# 1 Introduction

The first project is to implement a scanner (or lexer or tokenizer) for LangF, which will convert an input stream of characters into a stream of tokens. While such components of a compiler are often best written using a *lexical analyzer generator* (e.g., Lex or Flex (for C), JLex (for Java), ML-Lex or ML-ULex (for Standard ML), Alex (for Haskell), SILex (for Scheme)), you will write a scanner by hand.

# 2 LangF Lexical Conventions

LangF has four classes of *tokens*: keywords, delimiter and operator symbols, identifiers, and integer and string literals. Tokens can be separated by *whitespace* and/or *comments*.

LangF has the following set of keywords:

| | | | |
|---|---|---|---|
| **and** | **andalso** | **case** | **datatype** |
| **else** | **end** | **fn** | **fun** |
| **if** | **in** | **let** | **of** |
| **orelse** | **then** | **type** | **val** |

A keyword may not be used as an identifier.

LangF also has the following collection of delimiter and operator symbols:

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **+** | **−** | **\*** | **/** | **%** | **~** | **==** | **<>** | **<=** |
| **<** | **>=** | **>** | **^** | **(** | **)** | **[** | **]** | **{** |
| **}** | **−>** | **=>** | **=** | **:** | **,** | **;** | **\|** | **_** |

LangF has three classes of *identifiers*: type variables, constructors, and variables. An identifier in LangF can be any string of letters, digits, underscores, and (single) quote marks, with the following restrictions:

- a type variable must begin with a single quote character and a lower-case letter (e.g., `'alpha`)

- a constructor must begin with an upper-case letter (e.g., `List`)

- a variable must begin with a lower-case letter (e.g., `length`)

Identifiers are case-sensitive (e.g., `treemap` is different from `treeMap`).

Integer literals in LangF are written using decimal notation, with an optional negation symbol "**~**".

String literals in LangF start with a double quote """", followed by zero or more printable characters (ASCII codes 32 – 126), and terminated with a double quote. A string literal can contain the following C-like escape sequences:

| | | |
|---|---|---|
| \a | — | bell (ASCII code 7) |
| \b | — | backspace (ASCII code 8) |
| \f | — | form feed (ASCII code 12) |
| \n | — | newline (ASCII code 10) |
| \r | — | carriage return (ASCII code 13) |
| \t | — | horizontal tab (ASCII code 9) |
| \v | — | vertical tab (ASCII code 11) |
| \\ | — | backslash |
| \" | — | double quote |

A character in a string literal may also be specified by its numerical value (ASCII code) using the escape sequence "\$ddd$", where $ddd$ is a sequence of exactly three decimal digits. Strings in LangF may contain any 8-bit character, including embedded zeros, which can be specified as "\000".

Comments in LangF start anywhere outside a string literal with a " (*" and are terminated with a matching "*) ". Comments in LangF may be nested.

Whitespace in LangF is any non-empty sequence of spaces (ASCII code 32), horizontal tabs, newlines, vertical tabs, form feeds, or carriage returns.

# 3 Requirements

Your implementation should include (at least) the following module:

```
structure LangFScanner : LANGF_SCANNER
```

where the LANGF_SCANNER signature is as follows:

```
signature LANGF_SCANNER =
sig
   val scan : {reportError: string -> unit} ->
              (char, 'strm) StringCvt.reader ->
              (Tokens.token, 'strm) StringCvt.reader
end
```

The StringCvt.reader type is defined in the SML Basis Library as follows:

```
type ('item,'strm) reader = 'strm -> ('item * 'strm) option
```

A reader is a function that takes a stream and returns either SOME pair of the next item and the rest of the stream, or NONE when the end of the stream is reached. Thus scan is a function that takes a character reader and returns a token reader.

The **structure** Tokens : TOKENS module will be provided; the TOKENS signature is as follows:

```
signature TOKENS =
sig
    datatype token =
        KW_and
      | KW_andalso
      | ...

      | KW_type
      | KW_val
      | PLUS | MINUS | ASTERISK | SLASH | PERCENT | TILDE
      | EQEQ | LTGT | LTEQ | LT | GTEQ | GT
      | CARET
      | LPAREN (* ( *) | RPAREN (* ) *)
      | LSBRACK (* [ *) | RSBRACK (* ] *)
      | LCBRACK (* { *) | RCBRACK (* } *)
      | MINUS_ARROW (* -> *) | EQ_ARROW (* => *)
      | EQ | COLON | COMMA | SEMI | VBAR (* / *) | UNDERSCORE
      | TYVAR_NAME of Atom.atom
      | CON_NAME of Atom.atom
      | VAR_NAME of Atom.atom
      | INTEGER of IntInf.int
      | STRING of String.string

    val toString : token -> string
end
```

The tokens correspond to the various keywords, delimiter and operator symbols, identifiers, and literals. The identifier tokens (TYVAR_NAME, CON_NAME, and VAR_NAME) carry a unique string representation of the identifier. The **structure** Atom : ATOM module is provided by the SML/NJ Library[1]; you can view the ATOM signature and the Atom structure implementation via links in the HTML version of this document.

## 3.1 Errors

To support error reporting, the LangFScanner.scan function takes an initial argument of the type {reportError: string -> unit}. Characters in the input stream that cannot be recognized as part of a token or whitespace or part of a comment should be reported and discarded. For example, if the input program is the (almost) expression:

```
1 + @ + 3
```

then the scanner should report an error (for example Error: bad character '@') *and* return the following stream of tokens:

```
INTEGER (1)
+
+
INTEGER (3)
```

---

[1]The SML/NJ Library a collection of utility libraries that are distributed with SML/NJ (and available in some other SML implementations). To access the library, use $/smlnj-lib.cm in a CM file.

Any non-printable character that is not whitespace should be reported as an error and discarded from the character stream. Similarly, any printable character that is not part of a token, is not whitespace, and is not within a comment should be reported as an error and discarded from the character stream. Invalid escape sequences in strings should be reported as an error and discarded from the character stream. Finally, string literals and comments that are unterminated at the end of the character stream (that is, when the input character reader returns `NONE`) should be reported as errors.

## 4   GForge and Submission

We have set up GForge accounts (using your CNetId) and an SVN repository for each student on the GForge server. To *checkout* a copy of your repository, run the command:

```
svn checkout --username cnetid https://cs22610.cs.uchicago.edu/svn/cnetid-proj
```

On your first checkout, you should be prompted for your GForge password. Upon a successful checkout, a directory called `cnetid-proj` will be created in the current directory. Over the course of the quarter, each of the four projects will be stored as sub-directories within the `cnetid-proj` directory. Sources for Project 1 have been (or will shortly be) committed to your repository in the `project1` sub-directory. If you checkout your repository before the Project 1 sources have been committed, then you will need to *update* your local copy, by running the command:

```
svn update
```

from the `cnetid-proj` directory.

We will collect projects from the SVN repositories at 10pm on Friday, January 23; make sure that you have committed your final version before then. To do so, run the command:

```
svn commit
```

from the `cnetid-proj` directory.

## 5   Hints

- To complete the assignment, you should only need to make changes to the `cnetid-proj/project1/langfc-src/scanner/langf-scanner.sml` file.

- Work on error reporting last. Detecting errors and producing good error messages can be difficult; it is more important for your scanner to work on good programs than for it to "work" on bad programs.

- Executing the compiler (from the `cnetid-proj/project1` directory) with the command

  ```
  ./bin/langfc -Ckeep-scan=true file.lgf
  ```

  will produce a `file.scan.toks` file that contains the sequence of tokens returned by the scanner. Use this control and its output to check that your scanner is working as expected. The `tests/scanner` directory includes a number of tests (of increasing complexity); for each `testNN.lgf` file, there is a `testNN.out` file containing the sequence of tokens to be returned by the scanner and, if the test has lexical errors, a `testNN.err` file containing sample error messages to be reported by the scanner.

## Document History

**January 9, 2009**

Section 3: NAME **of** Atom.atom ⇒ VAR_NAME **of** Atom.atom

Section 3.1: Any non-printable character should be reported as an error and discarded from the character stream. $\Rightarrow$ Any non-printable character that is not whitespace should be reported as an error and discarded from the character stream.

**January 8, 2009**  Original version