The project for the course is to write a compiler for a language called Decaf. Decaf is a simple subset of C. To make debugging easier, you can directly compile Decaf source code using GCC to see what GCC generates.

# 1   Lexical Considerations

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

The reserved words are:

**bool break extern continue else false for while if int return sizeof true void**

Comments are started by // or # and are terminated by the end of the line. Text after a # is treated as a comment to make Decaf compatible with C; for GCC to compile Decaf programs directly, we have to use C preprocessor directives (starting with #) to tweak the programs a little bit. However, the Decaf compiler should ignore these C preprocessor directives. **Note: In any Decaf program that you want to compile with GCC, DO NOT USE # FOR COMMENTS. Although your Decaf compiler should ignore any text with a # in front of it, GCC will refuse to compile your program if it encounters any text starting with # that is not a valid C preprocessor directive.** We will give you a script for building your Decaf programs with GCC. This script will add C preprocessor directives to your Decaf code; you probably shouldn't modify these yourself unless you're sure you know what you're doing.

White space may appear between any lexical tokens. White space is defined as one or more spaces, tabs, line-break 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 64 bit signed. That is, decimal values between $-9223372036854775808$ and $9223372036854775807$. 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. For any integer literal, there may be an optional suffix "ll" following the last digit immediately without any whitespaces, and this suffix is purely for C compatibility and should be ignored in Decaf. Note that range checking is not performed during parsing. A long sequence of digits (e.g. 123456789123456789123) is still scanned as a single token.

A ⟨char⟩ is any printable ASCII character (ASCII values between decimal value 32 and 126) 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 literal tab, or "\n" to denote newline.

# 2 Reference Grammar

Meta-notation:

| | |
|---|---|
| $\langle foo \rangle$ | 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. |
| $\big[x\big]$ | means zero or one occurrence of $x$, *i.e.*, $x$ is optional; |
| | note that brackets in quotes $'\texttt{[}'\ '\texttt{]}'$ are terminals. |
| $x^*$ | means zero or more occurrences of $x$. |
| $x^+,$ | a comma-separated list of one or more x's. |
| | note that there is no comma following the last of $x$. |
| $\big\{\ \big\}$ | large braces are used for grouping; |
| | note that braces in quotes $'\texttt{\{}'\ '\texttt{\}}'$ are terminals. |
| $\vert$ | separates alternatives. |

$$\langle program \rangle \quad \rightarrow \quad \langle extern\_decl \rangle^* \ \langle field\_decl \rangle^* \ \langle method\_decl \rangle^*$$

$$\langle extern\_decl \rangle \quad \rightarrow \quad \textbf{extern} \ \langle id \rangle \ \textbf{()} \quad \textbf{;}$$

$$\langle field\_decl \rangle \quad \rightarrow \quad \langle type \rangle \ \big\{ \langle id \rangle \ \vert \ \langle id \rangle \ '\texttt{[}' \ \langle int\_literal \rangle \ '\texttt{]}' \ \big\}^+, \quad \textbf{;}$$

$$\langle method\_decl \rangle \quad \rightarrow \quad \big\{ \langle type \rangle \ \vert \ \textbf{void} \big\} \ \langle id \rangle \ \textbf{(} \ \Big[ \big\{ \langle type \rangle \ \langle id \rangle \big\}^+, \Big] \ \textbf{)} \ \langle block \rangle$$

$$\langle block \rangle \quad \rightarrow \quad '\texttt{\{}' \ \langle field\_decl \rangle^* \ \langle statement \rangle^* \ '\texttt{\}}'$$

$$\langle type \rangle \quad \rightarrow \quad \textbf{int} \ \vert \ \textbf{bool}$$

$$
\begin{aligned}
\langle statement \rangle \quad \rightarrow \quad & \langle location \rangle \ \langle assign\_op \rangle \ \langle expr \rangle \quad \textbf{;} \\
\vert \quad & \langle method\_call \rangle \quad \textbf{;} \\
\vert \quad & \textbf{if} \ \textbf{(} \ \langle expr \rangle \ \textbf{)} \ \langle block \rangle \ \big[ \textbf{else} \ \langle block \rangle \big] \\
\vert \quad & \textbf{for} \ \textbf{(} \ \langle id \rangle \ \textbf{=} \ \langle expr \rangle \ \textbf{;} \ \langle expr \rangle \ \textbf{;} \ \langle id \rangle \ \langle compound\_assign\_op \rangle \ \langle expr \rangle \ \textbf{)} \ \langle block \rangle \\
\vert \quad & \textbf{while} \ \textbf{(} \ \langle expr \rangle \ \textbf{)} \ \langle block \rangle \\
\vert \quad & \textbf{return} \ \big[ \langle expr \rangle \big] \quad \textbf{;} \\
\vert \quad & \textbf{break} \quad \textbf{;} \\
\vert \quad & \textbf{continue} \quad \textbf{;}
\end{aligned}
$$

$$
\begin{aligned}
\langle assign\_op \rangle \quad \rightarrow \quad & \textbf{=} \\
\vert \quad & \langle compound\_assign\_op \rangle
\end{aligned}
$$

$$
\begin{aligned}
\langle compound\_assign\_op \rangle \quad \rightarrow \quad & \textbf{+=} \\
\vert \quad & \textbf{-=}
\end{aligned}
$$

$$
\begin{aligned}
\langle method\_call \rangle \quad \rightarrow \quad & \langle method\_name \rangle \ \textbf{(} \ \big[ \langle expr \rangle^+, \big] \ \textbf{)} \\
\vert \quad & \langle method\_name \rangle \ \textbf{(} \ \big[ \langle extern\_arg \rangle^+, \big] \ \textbf{)}
\end{aligned}
$$

$$\langle method\_name \rangle \quad \rightarrow \quad \langle id \rangle$$

$$\langle location \rangle \quad \rightarrow \quad \langle id \rangle$$
$$| \quad \langle id \rangle \; '[' \; \langle expr \rangle \; ']'$$

$$\langle expr \rangle \quad \rightarrow \quad \langle location \rangle$$
$$| \quad \langle method\_call \rangle$$
$$| \quad \langle literal \rangle$$
$$| \quad \textbf{sizeof} \; ( \; \langle id \rangle \; )$$
$$| \quad \textbf{sizeof} \; ( \; \langle type \rangle \; )$$
$$| \quad \langle expr \rangle \; \langle bin\_op \rangle \; \langle expr \rangle$$
$$| \quad - \; \langle expr \rangle$$
$$| \quad ! \; \langle expr \rangle$$
$$| \quad ( \; \langle expr \rangle \; )$$

$$\langle extern\_arg \rangle \quad \rightarrow \quad \langle expr \rangle \; | \; \langle string\_literal \rangle$$

$$\langle bin\_op \rangle \quad \rightarrow \quad \langle arith\_op \rangle \; | \; \langle rel\_op \rangle \; | \; \langle eq\_op \rangle \; | \; \langle cond\_op \rangle$$

$$\langle arith\_op \rangle \quad \rightarrow \quad + \; | \; - \; | \; * \; | \; / \; | \; \%$$
$$\langle rel\_op \rangle \quad \rightarrow \quad < \; | \; > \; | \; <= \; | \; >=$$

$$\langle eq\_op \rangle \quad \rightarrow \quad == \; | \; !=$$

$$\langle cond\_op \rangle \quad \rightarrow \quad \&\& \; | \; ||$$

$$\langle literal \rangle \quad \rightarrow \quad \langle int\_literal \rangle \; [\textbf{ll}] \; | \; \langle char\_literal \rangle \; | \; \langle bool\_literal \rangle$$

$$\langle id \rangle \quad \rightarrow \quad \langle alpha \rangle \; \langle alpha\_num \rangle^*$$

$$\langle alpha\_num \rangle \quad \rightarrow \quad \langle alpha \rangle \; | \; \langle digit \rangle$$

$$\langle alpha \rangle \quad \rightarrow \quad a \; | \; b \; | \; \ldots \; | \; z \; | \; A \; | \; B \; | \; \ldots \; | \; Z \; | \; \_$$

$$\langle digit \rangle \quad \rightarrow \quad 0 \; | \; 1 \; | \; 2 \; | \; \ldots \; | \; 9$$

$$\langle hex\_digit \rangle \quad \rightarrow \quad \langle digit \rangle \; | \; a \; | \; b \; | \; c \; | \; d \; | \; e \; | \; f \; | \; A \; | \; B \; | \; C \; | \; D \; | \; E \; | \; F$$

$$\langle int\_literal \rangle \quad \rightarrow \quad \langle decimal\_literal \rangle \; | \; \langle hex\_literal \rangle$$

$$\langle decimal\_literal \rangle \quad \rightarrow \quad \langle digit \rangle \; \langle digit \rangle^*$$

$$\langle hex\_literal \rangle \quad \rightarrow \quad 0x \; \langle hex\_digit \rangle \; \langle hex\_digit \rangle^*$$

$$\langle bool\_literal \rangle \quad \rightarrow \quad \textbf{true} \; | \; \textbf{false}$$

$$\langle char\_literal \rangle \quad \rightarrow \quad ' \; \langle char \rangle \; '$$

$$\langle string\_literal \rangle \quad \rightarrow \quad " \; \langle char \rangle^* \; "$$

# 3 Semantics

A Decaf program consists of a single file. A program consists of extern declarations, 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 and returns **void**.

Execution of a Decaf program starts at method **main**.

## 3.1 Types

There are two basic types in Decaf — **int** and **bool**. In addition, there are single-dimensional arrays of integers (`int [ N ]`) and arrays of bools (`bool [ N ]`).

Arrays may be declared in any 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. Arrays are indexed by the usual bracket notation $a[i]$.

We use the size-of operator **sizeof** (Note that this is not a function) to evaluate the size of a variable (to an **int**) at compile time. The size of an **int** variable is 8, and the size of a **bool** variable is 1. The size of an array variable is equal to the number of elements in the array times the size of an array element. The **sizeof** operator can also take a **type**, which also evaluates to an **int**. **sizeof(int)** is 8, and **sizeof(bool)** is 1.

## 3.2 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 callouts, fields, and methods introduced in the top level of the program. The method scope consists of names of variables and formal parameters introduced in a method declaration. Additional local scopes exist within each ⟨block⟩ in the code; these can come after **if**, **while** and **for** statements. 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 or callout 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.

## 3.3 Locations

Decaf has two kinds of locations: local/global scalar variables and local/global array elements. Each location has a type. Locations of types **int** and **bool** contain integer values and bool values, respectively. Locations of types `int [ N ]` and `bool [ N ]` denote array elements. Since arrays

are statically sized in Decaf, global arrays may be allocated in the static data space of a program and need not be allocated on the heap. Local arrays may be dynamically allocated on the stack or statically allocated on the heap when appropriate.

Each location is initialized to a default value when it is declared. Integers have a default value of zero, and bools have a default value of **false**. Local variables must be initialized when the declaring scope is entered. Each element of a global array is initialized when the program starts. Each element of a local array is initialized when execution of the program enters the array's scope. In general, each time execution enters the scope of an array, its values must be reset to their defaults.

## 3.4   Assignment

Assignment is only permitted for scalar values. For the types **int** and **bool**, Decaf uses value-copy semantics, and the assignment ⟨location⟩ = ⟨expr⟩ copies the value resulting from the evaluation of ⟨expr⟩ into ⟨location⟩. The ⟨location⟩ += ⟨expr⟩ assignment increments the value stored in ⟨location⟩ by ⟨expr⟩, and is only valid for both ⟨location⟩ and ⟨expr⟩ of type **int**. The ⟨location⟩ -= ⟨expr⟩ assignment decrements the value stored in ⟨location⟩ by ⟨expr⟩, and is only valid for both ⟨location⟩ and ⟨expr⟩ of type **int**.

The ⟨location⟩ and the ⟨expr⟩ in an assignment must have the same type. For array types, ⟨location⟩ and ⟨expr⟩ must refer to a single array element which is also a scalar value.

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

## 3.5   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 one 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.

## 3.6   Control Statements

### 3.6.1   if

The **if** statement has the following semantics. First, the ⟨expr⟩ is evaluated. If the result is **true**, the **true** block is executed. Otherwise, the **else** block is executed, if it exists. Since Decaf requires

6

that the **true** and **false** blocks be enclosed in braces, there is no ambiguity in matching an **else** block with its corresponding **if** statement.

### 3.6.2 while

The **while** statement has the usual semantics. First, the ⟨expr⟩ is evaluated. If the result is **false**, control exits the loop. Otherwise, the loop body is executed. If control reaches the end of the loop body, the **while** statement is executed again.

### 3.6.3 for

The **for** statement is similar to C. The ⟨id⟩ is the loop index variable and must have been declared as an integer variable in the current scope or an outer scope. Because it must be an identifier, this means that array locations are not valid loop index variables. Before entering the loop body, it is assigned the value of the first ⟨expr⟩. The first expression is evaluated once, just prior to reaching the loop for the first time.

The second ⟨expr⟩ is the ending condition of the loop, which must be evaluated to **bool**, and it is evaluated every time before the loop body is executed, note that it may be an expression with side effects, for example, a function call.

After an execution of the loop body, the third part of the **for** statement "⟨id⟩⟨compound_assign_op⟩⟨expr⟩" is executed, which either increments or decrements ⟨id⟩ by the resulting value of ⟨expr⟩ (the third ⟨expr⟩ in the **for** statement). This ⟨expr⟩ must be evaluated every time, and it may have side effects as well. There is no restriction on what the ⟨id⟩ to be incremented or decremented should be, as it may be different from the first ⟨id⟩ in the **for** statement.

After incrementing or decrementing, the ending condition is checked again, and the loop repeats if it evaluates to **true**, and terminates if it evaluates to **false**.

## 3.7 Expressions

Expressions follow the normal rules for evaluation. In the absence of other constraints, operators (except ternary 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 *External Library*.

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**.)

We discussed the array length operator in Section 3.1.

The arithmetic operators (⟨arith_op⟩ and unary minus) have their usual meaning, as do the relational operators (⟨rel_op⟩). `%` computes the remainder of dividing its operands.

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

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

The boolean connectives `&&` and `||` are interpreted using short circuit evaluation as in Java. No side-effects of the second operand are 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:

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

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

## 3.8 External Library

Decaf includes a method for calling external functions similar to the C language. **extern** must be predeclared at the top of the file. The syntax (as specified in the grammar) is:

**extern** ⟨id⟩ () ;

All external functions are treated as if they return **int**. Once externs have been declared, they may be called similar to any function. The exceptions to this are that arguments to external functions may contain string literals and variables of array type. **This is the only use of the string literal in the decaf language**. Normal decaf methods may not contain string literals as arguments.

### 3.8.1 Extern Arguments

Expressions of bool or integer type are passed as integers; a string literal is passed as the memory address of its first character; an array variable is passed as the memory address of its first element. The return value of the function is passed back as an integer. The user of a **extern** 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. **The compiler is not responsible for verifying that externs have the correct number or type of arguments.**

### 3.8.2 External I/O Function

In addition to accessing the standard C library using **extern**, 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 **extern** mechanism.

# 4   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. This includes **extern** identifiers, which exist in the global 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 and returns **void** (note that since execution starts at method **main**, any methods defined after main will never be executed).

4. The ⟨int_literal⟩ in an array declaration must be greater than 0.

5. The number and types of arguments in a method call (non-extern) 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. String literals and array variables may not be used as arguments to non-extern methods. **Note: a[5] is not an array variable, it is an array location**

8. 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.

9. The expression in a **return** statement must have the same type as the declared result type of the enclosing method definition.

10. An ⟨id⟩ used as a ⟨location⟩ must name a declared local/global variable or formal parameter.

11. For all locations of the form ⟨id⟩[⟨expr⟩]

    (a) ⟨id⟩ must be an **array** variable, and
    (b) the type of ⟨expr⟩ must be **int**.

12. The argument of the **sizeof** operator must be a variable

13. The ⟨expr⟩ in an **if** or a **while** statement must have type **bool**.

14. The operands of ⟨arith_op⟩s and ⟨rel_op⟩s must have type **int**.

15. The operands of ⟨eq_op⟩s must have the same type, either **int** or **bool**.

16. The operands of ⟨cond_op⟩s and the operand of logical not (!) must have type **bool**.

17. The ⟨location⟩ and the ⟨expr⟩ in an assignment, ⟨location⟩ = ⟨expr⟩, must have the same type.

18. The ⟨location⟩ and the ⟨expr⟩ in an incrementing/decrementing assignment, ⟨location⟩ += ⟨expr⟩ and ⟨location⟩ -= ⟨expr⟩, must be of type **int**.

19. All **break** and **continue** statements must be contained within the body of a **for** or a **while**.

# 5  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. If the subscript of an array is found to be out of bounds, the program must terminate with exit value $-1$. If control falls off the end of a method that is declared to return a result, the program must terminate with exit value $-2$. The error messages output should be helpful to the programmer trying to find the problem in the source program.