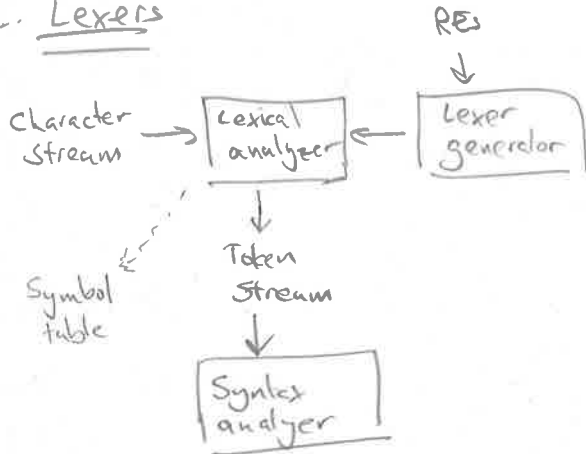


# Lexical analyzer (Pg. 1)

## Lecture topics:

- I. Lexers
- II. Regular expressions
- III. Finite automata

## I. Lexers



### Task of lexer:

- find tokens
- find token attributes
- strip whitespace and comments

### Why separate lexer from parser?

- more modular design
- may use more efficient algorithms
- may hide input format

### Some tools are scannerless

- do not separate lexer from parser
- attr. can be scannerless, but we won't use it that way

## II. Regular expressions

(RE example:  $(ab)^*abb$   
accepted strings:  $abb, cabb, abbabb, \dots$ )

	Regular Expressions	Context free grammars
Basic constructs	sequence (e.g. $ab$ ) alternation (e.g. $a b$ ) empty string ( $\epsilon$ ) grouping ( $(...)$ ) closure ( $...^*$ ) a.k.a. Kleene star, repetition	sequence alternation empty string ( $\epsilon$ ) recursion
Extended constructs	character classes e.g. $[_\wedge\wedge]$ $[a-zA-Z0-9]$ helper rules, but no recursion $+, ?, -, !,  $	grouping closure $+, ?$

Can emulate every RE with a CFG, e.g.

RE:	CFG:
$(ab)^*abb$	$X \rightarrow Yabb$ $Y \rightarrow YZ   \epsilon$ $Z \rightarrow a b$

$\Rightarrow$  CFGs are at least as powerful as REs.

Cannot evaluate every CFG with an RE, e.g.:

CFG:  
 $X \rightarrow \{x\}^n | y$   
 accepted strings:  $y, \{y\}, \{y\}^2, \{y\}^3, \dots, \{y\}^n$

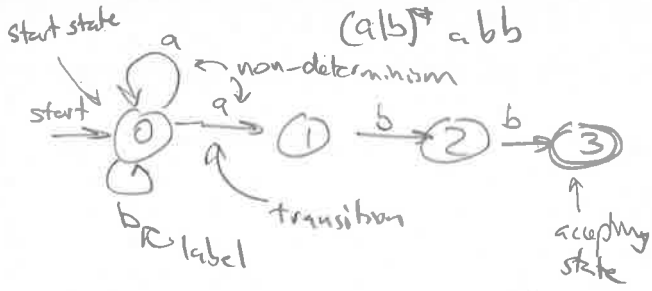
RE: cannot "count" parentheses

$\Rightarrow$  CFGs are more powerful than REs.

# Lexical analyzer (pg 2)

## III Finite automata

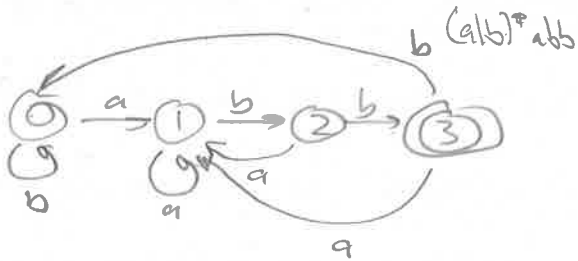
NFA = non-deterministic finite automaton



Transition table:

state	a	b	$\epsilon$
0	{0, 1}	{0}	$\emptyset$
1	$\emptyset$	{2}	$\emptyset$
2	$\emptyset$	{3}	$\emptyset$
3	$\emptyset$	$\emptyset$	$\emptyset$

DFA = deterministic finite automaton



Transition table

State	a	b
0	1	0
1	1	2
2	1	3
3	1	0

	DFA	NFA
labels from same state	unique complete	may be same may be incomplete
$\epsilon$ -labels	no	yes
transition table	states	state sets
expressiveness	same as REs	same as REs
tradeoffs	easy to simulate	easy to come up with