Compilers INF-400

Burak Arslan ext-inf400@burakarslan.com

Galatasaray Üniversitesi

Lecture II 2023-10-12

Course website

burakarslan.com/inf400

Lexing

Lexing;

- **1.** is the first step towards the **definition** of a language
- 2. is a left-to-right scan operation with some lookahead
- **3.** defines the set of token classes and the set of character sequences that belong to these classes
- 4. needs a set of characters to operate on called the alphabet



- 1. $\Sigma = \text{English characters}$ Language = English sentences
- 2. Not every string of English characters is an English sentence
- $\blacktriangleright \Sigma = \mathsf{ASCII}$
 - $\mathsf{Language} = \mathsf{C} \ \mathsf{programs}$
- Note: ASCII character set is different from English character set



```
import io
func main() -> int32 {
    io.print("Hello world!\n");
    return 0;
}
```

Lexing

```
KW_IMPORT IDENTIFIER(io) OP_NEWLINE
KW_FUNC IDENTIFIER(main) OP_LPAREN OP_RPAREN OP_RETURNS
IDENTIFIER(int32) OP_LBRACE OP_NEWLINE
IDENTIFIER(io) OP_DOT IDENTIFIER(print) OP_LPAREN
L_STRING("Hello world!\n") OP_RPAREN OP_SCOLON OP_NEWLINE
IDENTIFIER(return) L_INTEGER(0, 10) OP_SCOLON OP_NEWLINE
OP_RBRACE OP_NEWLINE
```

* Indented lines are continuations of the previous ones. There is actually only 5 lines in the above output



Language definitions start with tokens types:

Туре	Examples
IDENTIFIERs	foo bar main
INTEGERs	123 -456 42 0xFF
FLOATs	12.34 +56.78
KEYWORDs	if while import
OPERATORs	, . >= ==



A token type specifies a set of acceptable strings



A token type specifies a set of acceptable strings

Туре	Examples
INTEGERs	123 -456 42 0×FF

... but INTEGER set is infinite (?)



So we turn to regular languages. More specifically;

 Regular Expressions: for defining the sets
 Deterministic Finite (State) Automata (DFA): for testing set membership



Туре	Regex
IF	
\pm INTEGER (dec)	
+INTEGER (hex)	
IDENTIFIER	
KEYWORD	



Туре	Regex
IF	if
\pm INTEGER (dec)	
+INTEGER (hex)	
IDENTIFIER	
KEYWORD	



Туре	Regex
IF	if
\pm INTEGER (dec)	[0-9]+
+INTEGER (hex)	
IDENTIFIER	
KEYWORD	



Туре	Regex
IF	if
\pm INTEGER (dec)	-?[0-9]+
+INTEGER (hex)	
IDENTIFIER	
KEYWORD	



Туре	Regex
IF	if
\pm INTEGER (dec)	(+ -)?[0-9]+
+INTEGER (hex)	
IDENTIFIER	
KEYWORD	



Туре	Regex
IF	if
\pm INTEGER (dec)	(+ -)?[0-9]+
+INTEGER (hex)	0x[0-9a-fA-F]+
IDENTIFIER	
KEYWORD	



Туре	Regex
IF	if
\pm INTEGER (dec)	(+ -)?[0-9]+
+INTEGER (hex)	0x[0-9a-fA-F]+
IDENTIFIER	[a-zA-Z][a-zA-Z0-9]*
KEYWORD	



Туре	Regex
IF	if
\pm INTEGER (dec)	(+ -)?[0-9]+
+INTEGER (hex)	0x[0-9a-fA-F]+
IDENTIFIER	[a-zA-Z][a-zA-Z0-9]*
KEYWORD	(if while import func)



Let's classify the following token: if

Туре	Regex
IF	if
\pm INTEGER (dec)	(+ -)?[0-9]+
+INTEGER (hex)	0x[0-9a-f-A-F]+
IDENTIFIER	[a-zA-Z][a-zA-Z0-9]*
KEYWORD	(if while import func)



Let's classify the following token: if \Rightarrow IF

Туре	Regex
IF	if
\pm INTEGER (dec)	(+ -)?[0-9]+
+INTEGER (hex)	0x[0-9a-f-A-F]+
IDENTIFIER	[a-zA-Z][a-zA-Z0-9]*
KEYWORD	(if while import func)



Let's classify the following token: 1234

Туре	Regex
IF	if
\pm INTEGER (dec)	(+ -)?[0-9]+
+INTEGER (hex)	0x[0-9a-f-A-F]+
IDENTIFIER	[a-zA-Z][a-zA-Z0-9]*
KEYWORD	(if while import func)



Let's classify the following token: $1234 \Rightarrow \texttt{INTEGER}$

Туре	Regex
IF	if
\pm INTEGER (dec)	(+ -)?[0-9]+
+INTEGER (hex)	0x[0-9a-f-A-F]+
IDENTIFIER	[a-zA-Z][a-zA-Z0-9]*
KEYWORD	(if while import func)



Let's classify the following token: 1234A

Туре	Regex
IF	if
\pm INTEGER (dec)	(+ -)?[0-9]+
+INTEGER (hex)	0x[0-9a-f-A-F]+
IDENTIFIER	[a-zA-Z][a-zA-Z0-9]*
KEYWORD	(if while import func)



Let's classify the following token: $1234A \Rightarrow REJECT$

Туре	Regex
IF	if
\pm INTEGER (dec)	(+ -)?[0-9]+
+INTEGER (hex)	0x[0-9a-f-A-F]+
IDENTIFIER	[a-zA-Z][a-zA-Z0-9]*
KEYWORD	(if while import func)



Let's classify the following token: **import**

Туре	Regex
IF	if
\pm INTEGER (dec)	(+ -)?[0-9]+
+INTEGER (hex)	0x[0-9a-f-A-F]+
IDENTIFIER	[a-zA-Z][a-zA-Z0-9]*
KEYWORD	(if while import func)



Let's classify the following token: **import** \Rightarrow **KEYWORD**

Туре	Regex
IF	if
\pm INTEGER (dec)	(+ -)?[0-9]+
+INTEGER (hex)	0x[0-9a-f-A-F]+
IDENTIFIER	[a-zA-Z][a-zA-Z0-9]*
KEYWORD	(if while import func)



Let's classify the following token: **importer**

Туре	Regex
IF	if
\pm INTEGER (dec)	(+ -)?[0-9]+
+INTEGER (hex)	0x[0-9a-f-A-F]+
IDENTIFIER	[a-zA-Z][a-zA-Z0-9]*
KEYWORD	(if while import func)



Let's classify the following token: **importer** \Rightarrow **IDENTIFIER**

Туре	Regex
IF	if
\pm INTEGER (dec)	(+ -)?[0-9]+
+INTEGER (hex)	0x[0-9a-f-A-F]+
IDENTIFIER	[a-zA-Z][a-zA-Z0-9]*
KEYWORD	(if while import func)



Let's classify the following token: **importer** \Rightarrow **IDENTIFIER** ???

Туре	Regex
IF	if
\pm INTEGER (dec)	(+ -)?[0-9]+
+INTEGER (hex)	0x[0-9a-f-A-F]+
IDENTIFIER	[a-zA-Z][a-zA-Z0-9]*
KEYWORD	(if while import func)



The lookahead problem

- **1.** import \Rightarrow KW_IMPORT
- **2.** importer \Rightarrow IDEN(importer)



Ambiguities in grammars are practically impossible to avoid. C++ Examples:

- ► A < B > c;
 - ▶ if A and B are types, this is a variable definition;
 - if A and B are variables, this is a (pointless) comparison;
- cin >> var;
 - operator>>
- vector<unique_ptr<string>> ;
 - Right angle bracket



In the face of ambiguities;

Longest match is preferred

Order of rules matter



DFA that decides whether an input matches regexp: [0-9]+





DFA that decides whether an input matches regexp: if

start
$$\longrightarrow$$
 $s_0 \xrightarrow{i} s_1 \xrightarrow{f} s_2$



DFA ... regexp: ...?





The combined DFA becomes the lexical analyzer



Lexing

Deterministic Finite State Automata

The **lexer generator**'s job is to combine all regexp to a coherent **DFA**





The lexgen is used to assign **callbacks** to accepting states





Regular Expressions

Remember: Regular expressions specify sets of strings.

 $\forall A, B \in \textcircled{S}$, sets of strings over alphabet Σ ;

▶ Neutral: {""}
$$\implies \varepsilon^1 \neq \varnothing$$

- Union: $A \cup B \Longrightarrow (A|B)$
- ▶ Concatenation: ${s_1s_2 | s_1 \in A \land s_2 \in B} \Longrightarrow AB$
- ► Range: $\{$ "a", "b", ..., "z" $\} \implies [a z]$
- ▶ Range Excl.: (S) $\{ "a", "b", ..., "z" \} \implies [^a-z]$

¹singleton with an empty string

Regular Expressions

Repetitions: Let $A \in \mathfrak{S}$, sets of strings over alphabet Σ , $A^n = \underbrace{AA \dots A}_n$

• Optional:
$$A + \varepsilon \Longrightarrow A$$
?

• Zero or more:
$$\bigcup_{i\geq 0} A^i \Longrightarrow A^*$$

• One or more:
$$\bigcup_{i>0} A^i \Longrightarrow A^+$$

► Explicit:

Regular expressions are implemented using Finite State Automata

- ► A finite automaton consists of
- \blacktriangleright An input alphabet Σ
- A set of states S
- A start state n
- A set of accepting states $F \subseteq S$
- \blacktriangleright A set of transitions state \rightarrow input state

Finite State Automata

Two types:

- **DFA**: Deterministic Finite Automata
- ▶ NFA: Nondeterministic Finite Automata

Finite State Automata

Two types:

DFA:

- No more than one move per input
- \triangleright ε moves are forbidden
- ► NFA:
 - Zero or more moves per input
 - \triangleright ε moves are allowed

Finite State Automata

Regular expressions have direct NFA representations. Eg. (A|B):



Finite State Automata

Regular expressions have direct NFA representations. Eg. AB:



Finite State Automata

NFA and DFA are equivalent and recognize both regular languages.

DFAs are faster to execute

Finite State Automata

Conversion algorithm: **Simulation** / **Tracing** (recording execution)

- Start state of the DFA = States reachable from the start state of the NFA through ε input
- Let n, n', n'', \ldots states from NFA,
- ▶ Let d, d', d'', \ldots states from DFA,
- ▶ Add a new state $d \xrightarrow[a]{} d'$ if and only if n' is reachable from n, including ε input

Finite State Automata



Implementation

We use tables / grids / 2D arrays:



Implementation

All in all, lexer generators' job boil down to:

- Unify all regular expressions into a single NFA
- Perform NFA \Rightarrow DFA conversion
- Create DFA grid
- Execute DFA