

Lexical Analysis - An Introduction



The purpose of the front end is to deal with the input language

Errors

- Perform a membership test: code ∈ source language?
- Is the program well-formed (semantically)?
- Build an IR version of the code for the rest of the compiler

The front end is not monolithic





- Its part of speech (or syntactic category) is called its token type
- Scanner discards white space & (often) comments





- Checks stream of classified words (*parts of speech*) for grammatical correctness
- Determines if code is syntactically well-formed
- Guides checking at deeper levels than syntax
- Builds an IR representation of the code

We'll come back to parsing in a couple of lectures



Why study lexical analysis?

- We want to avoid writing scanners by hand
- We want to harness the theory from classes like CISC 303 Goals:
 - \rightarrow To simplify specification & implementation of scanners
 - \rightarrow To understand the underlying techniques and technologies



Where is Lexical Analysis used?

- For traditional languages but where else...
- Web page "compilation"
 - Lexical Analysis of HTML, XML, etc.
- Natural Language Processing
- Game Scripting Engines
- OS Shell Command Line
- Find
- Prototyping high-level languages
 - JavaScript, Perl, Python





 Finite Automaton (FA) – recognizers that can scan a stream of symbols to find lexemes (words)



Recognizer for *Register*

- An FA is a five-tuple $(S, \Sigma, \partial, S_o, S_F)$ where
 - S is the set of states
 - $\boldsymbol{\Sigma}$ is the alphabet
 - δ a set of transition functions where each takes a state and a character and returns another state
 - s_o is the start state
 - S_F is the set of final states

UNIVERSITY OF ELAWARE

Regular expressions (REs) describe regular languages

Regular Expression (over alphabet Σ)

- ε is a RE denoting the set {ε}
- If \underline{a} is in Σ , then \underline{a} is a RE denoting $\{\underline{a}\}$
- If x and y are REs denoting L(x) and L(y) then
 - \rightarrow Alternation : x | y is an RE denoting $L(x) \cup L(y)$
 - \rightarrow Concatenation: xy is an RE denoting L(x)L(y)
 - \rightarrow Closure: x^* is an RE denoting $L(x)^*$

<u>Precedence</u> is *closure*, then *concatenation*, then *alternation*

Examples of Regular Expressions

Identifiers:

Numbers:

Numbers can get much more complicated!



(the point)



Regular expressions can be used to specify the words to be translated to parts of speech by a lexical analyzer

Using results from automata theory and theory of algorithms, we can automatically build recognizers from regular expressions

Some of you may have seen this construction for string pattern matching

⇒ We study REs and associated theory to automate scanner construction !



Consider the problem of recognizing register names

 $Register \rightarrow r (\underline{0|1|2|} \dots | \underline{9}) (\underline{0|1|2|} \dots | \underline{9})^{*}$

- Allows registers of arbitrary number
- Requires at least one digit



Recognizer for *Register*

Transitions on other inputs go to an error state, s_e

Example

(continued)



DFA operation

- Start in state S_0 & take transitions on each input character
- DFA accepts a word <u>x iff x</u> leaves it in a final state (S_2)



Recognizer for Register

So,

- <u>r17</u> takes it through s_0 , s_1 , s_2 and accepts
- <u>r</u> takes it through s₀, s₁ and fails
- <u>a</u> takes it straight to s_e

Example

(continued)



To be useful, recognizer must turn into code



δ	r	0,1,2,3,4, 5,6,7,8,9	All others	
s ₀	\$ ₁	S _e	S _e	
S 1	S _e	s ₂	S _e S _e	
s ₂	S _e	s ₂		
S _e	S _e	S _e	S _e	

Skeleton recognizer

Table encoding RE

What if we need a tighter specification?



r (0|1|2| ... | 9) (0|1|2| ... | 9)* allows arbitrary numbers

- Accepts <u>r00000</u>
- Accepts <u>r99999</u>
- What if we want to limit it to <u>rO</u> through <u>r31</u>?

Write a tighter regular expression

- → Register → <u>r</u> ((0|1|2) ([0...9] | ε) | (4|5|6|7|8|9) | (3(0|1| ε)))
- $\rightarrow Register \rightarrow \underline{r0|r1|r2|} \dots |\underline{r31|r00|r01|r02|} \dots |\underline{r09}|$

Produces a more complex DFA

- Has more states
- Same cost per transition
- Same basic implementation

(continued)



The DFA for

Register → <u>r</u> ((0|1|2) ([0...9] |ε) | (4|5|6|7|8|9) | (3(0|1|ε)))



- Accepts a more constrained set of registers
- Same set of actions, more states

•	Tighte	ed)						
	δ	r	0,1	2	3	4-9	All others	
	s 0	S 1	S _e	S _e	S _e	S _e	S _e	
	S 1	s _e	S 2	S 2	S 5	S 4	S _e	
	s 2	S _e	S 3	S 3	S 3	S 3	S _e	
	S 3	S _e	S _e	s _e	s _e	s _e	S _e	← Runs in the
	S 4	S _e	S _e	S _e	S _e	S _e	S _e	same
	S 5	S _e	s ₆	S _e	s _e	s _e	S _e	recognizer
	S 6	S _e	S _e	S _e	S _e	S _e	S _e	
	S _e	S _e	S _e	S _e	S _e	S _e	S _e	

Table encoding RE for the tighter register specification