains as a way of structuring such proofs, e.g. [23]. Such proofs were complicated because the relations were formulated between elements of domains, and inverse-limit constructions acted as a surrogate for operational reasoning.

Instead, relying on the development of natural semantics [20], our predicates relate pairs of trees, thus avoiding most of these complications. Our definitions are patterned on those of Hannan [10]. In [28], we introduced these ideas, showing how the relations of [10] could be generalized and the proofs simplified.

## 7. The Stored-Program Machine

The stored-program machine is identical to the linear-data machine, except that the code is represented in cells of a linear instruction store $M$, much like the cells of the data store.

We use storage-layout relations to specify the correspondence between the two machines. We first show how linear-data machine code is represented in the linear program store. For example the representation of bytecode $(fetch_l \ i \ q_0)$ at location $p$ in $M$ is:

$$p : \begin{array}{|c|} \hline fetch_l \\ \hline i \\ \hline \phantom{i} \\ \hline \phantom{i} \\ \hline \end{array}$$

where $p + 2$ represents $q_0$. This notion is formalized by the storage-layout relation $\models_P$:

$$
\begin{aligned}
M \models_P p \simeq q \iff \\
\text{either} \quad & M(p) = halt \text{ and } q = (halt) \\
\text{or} \quad & M(p) = goto \text{ and } q = (goto \ q_0) \\
& \text{and } M \models_P M(p+1) \simeq q_0 \\
\text{or} \quad & q = (label \ q_0) \\
& \text{and } M \models_P p \simeq q_0 \\
\text{or} \quad & M(p) = fetch_l \text{ and } q = (fetch_l \ i \ q_0) \\
& \text{and } M(p+1) = i \text{ and } M \models_P p+2 \simeq q_0 \\
\text{or} \quad & M(p) = brf \text{ and } q = (brf \ q_0 \ q_1) \\
& \text{and } M \models_P p+2 \simeq q_0 \text{ and } M \models_P M(p+1) \simeq q_1 \\
\text{etc.} \quad &
\end{aligned}
$$

This change propagates through the other classes of data that are represented in the machine. Values correspond if they are alike except for the representation of programs. Data stores correspond up to pointer $p$ if they correspond location by location for all smaller locations. Heaps correspond if they both have the same number of used cells and the values of corresponding cells are related. All these can be written easily using the language of storage-layout relations in the style of Section 6.

Given these relations, we can define the correspondence between states of the linear-data machine and fetch states of the stored-program machine. For correspondence, the stored-program machine's instruction pointer must correspond to the linear-data machine's program. Furthermore the dedicated registers must have the same values in both machines and data in stores and heaps must correspond.

**Definition 8** *We say that under linear code store $M$, stored-program machine state $F(ip, up, sp, s, h)$ corresponds to a linear-data machine state $\langle q, up', sp', s', h' \rangle$ (written $M \models F(ip, up, sp, s, h) \simeq \langle q, up', sp', s', h' \rangle$) if and only if:*

1. $M \models_P ip \simeq q$

2. $up = up'$

3. $sp = sp'$

4. $M \models_{sp} s \simeq s'$

5. $M \models_H h \simeq h'$

The formal definition of the stored-program machine appears in figure 17. The program resides in the linear program store $M$, which is indexed by register $ip$. The machine uses a fetch-execute cycle. In a fetch state $F(ip, up, sp, s, h)$, the machine retrieves the contents of the location pointed to by $ip$ and goes into the execute state $E(ip, M(ip), up, sp, s, h)$. In the execute state, the machine decodes the current instruction and behaves accordingly, going into a fetch state. The machine's goal is to preserve correspondence with the behavior of the linear-data machine. The program store is treated as a global since it remains constant throughout execution.

The simulation theorem says that machines in corresponding states remain in corresponding states as they execute. Therefore, upon termination they will return corresponding values.

**Theorem 9** *Let $L_1$ and $L_2$ be states of the linear-data machine and $P_1$ and $P_2$ be fetch states of the stored-program machine. If $M \models P_1 \simeq L_1$, $L_1 \Longrightarrow L_2$, and $P_1 \Longrightarrow P_2$, then $M \models P_2 \simeq L_2$.*

**Proof:**
By analysis of each instruction. ∎

## 8. The Linearizer and its Correctness

In Theorem 9, we showed that if $L$ and $P$ were corresponding states of the linear-data machine and the stored-program machine, then $L$ and $P$ would always compute in corresponding states. In particular, if $L$ halted with an integer $n$, then so would $P$. Theorem 8 specified how to get from a Pure PreScheme program P to a suitable linear-data machine state by using `cp`

Fetch state :
$$F(ip, up, sp, s, h) \qquad\qquad \Longrightarrow \quad E(ip, M(ip), up, sp, s, h)$$

Execute state :
$$E(ip, halt, up, sp, s, h) \qquad\qquad \Longrightarrow \quad \langle \mathbf{ok}, s(sp) \rangle$$
$$E(ip, goto, up, sp, s, h) \qquad\qquad \Longrightarrow \quad F(M(ip+1), up, sp, s, h)$$
$$E(ip, fetch_l, up, sp, s, h) \qquad\qquad \Longrightarrow \quad F(ip+2, up, sp+1, s_1, h)$$
$$\text{where } s_1 = s[(\mathtt{lookup}\ M(ip+1)\ up\ s)/(sp+1)]$$
$$E(ip, brf, up, sp, s, h) \qquad\qquad \Longrightarrow \quad F(ip_1, up, sp-1, s, h)$$
$$\text{where } ip_1 = (\mathtt{truish}(s(sp)) \to (ip+2), M(ip+1))$$
$$E(ip, save\text{-}env, up, sp, s, h) \qquad \Longrightarrow \quad F(ip+1, up, sp+1, s[up/sp+1], h)$$
$$E(ip, restore\text{-}env, up, sp, s, h) \qquad \Longrightarrow \quad F(ip+1, s(sp-1), sp-1, s_1, h)$$
$$\text{where } s_1 = s[s(sp)/sp-1]$$
$$E(ip, update\text{-}store, up, sp, s, h) \quad \Longrightarrow \quad F(ip+2, up, sp, s, h_1)$$
$$\text{where } h_1 = (\mathtt{update}\ s(up-M(ip+1))\ s(sp)\ h)$$
$$E(ip, prim\text{-}apply, up, sp, s, h) \qquad \Longrightarrow \quad F(ip+3, up, sp_1, s_1, h_1)$$
$$\text{where } sp_1 = sp-(M(ip+1)-1)$$
$$\text{and } (w_1, h_1) = (\mathtt{apply\text{-}prim}\ M(ip+2)\ (\mathtt{collect}\ s\ sp\ M(ip+1))\ h)$$
$$\text{and } s_1 = s[w_1/sp-(M(ip+1)-1)]$$
$$E(ip, update\text{-}store_i, up, sp, s, h) \quad \Longrightarrow \quad F(ip+2, up, sp-1, s, h_1)$$
$$\text{where } h_1 = (\mathtt{update}\ s(up-M(ip+1))\ s(sp)\ h)$$
$$E(ip, prim\text{-}apply_i, up, sp, s, h) \qquad \Longrightarrow \quad F(ip+3, up, sp-M(ip+1), s, h_1)$$
$$\text{where } (w_1, h_1) = (\mathtt{apply\text{-}prim}\ M(ip+2)\ (\mathtt{collect}\ s\ sp\ M(ip+1))\ h)$$
$$E(ip, add\text{-}to\text{-}env, up, sp, s, h) \qquad \Longrightarrow \quad F(ip+2, up_1, sp_1, s_1, h)$$
$$\text{where } up_1 = (up+M(ip+1))$$
$$\text{and } sp_1 = (sp-M(ip+1))$$
$$\text{and } s_1 = (\mathtt{copy}\ s\ M(ip+1)\ (up+1)\ sp)$$
$$E(ip, add\text{-}to\text{-}env_g, up, sp, s, h) \qquad \Longrightarrow \quad F(ip+1, up+1, sp-1, s_1, h_1)$$
$$\text{where } s_1 = (s[(\mathtt{new}\ h)/(up+1)])$$
$$\text{and } h_1 = (\mathtt{update}\ (\mathtt{new}\ h)\ s(sp)\ h)$$
$$E(ip, closerecs, up, sp, s, h) \qquad\qquad \Longrightarrow \quad F(M(ip+1), up+M(ip+2), sp, s_1, h)$$
$$\text{where } s_1 = (\mathtt{spread}\ s(up+M(ip+2))\ (ip+3 \ldots ip+M(ip+2)+2))$$
$$E(ip, tail\text{-}call, up, sp, s, h) \qquad\qquad \Longrightarrow \quad F(s(sp), N_0+n, (N-1), s_1, h)$$
$$\text{where } n = sp-N-1$$
$$\text{and } s_1 = (\mathtt{copy}\ s\ M(ip+1)\ (N_0+1)\ (sp-1))$$
$$\text{otherwise} \qquad\qquad\qquad \Longrightarrow \quad \mathbf{error}$$

Figure 17: The fetch/execute cycle of the stored-program machine.

and `conc`, but Theorem 9 gives no corresponding specification for how to get to a stored-program machine state.

In this section, we describe an algorithm that takes the combinator code for a program $q$ and an initial code address $p$, and produces a linear instruction store $M$ such that $M \models_P p \simeq q$, that is, $M$ contains a program, starting at $p$, that corresponds to $q$.

We call this algorithm the *linearizer*. This would be a straightforward tree traversal, except that splits and joins must be accounted for. Splits occur at conditionals. For example, consider the combinator code emitted for a simple conditional expression. The syntax-directed compiler gives:

$$\texttt{cse}\,[\![(\texttt{if}\ S_0\ S_1\ S_2)]\!]gq = \texttt{cse}\,[\![S_0]\!]g(\mathit{brf}\,(\texttt{cse}\,[\![S_1]\!]g(\mathit{label}\ q))$$
$$(\texttt{cse}\,[\![S_2]\!]g(\mathit{goto}\ q)))$$

The linearizer must avoid making two copies of $q$. It does this by relying on the syntax-directed compiler to mark the join point with *label* and *goto* instructions. More precisely, the linearizer assumes that the left path into a join is marked by a *label* and the right path is marked by a *goto*. This requirement is formalized by the predicate $\mathit{legal}(q)$ in Figure 18, which also defines the join point $\mathit{join}(q)$ of a combinator program $q$. It is easy to confirm that code generated by the compiler satisfies these restrictions:

**Theorem 10** *For any program* P*,* $\texttt{conc}(\texttt{cp}\,[\![P]\!])$ *is legal.*

      **Proof:**  Routine induction. ∎

The linearizer is a set of three valuations on combinator-code programs. We present it in a style similar to that employed for the compiler. For each valuation, we present an informal description of its intended purpose, its formal specification, and some excerpts from its definition. As in Section 4, formulating the specification was the most difficult part of the process; once that was done, writing the linearizer and its proof was relatively straightforward. We therefore give the specifications in some detail, and the algorithm and proof in less detail.

We begin with some preliminary definitions needed for the specification.

**Definition 9** *For code stores* $M$ *and* $M'$*,* $M =^{p'}_p M'$ *if and only if* $(\forall i : p \le i < p')\ M(i) = M'(i)$.

**Definition 10** *The instructions tail-call and halt are called* leaves *and* leaf $(q)$ *is true if and only if* $q$ *is a leaf.*

| $q$ | $join(q)$ | $legal(q)$ |
|---|---|---|
| $(halt)$ | $(halt)$ | $true$ |
| $(goto\ q_0)$ | $(goto\ q_0)$ | $legal(q_0)$ |
| $(label\ q_0)$ | $(label\ q_0)$ | $legal(q_0)$ |
| $(fetch_l\ i\ q_0)$ | $join(q_0)$ | $legal(q_0)$ |
| $(brf\ q_0\ q_1)$ | $join(join(q_0))$ | $join(q_0) = (label\ q')$ |
| | | and $join(q_1) = (goto\ q')$ |
| | | and $legal(q')$ |
| $(save\text{-}env\ q_0)$ | $join(q_0)$ | $legal(q_0)$ |
| $(restore\text{-}env\ q_0)$ | $join(q_0)$ | $legal(q_0)$ |
| $(update\text{-}store\ i\ q_0)$ | $join(q_0)$ | $legal(q_0)$ |
| $(prim\text{-}apply\ n\ p\ q_0)$ | $join(q_0)$ | $legal(q_0)$ |
| $(update\text{-}store_i\ i\ q_0)$ | $join(q_0)$ | $legal(q_0)$ |
| $(prim\text{-}apply_i\ n\ p\ q_0)$ | $join(q_0)$ | $legal(q_0)$ |
| $(add\text{-}to\text{-}env\ nu\ q_0)$ | $join(q_0)$ | $legal(q_0)$ |
| $(add\text{-}to\text{-}env_g\ q_0)$ | $join(q_0)$ | $legal(q_0)$ |
| $(closerecs\ (q_1 \ldots q_n)\ q_0)$ | $join(q_0)$ | $(\forall i : 0 \leq i \leq n)$ |
| | | $legal(q_i)$ and $leaf\ (join(q_i))$ |
| $(tail\text{-}call\ n)$ | $(tail\text{-}call\ n)$ | $true$ |

Figure 18: Specification of legal combinator code.

We are now ready to describe the linearizer.

$\mathcal{L}$ is the initial linearization function. It takes as input a combinator-code program $q$ representing an entire program or procedure, a code store $M$, and a pointer $p$ to the next free position in $M$. It returns two values $\langle M', p' \rangle$, where $M'$ is a code store just like $M$ except that $M' \models_P p \simeq q$, and $p'$ is the next free position in $M'$.

This informal specification must be refined to clarify the meaning of phrases like "just like" and "next free position." This must be done in a way that will support the eventual induction proof. After considerable experimentation, we arrived at the following formal specification of $\mathcal{L}$:

If $q$ is legal and the join point of $q$ is a leaf then $\mathcal{L}[\![q]\!]Mp = \langle M', p' \rangle$ such that

1. $M =_0^p M'$

2. $p' \geq p$

3. $(\forall\ M'')$ if $M'' =_p^{p'} M'$, then $M \models_P p \simeq q$

The first condition guarantees that the linearizer does not change the locations in $M$ below $p$, corresponding to the informal statement that $p$ is the first free location in $M$ and $M'$ is "just like" $M$. The third condition guarantees that the linearization of $q$ lies between locations $p$ and $p' - 1$.

$\mathcal{L}$ is defined by:

$$\mathcal{L}[\![(halt)]\!]Mp = \langle M[halt/p], p + 1 \rangle$$

$$\mathcal{L}[\![(tail\text{-}call\ n)]\!]Mp = \langle M[tail\text{-}call/p, n/p + 1], p + 2 \rangle$$

otherwise
$$\mathcal{L}[\![q]\!]Mp = \texttt{let}\ \langle M', p', q' \rangle = \mathcal{L}'[\![q]\!]M(p + 2)$$
$$\qquad\qquad \texttt{in}\ \mathcal{L}[\![q']\!]M'p'$$

The idiom $M[tail\text{-}call/p, n/p + 1]$ corresponds to the emission of code: Here location $p$ is loaded with $tail\text{-}call$ and location $p + 1$ is loaded with the number $n$.

$\mathcal{L}'$ does the major work. It takes as input a combinator-code program $q$, a code store $M$, and a code pointer $p$ (the next free position in $M$). It returns three values $\langle M', p', q' \rangle$, where $M'$ is a modified code store, $p'$ is a pointer into the code store, and $q'$ is a combinator-code program.

The formal specification for $\mathcal{L}'$ is: If $q$ is legal then $\mathcal{L}'[\![q]\!]Mp = \langle M', p', q' \rangle$ such that

1. $M =_0^p M'$

2. $p' \geq p$

3. $q' = join(q)$

4. $(\forall M'')$ if $M'' =_p^{p'} M'$ and $M'' \models_p p' \simeq q'$, $M'' \models_P p \simeq q$.

The first and second conditions are similar to those for $\mathcal{L}$. The third condition states that $\mathcal{L}'$ computes $join(q)$. The fourth condition says that if $M''$ has a linearization of $q'$ beginning at $p'$, and is identical with $M'$ for positions between $p$ and $p'$, then $M''$ has a linearization of $q$ starting at $p$. So $p'$ is a position in the code store at which the linearizer expects an image of $q'$ to be placed.

When both sides of the conditional are linearized, the linearizer chooses which copy of $q'$ to linearize (they are guaranteed to be the same because $q$ is legal), and it inserts a goto at the place in the code store where the other branch wanted to put $q$. This is illustrated in the clause for linearizing a $brf$:

$$\mathcal{L}'[\![(brf\, q_0\, q_1)]\!]Mp$$
$$= \texttt{let}^* \; \langle M', p', (goto\, q') \rangle = \mathcal{L}'[\![q_0]\!]M(p+2)$$
$$\langle M'', p'', (label\, q') \rangle = \mathcal{L}'[\![q_1]\!]M'(p'+2)$$
$$\texttt{in}\; \mathcal{L}'[\![q']\!]M''[brf/p, (p'+2)/(p+1), goto\,/p', p''/(p'+1)]p''$$

Here $q'$ represents the common code after $q_0$ and $q_1$ rejoin each other. Though $q'$ is bound twice in this definition, the two bindings are guaranteed to coincide because of the assumption that $(brf\, q_0\, q_1)$ is legal. If $\mathcal{L}'[\![(brf\, q_0\, q_1)]\!]Mp = \langle M''', p''', q'' \rangle$, then $M'''$ will look like

| | |
|---|---|
| $p:$ | $brf$ |
| | $p'+2$ |
| | $\dots$   linearization of $q_0$ |
| $p':$ | $goto$ |
| | $p''$ |
| | $\dots$   linearization of $q_1$ |
| $p'':$ | $\dots$   linearization of $q'$ |
| $p''':$ | starting address for $q''$ |

We give some additional excerpts from the definition of $\mathcal{L}'$ to suggest the other cases.

$$\mathcal{L}'[\![(halt)]\!]Mp = \langle M, p, (halt) \rangle$$

$\mathcal{L}'[\![(goto\ q_0)]\!]Mp = \langle M, p, (goto\ q_0)\rangle$

$\mathcal{L}'[\![(label\ q_0)]\!]Mp = \langle M, p, (label\ q_0)\rangle$

$\mathcal{L}'[\![(fetch_l\ i\ q_0)]\!]Mp = \mathcal{L}'[\![q_0]\!]M[fetch_l\ /p, i/p+1](p+2)$

$\mathcal{L}'[\![(closerecs\ (q_1 \ldots q_n)\ q_0)]\!]Mp =$
   $\mathtt{let}\ \langle M', (p_1 \ldots p_n), p'\rangle = \mathcal{L}^*[\![(q_1 \ldots q_n)]\!]M(p+n+3)$
     $\mathtt{in}\ \mathcal{L}'[\![q_0]\!]r'[closerecs\ /p, p'/p+1, n/p+2,$
               $p_1/p+3, \ldots, p_n/(p+n+2)]p'$

$\mathcal{L}'[\![(tail\text{-}call\ n)]\!]Mp = \langle M, p, (tail\text{-}call\ n)\rangle$

The last valuation, $\mathcal{L}^*$, is used for linearizing a sequence of combinator-code programs. It takes as input a list of combinator-code programs $(q_1 \ldots q_n)$, a code store $M$, and a code pointer $p$ (the next free position in $M$) and returns three values $\langle M', (p_1 \ldots p_n), p'\rangle$. $M'$ is a code store just like $M$ except that $(\forall i : 1 \leq i \leq n)\ M' \models_P p_i \simeq q_i$. Again, $p'$ is the next free position in $M'$. $\mathcal{L}^*$ is used to linearize the procedure definitions in a *closerecs*.

The formal specification of $\mathcal{L}^*$ is: if $(\forall i : 1 \leq i \leq n)\ q_i$ is legal and the join point of $q_i$ is a leaf then $\mathcal{L}^*[\![(q_1 \ldots q_n)]\!]Mp = \langle M', (p_1 \ldots p_n), p'\rangle$ such that

1. $M =_0^p M'$

2. $p' \geq p$

3. $(\forall\ M : M =_p^{p'} M'),\ (\forall i : 1 \leq i \leq n)\ M \models_P p_i \simeq q_i$

$\mathcal{L}^*$ is defined by:

$\mathcal{L}^*[\![\langle\rangle]\!]Mp = \langle M, \langle\rangle, p\rangle$

$\mathcal{L}^*[\![q :: (q_1 \ldots q_n)]\!]Mp = \mathtt{let}^*\ \langle M', p'\rangle = \mathcal{L}[\![q]\!]Mp$
                                 $\langle r'', (p_1 \ldots p_n), p''\rangle = \mathcal{L}^*[\![(q_1 \ldots q_n)]\!]M'p'$
                     $\mathtt{in}\ \langle r'', p :: (p_1 \ldots p_n), p''\rangle$

The correctness of the linearizer is given as follows:

**Theorem 11** $\mathcal{L}$, $\mathcal{L}'$ and $\mathcal{L}^*$ satisfy their specifications.

    **Proof:**  By induction on the definitions of $\mathcal{L}$, $\mathcal{L}'$ and $\mathcal{L}^*$.

    ■

## 9. Implementation

The VLISP PreScheme compiler was implemented as two separate programs. The front end translates VLISP PreScheme into Pure PreScheme and the back end translates Pure PreScheme into code for the Motorola 68000.[3] The front end was implemented using about 4300 lines of code and the back end used 2300 lines. To facilitate bootstrapping as described in [8, Section 8], they were both written in the dialect of Scheme accepted by the VLISP Scheme byte-code compiler.

The back end used the 68000 as a target because it was the architecture most readily available to us at the start of the project. Later, the compiler was retargeted to produce SPARC code with very little effort.

Representing our abstract machines with the 68000 requires a mapping similar to the ones just detailed. This section informally describes that mapping and then describes the code generated for the VLISP Byte-Code Interpreter by the compiler.

Immediate data is modeled via a 32 bit quantity (4 bytes). Thus immediate data can be stored directly in any of the 68000 registers (`D0`–`D7` and `A0`–`A7`) or in 4 consecutive bytes of user memory.

The stack $s$ is modeled by the 68000 with the user stack, and its stack pointer, $sp$, is modeled by the 68000 stack pointer `SP` (address register `A7`). The environment register, $up$, is modeled by the address register `A1`. For efficiency, the initial environment pointer, $N_0$, and initial stack pointer, $N$, are also kept in address registers for quick reset on tail call.

Two details must be noted in using the 68000 user stack as described. First, the 68000 user stack grows downward rather than upward. Furthermore, one unit of immediate data is represented by four bytes of 68000 memory. Therefore the stack pointer manipulations and the linearization algorithm needed to be altered accordingly.

Procedures are represented by their effective addresses. When the initial environment $u_0$ is being created, these effective addresses (which are 32 bit quantities) are pushed on the stack. On tail-call, the effective address is loaded from the stack into the program counter.

The environment representation is almost completely analogous to that of the stored-program machine, with environment values kept in the first $4 * N$ bytes of the stack. The local stack is handled almost completely analogously as well. The only difference is that for efficiency, the top 8 values of the local stack are cached in in the data registers. The correctness of such

[3]The front end and the back end were developed at two different sites. The two programs were not merged solely because we did not take the time to resolve conflicting usages of Scheme program variables.

caching strategies can be performed as part of the linearization phase; see [17] for details. This simple minded register allocation technique was quite effective. Program execution speed increased by 35% when the technique was implemented.

Representation of the stored-program machine's heap is done differently for the different components of the heap. The first component of the heap, the store for mutable variables, is kept directly in the first $4 * j$ bytes of the user stack (where $j$ is the number of global variables in the program being compiled). The second component of the heap, which tells how many objects are stored in the first component, remains constant during the execution of a given Pure PreScheme program after the initial environment has been established (in particular, it is also $j$). It is modeled implicitly through the correctness of the instructions manipulating mutable data. The third component of the heap, that which the primitives manipulate, is simulated by the 68000 heap, primarily via operating system calls.

We used a program that computes Fibonacci numbers to give an initial benchmark for the back end of the compiler on an unloaded Sun 3/60. The Pure PreScheme program uses an explicit array to model the recursion stack. It was compared with an implementation written in C and compiled with gcc version 1.37 using the -O switch. In the C version, tail-recursion is replaced by explicit goto's, and arguments are passed using automatic variables, which gcc allocates in machine registers. The C version also uses `malloc` to allocate the array dynamically, just as the Pure PreScheme version does. We found that that our code used about three times as much CPU time as the C version.

A more realistic test of the VLISP software was constructed as follows. One of the authors wrote a C version of the VLISP Byte-Code Interpreter which was designed only for speed and not structured for verification. It was compiled on a Sun 4 using gcc version 2.1 with the -O2 switch, and the generated assembly code was studied. Any inefficiency observed in the code resulted in a modification to the C source. When the process was completed, all the variables used to describe the state of the byte-code interpreter were placed into machine registers, and the code sequences for the most used byte-code instructions appeared to be optimal.

Tests showed that it took 3.9 times longer to run programs on a Sun 3/60 using the verified PreScheme version of the VLISP Byte-Code Interpreter as compared with the C version compiled with gcc version 1.37 using the -O switch. We consider this a significant accomplishment given the extensive and intricate optimization techniques used by the GNU C compiler.

## 10. Conclusions

We have presented a verified compiler for PreScheme. Our architecture divided the compiler into three components:

1. A transformational front end that translates source text into a core language called Pure PreScheme.

2. A syntax-directed compiler that translates Pure PreScheme into a combinator-based tree-manipulation language. We call this language *combinator code.*

3. An assembler that translates combinator code into code for an abstract stored-program machine with linear memory for both data and code.

This factorization proved to be a successful decomposition. It allowed us to use both denotational and operational reasoning to their best advantage. Furthermore, each portion of the proof gave us confidence in a different part of the implementation.

The proof of correctness of the run-time structure is operational, rather than denotational, in its structure: that is, it proceeds by induction on the number of steps taken by the machine, rather than on the size of the program. We believe this is a fundamental improvement. Previous proofs of compiler correctness, such as those in [13, 23], were highly complex because they used induction on the construction of reflexive domains as a denotational analog of induction on length of computation. Our proof would likely be just as complex had we adopted this strategy (see [29] for an attempt to do a much smaller problem in a purely denotational style). By using induction on computation length directly, we avoid this indirection. Furthermore, since we deal directly with terms (trees) rather than with their denotations, we avoid the complication of "inclusive predicates," with their attendant complexity.

The implementation experience showed that having a validated compiler and run-time structure eliminated most bugs in the areas covered by the proofs. The various delivered versions of the compilers had a number of bugs, but these were almost entirely in one of two categories:

1. Incompatibilities between the various phases of the compiler, including errors in syntax, etc. These were artifacts of the circumstance that the various portions of the compiler, including the interface between the front and back ends, were developed concurrently.

2. Problems with the assembly code sequences generated for the concrete-machine instructions. These were mostly minor in nature (registers not being saved across routine calls, etc.), and were below the grain of the proof. Extending the proof to reach this level would require an extremely detailed model of the behavior of the machine and operating system (see, e.g. [16, 5]). In practice, however, the concrete machine was at a sufficiently low level that implementation of the primitives was easy.

The method of formalizing storage layout relations seems to be flexible enough to model standard representation strategies. More of these are presented in [28].

The proofs of the different stages had rather different styles. The proofs of the transformational front end (Section 3) were most like "traditional" semantics proofs, such as those in [13]. The proofs of the syntax-directed translator (Section 4.6) tended to be more like exercises in $\lambda$-calculus. We believe that a fairly simple simplification-based theorem prover could automate most of the cases in this induction. The back-end proofs (Sections 6.6–8) resemble traditional Hoare-style verification proofs, but seem quite stylized and may be amenable to mechanization by systems like the Boyer-Moore theorem prover [4]. Most of the individual proof cases were routine calculations. We believe this is a mark of success for this effort.

Use of a verified compiler need not entail catastrophic performance penalties. On a realistic example, we showed that a verified program compiled by our verified compiler used only about four times more CPU time compared with an unverified version compiled by the optimizing gcc.

We believe that this proof architecture is applicable to other languages as well. There are, however, a number of ways in which this might be done. One might use a modification of the transformational front end to translate the source language into a variant of Pure PreScheme. One could then modify the Pure PreScheme compiler appropriately. Another approach would be to modify the syntax-directed compiler of Sections 4–5 to use the denotational semantics of the source language instead of the semantics of Section 2. In this approach, the front end would be replaced by a source-to-source optimizer.

A key issue in either of these approaches is the incorporation of optimization into the compiler. In real compilers, optimization is performed on many levels. From the point of view of a compiler proof, the job of an optimizer is to derive some non-local invariants that allow the program to be transformed safely. For example, we needed to rely on the invariance of $N_0 \leq |dom(u)| \leq N$, which is not possible to verify locally without additional information. The operational type-soundness of the program is

another such invariant. We believe that developing a systematic way of incorporating such information is the most important outstanding problem in semantics-based compiler correctness.

# References

1. IEEE Std 1178-1990. *IEEE Standard for the Scheme Programming Language.* Institute of Electrical and Electronic Engineers, Inc., New York, NY (1991).

2. Appel, A. W. *Compiling with Continuations.* Cambridge University Press, Cambridge and New York (1992).

3. Barendregt, Henk P. *The Lambda Calculus: Its Syntax and Semantics.* North-Holland, Amsterdam (1981).

4. Boyer, R. and Moore, J. *A Computational Logic.* Academic Press (1979).

5. Boyer, R. S. and Yu, Y. Automated correctness proofs of machine code programs for a commercial microprocessor. In Kapur, D., editor, *Automated Deduction — CADE-11*, 11th International Conference on Automated Deduction, Springer Verlag (1992) 416–430.

6. Clinger, Will. The Scheme 311 compiler: An exercise in denotational semantics. In *1984 ACM Symposium on Lisp and Functional Programming*, The Association for Computing Machinery, Inc., New York (August 1984) 356–364.

7. Courcelle, Bruno. Fundamental properties of infinite trees. *Theoretical Computer Science*, 25 (1983) 95–169.

8. Guttman, J. D. and Swarup, V. VLISP: A verified implementation of scheme. Submitted to LASC.

9. Guttman, Joshua D., Wand, Mitchell, and Ramsdell, John D. Results and conclusions from the vlisp project. *Lisp and Symbolic Computation* (??).

10. Hannan, John. Making abstract machines less abstract. In Hughes, J., editor, *Functional Programming Languages and Computer Architecture, 5th ACM Conference*, Springer-Verlag, Berlin, Heidelberg, and New York (1991) 618–635.

11. Kelsey, Richard A. Realistic compilation by program transformation. In *Conf. Rec. 16th Ann. ACM Symp. on Principles of Programming Languages*, ACM (1989).

12. Kranz, D., Kelsey, R., Rees, J., Hudak, P., Philbin, J., and Adams, N. Orbit: An optimizing compiler for Scheme. In *SIGPLAN Notices* (1986) 219–233. Proceedings of the '86 Symposium on Compiler Construction.

13. Milne, R. and Strachey, Christopher. *A Theory of Programming Language Semantics*. Chapman and Hall, London (1976). Also Wiley, New York.

14. Milner, Robin. A theory of type polymorphism in programming. *Journal of Computer and System Sciences*, 17 (1978) 348–375.

15. Milner, Robin, Tofte, Mads, and Harper, Robert. *The Definition of Standard ML*. The MIT Press, Cambridge, MA (1990).

16. Moore, J S. *Piton: A Verified Assembly-Level Language*. Technical Report 22, Computational Logic, Inc., Austin, Texas (1988).

17. Oliva, Dino P. *Advice on Structuring Compiler Back Ends and Proving Them Correct*. PhD thesis, Northeastern University (1993).

18. Oliva, Dino P. and Wand, Mitchell. *A Verified Run-Time Structure for Pure PreScheme*. Technical Report NU-CCS-92-27, Northeastern University College of Computer Science (September 1992).

19. Paulson, Laurence C. *ML for the Working Programmer*. Cambridge University Press, Cambridge, Great Britain (1991).

20. Plotkin, Gordon D. *A Structural Approach to Operational Semantics*. Technical Report DAIMI FN-19, Aarhus University (1981).

21. Ramsdell, J. D., Farmer, W. M., Guttman, J. D., Monk, L. G., and Swarup, V. *The VLISP PreScheme Front End*. M 92B098, The MITRE Corporation (September 1992).

22. Steele, Guy L. *Rabbit: A Compiler for Scheme*. Technical Report 474, MIT AI Laboratory (1978).

23. Stoy, Joseph E. *Denotational Semantics: The Scott-Strachey Approach to Programming Language Theory*. MIT Press, Cambridge, MA (1977).

24. Wand, Mitchell. Deriving target code as a representation of continuation semantics. *ACM Transactions on Programming Languages and Systems*, 4, 3 (July 1982) 496–517.

25. Wand, Mitchell. Semantics-directed machine architecture. In *Conf. Rec. 9th ACM Symposium on Principles of Programming Languages* (1982) 234–241.

26. Wand, Mitchell. Loops in combinator-based compilers. *Information and Control*, 57, 2–3 (May/June 1983) 148–164.

27. Wand, Mitchell. Correctness of procedure representations in higher-order assembly language. In Brookes, S., editor, *Proceedings Mathematical Foundations of Programming Semantics '91*, Springer-Verlag, Berlin, Heidelberg, and New York (1992) 294–311.

28. Wand, M. and Oliva, D. P. Proving the correctness of storage representations. In *Proceedings of the 1992 ACM Conference on Lisp and Functional Programming*, ACM Press, New York (1992) 151–160.

29. Wand, Mitchell and Wang, Zheng-Yu. Conditional lambda-theories and the verification of static properties of programs. In *Proc. 5th IEEE Symposium on Logic in Computer Science* (1990) 321–332.