

15-411 Compiler Design: Lab 4

Fall 2009

Instructor: Frank Pfenning
TAs: Ruy Ley-Wild and Miguel Silva

Test Programs Due: 11:59pm, Thursday, October 22, 2009

Compilers Due: 11:59pm, Thursday, October 29, 2009

Revision 1: Oct 25, 2009 (see last page)

1 Introduction

The goal of the lab is to implement a complete compiler for the language L_4 . This language extends L_3 by structs (in the sense of C), pointers, and arrays. This means the main change from the third lab is that you will have to deal with *memory references*. Because pointers on the x86-64 architecture are 64 bits wide, this also means you will have to deal with differently sized data. Again, performance of both the compiler and the compiled code are only minor considerations at this stage. However, the threshold for compilation time will be slightly tighter, so you may need to start to pay attention to some performance issues.

2 Requirements

As for Labs 1, 2, and 3, you are required to hand in test programs as well as a complete working compiler that translates L_4 source programs into correct target programs written in x86-64 assembly language. When encountering an error in the input program (which can be a lexical, grammatical, or static semantics error) the compiler should terminate with a non-zero exit code and print a helpful error message. To test the target programs, we will assemble and link them using `gcc` on the lab machines and run them under fixed but generous time limits.

3 L_4 Syntax

The syntax of L_4 is defined by the context-free grammar in Figure 1. This is an extension of the grammar for L_3 in the sense that a correct L_3 program should still parse correctly if it does not use any of the new keywords. Ambiguities in this grammar are resolved according to the operator precedence table in Figure 2 and the rule that an `else` provides the alternative for the most recent eligible `if`. Whitespace and comments are as in L_3 .

$\langle \text{program} \rangle ::= \langle \text{gdecl} \rangle^* \langle \text{function} \rangle^*$
 $\langle \text{gdecl} \rangle ::= \mathbf{extern} \langle \text{type} \rangle \langle \text{ident} \rangle (\langle \text{paramlist} \rangle) ; |$
 $\quad \mathbf{struct} \langle \text{ident} \rangle \{ [\langle \text{ident} \rangle : \langle \text{type} \rangle ;]^* \} ;$
 $\langle \text{function} \rangle ::= \langle \text{type} \rangle \langle \text{ident} \rangle (\langle \text{paramlist} \rangle) \langle \text{body} \rangle$
 $\langle \text{paramlist} \rangle ::= \varepsilon | \langle \text{ident} \rangle : \langle \text{type} \rangle [, \langle \text{ident} \rangle : \langle \text{type} \rangle]^*$
 $\langle \text{body} \rangle ::= \{ \langle \text{decl} \rangle^* \langle \text{stmt} \rangle^* \}$
 $\langle \text{decl} \rangle ::= \mathbf{var} \langle \text{ident} \rangle [, \langle \text{ident} \rangle]^* : \langle \text{type} \rangle ;$
 $\langle \text{type} \rangle ::= \mathbf{int} | \mathbf{bool} | \langle \text{ident} \rangle | \langle \text{type} \rangle * | \langle \text{type} \rangle []$
 $\langle \text{stmt} \rangle ::= \langle \text{simp} \rangle ; | \langle \text{control} \rangle | ;$
 $\langle \text{simp} \rangle ::= \langle \text{exp} \rangle \langle \text{asop} \rangle \langle \text{exp} \rangle | \langle \text{exp} \rangle$
 $\langle \text{control} \rangle ::= \mathbf{if} (\langle \text{exp} \rangle) \langle \text{block} \rangle [\mathbf{else} \langle \text{block} \rangle] |$
 $\quad \mathbf{while} (\langle \text{exp} \rangle) \langle \text{block} \rangle | \mathbf{for} ([\langle \text{simp} \rangle] ; \langle \text{exp} \rangle ; [\langle \text{simp} \rangle]) \langle \text{block} \rangle |$
 $\quad \mathbf{continue} ; | \mathbf{break} ; | \mathbf{return} \langle \text{exp} \rangle ;$
 $\langle \text{block} \rangle ::= \langle \text{stmt} \rangle | \{ \langle \text{stmt} \rangle^* \}$
 $\langle \text{exp} \rangle ::= (\langle \text{exp} \rangle) | \langle \text{intconst} \rangle | \mathbf{true} | \mathbf{false} | \mathbf{NULL} |$
 $\quad \langle \text{ident} \rangle | * \langle \text{exp} \rangle | \langle \text{exp} \rangle . \langle \text{ident} \rangle | \langle \text{exp} \rangle [\langle \text{exp} \rangle] | \langle \text{exp} \rangle \rightarrow \langle \text{ident} \rangle |$
 $\quad \langle \text{unop} \rangle \langle \text{exp} \rangle | \langle \text{exp} \rangle \langle \text{binop} \rangle \langle \text{exp} \rangle | \langle \text{exp} \rangle ? \langle \text{exp} \rangle : \langle \text{exp} \rangle |$
 $\quad \langle \text{ident} \rangle ([\langle \text{exp} \rangle [, \langle \text{exp} \rangle]^*]) |$
 $\quad \mathbf{new} (\langle \text{type} \rangle [[\langle \text{exp} \rangle]])$
 $\langle \text{ident} \rangle ::= [\mathbf{A-Z_a-z}][\mathbf{0-9A-Z_a-z}]^*$
 $\langle \text{intconst} \rangle ::= [\mathbf{0-9}][\mathbf{0-9}]^* \quad (\text{in the range } 0 \leq \text{intconst} < 2^{32})$
 $\langle \text{asop} \rangle ::= = | += | -= | *= | /= | \% = | \& = | \wedge = | | = | \ll = | \gg =$
 $\langle \text{binop} \rangle ::= + | - | * | / | \% | < | \leq | > | \geq | == | != |$
 $\quad \&\& | || | \& | \wedge | | | \ll | \gg$
 $\langle \text{unop} \rangle ::= ! | \sim | -$

The precedence of unary and binary operators is given in Figure 2. Non-terminals are in $\langle \text{angle brackets} \rangle$, optional constituents in $[\text{brackets}]$. Terminals are in **bold**.

Figure 1: Grammar of L_4

Operator	Associates	Meaning
() [] -> .	left	parens, subscript, field dereference, field select
! ~ - *	right	logical not, bitwise not, unary minus, dereference
* / %	left	integer times, divide, modulo
+ -	left	plus, minus
<< >>	left	(arithmetic) shift left, right
< <= > >=	left	comparison
== !=	left	equality, disequality (integer or pointer)
&	left	bitwise and
^	left	bitwise exclusive or
	left	bitwise or
&&	left	logical and
	left	logical or
? :	right	conditional expression
= += -= *= /= %=		
&= ^= = <<= >>=	right	assignment operators

Figure 2: Precedence of operators, from highest to lowest

4 L_4 Runtime Objects

The L_4 language has four kinds of identifiers: those standing for functions (g), those standing for variables (x), those standing for structs (s), and those standing for field names (f). These are in separate name spaces, so a function, a variable, a struct, and a field may have the same name without conflict. Reserved words of the grammar (**extern**, **struct**, **var**, **int**, **bool**, **if**, **else**, **while**, **for**, **continue**, **break**, **return**, **true**, **false**, **NULL**, **new**) cannot be used as any kind of name.

The L_4 language has some complexity in its new constructs, dealing with memory (structs, pointers, and arrays). In order to describe these concisely we use an abstract syntax notation. Eliding some expressions which are as in L_3 :

$$\begin{array}{ll} \text{Types} & \tau ::= \mathbf{int} \mid \mathbf{bool} \mid s \mid \tau^* \mid \tau[] \\ \text{Expressions} & e ::= c \mid \mathbf{null} \mid x \mid *e \mid e.f \mid e[e] \mid e \rightarrow e \mid \mathbf{new} \tau \mid \mathbf{new} \tau[e] \mid \dots \end{array}$$

We also have to consider the runtime objects created during the execution of a program. We divide these into two classes: small values and large values. Small values can be held in registers, large values must be in memory. We use a for an address of an object in memory and n for an integer, and b for a boolean.

$$\begin{array}{ll} \text{Small Values} & w ::= a \mid n \mid b \\ \text{Large Values} & W ::= \{f_1=V_1; \dots f_n=V_n; \} \mid [V_0, \dots, V_{n-1}] \\ \text{Values} & V ::= w \mid W \end{array}$$

The first kind of large value is a struct with fields f_1, \dots, f_n , the second an array with n elements. On occasion, we need to know the exact size of objects of a given type, $|\tau|$. We compute this as follows from the type:

$$\begin{array}{ll} |\mathbf{int}| & = 4 \\ |\mathbf{bool}| & = 4 \\ |\tau^*| & = 8 \\ |\tau[]| & = 8 \\ |s| & = \|f_1:\tau_1; \dots f_n:\tau_n; \| \end{array}$$

The size of a struct, written here as $\|f_1:\tau_1; \dots f_n:\tau_n; \|$ is computed by traversing the fields from left to right, adding padding as necessary so that alignment restrictions on the fields are satisfied. **int**'s and **bool**'s are aligned at $0 \bmod 4$, small values of type τ^* and $\tau[]$ are aligned at $0 \bmod 8$, and structs are aligned according to their most strictly aligned field. Padding may need to be added at the end of a struct so that its total size is a multiple of its most strictly aligned field.

Arrays $[V_0, \dots, V_{n-1}]$ cannot be directly embedded in a struct or other array, since their size is unknown at compile time. This is the reason that a value of type $\tau[]$ is *small*: it is merely the starting address of the array in memory. Arrays must be aligned according to the requirement of their elements. The start of an array is returned by a call to the **malloc** or **calloc** library function which guarantees an address $0 \bmod 8$, which is the strictest alignment required in our language.

Addresses are represented as 8 byte (unsigned) integers in the machine. They are obtained from the runtime system which allocates new memory, or by addition (address arithmetic). There is a special address 0 which is never returned by the runtime system and will be used to denote the null pointer. Memory access is represented as $M[a]$, which reads from memory when used as a value

and writes to memory when used on the left-hand side of an assignment $M[a] \leftarrow w$. The number of bytes read or written depends on the size of the small value w .

Along similar lines, we write $V[x]$ for the value of a variable x . On the left hand side $V[x] \leftarrow w$ assigns to x ; on the right hand side $V[x]$ denotes the current value of x .

The most important judgments we make are typing, written $e : \tau$, and evaluation, written $e \Rightarrow w$. The first is part of the *static semantics*, which must satisfy some additional conditions that are stated informally. The latter constitutes the *dynamic semantics*, but is also only partially formal since evaluating expressions has effects that we only describe informally. As a general convention, expressions must be evaluated from left to right so that any side effects happen in a deterministic order.

We now present both the static and dynamic semantics for each new construct (booleans, structs, pointers, arrays) in turn.

Booleans

In *L3*, there was no separate type of booleans—integers were used directly as booleans. In *L4*, we distinguish booleans from integers in the source language, although their representation in the target is entirely straightforward, so not much of your compiler code will have to change. We have two constructors, `true` and `false`, and two destructors, conditional statements and expressions, for booleans.

$$\frac{}{\mathbf{true} : \mathbf{bool}} \qquad \frac{}{\mathbf{false} : \mathbf{bool}}$$

$$\frac{e_1 : \mathbf{bool} \quad s_2 \text{ valid} \quad s_3 \text{ valid}}{\mathbf{if } e_1 \mathbf{ then } s_2 \mathbf{ else } s_3 \text{ valid}} \qquad \frac{e_1 : \mathbf{bool} \quad e_2 : \tau \quad e_3 : \tau}{e_1 ? e_2 : e_3}$$

The exit condition in a while- or for-loop also must be a boolean. We also have a number of operators related to booleans.

$$\frac{e_1 : \mathbf{int} \quad e_2 : \mathbf{int} \quad \text{relop} \in \{<, \leq, >, \geq\}}{e_1 \text{ relop } e_2 : \mathbf{bool}} \qquad \frac{e_1 : \mathbf{bool} \quad e_2 : \mathbf{bool} \quad \text{logop} \in \{\&\&, ||\}}{e_1 \text{ logop } e_2 : \mathbf{bool}} \qquad \frac{e : \mathbf{bool}}{! e : \mathbf{bool}}$$

As for equality and disequality, please see the section on *Overloaded Operators* below.

Booleans should be represented canonically by 1 for `true` and 0 for `false` for interoperability with C (and the compiler for *L3*).

Structs

Statically, several properties must be checked for struct declarations.

- Every type identifier s used in a struct declaration or function must be declared with an explicit `struct s { . . . }`; somewhere in the file. The order of `struct` declarations is irrelevant, but the syntax restricts the file so that they can be mixed with `extern` declarations but must precede all function definitions.
- All struct declarations in a file must have distinct names s .
- Field names in a struct declaration must all be distinct. However, different struct declarations may reuse field names.

- Fields of structs can refer to other structs which can in turn refer to other structs and so on. However, any circular reference from a struct to itself following such a chain must go through a pointer or array type.

The main expression construct for structs is $e.f$. It is easily type-checked.

$$\frac{e : s \quad \mathbf{struct} \ s \ \{ \dots f : \tau ; \dots \}}{e.f : \tau}$$

We write $\text{offset}(s, f)$ for the offset of field f in structure s , in bytes. We suggest to compute this information early and store it in a table so that the compiler can access it easily.

The evaluation rule is slightly tricky, and carries a recurring theme. In several circumstances evaluating an expression yields an address of some object in memory. When this object is small, we can return it directly as a small value for further computation. When the object is large, we cannot directly operate on the object, but must return its address instead.

$$\begin{aligned} e.f &\Rightarrow M[a + k] && \text{if } e \Rightarrow a \text{ and } \text{offset}(s, f) = k \\ &&& \text{where } e : s \text{ and } \mathbf{struct} \ s \ \{ \dots f : \tau ; \dots \} \text{ and } \tau \text{ small} \\ &\Rightarrow a + k && \text{as above, except that } \tau \text{ large} \end{aligned}$$

We also have the derived form $e \rightarrow f$, which can be desugared into $(*e).f$.

Pointers

The typing for pointer operations are straightforward.

$$\frac{e : \tau*}{*e : \tau} \qquad \frac{}{\mathbf{new} \ \tau : \tau*} \qquad \frac{}{\mathbf{null} : \tau*}$$

Unfortunately, the typing of the **null** pointer present a challenge, since it potentially has infinitely many types. This is necessary so that expressions such as $p == \mathbf{NULL}$ or statements such as $p \rightarrow f = \mathbf{NULL}$ where the left-hand side has some type $\tau*$ are all well-typed.

We suggest to create a new type **any** used only internally. In key places during type-checking where some types are required to be equal, we can allow **any** on either side or both sides. In programming language terms we say that **null** is *polymorphic*; fortunately it is the only polymorphic construct in the language so you can still implement a relatively simple type checker. The main places where type comparisons take place is in compiling equality ($==$), disequality ($!=$), assignments, and function calls. To avoid pathologies in type-checking, we disallow $*\mathbf{null}$ because its type is ambiguous and may not be resolved by context. You can still reliably dereference the null pointer, should you be so inclined, for example with `var null : int*; null = NULL; *null;`

Assume a function g has prototype

$$\tau \ g(x_1 : \tau_1, \dots, x_n : \tau_n);$$

Then we have

$$\frac{e_i : \tau_i \quad (1 \leq i \leq n)}{g(e_1, \dots, e_n) : \tau}$$

$$\frac{e : \tau \quad (\mathbf{return} \text{ in body of } g)}{\mathbf{return} \ e \ \mathit{valid}}$$

where s *valid* indicates that the statement s is valid.

The operational semantics for $*e$ evaluates e to an address and then returns the memory contents at that address. Dereferencing the null pointer must raise the **SIGSEGV** exception. In an implementation this can be accomplished without any checks, because the operating system will prevent read access to address 0 and raise the appropriate exception.

null	\Rightarrow	0	
$*e$	\Rightarrow	$M[a]$	if $e \Rightarrow a$ with $a \neq 0$, $e : \tau*$ and τ small
	\Rightarrow	a	as above, except that τ large
		raises SIGSEGV	if $a = 0$
new τ	\Rightarrow	a	where $M[a], \dots, M[a + \tau - 1]$ are freshly allocated locations

The values stored in freshly allocated location must be all 0. This can be achieved with `calloc()` and means that values of type **int** are simply 0, value of type **bool** are **false**, values of type τ' are null pointers, all fields of structs are recursively set to 0, and values of array type have address 0 which is akin to a null array reference.

Arrays

Arrays are similar to pointers.

	$e_1 : \tau[] \quad e_2 : \mathbf{int}$	$e : \mathbf{int}$	
	$e_1[e_2] : \tau$	new $\tau[e] : \tau[]$	
$e_1[e_2]$	\Rightarrow	$M[a + n \tau]$	if $e_1 \Rightarrow a$ and $e_2 \Rightarrow n$ with $e_1 : \tau[], \tau$ small, $a \neq 0$, $M[a + n \tau]$ allocated
	\Rightarrow	$a + n \tau $	as above, except that τ large
		undefined	otherwise
new $\tau[e]$	\Rightarrow	a	if $e \Rightarrow n$ and $M[a], \dots, M[a + (n - 1) \tau]$ are freshly allocated

Values in a freshly allocated array are all initialized to 0, as for pointers. The value of an out-of-bounds array access is *undefined* in this version of the language.

Assignment

Assignment is now more general, because the left-hand side need not be a variable, but could be a more complicated expression so we can write, for example, `*p = *q` to copy the contents of q to the location denoted by p . The left-hand side of an assignment must be a so-called *lvalue* (“left value”), which is a certain kind of expression. We define lvalues, assuming that the expression is already known to be well-typed.

$$\text{Lvalues } v ::= x \mid *e \mid e.f \mid e_1[e_2]$$

For an assignment to be a valid statment, the lvalue on the left-hand side must have a small type.

$$\frac{v : \tau \quad e : \tau \quad \tau \text{ small}}{v = e \text{ valid}}$$

We also have a new statement which is just an expression e which is evaluated for possible effects. e is also required to have small type and its computed value is discarded so that it can be seen as shorthand for $x = e$ where x is a fresh variable not used anywhere else.

An lvalue denotes a *destination* for the assignment, which is either a variable x or an address a . We write $v \mapsto d$ if v denotes destination d . This is defined as follows:

$$\begin{aligned} x &\mapsto x \\ *e &\mapsto a && \text{if } e \Rightarrow a \\ e.f &\mapsto a + k && \text{if } e \Rightarrow a \text{ and } \text{offset}(s, f) = k \text{ for } e : s \\ e_1[e_2] &\mapsto a + n|\tau| && \text{if } e_1 \Rightarrow a \text{ and } e_2 \Rightarrow n \text{ with } e_1 : \tau[] \end{aligned}$$

Now an assignment $v = e$ has the following effect

$$\begin{aligned} v = e &\text{ has the effect of } V[x] \leftarrow w \\ &\text{if } v \mapsto x \text{ and } e \Rightarrow w \\ &\text{has the effect of } M[a] \leftarrow w \\ &\text{if } v \mapsto a \text{ and } e \Rightarrow w \text{ and } M[a] \text{ has been allocated} \\ &\text{has the effect of raising exception SIGSEGV} \\ &\text{if } v \mapsto a \text{ and } e \Rightarrow w \text{ and } a = 0 \\ &\text{is undefined otherwise} \end{aligned}$$

Because we check that v is an lvalue and that the type of e is small, this should cover all possible cases. Note that v must be evaluated to a destination before e is evaluated to a value. In the next lab, we will check for undefined memory access; in this lab, your program need not raise an exception the address is 0. All of the test cases should have well-defined behavior.

The meaning of complex assignment operators now also changes and is no longer a syntactic expansion because expressions can have side effects. An assignment

$$v \text{ op} = e$$

should execute as

$$\begin{aligned} V[x] &\leftarrow V[x] \text{ op } w && \text{if } v \mapsto x \text{ and } e \Rightarrow w \\ M[a] &\leftarrow M[a] \text{ op } w && \text{if } v \mapsto a \text{ and } e \Rightarrow w \end{aligned}$$

Functions

Regarding functions, several properties must be checked.

- Every function that is called must be declared somewhere in the file, either using **extern** or with a definition. The order of the functions in the file is irrelevant, but a function may not be declared more than once.
- Functions must be called with the correct number of arguments of the correct type.
- All function arguments and function results must be of small type.
- Functions must terminate with an explicit **return** statement. This is checked by verifying that each finite control flow path originating at the top of a function ends in a **return** statement. See a more detailed explanation in the description of *L2*.

- There must be a function `main()` which returns an integer as the result of the overall computation (and not an exit status).
- Each `break` or `continue` statement must occur inside a `for` or `while` loop.

Variables

Regarding variables, we need to check three properties, the last of which is new for this language.

- Every variable must be declared, either as a function parameter or an explicit local variable with a `var` declaration. Function parameters and local variables are local to a function and have nothing to do with parameters or local variables in other functions. Variables may not be redeclared, that is, names of parameters to a given function and its local variables must all be pairwise distinct. There are no global variables.
- Each local variable must be defined by an assignment before it is used. This is checked by verifying that along all control flow paths originating at the beginning of the function, each local variable is assigned to before it is used. This ensures that there will be no references to uninitialized variables. A precise specification of this condition is given in Lab 2.
- Any variable, either a function parameter or locally declared variable, must have a small type. Variables of large type are not permitted. This simplifies parameter passing considerably.

The lack of bounds check and this language means that, unlike *L3*, there are programs in *L4* whose result is not defined. Such programs are considered to be “erroneous”, although the compiler must permit them. For your test cases, you may not submit programs whose behavior is undefined.

Overloaded Operators

Equality (`==`) and disequality (`!=`) are overloaded, applying to integers, booleans, and pointers, always returning a boolean. See the discussion above on pointers. Operationally, the compiler will have to choose a 4 byte or 8 byte comparison for the two versions, which suggests after type-checking and during translation to intermediate code there should be two explicitly different operators.

5 Project Requirements

For this project, you are required to hand in test cases and a complete working compiler for *L4* that produces correct target programs written in Intel x86-64 assembly language. When we grade your work, we will use the `gcc` compiler to assemble and link the code you generate into executables using the provided runtime environment on the lab machines.

Test Files

Test files should have extension `.14` and start with one of the following lines

<code>#test return <i>i</i></code>	program must execute correctly and return <i>i</i>
<code>#test exception <i>n</i></code>	program must compile but raise runtime exception <i>n</i>
<code>#test error</code>	program must fail to compile due to an <i>L4</i> source error

followed by the program text. If an exception number is missing, any exception is accepted. Defined exceptions are at least SIGFPE (8), SIGSEGV (11), and SIGALRM (14). These are raised by division by 0 or division overflow (8), dereference of 0, stack overflow, or out-of-memory situations (11), or by a time-out (14). If n is omitted or 0, any exception is valid behavior. That's useful for program that might either run out of memory or time, but it is system-dependent which one is the case.

All test files should be collected into a directory `test/` (containing no other files) and submitted via the Autolab server.

Your test programs must test the new features of $L4$: booleans, structs, pointers, and arrays. Since the language is mostly backward compatible, you may also use the $L3$ test cases for regression testing. We cannot test programs whose results are undefined. It is therefore critical that your test programs do not perform any operations whose meaning is undefined. This includes accessing an array out of bounds or accessing a null array.

While your compiler must support arbitrary extern declarations, the test cases you submit can only refer to the three functions `printchar`, `printint`, `printhex`, and `printbool` provided by the runtime system and any other function explicitly defined in `l4rt.c` (see description below).

We would like some fraction of your test programs to compute “interesting” functions; please briefly describe such examples in a comment in the file. Disallowed are sorting programs.

Compiler Files

The files comprising the compiler itself should be collected in a directory `compiler/` which should contain a `Makefile`. **Important:** You should also update the `README` file and insert a description of your code and algorithms used at the beginning of this file. This will be a crucial guide for the grader.

Issuing the shell command

```
% make l4c
```

should generate the appropriate files so that

```
% bin/l4c <args>
```

will run your $L4$ compiler. The command

```
% make clean
```

should remove all binaries, heaps, and other generated files.

Using the Subversion Repository

Handin and handout of material is via the course subversion repository.

The handout files for this course can be checked out from our subversion repository via

```
% svn checkout https://cvs.concert.cs.cmu.edu/15411-f09/<team>
```

where `<team>` is the name of your team. You will find materials for this lab in the `lab4` subdirectory. Or, if you have checked out `15411-f09/<team>` directory before, you can issue the command `svn update` in that directory.

After first adding (with `svn add` or `svn copy` from a previous lab) and committing your handin directory (with `svn commit`) to the repository you can hand in your tests or compiler by selecting

S3 - Autograde your code in svn repository

from the Autolab server menu. It will perform one of

```
% svn checkout https://cvs.concert.cs.cmu.edu/15411-f09/<team>/lab4/tests
% svn checkout https://cvs.concert.cs.cmu.edu/15411-f09/<team>/lab4/compiler
```

to obtain the files directories to autograde, depending on whether you are handing in your test files or your compiler.

If you are submitting multiple versions, please remember to commit your changes to the repository before asking the Autolab server to grade them! And please do not include an compiled files or binaries in the repository!

What to Turn In

Hand-in on the Autolab server:

- At least 20 test cases, at least two of which generate an error and at least two others raise a runtime exception. The directory `tests/` should only contain your test files and be submitted via subversion or as a tar file as described above. The server will test your test files and notify you if there is a discrepancy between your answer and the outcome of the reference implementation. You may hand in as many times as you like before the deadline without penalty. If you feel the reference implementation is in error, please notify the instructors. The compiled binary for each test case should run in 2 seconds with the reference compiler on the lab machines; we will use a 5 second limit for testing compilers.

Test cases are due **11:59pm on Thu Oct 22, 2009**.

- The complete compiler. The directory `compiler/` should contain only the sources for your compiler and be submitted via subversion or as a tar file as described above. The Autolab server will build your compiler, run it on all extant test files, link the resulting assembly files against our runtime system (if compilation is successful), execute the binaries (each with a 5 second time limit), and finally compare the actual with the expected results. You may hand in as many times as you like before the deadline without penalty.

Compilers are due **11:59pm on Thu Oct 29, 2009**.

6 Notes and Hints

There is no address-of operator (`&`), and variables (be it local variables or procedure parameters) cannot hold large values. The suggested implementation therefore is to allocate structs and arrays on the heap. Our runtime system provides the function `void* calloc(size_t nobj, size_t nbytes)`; for this purpose, which allocates an array of `nobj` objects of size `nbytes` all initialized to 0. Your assembly code should call this function as necessary, assuming `size_t` is `unsigned int` (4 bytes).

There is no corresponding `free`, since we use a conservative garbage collector to reclaim storage.

To check that struct declarations are well-founded, we suggest a simple cycle-checking depth-first search which is in any case necessary to compute field offsets. Global tables to keep struct sizes, field offsets, and the types of functions seem like an appropriate strategy.

Data now can have different sizes, and you need to track this information throughout the various phases of compilation. We suggest you read Section 4 of the Bryant/O’Hallaron note on *x86-64 Machine-Level Programming* available from the Resources page, especially the paragraph on move instructions and the effects of 32 bit operators in the upper 32 of the 64 bit registers.

Your code must strictly adhere to the struct alignment requirements explained above. You may also read the Section 3.1.2 in the *Application Binary Interface* description available from the Resources page.

Also, the C calling conventions detailed in Lab 3 must still be obeyed. Since no large values can be passed, this remains relatively straightforward. Also as in Lab 3, functions declared with `extern` must retain their name, while all defined functions g must be exported as symbols `_14_g`. This prevents any conflicts between a C standard library function and the $L4$ name of a function. When a function `_14_g` is declared with `extern`, there may not be a local function g to avoid ambiguity.

Run-time Environment

The tests will be run in the standard Linux environment on the lab machines; the produced assembly code must conform to those standards. We recommend the use of `gcc -S` to produce assembly files from C sources which can provide template code and assembly language examples.

We will link your code against some library functions you can find in the files `runtime/14rt.h` and `runtime/14rt.c`. This will allow us to test your adherence to the C calling conventions on the x86-64 architecture, and it will allow you to call some standard library functions, say, to print messages, painful though it may be in a language without strings. The functions we explicitly define can be declared in $L4$ as

```
extern int printchar(c:int); /* print c as ASCII character */
extern int printint(n:int); /* print n in decimal */
extern int printhex(n:int); /* print n in hexadecimal */
extern int printbool(b:bool); /* print b as true or false */
```

More may be available as the code is released and updated.

If your compiler detects any (compile-time) errors in the source program, it should exit with a non-zero return code. If compilation succeeds and target code is generated, the compiler should then exit with a return code of 0.

7 Updates

Revision 1. Clarified how to evaluate an assignment $v = e$ by specifying that an lvalue v denotes a destination (to be computed first), the right-hand side e denotes a small value (to be computed second), followed by the assignment (which must raise a `SIGSEGV` exception if the destination is 0). See the description of Assignment on pages 7–8.