

Constructing Lexers with Ragel & Objective-C

Jim Rea
ProVUE Development

Typical Applications

1) Formula evaluation

- Spreadsheets
- Databases
- Programming languages

2) Other Text Based Protocols

- CSV
- HTTP

Database Calculations

The screenshot shows a database application window titled "Temperatures (Panorama 6) — Edited". The window contains a table with the following data:

Date	Location	Fahrenheit	Celsius
11/04/13	Huntington Beach, CA	64	18
11/04/13	Charlottesville, VA	48	9
11/04/13	Lahaina, HI	82	28
11/04/13	Saskatoon, SK	28	-2
11/05/13	Huntington Beach, CA	73	23
11/05/13	Charlottesville, VA	57	14
11/05/13	Lahaina, HI	81	27
11/05/13	Saskatoon, SK	21	-6
11/06/13	Huntington Beach, CA	75	24
11/06/13	Charlottesville, VA	64	18
11/06/13	Lahaina, HI	82	28
11/06/13	Saskatoon, SK	30	-1
11/07/13	Huntington Beach, CA	79	26
11/07/13	Charlottesville, VA	66	19
11/07/13	Lahaina, HI	81	27
11/07/13	Saskatoon, SK	32	0

A red arrow points from the "Fahrenheit" column of the first row to a "Formula:" dialog box. The dialog box contains a text input field with the formula $(\text{Fahrenheit} - 32) * 5 / 9$ and "Cancel" and "OK" buttons.

Lexical Analysis

Convert free flow text into a sequence of tokens

(Fahrenheit - 32) * 5 / 9

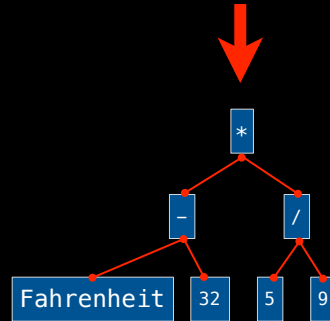


(Fahrenheit - 32) * 5 / 9

Parsing

Converts sequential data into a data structure

(Fahrenheit - 32) * 5 / 9



Implementing a Tokenizer by Hand

- 1) Character by character
- 2) Using NSScanner
- 3) Using NSRegularExpression

all of these techniques can be
tedious and error prone :(

Automated Tools

- | | |
|----------|----------|
| 1) lex | 4) flex |
| 2) yacc | 5) cup |
| 3) bison | 6) lemon |

and dozens more...

http://en.wikipedia.org/wiki/Comparison_of_parser_generators

Ragel

<http://www.complang.org/ragel/>

- Creates state machines for tokenizing text
 - Targets C, C++, Objective-C, D, Java and Ruby
- Written by Dr. Adrian Thurston
- First released January 2002
 - Latest Update February 2013

Why Ragel?

- Targets C/Objective-C
- Supports 16 bit characters (Unicode)
- No external dependencies
(no external library to link to)
- Flexible licensing
(no restrictions on generated code)
- Fast runtime execution
- Can be integrated into XCode
- 12 year track record, active community

only package I could find with this combination -- other unicode savvy tools required linking with a GPL'd library

Installing Ragel

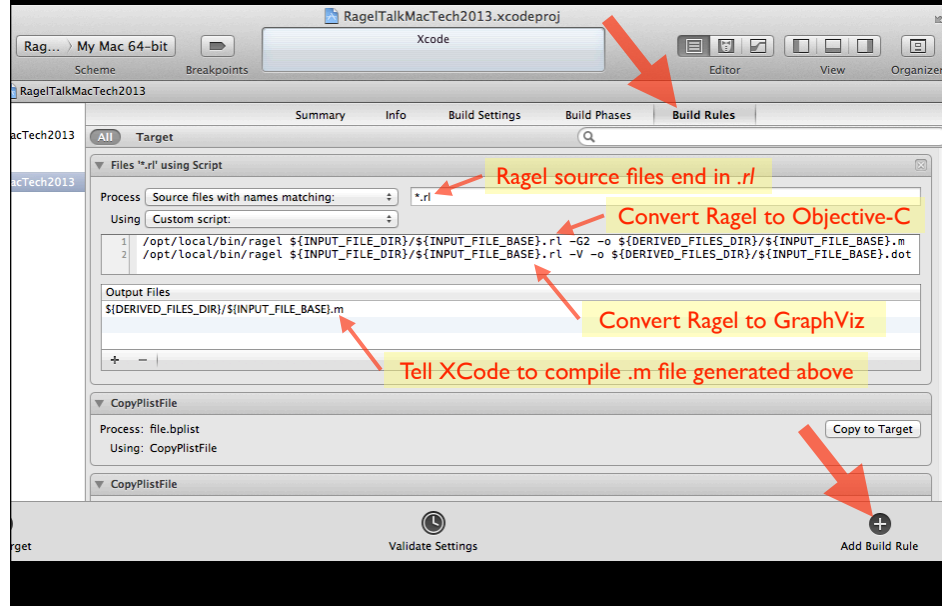
1) Install MacPorts

<http://www.macports.org/install.php>

2) Install Ragel (using Terminal.app)

```
sudo port install ragel
```

Adding Ragel to your project



More about Xcode Build Rules

<http://www.cocoawithlove.com/2010/02/custom-build-rules-generated-tables-and.html>



Adding a Ragel Source File

- 1) File > New > File...
- 2) Template: Other > Empty
- 3) Filename ends with .rl

```
test.rl  
formulas.rl  
parser.rl
```

Empty Ragel State Machine

```
#import <Foundation/Foundation.h>
#include <stdio.h>
#include <string.h>
#include <stdlib.h>
```

Everything else is host code
(in this case Objective-C)

```
%%{
```

```
    machine test_lexer;
```

```
}%%
```

Ragel code is inside %%{ ... }%% blocks

```
%% write data;
```

Ragel directives on lines beginning with %%

Anatomy of a Ragel Source File

imports/includes

state machine

regex action

regex action

regex action

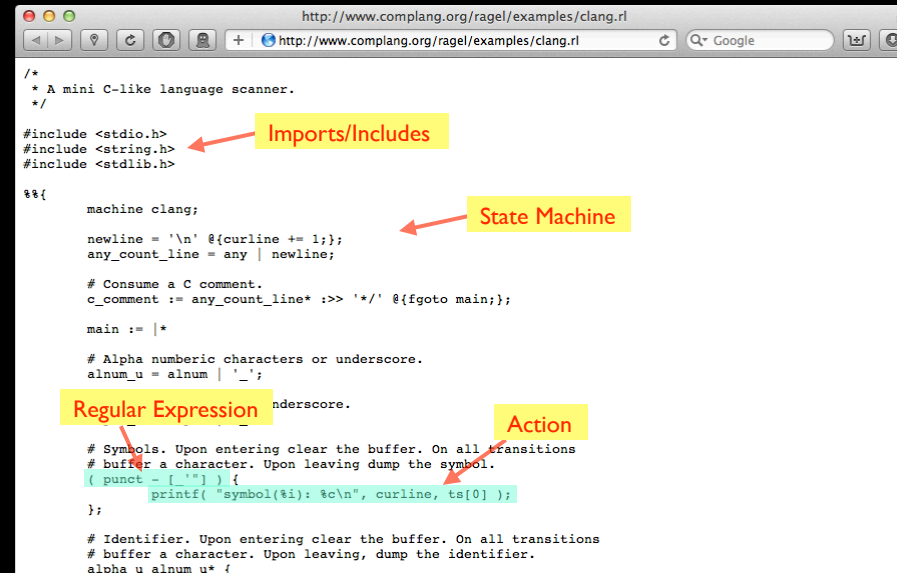
regex action

regex action

regex action

scanner function

Sample “Mini-C” Language Scanner



The image shows a web browser window displaying the source code of a Mini-C language scanner. The browser's address bar shows the URL `http://www.complang.org/ragel/examples/clang.rl`. The code is written in Ragel and includes several annotations with red arrows pointing to specific parts:

- Imports/Includes:** Points to the `#include <stdio.h>`, `#include <string.h>`, and `#include <stdlib.h>` lines.
- State Machine:** Points to the `machine clang;` line.
- Regular Expression:** Points to the `alnum_u = alnum | '_';` line.
- Action:** Points to the `printf("symbol(%i): %c\n", curline, ts[0]);` line.

```
/*
 * A mini C-like language scanner.
 */

#include <stdio.h>
#include <string.h>
#include <stdlib.h>

%%{
    machine clang;

    newline = '\n' @{curline += 1;};
    any_count_line = any | newline;

    # Consume a C comment.
    c_comment := any_count_line* :>> '*'/{fgoto main;};

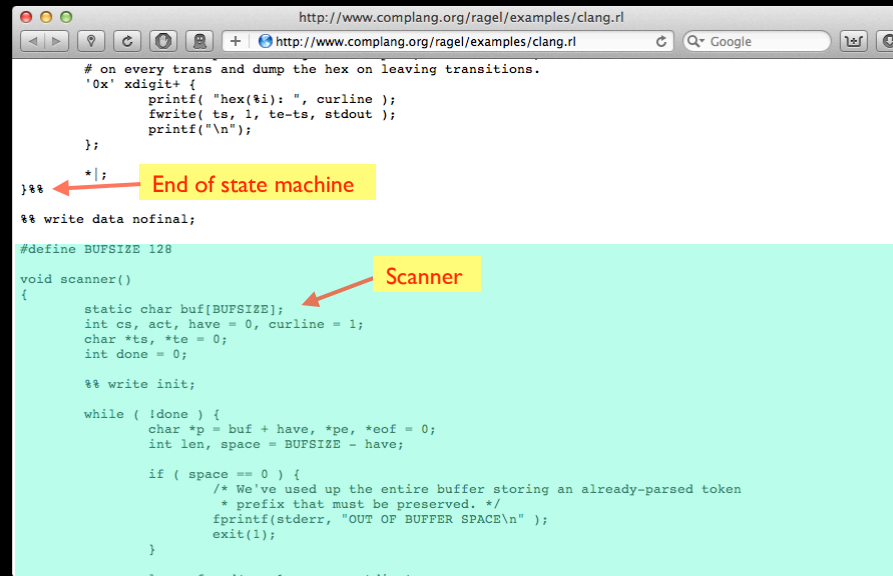
    main := |*

    # Alpha numeric characters or underscore.
    alnum_u = alnum | '_';

    # Symbols. Upon entering clear the buffer. On all transitions
    # buffer a character. Upon leaving dump the symbol.
    ( punct = ['.*'] ) {
        printf("symbol(%i): %c\n", curline, ts[0]);
    };

    # Identifier. Upon entering clear the buffer. On all transitions
    # buffer a character. Upon leaving, dump the identifier.
    alpha_u alnum_u* {
```


Sample “Mini-C” Language Scanner



```
http://www.complang.org/ragel/examples/clang.rl
http://www.complang.org/ragel/examples/clang.rl
# on every trans and dump the hex on leaving transitions.
'0x' xdigit+ {
    printf( "hex(%i): ", curline );
    fwrite( ts, 1, te-ts, stdout );
    printf("\n");
};
};
%% write data nofinal;

#define BUFSIZE 128

void scanner()
{
    static char buf[BUFSIZE];
    int cs, act, have = 0, curline = 1;
    char *ts, *te = 0;
    int done = 0;

    %% write init;

    while ( !done ) {
        char *p = buf + have, *pe, *eof = 0;
        int len, space = BUFSIZE - have;

        if ( space == 0 ) {
            /* We've used up the entire buffer storing an already-parsed token
             * prefix that must be preserved. */
            fprintf(stderr, "OUT OF BUFFER SPACE\n" );
            exit(1);
        }
    }
}
```

Simplified scanner() function

```
void scanner()
{
    // internal scanner( variables
    int len, space = BUFSIZE, curline = 1;
    static char buf[BUFSIZE]; // assume big enough for any input

    // Variables used for communicating with Ragel
    char *p = buf;           // pointer to the character data to process
    char *pe;                // pointer to the end of the data to process
    char *eof;               // same as pe unless using buffered input
    char *ts, *te = 0;       // start of token, end of token
    // Variables used internally by Ragel
    int cs = 0;              // current state
    int act = 0;             // used internally by Ragel

    %% write init;           // initialize state machine

    len = fread( p, 1, space, stdin);
    pe = p + len;
    eof = pe;

    %% write exec;           // run state machine
}
```

Simplified scanner() function

```
void scanner()
{
    // internal scanner( variables
    int len, space = BUFSIZE, curline = 1;
    static char buf[BUFSIZE]; // assume big enough for any input

    // Variables used for communicating with Ragel

    char *p = buf;           // pointer to the character data to process
    char *pe;                // pointer to the end of the data to process
    char *eof;               // same as pe unless using buffered input
    char *t;                 // temporary pointer
    // Variables used for state machine
    int cs;                  // current state
    int act;                 // action

    %% write init;           // initialize state machine

    len = fread( p, 1, space, stdin);
    pe = p + len;
    eof = pe;

    %% write exec;           // run state machine
}
```

Data Pointer – In C code this variable is expected to be a pointer to the character data to process. It should be initialized to the beginning of the data block on every run of the state machine.

Simplified scanner() function

```
void scanner()
{
    // internal scanner( variables
    int len, space = BUFSIZE, curline = 1;
    static char buf[BUFSIZE]; // assume big enough for any input

    // Variables used for communicating with Ragel
    char *p = buf; // pointer to the character data to process
    char *pe; // pointer to the end of the data to process
    char *eof; // same as pe unless using buffered input
    char *ts, *te = 0; // start of token, end of token
    // Variables for state machine
    int cs = 0;
    int act;

    %% write init; // initialize state machine

    len = fread( p, 1, space, stdin);
    pe = p + len;
    eof = pe;

    %% write exec; // run state machine
}
```

Data End Pointer – This should be initialized to *p* plus the data length on every run of the state machine.

Simplified scanner() function

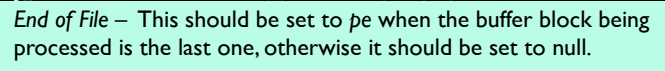
```
void scanner()
{
    // internal scanner( variables
    int len, space = BUFSIZE, curline = 1;
    static char buf[BUFSIZE]; // assume big enough for any input

    // Variables used for communicating with Ragel
    char *p = buf;           // pointer to the character data to process
    char *pe;                // pointer to the end of the data to process
    char *eof;               // same as pe unless using buffered input
    char *t, *te = 0;        // start of token, end of token
    // Variables used internally by Ragel
    int cs;
    int act;

    %% write init;           // initialize state machine

    len = fread( p, 1, space, stdin);
    pe = p + len;
    eof = pe;

    %% write exec;           // run state machine
}
```



Simplified scanner() function

```
void scanner()
{
    // internal scanner( variables
    int len, space = BUFSIZE, curline = 1;
    static char buf[BUFSIZE]; // assume big enough for any input

    // Variables used for communicating with Ragel
    char *p = buf;           // pointer to the character data to process
    char *pe;                // pointer to the end of the data to process
    char *eof;               // same as pe unless using buffered input
    char *ts, *te = 0;       // start of token, end of token
    // Variables used internally by Ragel
    int cs = 0;              // current state
    int act;

    %% write

    len = fread( p, 1, space, stdin);
    pe = p + len;
    eof = pe;

    %% write exec;           // run state machine
}
```

Token Start,Token End – These are set by Ragel as the input is analyzed, they will be used by our custom state machine actions.

Simplified scanner() function

```
void scanner()
{
    // internal scanner( variables
    int len, space = BUFSIZE, curline = 1;
    static char buf[BUFSIZE]; // assume big enough for any input

    // Variables used for communicating with Ragel

    char *p;
    char *p;
    char *e;
    char *t;
    // Variables used for communicating with Ragel
    int cs = 0; // current state
    int act = 0; // used internally by Ragel

    %% write init; // initialize state machine

    len = fread( p, 1, space, stdin);
    pe = p + len;
    eof = pe;

    %% write exec; // run state machine
}
```

The `write init` directive causes Ragel to emit initialization code. This should be executed once before the machine is started. At a very minimum this sets the current state to the start state. If other variables are needed by the generated code, they are also initialized here.

process
of process
input

Simplified scanner() function

```
void scanner()
{
    // internal scanner( variables
    int len, space = BUFSIZE, curline = 1;
    static char buf[BUFSIZE]; // assume big enough for any input

    // Variables used for communicating with Ragel
    char *p = buf;           // pointer to the character data to process
    char *pe;                // pointer to the end of the data to process
    char *eof;               // same as pe unless using buffered input
    char *ts, *te = 0;       // start of token, end of token
    // Variables used internally by Ragel
    int cs = 0;              // current state
    int act = 0;             // used internally by Ragel

    %% write init;           // initialize state machine

    len = fread( p, 1, space, stdin);
    pe = p + len;
    eof = pe;

    %% write exec;           // run state machine
}
```


Simplified scanner() function

```
void scanner()
{
    // internal scanner( variables
    int len, space = BUFSIZE, curline = 1;
    static char buf[BUFSIZE]; // assume big enough for any input

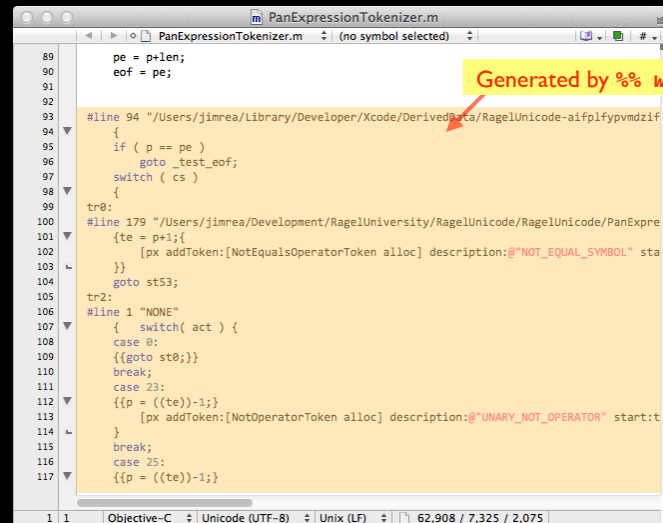
    // Variables used for communicating with Ragel
    char *p = buf;           // pointer to the character data to process
    char *pe;                // pointer to the end of the data to process
    char *eof;               // same as pe unless using buffered input
    char *ts, *te = 0;       // start of token, end of token
    // Variables used internally by Ragel
    int cs = 0;              // current state
    int act = 0;             // used internally by Ragel

    %% write
    len =
    pe = p + len;
    eof = pe;

    %% write exec;          // run state machine
}
```

The write exec directive causes Ragel to emit the state machine's execution code, which may be hundreds or thousands of lines of inline C code.

Code Generation Sample



```
89  pe = p+len;
90  eof = pe;
91
92
93
94  #line 94 "/Users/jimrea/Library/Developer/Xcode/DerivedData/RagelUnicode-aifplfypvmdzif
95  {
96  if ( p == pe )
97      goto _test_eof;
98  switch ( cs )
99  {
100  tr0:
101  #line 179 "/Users/jimrea/Development/RagelUniversity/RagelUnicode/RagelUnicode/PanExpre
102  {te = p+1;{
103  [px addToken:[NotEqualsOperatorToken alloc] description:@"NOT_EQUAL_SYMBOL" sta
104  }}
105  goto st53;
106  tr2:
107  #line 1 "NONE"
108  { switch( act ) {
109  case 0:
110  {{goto st0;}}
111  break;
112  case 23:
113  {{p = ((te))-1;}}
114  [px addToken:[NotOperatorToken alloc] description:@"UNARY_NOT_OPERATOR" start:t
115  }
116  break;
117  case 25:
118  {{p = ((te))-1;}}
```

Generated by %% write exec

1 1 Objective-C ⚡ Unicode (UTF-8) ⚡ Unix (LF) ⚡ 62,908 / 7,325 / 2,075

Setting up the State Machine

imports/includes

state machine

regex action

regex action

regex action

regex action

regex action

regex action

scanner function

Recognizing an Identifier Token

- An identifier token:

Starts with an alphabetic character or underscore,
followed by zero or more alphanumeric or
underscore characters

- Examples:

Name

Zarlon63

_zodiac_kayak

Recognizing an Identifier Token


```
# Alpha numeric characters or underscore.
alnum_u = alnum | '_';

# Alpha characters or underscore.
alpha_u = alpha | '_';

# Identifier
alpha_u alnum_u* {
    printf( "ident(%i): ", curline );
    fwrite( ts, 1, te-ts, stdout );
    printf("\n");
};
```

Recognizing an Identifier Token

```
# Alpha numeric characters or underscore.  
alnum_u = alnum | '_';  
  
# Alpha characters or underscore.  
alpha_u = alpha | '_';  
  
# Identifier  
alpha_u alnum_u* {  
    printf( "ident(%i): ", curline );  
    fwrite( ts, 1, te-ts, stdout );  
    printf("\n");  
};
```



First character of an identifier must be alphabetic or an underscore

Recognizing an Identifier Token

```
# Alpha numeric characters or underscore.  
alnum_u = alnum | '_';  
  
# Alpha characters or underscore.  
alpha_u = alpha | '_';  
  
# Identifier  
alpha_u alnum_u* {  
    printf("ident(%i): ", curline );  
    fwrite( ts, 1, te-ts, stdout );  
    printf("\n");  
};
```

Repeat 0 or more times


Next character must be alphanumeric or an underscore

Processing a Token

```
# Alpha numeric characters or underscore.
alnum_u = alnum | '_';

# Alpha characters or underscore.
alpha_u = alpha | '_';

# Identifier
alpha_u alnum_u* {
    printf( "ident(%i): ", curline );
    fwrite( ts, 1, te-ts, stdout );
    printf("\n");
};
```




Action to take for this type of token

Processing a Token

```
# Alpha numeric characters or underscore.
alnum_u = alnum | '_';

# Alpha characters or underscore.
alpha_u = alpha | '_';

# Identifier
alpha_u alnum_u* {
    printf( "ident(%i): ", curline );
    fwrite( ts, 1, te-ts, stdout );
    printf("\n");
};
```



Print "ident" + line number

Processing a Token

```
# Alpha numeric characters or underscore.
alnum_u = alnum | '_';

# Alpha characters or underscore.
alpha_u = alpha | '_';

# Identifier
alpha_u | alnum_u* {
    printf( "ident(%i): ", curline );
    fwrite( ts, 1, te-ts, stdout );
    printf( "\n" );
};
```

Print identifier token

Length of token

Pointer to start of token

Mini-C Scanner Operation

Input:

```
( Fahrenheit - 32 ) * 5 / 9
```

Output:

```
symbol(1):(
ident(1):Fahrenheit
symbol(1):-
int(1):32
symbol(1):)
symbol(1):*
int(1):5
symbol(1):/
int(1):9
```

Built in Regex Expressions

```
# Alpha numeric characters or underscore.  
alnum_u = alnum | '_';
```

```
# Alpha characters or underscore.  
alpha_u = [A-Za-z] | '_';
```

You can use standard Regex terms instead

Built-in “machine” for alphabetic characters

Built in Regex Expressions

any	<i>any character in the alphabet</i>	
alpha	<i>alphabetic characters</i>	[A-Za-z]
digit	<i>numeric digits</i>	[0-9]
alnum	<i>alphanumeric characters</i>	[A-Za-z0-9]
lower	<i>lower case alphabetic characters</i>	[a-z]
upper	<i>upper case alphabetic characters</i>	[A-Z]
xdigit	<i>hexadecimal digits</i>	[0-9A-Fa-f]
cntrl	<i>control characters</i>	0..31
graph	<i>graphical characters</i>	[!~]
print	<i>printable characters</i>	[-~]
punct	<i>punctuation characters</i>	[!-/:-@[~]
space	<i>whitespace</i>	[\t\v\f\n\r]
zlen	<i>zero length</i>	""
empty	<i>empty set (matches nothing)</i>	^any

Terminology

Ragel calls regular expressions “machines”

Regex Learning Resources

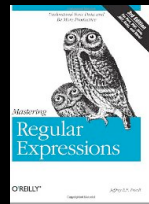
1) Online

<http://www.regular-expressions.info/tutorial.html>

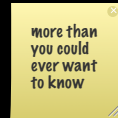
http://en.wikipedia.org/wiki/Regular_expression

<http://perldoc.perl.org/perlretut.html>

2) Books



Mastering
Regular
Expressions
Jeffery E.F. Friedl



Regular
Expressions
Cookbook
*Jan Goyvaerts,
Steven Levithan*

Ragel, Unicode & Objective-C

Classes

- RGEExpression
- ExpressionToken

RGExpression Class

Instance Variables

- Text
- Tokens (NSMutableArray)

Methods

- Convert text to tokens (calls Ragel)
- Add token to array (called from Ragel)
- Evaluate expression value

ExpressionToken Class

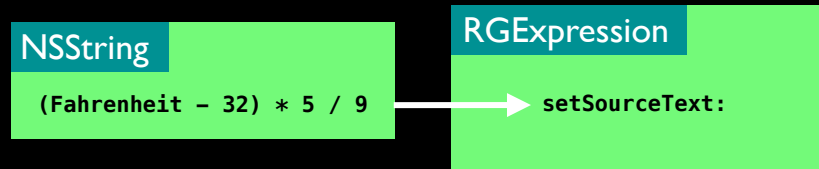
Instance Variables

- Text
- Value (type depends on subclass)

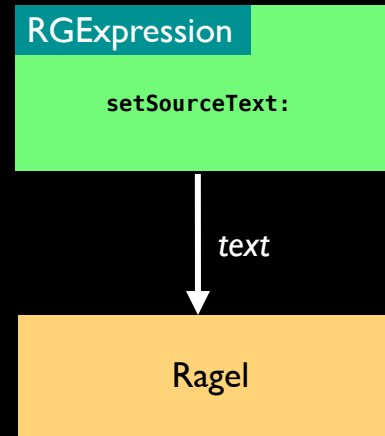
Subclasses

- IdentifierToken
- IntegerToken
- StringToken
- FloatToken

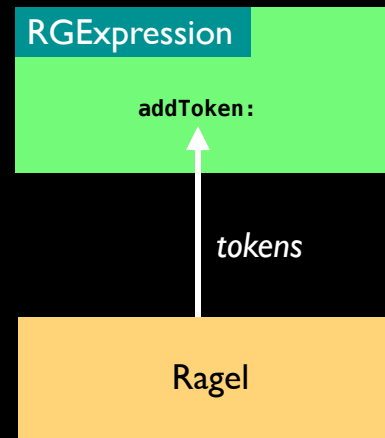
Ragel, Unicode & Objective-C



Ragel, Unicode & Objective-C



Ragel, Unicode & Objective-C



RGExpression Class

```
@interface RGExpression : NSObject {
    NSString * sourceText;
    NSData * sourceUnicode;
    NSMutableArray * tokens;
}

@implementation RGExpression

- (void) setSourceText:(NSString *) text
{
    sourceText = [text copy];
    sourceUnicode =
        [sourceText dataUsingEncoding:NSUTF16LittleEndianStringEncoding];
    short * sourceBase = (short *)[sourceUnicode bytes];
    unsigned long sourceLen = [sourceUnicode length]/2;
    tokens = [[NSMutableArray alloc] init];
    RGExpressionScanner(self, sourceBase, sourceLen);
}
```

RGExpression Class

```
@interface RGExpression : NSObject {
    NSString * sourceText;
    NSData * sourceUnicode;
    NSMutableArray * tokens;
}

@implementation RGExpression

- (void) setSourceText:(NSString *) text
{
    sourceText = [text copy];
    sourceUnicode =
        [sourceText dataUsingEncoding:NSUTF16LittleEndianStringEncoding];
    short * sourceBase = (short *)[sourceUnicode bytes];
    unsigned long sourceLen = [sourceUnicode length]/2;
    tokens = [[NSMutableArray alloc] init];
    RGExpressionScanner(self, sourceBase, sourceLen);
}
}
```


RGExpression Class

```
@interface RGExpression : NSObject {
    NSString * sourceText;
    NSData * sourceUnicode;
    NSMutableArray * tokens;
}

@implementation RGExpression

- (void) setSourceText:(NSString *) text
{
    sourceText = [text copy];
    sourceUnicode =
        [sourceText dataUsingEncoding:NSUTF16LittleEndianStringEncoding];
    short * sourceBase = (short *)[sourceUnicode bytes];
    unsigned long sourceLen = [sourceUnicode length]/2;
    tokens = [[NSMutableArray alloc] init];
    RGExpressionScanner(self, sourceBase, sourceLen);
}
```

RGExpression Class

```
@interface RGExpression : NSObject {
    NSString * sourceText;
    NSData * sourceUnicode;
    NSMutableArray * tokens;
}

@implementation RGExpression

- (void) setSourceText:(NSString *) text
{
    sourceText = [text copy];
    sourceUnicode =
        [sourceText dataUsingEncoding:NSUTF16LittleEndianStringEncoding];
    short * sourceBase = (short *)[sourceUnicode bytes];
    unsigned long sourceLen = [sourceUnicode length]/2;
    tokens = [[NSMutableArray alloc] init];
    RGExpressionScanner(self, sourceBase, sourceLen);
}
```

RGExpression Class

```
@interface RGExpression : NSObject {
    NSString * sourceText;
    NSData * sourceUnicode;
    NSMutableArray * tokens;
}

@implementation RGExpression

- (void) setSourceText:(NSString *) text
{
    sourceText = [text copy];
    sourceUnicode =
        [sourceText dataUsingEncoding:NSUTF16LittleEndianStringEncoding];
    short * sourceBase = (short *) [sourceUnicode bytes];
    unsigned long sourceLen = [sourceUnicode length]/2;
    tokens = [[NSMutableArray alloc] init];
    RGExpressionScanner(self, sourceBase, sourceLen);
}
```

For callbacks

Pointer to UTF-16

Number of characters

Ragel Scanner

```
%%{  
  machine clang;  
  alphatype short;  
  
  ... state machine definitions ...  
}%  
  
%% write data nofinal;  
  
void RgExpressionScanner(PanExpression *px, short *p, unsigned long len);  
void RgExpressionScanner(PanExpression *px, short *p, unsigned long len)  
{  
  short *pe, *eof;  
  int cs, act, curline = 1;  
  short *ts, *te = 0;  
  
  %% write init;  
  
  pe = p+len;  
  eof = pe;  
  
  %% write exec;  
}
```

Ragel Scanner

```
%%{  
  machine clang;  
  alphatype short;  
  ... state machine definitions ...  
}%  
  
%% write data nofinal;  
  
void RgExpressionScanner(PanExpression *px, short *p, unsigned long len);  
void RgExpressionScanner(PanExpression *px, short *p, unsigned long len)  
{  
  short *pe, *eof;  
  int cs, act, curline = 1;  
  short *ts, *te = 0;  
  
  %% write init;  
  
  pe = p+len;  
  eof = pe;  
  
  %% write exec;  
}
```

Tells Ragel to use UTF-16 instead of ASCII

Ragel Scanner

```
%%{  
  machine clang;  
  alphatype short;  
  
  ... state machine definitions ...  
}%  
  
%% write data nofinal;  
  
void RgExpressionScanner(PanExpression *x, short *p, unsigned long len);  
void RgExpressionScanner(PanExpression *px, short *p, unsigned long len)  
{  
  short *pe, *eof;  
  int cs, act, curline = 1;  
  short *ts, *te = 0;  
  
  %% write init;  
  
  pe = p+len;  
  eof = pe;  
  
  %% write exec;  
}
```

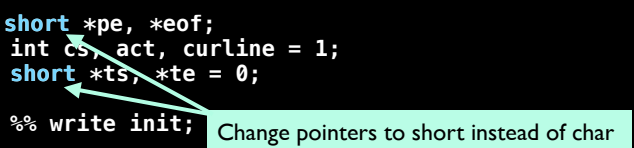
For callbacks

Number of characters

Pointer to UTF-16

Ragel Scanner

```
%%{  
  machine clang;  
  alphatype short;  
  
  ... state machine definitions ...  
}%  
  
%% write data nofinal;  
  
void RgExpressionScanner(PanExpression *px, short *p, unsigned long len);  
void RgExpressionScanner(PanExpression *px, short *p, unsigned long len)  
{  
  short *pe, *eof;  
  int cs, act, curline = 1;  
  short *ts, *te = 0;  
  %% write init;  
  pe = p+len;  
  eof = pe;  
  %% write exec;  
}
```



A diagram with a light blue rectangular box containing the text "Change pointers to short instead of char". Two green arrows originate from this box. One arrow points to the line "short *ts, *te = 0;" in the code. The other arrow points to the line "short *pe, *eof;" in the code.

Add Token to RGExpression

```
# Alpha numeric characters or underscore.
alnum_u = alnum | '_';

# Alpha characters or underscore.
alpha_u = alpha | '_';

# Identifier
alpha_u alnum_u* {
  [px addToken:[IdentifierToken alloc] start:ts end: te];
};
```


Add Token to RGExpression

```
# Alpha numeric characters or underscore.  
alnum_u = alnum | '_';
```

```
# Alpha characters or underscore.  
alpha_u = alpha | '_';
```

```
# Identifier  
alpha_u alnum_u* {  
  [px addToken:[IdentifierToken alloc] start:ts end:te];  
};
```

Start of token (in UTF-16 source)

End of token (in UTF-16 source)

Create uninitialized token object

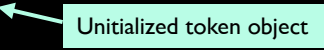
Callback to RGExpression object

Add Token to RGExpression

```
- (void) addToken:(ExpressionToken *)token
    start:(short *)ts
    end:(short *)te
{
    short * sourceBase = (short *)[sourceUnicode bytes];
    long tokenSpot = ts-sourceBase;
    long tokenLength = te-ts;
    NSString * xtext = [[NSString alloc]
        initWithBytes:(const void *)ts
        length: tokenLength*2
        encoding:NSUTF16LittleEndianStringEncoding];
    token = [token initWithExpression:self
        Token:xtext
        Start:tokenSpot
        Len:tokenLength];
    [tokens addObject:token];
}
```

Add Token to RGExpression

```
- (void) addToken:(ExpressionToken *)token
    start:(short *)ts
    end:(short *)te
{
    short * sourceBase = (short *)[sourceUnicode bytes];
    long tokenSpot = ts-sourceBase;
    long tokenLength = te-ts;
    NSString * xtext = [[NSString alloc]
        initWithBytes:(const void *)ts
        length: tokenLength*2
        encoding:NSUTF16LittleEndianStringEncoding];
    token = [token initWithExpression:self
        Token:xtext
        Start:tokenSpot
        Len:tokenLength];
    [tokens addObject:token];
}
```



Unitialized token object

Add Token to RGExpression


```
- (void) addToken:(ExpressionToken *)token
    start:(short *)ts
    end:(short *)te
{
    short * sourceBase = (short *)[sourceUnicode bytes];
    long tokenSpot = ts-sourceBase;
    long tokenLength = te-ts;
    NSString * xtext = [[NSString alloc]
        initWithBytes:(const void *)ts
        length: tokenLength*2
        encoding:NSUTF16LittleEndianStringEncoding];
    token = [token initWithExpression:self
        Token:xtext
        Start:tokenSpot
        Len:tokenLength];
    [tokens addObject:token];
}
```

Start of token (in UTF-16 source)

End of token (in UTF-16 source)

Add Token to RGExpression

```
- (void) addToken:(ExpressionToken *)token
    start:(short *)ts
    end:(short *)te
{
    short * sourceBase = (short *)[sourceUnicode bytes];
    long tokenSpot = ts-sourceBase;
    long tokenLength = te-ts;
    NSString * xtext = [[NSString alloc]
                        initWithBytes:(const void *)ts
                        length: tokenLength*2
                        encoding:NSUTF16LittleEndianStringEncoding];
    token = [token initWithExpression:self
                                Token:xtext
                                start:tokenSpot
                                Len:tokenLength];
    [tokens addObject:token];
}
```




String value of token

Add Token to RGExpression

```
- (void) addToken:(ExpressionToken *)token
    start:(short *)ts
    end:(short *)te
{
    short * sourceBase = (short *)[sourceUnicode bytes];
    long tokenSpot = ts-sourceBase;
    long tokenLength = te-ts;
    NSString * xtext = [[NSString alloc]
        initWithBytes:(const void *)ts
        length: tokenLength*2
        encoding:NSUTF16LittleEndianStringEncoding];
    token = [token initWithExpression:self
        Token:xtext
        Start:tokenSpot
        Len:tokenLength];
    [tokens addObject:token];
}
```

Initialize the token



Add Token to RGExpression

```
- (void) addToken:(ExpressionToken *)token
    start:(short *)ts
    end:(short *)te
{
    short * sourceBase = (short *)[sourceUnicode bytes];
    long tokenSpot = ts-sourceBase;
    long tokenLength = te-ts;
    NSString * xtext = [[NSString alloc]
        initWithBytes:(const void *)ts
        length: tokenLength*2
        encoding:NSUTF16LittleEndianStringEncoding];
    token = [token initWithExpression:self
        Token:xtext
        Start:tokenSpot
        Len:tokenLength];

    [tokens addObject:token];
}
```



Add token to array

Typical Operation

Input (NSString):

(Fahrenheit - 32) * 5 / 9

Output (NSMutableArray):

Symbol (Identifier Fahrenheit	Symbol -	Integer 32	Symbol)	Symbol *	Integer 5	Symbol /	Integer 9
-------------	--------------------------	-------------	---------------	-------------	-------------	--------------	-------------	--------------

Scanning for Unicode Characters

```
# Pi (Option-P)  
 $\pi$  | [] {  
  [px addToken:[FloatOperandToken alloc] start:ts end: te];  
};
```

Won't work, Ragel does not allow Unicode in source

Scanning for Unicode Characters

```
# Pi (Option-P)
0x3C0 | 0x220F {
    [px addToken:[FloatOperandToken alloc] start:ts end: te];
};
```

Scanning for Unicode Characters

π Π

```
# Pi (Option-P)  
0x3C0 | 0x220F {  
    [px addToken:[FloatOperandToken alloc] start:ts end: te];  
};
```

Scanning for Unicode Characters

```
# Pi (Option-P)
0x3C0 | 0x220F {
    [px addToken:[FloatOperandToken alloc] start:ts end: te];
};
```

Token Initialization

```
@implementation FloatOperandToken
@synthesize value;

- (id) initWithExpression:(PanExpression *)xpr
    Token:(NSString *)token
    Start:(long)start
    Len:(long)len;
{
    if (self = [super initWithExpression:xpr
        Token:token Start:start Len:len]) {
        dataType = FLOAT;
        NSString *text =
            [[expression getSourceText] substringWithRange:tokenStartEnd];
        if ([text isEqualToString:@"π"] || [text isEqualToString:@"[]"]) {
            value = M_PI;
        } else {
            value = [text doubleValue];
        }
    }
    return self;
}
```

Token Initialization

```
@implementation FloatOperandToken
@synthesize value;

- (id) initWithExpression:(PanExpression *)xpr
              Token:(NSString *)token
              Start:(long)start
              Len:(long)len;
{
    if (self = [super initWithExpression:xpr
                                      Token:token Start:start Len:len]) {
        dataType = FLOAT;
        NSString *text =
            [[expression getSourceText] substringWithRange:tokenStartEnd];
        if ([text isEqualToString:@"π"] || [text isEqualToString:@"[]"]) {
            value = M_PI;
        } else {
            value = [text doubleValue];
        }
    }
    return self;
}
```

Token Initialization

```
@implementation FloatOperandToken
@synthesize value;

- (id) initWithExpression:(PanExpression *)xpr
              Token:(NSString *)token
              Start:(long)start
              Len:(long)len;
{
    if (self = [super initWithExpression:xpr
                                      Token:token Start:start Len:len]) {
        dataType = FLOAT;
        NSString *text =
            [[expression getSourceText] substringWithRange:tokenStartEnd];
        if ([text isEqualToString:@"π"] || [text isEqualToString:@"[]"]) {
            value = M_PI;
        } else {
            value = [text doubleValue];
        }
    }
    return self;
}
```

Token Initialization

```
@implementation FloatOperandToken
@synthesize value;

- (id) initWithExpression:(PanExpression *)xpr
                Token:(NSString *)token
                Start:(long)start
                Len:(long)len;
{
    if (self = [super initWithExpression:xpr
                                Token:token Start:start Len:len]) {
        dataType = FLOAT;
        NSString *text =
            [[expression getSourceText] substringWithRange:tokenStartEnd];
        if ([text isEqualToString:@"π"] || [text isEqualToString:@"Π"]) {
            value = M_PI;
        } else {
            value = [text doubleValue];
        }
    }
    return self;
}
```


Token Initialization

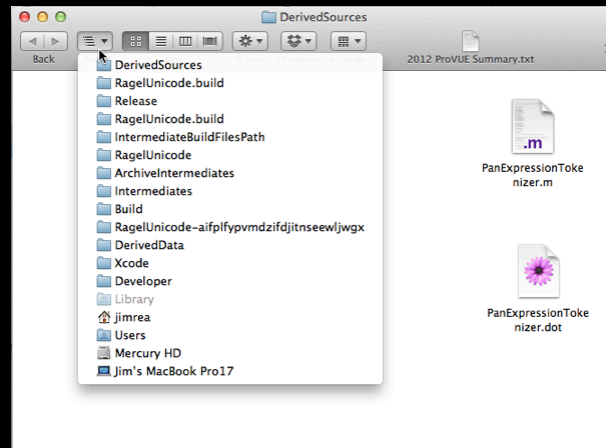
```
@implementation FloatOperandToken
@synthesize value;

- (id) initWithExpression:(PanExpression *)xpr
    Token:(NSString *)token
    Start:(long)start
    Len:(long)len;
{
    if (self = [super initWithExpression:xpr
        Token:token Start:start Len:len]) {
        dataType = FLOAT;
        NSString *text =
            [[expression getSourceText] substringWithRange:tokenStartEnd];
        if ([text isEqualToString:@"π"] || [text isEqualToString:@"[]"]) {
            value = M_PI;
        } else {
            value = [text doubleValue];
        }
    }
    return self;
}
```

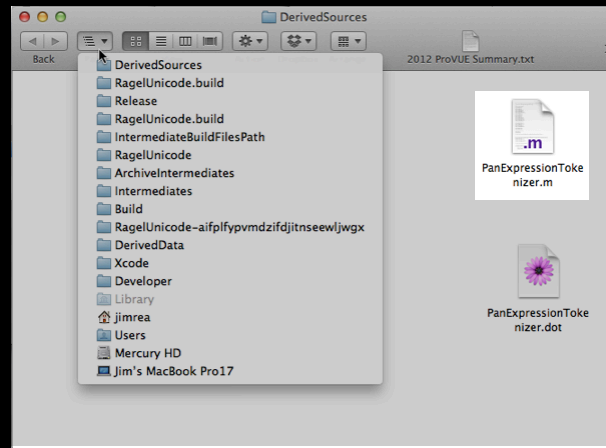
Floating Point Constant

```
# Floating Point Number
(digit+ '.' digit+ ('e' ('-'|'+')? digit+)? | (digit+ 'e' ('-'|'+')?
digit+ ) {
    [px addToken:[FloatOperandToken alloc] start:ts end: te];
};
```

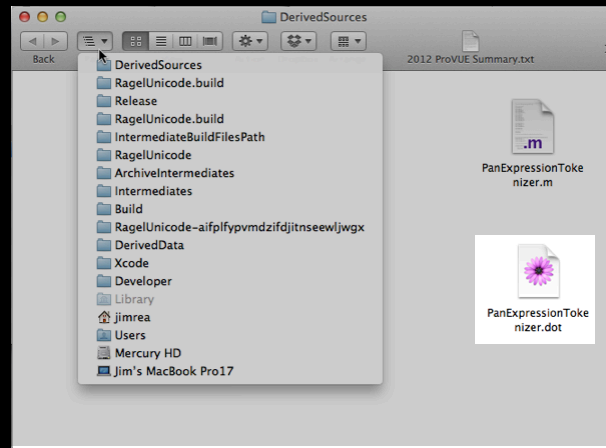
Ragel Outputs



Ragel Outputs



Ragel Outputs



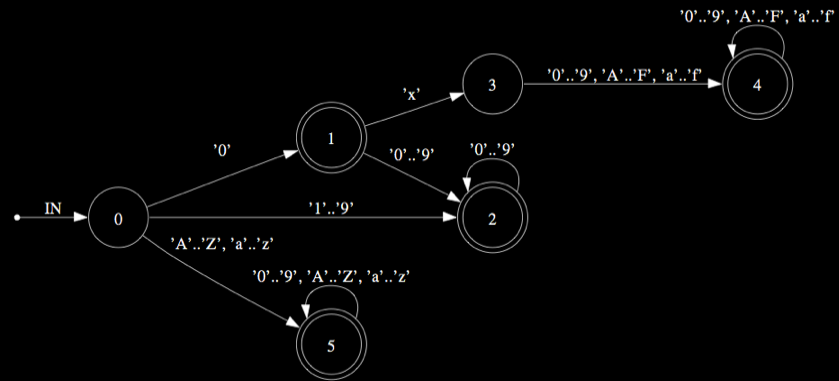
GraphViz

Open source graph visualization software

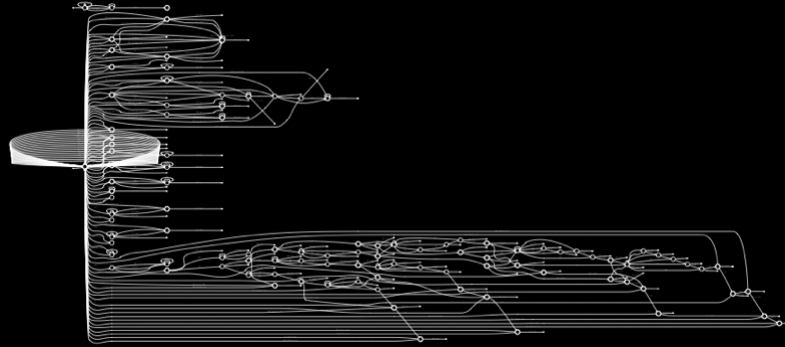
<http://www.graphviz.org>



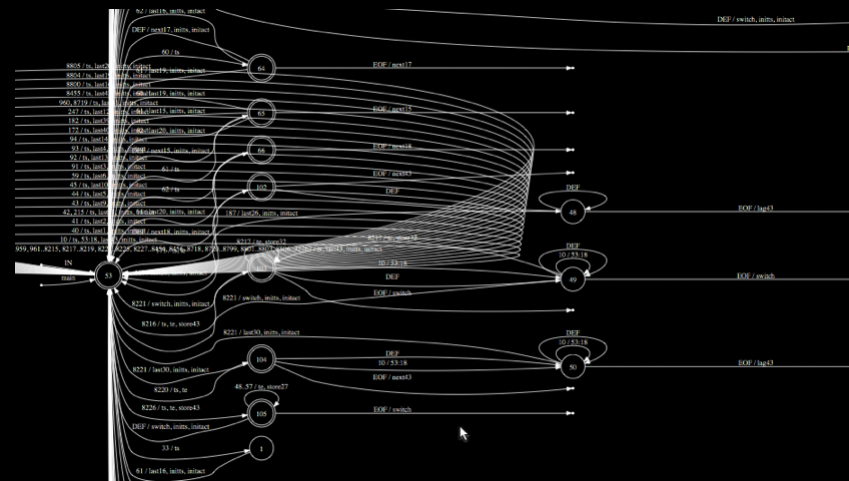
GraphViz



GraphViz



GraphViz



Resources

1) Ragel Home Page

<http://www.complang.org/ragel/>

2) Ragel User Guide

<http://www.complang.org/ragel/ragel-guide-6.6.pdf>

3) My Google+ post on using Ragel with Xcode

<https://plus.google.com/109028627294998653069/posts/C6d1YkcY8kW>

Credits

Ragel was written by Dr. Adrian Thurston, a software researcher at a Canadian network security firm. It was originally developed in early 2000 and was first released January 2002. Many people have contributed feedback, ideas and code.