



BAWK, A BINARY AWK

Kevin M. Graney
kmg2165@columbia.edu

Programming Languages and Translators
COMS W4115, Spring 2013, Professor Stephen Edwards

Contents

1	Introduction	3
1.1	Overview	3
1.2	Sample program	3
2	Language tutorial	4
2.1	A sample program	4
2.2	Executing bawk files	4
2.3	Walking through the example	6
3	Language reference manual	8
3.1	Lexical Conventions	8
3.1.1	Comments	8
3.1.2	Identifiers	8
3.1.3	Keywords	8
3.1.4	Constants	8
3.1.5	String Literals	8
3.2	Meaning of Identifiers	9
3.2.1	Function Name	9
3.2.2	Variable Name	9
3.2.3	Pattern Binding Variable Name	9
3.2.4	Special Identifiers	9
3.3	Statements	9
3.3.1	Block Statements	10
3.3.2	Pattern Statements	10
3.3.3	Expression Statements	10
3.3.4	Function Declaration	10
3.3.5	If statement	10
3.3.6	If/else statement	11
3.4	Expressions	11
3.4.1	Primary Expression	11
3.4.2	Function Calls	11
3.4.3	Multiplicative Operators	11
3.4.4	Additive Operators	12
3.4.5	Relational Operators	12
3.4.6	Equality Operators	12

3.4.7	Assignment Expressions	12
3.5	Pattern Expressions	12
3.5.1	Constants	13
3.5.2	Bindings	13
3.5.3	Pattern Variables	14
3.5.4	String Literals	14
4	Project plan	15
4.1	Project log	16
5	Architectual design	48
5.1	Design decisions	48
5.2	Execution phases	48
5.2.1	Bytecode compiler	48
5.2.2	Bytecode operations	49
5.2.3	Bytecode runtime	50
5.3	Limitations	50
6	Test plan	51
7	Lessons learned	52
7.1	Language design	52
7.2	Scanner and parser generators	52
7.2.1	GraphViz output	52
7.3	Compiler and bytecode design	53
7.3.1	Return statement problem	53
7.3.2	Lack of full support for string types	53
8	Source code listing	54
8.1	Compiler and bytecode interpreter	54
8.1.1	Scanner & parser definitions	54
8.1.2	Interfaces	58
8.1.3	Source Files	63
8.2	Test files	77
8.2.1	test-arith.bawk	77
8.2.2	test-cond1.bawk	77
8.2.3	test-func1.bawk	78
8.2.4	test-func2.bawk	78
8.2.5	test-func3.bawk	78
8.2.6	test-pat1.bawk	79
8.2.7	test-pat2.bawk	79
8.2.8	test-pat3.bawk	80
8.2.9	test-pat4.bawk	80
8.2.10	test-pat5.bawk	81
8.2.11	test-pat6.bawk	81
8.2.12	test-pat7.bawk	82
8.2.13	test-pat8.bawk	82

8.2.14	test-pat9.bawk	83
8.2.15	test-pat10.bawk	84
8.2.16	test-pat11.bawk	84
8.2.17	test-scope1.bawk	85
8.2.18	test-scope2.bawk	85
8.2.19	test-useful1.bawk	86
8.2.20	test-var1.bawk	86
8.3	Misc. Files	86

Listings

1.1	A simple bawk program	3
2.1	Showing how RP works	6
4.1	Git log for project development	16
5.1	Possibly ambiguous code	48
8.1	scanner.mll	54
8.2	parser.mly	56
8.3	parser_help.mli	58
8.4	ast_types.mli	59
8.5	bytecode_types.mli	59
8.6	ast.mli	60
8.7	compile.mli	61
8.8	bytecode.mli	62
8.9	reader.mli	62
8.10	utile.mli	63
8.11	bawk.ml	63
8.12	parser_help.ml	64
8.13	ast.ml	65
8.14	compile.ml	68
8.15	bytecode.ml	72
8.16	reader.ml	75
8.17	utile.ml	76
8.18	Test program (test-arith.bawk)	77
8.19	Expected output (test-arith.bawk.out)	77
8.20	Test program (test-cond1.bawk)	77
8.21	Expected output (test-cond1.bawk.out)	77
8.22	Test program (test-func1.bawk)	78
8.23	Expected output (test-func1.bawk.out)	78
8.24	Test program (test-func2.bawk)	78
8.25	Expected output (test-func2.bawk.out)	78
8.26	Test program (test-func3.bawk)	78
8.27	Expected output (test-func3.bawk.out)	79
8.28	Test program (test-pat1.bawk)	79
8.29	Expected output (test-pat1.bawk.out)	79
8.30	Test program (test-pat2.bawk)	79
8.31	Expected output (test-pat2.bawk.out)	80

8.32	Test program (test-pat3.bawk)	80
8.33	Expected output (test-pat3.bawk.out)	80
8.34	Test program (test-pat4.bawk)	80
8.35	Expected output (test-pat4.bawk.out)	81
8.36	Test program (test-pat5.bawk)	81
8.37	Expected output (test-pat5.bawk.out)	81
8.38	Test program (test-pat6.bawk)	81
8.39	Expected output (test-pat6.bawk.out)	81
8.40	Test program (test-pat7.bawk)	82
8.41	Expected output (test-pat7.bawk.out)	82
8.42	Test program (test-pat8.bawk)	82
8.43	Expected output (test-pat8.bawk.out)	83
8.44	Test program (test-pat9.bawk)	83
8.45	Expected output (test-pat9.bawk.out)	83
8.46	Test program (test-pat10.bawk)	84
8.47	Expected output (test-pat10.bawk.out)	84
8.48	Test program (test-pat11.bawk)	84
8.49	Expected output (test-pat11.bawk.out)	84
8.50	Test program (test-scope1.bawk)	85
8.51	Expected output (test-scope1.bawk.out)	85
8.52	Test program (test-scope2.bawk)	85
8.53	Expected output (test-scope2.bawk.out)	85
8.54	Test program (test-useful1.bawk)	86
8.55	Expected output (test-useful1.bawk.out)	86
8.56	Test program (test-var1.bawk)	86
8.57	Expected output (test-var1.bawk.out)	86
8.58	Makefile	86
8.59	run_tests.sh	87

Chapter 1

Introduction

1.1 Overview

The bawk language, whose name is derived from ‘binary awk’, is intended to be a small, special-purpose language for the parsing of formatted binary files. In the spirit of awk, bawk will match rules in binary files and extract data based on these rules. Bawk will solve the problem of quickly decoding a binary file of a known format to extract information much the way awk solves this same problem for formatted text files.

The bawk interpreter will take two inputs for execution: (1) text in the bawk language, and (2) binary data to run the bawk program on. This workflow will be similar to that of a basic awk workflow, and the language itself is inspired by awk, with the same basic program construct. That is, a program will be a series of patterns and statements, with the statements executing when their respective pattern matches.

1.2 Sample program

A representative bawk program is shown below. This program operates on PNG format image files and will print two numbers, newline separated. The first number is the width of the image and the second number is the height.

Listing 1.1: A simple bawk program

```
1 /* print the size of a PNG file */
2
3 /89 "PNG" 0d 0a 1a 0a/ {
4     /length:uint4 "IHDR"/ {
5         /width:uint4 height:uint4/ {
6             print(width);
7             print(height);
8         }
9     }
10 }
```

Chapter 2

Language tutorial

The bawk language is based on pattern matching bytes and extracting values from binary format data files.

2.1 A sample program

The sample program shown in §1.2 illustrates how to print the width and height of a given PNG image file using bawk. This is a perfect example of an application of the bawk language: a simple binary format file for which the extraction of integer data is desired. We will use this sample program in the next section to show how bawk is executed.

2.2 Executing bawk files

To execute a bawk file, call bawk with a `-e` flag and specify the file on which to operate. The bawk code is then to be provided via standard input.

```
./bawk -e [binary filename] < [bawk filename]
```

This command will compile and run the bawk program over the given binary file. The `-e` flag is optional and may be omitted. (The default operation is to execute the program if no flags are given.)

To execute the sample program, named `tests/test-util1.bawk`, we run bawk on a sample file, `tests/lichtenstein.png`, and see that the height and width of this file are both 512 pixels.

```
$ ./bawk -e tests/lichtenstein.png < tests/test-useful1.bawk
512
512
```

To see the bytecode instructions for a given program use instead the `-c` flag. With this flag, the bytecode instructions will be printed to standard output.

```
./bawk -c < [bawk filename]
```

If we compile our sample program for PNG files, we can see the bytecode generated, which is shown below.


```

$ ./bawk -c < tests/test-useful1.bawk
0: Ldp
1: Ldp
2: Beo 124
3: Lod 0
4: Bne 10
5: Rdb 1
6: Lit 137
7: Bin Sub
8: Bne 119
9: Bra 14
10: Rdb 1
11: Lit 137
12: Bin Sub
13: Bne 119
14: Rdb 1
15: Lit 80
16: Bin Sub
17: Bne 119
18: Rdb 1
19: Lit 78
20: Bin Sub
21: Bne 119
22: Rdb 1
23: Lit 71
24: Bin Sub
25: Bne 119
26: Lod 0
27: Bne 33
28: Rdb 1
29: Lit 13
30: Bin Sub
31: Bne 119
32: Bra 37
33: Rdb 1
34: Lit 13
35: Bin Sub
36: Bne 119
37: Lod 0
38: Bne 44
39: Rdb 1
40: Lit 10
41: Bin Sub
42: Bne 119
43: Bra 48
44: Rdb 1
45: Lit 10
46: Bin Sub
47: Bne 119
48: Lod 0
49: Bne 55
50: Rdb 1
51: Lit 26
52: Bin Sub
53: Bne 119
54: Bra 59
55: Rdb 1
56: Lit 26
57: Bin Sub
58: Bne 119
59: Lod 0
60: Bne 66
61: Rdb 1
62: Lit 10
63: Bin Sub
64: Bne 119
65: Bra 70
66: Rdb 1
67: Lit 10
68: Bin Sub
69: Bne 119
70: Ldp
71: Ldp
72: Beo 116
73: Rdb 4
74: Str 2
75: Rdb 1
76: Lit 73
77: Bin Sub
78: Bne 111
79: Rdb 1
80: Lit 72
81: Bin Sub
82: Bne 111
83: Rdb 1
84: Lit 68
85: Bin Sub
86: Bne 111
87: Rdb 1
88: Lit 82

```

89: Bin Sub	108: Skp
90: Bne 111	109: Skp
91: Ldp	110: Bra 116
92: Ldp	111: Skp
93: Beo 108	112: Beo 116
94: Rdb 4	113: Rdb 1
95: Str 3	114: Drp
96: Rdb 4	115: Bra 71
97: Str 4	116: Skp
98: Lod 3	117: Skp
99: Jsr -1	118: Bra 124
100: Lod 4	119: Skp
101: Jsr -1	120: Beo 124
102: Bra 108	121: Rdb 1
103: Skp	122: Drp
104: Beo 108	123: Bra 1
105: Rdb 1	124: Skp
106: Drp	125: Skp
107: Bra 92	126: Hlt

2.3 Walking through the example

The first line of the sample program in §1.2 matches the pattern 89 "PNG" 0d0a1a0a, where the string "PNG" is translated into ASCII bytes 50 4e 47. Bawk will scan through the input file until finding this pattern. On the first occurrence, the block statement following will be executed, and the nested patterns will attempt to match starting at the end of the match from the higher-level pattern. When a block exits, the pattern matching returns to executing at the position where it last began searching.

There is a special variable in bawk, RP, that can be used to retrieve or set the *read pointer* over the binary file. The value of RP represents the number of bytes from the beginning of the file. If we modify the example program to show us RP at various execution points it should become clear how the variable behaves.

Listing 2.1: Showing how RP works

```

1 print (RP);
2 /89 "PNG" 0d 0a 1a 0a/ {
3   print (RP);
4   /length:uint4 "IHDR"/ {
5     print (RP);
6     /width:uint4 height:uint4/ {
7       print (RP);
8     }
9     print (RP);
10  }
11  print (RP);
12 }
13 print (RP);
14 RP + 100;

```

```
15 print (RP);  
16 /ffee/ { print (RP); }  
17 print (RP);
```

Running this code produces the following output on `lichtenstein.png`.

```
0  
8  
16  
24  
16  
8  
0  
0  
15705  
0
```

As you can see, the RP pointer resets at the end of each pattern matching block, and can also be manually advanced. Pattern matching always searches forward from RP.

In addition to pattern matching `bawk` supports dynamic variable scoping, functions, conditional statements, and mixed endianness. See the language reference manual in Chapter 3 for details.

Chapter 3

Language reference manual

3.1 Lexical Conventions

A program consists of a single bawk language character file. There are six classes of tokens: identifiers, keywords, constants, string literals, operators, and separators. White space characters are ignored except as they separate tokens or when they appear in string literals.

3.1.1 Comments

Comments follow the ANSI C style, beginning with `/*` and ending with `*/`. They do not nest, and they cannot be present in quoted strings.

3.1.2 Identifiers

An identifier is a sequence of uppercase and lowercase letters A–Z, the numerals 0–9, and the underscore character. The first character of an identifier can not be a numeral.

```
<identifier> ::= <letter>  
             | <identifier> (<letter> | <digit>)
```

3.1.3 Keywords

The following identifiers are reserved for use as keywords, and should not be used otherwise: `if`, `else`, `return`, `def`, `int1`, `int2`, `int4`, `uint1`, `uint2`, `uint4`.

3.1.4 Constants

Bawk only supports a single kind of constant: the integer constant. The semantic meaning of integer constants depends on the context in which they occur. Except for inside of pattern expressions (see §3.5.1), the following is true. An integer constant consisting of a sequence of digits is taken to be in decimal. A sequence of hexadecimal digits prefaced by `0x` is taken to be in hexadecimal (base 16).

3.1.5 String Literals

A string is a sequence of characters enclosed in double quotes. The only valid place for strings is inside of pattern expressions (see §3.5.4).

3.2 Meaning of Identifiers

Identifiers can have several different meanings based on how they are first used. Once an identifier is associated with a particular meaning it cannot be disassociated with that meaning.

3.2.1 Function Name

Previously unused identifiers become function names once a function is declared with a given name.

3.2.2 Variable Name

Previously unused identifiers become variable names once an assignment expression (see §3.4.7) executes with the new identifier on the left hand side. Variables have scope where they are first created in this way. Scoping is dynamic.

3.2.3 Pattern Binding Variable Name

Previously unused identifiers become binded variable names once a pattern expression with a binding pattern token (see §3.5.2) matches.

3.2.4 Special Identifiers

There are certain special identifiers that behave like variables, but whose meaning affects the program execution.

- RP – The *file pointer* variable indicates where in the file pattern statements begin matching (see §3.3.2). The value contains the number of bytes from the beginning of the file. At the start of the bawk program this value defaults to 0.
- LE – Force *little-endian* interpretation of integers and strings while pattern matching in the file. A value of 0 indicates big-endian interpretation and a value of 1 indicates little-endian interpretation. The default value is 0.

3.3 Statements

The expressions described in §3.4 are a specific form of statement in the bawk language. There are three different types of statements, and a list of statements forms a bawk program. The three different types of statements are: expressions (described in §3.4), block statements, and pattern statements.

$\langle \text{program} \rangle ::= \langle \text{statement-list} \rangle$

$\langle \text{statement-list} \rangle ::= e$
| $\langle \text{statement} \rangle$
| $\langle \text{statement-list} \rangle \langle \text{statement} \rangle$

$\langle \text{statement} \rangle ::= \langle \text{expression-statement} \rangle$
 $\quad | \langle \text{block-statement} \rangle$
 $\quad | \langle \text{pattern-statement} \rangle$
 $\quad | \langle \text{function-declaration} \rangle$

3.3.1 Block Statements

Block statements are used to combine a sequence of statements into a single statement. The sequence of statements inside the block is executed in order when the block itself is executed. Variables first used inside the block are locally scoped, but names from outside the block are also available.

$\langle \text{block-statement} \rangle ::= \{ \langle \text{statement-list} \rangle \}$

3.3.2 Pattern Statements

Pattern statements allow a type of pattern matching to be performed on the binary file. Pattern statements consist of a pattern expression (see §3.5) and a statement. If the pattern expression matches the data file at the location of RP at the time of evaluation of the pattern statement, then the corresponding statement is executed. If the pattern does not match, then the statement is not executed.

$\langle \text{pattern-statement} \rangle ::= / \langle \text{pattern-expression} \rangle / \langle \text{statement} \rangle$

3.3.3 Expression Statements

An expression statement is an expression followed by a semicolon. Expressions are described in §3.4.

$\langle \text{expression-statement} \rangle ::= \langle \text{expression} \rangle ;$

3.3.4 Function Declaration

A function declaration is a statement declaring a new function. Function declarations consist of an identifier to serve as the function name, followed by an open parenthesis, an optional comma separated list of identifiers to serve as function parameter names, a closing parenthesis, and a statement to serve as the function body.

$\langle \text{function-declaration} \rangle ::= \text{def} \langle \text{identifier} \rangle (\langle \text{function-decl-params} \rangle) \langle \text{statement} \rangle$

$\langle \text{function-decl-params} \rangle ::= \epsilon$
 $\quad | \langle \text{identifier} \rangle$
 $\quad | \langle \text{function-decl-params} \rangle , \langle \text{identifier} \rangle$

3.3.5 If statement

If the expression provided to the if statement evaluates to a non-zero value then the subsequent statement is executed. If the expression evaluates to a value of zero, then the statement is not executed.

$\langle \text{if-statement} \rangle ::= \text{if} (\langle \text{expression} \rangle) \langle \text{statement} \rangle$

3.3.6 If/else statement

If/else statements work like if statements, with an additional second statement. If the expression evaluates to a non-zero value then the first statement is executed and the second is not. If the expression evaluates to zero then the first statement is not executed and the second is executed.

$\langle \text{if-statement} \rangle ::= \text{'if' ' (' } \langle \text{expression} \rangle \text{')' } \langle \text{statement} \rangle \text{' else' } \langle \text{statement} \rangle$

3.4 Expressions

3.4.1 Primary Expression

$\langle \text{primary-expression} \rangle ::= \langle \text{constant} \rangle$
| $\langle \text{identifier} \rangle$
| $\langle \text{string-literal} \rangle$
| $\langle \text{function-call} \rangle$
| $\text{' (' } \langle \text{expression} \rangle \text{')'}$

$\langle \text{expression} \rangle ::= \langle \text{assignment-expression} \rangle$
| $\langle \text{additive-expression} \rangle$

3.4.2 Function Calls

A function call is an identifier, known as the function name, followed by a pair of parenthesis containing a, possibly empty, comma-separated list of expressions. These expressions constitute arguments to the function. When functions are called a copy of each argument is made and used within the function, that is, all function calls pass arguments by value.

$\langle \text{function-call} \rangle ::= \langle \text{identifier} \rangle \text{' (' } \langle \text{function-call-params} \rangle \text{')'}$

$\langle \text{function-call-params} \rangle ::= \epsilon$
| $\langle \text{expression} \rangle$
| $\langle \text{function-call-params} \rangle \text{' , ' } \langle \text{expression} \rangle$

3.4.3 Multiplicative Operators

The multiplicative operators, * and /, are left associative. The result of * is multiplication of the two operands, and the result of / is the quotient of the two operands. The quotient with a second operand of zero yields an undefined result.

$\langle \text{multiplicative-expression} \rangle ::= \langle \text{primary-expression} \rangle$
| $\langle \text{multiplicative-expression} \rangle \text{' * ' } \langle \text{primary-expression} \rangle$
| $\langle \text{multiplicative-expression} \rangle \text{' / ' } \langle \text{primary-expression} \rangle$

3.4.4 Additive Operators

The two additive operators, + and -, are left associative. The expected arithmetic operation is performed on integers. The result of the + operator is the sum of the operands, and the result of the - operator is the difference of the operands.

```
<additive-expression> ::= <multiplicative-expression>  
| <additive-expression> '+' <multiplicative-expression>  
| <additive-expression> '-' <multiplicative-expression>
```

3.4.5 Relational Operators

The relational operators are all left associative. The operators all produce 0 if the specified relation is false and 1 if the specified relation is true. The operators are < (less than), > (greater than), <= (less than or equal to), >= (greater than or equal to).

```
<relational-expression> ::= <additive-expression>  
| <relational-expression> '<' <additive-expression>  
| <relational-expression> '>' <additive-expression>  
| <relational-expression> '<=' <additive-expression>  
| <relational-expression> '>=' <additive-expression>
```

3.4.6 Equality Operators

The == (equal to) and != (not equal to) operators compare operands for equality. Like the relational operators, the equality operators produce 0 if the relationship is false and 1 if the relationship is true. The equality operators have a lower precedence than the relational operators, and are also left associative.

```
<equality-expression> ::= <relational-expression>  
| <equality-expression> '==' <relational-expression>  
| <equality-expression> '!=' <relational-expression>
```

3.4.7 Assignment Expressions

All assignment expressions require a modifiable identifier on the left hand side. The only type of modifiable identifier is a variable (see §3.2.2). The assignment operator, =, is right-associative and the return value of assignment is the value being assigned.

```
<assignment-expression> ::= <equality-expression>  
| <identifier> '=' <assignment-expression>
```

3.5 Pattern Expressions

A pattern expression matches content in the binary data file. Pattern expressions are distinct in syntax from the expressions described in §3.4.


```

<pattern-expression> ::= c
    | <pattern-token>
    | <pattern-expression> <pattern-token>

```

```

<pattern-token> ::= <pattern-constant>
    | <pattern-binding>
    | <pattern-variable>
    | <string-literal>

```

Pattern expressions can contain constants or typed variable bindings. Each constant or binding is known as a pattern term, and pattern terms are whitespace delimited. Already bound pattern variables can also be used in pattern statements.

3.5.1 Constants

Inside pattern expressions, constants work differently than those described in §3.1.4. All constants in pattern terms are implicitly expressed in hexadecimal form, and the 0x prefix required to express hexadecimal constants in §3.1.4 must be omitted.

Constants are read in hexadecimal as bytes. For this reason leading zeros on a value have semantic meaning. For example /0000abcd/ matches the four byte pattern 00 00 ab cd in the binary file, while /abcd/ matches the two byte pattern ab cd.

The byte order expected in the file is dependent on the run-time value of the LE special variable (see §3.2.4). If this variable is non-zero, then constants are interpreted in *little endian* order. If the variable is zero, then the interpretation is *big endian*. The default value is zero. As an example, if LE != 0 the pattern /aabb/ matches the file with two bytes in the order bb aa, and if LE == 0 then the pattern matches bytes aa bb in that order.

Whitespace does make a difference in the interpretation of pattern constants. When LE != 0, the patterns /aabb/ and /aa bb/ are distinctly different. If LE == 0, these two patterns are the same. Whitespace added between byte boundaries also changes the semantic meaning; /01ab/ is *not* identical to /0 1 a b/. If a constant has an odd number of hexadecimal digits, an implied 0 is added to the left-most side of the value. That is /123/ is the same as /0123/ (is the same as /01 23/ if LE == 0). For this reason the example of /0 1 a b/ is identical to /00 01 0a 0b/ (and /00010a0b/ if LE == 0).

3.5.2 Bindings

A bind pattern term consists of an identifier, a colon (:), and a bind type. Valid bind types are: int1, int2, int4, uint2, uint4. The int1 type matches a one-byte integer, the uint4 type matches an unsigned 4-byte integer. Identifiers have the same naming rules as §3.1.2.

```

<pattern-binding> ::= <identifier> ':' <bind-type>

```

```

<bind-type> ::= 'int1'
    | 'int2'
    | 'int4'
    | 'uint1'

```

```
| 'uint2'  
| 'uint4'
```

3.5.3 Pattern Variables

Previously bound pattern variables can be used in a pattern expression to match the previously bound value.

⟨pattern-variable⟩ ::= ⟨identifier⟩

3.5.4 String Literals

String literals are also supported in pattern expressions. The syntax is exactly the same as strings in §3.1.5. The string literal is converted to and interpreted as a sequence of ASCII-encoded bytes, exactly how constants are (see §3.5.1).

Chapter 4

Project plan

The plan for this project was to execute the following basic steps.

1. *Determine objectives of the language* – the language was designed to be useful for extracting data quickly from binary files. It was modelled after an awk-style syntax, but was constrained to be simple enough to implement in a single semester.
2. *Design the language syntax* – the language syntax was modelled after awk and the microc language, and the concept was tested in parallel with writing a parser for it.
 - (a) *Write the scanner/parser* – This was critical to working through details of the language in parallel with designing it.
 - (b) *Design the AST types* – The abstract syntax tree (AST) types were designed to be as simple as possible, so that the translation to bytecode could be done easily.
3. *Write the logic to output the AST to GraphViz* – This was done mostly to help debug the AST structure, but also to make sure the AST was easy enough to traverse before attempting the bytecode translation. Translating the AST to GraphViz code enabled the creation of a visualization that could be generated for different inputs to make sure very early on that things like operator precedence were working correctly.
4. *Design the bytecode* – The plan was to take the bytecode from microc and add some basic instructions specific to the reading of binary files, such as an instruction to read a given number of bytes, branch conditionally on reaching end-of-file, etc. The stack-style design of the microc bytecode made it easier to implement than a register-based bytecode. One change I did intend to make was switching from relative branching to absolute branching.
5. *Implement the translation logic* – The translation logic would hopefully be a straightforward step once the previous were completed. The plan was to implement translation from AST structures into basic blocks, and implement the necessary bytecode instructions in parallel, while testing along the way.
6. *Project report* – The project report would be completed for required intermediate deadlines, with the remainder written after finishing the bulk of the coding.

7. *Test* – Finally, a conclusive round of testing would be done to ensure the language has no glaring bugs. Inefficiency and minor non-fatal flaws would be tolerated.

Development tools used included git and GitHub for version control, Make for generating both the bawk executable as well as class documents, and a few decent text editors. GraphViz was used for visualizing the AST, and \LaTeX was used to generate all the documents.

4.1 Project log

The project log, listed as git commits, is shown below. The full git repository is included in the project archive file.

Listing 4.1: Git log for project development

```

1 commit 9b2726855909c940c129ce02e3e81611238e5ed0
2 Author: Kevin Graney <nanonet@gmail.com>
3 Date:   Mon Aug 12 21:06:52 2013 -0400
4
5     Report work
6
7     Added the git log as a project log and finished up a few sections.
8
9     Makefile |      6 +-
10    plt_docs/gitlog.tex | 1138 ++++++
11    plt_docs/report-lessons.tex |    8 +-
12    plt_docs/report-plan.tex |    6 +-
13    plt_docs/report-test.tex |    1 -
14    5 files changed, 1155 insertions(+), 4 deletions(-)
15
16 commit 057a96571e6b99e2cfbf81126564acbeb4095bd8
17 Author: Kevin Graney <nanonet@gmail.com>
18 Date:   Mon Aug 12 20:17:55 2013 -0400
19
20     Report work
21
22    plt_docs/report-architecture.tex |    4 +-
23    plt_docs/report-introduction.tex |    1 +
24    plt_docs/report-lrm.tex |   16 +++-
25    plt_docs/report-plan.tex |    8 +-
26    plt_docs/report-tutorial.tex |   195 ++++++
27    plt_docs/report.tex |    7 +-
28    6 files changed, 224 insertions(+), 7 deletions(-)
29
30 commit 5d13044f761abc93542db2260c05d2c4087c3950
31 Author: Kevin Graney <nanonet@gmail.com>
32 Date:   Sun Aug 11 17:38:11 2013 -0400
33
34     Remove old test case
35
36    tests/simple.bawk |    6 -----
37    1 file changed, 6 deletions(-)
38
39 commit 6e9314853eef20ed8d3a856b0f786f1ba68fale1
40 Author: Kevin Graney <nanonet@gmail.com>

```

```

41 Date: Sun Aug 11 17:37:50 2013 -0400
42
43 Add test of int vs. uint types
44
45 tests/test-pat11.bawk | 8 ++++++++
46 tests/test-pat11.bawk.out | 2 ++
47 2 files changed, 10 insertions(+)
48
49 commit 5c58672c9c786b60cc336fe221d89d72318847f1
50 Author: Kevin Graney <nanonet@gmail.com>
51 Date: Sun Aug 11 17:30:25 2013 -0400
52
53 Report work
54
55 plt_docs/report-lrm.tex | 5 ++---
56 plt_docs/report-plan.tex | 27 ++++++++-----
57 2 files changed, 21 insertions(+), 11 deletions(-)
58
59 commit 820fc3560c2251d4421636ad13389da4666b1df3
60 Author: Kevin Graney <nanonet@gmail.com>
61 Date: Sun Aug 11 15:47:49 2013 -0400
62
63 Report work
64
65 plt_docs/report-architecture.tex | 5 +++-
66 plt_docs/report-lessons.tex | 8 ++++++-
67 2 files changed, 11 insertions(+), 2 deletions(-)
68
69 commit 9c4ec989cb2718b53c679b21e2521e53b137bff2
70 Author: Kevin Graney <nanonet@gmail.com>
71 Date: Sun Aug 11 15:47:38 2013 -0400
72
73 Line length enforcement for report
74
75 compile.ml | 3 +-
76 reader.mli | 16 ++++++-----
77 2 files changed, 11 insertions(+), 8 deletions(-)
78
79 commit 42d780a9e61394a6d229e5371b498a3c4bed9189
80 Author: Kevin Graney <nanonet@gmail.com>
81 Date: Sun Aug 11 14:56:07 2013 -0400
82
83 Work on final report
84
85 plt_docs/report-architecture.tex | 3 +++
86 plt_docs/report-plan.tex | 10 ++++++++
87 plt_docs/report-test.tex | 27 ++++++++
88 plt_docs/report.tex | 28 +++-----
89 4 files changed, 44 insertions(+), 24 deletions(-)
90
91 commit 03f47bf80f889638f51df5aefe62798dc5e715f1
92 Author: Kevin Graney <nanonet@gmail.com>
93 Date: Sun Aug 11 13:42:53 2013 -0400
94
95 Work on final report
96

```

```

97 Makefile | 2 +
98 plt_docs/columbia_university.eps | 6907 ++++++
99 plt_docs/report-architecture.tex | 47 +
100 plt_docs/report-lessons.tex | 9 +
101 plt_docs/report-lrm.tex | 16 +-
102 plt_docs/report.tex | 34 +-
103 6 files changed, 7007 insertions(+), 8 deletions(-)
104
105 commit b52b971c5ff914e23ea5b32b656fa837265ab03b
106 Author: Kevin Graney <nanonet@gmail.com>
107 Date: Sat Aug 10 14:46:23 2013 -0400
108
109 Fail on bad use of string literals
110
111 String literals can only be used in pattern expressions.
112
113 compile.ml | 4 +++
114 1 file changed, 4 insertions(+)
115
116 commit 7caaed970066afd42fccf9450191928f934db0a0
117 Author: Kevin Graney <nanonet@gmail.com>
118 Date: Sat Aug 10 14:46:04 2013 -0400
119
120 Add Beq instruction to bytecode interpreter
121
122 bytecode.ml | 3 +++
123 1 file changed, 3 insertions(+)
124
125 commit bb997604d2c1dacde12fbbd0633bb961e66e5949
126 Author: Kevin Graney <nanonet@gmail.com>
127 Date: Sat Aug 10 14:22:20 2013 -0400
128
129 Add if/else construct to AST GraphViz output
130
131 ast.ml | 6 +++++
132 1 file changed, 6 insertions(+)
133
134 commit 94dlb482f17b1f0bdaffcf0bc9f95b59a8ca47c6
135 Author: Kevin Graney <nanonet@gmail.com>
136 Date: Sat Aug 10 13:57:04 2013 -0400
137
138 Add test case for LE and a string constant
139
140 tests/test-pat10.bawk | 13 ++++++
141 tests/test-pat10.bawk.out | 4 +++
142 2 files changed, 17 insertions(+)
143
144 commit 09902de304da484f90227cbc8c4740f6807f7174
145 Author: Kevin Graney <nanonet@gmail.com>
146 Date: Sat Aug 10 13:11:19 2013 -0400
147
148 Add test case for LE
149
150 tests/test-pat9.bawk | 36 ++++++
151 tests/test-pat9.bawk.out | 10 +++++
152 2 files changed, 46 insertions(+)

```

```

153
154 commit eaa5733f996315216274479668d5267b01e38eed
155 Author: Kevin Graney <nanonet@gmail.com>
156 Date: Sat Aug 10 13:02:29 2013 -0400
157
158 Add support for little endian constants
159
160 Pattern constants (of hexadecimal integer form) are now supported
161 when LE!=0 to be interpreted as a little endian value.
162
163 compile.ml | 13 ++++++++--
164 1 file changed, 12 insertions(+), 1 deletion(-)
165
166 commit 1f92e256bc1632bc000b2a509f8bab3169808a90
167 Author: Kevin Graney <nanonet@gmail.com>
168 Date: Sat Aug 10 13:01:38 2013 -0400
169
170 Add support for previously bound pattern variables
171
172 Previously bound pattern variables can be used in future pattern
173 matches.
174
175 compile.ml | 18 ++++++++
176 parser.mly | 8 +++++--
177 tests/test-pat8.bawk | 10 ++++++++
178 tests/test-pat8.bawk.out | 4 ++++
179 4 files changed, 39 insertions(+), 1 deletion(-)
180
181 commit b03fc3blad8c2603175085bfbd8f14b7ec42b6d3
182 Author: Kevin Graney <nanonet@gmail.com>
183 Date: Sat Aug 10 12:12:40 2013 -0400
184
185 Test writing of RP
186
187 tests/test-pat7.bawk | 16 ++++++++
188 tests/test-pat7.bawk.out | 5 ++++
189 2 files changed, 21 insertions(+)
190
191 commit 8e005718abc113fdd2105bdd9ffbc70016482ba5
192 Author: Kevin Graney <nanonet@gmail.com>
193 Date: Sat Aug 10 10:55:02 2013 -0400
194
195 Add LE as a default variable
196
197 compile.ml | 12 ++++++----
198 1 file changed, 8 insertions(+), 4 deletions(-)
199
200 commit d2c07179ed76266c6aea84b8ed313ed31d9cddd4
201 Author: Kevin Graney <nanonet@gmail.com>
202 Date: Sat Aug 10 10:43:14 2013 -0400
203
204 Support assignment to RP
205
206 bytecode.ml | 3 +++
207 1 file changed, 3 insertions(+)
208

```

```

209 commit dd0b8d75c79828e90bb557b78738ae31e4b3f487
210 Author: Kevin Graney <nanonet@gmail.com>
211 Date: Sat Aug 10 10:38:10 2013 -0400
212
213 Add bind type names to reserved words
214
215 plt_docs/report-lrm.tex | 4 ++--
216 1 file changed, 2 insertions(+), 2 deletions(-)
217
218 commit b0f748857c35b232b3c0f7e6f3f1cd0021ba1c45
219 Author: Kevin Graney <nanonet@gmail.com>
220 Date: Sat Aug 10 10:37:01 2013 -0400
221
222 Change dynamic scoping of assignment
223
224 Changed the scoping of the assignment operator to be dynamic
225 and match the scoping of how variables are read.
226
227 compile.ml | 8 +++++---
228 tests/test-scope1.bawk.out | 2 +-
229 tests/test-scope2.bawk | 16 ++++++
230 tests/test-scope2.bawk.out | 5 +++++
231 4 files changed, 27 insertions(+), 4 deletions(-)
232
233 commit 22c9279c64a1d417a2fcee8b6fc8086b65de1e49
234 Author: Kevin Graney <nanonet@gmail.com>
235 Date: Sat Aug 10 10:36:07 2013 -0400
236
237 Add test case of if/else constructs
238
239 tests/test-cond1.bawk | 15 ++++++
240 tests/test-cond1.bawk.out | 2 ++
241 2 files changed, 17 insertions(+)
242
243 commit 49af898a754e376a4856817d9420b1e290345e1f
244 Author: Kevin Graney <nanonet@gmail.com>
245 Date: Thu Aug 8 19:51:14 2013 -0400
246
247 Change sample program to use uint types
248
249 tests/test-usefull.bawk | 8 +++-----
250 1 file changed, 3 insertions(+), 5 deletions(-)
251
252 commit e6bd2027d36914828584a5d4cca14479d2d1397c
253 Author: Kevin Graney <nanonet@gmail.com>
254 Date: Thu Aug 8 19:47:44 2013 -0400
255
256 Create final report skeleton
257
258 Created a skeleton of the final report based on the syllabus.
259
260 plt_docs/report-introduction.tex | 8 ++
261 plt_docs/report-lrm.tex | 220 ++++++
262 plt_docs/report.tex | 52 ++++++
263 3 files changed, 280 insertions(+)
264

```



```

265 commit 87e3d252b4ed368753869e6c7ea161b4728ef23d
266 Author: Kevin Graney <nanonet@gmail.com>
267 Date: Thu Aug 8 19:42:39 2013 -0400
268
269 Add if/else statement support
270
271 This will be required for supporting runtime resolution of RP
272 easily, so it seems trivial to just implement the if/else construct
273 in the language too.
274
275 ast_types.mli | 1 +
276 compile.ml | 18 ++++++-----
277 parser.mly | 8 ++++++-
278 scanner.mll | 3 +++
279 4 files changed, 23 insertions(+), 7 deletions(-)
280
281 commit 3cf57fd37d723fb04bdd7235ee61fa58989cdd3a
282 Author: Kevin Graney <nanonet@gmail.com>
283 Date: Sun Aug 4 15:57:47 2013 -0400
284
285 Remove unused variable
286
287 scanner.mll | 2 +-
288 1 file changed, 1 insertion(+), 1 deletion(-)
289
290 commit 4eea7d725f07304ab6c4c5f231632b55769f0b72
291 Author: Kevin Graney <nanonet@gmail.com>
292 Date: Sun Aug 4 15:50:55 2013 -0400
293
294 Add test case for bind types
295
296 tests/test-pat6.bawk | 23 ++++++
297 tests/test-pat6.bawk.out | 12 ++++++
298 2 files changed, 35 insertions(+)
299
300 commit 905b922d5078080bd0575690d148f68b501cb70f
301 Author: Kevin Graney <nanonet@gmail.com>
302 Date: Sun Aug 4 15:49:14 2013 -0400
303
304 Add uint1 type
305
306 ast.ml | 3 +++
307 ast_types.mli | 1 +
308 scanner.mll | 1 +
309 3 files changed, 5 insertions(+)
310
311 commit 24dd8c5b1b31fe60fb58003cb96b62ce33ecdc87
312 Author: Kevin Graney <nanonet@gmail.com>
313 Date: Sun Aug 4 15:27:56 2013 -0400
314
315 Add support for differentiating int/uint types
316
317 Added compiler support for treating integer types as signed. This
318 assumes that the interpretation stored in the file is twos
319 complement.
320

```

```

321 compile.ml |      7 +++++--
322 1 file changed, 5 insertions(+), 2 deletions(-)
323
324 commit 7a5b4836d5aea31cfef8db3d9a522400fc8bd10c
325 Author: Kevin Graney <nanonet@gmail.com>
326 Date:   Sun Aug 4 15:27:04 2013 -0400
327
328     Add instruction to convert to twos complement
329
330     Added a bytecode instruction to reinterpret the top of the stack
331     as a twos complement value with the given bit length.
332
333 bytecode.ml      |    10 ++++++++--
334 bytecode_types.mli |     1 +
335 utile.ml        |     6 ++++++
336 utile.mli       |     4 +++++
337 4 files changed, 20 insertions(+), 1 deletion(-)
338
339 commit 749e0bf50ad1a577adea51496b6466262f8a0274
340 Author: Kevin Graney <nanonet@gmail.com>
341 Date:   Sun Aug 4 15:25:37 2013 -0400
342
343     Add function to map binding types to their sizes
344
345 ast.ml |      7 ++++++
346 ast.mli |     4 +++++
347 2 files changed, 11 insertions(+)
348
349 commit 22658bec76ebd2d2cf0611d790582fcbf91d49ef
350 Author: Kevin Graney <nanonet@gmail.com>
351 Date:   Sun Aug 4 13:32:05 2013 -0400
352
353     Add pattern strings to AST GraphViz output
354
355 ast.ml |      7 +++++--
356 1 file changed, 5 insertions(+), 2 deletions(-)
357
358 commit ec915c6e62c8070c8b74f7f29f05ddb96af89d5d
359 Author: Kevin Graney <nanonet@gmail.com>
360 Date:   Sun Aug 4 13:23:56 2013 -0400
361
362     Add LaTeX file for final report
363
364 .gitignore      |     1 +
365 Makefile        |     2 +-
366 plt_docs/report.tex |  114 +++++
367 3 files changed, 116 insertions(+), 1 deletion(-)
368
369 commit 37be880f30d6d0e0839f76b2e7d995e310a81821
370 Author: Kevin Graney <nanonet@gmail.com>
371 Date:   Sun Aug 4 12:28:53 2013 -0400
372
373     Add test case for variable scoping
374
375     Show that variable scoping works inside of function calls
376

```

```

377 tests/test-func3.bawk      | 20 ++++++
378 tests/test-func3.bawk.out | 7 ++++++
379 2 files changed, 27 insertions(+)
380
381 commit bad13c86818b3be0eb077773a5326e0c67da40e0
382 Author: Kevin Graney <nanonet@gmail.com>
383 Date: Sat Aug 3 17:29:40 2013 -0400
384
385 Add support for function calls with arguments
386
387 Added support for function call arguments. Values are placed on
388 the stack at an offset from the frame pointer. The symbol table
389 contains values <-100 to represent these values differently from
390 global variables.
391
392 bytecode.ml      | 31 ++++++-----
393 bytecode_types.mli | 8 +----
394 compile.ml      | 19 ++++++----
395 3 files changed, 44 insertions(+), 14 deletions(-)
396
397 commit 96ec3a1b3a402498188f1ff289c6aee524d058ce
398 Author: Kevin Graney <nanonet@gmail.com>
399 Date: Sat Aug 3 17:28:20 2013 -0400
400
401 Add support for starting enumeration at a value
402
403 Added an optional argument to start the enumeration at a specific
404 value. The default remains 0.
405
406 utile.ml | 4 +--
407 utile.mli | 3 +--
408 2 files changed, 4 insertions(+), 3 deletions(-)
409
410 commit 5f3fa692057c26791269550c5f8f11da956dcb20
411 Author: Kevin Graney <nanonet@gmail.com>
412 Date: Fri Aug 2 20:55:34 2013 -0400
413
414 Make pattern matching work across entire file
415
416 Pattern matching now searches forward over the the entire file
417 instead of expecting a match at the current location of RP. This
418 is implemented in an incredibly inefficient way, but this should
419 be good enough for class.
420
421 bytecode.ml      | 6 +----
422 compile.ml      | 30 ++++++-----
423 tests/test-pat3.bawk | 2 +-
424 3 files changed, 28 insertions(+), 10 deletions(-)
425
426 commit 76e451ee8c908db8d8f1755cc239a4ac47230052
427 Author: Kevin Graney <nanonet@gmail.com>
428 Date: Fri Aug 2 20:53:27 2013 -0400
429
430 Add support for reading the RP variable
431
432 The variable is added as a special case in the bytecode interpreter.

```

```

433     Writing to this variable is not yet supported.
434
435 bytecode.ml |      3 +++
436 compile.ml  |      2 +-
437 2 files changed, 4 insertions(+), 1 deletion(-)
438
439 commit b6edc2975e4ee271ba7d8e0622f830c61ffefd59
440 Author: Kevin Graney <nanonet@gmail.com>
441 Date:   Fri Aug 2 20:51:32 2013 -0400
442
443     Added the Beo and Drp instructions to bytecode
444
445     Beo branches if the next byte reaches EOF and Drp drops a value
446     from the stack.
447
448 bytecode.ml      |      7 ++++++
449 bytecode_types.mli |      1 +
450 compile.ml       |      1 +
451 reader.ml        |     17 ++++++-----
452 reader.mli       |      4 +++-
453 5 files changed, 19 insertions(+), 11 deletions(-)
454
455 commit e794ccf368d43475081411e6b7b063aaa6274f1a
456 Author: Kevin Graney <nanonet@gmail.com>
457 Date:   Fri Aug 2 19:41:42 2013 -0400
458
459     Add test for scoping of variables
460
461 tests/test-scope1.bawk      |    15 ++++++
462 tests/test-scope1.bawk.out |      5 +++++
463 2 files changed, 20 insertions(+)
464
465 commit f5652c7d9dff21400ccb243a805d0148308afb87
466 Author: Kevin Graney <nanonet@gmail.com>
467 Date:   Fri Aug 2 19:37:59 2013 -0400
468
469     Merge symbol map into a single StringMap
470
471     Removed the function_map and variable_map objects and created a
472     single symbol_map StringMap to hold all the symbols. There still
473     exists a bindings map to contain the size of the binding variables.
474
475 bawk.ml      |      4 +-
476 compile.ml   |     60 ++++++-----
477 compile.mli  |     13 ++++++-----
478 3 files changed, 36 insertions(+), 41 deletions(-)
479
480 commit d238b3e53cb35d6285251e98d202ff6fafa6450b
481 Author: Kevin Graney <nanonet@gmail.com>
482 Date:   Fri Aug 2 18:34:12 2013 -0400
483
484     Add parent pointers to env
485
486     Added parent pointers to the environment records to support linked
487     lists. Create a simple resolution function for function calls so
488     things like print() (a built-in function) work from inside functions.

```

```

489
490 compile.ml | 22 ++++++-----
491 compile.mli | 1 +
492 2 files changed, 21 insertions(+), 2 deletions(-)
493
494 commit 2fa4bd4f81bbe1e9f4980029be7f7d03bd3a5ae0
495 Author: Kevin Graney <nanonet@gmail.com>
496 Date: Fri Aug 2 18:57:03 2013 -0400
497
498 Add Beq to label resolution
499
500 compile.ml | 1 +
501 1 file changed, 1 insertion(+)
502
503 commit dd63a8c9d9a101245eadef753d7812bf3fcf67c7
504 Author: Kevin Graney <nanonet@gmail.com>
505 Date: Fri Aug 2 13:48:32 2013 -0400
506
507 Added support for strings in pattern expressions
508
509 Added support for string literals inside of pattern expressions.
510 This allows the first useful test case to be added for finding
511 the width and height of a PNG file.
512
513 ast_types.mli | 1 +
514 compile.ml | 6 +++++-
515 parser.mly | 1 +
516 tests/test-usefull.bawk | 12 ++++++
517 tests/test-usefull.bawk.out | 2 ++
518 utile.ml | 3 +++
519 utile.mli | 4 ++++
520 7 files changed, 28 insertions(+), 1 deletion(-)
521
522 commit fdad37b643f96444bb64620717901e649dd514d1
523 Author: Kevin Graney <nanonet@gmail.com>
524 Date: Fri Aug 2 13:30:19 2013 -0400
525
526 Added explode and implode functions to Utile
527
528 These are functions defined in the OCaml docs as inefficient ways
529 to convert a string to a char list or a char list to a string
530 respectively.
531
532 utile.ml | 13 ++++++
533 utile.mli | 6 ++++++
534 2 files changed, 19 insertions(+)
535
536 commit 34e7d8295fcaf5f18103305609d5f583f95d3b38
537 Author: Kevin Graney <nanonet@gmail.com>
538 Date: Fri Aug 2 13:29:14 2013 -0400
539
540 Add boolean operators to bytecode interpreter.
541
542 These may not be of limited use since we don't have conditionals,
543 but they should work.
544

```

```

545 bytecode.ml |      7 ++++++
546 1 file changed, 7 insertions(+)
547
548 commit 4d96f677c8b0ce4c67b202f3b0fe886d0442eedf
549 Author: Kevin Graney <nanonet@gmail.com>
550 Date:   Fri Aug 2 13:00:15 2013 -0400
551
552     Remove error-prone branch
553
554     No need to have Bne instruction for binding pattern match, since
555     this will cause it to reject cases where the variable binds
556     to a value of zero.
557
558 compile.ml |      1 -
559 1 file changed, 1 deletion(-)
560
561 commit 79603df70433b8ed4b836236cc53f8d3b49c3531
562 Author: Kevin Graney <nanonet@gmail.com>
563 Date:   Fri Aug 2 13:00:03 2013 -0400
564
565     Add test case for int4 bind type
566
567 tests/test-pat5.bawk      |      5 +++++
568 tests/test-pat5.bawk.out |      1 +
569 2 files changed, 6 insertions(+)
570
571 commit 85c646aaf977d13568ecb04c4f741ef603f43937
572 Author: Kevin Graney <nanonet@gmail.com>
573 Date:   Fri Aug 2 12:13:38 2013 -0400
574
575     Add test case for int1 pattern binding
576
577 tests/test-pat4.bawk      |     12 ++++++
578 tests/test-pat4.bawk.out |      3 +++
579 2 files changed, 15 insertions(+)
580
581 commit 12cf98e1a573a16d463f25e2669b0f4642a7984a
582 Author: Kevin Graney <nanonet@gmail.com>
583 Date:   Wed Jul 31 20:25:40 2013 -0400
584
585     Initial cut at pattern variables
586
587     Added pattern variables.  The implementation is limited.  Only int1
588     type works at the moment, the namespace for variables is shared
589     with the globals, and the storage of the data is still messy.
590
591 bytecode.ml |      3 +++
592 compile.ml  |     31 ++++++
593 compile.mli |      2 +-
594 3 files changed, 32 insertions(+), 4 deletions(-)
595
596 commit c1efac5a1334654490cfbe2472498d633c8d31d8
597 Author: Kevin Graney <nanonet@gmail.com>
598 Date:   Sun Jul 28 22:13:01 2013 -0400
599
600     Add simple global variables

```

```

601
602   Added simple global variables.  These are implemented as an array
603   in the bytecode interpreter.  The compiler assigns each name a
604   unique entry in the array.  Currently all variables are considered
605   global.  Variables are not being set at an offset of FP on the stack
606   because they are dynamically created and not required to be declared.
607
608   bytecode.ml           |    9 ++++++++
609   compile.ml           |   33 ++++++++
610   compile.mli          |    7 ++++++++
611   tests/test-var1.bawk |    5 ++++++
612   tests/test-var1.bawk.out |    2 ++
613   5 files changed, 56 insertions(+)
614
615   commit 2ee86a40da8b5aed8e0e7844c1b4ccb7a3ef9d39
616   Author: Kevin Graney <nanonet@gmail.com>
617   Date:   Sat Jul 27 21:02:36 2013 -0400
618
619       Add basic function call test cases
620
621       Added test cases for function calls with no arguments.
622
623   tests/test-func1.bawk |   10 ++++++++
624   tests/test-func1.bawk.out |    5 ++++++
625   tests/test-func2.bawk |    7 ++++++++
626   tests/test-func2.bawk.out |    2 ++
627   4 files changed, 24 insertions(+)
628
629   commit d79330e9a40633c38a61a7760254a9a3756689f0
630   Author: Kevin Graney <nanonet@gmail.com>
631   Date:   Sat Jul 27 20:51:53 2013 -0400
632
633       Support function calls by name
634
635       Added support for function calls by name.  Still need to decide
636       how exactly the 'env' record is going to work.
637
638   bytecode.ml           |    6 ++++++
639   bytecode_types.mli   |    2 +-
640   compile.ml           |   15 ++++++++----
641   3 files changed, 18 insertions(+), 5 deletions(-)
642
643   commit bd94bb11eb0546ff2633e152df2e61759144d8ab
644   Author: Kevin Graney <nanonet@gmail.com>
645   Date:   Sat Jul 27 19:12:03 2013 -0400
646
647       Implement compilation of function declarations
648
649       Each function declaration is preceded by a Bra call to skip over
650       it in case the PC approaches it from something prior falling
651       through it (e.g. at the start of the program).  This is a quicker
652       hack than actually reobjecting all of them after the Hlt instruction
653       or something.
654
655   bytecode.ml           |    2 ++
656   bytecode_types.mli   |    9 +++++----

```

```

657 compile.ml | 8 ++++++-
658 3 files changed, 14 insertions(+), 5 deletions(-)
659
660 commit 32a42f8b1cf9abc0f90ab10ea98aaee9e83aa4f5
661 Author: Kevin Graney <nanonet@gmail.com>
662 Date: Sat Jul 27 17:35:04 2013 -0400
663
664 Don't compile code when printing out the AST
665
666 The AST is useful for debugging, often of code that doesn't compile.
667 Therefore, it is silly to compile the code (and abort) before
668 printing the AST.
669
670 bawk.ml | 7 ++++---
671 1 file changed, 4 insertions(+), 3 deletions(-)
672
673 commit 2b25b94e38e0de1f6dbe991c31460cffe3f04186
674 Author: Kevin Graney <nanonet@gmail.com>
675 Date: Sat Jul 27 17:33:46 2013 -0400
676
677 Adding function declarations to parser
678
679 To declare functions the syntax is changed so a "def" keyword must
680 precede the function identifier. This eliminates reduce/reduce
681 conflicts between expressions and the function declarations.
682
683 ast.ml | 19 ++++++-----
684 ast_types.mli | 8 ++++++-
685 parser.mly | 16 ++++++-----
686 scanner.mll | 1 +
687 4 files changed, 41 insertions(+), 3 deletions(-)
688
689 commit 821a574a5238d228136cb885a78e5bae6606f010
690 Author: Kevin Graney <nanonet@gmail.com>
691 Date: Sat Jul 27 16:06:49 2013 -0400
692
693 Fix bug in pattern match compilation
694
695 The file wasn't seeking back to the pre-pattern match location when
696 the pattern did not match.
697
698 compile.ml | 3 +--
699 1 file changed, 1 insertion(+), 2 deletions(-)
700
701 commit a36cc2c7e1545108511263a18654f9dfcce6643a
702 Author: Kevin Graney <nanonet@gmail.com>
703 Date: Sat Jul 27 16:06:12 2013 -0400
704
705 Add test cases for pattern matching
706
707 run_tests.sh | 2 +-
708 tests/lichtenstein.png | Bin 0 -> 365198 bytes
709 tests/test-pat1.bawk | 9 +++++++
710 tests/test-pat1.bawk.out | 2 ++
711 tests/test-pat2.bawk | 9 +++++++
712 tests/test-pat2.bawk.out | 3 +++

```



```

713 tests/test-pat3.bawk      |   13 ++++++++
714 tests/test-pat3.bawk.out |    3 +++
715 8 files changed, 40 insertions(+), 1 deletion(-)
716
717 commit 8314d555d1c85be1bc9113a7d9fc6bd53ec5357a
718 Author: Kevin Graney <nanonet@gmail.com>
719 Date:   Sat Jul 27 16:04:21 2013 -0400
720
721     Fix interpreter bug in Bne
722
723     Top of stack (compare-to value) was not being popped.
724
725 bytecode.ml |    2 +-
726 1 file changed, 1 insertion(+), 1 deletion(-)
727
728 commit abfecc6c182861783762343b579e344382912429
729 Author: Kevin Graney <nanonet@gmail.com>
730 Date:   Sat Jul 27 15:50:13 2013 -0400
731
732     Add instructions to store file position on stack
733
734     Added instructions to store the file position on the stack and
735     read a file position to seek to from the stack. This enables the
736     bytecode interpreter to return to the position in the file at
737     the conclusion of a pattern statement.
738
739 bytecode.ml      |   12 ++++++++
740 bytecode_types.mli |    2 ++
741 compile.ml      |    2 ++
742 reader.ml       |    6 +++++
743 reader.mli      |    6 +++++
744 5 files changed, 27 insertions(+), 1 deletion(-)
745
746 commit 3753438ca40393e69f64184a9f1795224c8d2d18
747 Author: Kevin Graney <nanonet@gmail.com>
748 Date:   Sat Jul 27 15:49:06 2013 -0400
749
750     Fix error in built-in print function
751
752     The built-in print function didn't pop it's argument from the stack.
753
754 bytecode.ml |    2 +-
755 1 file changed, 1 insertion(+), 1 deletion(-)
756
757 commit 0e87833401824d6d87596a7f107e53d9a837a785
758 Author: Kevin Graney <nanonet@gmail.com>
759 Date:   Sat Jul 27 15:28:56 2013 -0400
760
761     Add branch rewriting on second pass
762
763 bytecode.ml |    9 ++++++++
764 bytecode.mli |    7 ++++++
765 compile.ml  |   27 ++++++++
766 3 files changed, 42 insertions(+), 1 deletion(-)
767
768 commit 2bf1ec02de8c3fe35d3b154b68d56eb7b127ae45

```

```

769 Author: Kevin Graney <nanonet@gmail.com>
770 Date: Sat Jul 27 15:19:07 2013 -0400
771
772 Add Reader module and support for branching
773
774 Added the Reader module and support for resolving branch addresses
775 correctly using the 'Label' pseudo instruction and two passes.
776
777 Makefile | 1 +
778 bytecode.ml | 4 ++++
779 reader.ml | 50 ++++++
780 reader.mli | 26 ++++++
781 4 files changed, 81 insertions(+)
782
783 commit df48722be79dd90290f9ff2bf6e394c62dac47cc
784 Author: Kevin Graney <nanonet@gmail.com>
785 Date: Sat Jul 27 15:16:20 2013 -0400
786
787 Don't compile inside Bytecode module
788
789 This removes the dependency of the Bytecode module on the
790 Compile module. OCaml appears somewhat picky about enforcing
791 non-cyclic dependencies.
792
793 bawk.ml | 5 +++--
794 bytecode.ml | 8 +++-----
795 bytecode.mli | 4 +++--
796 3 files changed, 8 insertions(+), 9 deletions(-)
797
798 commit 6092ae879e5ed00e52f74100f2435643a21f66ae
799 Author: Kevin Graney <nanonet@gmail.com>
800 Date: Sat Jul 27 14:56:06 2013 -0400
801
802 Use ocamldep to find linking order
803
804 Use 'ocamldep -sort' to determine the order the .cmo files get
805 linked.
806
807 Makefile | 2 +-
808 1 file changed, 1 insertion(+), 1 deletion(-)
809
810 commit ad4836c8a8e890289a5d77a2998242ff560f05e0
811 Author: Kevin Graney <nanonet@gmail.com>
812 Date: Sat Jul 27 14:45:42 2013 -0400
813
814 Allow enumerate to do some primitive introspection
815
816 Changed the step function for enumerate to allow it to behave
817 differently based on the value of the head of the content list.
818 This is useful for enumerating bytecode with pseudo instructions,
819 where we want to no give the pseudo instructions their own address.
820
821 compile.ml | 2 +-
822 utile.ml | 4 +++--
823 utile.mli | 2 +-
824 3 files changed, 4 insertions(+), 4 deletions(-)

```

```

825
826 commit 3e153a4adf5f5be93ffa7a127487573a38b51a3b
827 Author: Kevin Graney <nanonet@gmail.com>
828 Date: Sat Jul 27 14:19:55 2013 -0400
829
830 Add bytecode support for branches
831
832 Added the Label pseudo-instruction to the bytecode to allow the
833 branching addresses to be resolved. Basically labels are entered
834 by the one stage of the compilation and then resolved to actual
835 addresses by a later stage.
836
837 bytecode.ml | 5 +++++
838 bytecode_types.mli | 4 +++-
839 compile.ml | 15 ++++++++-----
840 3 files changed, 19 insertions(+), 5 deletions(-)
841
842 commit 41b603c93a75d06b673989533ed4541735910d63
843 Author: Kevin Graney <nanonet@gmail.com>
844 Date: Sat Jul 27 13:45:43 2013 -0400
845
846 Support passing a binary filename in argv
847
848 Added support for reading argv[2] as the filename of a binary file
849 to process. Argument parsing is still very primitive, but this
850 is enough to build something that works.
851
852 bawk.ml | 10 ++++++++--
853 bytecode.ml | 6 +++---
854 bytecode.mli | 4 +++-
855 3 files changed, 14 insertions(+), 6 deletions(-)
856
857 commit 747889f4963cf7c90a2dlf6219c5b8ccf59b7663
858 Author: Kevin Graney <nanonet@gmail.com>
859 Date: Thu Jul 25 22:00:04 2013 -0400
860
861 Begin implementing pattern matching compilation
862
863 A start on the implementation of compiling pattern expressions into
864 bytecode.
865
866 compile.ml | 19 ++++++++-----
867 1 file changed, 18 insertions(+), 1 deletion(-)
868
869 commit a74c6bd36d904a37206f7e2b5dc6bfff9c74ca91
870 Author: Kevin Graney <nanonet@gmail.com>
871 Date: Thu Jul 25 21:59:37 2013 -0400
872
873 Add more instructions to string_of_instruction
874
875 bytecode.ml | 4 +++-
876 1 file changed, 3 insertions(+), 1 deletion(-)
877
878 commit 368e05b6506f7ea33dfc16c5fb91694195c0ca92
879 Author: Kevin Graney <nanonet@gmail.com>
880 Date: Thu Jul 25 21:56:25 2013 -0400

```

```

881
882   Add size_of_bind_type function
883
884   Added a function to return the size in bytes of a given bind_type.
885
886   ast.ml |      7 ++++++
887   ast.mli |     4 ++++
888   2 files changed, 11 insertions(+)
889
890   commit 3057ec54d1eedc46d4fcf6a549c92c05fb0350a0
891   Author: Kevin Graney <nanonet@gmail.com>
892   Date:   Thu Jul 25 21:33:15 2013 -0400
893
894   Reverse parsing order of pattern tokens
895
896   Pattern tokens were stored in the AST in the reverse of the desired
897   order. This switches the order.
898
899   parser.mly |    5 ++++-
900   1 file changed, 4 insertions(+), 1 deletion(-)
901
902   commit bfc73199e14bf37ca708d6c2e47258bableaf02d
903   Author: Kevin Graney <nanonet@gmail.com>
904   Date:   Thu Jul 25 21:29:51 2013 -0400
905
906   Add PC values to bytecode output
907
908   Added value of program counter address to instructions when printing
909   the bytecode output.
910
911   bytecode.ml |    4 +++-
912   1 file changed, 3 insertions(+), 1 deletion(-)
913
914   commit ede71a47f4d8fa8e1f2b40893044358bd9046007
915   Author: Kevin Graney <nanonet@gmail.com>
916   Date:   Thu Jul 25 21:18:37 2013 -0400
917
918   Add parsing of hex constants in pattern expression
919
920   Added parsing of hexadecimal constants in the pattern expression,
921   as defined in the LRM. Big/little endian conversion should
922   probably be handled by the bytecode at runtime.
923
924   parser_help.ml |  18 ++++++
925   1 file changed, 17 insertions(+), 1 deletion(-)
926
927   commit f8dbc43c280dc1264938e567664ee38994e95368
928   Author: Kevin Graney <nanonet@gmail.com>
929   Date:   Wed Jul 24 20:03:51 2013 -0400
930
931   Add regression tests
932
933   Started the regression test suite with a very simple test case.
934
935   run_tests.sh | 123 ++++++
936   tests/test-arith.bawk | 1 +

```

```

937 tests/test-arith.bawk.out | 1 +
938 3 files changed, 125 insertions(+)
939
940 commit 35794eb2cc8055daffff965986bc070b5868e1142
941 Author: Kevin Graney <nanonet@gmail.com>
942 Date: Wed Jul 24 19:40:07 2013 -0400
943
944 Halt at conclusion of program
945
946 Added the Hlt instruction to the end of compiled programs, and
947 implemented halt functionality in the interpreter. Previously
948 we just kept reading through to the next instruction in memory.
949
950 bytecode.ml | 2 ++
951 compile.ml | 2 +-
952 2 files changed, 3 insertions(+), 1 deletion(-)
953
954 commit 1b3bb14eb3a13beede68a3c66d9e389f6ca2d63e
955 Author: Kevin Graney <nanonet@gmail.com>
956 Date: Tue Jul 23 20:30:03 2013 -0400
957
958 Add basic arithmetic operators to interpreter
959
960 Added basic arithmetic instructions to bytecode interpreter. It's
961 still missing a lot of things, but this is a start.
962
963 bytecode.ml | 23 ++++++-----
964 bytecode.mli | 6 +++++-
965 2 files changed, 27 insertions(+), 2 deletions(-)
966
967 commit f9aaa15acb4fee676f7ba5bc4468168991d6d5bc
968 Author: Kevin Graney <nanonet@gmail.com>
969 Date: Tue Jul 23 19:31:00 2013 -0400
970
971 Add stub for bytecode execution flag
972
973 Added the -e flag to compile a program into bytecode and then
974 execute that bytecode.
975
976 bawk.ml | 7 ++++---
977 bytecode.ml | 4 ++++
978 bytecode.mli | 3 +++
979 3 files changed, 11 insertions(+), 3 deletions(-)
980
981 commit 878dd0aa65efeea16fladece43e9b8b1683b1c80
982 Author: Kevin Graney <nanonet@gmail.com>
983 Date: Mon Jul 22 23:03:40 2013 -0400
984
985 Fix bug where built-in function addresses increase
986
987 Fixed a bug where built-in function addresses increased as the
988 name list is read instead of decreased. (Negative addresses are
989 used to identify built-in functions.)
990
991 compile.ml | 3 +-
992 1 file changed, 2 insertions(+), 1 deletion(-)

```

```

993
994 commit 51109abddc9345b4147b6771c0f2c3c0eacff4b0
995 Author: Kevin Graney <nanonet@gmail.com>
996 Date: Mon Jul 22 23:02:41 2013 -0400
997
998 Add step function to enumerate
999
1000 Added an optional argument to enumerate to allow a step function
1001 to be specified.
1002
1003 utile.ml | 4 ++--
1004 utile.mli | 2 +-
1005 2 files changed, 3 insertions(+), 3 deletions(-)
1006
1007 commit c5f9a558bbec35da21b91546ebc287afc0b4b510
1008 Author: Kevin Graney <nanonet@gmail.com>
1009 Date: Mon Jul 22 22:53:57 2013 -0400
1010
1011 Populate starting environment with built-ins
1012
1013 Populated the default environment with built-in functions. They
1014 are read from an array, enumerated.
1015
1016 bawk.ml | 13 ++++++++
1017 compile.ml | 14 ++++++-----
1018 compile.mli | 38 ++++++
1019 3 files changed, 60 insertions(+), 5 deletions(-)
1020
1021 commit b970704fb337247656f959c79fc90415f68566f5
1022 Author: Kevin Graney <nanonet@gmail.com>
1023 Date: Mon Jul 22 22:53:13 2013 -0400
1024
1025 Add enumerate function
1026
1027 Added enumerate function to Utile library.
1028
1029 utile.ml | 5 +++++
1030 utile.mli | 6 +++++-
1031 2 files changed, 10 insertions(+), 1 deletion(-)
1032
1033 commit 1b added 68335ca43bafa27186611dc21a3edca399f
1034 Author: Kevin Graney <nanonet@gmail.com>
1035 Date: Mon Jul 22 21:42:09 2013 -0400
1036
1037 Add function call stubs
1038
1039 Added preliminary logic for generating bytecode for function
1040 calls. Still stub code. Also added basic environment to support
1041 things like function names (and eventually scoping).
1042
1043 bytecode.ml | 3 ++-
1044 compile.ml | 35 ++++++
1045 compile.mli | 4 ++--
1046 3 files changed, 33 insertions(+), 9 deletions(-)
1047
1048 commit 19dc1c220040c645ab88bcbc946062360e2a53e0

```

```

1049 Author: Kevin Graney <nanonet@gmail.com>
1050 Date:   Fri Jul 19 09:50:16 2013 -0400
1051
1052     Ignore autogenerated makefiles
1053
1054     .gitignore |      1 +
1055     1 file changed, 1 insertion(+)
1056
1057 commit 5fed5f231c68a46dc7208fbb4ae8e267d0ebf436
1058 Author: Kevin Graney <nanonet@gmail.com>
1059 Date:   Fri Jul 19 09:49:50 2013 -0400
1060
1061     Add more stub code to the translator
1062
1063     bytecode.ml          |      1 +
1064     bytecode_types.mli  |      3 ++-
1065     compile.ml          |      6 ++++-
1066     3 files changed, 8 insertions(+), 2 deletions(-)
1067
1068 commit 74d7f39a9d0a6f80c5e207947fc9f1e009a9a440
1069 Author: Kevin Graney <nanonet@gmail.com>
1070 Date:   Wed Jul 17 22:17:26 2013 -0400
1071
1072     Remove stub value and replace with actual.
1073
1074     compile.ml |      2 +-
1075     1 file changed, 1 insertion(+), 1 deletion(-)
1076
1077 commit 8468b5b5886bdf32c26d618185ca71ab7de4bfbf
1078 Author: Kevin Graney <nanonet@gmail.com>
1079 Date:   Wed Jul 17 22:14:36 2013 -0400
1080
1081     Add bytecode creation support for binary operators.
1082
1083     Added bytecode generation for binary operators.
1084
1085     bytecode.ml          |      1 +
1086     compile.ml          |      4 +++-
1087     tests/simple.bawk  |      2 +-
1088     3 files changed, 5 insertions(+), 2 deletions(-)
1089
1090 commit f9939d8c691397880f4a5e47ede0eed287601714
1091 Author: Kevin Graney <nanonet@gmail.com>
1092 Date:   Wed Jul 17 22:12:20 2013 -0400
1093
1094     Change string_of_operator to produce bytecode
1095
1096     Changed string_of_operator from producing human readable strings
1097     to bytecode operator strings.  This benefits the bytecode generation
1098     and makes the AST graph more concise.
1099
1100     ast.ml |      20 ++++++++-----
1101     ast.mli |       4 ++++
1102     2 files changed, 14 insertions(+), 10 deletions(-)
1103
1104 commit d8f9b1658ca62d6920c2a69e9d7535f420b1c62f

```

```

1105 Author: Kevin Graney <nanonet@gmail.com>
1106 Date: Wed Jul 17 21:52:47 2013 -0400
1107
1108 Add primitive bytecode creation
1109
1110 Added Bytecode_types module to mimic pattern used for Ast_types.
1111 Also got simple translation working, so the -c flag now produces
1112 the bytecode output. Currently only works for primitive programs.
1113
1114 Makefile | 2 +-
1115 bytecode.ml | 24 ++++++-----
1116 bytecode.mli | 10 ++++++---
1117 bytecode_types.mli | 16 ++++++-----
1118 compile.ml | 18 ++++++-----
1119 compile.mli | 5 +++++
1120 tests/simple.bawk | 6 +++++
1121 7 files changed, 55 insertions(+), 26 deletions(-)
1122
1123 commit 3277f84748d4b697c71bc01b8fee79b6a9a91200
1124 Author: Kevin Graney <nanonet@gmail.com>
1125 Date: Wed Jul 17 21:50:48 2013 -0400
1126
1127 Produce ocaml doc by default
1128
1129 Changed Makefile to produce design_docs target by default. (This
1130 is the ocaml doc output.)
1131
1132 Makefile | 2 +-
1133 1 file changed, 1 insertion(+), 1 deletion(-)
1134
1135 commit c269dc76c098f6f880e395abf334a6262a7cae20
1136 Author: Kevin Graney <nanonet@gmail.com>
1137 Date: Wed Jul 17 19:34:53 2013 -0400
1138
1139 Automatic dependency generation
1140
1141 Makefile now generates dependencies from ocamldep automatically. It's
1142 a little ugly in that it makes the dependencies regardless of what the
1143 target is (including clean) but this is still better than manually
1144 updating the Makefile all the time.
1145
1146 Makefile | 47 ++++++-----
1147 1 file changed, 19 insertions(+), 28 deletions(-)
1148
1149 commit 21e268ef44da666feaeabcd77926a3884f25f87f
1150 Author: Kevin Graney <nanonet@gmail.com>
1151 Date: Wed Jul 17 19:15:14 2013 -0400
1152
1153 Add bytecode and compile modules
1154
1155 Add primitive bytecode and compile modules.
1156
1157 Makefile | 11 ++++++---
1158 bawk.ml | 2 +-
1159 bytecode.ml | 21 ++++++-----
1160 bytecode.mli | 2 ++

```



```

1161 compile.ml | 10 ++++++++
1162 5 files changed, 42 insertions(+), 4 deletions(-)
1163
1164 commit 75444530c89124b5ba77072bacbfd4d0a4be50db
1165 Author: Kevin Graney <nanonet@gmail.com>
1166 Date: Sun Jul 14 14:38:29 2013 -0400
1167
1168 Add support for command line arguments
1169
1170 Added primitive support for command line arguments to distinguish
1171 between actions such as compile, output the AST, output the Bytecode,
1172 etc.
1173
1174 bawk.ml | 35 ++++++++-----
1175 1 file changed, 23 insertions(+), 12 deletions(-)
1176
1177 commit 7c73eafb8a0d38620081c90e417c3110561c7302
1178 Author: Kevin Graney <nanonet@gmail.com>
1179 Date: Sun Jul 14 14:01:17 2013 -0400
1180
1181 Expand definition of bind patterns
1182
1183 Work on LRM to expand the definition of bind patterns.
1184
1185 plt_docs/lrm.tex | 20 ++++++++--
1186 1 file changed, 18 insertions(+), 2 deletions(-)
1187
1188 commit 244a3faaee0a96b850562e183fc94103800e6598
1189 Author: Kevin Graney <nanonet@gmail.com>
1190 Date: Thu Jun 27 14:39:18 2013 -0400
1191
1192 LRM work
1193
1194 plt_docs/lrm.tex | 149 ++++++++-----
1195 1 file changed, 108 insertions(+), 41 deletions(-)
1196
1197 commit c52dd031082c4562d9030a5784bc19ec1994fa94
1198 Author: Kevin Graney <nanonet@gmail.com>
1199 Date: Thu Jun 27 10:43:05 2013 -0400
1200
1201 Add assignment operator
1202
1203 Added right-associative assignment operator.
1204
1205 ast.ml | 4 ++++
1206 ast_types.mli | 1 +
1207 parser.mly | 2 ++
1208 tests/prototype.bawk | 6 ++++++
1209 4 files changed, 13 insertions(+)
1210
1211 commit 0dc571931abd863381854157cd1a7c67d17efa98
1212 Author: Kevin Graney <nanonet@gmail.com>
1213 Date: Thu Jun 27 10:08:00 2013 -0400
1214
1215 Add preliminary pattern expr constant parsing
1216

```

```

1217   Added preliminary code to differentiate between integer constants
1218   in pattern expressions and those in plain expressions so they can
1219   be treated differently. Pattern expression constants will always
1220   be interpreted in hexadecimal as a series of bytes, and plain
1221   expression constants will be more C-like in their interpretation.
1222
1223   Makefile          |   16 ++++++++-----
1224   ast.ml            |   11 ++++++-----
1225   ast_types.mli     |    3 +-+
1226   parser.mly        |   10 ++++++-----
1227   parser_help.ml    |    5 +++++
1228   parser_help.mli   |    6 +++++
1229   scanner.mll       |    6 +++---
1230   7 files changed, 38 insertions(+), 19 deletions(-)
1231
1232   commit 03c361378ae507fae887147f1b22f85e51776cff
1233   Author: Kevin Graney <nanonet@gmail.com>
1234   Date:   Thu Jun 27 09:31:13 2013 -0400
1235
1236       Add Sublime Text 2 project file
1237
1238       Also update .gitignore to ignore other Sublime Text files.
1239
1240   .gitignore        |    2 ++
1241   bawk.sublime-project |   21 ++++++
1242   2 files changed, 23 insertions(+)
1243
1244   commit a345dd50a587d0b70d2709cf64a0908521ad2c4a
1245   Author: Kevin Graney <nanonet@gmail.com>
1246   Date:   Thu Jun 27 09:16:11 2013 -0400
1247
1248       Remove recursive makefiles
1249
1250       Changing to a flat structure for simplicity.
1251
1252   Makefile          |    9 ++++++--
1253   plt_docs/Makefile |    9 -----
1254   2 files changed, 8 insertions(+), 10 deletions(-)
1255
1256   commit b950f353c6c05e0c8400cf63c0114da0b9d0142a
1257   Author: Kevin Graney <nanonet@gmail.com>
1258   Date:   Tue Jun 25 22:34:34 2013 -0400
1259
1260       Work on LRM
1261
1262   plt_docs/lrm.tex |   16 ++++++
1263   1 file changed, 13 insertions(+), 3 deletions(-)
1264
1265   commit 0f3edccfc005cc678f3e12d4a201f2f926a86ead
1266   Author: Kevin Graney <nanonet@gmail.com>
1267   Date:   Tue Jun 25 20:17:15 2013 -0400
1268
1269       Work on LRM
1270
1271   plt_docs/lrm.tex |   70 ++++++
1272   1 file changed, 42 insertions(+), 28 deletions(-)

```

```

1273
1274 commit a05a96b731bd0492b2efe2cbfe709b7fcd72786e
1275 Author: Kevin Graney <nanonet@gmail.com>
1276 Date: Mon Jun 24 22:46:32 2013 -0400
1277
1278     Added to prototype.bawk test cases
1279
1280 tests/prototype.bawk |    6 ++++++
1281 1 file changed, 6 insertions(+)
1282
1283 commit 9f52f5177a5b24b421b5e63fab59532a772e5278
1284 Author: Kevin Graney <nanonet@gmail.com>
1285 Date: Mon Jun 24 22:46:09 2013 -0400
1286
1287     LRM work
1288
1289 plt_docs/lrm.tex |   45 ++++++-----
1290 1 file changed, 39 insertions(+), 6 deletions(-)
1291
1292 commit 01cbe8ea67897d83343bc32614160a99d8abe291
1293 Author: Kevin Graney <nanonet@gmail.com>
1294 Date: Mon Jun 24 22:45:16 2013 -0400
1295
1296     Add string literals
1297
1298     Added lexer/parser support for string literals. Parser support
1299     may need to be firmed up still (assignment to a string literal
1300     should not be allowed, so we need some kind of lvalue/rvalue
1301     distinction).
1302
1303 ast.ml           |    3 +++
1304 ast_types.mli   |    1 +
1305 parser.mly      |    2 ++
1306 scanner.mll     |   72 ++++++-----
1307 4 files changed, 77 insertions(+), 1 deletion(-)
1308
1309 commit 840aa273da4bf579668a4cfdfae7ef5c40e0d825
1310 Author: Kevin Graney <nanonet@gmail.com>
1311 Date: Sun Jun 23 09:47:36 2013 -0400
1312
1313     Fix warning about non-unit type
1314
1315     Changed to use ignore in the print out of the GraphViz DOT file
1316     where a function with a side effect is called that returns a value
1317     we don't care about.
1318
1319 ast.ml |    2 +-
1320 1 file changed, 1 insertion(+), 1 deletion(-)
1321
1322 commit 16b81c4478d499ff5d68e9baf0fb2a79762df0af
1323 Author: Kevin Graney <nanonet@gmail.com>
1324 Date: Sun Jun 23 09:46:17 2013 -0400
1325
1326     Add precedence and associativity for comparators
1327
1328     Resolves a bunch of shift/reduce conflicts in the parser generator.

```

```

1329     The precedence ordering is taken from the microc example.
1330
1331 parser.mly |      2 ++
1332 1 file changed, 2 insertions(+)
1333
1334 commit f8eb3945c676eb2bdd89424824854d6ec753173c
1335 Author: Kevin Graney <nanonet@gmail.com>
1336 Date:   Wed Jun 19 21:13:49 2013 -0400
1337
1338     Initial test case
1339
1340     A simple test case to test several different parts of the parser.
1341     Useful for looking at the AST GraphViz DOT output of.
1342
1343 tests/prototype.bawk |  23 ++++++
1344 1 file changed, 23 insertions(+)
1345
1346 commit 9a51cf106b5cfa45da43624a578ae26b9454d415
1347 Author: Kevin Graney <nanonet@gmail.com>
1348 Date:   Wed Jun 19 21:13:30 2013 -0400
1349
1350     Work on language reference manual
1351
1352 plt_docs/lrm.tex |    92 ++++++-----
1353 1 file changed, 86 insertions(+), 6 deletions(-)
1354
1355 commit a173295ba72fd0e3b8561a1f144af84f28ec5616
1356 Author: Kevin Graney <nanonet@gmail.com>
1357 Date:   Wed Jun 19 21:11:59 2013 -0400
1358
1359     Let identifiers start with underscore
1360
1361 scanner.mll |      2 +-
1362 1 file changed, 1 insertion(+), 1 deletion(-)
1363
1364 commit 7270f30448c43458653804ee80195b9145eecb92
1365 Author: Kevin Graney <nanonet@gmail.com>
1366 Date:   Sun Jun 16 19:58:26 2013 -0400
1367
1368     Add comparator operators
1369
1370     Added equality comparison along with >, <, >=, <=, !=.
1371
1372 ast.ml          |      6 ++++++
1373 ast_types.mli  |      2 ++
1374 parser.mly     |      7 ++++++
1375 3 files changed, 15 insertions(+)
1376
1377 commit 93514726fb63d927360c2b0682a9648b276ad526
1378 Author: Kevin Graney <nanonet@gmail.com>
1379 Date:   Sat Jun 15 23:34:11 2013 -0400
1380
1381     Allow function calls with no arguments
1382
1383 parser.mly |      1 +
1384 1 file changed, 1 insertion(+)

```

```

1385
1386 commit 42dlfff15e1ed4a5223c6bcf03c95d51077b5c1a5
1387 Author: Kevin Graney <nanonet@gmail.com>
1388 Date: Sat Jun 15 23:27:02 2013 -0400
1389
1390 Make edge insertion consistent
1391
1392 Some of the edge insertion for the GraphViz DOT generation of the
1393 AST was inconsistent. Some functions assumed the caller made the
1394 edge and some assumed the callee, resulting in a few double or
1395 missing edges. Now the caller is always required to create a
1396 node and a link to its parent.
1397
1398 ast.ml | 16 ++++++-----
1399 1 file changed, 7 insertions(+), 9 deletions(-)
1400
1401 commit 5ce7f0fc576ff7b9a272c2363237ae9d49a4c0c3
1402 Author: Kevin Graney <nanonet@gmail.com>
1403 Date: Sat Jun 15 23:11:21 2013 -0400
1404
1405 Add function calls to AST
1406
1407 Added parsing of function calls to the AST. Identified an
1408 inconsistency in how parent edges are drawn in the GraphViz output,
1409 so that bug needs to be fixed in the future.
1410
1411 ast.ml | 7 ++++++
1412 ast_types.mli | 1 +
1413 parser.mly | 10 ++++++++
1414 scanner.mll | 1 +
1415 4 files changed, 18 insertions(+), 1 deletion(-)
1416
1417 commit 386c4d95c38f1975579fce7e6152e568960da326
1418 Author: Kevin Graney <nanonet@gmail.com>
1419 Date: Sat Jun 15 21:14:55 2013 -0400
1420
1421 Fix bug with reverse ordered statement_list
1422
1423 Fixed a bug in the parser where the statement_list was in reverse
1424 order, and the GraphViz DOT generating function was reversing it
1425 to hide the problem. Should be correct now, with a Yacc pattern
1426 that reverses it in the parser.
1427
1428 ast.ml | 5 +---
1429 parser.mly | 5 +++-
1430 2 files changed, 6 insertions(+), 4 deletions(-)
1431
1432 commit 0f627538861edd0c9ecaf3f42fdae3426a521d27
1433 Author: Kevin Graney <nanonet@gmail.com>
1434 Date: Sat Jun 15 20:16:03 2013 -0400
1435
1436 Add support for decimal literals
1437
1438 Added support for base 10 decimal literals.
1439
1440 scanner.mll | 2 ++

```

```

1441 1 file changed, 2 insertions(+)
1442
1443 commit 8b7ed5b0dd1cd37eb6888f807ecd3c1c22c763ea
1444 Author: Kevin Graney <nanonet@gmail.com>
1445 Date: Sat Jun 15 20:08:01 2013 -0400
1446
1447 Add arithmetic binary operators
1448
1449 Added addition, subtraction, multiplication, and division. Had to
1450 define division in terms of the existing FSLASH token. This makes
1451 the code a little unclear, and might be changed for clarity
1452 reasons in the future.
1453
1454 ast.ml | 11 ++++++++
1455 ast_types.mli | 3 +++
1456 parser.mly | 10 ++++++--
1457 scanner.mll | 2 +-
1458 4 files changed, 24 insertions(+), 2 deletions(-)
1459
1460 commit f5ed9c83e92888500ae3c5ca99f942cb67f3f9ed
1461 Author: Kevin Graney <nanonet@gmail.com>
1462 Date: Sat Jun 15 19:36:03 2013 -0400
1463
1464 Make output graph read left-to-right with input
1465
1466 The ouput GraphViz DOT tree for the AST now reads left-to-right in
1467 order with the input lines of code. Before it read right-to-left
1468 because the lists were not reversed before being passed to
1469 List.fold_left.
1470
1471 ast.ml | 9 +++++----
1472 1 file changed, 5 insertions(+), 4 deletions(-)
1473
1474 commit 3fd147d56da971257ce3df971ba2c136ce3b555c
1475 Author: Kevin Graney <nanonet@gmail.com>
1476 Date: Sat Jun 15 19:14:54 2013 -0400
1477
1478 Add interfaces for ast.ml
1479
1480 Added interface definitions and documentation for some of the
1481 functions in ast.ml.
1482
1483 ast.mli | 26 ++++++
1484 1 file changed, 26 insertions(+)
1485
1486 commit 901689776b1286fdc994959ddf8c32c6f4e838cd
1487 Author: Kevin Graney <nanonet@gmail.com>
1488 Date: Sat Jun 15 18:59:58 2013 -0400
1489
1490 Add design_docs target to top-level Makefile
1491
1492 Makefile | 6 +++++-
1493 1 file changed, 5 insertions(+), 1 deletion(-)
1494
1495 commit c49ed3e989557b9d5b32638972495963074f9122
1496 Author: Kevin Graney <nanonet@gmail.com>

```

```

1497 Date: Sat Jun 15 18:59:31 2013 -0400
1498
1499 Add pattern expressions to DOT output.
1500
1501 Added pattern expression to the GraphViz DOT output of the AST.
1502
1503 ast.ml | 28 ++++++++-----
1504 1 file changed, 27 insertions(+), 1 deletion(-)
1505
1506 commit 16fe0e54bba34ac41e52c74382a6f8ed32fd4ff3
1507 Author: Kevin Graney <nanonet@gmail.com>
1508 Date: Sat Jun 15 18:58:24 2013 -0400
1509
1510 Add expressions to DOT output
1511
1512 Added the expressions to the GraphViz DOT output of the AST.
1513
1514 ast.ml | 13 ++++++++
1515 1 file changed, 13 insertions(+)
1516
1517 commit 569b897a0d9b42356f72fc51603f2bd1d86c1d3f
1518 Author: Kevin Graney <nanonet@gmail.com>
1519 Date: Sat Jun 15 18:23:39 2013 -0400
1520
1521 Add support for non-terminals to DOT output
1522
1523 Added support for non-terminal nodes to the GraphViz dot output.
1524
1525 ast.ml | 11 ++++++---
1526 1 file changed, 9 insertions(+), 2 deletions(-)
1527
1528 commit b4b6df9eff32b4111f25569960620560b82487c0
1529 Author: Kevin Graney <nanonet@gmail.com>
1530 Date: Sat Jun 15 17:10:09 2013 -0400
1531
1532 Add additional punctuation to lexer
1533
1534 Added additional punctuation symbols to lexer for binary operators
1535 as well as parenthesis. Changed strings of length 1 to characters
1536 in pattern matching.
1537
1538 parser.mly | 9 ++++++---
1539 scanner.mll | 23 ++++++++-----
1540 2 files changed, 26 insertions(+), 6 deletions(-)
1541
1542 commit 466ff9a815d9c8b79df2216b9eae1345664259cc
1543 Author: Kevin Graney <nanonet@gmail.com>
1544 Date: Sat Jun 15 16:18:24 2013 -0400
1545
1546 Allow underscore in literal names
1547
1548 Allow an underscore as the second and subsequent character in
1549 literal names. In the future it may also be allowed as the first
1550 character.
1551
1552 scanner.mll | 2 +-

```

```

1553 1 file changed, 1 insertion(+), 1 deletion(-)
1554
1555 commit 4434b3157a6093affd36ffed17364181749eff84
1556 Author: Kevin Graney <nanonet@gmail.com>
1557 Date: Sat Jun 15 16:18:03 2013 -0400
1558
1559 Add Literals to the AST
1560
1561 ast_types.mli | 1 +
1562 parser.mly | 1 +
1563 2 files changed, 2 insertions(+)
1564
1565 commit 30fe72378962e036fbff55887fec066fd286b21a
1566 Author: Kevin Graney <nanonet@gmail.com>
1567 Date: Sat Jun 15 16:14:31 2013 -0400
1568
1569 Remove extraneous make targets in plt_docs
1570
1571 Removed the lrm and proposal targets since it's just as easy
1572 to type "make lrm.pdf" and "make proposal.pdf".
1573
1574 plt_docs/Makefile | 5 +----
1575 1 file changed, 1 insertion(+), 4 deletions(-)
1576
1577 commit 704d83a3797db42f08b46b00dd6cb2a3c31a8907
1578 Author: Kevin Graney <nanonet@gmail.com>
1579 Date: Wed Jun 12 21:15:26 2013 -0400
1580
1581 Add PLT documents
1582
1583 Added documents for PLT (COMS W4115) to the repository.
1584
1585 .gitignore | 7 ++++
1586 plt_docs/Makefile | 12 ++++++
1587 plt_docs/lrm.tex | 27 ++++++++
1588 plt_docs/proposal.tex | 109 ++++++++
1589 4 files changed, 155 insertions(+)
1590
1591 commit e500cbd1e47b791227ff13233cb0c38bf9940674
1592 Author: Kevin Graney <nanonet@gmail.com>
1593 Date: Wed Jun 12 21:11:38 2013 -0400
1594
1595 Renamed docs folder
1596
1597 Renamed the 'docs' folder since it's a little too generic, and for
1598 this project several types of documentation are expected. The
1599 ocaml doc documentation is now in a folder named 'design_docs'.
1600
1601 .gitignore | 4 +-
1602 design_docs/Makefile | 2 ++
1603 docs/Makefile | 2 --
1604 3 files changed, 4 insertions(+), 4 deletions(-)
1605
1606 commit 0158e187d22268255cca67ab6683321535cae3cb
1607 Author: Kevin Graney <nanonet@gmail.com>
1608 Date: Sun Jun 9 23:15:40 2013 -0400

```



```

1609
1610     Clean up the parse tree dot code generation
1611
1612     Cleaned up the print_tree function to abstract away some of the
1613     dot code generation into a couple routines.  Still need to expand
1614     the function to more types in the AST.
1615
1616     ast.ml | 47 ++++++-----
1617     1 file changed, 28 insertions(+), 19 deletions(-)
1618
1619     commit 691c33954f98ba711216915f3833f686eb98e852
1620     Author: Kevin Graney <nanonet@gmail.com>
1621     Date: Sun Jun 9 21:28:57 2013 -0400
1622
1623     Change print_tree to print AST in dot form
1624
1625     Changed the print_tree method to print out graphviz dot code to
1626     draw the tree.  This is still very crude, but it's much easier to
1627     visualize how the tree is being formed.
1628
1629     To support this change the Ast_types.program type was removed, and
1630     now the root program is represented as a Block statement instead of
1631     a list of statements.  This should be a nicer overall
1632     representation anyway.  The parser handles the type conversion.
1633
1634     ast.ml | 37 ++++++-----
1635     ast.mli | 2 +-
1636     ast_types.mli | 4 +--
1637     parser.mly | 4 +--
1638     4 files changed, 38 insertions(+), 9 deletions(-)
1639
1640     commit d86ec95db986452cdc51e728b06e1e25a1a2c858
1641     Author: Kevin Graney <nanonet@gmail.com>
1642     Date: Sun Jun 9 19:32:45 2013 -0400
1643
1644     Add variable bindings to pattern
1645
1646     Added support for variable bindings to the pattern matching.
1647
1648     ast_types.mli | 14 ++++++--
1649     parser.mly | 10 +++++--
1650     scanner.mll | 13 +++++-----
1651     3 files changed, 28 insertions(+), 9 deletions(-)
1652
1653     commit 39e981dd9607dd6354adec46ceef31885dfd7024
1654     Author: Kevin Graney <nanonet@gmail.com>
1655     Date: Sun Jun 9 19:28:19 2013 -0400
1656
1657     Change hex literals to require 0x prefix
1658
1659     Changed hex literals to require 0x prefix.  This avoids ambiguity
1660     with variable name literals, which are required to start with a
1661     letter.  (e.g. is 'ff' a variable or 255?)
1662
1663     parser.mly | 4 +--
1664     scanner.mll | 4 +--

```

```

1665 2 files changed, 4 insertions(+), 4 deletions(-)
1666
1667 commit ee673b1bcb29a7986af13e8050e096ed697b067f
1668 Author: Kevin Graney <nanonet@gmail.com>
1669 Date: Sun Jun 9 18:19:43 2013 -0400
1670
1671 Add pattern syntax to parser
1672
1673 First cut at getting the patten matching syntax into the parser.
1674
1675 parser.mly | 4 +---
1676 1 file changed, 1 insertion(+), 3 deletions(-)
1677
1678 commit 04ff7d64bb1eae6be83abe054d3a5dbe0596d710
1679 Author: Kevin Graney <nanonet@gmail.com>
1680 Date: Sun Jun 9 18:09:17 2013 -0400
1681
1682 Ignore auto-generated documentation
1683
1684 .gitignore | 2 ++
1685 1 file changed, 2 insertions(+)
1686
1687 commit 128bd7649ab921477918371202fe9900f124ef38
1688 Author: Kevin Graney <nanonet@gmail.com>
1689 Date: Sun Jun 9 12:58:16 2013 -0400
1690
1691 Update ocamldep dependencies
1692
1693 Updated ocamldep dependencies. These remain static for now.
1694
1695 Makefile | 26 ++++++++-----
1696 1 file changed, 16 insertions(+), 10 deletions(-)
1697
1698 commit 49ac872a829372682f68b1a1ca7bbb57ebde926c
1699 Author: Kevin Graney <nanonet@gmail.com>
1700 Date: Sun Jun 9 12:56:52 2013 -0400
1701
1702 Add block statements to grammar
1703
1704 Added support for block statements to the grammar, and began
1705 adding some scanner ability for the patterns. The bawk main loop
1706 is just dummy logic that prints 'item' for each statement it
1707 encounters.
1708
1709 ast.ml | 7 ++++++
1710 ast.mli | 4 +---
1711 ast_types.mli | 12 ++++++++
1712 bawk.ml | 16 ++++++++-----
1713 parser.mly | 25 ++++++++-----
1714 scanner.mll | 2 +-
1715 6 files changed, 54 insertions(+), 12 deletions(-)
1716
1717 commit 916197c3490919ccef12d77e9701a88a0b87c3b8
1718 Author: Kevin Graney <nanonet@gmail.com>
1719 Date: Sun Jun 9 12:09:11 2013 -0400
1720

```

```
1721     Add .gitignore
1722
1723     Ignore OCaml build files
1724
1725     .gitignore |    11 ++++++++
1726     1 file changed, 11 insertions(+)
1727
1728     commit 56c8597a3730aeb0a4d92c21379ce09974b9f418
1729     Author: Kevin Graney <nanonet@gmail.com>
1730     Date:   Sat Jun 8 22:45:30 2013 -0400
1731
1732     Initial commit
1733
1734     Basic parser for reading in a hexadecimal string and printing the
1735     value in hex.
1736
1737     Makefile      |   32 ++++++++
1738     ast.mli       |    4 ++++
1739     bawk.ml       |    9 ++++++++
1740     docs/Makefile |    2 ++
1741     parser.mly    |   18 ++++++++
1742     scanner.mll   |   30 ++++++++
1743     utile.ml      |    5 +++++
1744     utile.mli     |    5 +++++
1745     8 files changed, 105 insertions(+)
```

Chapter 5

Architectural design

5.1 Design decisions

The design decisions for this project were not made with considerations towards performance or scalability. Design decisions were based around ease of implementation and demonstration of compiler concepts. For this reason some features of a practical language were excluded because they cost significant effort for minimal benefit to the experience of writing a compiler.

5.2 Execution phases

The bawk interpreter consists of a bytecode compiler and a bytecode runtime environment. The two are seamlessly integrated into a single executable program, bawk. By calling this program with a `-c` option the results of the bytecode compiler are output, and by calling it with no option (or an explicit `-e`) the the given program is compiled internally to bytecode and then immediately executed in the bytecode interpreter.

5.2.1 Bytecode compiler

The bytecode compiler consists of a scanner (generated by OCamllex), a parser (generated by OCamllyacc), and the compiler itself. The scanner parses a sequence of input characters that form a given bawk language program and returns a list of tokens. The token types include punctuation that has meaning in the language, the C-style comments (`/* . . . */`), identifiers, keywords, constants, and string literals. Whitespace is consumed by the scanner, but serves no purpose other than to separate tokens.

There is at least one instance of ambiguity in the scanner. Identifier names, as defined in §3.1.2, consist of a letter or underscore followed by letters, underscores, and/or numerals. However, constants in pattern expressions, as defined in §3.5.1, consist of a sequence of hexadecimal digits, possibly including the letters 'a'-'f' in uppercase or lowercase. This leads to possible confusion in the code chunk listed below. On line 2 it is not known if `cafe123` is a variable name or a hexadecimal number. This ambiguity is actually resolved by the parser, which treats ambiguous constants as constants if they are valid hexadecimal.

Listing 5.1: Possibly ambiguous code

```
1 /cafe123:int1/ {
2     /cafe123/ {
3         /* do_something */
4     }
5 }
```

The parser converts the sequence of tokens into an abstract syntax tree (AST). The types used in the AST are defined in `ast_types.mli`, and rely on the OCaml tagged union types to represent the tree. The OCaml yacc specification for the parser is defined in the `parser.mly` file. The parser definition is fairly straight forward. There is some conversion done at this stage from strings to integers, but primarily it simply builds the AST using the tagged union types defined in the `Ast_types` module. The root of the generated AST is a statement, specifically a `Ast_types.Block` type, which contains a list of statements.

The bawk compiler can output the AST in a GraphViz format, which is useful for visualizing and debugging the AST. The `-ast` flag to the compiler will output GraphViz code. This can be turned into a PDF by piping it to GraphViz using the command

```
bawk -ast [binary filename] < [source filename] | dot -Tpdf -o out-file.pdf.
```

The logic to generate the GraphViz code is the majority of the `ast.ml` file. Other functions in that file serve to provide a few helper functions used when parsing the AST.

In the translation phase of the compiler, most of the heavy lifting is done by translating AST structures into basic blocks. In addition, the compiler keeps track of a symbol table in the `env` type in `compile.ml`. The structure forms a linked list, with a new node added each time a new, nested scope, is entered. Each node has a pointer to the parent scope, and a `StringMap` from symbols to either (1) the jump address for a subroutine, or (2) the index into globals. The environment also contains additional information for binding variables, which includes the size of the binding.

5.2.2 Bytecode operations

The compiler takes the AST created by the parser portion and translates it to a custom stack-based bytecode. The instructions for the bytecode are defined in `bytecode_types.mli`. Most of the instructions are related directly to the stack, but several are designed specifically for the application of this language: reading binary files.

The translation phase from the AST to bytecode is really a two phase translation process. The first does the basic translation from the AST to basic blocks of bytecode. In this initial translation the various branch instructions are given a unique value referencing a `Label` instruction. The `Label` instruction is a special pseudo instruction that is removed in the second phase of the translation. When the `Label` instructions are removed, their address, or position in the instruction list, is noted. These addresses then replace the values in all the branch instructions, converting them from the value referencing a label to the address of that label.

All the branch instructions in the bytecode operate on absolute addresses. There is no relative addressing in the bytecode by design, simply for simplicity. After the second phase of the translation is complete, the `Label` instructions removed and branch instructions pointing to addresses, the bytecode is ready for the runtime interpreter.

5.2.3 Bytecode runtime

The runtime interpreter is fairly simple, and defined primarily in `bytecode.ml`. The runtime executes the stack-based bytecode in a virtual machine environment. For simplicity the bytecode instructions composing the program, the execution stack, and global variable values are stored in three different arrays: `instructions`, `stack`, and `globals` respectively. Execution is performed by the recursive `exec` function, which takes a frame pointer, a stack pointer, and a program counter.

The program counter is simply an index into the `instructions` array for the instruction being executed in the current cycle. When the recursive call to `exec` is made, to execute the next instruction, the program counter is either incremented by one to run the next sequential instruction, or, in the case of a branch or function return instruction, called with some absolute address from the stack or an instruction argument. Since all branching is to absolute addresses the bytecode interpreter never computes relative offsets from the current program counter.

The stack grows from an index of zero in the `stack` array upwards, with items pushed and popped as individual instructions are executed. The argument to `exec` is a pointer to where the next value can be pushed on the stack. The stack only holds integer values, which is the primary motivation for not supporting functions, such as a `printf`-like function, that take string arguments.

Function arguments are placed on the stack, and the frame pointer references the location of these arguments. The previous frame pointer is also placed onto the stack and reset when the function returns, allowing nested function calls with arguments to work properly.

Global variables are stored in the `globals` array, with each symbol, other than function arguments, assigned an offset into this array. The name ‘globals’ is a bit of a misnomer, since variables local to blocks or functions are also stored in this same array. The slots in the `globals` array are not re-used; they last the entire duration of the program and are statically assigned to symbols by the compiler.

The bytecode interpreter stack is limited to 1024 values, and the `globals` array is limited to this same size of 1024 values. These limits are *not* enforced by the compiler, but the bytecode interpreter will crash if they are exceeded. The length of the `instructions` array is not explicitly limited.

5.3 Limitations

One major limitation of the language design is the lack of any type of return statement or return value from a function. This prevents the use of many types of useful recursion, and really limits the utility of the language. See §7.3.1 for more details on this limitation.

Chapter 6

Test plan

The testing for this project was performed using a variety of very tiny bawk language programs designed to exercise various features of the language. These programs are provided with the source code in the tests directory. Each program is named test-<name>.bawk and its output is expected to be identical to the test-<name>.bawk.out file. Both the source and expected output files are listed in §8.2.

The tests are all designed to operate on the binary file tests/lichtenstein.png, which serves as a modestly sized sample file. The beginning of the lichtenstein.png file is shown below for convenience, and the file itself can be found packaged with the bawk source.

```
0000000: 8950 4e47 0d0a 1a0a 0000 000d 4948 4452 .PNG.....IHDR
0000010: 0000 0200 0000 0200 0802 0000 007b 1a43 .....{.C
0000020: ad00 0000 0970 4859 7300 000a f000 000a .....pHYs.....
0000030: f001 42ac 3498 0000 0007 7449 4d45 07d7 ..B.4.....tIME..
0000040: 0511 0921 0919 38af 7500 0020 0049 4441 ...!.8.u.. .IDA
0000050: 5478 dab4 bd5b ac6d 5b76 1dd4 5aef 63ae Tx...[.m[v..Z.c.
0000060: bdcf bdb7 1ed7 a9b2 2b71 d995 38b1 63f2 .....+q..8.c.
0000070: b021 4fe7 0181 8404 a280 1421 10e2 2522 .!0.....!..%"
0000080: 7ef8 46fc 874f 7ef9 e40f 81f8 43e1 2bca ~.F..0~.....C.+
0000090: 4f82 2148 0ea0 5809 2838 ca8b d871 1c95 0.!H..X.(8...q..
00000a0: 53b6 eb75 5fe7 ec35 476f 7cf4 3ec6 1c73 S..u_..5Go|.>..s
00000b0: edb5 cfbf 65c8 aeab 5be7 eeb3 f75a 73cd ....e...[....Zs.
00000c0: 3946 1fbd b7d6 7aeb fcbd fff6 7f2b 7a17 9F....z.....+z.
00000d0: 4220 626b 6c06 2a48 9a39 e980 0114 0420 B bkl.*H.9.....
00000e0: 22f6 d09b c075 8fab 60f4 2638 f317 ed72 "...u..'&8...r
00000f0: d1e6 7cd8 daa5 b953 4e38 e194 2481 8289 ..|....SN8..$...
...
0059270: 3122 4a04 2150 a6c8 ff1f 00c1 c91a 2379 1"J.!P.....#y
0059280: 1a27 0000 0000 4945 4e44 ae42 6082 .'. ....IEND.B'.
```

The test programs are executed by a shell script, run_tests.sh, which runs each program and compares its output to the expected values.

Chapter 7

Lessons learned

This project provided a number of valuable software engineering and language design lessons. The language design lessons were invaluable, as this was my first attempt at writing a compiler from scratch that had any level of sophistication beyond regular expression matching on text files.

7.1 Language design

The design of the language was critical to being successful with this project. It was crucial to really think through the specifics of the language before even thinking about writing the compiler's translation logic. It was critical to write the scanner and parser in parallel with hashing out the language details. Writing these helped force the resolution of many shift/reduce ambiguities and couple reduce/reduce ambiguities in the original, pre-parser, design.

One visible addition to the language that came about from writing the parser is the `def` keyword before function definitions. This was added to resolve a conflict in the parser. Similarly, when adding `if` statements I realized that parentheses were required around the conditional expression when a parser conflict was resolved after adding them.

7.2 Scanner and parser generators

As described in §5.2.1, there is some ambiguity in the scanner distinguishing between hexadecimal values and identifier names. I didn't realize this ambiguity in the initial language design until testing the compiler with a hexadecimal constant beginning with a letter. This ambiguity was obviously not detected as a shift/reduce or reduce/reduce conflict by the parser generator because it was a collision of definitions in the scanner. Luckily it was possible to change the language definition slightly and add some logic in the parser to resolve the ambiguity.

7.2.1 GraphViz output

The GraphViz output was very useful in the early stages of writing the scanner and parser, but as it became clear that the basic design was solid the utility of the visualization decreased. Despite the decreasing returns on invested time, it wasn't a whole lot more difficult than implementing an ASCII format of AST output.

7.3 Compiler and bytecode design

7.3.1 Return statement problem

Late in the project I realized that I did not implement any type of return statement in the language. This was an oversight, but unfortunately the design of the bytecode did not make it easy to add at the last minute. The bytecode instruction `Rts` returns from a subroutine by popping the arguments to the function off the stack. Unfortunately this instruction requires a hard-coded constant for the number of instructions to pop. Without changing the bytecode and its interpreter to support reading this value from the stack it is difficult to return from an arbitrary point in the function. It could likely still be accomplished without a bytecode interpreter change, but it would require some redesign of the `translate` function in `compile.ml` to pass around additional information, namely a `Label` number for the end of the function. This change was just too large to risk making in the last week of the semester, so the language stands without a return statement.

7.3.2 Lack of full support for string types

One way to improve the `bawk` language would be to add support for strings inside of expressions. This would facilitate adding support for functions such as `printf` and also the reading of a fixed-length or null-terminated string from a file. Support for strings in this way was one of the initial conceptual ideas behind the language, but late in the semester it became apparent that the bytecode interpreter would have to be overhauled to support both integers and strings on the stack and in the `globals` array.

Chapter 8

Source code listing

8.1 Compiler and bytecode interpreter

8.1.1 Scanner & parser definitions

scanner.mll

The scanner.mll file contains Ocamllex specifications for the bawk scanner.

Listing 8.1: scanner.mll

```
1 { open Parser
2
3 let initial_string_buffer = String.create 256
4 let string_buff = ref initial_string_buffer
5 let string_index = ref 0
6
7 let is_in_string = ref false
8 let in_string () = !is_in_string
9
10 let reset_string_buffer () =
11     string_buff := initial_string_buffer;
12     string_index := 0
13
14 let store_string_char c =
15     if !string_index >= String.length (!string_buff) then
16     begin
17         let new_buff = String.create (String.length (!string_buff) * 2) in
18             String.blit (!string_buff) 0 new_buff 0 (String.length (!string_buff));
19             string_buff := new_buff
20     end;
21     String.unsafe_set (!string_buff) (!string_index) c;
22     incr string_index
23
24 let store_lexeme lexbuf =
25     let s = Lexing.lexeme lexbuf in
26     for i = 0 to String.length s - 1 do
27         store_string_char s.[i];
28     done
29
30 let get_stored_string () =
```

```

31 let s = String.sub (!string_buff) 0 (!string_index) in
32 string_buff := initial_string_buffer;
33 s
34
35 let char_for_backslash = function
36 | 'n' -> '\010'
37 | 'r' -> '\013'
38 | 'b' -> '\008'
39 | 't' -> '\009'
40 | c -> c
41
42 }
43
44 let hex_character = ['0'-'9' 'A'-'F' 'a'-'f']
45 let hex_sequence = "0x" hex_character+
46 let hexdigit_sequence = hex_character+
47 let literal = ['A'-'Z' 'a'-'z' '_' ]['0'-'9' 'A'-'Z' 'a'-'z' '_' ]*
48
49 let newline = ('\010' | "\013\010")
50
51 rule token = parse
52 | [' ' '\t' '\r' '\n'] { token lexbuf }
53 | "/*" { comment lexbuf }
54
55 (* language semantics *)
56 | ':' { COLON }
57 | ';' { SEMICOLON }
58 | ',' { COMMA }
59 | '/' { FSLASH }
60 | '{' { LBRACE }
61 | '}' { RBRACE }
62 | '(' { LPAREN }
63 | ')' { RPAREN }
64 | '+' { PLUS }
65 | '-' { MINUS }
66 | '*' { TIMES }
67 (* DIVIDE is the same as FSLASH *)
68 | '=' { ASSIGN }
69 | "==" { EQ }
70 | "!=" { NEQ }
71 | '<' { LT }
72 | "<=" { LEQ }
73 | '>' { GT }
74 | ">=" { GEQ }
75
76 (* type definitions *)
77 | "int1" { BIND_TYPE(Ast_types.Int_1_byte) }
78 | "int2" { BIND_TYPE(Ast_types.Int_2_bytes) }
79 | "int4" { BIND_TYPE(Ast_types.Int_4_bytes) }
80 | "uint1" { BIND_TYPE(Ast_types.UInt_1_byte) }
81 | "uint2" { BIND_TYPE(Ast_types.UInt_2_bytes) }
82 | "uint4" { BIND_TYPE(Ast_types.UInt_4_bytes) }
83
84 | "if" { IF }
85 | "else" { ELSE }
86

```

```

87 | "def" { DEF }
88 | literal as lit { LITERAL(lit) }
89
90 | hexdigit_sequence as lit { INT_LITERAL(lit)}
91 | hex_sequence as lit { INT_LITERAL(lit) }
92
93 (* string literals *)
94 | "\"" { reset_string_buffer();
95         is_in_string := true;
96         let string_start = lexbuf.Lexing.lex_start_p in
97         string lexbuf;
98         is_in_string := false;
99         lexbuf.Lexing.lex_start_p <- string_start;
100        STRING_LITERAL(get_stored_string()) }
101
102 | eof { EOF }
103
104
105 and comment = parse
106 | "*/" { token lexbuf }
107 | _ { comment lexbuf }
108
109 and string = parse
110   "' '
111   { () }
112 | '\\' newline ([' ' '\t'])
113   { string lexbuf }
114 | '\\' ['\\' '\n' '\t' '\b' '\r' ' ' ]
115   { store_string_char(char_for_backslash(Lexing.lexeme_char lexbuf 1));
116     string lexbuf }
117 | eof
118   { is_in_string := false;
119     assert false
120     (*raise (Error (Unterminated_string, 0))* ) }
121 | _
122   { store_string_char(Lexing.lexeme_char lexbuf 0);
123     string lexbuf }

```

parser.mly

The parser.mly file contains Ocaml yacc specifications for the bawk parser.

Listing 8.2: parser.mly

```

1 %{
2 open Parser_help
3 %}
4
5 %token <Ast_types.bind_type> BIND_TYPE
6
7 /* punctuation */
8 %token COLON SEMICOLON COMMA FSLASH RBRACE LBRACE RPAREN LPAREN DEF IF ELSE
9
10 /* operators */
11 %token PLUS MINUS TIMES ASSIGN

```

```

12
13 /* comparators */
14 %token EQ NEQ LT LEQ GT GEQ
15
16 %token <string> STRING_LITERAL
17 %token <string> LITERAL
18 %token <string> INT_LITERAL
19 %token EOF
20
21 /* associativity and precedence */
22 %nonassoc NOELSE
23 %nonassoc ELSE
24 %right ASSIGN
25 %left EQ NEQ
26 %left LT GT LEQ GEQ
27 %left PLUS MINUS
28 %left TIMES FSLASH
29
30 %start program
31 %type <Ast_types.statement> program
32
33 %%
34
35 pat_token:
36     | INT_LITERAL { parse_pattern_const $1 }
37     | LITERAL COLON BIND_TYPE { Ast_types.Binding($1, $3) }
38     | LITERAL {
39         (* Resolve a little ambiguity from the scanning phase, basically
40          binding variables can't be valid hex numbers or we don't know if it's
41          a hex number or a variable. Here it's resolved to be a number. *)
42         try parse_pattern_const $1
43         with Failure(_) -> Ast_types.Literal($1)
44         }
45     | STRING_LITERAL { Ast_types.PatString($1) }
46
47 pat_expr:
48     rev_pat_expr { List.rev $1 }
49
50 rev_pat_expr:
51     | { [] }
52     | rev_pat_expr pat_token { $2 :: $1 }
53
54 statement:
55     | expr SEMICOLON { Ast_types.Expr($1) }
56     | LBRACE statement_list RBRACE { Ast_types.Block($2) }
57     | FSLASH pat_expr FSLASH statement { Ast_types.Pattern($2,$4) }
58     | DEF LITERAL LPAREN lit_list RPAREN LBRACE statement_list RBRACE
59         { Ast_types.FunctionDecl({
60             Ast_types.fname = $2;
61             Ast_types.arguments = $4;
62             Ast_types.body = Ast_types.Block($7)}) }
63     }
64     | IF LPAREN expr RPAREN statement %prec NOELSE
65         { Ast_types.If($3, $5, Ast_types.Block([])) }
66     | IF LPAREN expr RPAREN statement ELSE statement
67         { Ast_types.If($3, $5, $7) }

```

```

68
69 expr:
70   | INT_LITERAL { Ast_types.LitInt(int_of_string $1) }
71   | expr PLUS   expr { Ast_types.Binopt($1, Ast_types.Add, $3) }
72   | expr MINUS  expr { Ast_types.Binopt($1, Ast_types.Subtract, $3) }
73   | expr TIMES  expr { Ast_types.Binopt($1, Ast_types.Multiply, $3) }
74   | expr FSLASH expr { Ast_types.Binopt($1, Ast_types.Divide, $3) }
75   | expr EQ     expr { Ast_types.Binopt($1, Ast_types.Equal, $3) }
76   | expr NEQ    expr { Ast_types.Binopt($1, Ast_types.Neq, $3) }
77   | expr LT     expr { Ast_types.Binopt($1, Ast_types.Less, $3) }
78   | expr LEQ    expr { Ast_types.Binopt($1, Ast_types.Leq, $3) }
79   | expr GT     expr { Ast_types.Binopt($1, Ast_types.Greater, $3) }
80   | expr GEQ    expr { Ast_types.Binopt($1, Ast_types.Geq, $3) }
81   | LITERAL ASSIGN expr { Ast_types.Assign($1, $3) }
82   | LITERAL { Ast_types.ExprLiteral($1) }
83   | LITERAL LPAREN expr_list RPAREN { Ast_types.Call($1, $3) }
84   | STRING_LITERAL { Ast_types.LitString($1) }
85
86 expr_list:
87   | rev_expr_list { List.rev $1 }
88
89 rev_expr_list:
90   | { [] }
91   | expr { [$1] }
92   | rev_expr_list COMMA expr { $3 :: $1 }
93
94 statement_list:
95   | rev_statement_list { List.rev $1 }
96
97 rev_statement_list:
98   | { [] }
99   | rev_statement_list statement { $2 :: $1 }
100
101 lit_list:
102   | rev_lit_list { List.rev $1 }
103
104 rev_lit_list:
105   | { [] }
106   | LITERAL { [$1] }
107   | rev_lit_list COMMA LITERAL { $3 :: $1 }
108
109 program:
110   statement_list { Ast_types.Block($1) }

```

8.1.2 Interfaces

Listing 8.3: parser_help.mli

```

1
2
3 (** [parse_pattern_const] converts a pattern constant into a [PatternBytes]
4   list of integers for the AST. This method interprets leading zeros as
5   significant among other things. *)
6 val parse_pattern_const: string -> Ast_types.pat_token

```

Listing 8.4: ast_types.mli

```
1 (** Types used by the AST *)
2
3 type operators = Add | Subtract | Multiply | Divide
4   | Equal | Neq | Less | Leq | Greater | Geq
5
6 type bind_type =
7   | Int_1_byte
8   | Int_2_bytes
9   | Int_4_bytes
10  | UInt_1_byte
11  | UInt_2_bytes
12  | UInt_4_bytes
13
14 type literal = string
15
16 type expr =
17   | LitInt of int
18   | ExprLiteral of string
19   | Binopt of expr * operators * expr
20   | Call of string * expr list
21   | LitString of string
22   | Assign of string * expr
23
24 type pat_token =
25   | PatternByte of int
26   | PatternBytes of pat_token list
27   | Binding of literal * bind_type
28   | Literal of literal
29   | PatString of string
30
31 type pat_expr = pat_token list
32
33 type func_decl = {
34   fname: string;
35   arguments: string list;
36   body: statement;
37 }
38 and statement =
39   | Pattern of pat_expr * statement
40   | Block of statement list
41   | Expr of expr
42   | FunctionDecl of func_decl
43   | If of expr * statement * statement
```

Listing 8.5: bytecode_types.mli

```
1
2 type instruction =
3   Lit of int      (* Push a literal *)
4   | Drp          (* Discard a value *)
5   | Bin of Ast_types.operators (* Perform arithmetic on top of stack *)
6   | Two of int   (* Convert top of stack to two's complement *)
7   | Lod of int  (* Fetch global variable *)
8   | Str of int  (* Store global variable *)
9   | Lfp        (* Load frame pointer from the stack *)
```

```

10 | Sfp          (* Store frame pointer to the stack *)
11 | Jsrf of int  (* Call function by absolute address *)
12 | Ent         (* Push FP, FP -> SP, SP += i *)
13 | Rts of int  (* Restore FP, SP, consume formals, push result *)
14 | Beq of int  (* Branch absolute if top-of-stack is zero *)
15 | Bne of int  (* Branch absolute if top-of-stack is non-zero *)
16 | Bra of int  (* Branch absolute *)
17 | Beo of int  (* Branch if end of file *)
18 | Hlt        (* Terminate *)
19 | Rdb of int  (* Read a number of bytes from the file *)
20 | Ldp        (* Load the position in the file onto the stack *)
21 | Skp        (* Seek the file to the position on the stack *)
22
23 (* Pseudo bytecode instructions used by compiler *)
24 | Label of int (* Used for branching jumps *)

```

Listing 8.6: ast.mli

```

1
2
3 (** [print_tree] outputs the AST in GraphViz DOT format. This is useful for
4     visualizing how the parser is constructing the tree *)
5 val print_tree: Ast_types.statement -> unit
6
7 (** [string_of_bind_type] converts an [Ast_types.bind_type] to a [string] *)
8 val string_of_bind_type: Ast_types.bind_type -> string
9
10 (** [size_of_bind_type] returns the size, in bytes, of an [Ast_types.bind_type]
11     *)
12 val size_of_bind_type: Ast_types.bind_type -> int
13
14 (** [is_signed_type] returns true if the bind type is signed and false
15     otherwise *)
16 val is_signed_type: Ast_types.bind_type -> bool
17
18 (** [make_terminal] creates a terminal node for the node with an id of [this],
19     and labels it with the string argument. If the optional [parent] argument
20     is provided a link is also created between this node and the parent. *)
21 val make_terminal: this:int -> ?parent:int -> string -> unit
22
23 (** [make_nonterminal] creates a non-terminal node for the node with an id of
24     [this], and labels it with the string argument. If the optional [parent]
25     argument is provided a link is also created between this node and the
26     parent. *)
27 val make_nonterminal: this:int -> ?parent:int -> string -> unit
28
29 (** [make_link] creates a link between two of the GraphViz nodes *)
30 val make_link: int -> int -> unit
31
32 (** [folded_printer] returns a function that can be used with [List.fold_left]
33     to traverse a list of objects that should be children of a common parent.
34     The parent's ID is the second argument and the print function is the
35     first argument. *)
36 val folded_printer: ('a -> int -> int -> 'b) -> int -> int -> 'a -> 'b
37
38 (** [string_of_operator] returns a string of the bytecode instruction for

```



```

39     the given operator. *)
40 val string_of_operator: Ast_types.operators -> string

```

Listing 8.7: compile.mli

```

1  (** Compilation of bawk code *)
2
3  module StringMap : (* TODO: figure out how to get rid of this monstrosity! *)
4    sig
5      type key = String.t
6      type 'a t = 'a Map.Make(String).t
7      val empty : 'a t
8      val is_empty : 'a t -> bool
9      val mem : key -> 'a t -> bool
10     val add : key -> 'a -> 'a t -> 'a t
11     val singleton : key -> 'a -> 'a t
12     val remove : key -> 'a t -> 'a t
13     val merge :
14       (key -> 'a option -> 'b option -> 'c option) -> 'a t -> 'b t -> 'c t
15     val compare : ('a -> 'a -> int) -> 'a t -> 'a t -> int
16     val equal : ('a -> 'a -> bool) -> 'a t -> 'a t -> bool
17     val iter : (key -> 'a -> unit) -> 'a t -> unit
18     val fold : (key -> 'a -> 'b -> 'b) -> 'a t -> 'b -> 'b
19     val for_all : (key -> 'a -> bool) -> 'a t -> bool
20     val exists : (key -> 'a -> bool) -> 'a t -> bool
21     val filter : (key -> 'a -> bool) -> 'a t -> 'a t
22     val partition : (key -> 'a -> bool) -> 'a t -> 'a t * 'a t
23     val cardinal : 'a t -> int
24     val bindings : 'a t -> (key * 'a) list
25     val min_binding : 'a t -> key * 'a
26     val max_binding : 'a t -> key * 'a
27     val choose : 'a t -> key * 'a
28     val split : key -> 'a t -> 'a t * 'a option * 'a t
29     val find : key -> 'a t -> 'a
30     val map : ('a -> 'b) -> 'a t -> 'b t
31     val mapi : (key -> 'a -> 'b) -> 'a t -> 'b t
32   end
33
34
35   type pattern_binding = {
36     loc: int;
37     size: int;
38   }
39
40   type env = {
41     symbol_map: int StringMap.t;
42     bindings: pattern_binding StringMap.t;
43     parent: env ref option;
44   }
45
46   (** [translate_program] takes an AST statement and returns a list of bytecode
47       instructions implementing the statement *)
48   val translate_program: Ast_types.statement -> Bytecode_types.instruction list
49
50   val clean_environment: env

```

Listing 8.8: bytecode.mli

```
1 (** Bytecode generation *)
2
3 (** [execute_bytecode] is a function called from the main argument parser. It
4     generates and then executes the bytecode instructions. *)
5 val execute_bytecode: Bytecode_types.instruction list -> in_channel -> unit
6
7 (** [print_bytecode] is a function called from the main argument parser. It
8     generates and then prints out the bytecode instructions. *)
9 val print_bytecode: Bytecode_types.instruction list -> unit
10
11 (** [string_of_instruction] converts a bytecode instruction into a string
12     representation *)
13 val string_of_instruction: Bytecode_types.instruction -> string
14
15 (** [execute_instructions] interprets an array of bytecode instructions
16     as a program, executing said instructions. *)
17 val execute_instructions: Bytecode_types.instruction array -> in_channel -> unit
18
19 (** [is_pseudo] returns true if the given bytecode instruction is a pseudo
20     instruction and false otherwise *)
21 val is_pseudo: Bytecode_types.instruction -> bool
22
23 val enumerate_instructions: Bytecode_types.instruction list ->
24     (int * Bytecode_types.instruction) list
```

Listing 8.9: reader.mli

```
1 (** The Reader module contains functions for extracting and decoding bits of
2     information from an [in_channel] for a given binary file. *)
3
4
5 (** [read_byte] reads a single byte from the input channel and returns an
6     integer representing that byte *)
7 val read_byte: in_channel -> int
8
9 (** [read_bytes] reads the next [n] bytes from the input channel into a list of
10    integers *)
11 val read_bytes: in_channel -> int -> int list
12
13 (** [read_string_null] reads a null-terminated string from the input channel *)
14 val read_string_null: in_channel -> string
15
16 (** [read_string_fixed] reads a fixed-length string from the input channel *)
17 val read_string_fixed: in_channel -> int -> string
18
19 (** [read_unsigned] reads the next [n] bytes as an unsigned integer
20
21     The algorithm is: read bytes into a list with [read_bytes], traverse the
22     list in reverse order, and accumulate a sum of every byte shifted by it's
23     position in the list. *)
24 val read_unsigned: in_channel -> int -> int
25
26 (** [get_pos] returns the position in the file *)
27 val get_pos: in_channel -> int
28
```

```

29 (** [set_pos] set the position in the file *)
30 val set_pos: in_channel -> int -> unit
31
32 val advance: in_channel -> int -> unit
33
34 (** [is_eof] returns true if the next byte of the file reaches EOF and false
35     otherwise *)
36 val is_eof: in_channel -> bool

```

Listing 8.10: utile.mli

```

1 (** General purpose utility functions *)
2
3
4 (** [int_of_hex] converts a string of hexadecimal digits into an integer *)
5 val int_of_hex: string -> int
6
7 (** [enumerate] converts a list of values in a list of tuples of an integer
8     followed by that value. Integers begin at zero and count upwards *)
9 val enumerate: ?step:('a -> int -> int) -> ?start:(int) -> 'a list
10    -> (int * 'a) list
11
12 (** [explode] converts a string into a list of characters *)
13 val explode: string -> char list
14
15 (** [implode] converts a list of characters into a string *)
16 val implode: char list -> string
17
18 (** [bytes_of_string] converts a string into a list of integers representing
19     the byte values for each character *)
20 val bytes_of_string: string -> int list
21
22 (** [signed_of_unsigned] converts an integer into its twos' complement
23     interpretation *)
24 val signed_of_unsigned: int -> int -> int

```

8.1.3 Source Files

Listing 8.11: bawk.ml

```

1 (** The main function of the bawk compiler *)
2 open Compile
3
4 module StringMap = Map.Make(String) (* TODO: can we use the other def? *)
5
6 (** Possible actions for the compiler to take *)
7 type action = Ast | Compile | Execute
8             | D_print_clean_env
9
10 (** [decode_action] reads the [argv] array of command line arguments and returns
11     the action that the compiler is being asked to take. This is fairly
12     primitive, and only the first argument is used to make the decision. *)
13 let decode_action argv =
14     if Array.length argv > 1 then

```

```

15     List.assoc argv.(1) [
16         ("-ast", Ast);
17         ("-c", Compile);
18         ("-e", Execute);
19         ("-D-print-clean-env", D_print_clean_env)
20     ]
21     else Execute;;
22
23 let decode_in_file argv =
24     if Array.length argv > 2 then
25         open_in_bin argv.(2)
26     else
27         open_in_bin "/dev/random";;
28
29
30 (** [parse_channel] runs an input channel through the lexer and parser phases
31     of the compiler returning an AST. *)
32 let parse_channel channel =
33     let lexbuf = Lexing.from_channel channel in
34     let program = Parser.program Scanner.token lexbuf in
35     program;;
36
37 let _ =
38     let action = decode_action Sys.argv in
39     let in_file = decode_in_file Sys.argv in
40     match action with
41     | Ast -> Ast.print_tree (parse_channel stdin)
42     | Compile -> Bytecode.print_bytecode
43         (Compile.translate_program (parse_channel stdin))
44     | Execute -> Bytecode.execute_bytecode
45         (Compile.translate_program (parse_channel stdin)) in_file
46
47     (* Debug actions *)
48     | D_print_clean_env ->
49         let env = Compile.clean_environment in
50         Printf.printf "Starting symbol table:\n";
51         StringMap.iter (fun s i -> Printf.printf "\t%d %s\n" i s)
52             env.symbol_map;;

```

Listing 8.12: parser_help.ml

```

1
2 let split_bytes s =
3     let split s =
4         let num_bytes = (String.length s) / 2 in
5         let arr = Array.make num_bytes (String.create 2) in
6         for i = 0 to num_bytes - 1 do
7             arr.(i) <- String.create 2;
8             String.blit s (2*i) arr.(i) 0 2
9         done;
10        Array.to_list arr
11    in
12    let pieces =
13        if (String.length s) mod 2 == 0 then split s
14        else split (Printf.sprintf "0%s" s)
15    in

```

```

16 List.map (fun x -> int_of_string ("0x" ^ x)) pieces;;
17
18 let parse_pattern_const s =
19   let lst = split_bytes s in
20   let bytes = List.map (fun x -> Ast_types.PatternByte x) lst in
21   Ast_types.PatternBytes bytes

```

Listing 8.13: ast.ml

```

1 open Ast_types
2 open Printf
3
4 let make_link n1 n2 =
5   printf "%d -> %d\n" n1 n2;;
6
7 let conditional_link_with_parent this_id parent_id =
8   if parent_id != -1 then make_link parent_id this_id
9   else ();;
10
11 let make_terminal ~this:this_id ?parent:(parent_id=(-1)) label =
12   printf "%d [label=\"%s\" shape=hexagon style=filled]\n" this_id label;
13   conditional_link_with_parent this_id parent_id;;
14
15 let make_nonterminal ~this:this_id ?parent:(parent_id=(-1)) label =
16   printf "%d [label=\"%s\" shape=plaintext]\n" this_id label;
17   conditional_link_with_parent this_id parent_id;;
18
19 let folded_printer func parent_id =
20   (fun id x -> let this_id = id + 1 in
21     func x this_id parent_id );;
22
23 let print_plain str id parent =
24   make_nonterminal str ~this:id ~parent:parent;
25   id + 1
26
27 let size_of_bind_type = function
28   | Int_1_byte -> 1
29   | Int_2_bytes -> 2
30   | Int_4_bytes -> 4
31   | UInt_1_byte -> 1
32   | UInt_2_bytes -> 2
33   | UInt_4_bytes -> 4
34
35 let string_of_bind_type = function
36   | Int_1_byte -> "int1"
37   | Int_2_bytes -> "int2"
38   | Int_4_bytes -> "int4"
39   | UInt_1_byte -> "uint1"
40   | UInt_2_bytes -> "uint2"
41   | UInt_4_bytes -> "uint4"
42
43 let is_signed_type = function
44   | Int_1_byte -> true
45   | Int_2_bytes -> true
46   | Int_4_bytes -> true
47   | UInt_1_byte -> false

```

```

48 | UInt_2_bytes -> false
49 | UInt_4_bytes -> false
50
51 let string_of_operator = function
52 | Add -> "Add"
53 | Subtract -> "Sub"
54 | Multiply -> "Mul"
55 | Divide -> "Div"
56 | Equal -> "Eq"
57 | Neq -> "Neq"
58 | Less -> "Lt"
59 | Leq -> "Leq"
60 | Greater -> "Gt"
61 | Geq -> "Geq"
62
63 (* each *_print function returns the next id available for use *)
64 let print_tree prog =
65   let rec stmt_print root id parent =
66     match root with
67     | Pattern(pat_expr, stmt) ->
68       let pattern_id = id + 1 in
69       let stmt_id = pat_expr_print pat_expr pattern_id id in
70       let consumed_ids = stmt_print stmt (stmt_id + 1) stmt_id in
71       make_nonterminal "pattern" ~this:id ~parent:parent;
72       make_nonterminal "statement" ~this:stmt_id ~parent:id;
73       consumed_ids
74
75     | Block(statement_lst) ->
76       make_nonterminal "block" ~this:id ~parent:parent;
77       List.fold_left (folded_printer stmt_print id) id statement_lst
78
79     | Expr(expr) ->
80       make_nonterminal "expression" ~this:id ~parent:parent;
81       expr_print expr (id + 1) id
82
83     | FunctionDecl(decl) ->
84       let stmt_id = stmt_print decl.body (id + 1) id in
85       let consumed_ids = List.fold_left
86         (folded_printer print_plain stmt_id) (stmt_id + 1)
87         decl.arguments
88       in
89       make_nonterminal "arguments" ~this:stmt_id ~parent:id;
90       make_nonterminal (sprintf "fdef:%s" decl.fname) ~this:id
91         ~parent:parent;
92       consumed_ids
93
94     | If(cond, stmt1, stmt2) ->
95       let if_id = expr_print cond (id + 1) id in
96       let ifstmt = stmt_print stmt1 if_id id in
97       let consumed_ids = stmt_print stmt2 ifstmt id in
98       make_nonterminal "if" ~this:id ~parent:parent;
99       consumed_ids
100
101     (* | _ ->
102       make_nonterminal "other_stmt" ~this:id ~parent:parent;
103       id + 1 *)

```

```

104 and expr_print expr id parent =
105   match expr with
106   | LitInt(value) ->
107     make_terminal (string_of_int value) ~this:id ~parent:parent;
108     id + 1
109   | Binopt(e1, op, e2) ->
110     let e2_id = expr_print e1 (id + 1) id in
111     let consumed_ids = expr_print e2 e2_id id in
112     make_nonterminal (string_of_operator op) ~this:id ~parent:parent;
113     consumed_ids
114   | Call(name, arg_list) ->
115     let name_id = id + 1 in
116     let arg_id = name_id + 1 in
117     make_nonterminal "call" ~this:id ~parent:parent;
118     make_terminal name ~this:name_id ~parent:id;
119     make_nonterminal "arguments" ~this:arg_id ~parent:id;
120     List.fold_left (folded_printer expr_print arg_id) arg_id arg_list
121   | LitString(str) ->
122     make_terminal str ~this:id ~parent:parent;
123     id + 1;
124   | Assign(name, expr) ->
125     make_nonterminal "assign" ~this:id ~parent:parent;
126     make_nonterminal name ~this:(id + 1) ~parent:id;
127     expr_print expr (id + 2) id
128   | _ ->
129     make_nonterminal "other_expr" ~this:id ~parent:parent;
130     id + 1
131
132 and pat_expr_print lst id parent =
133   make_nonterminal "pat_expr" ~this:id ~parent:parent;
134   List.fold_left (folded_printer pat_token_print id) id lst
135
136 and pat_token_print token id parent =
137   match token with
138   | Binding(literal, bind_type) ->
139     let literal_id = id + 1 in
140     let type_id = id + 2 in
141     make_nonterminal "binding" ~this:id ~parent:parent;
142     make_terminal literal ~this:literal_id ~parent:id;
143     make_terminal (string_of_bind_type bind_type)
144       ~this:type_id ~parent:id;
145     id + 3
146   | Literal(literal) ->
147     make_nonterminal "literal" ~this:id ~parent:parent;
148     make_terminal literal ~this:(id + 1) ~parent:id;
149     id + 2
150   | PatternBytes(values) ->
151     make_nonterminal "bytes" ~this:id ~parent:parent;
152     List.fold_left (folded_printer pat_token_print id) id values
153   | PatternByte(value) ->
154     make_terminal (string_of_int value) ~this:id ~parent:parent;
155     id + 1
156   | PatString(str) ->
157     make_terminal (sprintf "%s" str) ~this:id ~parent:parent;
158     id + 1
159   (*| _ ->

```

```

160         make_nonterminal "foo" ~this:(id+1) ~parent:id;
161         id + 2*)
162     in
163     print_string "digraph AST {\n";
164     print_string "ordering = out;\n";
165     ignore (stmt_print prog 0 0);
166     print_string "}";;

```

Listing 8.14: compile.ml

```

1 open Ast_types
2 open Bytecode_types
3
4 module StringMap = Map.Make(String)
5 exception Compile_error of string;;
6
7 type pattern_binding = {
8     loc: int;
9     size: int;
10 }
11
12 type env = {
13     symbol_map: int StringMap.t;
14     bindings: pattern_binding StringMap.t;
15     parent: env ref option;
16 }
17
18 let built_in_functions = ["print"; "RP"];;
19 let special_vars = ["LE"];;
20
21 let clean_environment =
22     let preset_binding_list =
23         (Utile.enumerate ~step:(fun x y -> y - 1) ~start:(-1)
24          built_in_functions)
25         @ Utile.enumerate special_vars in
26     let built_ins =
27         List.fold_left (fun map item ->
28             let (value, key) = item in
29             StringMap.add key (value) map
30         ) StringMap.empty preset_binding_list in
31     {
32         symbol_map = built_ins;
33         bindings = StringMap.empty;
34         parent = None;
35     }
36
37 let new_environment parent_env =
38     ref {
39         symbol_map = StringMap.empty;
40         bindings = StringMap.empty;
41         parent = Some parent_env;
42     };;
43
44 let rec resolve_symbol env name =
45     try StringMap.find name !env.symbol_map
46     with Not_found ->

```



```

47     match !env.parent with
48     None -> let error_msg = Printf.sprintf "No such symbol: %s" name in
49         raise (Compile_error error_msg)
50     | Some(env) -> resolve_symbol env name;;
51
52 let rec resolve_binding env name =
53     try StringMap.find name !env.bindings
54     with Not_found ->
55         match !env.parent with
56         None -> let error_msg = Printf.sprintf "No such symbol: %s" name in
57             raise (Compile_error error_msg)
58         | Some(env) -> resolve_binding env name;;
59
60 let add_function env fname addr =
61     let symbol_map_new = StringMap.add fname addr env.symbol_map in
62     {env with symbol_map = symbol_map_new}
63
64 let global_counter = ref (List.length special_vars);;
65 let get_next_global () =
66     global_counter := !global_counter + 1;
67     !global_counter;;
68
69 let add_variable_force env vname id =
70     let symbol_map_new = StringMap.add vname id
71         env.symbol_map in
72     {env with symbol_map = symbol_map_new}
73
74 let add_variable env vname =
75     let symbol_map_new = StringMap.add vname (get_next_global ())
76         env.symbol_map in
77     {env with symbol_map = symbol_map_new}
78
79 let rec get_var_address env vname =
80     try StringMap.find vname !env.symbol_map;
81     with Not_found ->
82         match !env.parent with
83         None -> env := add_variable !env vname;
84             StringMap.find vname !env.symbol_map
85         | Some(env) -> get_var_address env vname;;
86
87 let add_binding env bname size =
88     let vaddr = get_var_address env bname in
89     (* let symbol_map_new = StringMap.add bname vaddr env.symbol_map in *)
90     let bindings_new = StringMap.add bname { loc = vaddr; size = size}
91         !env.bindings in
92     {!env with bindings = bindings_new}
93
94 exception Bad_binding of string;;
95
96 let get_binding_address env bname size =
97     let bind_info =
98     try StringMap.find bname !env.bindings;
99     with Not_found ->
100         env := add_binding env bname size;
101         StringMap.find bname !env.bindings
102     in

```

```

103   if bind_info.size != size then
104       raise (Bad_binding "size mismatch")
105   else bind_info.loc;;
106
107   let label_counter = ref 0;;
108   let get_new_label () =
109       label_counter := !label_counter + 1;
110       !label_counter;;
111
112   let form_label_map instructions =
113       let arr = Array.make (!label_counter + 1) 0 in
114       let enumerated = Bytecode.enumerate_instructions instructions in
115       List.iter (fun x ->
116           let (addr, instr) = x in
117           match instr with
118               Label (id) -> arr.(id) <- addr; ()
119           | _ -> ()
120       ) enumerated;
121       arr;;
122
123   let resolve_labels instructions =
124       let label_map = form_label_map instructions in
125       let rec emit_resolution = function
126           [] -> []
127           (* rewrite branches *)
128           | Bne(id)::t -> Bne label_map.(id) :: emit_resolution t
129           | Bra(id)::t -> Bra label_map.(id) :: emit_resolution t
130           | Beq(id)::t -> Beq label_map.(id) :: emit_resolution t
131           | Beo(id)::t -> Beo label_map.(id) :: emit_resolution t
132           | Jsr(id)::t when id >= 0 -> Jsr label_map.(id) :: emit_resolution t
133
134           (* remove the pseudo instructions *)
135           | Label(id)::t -> emit_resolution t
136
137           (* all other instructions *)
138           | h::t -> h :: emit_resolution t
139       in emit_resolution instructions;;
140
141   let rec translate_expr env expr =
142       let recurse = translate_expr env in
143       match expr with
144           LitInt(integer) ->
145               [Bytecode_types.Lit(integer)]
146       | ExprLiteral(var_name) ->
147           let vaddr = resolve_symbol env var_name in
148           [Lod vaddr]
149       | Binopt(e1, op, e2) ->
150           recurse e1 @ recurse e2 @ [Bin op]
151       | Call(func_name, args) ->
152           let function_addr = resolve_symbol env func_name in
153           (List.concat (List.map recurse args)) @ [Jsr function_addr]
154       | Assign(var_name, expr) ->
155           let vaddr = get_var_address env var_name in
156           translate_expr env expr @ [
157               Str vaddr
158           ]

```

```

159 | LitString(str) ->
160 |   let error_msg = Printf.sprintf "Invalid use of string literal: %s" str
161 |   in raise (Compile_error error_msg)
162 |   []
163
164 and translated_pattern env expr fail_label =
165 | let rec check_item = function
166 |   PatternByte(value) -> [
167 |     Rdb 1;
168 |     Lit value;
169 |     Bin Subtract;
170 |     Bne fail_label (* branch to failed match *)
171 |   ]
172 | PatternBytes(bytes) ->
173 |   let big_endian = List.flatten (List.map check_item bytes) in
174 |   let little_endian = List.flatten
175 |     (List.map check_item (List.rev bytes)) in
176 |   let le_label = get_new_label () in
177 |   let end_label = get_new_label () in
178 |   [
179 |     Lod 0; (* Load LE *)
180 |     Bne le_label;
181 |   ]
182 |   @ big_endian
183 |   @ [ Bra end_label; Label le_label; ]
184 |   @ little_endian
185 |   @ [ Label end_label; ]
186 | Binding(literal, bind_type) ->
187 |   let num_bytes = Ast.size_of_bind_type bind_type in
188 |   let vaddr = get_binding_address env literal num_bytes in
189 |   [
190 |     (* TODO: handle EOF here? *)
191 |     Rdb num_bytes;
192 |   ]
193 |   @ (
194 |     if Ast.is_signed_type bind_type then [Two (num_bytes * 8);]
195 |     else []
196 |   )
197 |   @ [ Str vaddr; ]
198 | PatString(str) ->
199 |   let bytes = Utile.bytes_of_string str in
200 |   let ast_bytes = List.map (fun x -> PatternByte x) bytes in
201 |   translated_pattern env ast_bytes fail_label @ [];
202 | Literal(symbol) ->
203 |   let register = resolve_symbol env symbol in
204 |   let binding = resolve_binding env symbol in
205 |   let size = binding.size in
206 |   [
207 |     Lod register;
208 |     Rdb size;
209 |     Bin Subtract;
210 |     Bne fail_label
211 |   ]
212 | in
213 | List.flatten (List.map check_item expr)
214

```

```

215 and translate env stmt =
216   let recurse = translate env in
217   match stmt with
218     Block(stmt_list) ->
219       List.flatten (List.map recurse stmt_list)
220   | Expr(expr) ->
221     translate_expr env expr
222   | Pattern(pat_expr, stmt) ->
223     let start_label = get_new_label () in
224     let fail_label = get_new_label () in
225     let end_label = get_new_label () in
226
227     let new_env = new_environment env in
228
229     let pattern_code = translated_pattern new_env pat_expr fail_label in
230     [ Ldp; Label start_label; Ldp; Beo end_label; ] @ pattern_code @
231     translate new_env stmt @
232     [
233       Bra end_label;
234       Label fail_label;
235       Skp;
236       Beo end_label;
237       Rdb 1; Drp;
238       Bra start_label;
239       Label end_label;
240       Skp;
241       Skp;
242     ]
243   | FunctionDecl(decl) ->
244     let start_label = get_new_label () in
245     let end_label = get_new_label () in
246     let num_args = List.length decl.arguments in
247     let new_env = new_environment env in
248     List.iter (fun x ->
249       let (id,name) = x in
250       new_env := add_variable_force !new_env name id;
251       ());
252     (Utile.enumerate ~step:(fun x y -> y - 1) ~start:(-100)
253      (List.rev decl.arguments));
254     env := add_function !env decl.fname start_label;
255     [ Bra end_label; Label start_label; Ent; ] @
256     translate new_env decl.body @
257     [ Rts num_args; Label end_label ]
258   | If(condition, ifstmt, elsestmt) ->
259     let if_label = get_new_label () in
260     let end_label = get_new_label () in
261     translate_expr env condition
262     @ [ Bne if_label; ] @ translate env elsestmt
263     @ [ Bra end_label; Label if_label; ]
264     @ translate env ifstmt @ [ Label end_label; ];;
265
266 let translate_program stmt =
267   let env = ref clean_environment in
268   resolve_labels (translate env stmt @ [Hlt]);;

```

Listing 8.15: bytecode.ml

```
1 open Printf
2 open Bytecode_types
3 open Ast_types
4
5 let string_of_instruction = function
6   Lit(integer) -> sprintf "Lit %d" integer
7   | Bin(operator) -> "Bin " ^ Ast.string_of_operator operator
8   | Two(num) -> sprintf "Two %d" num
9   | Rdb(num) -> sprintf "Rdb %d" num
10  | Jsr(num) -> sprintf "Jsr %d" num
11  | Rts(num) -> sprintf "Rts %d" num
12  | Lod(num) -> sprintf "Lod %d" num
13  | Str(num) -> sprintf "Str %d" num
14  | Bra(addr) -> sprintf "Bra %d" addr
15  | Beq(addr) -> sprintf "Beq %d" addr
16  | Bne(addr) -> sprintf "Bne %d" addr
17  | Beo(addr) -> sprintf "Beo %d" addr
18  | Drp -> "Drp"
19  | Ldp -> "Ldp"
20  | Skp -> "Skp"
21  | Lfp -> "Lfp"
22  | Sfp -> "Sfp"
23  | Ent -> "Ent"
24  | Hlt -> "Hlt"
25
26  | Label(id) -> sprintf "Label %d (PSEUDO)" id
27
28 let is_pseudo = function
29   Label(id) -> true
30   | _ -> false
31
32 let enumerate_instructions lst =
33   Utile.enumerate ~step:(fun x y ->
34     if is_pseudo x then y
35     else y + 1) lst;;
36
37 let print_bytecode instructions =
38   List.iter (fun x -> let (i, ins) = x in
39     printf "%9d: %s\n" i (string_of_instruction ins))
40     (enumerate_instructions instructions);;
41
42 let execute_instructions instructions on_file =
43   let stack = Array.make 1024 0 in
44   let globals = Array.make 1024 0 in
45   let rec exec fp sp pc =
46     (* Printf.printf "(fp=%d sp=%d pc=%d) " fp sp pc;
47     Array.iter (fun d -> Printf.printf "%d " d) stack;
48     Printf.printf "\n"; *)
49     match instructions.(pc) with
50     Lit i -> stack.(sp) <- i;
51     exec fp (sp + 1) (pc + 1)
52   | Bin op ->
53     let op1 = stack.(sp - 2) and op2 = stack.(sp - 1) in
54     stack.(sp - 2) <- (
55       let boolean i = if i then 1 else 0 in
```

```

56         match op with
57             Add -> op1 + op2
58             | Subtract -> op1 - op2
59             | Multiply -> op1 * op2
60             | Divide -> op1 / op2
61             | Equal -> boolean (op1 == op2)
62             | Neq -> boolean (op1 != op2)
63             | Less -> boolean (op1 < op2)
64             | Leq -> boolean (op1 <= op2)
65             | Greater -> boolean (op1 > op2)
66             | Geq -> boolean (op1 >= op2)
67         );
68         exec fp (sp - 1) (pc + 1)
69     | Two size ->
70         let operand = stack.(sp - 1) in
71         stack.(sp - 1) <- (
72             Utile.signed_of_unsigned operand size
73         );
74         exec fp sp (pc + 1)
75     | Jsr (-1) -> (* print function *)
76         print_endline (string_of_int stack.(sp - 1));
77         exec fp (sp - 1) (pc + 1)
78     | Rdb (1) ->
79         stack.(sp) <- Reader.read_byte on_file;
80         (* Printf.printf "Read byte: %x\n" stack.(sp); *)
81         exec fp (sp + 1) (pc + 1)
82     | Rdb (n) ->
83         stack.(sp) <- Reader.read_unsigned on_file n;
84         exec fp (sp + 1) (pc + 1)
85     | Ldp ->
86         stack.(sp) <- Reader.get_pos on_file;
87         (* Printf.printf "Store position: %d at %d\n" stack.(sp) sp; *)
88         exec fp (sp + 1) (pc + 1)
89     | Skp ->
90         Reader.set_pos on_file stack.(sp - 1);
91         (* Printf.printf "Set position: %d\n" stack.(sp - 1); *)
92         exec fp (sp - 1) (pc + 1)
93     | Bne addr -> if stack.(sp - 1) == 0 then
94         exec fp (sp - 1) (pc + 1) else
95         exec fp (sp - 1) addr
96     | Beq addr -> if stack.(sp - 1) != 0 then
97         exec fp (sp - 1) (pc + 1) else
98         exec fp (sp - 1) addr
99     | Bra addr -> exec fp sp addr
100    | Beo addr ->
101        if Reader.is_eof on_file then exec fp sp addr
102        else exec fp sp (pc + 1)
103    | Jsr addr ->
104        stack.(sp) <- pc + 1;
105        exec fp (sp + 1) addr
106    | Rts(num) ->
107        (* Printf.printf "\nReturn to %d from %d (fp=%d, num=%d, sp=%d)\n"
108            stack.(fp - 1) (fp - 1) stack.(sp - 1) num sp; *)
109        exec stack.(fp) (fp - num) stack.(fp - 1)
110    | Lod(-2) -> (* RP variable *)
111        stack.(sp) <- Reader.get_pos on_file;

```

```

111         exec fp (sp + 1) (pc + 1)
112     | Lod index ->
113         (
114             if (index > 0) then stack.(sp) <- globals.(index)
115             else stack.(sp) <- stack.(fp + index + 100 - 2)
116         );
117         exec fp (sp + 1) (pc + 1)
118
119     | Str(-2) -> (* RP variable *)
120         Reader.set_pos on_file stack.(sp);
121         exec fp (sp - 1) (pc + 1)
122     | Str index ->
123         (
124             if (index > 0) then globals.(index) <- stack.(sp - 1)
125             else stack.(fp + index + 100 - 2) <- stack.(sp - 1)
126         );
127         exec fp (sp - 1) (pc + 1)
128     | Ent ->
129         stack.(sp) <- fp;
130         exec sp (sp + 1) (pc + 1)
131     | Lfp ->
132         exec stack.(sp) (sp - 1) (pc + 1)
133     | Drp ->
134         exec fp (sp - 1) (pc + 1)
135     | Hlt -> ()
136 in exec 0 0 0
137
138 let execute_bytecode instructions on_file =
139     execute_instructions (Array.of_list instructions) on_file;;

```

Listing 8.16: reader.ml

```

1
2 let read_byte ic =
3     let byte = input_char ic in
4     int_of_char byte;;
5
6 let read_bytes ic n =
7     let rec read c lst =
8         if c == 0 then List.rev lst
9         else read (c - 1) (read_byte ic :: lst)
10    in read n [];;
11
12 let read_string_null ic =
13     let rec scan ic str =
14         let c = input_char ic in
15         match c with
16         | '\000' -> Utile.implode (List.rev str)
17         | _ -> scan ic (c :: str)
18    in scan ic [];;
19
20 let read_string_fixed ic n =
21     let rec scan ic str n =
22         if n == 0 then Utile.implode (List.rev str)
23         else
24             let c = input_char ic in

```

```

25         scan ic (c :: str) (n - 1)
26     in scan ic [] n;;
27
28 let read_unsigned ic n =
29     let bytes = read_bytes ic n in
30     let rbytes = List.rev bytes in
31     let rec form_int n shift bytes =
32         match bytes with
33         | [] -> n
34         | h :: t ->
35             let next_n = n + (h lsl shift) in
36             form_int next_n (shift + 8) t
37     in form_int 0 0 rbytes;;
38
39 let advance ic n =
40     let pos = pos_in ic in
41     let new_pos = pos + n in
42     seek_in ic new_pos;;
43
44 let get_pos ic =
45     pos_in ic;;
46
47 let set_pos ic n =
48     seek_in ic n;;
49
50 let is_eof ic =
51     let pos = get_pos ic in
52     let result = try input_char ic; false with End_of_file -> true in
53     set_pos ic pos; result;;

```

Listing 8.17: utile.ml

```

1
2
3 let int_of_hex str =
4     int_of_string ("0x" ^ str);;
5
6 let enumerate ?step:(step=(fun x y -> y + 1)) ?start:(start=0) lst =
7     let rec enum count = function
8         [] -> []
9         | h :: t -> (count, h) :: (enum (step h count) t)
10    in enum start lst;;
11
12 let explode s =
13     let rec exp i l =
14         if i < 0 then l else exp (i - 1) (s.[i] :: l) in
15     exp (String.length s - 1) [];;
16
17 let implode l =
18     let res = String.create (List.length l) in
19     let rec imp i = function
20         | [] -> res
21         | c :: l -> res.[i] <- c; imp (i + 1) l in
22     imp 0 l;;
23
24 let bytes_of_string str =

```



```

25   let chars = explode str in
26   List.map int_of_char chars;;
27
28 let signed_of_unsigned value num_bits =
29   let shift = num_bits - 1 in
30   let mask = 1 lsl shift in
31   if mask land value > 0 then (value land (lnot mask)) - mask
32   else value;;

```

8.2 Test files

All test files are designed to be run on the input file `lichtenstein.png` unless specified otherwise.

8.2.1 test-arith.bawk

Listing 8.18: Test program (test-arith.bawk)

```

1 print(0 + 10 - 9 + 3);

```

Listing 8.19: Expected output (test-arith.bawk.out)

```

1 4

```

8.2.2 test-cond1.bawk

Listing 8.20: Test program (test-cond1.bawk)

```

1
2 if (100 > 100*2) {
3   print(1);
4 } else {
5   print(0);
6 }
7
8
9 if (100) {
10  print (100);
11 }
12
13 if (0) {
14   print (0);
15 }

```

Listing 8.21: Expected output (test-cond1.bawk.out)

```

1 0
2 100

```

8.2.3 test-func1.bawk

Listing 8.22: Test program (test-func1.bawk)

```
1
2 def food () {
3     print(100);
4 }
5
6 print(1);
7 food();
8 print(2);
9 food();
10 print(3);
```

Listing 8.23: Expected output (test-func1.bawk.out)

```
1 1
2 100
3 2
4 100
5 3
```

8.2.4 test-func2.bawk

Listing 8.24: Test program (test-func2.bawk)

```
1 print(1);
2
3 def food () {
4     print(100);
5 }
6
7 print(2);
```

Listing 8.25: Expected output (test-func2.bawk.out)

```
1 1
2 2
```

8.2.5 test-func3.bawk

Listing 8.26: Test program (test-func3.bawk)

```
1
2 a = 100;
3
4 def bar(c) {
5     print(c);
6 }
7
8 def foo(a, b, c) {
```

```

9   print (a);
10  print (b);
11  print (c);
12  a = 10;
13  bar(c * 2);
14  print (a);
15  print (c);
16  c = c * 2;
17  print (c);
18 }
19
20 foo(50+20, 10, a);

```

Listing 8.27: Expected output (test-func3.bawk.out)

```

1 70
2 10
3 100
4 200
5 10
6 100
7 200

```

8.2.6 test-pat1.bawk

Listing 8.28: Test program (test-pat1.bawk)

```

1 /89504e47/ { /* should match */
2   print(10);
3 }
4
5 /89504347/ { /* should not match */
6   print(20);
7 }
8
9 print(30);

```

Listing 8.29: Expected output (test-pat1.bawk.out)

```

1 10
2 30

```

8.2.7 test-pat2.bawk

Listing 8.30: Test program (test-pat2.bawk)

```

1 /89504e47/ { /* should match */
2   print(10);
3 }
4
5 /89504e47/ { /* should match */
6   print(20);
7 }

```

```
8
9 print(30);
```

Listing 8.31: Expected output (test-pat2.bawk.out)

```
1 10
2 20
3 30
```

8.2.8 test-pat3.bawk

Listing 8.32: Test program (test-pat3.bawk)

```
1 /8950/ { /* should match */
2     print(10);
3 }
4
5 /8951123456789/ { /* should not match */
6     print(20);
7 }
8
9 /8950/ { /* should match */
10    print(40);
11 }
12
13 print(30);
```

Listing 8.33: Expected output (test-pat3.bawk.out)

```
1 10
2 40
3 30
```

8.2.9 test-pat4.bawk

Listing 8.34: Test program (test-pat4.bawk)

```
1 a = 100;
2
3 /89 foo:int1 4e47/ { /* should match */
4     print(foo);
5 }
6
7 /89504e47/ { /* should match */
8     print(20);
9 }
10
11
12 print(30);
```

Listing 8.35: Expected output (test-pat4.bawk.out)

```
1 80
2 20
3 30
```

8.2.10 test-pat5.bawk

Listing 8.36: Test program (test-pat5.bawk)

```
1
2 /89 foo2:int4/ { /* should match */
3     print(foo2); /* foo = 0x504e470d */
4 }
```

Listing 8.37: Expected output (test-pat5.bawk.out)

```
1 1347307277
```

8.2.11 test-pat6.bawk

Listing 8.38: Test program (test-pat6.bawk)

```
1
2 /* these values should be different */
3
4 /000a value:int4/ { print(value); }
5 /000a value:uint4/ { print(value); }
6
7 /000a value:int2/ { print (value); }
8 /000a value:uint2/ { print (value); }
9
10 /000a value:int1/ { print (value); }
11 /000a value:uint1/ { print (value); }
12
13
14 /* these values should be the same */
15
16 /000d value:int4/ { print(value); }
17 /000d value:uint4/ { print(value); }
18
19 /000d value:int2/ { print(value); }
20 /000d value:uint2/ { print(value); }
21
22 /000d value:int1/ { print(value); }
23 /000d value:uint1/ { print(value); }
```

Listing 8.39: Expected output (test-pat6.bawk.out)

```
1 -268435446
2 4026531850
3 -4096
4 61440
```

```
5 -16
6 240
7 1229472850
8 1229472850
9 18760
10 18760
11 73
12 73
```

8.2.12 test-pat7.bawk

Listing 8.40: Test program (test-pat7.bawk)

```
1
2 my_RP = RP;
3 print(my_RP);
4
5 / foo:int1 19 38af / {
6     print(foo);
7     my_RP = RP;
8     print(my_RP);
9 }
10
11 print(my_RP);
12 RP = my_RP;
13
14 / foo:int1 00 / {
15     print(foo);
16 }
```

Listing 8.41: Expected output (test-pat7.bawk.out)

```
1 0
2 9
3 72
4 72
5 117
```

8.2.13 test-pat8.bawk

Listing 8.42: Test program (test-pat8.bawk)

```
1
2 /4e47 value:int1/ {
3     print(value);
4     print(RP);
5     /value value2:int2/ {
6         print(value2);
7         print(RP);
8     }
9 }
```

Listing 8.43: Expected output (test-pat8.bawk.out)

```
1 13
2 5
3 18760
4 14
```

8.2.14 test-pat9.bawk**Listing 8.44: Test program (test-pat9.bawk)**

```
1 /* LE defaults to false (0) */
2
3 / 091938af / {
4     print(100);
5     print(RP);
6 }
7
8 LE = 1;
9
10 / af381909 / {
11     print(200);
12     print(RP);
13 }
14
15 / 09 19 38 af / {
16     print(300);
17     print(RP);
18 }
19
20 / af 38 19 09 / {
21     print(401);
22     print(RP);
23 }
24
25 LE = 0;
26
27 / 091938af / {
28     print(500);
29     print(RP);
30 }
31
32
33 / 09 19 38 af / {
34     print(600);
35     print(RP);
36 }
```

Listing 8.45: Expected output (test-pat9.bawk.out)

```
1 100
2 72
3 200
4 72
5 300
```

```
6 72
7 500
8 72
9 600
10 72
```

8.2.15 test-pat10.bawk

Listing 8.46: Test program (test-pat10.bawk)

```
1
2 LE = 0;
3 / "IHDR" / {
4     print(100);
5     print(RP);
6 }
7
8 LE = 1;
9 / "IHDR" / {
10    print(200);
11    print(RP);
12 }
```

Listing 8.47: Expected output (test-pat10.bawk.out)

```
1 100
2 16
3 200
4 16
```

8.2.16 test-pat11.bawk

Listing 8.48: Test program (test-pat11.bawk)

```
1
2 /"IEND" end:uint2/ {
3     print(end);
4 }
5
6 /"IEND" end:int2/ {
7     print(end);
8 }
```

Listing 8.49: Expected output (test-pat11.bawk.out)

```
1 44610
2 -20926
```


8.2.17 test-scope1.bawk

Listing 8.50: Test program (test-scope1.bawk)

```
1
2 a = 1000;
3 b = 200;
4
5 def foo() {
6     a = 10 + 100;
7     print (a);
8     print (b);
9 }
10
11 print (a);
12 foo();
13 print (a);
14 print (b);
```

Listing 8.51: Expected output (test-scope1.bawk.out)

```
1 1000
2 110
3 200
4 110
5 200
```

8.2.18 test-scope2.bawk

Listing 8.52: Test program (test-scope2.bawk)

```
1
2 b = 200;
3
4 def foo() {
5     a = 10 + 100;
6     print (a);
7     print (b);
8 }
9
10 a = 1000;
11
12 print (a);
13 foo();
14 print (a);
15 print (b);
```

Listing 8.53: Expected output (test-scope2.bawk.out)

```
1 1000
2 110
3 200
4 110
5 200
```

8.2.19 test-useful1.bawk

Listing 8.54: Test program (test-useful1.bawk)

```
1 /* print the size of a PNG file */
2
3 /89 "PNG" 0d 0a 1a 0a/ {
4   /length:uint4 "IHDR"/ {
5     /width:uint4 height:uint4/ {
6       print(width);
7       print(height);
8     }
9   }
10 }
```

Listing 8.55: Expected output (test-useful1.bawk.out)

```
1 512
2 512
```

8.2.20 test-var1.bawk

Listing 8.56: Test program (test-var1.bawk)

```
1
2 a = 10 + 20 * 20;
3 b = 20 * 20;
4 print(a);
5 print(b);
```

Listing 8.57: Expected output (test-var1.bawk.out)

```
1 410
2 400
```

8.3 Misc. Files

Listing 8.58: Makefile

```
1 DEP_FILENAME=.ocamldeps.mk
2
3 OBJS = utile.cmo \
4   reader.cmo \
5   scanner.cmo \
6   parser_help.cmo \
7   parser.cmo \
8   ast.cmo \
9   compile.cmo \
10  bytecode.cmo \
11  bawk.cmo
12
13 default: bawk plt_docs/lrm.pdf plt_docs/proposal.pdf design_docs
```

```

14
15 plt_docs/report.pdf: $(wildcard plt_docs/report-*.tex) plt_docs/gitlog.tex
16
17 plt_docs/gitlog.tex:
18     git log --stat > plt_docs/gitlog.tex
19
20 %.pdf: %.tex
21     cd $(shell dirname $@); pdflatex $(shell basename $<)
22     cd $(shell dirname $@); pdflatex $(shell basename $<)
23
24 bawk: $(OBJS)
25     ocamlc -o bawk $(shell ocamldep -sort *.ml | sed 's/\.ml/\.cmo/g')
26
27 scanner.ml: scanner.mll
28     ocamllex scanner.mll
29
30 parser.ml parser.mli: parser.mly
31     ocamlyacc parser.mly
32
33 %.cmo: %.ml | $(DEP_FILENAME)
34     ocamlc -c $<
35
36 %.cmi: %.mli | $(DEP_FILENAME)
37     ocamlc -c $<
38
39 clean:
40     rm -f bawk *.cmi *.cmo scanner.ml parser.ml parser.mli
41     rm -f design_docs/*.html design_docs/*.css
42     rm -f plt_docs/*.pdf plt_docs/*.toc plt_docs/*.aux plt_docs/*.log plt_docs/*.
43         lol
44     rm -f $(DEP_FILENAME)
45     rm -f plt_docs/gitlog.tex
46
47 design_docs: bawk
48     $(MAKE) -C ./design_docs
49
50 include $(DEP_FILENAME)
51
52 $(DEP_FILENAME): parser.ml parser.mli scanner.ml
53     ocamldep *.ml *.mli > $(DEP_FILENAME)
54
55 .PHONY: clean default design_docs

```

Listing 8.59: run_tests.sh

```

1 #!/bin/sh
2
3 BAWK="./bawk"
4
5 # Set time limit for all operations
6 ulimit -t 30
7
8 globallog=testall.log
9 rm -f $globallog
10 error=0
11 globalerror=0

```

```

12
13 keep=0
14
15 Usage() {
16     echo "Usage: testall.sh [options] [.mc files]"
17     echo "-k    Keep intermediate files"
18     echo "-h    Print this help"
19     exit 1
20 }
21
22 SignalError() {
23     if [ $error -eq 0 ] ; then
24         echo "FAILED"
25         error=1
26     fi
27     echo " $1"
28 }
29
30 # Compare <outfile> <reffile> <difffile>
31 # Compares the outfile with reffile. Differences, if any, written to difffile
32 Compare() {
33     generatedfiles="$generatedfiles $3"
34     echo diff -b $1 $2 ">" $3 1>&2
35     diff -b "$1" "$2" > "$3" 2>&1 || {
36         SignalError "$1 differs"
37         echo "FAILED $1 differs from $2" 1>&2
38     }
39 }
40
41 # Run <args>
42 # Report the command, run it, and report any errors
43 Run() {
44     echo $* 1>&2
45     eval $* || {
46         SignalError "$1 failed on $*"
47         return 1
48     }
49 }
50
51 Check() {
52     error=0
53     basename=`echo $1 | sed 's/.*\\\/\\\/
54                 s/.mc//'\`
55     reffile=`echo $1 | sed 's/.mc$//'\`
56     basedir="`echo $1 | sed 's/\[^\/\]*$//'\`/."
57
58     echo -n "$basename..."
59
60     echo 1>&2
61     echo "##### Testing $basename" 1>&2
62
63     generatedfiles=""
64
65     generatedfiles="$generatedfiles ${basename}.i.out" &&
66     Run "$BAWK" "-e tests/lichtenstein.png" "<" $1 ">" ${basename}.i.out &&
67     Compare ${basename}.i.out ${reffile}.out ${basename}.i.diff

```

```

68
69 #generatedfiles="$generatedfiles ${basename}.c.out" &&
70 #Run "$BAWK" "-c" "<" $1 ">" ${basename}.c.out &&
71 #Compare ${basename}.c.out ${reffile}.out ${basename}.c.diff
72
73 # Report the status and clean up the generated files
74
75 if [ $error -eq 0 ] ; then
76 if [ $keep -eq 0 ] ; then
77     rm -f $generatedfiles
78 fi
79 echo "OK"
80 echo "##### SUCCESS" 1>&2
81 else
82 echo "##### FAILED" 1>&2
83 globalerror=$error
84 fi
85 }
86
87 while getopts kdpsh c; do
88     case $c in
89         k) # Keep intermediate files
90             keep=1
91             ;;
92         h) # Help
93             Usage
94             ;;
95         esac
96 done
97
98 shift `expr $OPTIND - 1`
99
100 if [ $# -ge 1 ]
101 then
102     files=$@
103 else
104     files="tests/fail-*.bawk tests/test-*.bawk"
105 fi
106
107 for file in $files
108 do
109     case $file in
110         *test-*)
111             Check $file 2>> $globallog
112             ;;
113         *fail-*)
114             CheckFail $file 2>> $globallog
115             ;;
116         *)
117             echo "unknown file type $file"
118             globalerror=1
119             ;;
120     esac
121 done
122
123 exit $globalerror

```