

# Computing if a token can follow

**first**( $B_1 \dots B_p$ ) =  $\{a \in \Sigma \mid B_1 \dots B_p \Rightarrow \dots \Rightarrow aw\}$

**follow**( $X$ ) =  $\{a \in \Sigma \mid S \Rightarrow \dots \Rightarrow \dots Xa \dots\}$

There exists a derivation from the start symbol that produces a sequence of terminals and nonterminals of the form  $\dots Xa \dots$   
(the token  $a$  follows the non-terminal  $X$ )

# Rule for Computing Follow

Given  $X ::= YZ$  (for reachable  $X$ )

then  $\text{first}(Z) \subseteq \text{follow}(Y)$

and  $\text{follow}(X) \subseteq \text{follow}(Z)$

$S \Rightarrow Xa \Rightarrow YZa \Rightarrow Yba$

$S ::= Xa$

$X ::= YZ$

$Y ::= b$

$Z ::= c$

**Now take care of nullable ones as well:**

For each rule  $X ::= Y_1 \dots Y_p \dots Y_q \dots Y_r$

$\text{follow}(Y_p)$  should contain:

- $\text{first}(Y_{p+1}Y_{p+2}\dots Y_r)$
- also  $\text{follow}(X)$  if  $\text{nullable}(Y_{p+1}Y_{p+2}Y_r)$

# Compute nullable, first, follow

stmtList ::=  $\epsilon$  | stmt stmtList

stmt ::= assign | block

assign ::= **ID = ID ;**

block ::= **beginof ID stmtList ID ends**

Compute **follow** (for that we need **nullable,first**)

# Conclusion of the Solution

The grammar is not LL(1) because we have

- **nullable**(*stmtList*)
- **first**(*stmt*)  $\cap$  **follow**(*stmtList*) = {**ID**}
- If a recursive-descent parser sees **ID**, it does not know if it should
  - finish parsing *stmtList* or
  - parse another *stmt*

# LL(1) Grammar - good for building recursive descent parsers

- Grammar is LL(1) if for each nonterminal  $X$ 
  - first sets of different alternatives of  $X$  are disjoint
  - if nullable( $X$ ), first( $X$ ) must be disjoint from follow( $X$ )
- For each LL(1) grammar we can build recursive-descent parser
- Each LL(1) grammar is unambiguous
- If a grammar is not LL(1), we can sometimes transform it into equivalent LL(1) grammar

# Table for LL(1) Parser: Example

$S ::= B \text{ EOF}$   
(1)

$B ::= \varepsilon \mid B (B)$   
(1)      (2)

nullable: B

$\text{first}(S) = \{ ( \}$

$\text{follow}(S) = \{ \}$

$\text{first}(B) = \{ ( \}$

$\text{follow}(B) = \{ ), (, \text{EOF} \}$

empty entry:  
when parsing S,  
if we see ),  
report error

**Parsing table:**

	EOF	(	)
S	{1}	{1}	{ }
B	{1}	{1,2}	{1}

**parse conflict - choice ambiguity:  
grammar not LL(1)**

1 is in entry because ( is in follow(B)

2 is in entry because ( is in first(B(B))

# Table for LL(1) Parsing

Tells which alternative to take, given current token:

choice : Nonterminal x Token  $\rightarrow$  Set[Int]

$$\begin{array}{l} A ::= (1) B_1 \dots B_p \\ \quad | (2) C_1 \dots C_q \\ \quad | (3) D_1 \dots D_r \end{array}$$

if  $t \in \text{first}(C_1 \dots C_q)$  add 2  
to choice(A,t)  
if  $t \in \text{follow}(A)$  add K to choice(A,t)  
where K is nullable alternative

For example, when parsing A and seeing token t

choice(A,t) = {2} means: parse alternative 2 ( $C_1 \dots C_q$ )

choice(A,t) = {1} means: parse alternative 3 ( $D_1 \dots D_r$ )

choice(A,t) = {} means: report syntax error

choice(A,t) = {2,3} : not LL(1) grammar

# Transform Grammar for LL(1)

$S ::= B \text{ EOF}$

$B ::= \varepsilon \mid B (B)$   
(1)      (2)

Transform the grammar so that parsing table has no conflicts.

$S ::= B \text{ EOF}$

$B ::= \varepsilon \mid (B) B$   
(1)      (2)

Left recursion is bad for LL(1)

Old parsing table:

	EOF	(	)
S	{1}	{1}	{}
B	{1}	{1,2}	{1}

**conflict - choice ambiguity:  
grammar not LL(1)**

- 1 is in entry because ( is in follow(B)
- 2 is in entry because ( is in first(B(B))

	EOF	(	)
S			
B			

choice(A,t)



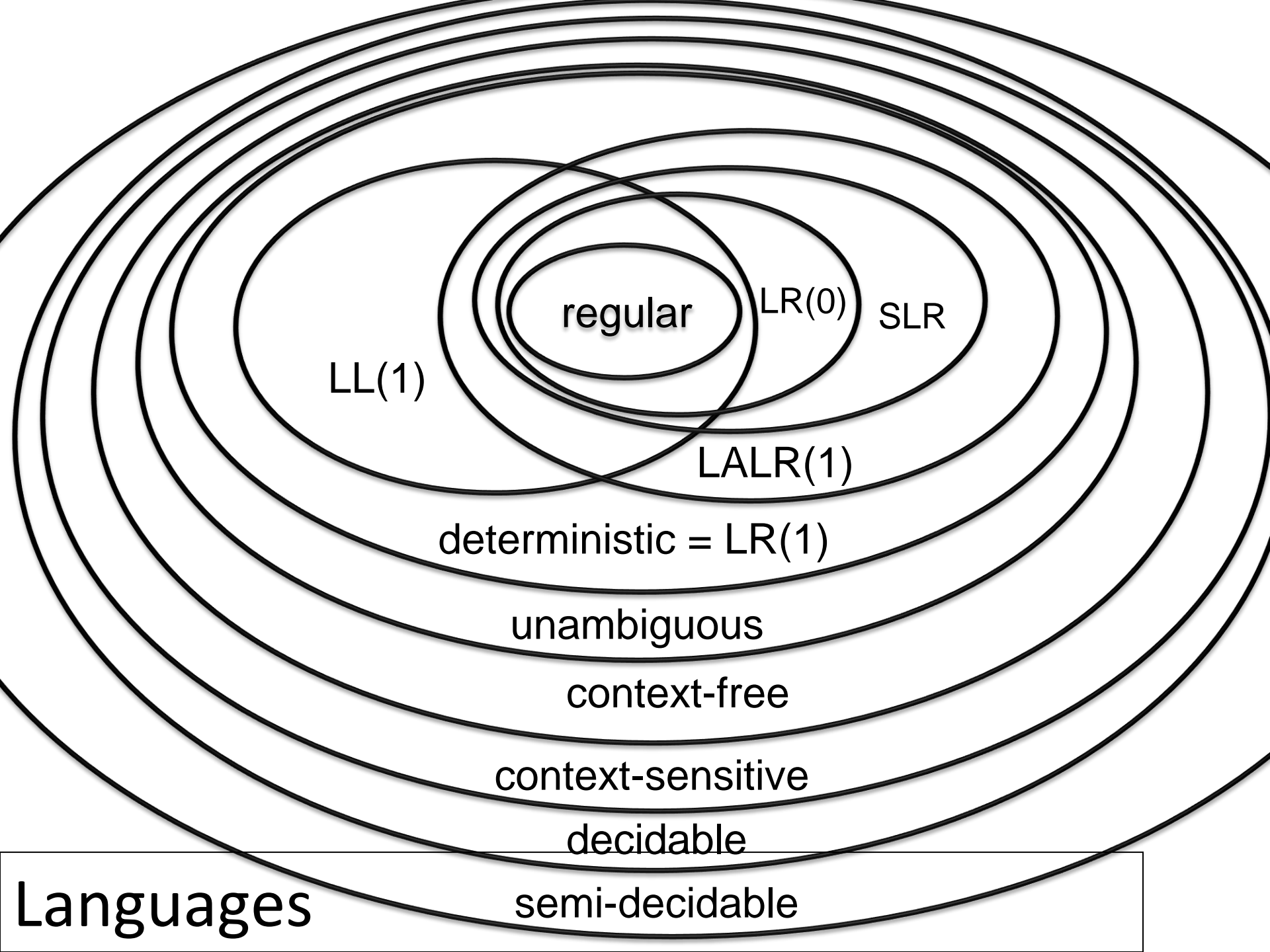
# Parse Table is Code for Generic Parser

```
var stack : Stack[GrammarSymbol] // terminal or non-terminal
stack.push(EOF);
stack.push(StartNonterminal);
var lex = new Lexer(inputFile)
while (true) {
  X = stack.pop
  t = lex.curent
  if (isTerminal(X))
    if (t==X) if (X==EOF) return success
    else lex.next // eat token t
  else parseError("Expected " + X)
else { // non-terminal
  cs = choice(X)(t) // look up parsing table
  cs match { // result is a set
  case {i} => { // exactly one choice
    rhs = p(X,i) // choose correct right-hand side
    stack.push(reverse(rhs)) }
  case {} => parseError("Parser expected an element of " + unionOfAll(choice(X)))
  case _ => crash("parse table with conflicts - grammar was not LL(1)")
  }
}
```

# What if we cannot transform the grammar into LL(1)?

1) Redesign your language

2) Use a more powerful parsing technique



regular

LR(0)

SLR

LL(1)

LALR(1)

deterministic = LR(1)

unambiguous

context-free

context-sensitive

decidable

semi-decidable

Languages

# Remark: Grammars and Languages

- Language  $S$  is a set of words
- For each language  $S$ , there can be multiple possible grammars  $G$  such that  $S=L(G)$
- Language  $S$  is
  - Non-ambiguous if there exists a non-ambiguous grammar for it
  - LL(1) if there is an LL(1) grammar for it
- Even if a language has ambiguous grammar, it can still be non-ambiguous if it also has a non-ambiguous grammar

# Parsing General Grammars: Why

- Can be difficult or impossible to make grammar unambiguous
- Some inputs are more complex than simple programming languages
  - mathematical formulas:  
 $x = y \wedge z$       ?       $(x=y) \wedge z$        $x = (y \wedge z)$
  - future programming languages
  - natural language:

*I saw the man with the telescope.*

# Ambiguity

1)



2)



*I saw the man with the telescope.*

# CYK Parsing Algorithm

C:

[John Cocke](#) and Jacob T. Schwartz (1970). Programming languages and their compilers: Preliminary notes. Technical report, [Courant Institute of Mathematical Sciences](#), [New York University](#).

Y:

Daniel H. **Younger** (1967). Recognition and parsing of context-free languages in time  $n^3$ . *Information and Control* 10(2): 189–208.

K:

[T. Kasami](#) (1965). An efficient recognition and syntax-analysis algorithm for context-free languages. Scientific report AFCRL-65-758, Air Force Cambridge Research Lab, [Bedford, MA](#).

# Two Steps in the Algorithm

1) Transform grammar to normal form  
called Chomsky Normal Form

(Noam Chomsky, mathematical linguist)

2) Parse input using transformed grammar  
dynamic programming algorithm

“a method for solving complex problems by breaking them down into simpler steps.

It is applicable to problems exhibiting the properties of overlapping subproblems” (>WP)



# Chomsky Normal Form

- Essentially, only binary rules
- Concretely, these kinds of rules:

$X ::= Y Z$	binary rule $X, Y, Z$ - non-terminals
$X ::= a$	non-terminal as a name for token
$S ::= \varepsilon$	only for top-level symbol $S$

# Balanced Parentheses Grammar

Original grammar G

$$S \rightarrow \varepsilon \mid (S) \mid SS$$

Modified grammar in Chomsky Normal Form:

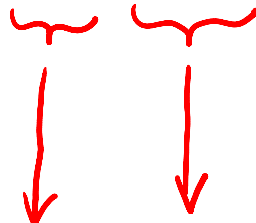
$$S \rightarrow \varepsilon \mid S' \quad \leftarrow \text{if } "" \in L(G)$$

$$\begin{array}{l}
 S' \rightarrow N_{(} N_{S)} \mid N_{(} N_{)} \mid S' S' \\
 N_{S)} \rightarrow S' N_{)} \\
 N_{(} \rightarrow ( \\
 N_{)} \rightarrow )
 \end{array}
 \left. \begin{array}{l}
 \text{Rules} \\
 \text{Rules}
 \end{array} \right\}
 \begin{array}{l}
 N \rightarrow N_1 N_2 \\
 \text{nonterminals} \\
 N \rightarrow t \\
 \begin{array}{l}
 \uparrow \text{nonterminal} \quad \uparrow \text{terminal}
 \end{array}
 \end{array}$$

- Terminals: ( )    Nonterminals: S S' N<sub>S)</sub> N<sub>)</sub> N<sub>(</sub>  
nonterminal with funny name

# Idea How We Obtained the Grammar

$$S \rightarrow ( S )$$



$$S' \rightarrow N_{(} N_{)} \mid N_{(} N_{)}$$

because  $S$  can be empty  
but  $S'$  cannot

$$N_{(} \rightarrow ($$

$$N_{)} \rightarrow S' N_{)}$$

$$N_{)} \rightarrow )$$

Chomsky Normal Form transformation  
can be done fully mechanically

# Transforming Grammars into Chomsky Normal Form

## Steps:

1. remove unproductive symbols
2. remove unreachable symbols
3. remove epsilons (no non-start nullable symbols)
4. remove single non-terminal productions  $X ::= Y$
5. transform productions w/ more than 3 on RHS
6. make terminals occur alone on right-hand side

## 4) Eliminating single productions

- Single production is of the form

$X ::= Y$

where  $X, Y$  are non-terminals

$\text{program} ::= \text{stmtSeq}$

$\text{stmtSeq} ::= \text{stmt}$

$\quad \quad \quad | \text{stmt} ; \text{stmtSeq}$

$\text{stmt} ::= \text{assignment} | \text{whileStmt}$

$\text{assignment} ::= \text{expr} = \text{expr}$

$\text{whileStmt} ::= \text{while} (\text{expr}) \text{stmt}$

## 4) Eliminate single productions - Result

- Generalizes removal of epsilon transitions from non-deterministic automata

program ::= expr = expr | while (expr) stmt  
          | stmt ; stmtSeq

stmtSeq ::= expr = expr | while (expr) stmt  
          | stmt ; stmtSeq

stmt ::= expr = expr | while (expr) stmt

assignment ::= expr = expr

whileStmt ::= while (expr) stmt

## 4) “Single Production Terminator”

- If there is single production  
 $X ::= Y$  put an edge  $(X, Y)$  into graph
- If there is a path from  $X$  to  $Z$  in the graph, and there is rule  $Z ::= s_1 s_2 \dots s_n$  then add rule  
 $X ::= s_1 s_2 \dots s_n$

At the end, remove all single productions.

$\text{program} ::= \text{expr} = \text{expr} \mid \text{while}(\text{expr}) \text{stmt}$   
 $\quad \mid \text{stmt} ; \text{stmtSeq}$

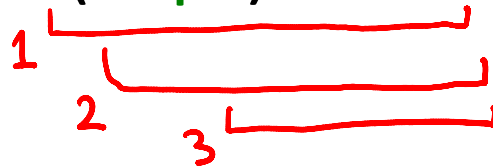
$\text{stmtSeq} ::= \text{expr} = \text{expr} \mid \text{while}(\text{expr}) \text{stmt}$   
 $\quad \mid \text{stmt} ; \text{stmtSeq}$

$\text{stmt} ::= \text{expr} = \text{expr} \mid \text{while}(\text{expr}) \text{stmt}$

## 5) No more than 2 symbols on RHS

$\text{stmt} ::= \text{while } (\text{expr}) \text{ stmt}$

becomes



$\text{stmt} ::= \text{while } \text{stmt}_1$

$\text{stmt}_1 ::= ( \text{stmt}_2$

$\text{stmt}_2 ::= \text{expr } \text{stmt}_3$

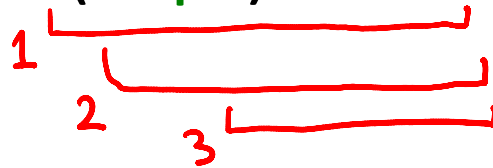
$\text{stmt}_3 ::= ) \text{stmt}$



## 6) A non-terminal for each terminal

$\text{stmt} ::= \text{while } (\text{expr}) \text{ stmt}$

becomes



$\text{stmt} ::= N_{\text{while}} \text{stmt}_1$

$\text{stmt}_1 ::= N_{(} \text{stmt}_2$

$\text{stmt}_2 ::= \text{expr} \text{stmt}_3$

$\text{stmt}_3 ::= N_{)} \text{stmt}$

$N_{\text{while}} ::= \text{while}$

$N_{(} ::= ($

$N_{)} ::= )$

# Parsing using CYK Algorithm

- Transform grammar into Chomsky Form:
  1. remove unproductive symbols
  2. remove unreachable symbols
  3. remove epsilons (no non-start nullable symbols)
  4. remove single non-terminal productions  $X ::= Y$
  5. transform productions of arity more than two
  6. make terminals occur alone on right-hand sideHave only rules  $X ::= Y Z$ ,  $X ::= t$ , and possibly  $S ::= \epsilon$
- Apply CYK dynamic programming algorithm

# Dynamic Programming to Parse Input

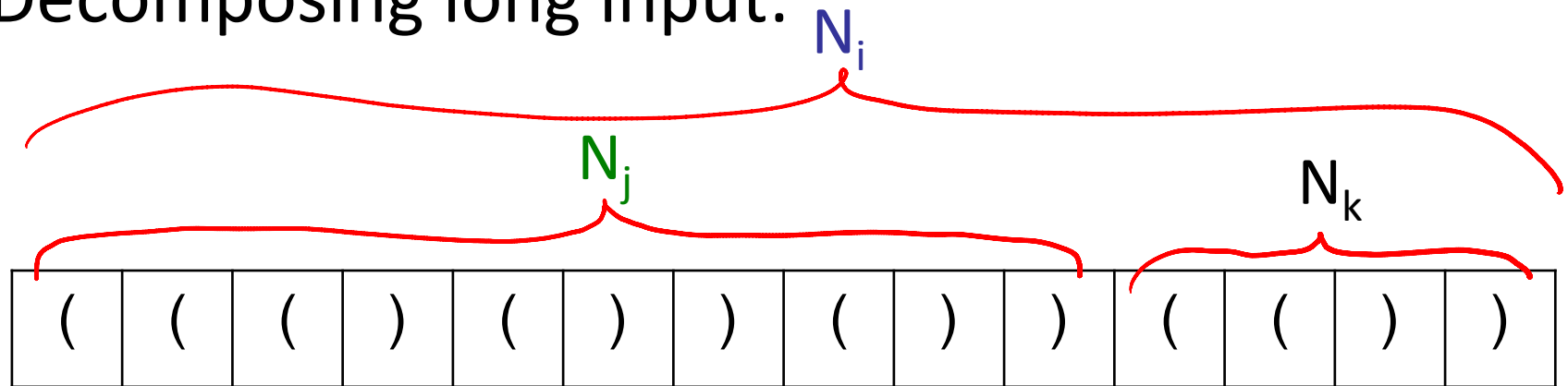
Assume Chomsky Normal Form, 3 types of rules:

$S \rightarrow \varepsilon \mid S'$  (only for the start non-terminal)

$N_j \rightarrow t$  (names for terminals)

$N_i \rightarrow N_j N_k$  (just **2** non-terminals on RHS)

Decomposing long input:



find all ways to parse substrings of length 1,2,3,...

# Parsing an Input

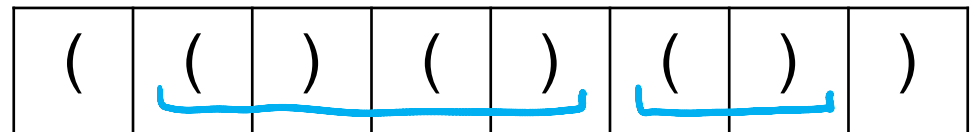
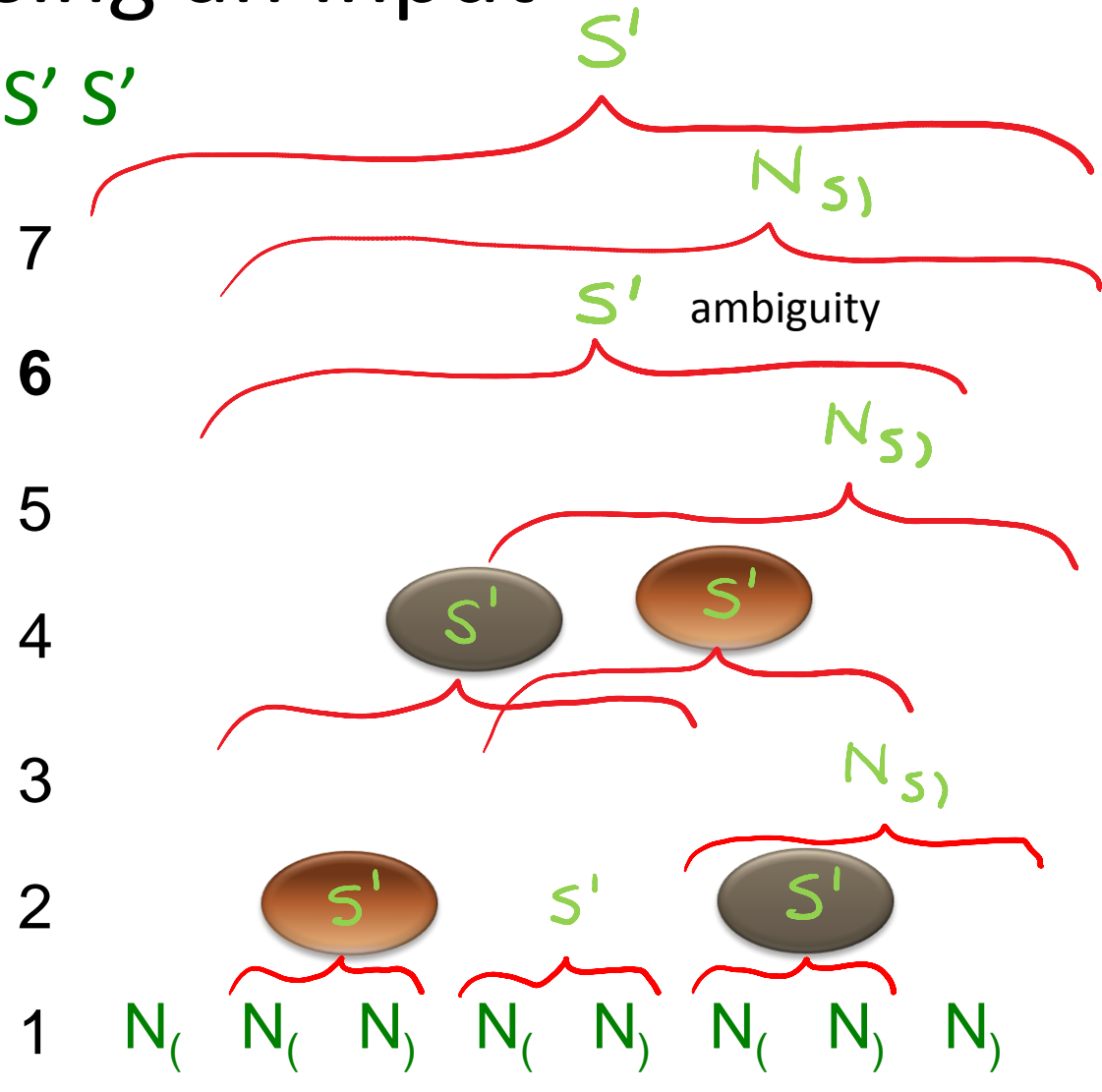
$$S' \rightarrow N_{(} N_{)} \mid N_{(} N_{)} \mid S' S'$$

$$N_{)} \rightarrow S' N_{(}$$

$$N_{(} \rightarrow ($$

$$N_{)} \rightarrow )$$

substring  
length



# Algorithm Idea

$$S' \rightarrow S' S'$$

$w_{pq}$  – substring from  $p$  to  $q$

$d_{pq}$  – all non-terminals that could expand to  $w_{pq}$

Initially  $d_{pp}$  has  $N_{w(p,p)}$

key step of the algorithm:

if  $X \rightarrow YZ$  is a rule,

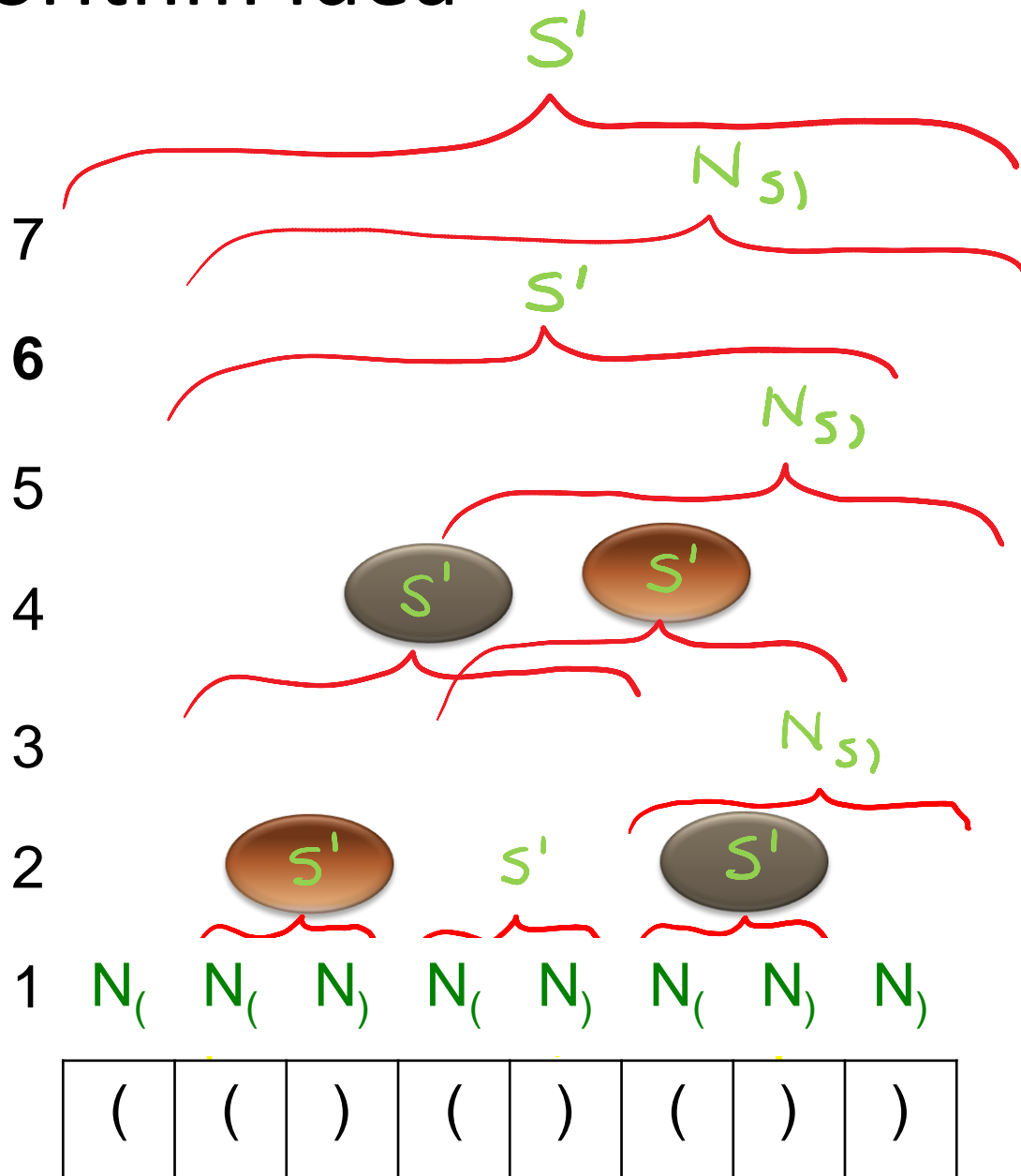
$Y$  is in  $d_{pr}$ , and

$Z$  is in  $d_{(r+1)q}$

then put  $X$  into  $d_{pq}$

( $p \leq r < q$ ),

in increasing value of  $(q-p)$



# Algorithm

INPUT: grammar  $G$  in Chomsky normal form  
word  $w$  to parse using  $G$

OUTPUT: true iff ( $w$  in  $L(G)$ )

$N = |w|$

var  $d$  : Array[ $N$ ][ $N$ ]

for  $p = 1$  to  $N$  {

$d(p)(p) = \{X \mid G \text{ contains } X \rightarrow w(p)\}$

for  $q$  in  $\{p + 1 .. N\}$   $d(p)(q) = \{\}$  }

for  $k = 2$  to  $N$  // substring length

for  $p = 0$  to  $N - k$  // initial position

for  $j = 1$  to  $k - 1$  // length of first half

val  $r = p + j - 1$ ; val  $q = p + k - 1$ ;

for  $(X ::= Y Z)$  in  $G$

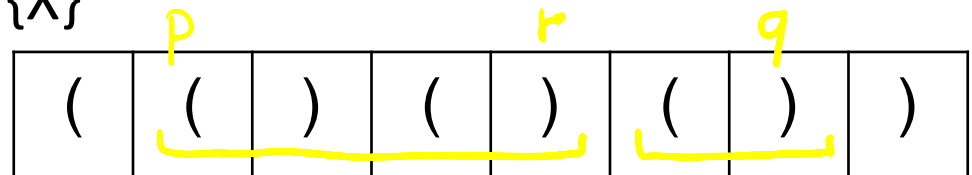
if  $Y$  in  $d(p)(r)$  and  $Z$  in  $d(r + 1)(q)$

$d(p)(q) = d(p)(q) \cup \{X\}$

return  $S$  in  $d(0)(N - 1)$

What is the running time  
as a function of grammar  
size and the size of input?

$O(\quad)$



# Parsing another Input

$S' \rightarrow N_{(} N_{S)} \mid N_{(} N_{)} \mid S' S'$

$N_{S)} \rightarrow S' N_{)}$

$N_{(} \rightarrow ($

$N_{)} \rightarrow )$

substring  
length

7

6

5

4

3

2

1

$N_{(}$   $N_{)}$   $N_{(}$   $N_{)}$   $N_{(}$   $N_{)}$   $N_{(}$   $N_{)}$

(	)	(	)	(	)	(	)
---	---	---	---	---	---	---	---

# Number of Parse Trees

- Let  $w$  denote word  $()()()$ 
  - it has two parse trees
- Give a lower bound on number of parse trees of the word  $w^n$  ( $n$  is positive integer)  
 $w^5$  is the word  
 $()()() ()()() ()()() ()()() ()()()$
- CYK represents all parse trees compactly
  - can re-run algorithm to extract first parse tree, or enumerate parse trees one by one



# Earley's Algorithm

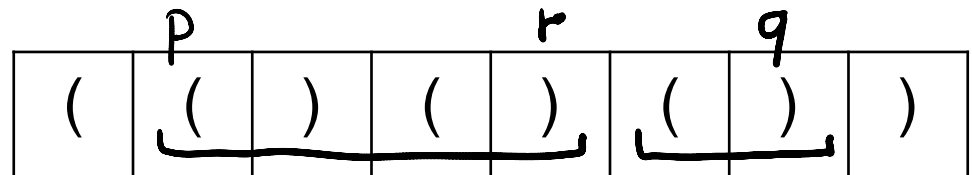
also parser arbitrary grammars

J. Earley, "An efficient context-free parsing algorithm", *Communications of the Association for Computing Machinery*, **13**:2:94-102, 1970.

# CYK vs Earley's Parser Comparison

$Z ::= X Y$                        $Z$  parses  $w_{pq}$

- CYK: if  $d_{pr}$  parses  $X$  and  $d_{(r+1)q}$  parses  $Y$ , then in  $d_{pq}$  stores symbol  $Z$
- Earley's parser:  
in set  $S_q$  stores *item* ( $Z ::= XY. , p$ )
- Move forward, similar to top-down parsers
- Use **dotted rules** to avoid binary rules



# Example: expressions

$D ::= e \text{ EOF}$

$e ::= \text{ID} \mid e - e \mid e == e$

## Rules with a dot inside

$D ::= . e \text{ EOF} \mid e . \text{ EOF} \mid e \text{ EOF} .$

$e ::= . \text{ID} \mid \text{ID} .$

$\mid . e - e \mid e . - e \mid e - . e \mid e - e .$

$\mid . e == e \mid e . == e \mid e == . e \mid e == e .$

		ID	-	ID	==	ID	EOF
	$\epsilon$	ID	ID-	ID-ID	ID-ID==	ID-ID==ID	
ID		$\epsilon$	-	-ID	-ID==	-ID==ID	
-			$\epsilon$	ID	ID==	ID==ID	
ID				$\epsilon$	==	==ID	
==					$\epsilon$	ID	
ID	$S ::= . e EOF \mid e . EOF \mid e EOF .$ $e ::= . ID \mid ID .$					$\epsilon$	
EOF	$\mid . e - e \mid e . - e \mid e - . e \mid e - e .$ $\mid . e == e \mid e . == e \mid e == . e \mid e == e .$						$\epsilon$

		ID <i>s<sub>1</sub></i>	- <i>s<sub>2</sub></i>	ID <i>s<sub>3</sub></i>	==	ID	EOF
	$\epsilon$ .e EOF .ID .e-e .e=e	ID ID. e, EOF e, -e e, =e	ID- e-.e	ID-ID e-e, e, EOF e, -e e, =e	ID-ID== e=.e	ID-ID==ID e=e. e-e.	e, EOF
ID		$\epsilon$	-	-ID	-ID==	-ID==ID	
-			$\epsilon$ .ID .e-e .e=e	ID ID. e, -e e, =e	ID== e=.e	ID==ID e=e.	
ID				$\epsilon$	==	==ID	
==					$\epsilon$ .ID .e-e .e=e	ID ID. e, -e e, =e	
ID	S ::= . e EOF   e . EOF   e EOF . e ::= . ID   ID .						$\epsilon$
EOF	. e - e   e . - e   e - . e   e - e .   . e == e   e . == e   e == . e   e == e .						$\epsilon$