

Democritus Language Final Report

Amy Xu
xx2152
Project Manager

Emily Pakulski
enp2111
Language Guru

Amarto Rajaram
aar2160
System Architect

Kyle Lee
kpl2111
Tester

May 11, 2016

Contents

1	Introduction	4
1.1	Motivation	4
1.2	Product Goals	4
	Native Concurrency and Atomicity	4
	Portability	4
	Flexibility	4
2	Language Tutorial	5
2.1	Setup and Installation	5
2.2	Compiling Your Code	5
2.3	Writing Code	6
2.4	Getting Started	6
	Declarations	6
	Types	6
	Primitives	6
	Strings	6
	Structs	7
	Operators	7
	Binary Operators:	7
	Unary Operators:	7
2.5	Pointers	8
2.6	Control Flow	8
	Conditional Branching	8
	Loops and Iteration	8
2.7	Multithreading and Atomicity	9
2.8	Miscellaneous	9
	Malloc	9
	File I/O	10
	Sockets API	10
3	Language Reference Manual	11
3.1	Introduction	11
3.2	Structure of a Democritus Program	11
3.3	Data types	13
	Primitive Types	13
	int	13
	float	13
	boolean	13
	pointer	13
	Complex Types	13

	string	13
	struct	13
3.4	Lexical Conventions	14
	Identifiers	14
	Reserved Keywords	14
	Literals	14
	Integer Literals	14
	Boolean Literals	14
	String Literals	14
	All Democritus Tokens	15
	Punctuation	16
	Semicolon	16
	Curly Brackets	16
	Parentheses	16
	Comments	16
3.5	Variable Declarations	16
	Variable Declaration	16
	Struct Declaration	17
3.6	Expressions and Operators	17
	Expressions	17
	Literal	17
	Identifier	17
	Binary Operation	17
	Unary Operation	17
	Struct Access	17
	Struct Assignment	17
	Function Call	17
	Variable Assignment	18
	Parenthisization	18
	Binary and Unary Operations	18
	Arithmetic Operations	19
	Addition and Subtraction	19
	Multiplication and Division	19
	Modulo	19
	Boolean Expressions	19
	Equality	19
	Negation	19
	Comparison	19
	Chained Expressions	19
	Parentheses	19
	Function Calls	20
	Pointers and References	20
	Operator Precedence and Associativity	20
3.7	Statements	20
	Expressions	21
	Return Statements	21
	Nested Blocks	21
	Conditional Statements	21
	Conditional Loops	21
3.8	Functions	22
	Overview and Grammar	22

	Calling and Recursion	22
3.9	Concurrency	23
	Overview	23
	Atomic Inline Declarations	23
	Atomic Parameter Declarations	23
	Spawning Threads	23
4	Project Plan	25
4.1	Planning	25
4.2	Workflow	25
4.3	Team Member Responsibilities	25
4.4	Git Logs	25
5	Architecture Overview	26
5.1	Compiler Overview	26
	The Scanner	26
	The Parser	26
	The Semantic Analyzer	26
	The Code Generator	26
6	Testing	27
6.1	Integration Testing	27
	Development and Testing Process	27
	Aside: Unit Testing	28
6.2	The Test Suite and Automated Regression Testing	28
7	Lessons Learned	29
7.1	Amy	29
7.2	Emily	29
7.3	Amarto	29
7.4	Kyle	29
8	Code Listing	30

1. Introduction

Democritus is a programming language with a static type system and native support for concurrent programming via its `atomic` keyword, with facilities for both imperative and functional programming. Democritus is compiled to the LLVM (Low Level Virtual Machine) intermediate form, which can then be optimized to machine-specific assembly code. Democritus' syntax draws inspiration from contemporary languages, aspiring to emulate Go and Python in terms of focusing on use cases familiar to the modern software engineer, emphasizing readability, and having “one – and preferably only one – obvious way to do it”¹.

1.1 Motivation

The main motivation behind Democritus was to create a lower level imperative language that supported concurrency ‘out-of-the-box’. The race condition arising from threading would be solved by the language’s `atomic` keyword, which provides a native locking system for variables and data. Users would then be able to write and run simple multi-threaded applications relatively quickly.

1.2 Product Goals

Native Concurrency and Atomicity

Users should be able to easily and quickly thread their program with minimal worry about race conditions. Their development process should not be hindered by the use of multithreading, nor should they have to define special threading classes as is common in some other languages.

Portability

Developed under the LLVM IR, code written in Democritus can be compiled and run on any machine that LLVM can run on. As an industrial-level compiler, LLVM offers robustness and portability as the compiler back-end of Democritus.

Flexibility

Though Democritus is not an object-oriented language, it seeks to grant users flexibility in functionality, supporting structures, standard primitive data types, and native string support.

¹<http://c2.com/cgi/wiki?PythonPhilosophy>

2. Language Tutorial

Democritus is a statically-typed, imperative language with standard methods for conditional blocks, iteration, variable assignment, and expression evaluation. In this chapter, we will cover environment configuration as well as utilizing both Democritus' basic and advanced features.

2.1 Setup and Installation

To set up the Democritus compiler, OCaml and LLVM must be installed. Testing and development was done in both native Ubuntu 15.04 and Ubuntu 14.04 running on a virtual machine.

```
sudo apt-get install m4 clang-3.7 clang clang-3.7-doc libclang-common-3.7-dev libclang-3.7-dev libclang1-3.7 libclang1-3.7-dbg libllvm-3.7-ocaml-dev libllvm3.7 libllvm3.7-dbg lldb-3.7 llvm-3.7 llvm-3.7-dev llvm-3.7-doc llvm-3.7-examples llvm-3.7-runtime clang-modernize-3.7 clang-format-3.7 python-clang-3.7 lldb-3.7-dev liblldb-3.7-dbg opam llvm-runtime
```

For Ubuntu 15.04, we need the matching LLVM 3.6 OCaml Library.

```
sudo apt-get install -y ocaml m4 llvm opam
opam init
opam install llvm.3.6 ocamlfind
eval `opam config env`
```

For Ubuntu 14.04:

```
sudo apt-get install m4 llvm software-properties-common

sudo add-apt-repository --yes ppa:avsm/ppa
sudo apt-get update -qq
sudo apt-get install -y opam
opam init

eval `opam config env`

opam install llvm.3.4 ocamlfind
```

After setting up the environment, clone the git repository into your desired installation directory:

```
git clone https://github.com/DemocritusLang/Democritus.git
```

2.2 Compiling Your Code

To build the compiler, cd into the Democritus repository, and run `make`.

If building fails, try running `eval 'opam config env'`, which should update your local environment use OPAM packages and compilers. It's recommended to add the above command to your shell's configuration file if you plan on developing with Democritus.

To compile code, simply run

```
./Democritus < filename.dem > outfile.lll
```

To run compiled code, call `lll` on the output:

```
lll outfile.lll
```

2.3 Writing Code

Code can be written in any text file, but Democritus source files should have the `.dem` extension by convention. Democritus programs consist of global function, struct, and variable declarations. Only the code inside `main()` will be executed at runtime. At this time, linking is not included in the Democritus compiler; all code should be written and compiled from a single `.dem` source file.

2.4 Getting Started

Declarations

Functions are declared with the `<function func_name(a type, b type) return_type>` syntax. Variables are declared with the `<let var_name var_type;>` syntax. Statements are terminated with the semicolon `;`. Note that all variable declarations must happen before statements (including assignments) in any given function.

```
function triangle_area(base int, height int) int{
    return base*height/2;
}
```

Types

Primitives

Primitive types in Democritus include booleans and integers. The void type is also used for functions.

Strings

Strings are built-in to Democritus. String literals are added to global static memory at runtime, and string variables point to the literals. These literals are automatically null-terminated.

```
function main() int{
    let s string;
    let foo int;
    let bar bool;

    bar = true;
    s = "Hello, World!";
    foo = 55;
    bar = false;

    return 0;
}
```

Structs

Structs are declared at the global level with the `<struct struct_type { named fields }>` syntax. Struct declarations may also be nested.

```
struct Person{
    let name string;
    let age int;
    let info struct Info;
}

struct Info{
    let education string;
    let salary int;
}

function main() int{
    let p struct Person;

    p.name = "Joe";
    p.age = 30;
    p.info.education = "Bachelor's";
    p.info.salary = 99999;

    print(p.name);
    print(" earns: ");
    print_int(p.info.salary);

    return 0;
}
```

Operators

Democritus includes the ‘standard’ set of operators, defined as follows:

Binary Operators:

- arithmetic: `+`, `-`, `*`, `/`, `%`
- logical: `==`, `!=`, `<`, `<=`, `>`, `>=`, `&&` (and), `||` (or)

Unary Operators:

- arithmetic: `-`
- logical: `!` (not)
- addressing: `&` (reference), `*` (reference)

Logical expressions return a boolean value.

The expressions on each side of a binary operation must be of the same type. The `&&`, `||`, and `!` operators must be called on boolean expressions.

References can only be called on addressable fields (such as variables, or struct fields). Dereferences can only be called on pointer types.

2.5 Pointers

2.6 Control Flow

As an imperative language, Democritus executes statements sequentially from the top of any given function to the bottom. Branching and iteration is done similarly to many other imperative languages.

Conditional Branching

Conditional branching is done with:

```
if(boolean expression)
{
    /* do something here */
}

else
{
    /* do alternative here */
}
```

Here is an example of conditional branching in Democritus:

```
struct Person{
    let education string;
    let name string;
    let age int;
    let working bool;
}

function main() int{
    let p Struct person;
    p.name = "Joe"
    p.education = "Bachelor's";
    p.age = 25;
    p.working = false;

    if(p.working){
        print(p.name);
        print(" works.\n");
    }else{
        print(p.name);
        print(" is looking for work.\n");
    }

    return 0;
}
```

This program prints “Joe is looking for work.”

Loops and Iteration

Iteration can be done either via a `for` loop. A `for(e1; e2; e3)` loop may take three expressions; `e1` is called prior to entering the loop, `e2` is a boolean conditional statement for the loop, and `e3` is called after

each iteration. Both **e1** and **e3** are optional; omitting both converts the **for** loop into the conditional **while** used in other languages.

```
function main() int{
    let i int;
    for(i = 0; i<42; i=i+1){
        i = i+1;
    }

    for(;false;){
        // This block will never be reached
    }

    for(true){
        print_int(i);
    }
    return 0;
}
```

This program will print 42 forever.

2.7 Multithreading and Atomicity

Democritus supports threading with the **thread()** function call, which then calls the underlying **pthread** function in C. Any defined function can be called with multiple threads. The calling syntax is as follows:

```
thread("functionname", <comma separated args>, #threads);
```

Multithreaded functions must take a **starvoid** type as input and return a **starvoid** to conform with C's calling convention. An example of a multithreaded program:

```
function multiprint(noop starvoid) starvoid{
    let x starvoid;
    print("Hello, World!\n");
    return x;
}

function main() int{
    thread("multiprint", 0, 6); /* "Hello, World!" will be printed six times. */
    return 0;
}
```

2.8 Miscellaneous

Besides threading, a couple of other functions from C have been bound to Democritus.

Malloc

malloc(size) may be called, returning a pointer to a newly heap-allocated block of *size* bytes. These pointers may be bound to strings, which themselves are pointers to string literals. An example utilizing malloc will be included with file I/O.

File I/O

Files may be opened with `open()`. This call returns an integer, which may be then bound as a file descriptor. C functions such as `write()`, `read()`, or `lseek()` may then be called on the file descriptor.

- `open(filename, fd, fd2)`: opens a file. `filename` is a string referring to a file to be opened. `fd` and `fd2` are file descriptors used for `open`.
- `write(fd, text, length)`: writes to a file. `fd` refers to the file descriptor of an open file. `text` is a string representing text to be written. `length` is an integer specifying the number of bytes to be written.
- `read(fd, buf, length)`: reads from a file. `fd` refers to the file descriptor of an open file. `buf` is a pointer to malloc'd or allocated space. `length` is an integer representing amount of data to be read. Buffers should be malloc'd before reading.
- `lseek(fd, offset, whence)`: sets a file descriptors cursor position. `fd` is the file descriptor to an open file. `offset` is an integer describing how many bytes the cursor should be offset by, and `whence` is an integer describing how `offset` should be applied: as an absolute location, relative location to the current cursor, or relative location to the end of the file.

For more detailed information on these calls, run `man function.name`.

```
function main() int{

    let fd int;
    let malloced string;

    fd = open("tests/HELLOOOOOO.txt", 66, 384);    /* Open this file */
    write(fd, "hellooo!\n", 10);                  /* Write these 10 bytes */

    malloced = malloc(10);                        /* Allocate space for the data and null terminator */
    lseek(fd, 0, 0);                              /* Jump to the front of the file */
    read(fd, malloced, 10);                       /* Read the data we just wrote into the buffer */

    print(malloced);                             /* Prints "hellooo!\n" */
    return 0;
}
```

Sockets API

3. Language Reference Manual

3.1 Introduction

In this language reference manual, Democritus, its syntax, and underlying operating mechanisms will be documented. In the grammars shown in this reference manual, all terminals are expressed in uppercase and all nonterminals are kept lowercase. The Lexical Conventions section will detail terminals (also known as tokens).

3.2 Structure of a Democritus Program

A basic Democritus program reduces to a list of global variable, struct, and function declarations. Code ‘to be executed’ should be written in functions. These declarations are accessible and usable from any scope in a Democritus program. At runtime, the function `main()` will be executed.

The full grammar of a program is as follows:

```
program:
    decls EOF

decls:
    /* nothing */
    | decls vdecl
    | decls fdecl
    | decls sdecl

fdecl:
    FUNCTION ID LPAREN formals_opt RPAREN typ LBRACE vdecl_list stmt_list RBRACE

formals_opt:
    /* nothing */
    | formal_list

formal_list:
    ID typ
    | formal_list COMMA ID typ

typ:
    INT
    | FLOAT
    | BOOL
    | VOID
```

```

| STRTYPE
| STRUCT ID
| VOIDSTAR
| STAR %prec POINTER typ

vdecl_list:
    /* nothing */
    | vdecl_list vdecl

vdecl:
    LET ID typ SEMI

sdecl:
    STRUCT ID LBRACE vdecl_list RBRACE

stmt_list:
    /* nothing */
    | stmt_list stmt

stmt:
    expr SEMI
    | RETURN SEMI
    | RETURN expr SEMI
    | LBRACE stmt_list RBRACE
    | IF LPAREN expr RPAREN stmt %prec NOELSE
    | IF LPAREN expr RPAREN stmt ELSE stmt
    | FOR LPAREN expr_opt SEMI expr SEMI expr_opt RPAREN stmt
    | FOR LPAREN expr RPAREN stmt

expr_opt:
    /* nothing */
    | expr

expr:
    LITERAL
    | FLOATLITERAL
    | TRUE
    | FALSE
    | ID
    | STRING
    | expr PLUS expr
    | expr MINUS expr
    | expr STAR expr
    | expr DIVIDE expr
    | expr MOD expr
    | expr EQ expr
    | expr NEQ expr
    | expr LT expr
    | expr LEQ expr
    | expr GT expr
    | expr GEQ expr

```

```

| expr AND      expr
| expr OR       expr
| expr DOT      ID
| expr DOT      ID ASSIGN expr
| MINUS expr     %prec NEG
| STAR expr     %prec Deref
| REF expr
| NOT expr
| ID ASSIGN expr
| ID LPAREN actuals_opt RPAREN
| LPAREN expr RPAREN

```

```

actuals_opt:
    /* nothing */
| actuals_list

```

```

actuals_list:
    expr
| actuals_list COMMA expr

```

3.3 Data types

Primitive Types

int

A standard 32-bit two's-complement signed integer. It can take any value in the inclusive range (-2147483648, 2147483647).

float

A 64-bit floating precision number, represented in the IEEE 754 format.

boolean

A 1-bit true or false value.

pointer

A 64-bit pointer that holds the value to a location in memory; pointers may be passed and dereferenced.

Complex Types

string

An immutable array of characters, implemented as a native data type in Democritus. String variables are 8-bit pointers to the location of the string literal in the global static memory.

struct

A struct is a simple user-defined data structure that holds various data types, such as primitives, other structs, or pointers.

3.4 Lexical Conventions

In this subsection, we will cover the standard lexical conventions for Democritus. Lexical elements are scanned as ‘tokens,’ which are then parsed into a valid Democritus program. Democritus is a free-format language, discarding all whitespace characters such as ‘ ’, `\t`, and `\n`.

Identifiers

Identifiers for Democritus will be defined as follows: any sequence of letters and numbers without whitespaces and not a keyword will be parsed as an identifier. Identifiers must start with a letter, but they may contain any lowercase or uppercase ASCII letter, numbers, and the underscore ‘_’. Identifiers are case-sensitive, so ‘var1’ and ‘Var1’ would be deemed separate and unique. Identifiers are used to identify named elements, such as variables, struct fields, and functions. Note that identifiers cannot begin with a number. The following is a regular expression for identifiers:

```
ID = "[ 'a'-'z' 'A'-'Z' ] [ 'a'-'z' 'A'-'Z' '0'-'9' '_' ] *"
```

Reserved Keywords

The following is a list of reserved Democritus keywords:

<code>if</code>	<code>else</code>	<code>for</code>	<code>return</code>	<code>int</code>
<code>bool</code>	<code>void</code>	<code>true</code>	<code>false</code>	<code>string</code>
<code>struct</code>	<code>*void</code>	<code>function</code>	<code>let</code>	

Literals

Literals are used to represent various values or constants within the language.

Integer Literals

Integer literals are simply a sequence of ASCII digits, represented in decimal.

```
INT = "[ '0'-'9' ] +"
```

Boolean Literals

Boolean literals represent the two possible values that boolean variables can take, `true` or `false`. These literals are represented in lowercase.

```
BOOLEAN = "true|false"
```

String Literals

String literals represent strings of characters, including escaped characters. String literals are automatically null-terminated. Strings are opened and closed with double quotations. A special OCaml `lexbuf` was used to parse string literals.

(Taken from <http://realworldocaml.org/v1/en/html/parsing-with-ocamllex-and-menhir.html>)

```
read_string buf = parse
```

```

| '"'          { STRING (Buffer.contents buf) }
| '\\\' '/'    { Buffer.add_char buf '/' ; read_string buf lexbuf }
| '\\\' '\\\'   { Buffer.add_char buf '\\\' ; read_string buf lexbuf }
| '\\\' 'b'     { Buffer.add_char buf 'b' ; read_string buf lexbuf }
| '\\\' 'f'     { Buffer.add_char buf '\012' ; read_string buf lexbuf }
| '\\\' 'n'     { Buffer.add_char buf 'n' ; read_string buf lexbuf }
| '\\\' 'r'     { Buffer.add_char buf 'r' ; read_string buf lexbuf }
| '\\\' 't'     { Buffer.add_char buf 't' ; read_string buf lexbuf }
| [^ '"' '\\\' ]+ { Buffer.add_string buf (Lexing.lexeme lexbuf);
                  read_string buf lexbuf
                }
}

```

All Democritus Tokens

The list of tokens used in Democritus are as follows:

```

| "/"          { comment lexbuf }
| "/*"        { multicomment lexbuf }
| '('         { LPAREN }
| ')'         { RPAREN }
| '{'         { LBRACE }
| '}'         { RBRACE }
| ';'         { SEMI }
| ','         { COMMA }
| '+'         { PLUS }
| '-'         { MINUS }
| '*'         { STAR }
| '&'         { REF }
| '.'         { DOT }
| '/'         { DIVIDE }
| '='         { ASSIGN }
| "=="        { EQ }
| "!="        { NEQ }
| '<'         { LT }
| "<="        { LEQ }
| ">"         { GT }
| ">="        { GEQ }
| "&&"        { AND }
| "||"        { OR }
| "!"         { NOT }
| "if"        { IF }
| "else"      { ELSE }
| "for"       { FOR }
| "return"    { RETURN }
| "int"       { INT }
| "bool"      { BOOL }
| "void"      { VOID }
| "true"      { TRUE }
| "string"    { STRTYPE }
| "struct"    { STRUCT }
| "*void"     { VOIDSTAR }

```



```

| "false"    { FALSE }
| "function" { FUNCTION }
| "let"      { LET }
| ['0'-'9']+ as lxm { LITERAL(int_of_string lxm) }
| ['a'-'z' 'A'-'Z'] ['a'-'z' 'A'-'Z' '0'-'9' '_' ]* as lxm { ID(lxm) }
| '"'       { read_string (Buffer.create 17) lexbuf } (* refer to read_string above *)
| eof { EOF }

```

Punctuation

Semicolon

The semicolon ‘;’ is required to terminate any statement in Democritus.

```
statement SEMI
```

Curly Brackets

In order to keep the language free-format, curly braces are used to delineate separate and nested blocks. These braces are required even for single-statement conditional and iteration loops.

```
LBACE statements RBACE
```

Parentheses

To assert precedence, expressions may be encapsulated within parentheses to guarantee order of operations.

```
LPAREN expression RPAREN
```

Comments

Comments may either be single-line, initialized with two backslashes, or multi-line, enclosed by `*` and `*\`.

```
COMMENT = ("// [^\n']* \n") | ("/*" [^ "*/"]* "*/")
```

3.5 Variable Declarations

In Democritus, local variables must be declared at the top of each function, before being later assigned.

Variable Declaration

Democritus requires all named variables to be declared with its type at the top of each function. Named variables are declared with the `let [ID] type` syntax. Assignment to these variables may then be done with `=`.

The grammar for variable declarations is as follows:

```

vdecl:
  LET ID typ SEMI

```

```

typ:
  INT

```

```
| BOOL
| VOID
| STRTYPE
| STRUCT ID
| VOIDSTAR
| STAR %prec POINTER typ
```

Struct Declaration

Structs are defined at the global scope, and can then be declared as variables. The global definitions are as follows:

```
sdecl:
    STRUCT ID LBRACE vdecl_list RBRACE

vdecl_list:
    | vdecl_list vdecl
```

3.6 Expressions and Operators

Expressions

Expressions may be any of the following:

Literal

A literal of any type, as detailed in the lexical conventions section.

Identifier

An identifier for a variable.

Binary Operation

A binary operation between an expression and another expression.

Unary Operation

A unary operation acting on the expression appearing on the immediate right of the operator.

Struct Access

An expression of a **struct** type accessing an identifier field with the **dot** (.) operator.

Struct Assignment

An expression of a **struct** type assigning a value to one of its fields (accessed with the **dot** (.) operator) using the **=** operator.

Function Call

A call to a function along with its formal arguments.

Variable Assignment

An identifier being assigned a value with the = operator.

Parenthisization

Another expression nested within parentheses.

The grammar for expressions is as follows:

```
expr:
  LITERAL
| FLOATLITERAL
| TRUE
| FALSE
| ID
| STRING
| expr PLUS expr          (* expr TERMINAL expr are binary operations *)
| expr MINUS expr
| expr STAR expr
| expr DIVIDE expr
| expr MOD expr
| expr EQ expr
| expr NEQ expr
| expr LT expr
| expr LEQ expr
| expr GT expr
| expr GEQ expr
| expr AND expr
| expr OR expr
| expr DOT ID             (* struct access *)
| expr DOT ID ASSIGN expr (* struct assign *)
| MINUS expr %prec NEG    (* unary arith negate *)
| STAR expr %prec Deref   (* unary deref *)
| REF expr                (* unary ref *)
| NOT expr                (* unary log negate *)
| ID ASSIGN expr
| ID LPAREN actuals_opt RPAREN (* function call *)
| LPAREN expr RPAREN        (* paren'd expr *)
```

Binary and Unary Operations

A binary operation operates on the two expressions on the left and right side of the operator. Binary operations may be:

- an addition, subtraction, mult., division, or modulo on two arithmetic expressions (+, -, *, /, %). Modulo only works on integer types.
- equality or inequality expression between boolean expressions (==, !=, <, <=, >, >=, &&, ||)

A unary operation operates on the expression on the operator's right side:

- a negation of an arithmetic expression (-)

- a dereference of a pointer type (*)
- an address reference of a variable or field within a struct (&)
- a negation of a boolean expression (!)

Arithmetic Operations

Democritus supports all the arithmetic operations standard to most general-purpose languages, documented below. Automatic casting has not been included in the language, and the compiler will throw an error in the case that arithmetic operations are performed between the same types of expressions.

Addition and Subtraction

Addition works with the + character, behaving as expected. Subtraction is called with -.

Multiplication and Division

Multiplication is called with *, and division with /. Division between integers discards the fractional part of the division.

Modulo

The remainder of an integer division operation can be computed via the modulo % operator.

Boolean Expressions

Democritus features all standard logical operators, utilizing ! for negation, and && and || for **and** and **or**, respectively. Each expression will return a boolean value of true or false.

Equality

Equality is tested with the == operator. Inequality is tested with !=. Equality may be tested on both boolean and arithmetic expressions.

Negation

Negation is done with !, a unary operation for boolean expressions.

Comparison

Democritus also features the <, <=, >, and >= operators. These represent less than, less than or equal to, greater than, and greater than or equal to, respectively. These operators are called on arithmetic expressions and return a boolean value.

Chained Expressions

Boolean expressions can be chained with && and ||, representing **and** and **or**. These operators have lower precedence than any of the other boolean operators described above. The **and** operator has a higher precedence than **or**.

Parentheses

Parentheses are used to group expressions together, since they have the highest order of precedence. Using parentheses will ensure that whatever is encapsulated within will be evaluated first.

Function Calls

Function calls are treated as expressions with a type equal to their return type. As an applicative-order language, Democritus evaluates function arguments first before passing them to the function. The grammar for function calls is as follows:

```
expr:
  .
  .
  .
  | ID LPAREN actuals_opt RPAREN    (* Function call *)

actuals_opt:
  /* nothing */
  | actuals_list

actuals_list:
  expr
  | actuals_list COMMA expr
```

Pointers and References

Referencing and dereferencing operations are used to manage memory and addressing in Democritus. The unary operator `&` gives a variable or struct field's address in memory, and the operator `*` dereferences a pointer type.

Operator Precedence and Associativity

Precedence	Operator	Description	Associativity
1	()	Parenthesis	Left-to-right
2	()	Function call	Left-to-right
3	*	Dereference	Right-to-left
	&	Address-of	
	!	Negation	
	-	Unary minus	
4	*	Multiplication	Left-to-right
	/	Division	
	%	Modulo	
5	+	Addition	Left-to-right
	-	Subtraction	
6	>> =	For relational > and ≥ respectively	Left-to-right
	<<=	For relational < and ≤ respectively	
7	== !=	For relational = and ≠ respectively	Left-to-right
8	&&	Logical and	Left-to-right
9		Logical or	Left-to-right
10	=	Assignment	Right-to-left

3.7 Statements

Statements written in a Democritus program are run from the top to bottom, sequentially. Statements can reduce to the following:

Expressions

An expression statement consists of an expression followed by a semicolon. Expressions in expression statements will be evaluated, with their values calculated.

Return Statements

A return statement is either a `RETURN SEMI` or `RETURN expr SEMI`. They are used as endpoints of a function, and control from a function returns to the original caller when a return statement is executed. Returns may be empty or return a type, though non-void functions must return an expression of their type.

Nested Blocks

A nested block is another statement list encapsulated within braces `{}`.

Conditional Statements

Conditional statements follow the `IF (boolean expr) stmt1 ELSE stmt2` format. When the `expr` evaluates to true, `stmt1` is run. Otherwise, if an `ELSE` and `stmt2` have been specified, `stmt2` is run.

Conditional Loops

A conditional loop is similar to a conditional statement, except in that it will loop or run repeatedly until its given boolean expression evaluates to `false`. In the case that the expression never evaluates to `false`, an infinite loop will occur.

Democritus eliminates the `while` keyword sometimes used in conditional iteration. Conditional loops follow the `FOR (expr1;boolean expr2;expr3) stmt` format where `expr1` and `expr3` are optional expressions to be evaluated prior to entering the `for` loop and upon each loop completion, respectively. Prior to entering or re-entering the loop, `expr2` is evaluated; control only transfers to `stmt` if this evaluation returns true. If both `expr1` and `expr3` are omitted, a simpler `for` loop can be written of the form `FOR (boolean expr) stmt`.

The full grammar for statements is as follows:

```
stmt_list:
    /* nothing */
    | stmt_list stmt

stmt:
    expr SEMI
    | RETURN SEMI
    | RETURN expr SEMI
    | LBRACE stmt_list RBRACE
    | IF LPAREN expr RPAREN stmt %prec NOELSE
    | IF LPAREN expr RPAREN stmt ELSE stmt
    | FOR LPAREN expr_opt SEMI expr SEMI expr_opt RPAREN stmt
    | FOR LPAREN expr RPAREN stmt
```

3.8 Functions

Overview and Grammar

Functions can be defined in Democritus to return one or no data type. Functions are evaluated via eager (applicative-order) evaluation and the function implementation must directly follow the function header. The grammar for function declarations is as follows:

```
fdecl:
    FUNCTION ID LPAREN formals_opt RPAREN typ LBRACE vdecl_list stmt_list RBRACE

formals_opt:
    /* nothing */
    | formal_list

formal_list:
    ID typ
    | formal_list COMMA ID typ
```

All functions require **return** statements at the end, and must return an expression of the same type as the function. **void** functions may simply terminate with an empty **return** statement.

```
    }
function function-name([formal-arg type, ... ]) type-r {

    [function implementation]
    return [variable of type type-r]
}
```

Calling and Recursion

Functions may be recursive and call themselves:

```
function recursive_func(i int) void {

    if (i < 0) {
        return;
    } else {
        print(  h i  );
        recursive_func(i-1);    // Call ourselves again.
    }
}
```

Functions may be called within other functions:

```
function main() void{
    recursive_func(3);
    return;                // Return nothing for void.
}
```

3.9 Concurrency

Overview

Democritus intends to cater to modern software engineering use cases. Developments in the field are steering us more and more towards highly concurrent programming as the scale at which software is used trends upward.

With this in mind, Democritus adds support for the `atomic` keyword, used as a modifier at the declaration step. The keyword can be used both in an inline declaration and when declaring a function's types. We also wrap the `pthread_t` datatype and related functions.

Atomic Inline Declarations

Under the hood, declaring a variable with the `atomic` keyword embeds a locking structure into the type, as well as exposing the `lock()` and `unlock()` functions. If the keyword is used with a standard data type, the compiler replaces the normal version of that type with a version that includes the lock and the functions described above.

```
}
x int;
x.lock(); x.unlock(); /* undefined! */

y atomic int;
y.lock(); y.unlock(); /* defined! */
}
```

Atomic Parameter Declarations

A function whose formal parameters are `atomic` will throw a compile-time error if a non-atomic type is passed in. The idea is to force the programmer to document which functions are safe to use in a multi-threaded context and which are not.

```
}

function [function name]([formal_arg atomic type]) atomic type {
    formal_arg.lock();
    /* do something */
    formal_arg.unlock();
    return formal_arg
}
```

Naturally, rather than calling `lock()` and `unlock()` manually, the programmer can implement atomic operations.

Spawning Threads

To spawn threads, Democritus uses a wrapper around the C-language `pthread` family of functions.

The `thread_t` data type wraps `pthread_t`.

To spawn a thread, the thread function takes a variable number of arguments where the first argument is a function and the remaining optional arguments are the arguments for that function. It returns an error code.

The `detach` boolean determines whether or not the parent thread will be able to join on the thread or not.


```
{  
  thread(f function, boolean detach, [arg1 type, arg2 type, ...]) int;  
}
```

4. Project Plan

4.1 Planning

Much of the planning was facilitated by our weekly meetings with David Watkins, our TA. He very clearly explained what the requirements of each milestone entailed, and helped keep expectations transparent. Since Professor Edwards had emphasized the need for vertical development of features instead of horizontal building of each compiler layer, we quickly identified the key features that would be required to enable the key functionality of our language. Two of the most important components were structs and threads.

Our initial plan, which we largely followed, was to complete the Language Reference Manual, experiment with the layers of the compiler and get `Hello, World!` working, and then use what we had learned to begin implementing the more crucial aspects of the language.

4.2 Workflow

Workflow was facilitated by Git and GitHub, which allowed for the team to easily work on multiple features simultaneously and (usually) merge together features without overlapping conflicts. The Git workflow reached an optimal point by the conclusion of the project; new features would be developed, tested, and finalized in separate branches, and the commits for that feature would be squashed down until a single commit representing the new feature would be merged into the master branch.

Features were developed individually or through paired programming, depending on the scope and complexity of the feature. GitHub and division of labor allowed for many team members to work independently, at different times of the day per their own schedule. Branching allowed for one person to quickly deploy buggy code to another in hopes of resolving the issue, without any modification or bad commits to the master branch. GroupMe was used extensively for inter-team communication.

4.3 Team Member Responsibilities

Team Member	Responsibilities	GitHub Handle
Amy Xu	structs, nested structs, pointers, malloc	axxu
Emily Pakulski	C-bindings, reformatting tests, atomicity	ohEmily
Amarto Rajaram	C-bindings, pthread, file I/O, malloc	Amarto
Kyle Lee	structs, LRM, Final Report, debug and assist Amy	kyle-lee

4.4 Git Logs

5. Architecture Overview

Democritus' compiler is built off of Professor Stephen Edwards' `MicroC` compiler.

5.1 Compiler Overview

Several files make up the source code of the compiler. These include:

- `scanner.mll`: the OCamllex scanner.
- `ast.ml`: the abstract syntax tree, summarizing the overall structure of a Democritus program.
- `parser.mly`: the Ocamlyacc parser. Tokens from the scanner are parsed into the abstract syntax tree in the parser.
- `semant.ml`: the semantic analyzer.
- `codegen.ml`: the LLVM IR code generator.
- `democritus.ml`: the overarching OCaml program that calls the four main steps of the compiler.
- `bindings.c`: a C file that provides key underlying C functionality, such as for threads, which is compiled to LLVM bytecode.

The Scanner

The scanner is simply a text scanner that parses text into various tokens, to then be interpreted by the parser. The regular expressions used by the scanner were listed above in the language reference chapter.

The Parser

The parser is a token scanner that converts the tokens read into a valid abstract syntax tree of the program. If the program follows valid syntax, it will be parsed accordingly. Otherwise, compilation of code will yield a parse error. The structure of the program is as follows:

The Semantic Analyzer

The Code Generator

6. Testing

As with any software project, extensive testing was required to verify that all the features being implemented were working properly.

6.1 Integration Testing

Development and Testing Process

Development of new features required making them pass through the scanner, parser, semantic analyzer, and then code generation, in that order. When envisioning or developing a new feature, the testing process would proceed as follows:

1. Write example code implementing and utilizing the desired feature. (E.g. writing a struct definition in a new test file).
2. Modify the scanner (if needed) to read new tokens required by the new feature.
3. Modify the parser (usually needed) to change the grammar of the program to accept the new feature and pass necessary information (E.g. struct field names) to the semantic analyzer.
4. Modify the example code and test it so that only the ‘correct’ implementation of the feature passes the parser. Modify the scanner and parser until this step passes.
5. Modify the semantic analyzer so that it detects possible semantic issues that could arise from utilization of the new feature (E.g. accessing an undefined field in a struct or an undefined struct).
6. Modify the example code and test it so that only the ‘correct’ implementation passes the semantic analyzer; try testing multiple cases that should cause the analyzer to raise an error. Modify the semantic analyzer until this step passes.
7. Modify the code generator so that it generates the appropriate LLVM IR representing your new feature (E.g., allocating the correct amount of memory for new structs, building a map of struct field indexes, calling `LLVM.build_struct_gep`, etc.).
8. Modify your example code to utilize your feature and produce some visible effect or output (E.g. assigning a struct field, doing arithmetic on it, then printing it).
9. Test the code and ensure that running the program produces the expected output or effect; continue working on code generation until it does.

The process of writing test code, compiling it, and observing its output after being run as LLVM IR was the integration testing method that the Democritus team utilized throughout development. It helped ensure that whole features were working properly, and that the language, built up from multiple features, was still functioning correctly. Integration testing was done on all new features added to the language, as well as the existing ones from MicroC (such as basic variable assignment, conditional iteration, etc).

Aside: Unit Testing

Unit testing was not overly utilized in this development process, besides for testing to ensure that new features could pass certain layers of the compiler while working towards a passing integration test. This is because unit tests can still pass, while whole features lose vertical integration in the process of building up a compiler. This is because new features may often conflict with each other and the successful introduction of one feature could very well mean the breaking of another. This leads us to the test suite and automated regression testing.

6.2 The Test Suite and Automated Regression Testing

Democritus' test suite was built upon MicroC's automated regression testing package. Within the `tests` directory, there are dozens of integration test files for various language features as well as their expected `stdout` output. Additionally, there are several 'fail' tests used for showing invalid Democritus code as well as their expected error outputs.

The automated regression testing suite was used to quickly test all major language features by compiling each test, writing the error thrown by compilation (if it was a failure) or output of running the LLVM file (if compilation was a success) to a temporary file, and comparing that output to the expected output of each test with `diff`. The automated test was a shell script, invoked with `./testall.sh` in the Democritus root directory.

The test suite was used frequently throughout development; while developing new features, team members would utilize the test suite to ensure that all major features of the language were still working. If a certain test in the suite failed, more verbose information about the test's failure could be accessed in the `testall.log` file generated by the testing suite. The automated regression testing was crucial in ensuring that the language stayed consistent and working, and that our master branch remained 'updated' and error-free.

7. Lessons Learned

7.1 Amy

7.2 Emily

7.3 Amarto

7.4 Kyle

8. Code Listing