

Massachusetts Institute of Technology
Department of Electrical Engineering and Computer Science

6.035, Fall 2006

Handout 5 — Decaf Language

Wednesday, September 6

The project for the course is to write a compiler for a language called Decaf. Decaf is a simple imperative language similar to C or Pascal.

Lexical Considerations

All Decaf keywords are lowercase. Keywords and identifiers are case-sensitive. For example, **if** is a keyword, but **IF** is a variable name; `foo` and `Foo` are two different names referring to two distinct variables.

The reserved words are:

boolean break callout class continue else false for forpar if int return true void

Note that **Program** (see below) is not a keyword, but an identifier with a special meaning in certain circumstances.

Comments are started by `//` and are terminated by the end of the line.

White space may appear between any lexical tokens. White space is defined as one or more spaces, tabs, page and line-breaking characters, and comments.

Keywords and identifiers must be separated by white space, or a token that is neither a keyword nor an identifier. For example, `thisfortrue` is a single identifier, not three distinct keywords. If a sequence begins with an alphabetic character or an underscore, then it, and the longest sequence of characters following it forms a token.

String literals are composed of `<char>`s enclosed in double quotes. A character literal consists of a `<char>` enclosed in single quotes.

Numbers in Decaf are 32 bit signed. That is, decimal values between -2147483648 and 2147483647. If a sequence begins with `0x`, then these first two characters and the longest sequence of characters drawn from `[0-9a-fA-F]` form a hexadecimal integer literal. If a sequence begins with a decimal digit (but not `0x`), then the longest prefix of decimal digits forms a decimal integer literal. Note that range checking is performed later. A long sequence of digits (e.g. 123456789123456789) is still scanned as a single token.

A `<char>` is any printable ASCII character (ASCII values between decimal value 32 and 126, or octal 40 and 176) other than quote (`"`), single quote (`'`), or backslash (`\`), plus the 2-character sequences `"\"` to denote quote, `"\'` to denote single quote, `"\"` to denote backslash, `"\t` to denote a literal tab, or `"\n` to denote newline.

Reference Grammar

Meta-notation:

$\langle \text{foo} \rangle$	means foo is a nonterminal.
foo	(in bold font) means that foo is a terminal i.e., a token or a part of a token.
$[x]$	means zero or one occurrence of x , <i>i.e.</i> , x is optional; note that brackets in quotes ' $[\]$ ' are terminals.
x^*	means zero or more occurrences of x .
$x^+,$	a comma-separated list of one or more x 's.
$\{ \}$	large braces are used for grouping; note that braces in quotes ' $\{ \}$ ' are terminals.
$ $	separates alternatives.

$\langle \text{program} \rangle \rightarrow \text{class Program } \{ \langle \text{field_decl} \rangle^* \langle \text{method_decl} \rangle^* \}$

$\langle \text{field_decl} \rangle \rightarrow \langle \text{type} \rangle \{ \langle \text{id} \rangle \mid \langle \text{id} \rangle \text{ '}' \langle \text{int_literal} \rangle \text{'}' }^+, ;$

$\langle \text{method_decl} \rangle \rightarrow \{ \langle \text{type} \rangle \mid \text{void} \} \langle \text{id} \rangle ([\{ \langle \text{type} \rangle \langle \text{id} \rangle \}^+,]) \langle \text{block} \rangle$

$\langle \text{block} \rangle \rightarrow \{ \langle \text{var_decl} \rangle^* \langle \text{statement} \rangle^* \}$

$\langle \text{var_decl} \rangle \rightarrow \langle \text{type} \rangle \langle \text{id} \rangle^+, ;$

$\langle \text{type} \rangle \rightarrow \text{int} \mid \text{boolean}$

$\langle \text{statement} \rangle \rightarrow \langle \text{location} \rangle = \langle \text{expr} \rangle ;$
 $| \langle \text{method_call} \rangle ;$
 $| \text{if } (\langle \text{expr} \rangle) \langle \text{block} \rangle [\text{else } \langle \text{block} \rangle]$
 $| \text{for } \langle \text{id} \rangle = \langle \text{expr} \rangle , \langle \text{expr} \rangle \langle \text{block} \rangle$
 $| \text{forpar } \langle \text{id} \rangle = \langle \text{expr} \rangle , \langle \text{expr} \rangle \langle \text{block} \rangle$
 $| \text{return } [\langle \text{expr} \rangle] ;$
 $| \text{break } ;$
 $| \text{continue } ;$
 $| \langle \text{block} \rangle$

$\langle \text{method_call} \rangle \rightarrow \langle \text{method_name} \rangle ([\langle \text{expr} \rangle^+,])$
 $| \text{callout } (\langle \text{string_literal} \rangle [, \langle \text{callout_arg} \rangle^+,])$

$\langle \text{method_name} \rangle \rightarrow \langle \text{id} \rangle$

$\langle \text{location} \rangle \rightarrow \langle \text{id} \rangle$
 $| \langle \text{id} \rangle \text{ '}' \langle \text{expr} \rangle \text{'}'$

$\langle \text{expr} \rangle \rightarrow \langle \text{location} \rangle$
 $\quad | \langle \text{method_call} \rangle$
 $\quad | \langle \text{literal} \rangle$
 $\quad | \langle \text{expr} \rangle \langle \text{bin_op} \rangle \langle \text{expr} \rangle$
 $\quad | - \langle \text{expr} \rangle$
 $\quad | ! \langle \text{expr} \rangle$
 $\quad | (\langle \text{expr} \rangle)$

$\langle \text{callout_arg} \rangle \rightarrow \langle \text{expr} \rangle | \langle \text{string_literal} \rangle$

$\langle \text{bin_op} \rangle \rightarrow \langle \text{arith_op} \rangle | \langle \text{rel_op} \rangle | \langle \text{eq_op} \rangle | \langle \text{cond_op} \rangle$

$\langle \text{arith_op} \rangle \rightarrow + | - | * | / | \% | \ll | \gg$

$\langle \text{rel_op} \rangle \rightarrow < | > | <= | >=$

$\langle \text{eq_op} \rangle \rightarrow == | !=$

$\langle \text{cond_op} \rangle \rightarrow \&\& | ||$

$\langle \text{literal} \rangle \rightarrow \langle \text{int_literal} \rangle | \langle \text{char_literal} \rangle | \langle \text{bool_literal} \rangle$

$\langle \text{id} \rangle \rightarrow \langle \text{alpha} \rangle \langle \text{alpha_num} \rangle^*$

$\langle \text{alpha_num} \rangle \rightarrow \langle \text{alpha} \rangle | \langle \text{digit} \rangle$

$\langle \text{alpha} \rangle \rightarrow \mathbf{a} | \mathbf{b} | \dots | \mathbf{z} | \mathbf{A} | \mathbf{B} | \dots | \mathbf{Z} | _ | \cdot$

$\langle \text{digit} \rangle \rightarrow 0 | 1 | 2 | \dots | 9$

$\langle \text{hex_digit} \rangle \rightarrow \langle \text{digit} \rangle | \mathbf{a} | \mathbf{b} | \mathbf{c} | \mathbf{d} | \mathbf{e} | \mathbf{f} | \mathbf{A} | \mathbf{B} | \mathbf{C} | \mathbf{D} | \mathbf{E} | \mathbf{F}$

$\langle \text{int_literal} \rangle \rightarrow \langle \text{decimal_literal} \rangle | \langle \text{hex_literal} \rangle$

$\langle \text{decimal_literal} \rangle \rightarrow \langle \text{digit} \rangle \langle \text{digit} \rangle^*$

$\langle \text{hex_literal} \rangle \rightarrow 0\mathbf{x} \langle \text{hex_digit} \rangle \langle \text{hex_digit} \rangle^*$

$\langle \text{bool_literal} \rangle \rightarrow \mathbf{true} | \mathbf{false}$

$\langle \text{char_literal} \rangle \rightarrow ' \langle \text{char} \rangle '$

$\langle \text{string_literal} \rangle \rightarrow " \langle \text{char} \rangle^* "$

Semantics

A Decaf program consists of a single class declaration for a class called **Program**. The class declaration consists of field declarations and method declarations. Field declarations introduce variables that can be accessed globally by all methods in the program. Method declarations introduce functions/procedures. The program must contain a declaration for a method called **main** that has no parameters. Execution of a Decaf program starts at method **main**.

Types

There are two basic types in Decaf — **int** and **boolean**. In addition, there are arrays of integers (**int** [N]) and arrays of booleans (**boolean** [N]).

Arrays may be declared only in the global (class declaration) scope. All arrays are one-dimensional and have a compile-time fixed size. Arrays are indexed from 0 to $N - 1$, where $N > 0$ is the size of the array. The usual bracket notation is used to index arrays. Since arrays have a compile-time fixed size and cannot be declared as parameters (or local variables), there is no facility for querying the length of an array variable in Decaf.

Scope Rules

Decaf has simple and quite restrictive scope rules. All identifiers must be defined (textually) before use. For example:

- a variable must be declared before it is used.
- a method can be called only by code appearing after its header. (Note that recursive methods are allowed.)

There are at least two valid scopes at any point in a Decaf program: the global scope, and the method scope. The global scope consists of names of fields and methods introduced in the (single) **Program** class declaration. The method scope consists of names of variables and formal parameters introduced in a method declaration. Additional local scopes exist within each \langle block \rangle in the code; these can come after **if**, **for**, and **forpar** statements, or inserted anywhere a \langle statement \rangle is legal. An identifier introduced in a method scope can shadow an identifier from the global scope. Similarly, identifiers introduced in local scopes shadow identifiers in less deeply nested scopes, the method scope, and the global scope.

Variable names defined in the method scope or a local scope may shadow method names in the global scope. In this case, the identifier may only be used as a variable until the variable leaves scope.

No identifier may be defined more than once in the same scope. Thus field and method names must all be distinct in the global scope, and local variable names and formal parameters names must be distinct in each local scope.

Locations

Decaf has two kinds of locations: local/global scalar variables and (global) array elements. Each location has a type. Locations of types **int** and **boolean** contain integer values and boolean values, respectively. Locations of types **int** [N] and **boolean** [N] denote array elements. Since arrays are statically sized in Decaf, arrays may be allocated in the static data space of a program and need not be allocated on the heap.

Each location is initialized to a default value when it is declared. Integers have a default value of zero, and booleans have a default value of **false**. Local variables must be initialized when the declaring scope is entered. Array elements are initialized when the program starts.

Assignment

Assignment is only permitted for scalar values. For the types **int** and **boolean**, Decaf uses value-copy semantics, and the assignment $\langle \text{location} \rangle = \langle \text{expr} \rangle$ copies the value resulting from the evaluation of $\langle \text{expr} \rangle$ into $\langle \text{location} \rangle$. For array types, $\langle \text{location} \rangle$ and $\langle \text{expr} \rangle$ must refer to a single array element which is also a scalar value.

The $\langle \text{location} \rangle$ and the $\langle \text{expr} \rangle$ in an assignment must have the same type.

It is legal to assign to a formal parameter variable within a method body. Such assignments affect only the method scope.

Method Invocation and Return

Method invocation involves (1) passing argument values from the caller to the callee, (2) executing the body of the callee, and (3) returning to the caller, possibly with a result.

Argument passing is defined in terms of assignment: the formal arguments of a method are considered to be like local variables of the method and are initialized, by assignment, to the values resulting from the evaluation of the argument expressions. The arguments are evaluated from left to right.

The body of the callee is then executed by executing the statements of its method body in sequence.

A method that has no declared result type can only be called as a statement, *i.e.*, it cannot be used in an expression. Such a method returns control to the caller when **return** is called (no result expression is allowed) or when the textual end of the callee is reached.

A method that returns a result may be called as part of an expression, in which case the result of the call is the result of evaluating the expression in the **return** statement when this statement is reached. It is illegal for control to reach the textual end of a method that returns a result.

A method that returns a result may also be called as a statement. In this case, the result is ignored.

Control Statements

if

The **if** statement has the usual semantics. First, the $\langle \text{expr} \rangle$ is evaluated. If the result is **true**, the **true** arm is executed. Otherwise, the **else** arm is executed, if it exists. Since Decaf requires that the **true** and **else** arms be enclosed in braces, there is no ambiguity in matching an **else** arm with its corresponding **if** statement.

for

The **for** statement is similar to a **do** loop in Fortran. The $\langle \text{id} \rangle$ is the loop index variable and it shadows a variable of the same name declared in an outer scope if one exists. The loop index variable declares an integer variable whose scope is limited to the body of the loop. The first $\langle \text{expr} \rangle$ is the initial value of the loop index variable and the second $\langle \text{expr} \rangle$ is the ending value of the loop index variable. Each of these expressions are evaluated once, just prior to reaching the loop for

the first time. Each expression must evaluate to an integer value. The loop body is executed if the current value of the index variable is less than the ending value. After an execution of the loop body, the index variable is incremented by 1, and the new value is compared to the ending value to decide if another iteration should execute.

break and **continue** have their usual meaning: **break** exits the innermost enclosing loop, resuming execution at the next statement; **continue** branches to the beginning of the body of the innermost enclosing loop, increments the loop index, and re-compares the index variable to the ending value.

forpar

The **forpar** statement is similar to the **for** statement but it denotes that all iterations of the loop body can execute parallel. The programmer is responsible for determining that the code specified in a parallel loop is valid and suitable for parallel processing. If inter-iteration dependencies exist in a loop, then the execution may result in erroneous results. The programmer must ensure that these cases do not arise.

Any variable defined inside of the loop is considered private to an iteration and any variable that reaches the loop (global or defined in outer scope) is considered to be shared across iterations. The sharing for each scalar location is undefined, i.e., there is no synchronization for a shared scalar location between parallel iterations. An array can be properly shared across iterations as long as at most one iteration writes to each element of the array.

The behavior of executing a **return** or a **break** in the body of a **forpar** is undefined.

Expressions

Expressions follow the normal rules for evaluation. In the absence of other constraints, operators with the same precedence are evaluated from left to right. Parentheses may be used to override normal precedence.

A location expression evaluates to the value contained by the location.

Method invocation expressions are discussed in *Method Invocation and Return*. Array operations are discussed in *Types*. I/O related expressions are discussed in *Library Callouts*.

Integer literals evaluate to their integer value. Character literals evaluate to their integer ASCII values, e.g., 'A' represents the integer 65. (The type of a character literal is **int**.)

The arithmetic operators (`<arith_op>` and unary minus) have their usual precedence and meaning, as do the relational operators (`<rel_op>`). `%` computes the remainder of dividing its operands. The `<<` is a left shift and the `>>` is an arithmetic right shift.

Relational operators are used to compare integer expressions. The equality operators, `==` and `!=` are defined for **int** and **boolean** types only, can be used to compare any two expressions having the same type. (`==` is “equal” and `!=` is “not equal”).

The result of a relational operator or equality operator has type **boolean**.

The boolean connectives `&&` and `||` are interpreted using short circuit evaluation as in Java. The side-effects of the second operand are not executed if the result of the first operand determines the value of the whole expression (i.e., if the result is false for `&&` or true for `||`).

Operator precedence, from highest to lowest:

<i>Operators</i>	<i>Comments</i>
-	unary minus
!	logical not
* / %	multiplication, division, remainder
+ -	addition, subtraction
<< >>	shifts
< <= >= >	relational
== !=	equality
&&	conditional and
	conditional or

Note that this precedence is not reflected in the reference grammar.

Library Callouts

Decaf includes a primitive method for calling functions provided in the runtime system, such as the standard C library or user-defined functions.

The primitive method for calling functions is:

int callout ($\langle\text{string_literal}\rangle$, $[\langle\text{callout_arg}\rangle^+,]$) — the function named by the initial string literal is called and the arguments supplied are passed to it. Expressions of boolean or integer type are passed as integers; string literals or expressions with array type are passed as pointers. The return value of the function is passed back as an integer. The user of **callout** is responsible for ensuring that the arguments given match the signature of the function, and that the return value is only used if the underlying library function actually returns a value of appropriate type. Arguments are passed to the function in the system’s standard calling convention.

In addition to accessing the standard C library using **callout**, an I/O function can be written in C (or any other language), compiled using standard tools, linked with the runtime system, and accessed by the **callout** mechanism.

Semantic Rules

These rules place additional constraints on the set of valid Decaf programs besides the constraints implied by the grammar. A program that is grammatically well-formed and does not violate any of the following rules is called a *legal* program. A robust compiler will explicitly check each of these rules, and will generate an error message describing each violation it is able to find. A robust compiler will generate at least one error message for each illegal program, but will generate no errors for a legal program.

1. No identifier is declared twice in the same scope.
2. No identifier is used before it is declared.

3. The program contains a definition for a method called **main** that has no parameters (note that since execution starts at method **main**, any methods defined after **main** will never be executed).
4. The $\langle \text{int_literal} \rangle$ in an array declaration must be greater than 0.
5. The number and types of arguments in a method call must be the same as the number and types of the formals i.e., the signatures must be identical.
6. If a method call is used as an expression, the method must return a result.
7. A **return** statement must not have a return value unless it appears in the body of a method that is declared to return a value.
8. The expression in a **return** statement must have the same type as the declared result type of the enclosing method definition.
9. An $\langle \text{id} \rangle$ used as a $\langle \text{location} \rangle$ must name a declared local/global variable or formal parameter.
10. For all locations of the form $\langle \text{id} \rangle [\langle \text{expr} \rangle]$
 - (a) $\langle \text{id} \rangle$ must be an **array** variable, and
 - (b) the type of $\langle \text{expr} \rangle$ must be **int**.
11. The $\langle \text{expr} \rangle$ in **if** statement must have type **boolean**.
12. The operands of $\langle \text{arith_op} \rangle$ s and $\langle \text{rel_op} \rangle$ s must have type **int**.
13. The operands of $\langle \text{eq_op} \rangle$ s must have the same type, either **int** or **boolean**.
14. The operands of $\langle \text{cond_op} \rangle$ s and the operand of logical not (!) must have type **boolean**.
15. The $\langle \text{location} \rangle$ and the $\langle \text{expr} \rangle$ in an assignment, $\langle \text{location} \rangle = \langle \text{expr} \rangle$, must have the same type.
16. The initial $\langle \text{expr} \rangle$ and the ending $\langle \text{expr} \rangle$ of **for** and **forpar** must have type **int**.
17. All **break** and **continue** statements must be contained within the body of a **for** or **forpar**.

Run Time Checking

In addition to the constraints described above, which are statically enforced by the compiler's semantic checker, the following constraints are enforced dynamically: the compiler's code generator must emit code to perform these checks; violations are discovered at run-time.

1. The subscript of an array must be in bounds.
2. Control must not fall off the end of a method that is declared to return a result.

When a run-time error occurs, an appropriate error message is output to the terminal and the program terminates. Such error messages should be helpful to the programmer trying to find the problem in the source program.