

Lexing

CS 132: Compiler Design

January 24, 2011

DEFINITIONS

- ✓ What is lexing (a.k.a. tokenizing)?

DEFINITIONS

- ✓ What is lexing (a.k.a. tokenizing)?
- ✓ What is parsing?

DEFINITIONS

- ✓ What is lexing (a.k.a. tokenizing)?
- ✓ What is parsing?
- ✓ What is abstract syntax?

CONCRETE SYNTAXES

Pascal vs. C vs. Python

```
if space_left <> 0 then
  begin
    processLine(line);
    getNextLine
  end
else
  begin
    getNextLine;
    skipped := 1
  end
```

```
if (space_left != 0)
  {
    processLine(line);
    getNextLine();
  }
else
  {
    getNextLine();
    skipped = 1;
  }
```

```
if (space_left != 0):
  processLine(line)
  getNextLine()
else:
  getNextLine()
  skipped = 1
```

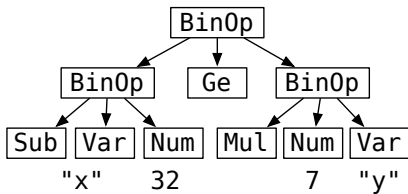
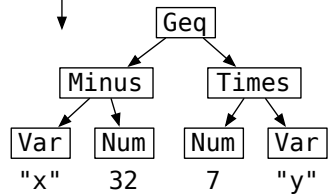
TRADITIONAL PARSING

(x - 3 2) > = 7 * y

LEXING

LPAREN ID DASH INT RPAREN GEQ NUM STAR ID
 "x" 32 7 "y"

PARSING



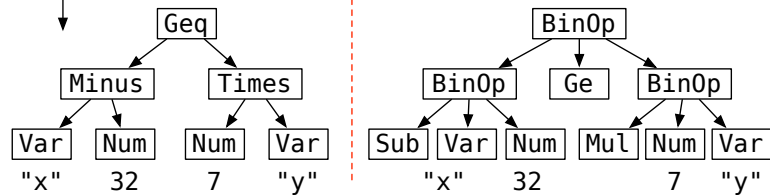
TRADITIONAL PARSING

(x - 3 2) > = 7 * y

LEXING

LPAREN ID DASH INT RPAREN GEQ NUM STAR ID
 "x" 32 7 "y"

PARSING



Question: Why bother with this two-step process?

REGULAR EXPRESSION REVIEW

- ✓ What are the usual "primitive" regular expressions?

M[ou]'?am+[ae]r .*([AEae]l[-])?[GKQ] h?[aeu]+([dtz][dhz]?)af[iy]

M[ou]'?am+[ae]r .*([AEae]l[-])?[GKQ] h?[aeu]+([dtz][dhz]?)+af[iy]

Muammar Qaddafi

Mo'amar Gadhafi

Muammar Kaddafi

Muammar Qadhafi

Moammar El Kadhafi

Muammar Gadafi

Mu'amar al-Qadafi

Moamer El Kazzafi

Moamar al-Gaddafi

Mu'amar Al Qathafi

Muammar Al Qathafi

Mo'amar el-Gadhafi

Moamar El Kadhafi

Muammar al-Qadhafi

Mu'amar al-Qadhdhafi

Mu'amar Qadafi

Moamar Gaddafi

Mu'amar Qadhdhafi

Muammar al-Khaddafi

Mu'amar al-Kadafi

Muammar Ghaddafy

Muammar Ghadafi

Muammar Ghaddafi

Muamar Kaddafi

Muammar Quathafi

Muammar Ghaddafi

Muamar Al-Kaddafi

Moammar Khadafy

Moammar Qudhafi

Mu'amar al-Qaddafi

Mu'amar Muhammad Abu Minyar al-Qadhafi

APPLYING REGULAR EXPRESSIONS

Give a regular expression for Ada identifiers, which

- ✓ Can contain letters, digits and underscores
- ✓ Begin with a letter
- ✓ Have no consecutive underscores or an underscore at the end

Give a regular expression for dollar amounts on US checks. (Note: many printed checks contain extra asterisks to prevent fraud, i.e., \$****1.00.)



FINITE AUTOMATA

Draw NFAs and DFAs for the following regular expressions:

- ✓ If keyword: `if`
- ✓ Case-insensitive if keyword: `[Ii][Ff]`
- ✓ Identifiers: `[a-zA-z_][a-zA-Z_0-9]*`

PRACTICAL ISSUES

- ✓ Keywords: `if` `int` `end`
 - ✓ Identifiers: `[a-zA-z_][a-zA-Z_0-9]*`
-

- ✓ Should `intend` be one token or two?

PRACTICAL ISSUES

- ✓ Keywords: `if` `int` `end`
 - ✓ Identifiers: `[a-zA-z_][a-zA-Z_0-9]*`
-

- ✓ Should `intend` be one token or two?
- ✓ Is `if` a keyword or an identifier?

PRACTICAL ISSUES

- ✓ Keywords: `if` `int` `end`
 - ✓ Identifiers: `[a-zA-z_][a-zA-Z_0-9]*`
-

- ✓ Should `intend` be one token or two?
- ✓ Is `if` a keyword or an identifier?
- ✓ How many tokens are there in `3-2` ?

PRACTICAL ISSUES

- ✓ Keywords: `if` `int` `end`
 - ✓ Identifiers: `[a-zA-z_][a-zA-Z_0-9]*`
-

- ✓ Should `intend` be one token or two?
- ✓ Is `if` a keyword or an identifier?
- ✓ How many tokens are there in `3-2` ?
- ✓ How many tokens are there in `3 + -2` ?

PRACTICAL ISSUES

- ✓ Keywords: `if` `int` `end`
 - ✓ Identifiers: `[a-zA-z_][a-zA-Z_0-9]*`
-

- ✓ Should `intend` be one token or two?
- ✓ Is `if` a keyword or an identifier?
- ✓ How many tokens are there in `3-2` ?
- ✓ How many tokens are there in `3 + -2` ?
- ✓ What about `3-2147483648` ?

NOTE ON RESERVED WORDS

Many languages have reserved words (*keywords*) that cannot be used as identifiers (`if`, `return`, etc.)

NOTE ON RESERVED WORDS

Many languages have reserved words (*keywords*) that cannot be used as identifiers (*if, return, etc.*)

On the other hand, there are languages like PL/I:

IF THEN THEN THEN = ELSE; ELSE ELSE = THEN;

BUILDING A LEXER

Pascal included tokens of the following form:

- ✓ Integers: $[0-9]^+$
- ✓ Real numbers: $[0-9]^+ \cdot [0-9]^+$
- ✓ Range marker: $..$

How do we combine them into a lexer?

BUILDING A LEXER

Pascal included tokens of the following form:

- ✓ Integers: $[0-9]^+$
- ✓ Real numbers: $[0-9]^+ \cdot [0-9]^+$
- ✓ Range marker: $..$

How do we combine them into a lexer?

What is the worst-case running time for a DFA-based lexer?

LEXING CHALLENGES: FORTRAN

- ✓ Pre-1977 Hollerith Constants

15HCOMPILER DESIGN

LEXING CHALLENGES: FORTRAN

- ✓ Pre-1977 Hollerith Constants

```
15HCOMPILER DESIGN
```

- ✓ Fixed-Column Layouts

```
CALL DOIT(X,Y,Z)
```

```
CALL DOIT(X,Y,Z)
```

LEXING CHALLENGES: FORTRAN

- ✓ Pre-1977 Hollerith Constants

```
15HCOMPILER DESIGN
```

- ✓ Fixed-Column Layouts

```
CALL DOIT(X,Y,Z)
```

```
CALL DOIT(X,Y,Z)
```

- ✓ Whitespace is *not* significant.

```
DO 10 I = 1,15
```

```
DO 10 I = 1.15
```

```
REAL X
```

```
REAL X = 3.5
```

```
INTEGER FUNCTION A(I)
```

PARSING CHALLENGES: C, C++, C#, JAVA, ...

- ✓ Often can be resolved by extra information from the lexer.

X * Y;

X (Y);

X < Y > Z;

X ? Y : Z;

X ? Y = Z;

(T)-(U);

PARSING CHALLENGES: C, C++, C#, JAVA, ...

- ✓ Often can be resolved by extra information from the lexer.

```
X * Y;
```

```
X (Y);
```

```
X < Y > Z;
```

```
X ? Y : Z;
```

```
X ? Y = Z;
```

```
(T)-(U);
```

- ✓ Not all challenges are this easy, though

```
int(x), y, *const z;    // variable declarations
```

```
int(x), y, new int;    // list of expressions
```

```
int(x), y, z = 0;      // variable declarations
```

IMPLEMENTATION

Once we have a DFA, how do we turn it into code?