

CS 4120

Introduction to Compilers

Andrew Myers
Cornell University

Lecture 6: Bottom-Up Parsing
9/9/09

Bottom-up parsing

- A more powerful parsing technology
- LR grammars -- more expressive than LL
 - can handle left-recursive grammars, virtually all programming languages
 - Easier to express programming language syntax
- Shift-reduce parsers
 - construct right-most derivation of program
 - automatic parser generators (e.g., yacc, CUP, ocamllyacc)
 - detect errors as soon as possible
 - allows better error recovery

Top-down parsing

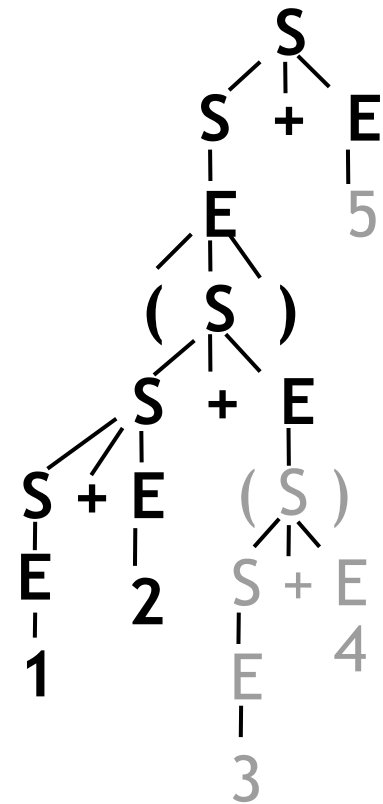
$(1+2+(3+4))+5$

$S \rightarrow S + E \rightarrow E + E \rightarrow (S) + E \rightarrow (S + E) + E \rightarrow (S + E + E) + E \rightarrow (E + E + E) + E \rightarrow (1 + E + E) + E \rightarrow (1 + 2 + E) + E \dots$

- In left-most derivation, entire tree above a token (2) has to be expanded when encountered
- Must be able to predict productions!

$$S \rightarrow S + E \mid E$$

$$E \rightarrow n \mid (S)$$



Bottom-up parsing

- Right-most derivation -- backward
 - Start with the tokens
 - End with the start symbol

$$\begin{array}{l} S \rightarrow S + E \mid E \\ E \rightarrow \text{number} \mid (S) \end{array}$$

$$\begin{aligned} (1+2+(3+4))+5 &\leftarrow (E+2+(3+4))+5 \leftarrow (S+2+(3+4))+5 \\ &\leftarrow (S+E+(3+4))+5 \leftarrow (S+(3+4))+5 \leftarrow (S+(E+4))+5 \\ &\leftarrow (S+(S+4))+5 \leftarrow (S+(S+E))+5 \leftarrow (S+(S))+5 \leftarrow (S \\ &+E)+5 \leftarrow (S)+5 \leftarrow E+5 \leftarrow S+E \leftarrow S \end{aligned}$$

Progress of bottom-up parsing

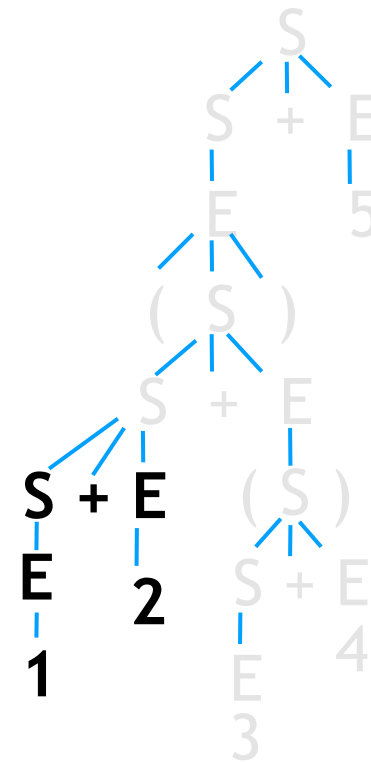
<div> <div>↑</div> <div>right-most derivation</div> <div>↓</div> </div>	$(1+2+(3+4))+5 \leftarrow$		$(1+2+(3+4))+5$
	$(\textcolor{blue}{E}+2+(3+4))+5 \leftarrow$	$(1$	$+2+(3+4))+5$
	$(\textcolor{blue}{S}+2+(3+4))+5 \leftarrow$	$(1$	$+2+(3+4))+5$
	$(S+\textcolor{blue}{E}+(3+4))+5 \leftarrow$	$(1+2$	$+ (3+4))+5$
	$(\textcolor{blue}{S}+(3+4))+5 \leftarrow$	$(1+2+(3$	$+4))+5$
	$(S+(\textcolor{blue}{E}+4))+5 \leftarrow$	$(1+2+(3$	$+4))+5$
	$(S+(\textcolor{blue}{S}+4))+5 \leftarrow$	$(1+2+(3$	$+4))+5$
	$(S+(S+\textcolor{blue}{E}))+5 \leftarrow$	$(1+2+(3+4$	$))+5$
	$(S+(\textcolor{blue}{S}))+5 \leftarrow$	$(1+2+(3+4$	$))+5$
	$(S+\textcolor{blue}{E}))+5 \leftarrow$	$(1+2+(3+4)$	$))+5$
	$(\textcolor{blue}{S}))+5 \leftarrow$	$(1+2+(3+4)$	$))+5$
	$\textcolor{blue}{E}+5 \leftarrow$	$(1+2+(3+4))$	$+5$
	$S+\textcolor{blue}{E} \leftarrow$	$(1+2+(3+4))+5$	
	$\textcolor{blue}{S}$	$(1+2+(3+4))+5$	

Bottom-up parsing

- $(1+2+(3+4))+5 \leftarrow (E+2+(3+4))+5 \leftarrow$
 $(S+2+(3+4))+5 \leftarrow (S+E+(3+4))+5 \dots$

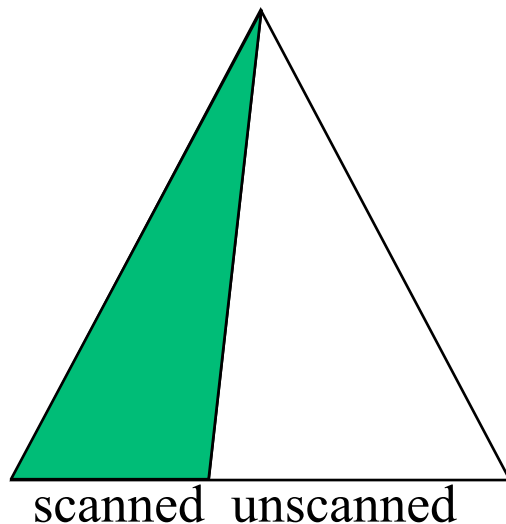
$S \rightarrow S + E \mid E$
 $E \rightarrow \text{number} \mid (S)$

- Advantage of bottom-up parsing:
select productions using more
information

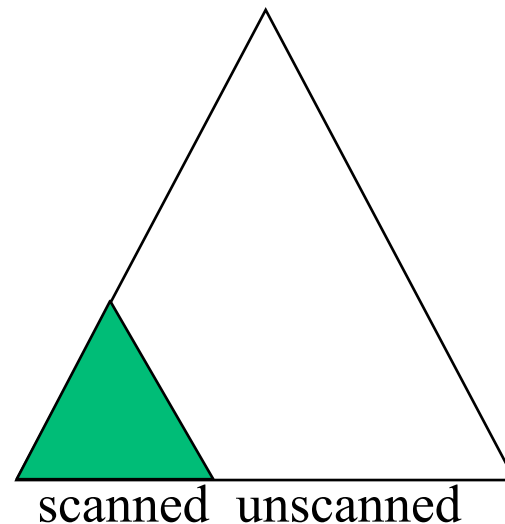


Top-down vs. Bottom-up

Bottom-up: Don't need to figure out as much of the parse tree for a given amount of input



Top-down



Bottom-up

Shift-reduce parsing

- Parsing is a sequence of *shift* and *reduce* operations
- Parser state is a stack of terminals and non-terminals (grows to the right)
- Unconsumed input is a string of terminals
- Current derivation step is always stack+input

Derivation step

(1+2+(3+4))+5 ←

(**E**+2+(3+4))+5 ←

(**S**+2+(3+4))+5 ←

(S+**E**+(3+4))+5 ←

stack

(E

(S

(S+E

unconsumed input

(1+2+(3+4))+5

+2+(3+4))+5

+2+(3+4))+5

+(3+4))+5

Shift-reduce parsing

- Parsing is a sequence of *shifts* and *reduces*

-
- Shift** : move lookahead token to stack. No effect on derivation.

stack	input	action
(1+2+(3+4))+5	<i>shift</i> 1
(1	+2+(3+4))+5	

- Reduce** : Replace symbols γ in top of stack with non-terminal symbol X , corresponding to production $X \rightarrow \gamma$ (pop γ , push X). Reduces rightmost nonterminal.

stack	input	action
(<u>S+E</u>	+(3+4))+5	<i>reduce</i> $S \rightarrow S+E$
(S	+(3+4))+5	

Shift-reduce parsing

$$S \rightarrow S + E \mid E$$

$$E \rightarrow \text{number} \mid (S)$$

derivation	stack	input stream	action
$(1+2+(3+4))+5 \leftarrow$		$(1+2+(3+4))+5$	<i>shift</i>
$(1+2+(3+4))+5 \leftarrow$	($1+2+(3+4))+5$	<i>shift</i>
$(1+2+(3+4))+5 \leftarrow$	(1	$+2+(3+4))+5$	<i>reduce</i> $E \rightarrow \text{num}$
$(E+2+(3+4))+5 \leftarrow$	(E	$+2+(3+4))+5$	<i>reduce</i> $S \rightarrow E$
$(S+2+(3+4))+5 \leftarrow$	(S	$+2+(3+4))+5$	<i>shift</i>
$(S+2+(3+4))+5 \leftarrow$	(S+	$2+(3+4))+5$	<i>shift</i>
$(S+2+(3+4))+5 \leftarrow$	(S+2	$+(3+4))+5$	<i>reduce</i> $E \rightarrow \text{num}$
$(S+E+(3+4))+5 \leftarrow$	(S+E	$+(3+4))+5$	<i>reduce</i> $S \rightarrow S+E$
$(S+(3+4))+5 \leftarrow$	(S	$+(3+4))+5$	<i>shift</i>
$(S+(3+4))+5 \leftarrow$	(S+	$(3+4))+5$	<i>shift</i>
$(S+(3+4))+5 \leftarrow$	(S+($3+4))+5$	<i>shift</i>
$(S+(3+4))+5 \leftarrow$	(S+(3	$+4))+5$	<i>reduce</i> $E \rightarrow \text{num}$

Problem

- How do we know which action to take -- whether to shift or reduce, and which production?
- Sometimes **can** reduce but **shouldn't**.
 - e.g., $X \rightarrow \varepsilon$ can *always* be reduced
- Sometimes can reduce in more than one way.

Action Selection Problem

- Given stack σ and look-ahead symbol b , should parser:
 - **shift** b onto the stack (making it σb)
 - **reduce** some production $X \rightarrow \gamma$ assuming that stack has the form $\alpha \gamma$ (making it αX)
- If stack has form $\alpha \gamma$, should apply reduction $X \rightarrow \gamma$ (or shift) depending on stack prefix α
 - α is different for different possible reductions, since γ 's have different length.
 - How to keep track of possible reductions?

Parser States

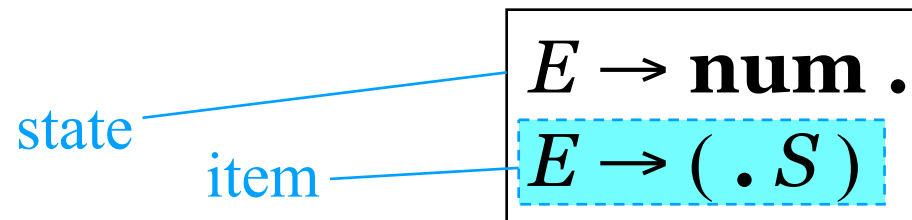
- Goal: know what reductions are legal at any given point.
- Idea: summarize all possible stacks σ (and prefixes α) as a finite parser **state**
 - Parser state is computed by a DFA that reads in the stack σ
 - Accept states of DFA: unique reduction!
- Summarizing discards information
 - affects what grammars parser handles
 - affects size of DFA (number of states)

LR(0) parser

- **L**eft-to-right scanning, **R**ight-most derivation, “**zero**” look-ahead characters
- Too weak to handle most language grammars (e.g., “sum” grammar)
- But will help us understand shift-reduce parsing...

LR(0) states

- A state is a set of *items* keeping track of progress on possible upcoming reductions
- An *LR(0) item* is a production from the language with a separator “.” somewhere in the RHS of the production



- Stuff before “.” is already on stack (beginnings of possible γ 's to be reduced)
- Stuff after “.” : what we might see next
- The prefixes α represented by state itself

An LR(0) grammar: non-empty lists

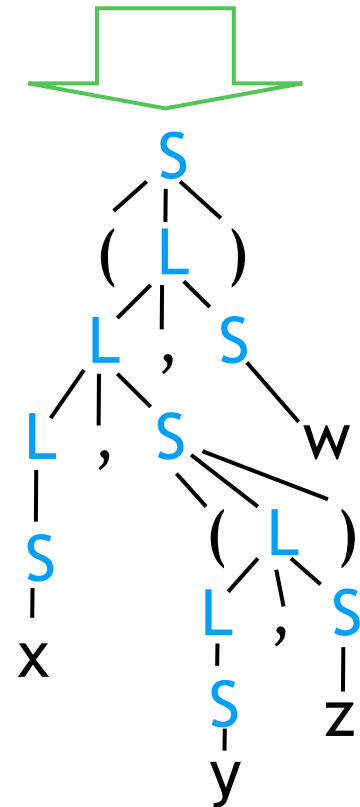
$$S \rightarrow (L) \mid id$$

$$L \rightarrow S \mid L , S$$

x (x,y) (x, (y,z), w)

(((x))) (x, (y, (z, w)))

(x, (y,z), w)



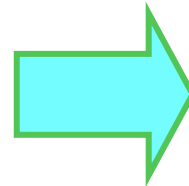
Start State & Closure

$$\begin{array}{l} S \rightarrow (L) \mid id \\ L \rightarrow S \mid L, S \end{array}$$

DFA start state

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

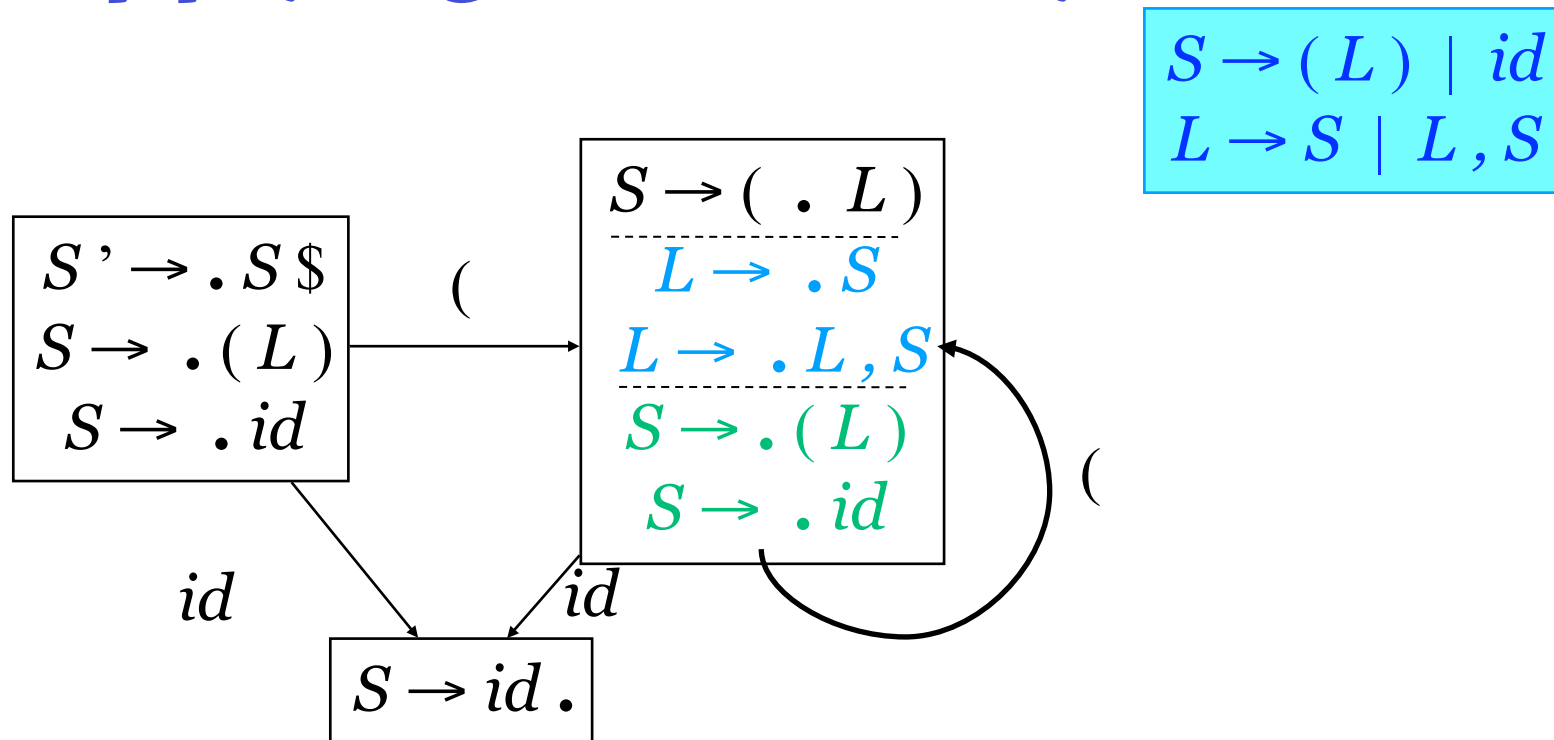
closure


$$\begin{array}{l} S' \rightarrow . S \$ \\ S \rightarrow . (L) \\ S \rightarrow . id \end{array}$$

Constructing a DFA to read stack:

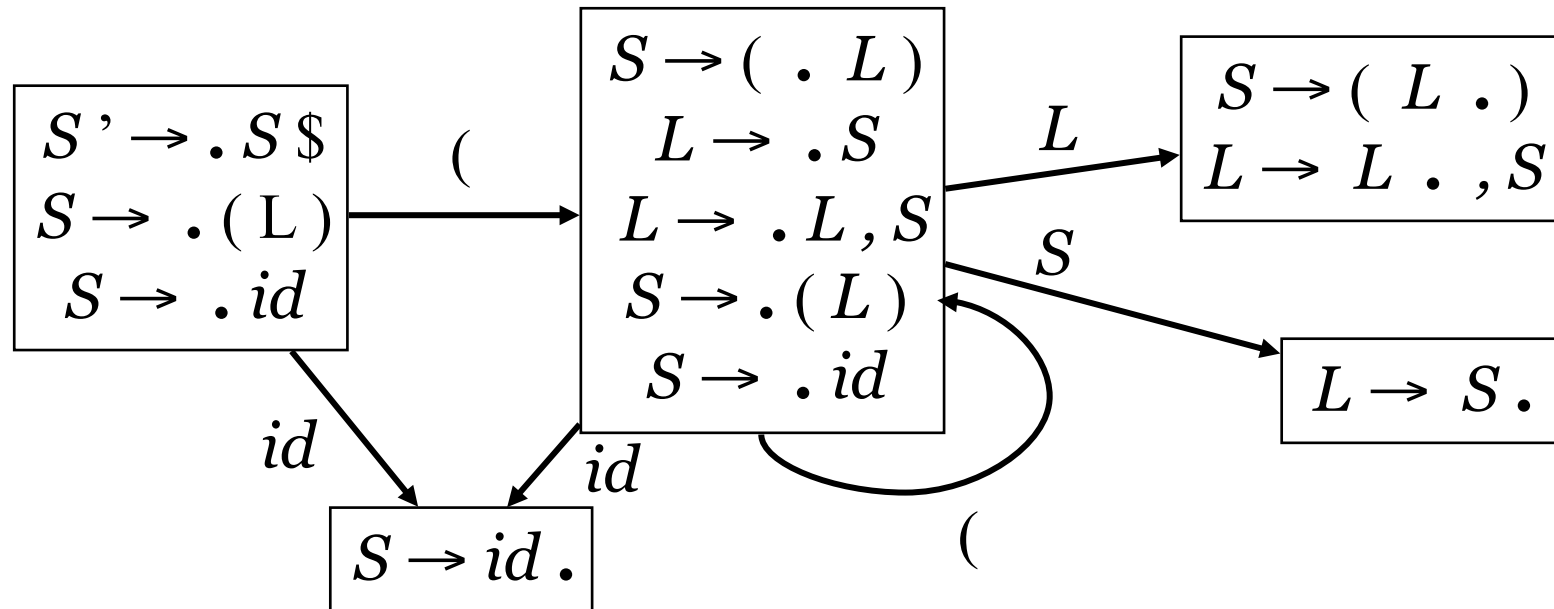
- First step: augment grammar with prod'n $S' \rightarrow S \$$
- Start state of DFA: empty stack = $S' \rightarrow . S \$$
- *Closure* of a state adds items for all productions whose LHS occurs in an item in the state, just after “.”
 - set of possible productions to be reduced next
 - Added items have the “.” located at the beginning: no symbols for these items on the stack yet

Applying terminal symbols



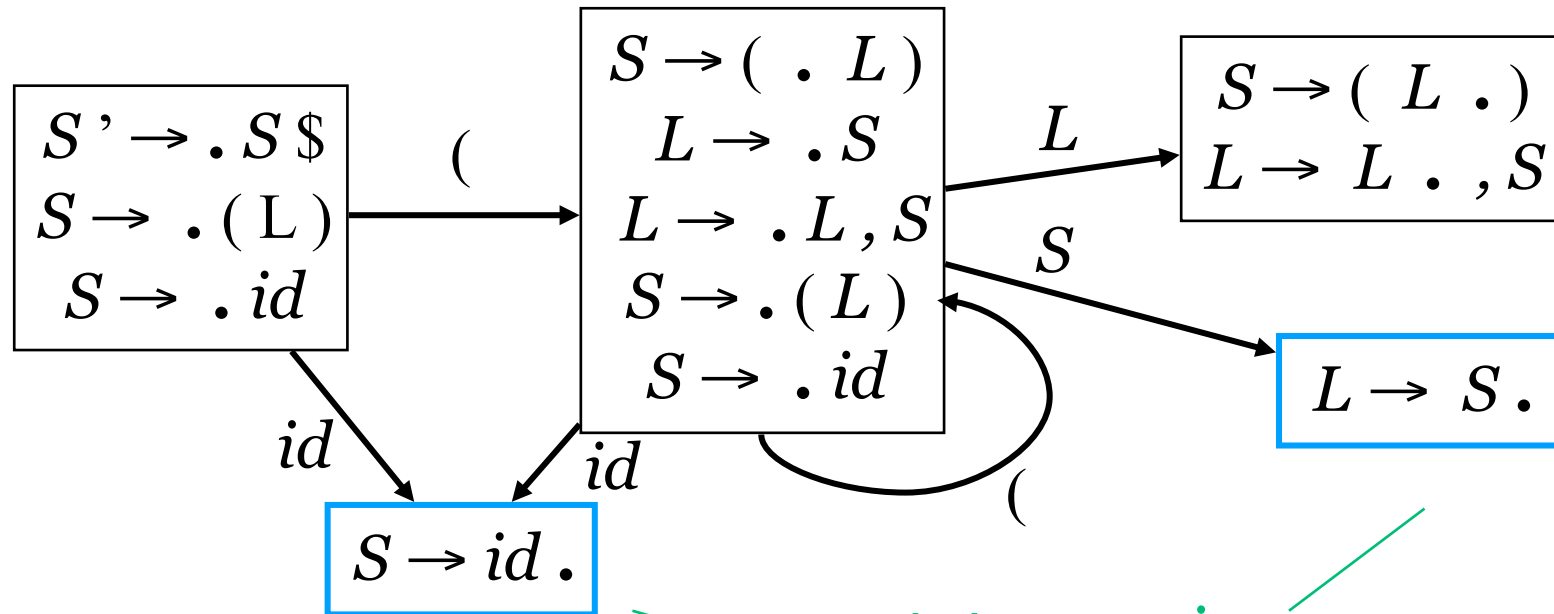
In new state, include all items that have appropriate input symbol just after dot, advance dot in those items, *and take closure*.

Applying non-terminals



- Non-terminals on stack treated just like terminals (but added by reductions)

Applying reduce actions

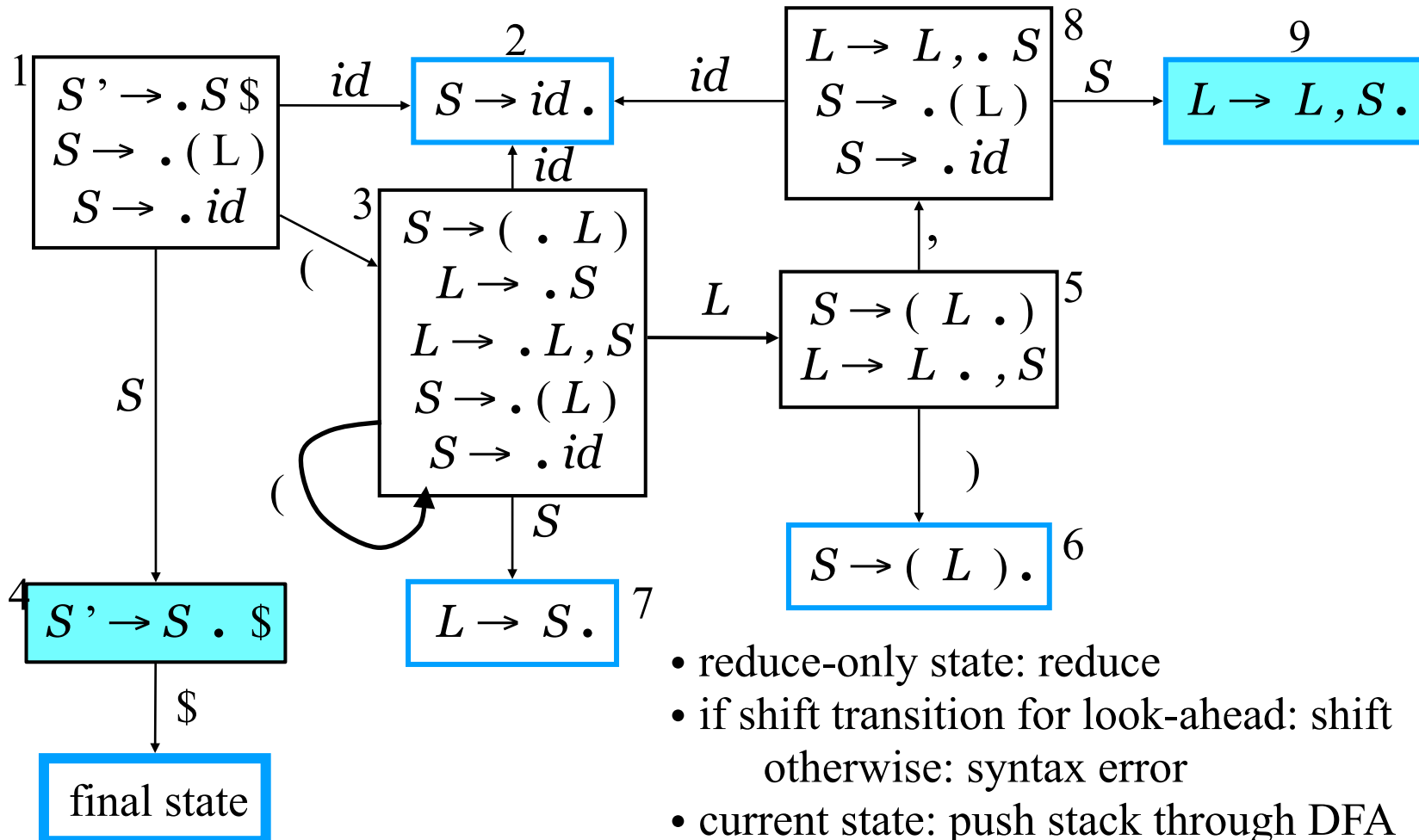


*states causing
reductions*

- Pop RHS off stack, replace with LHS X ($X \rightarrow \gamma$), rerun DFA (e.g. (x))

Full DFA (Appel)

$S \rightarrow (L) \mid id$
 $L \rightarrow S \mid L, S$



$$S \rightarrow (L) \mid id$$

$$L \rightarrow S \mid L, S$$

Parsing example: ((x),y)

derivation	stack	input	action
((x),y) ←	₁	((x),y)	shift, goto 3
((x),y) ←	₁ (₃	(x),y)	shift, goto 3
((x),y) ←	₁ (₃ (₃	x),y)	shift, goto 2
((x),y) ←	₁ (₃ (₃ x ₂),y)	reduce $S \rightarrow id$
((S),y) ←	₁ (₃ (₃ S ₇),y)	reduce $L \rightarrow S$
((L),y) ←	₁ (₃ (₃ L ₅),y)	shift, goto 6
((L),y) ←	₁ (₃ (₃ L ₅) ₆	,y)	reduce $S \rightarrow (L)$
(S ,y) ←	₁ (₃ S ₇	,y)	reduce $L \rightarrow S$
(L ,y) ←	₁ (₃ L ₅	,y)	shift, goto 8
(L,y) ←	₁ (₃ L ₅ , ₈	y)	shift, goto 9
(L,y) ←	₁ (₃ L ₅ , ₈ y ₂)	reduce $S \rightarrow id$
(L, S) ←	₁ (₃ L ₅ , ₈ S ₉)	reduce $L \rightarrow L, S$
(L) ←	₁ (₃ L ₅)	shift, goto 6
(L) ←	₁ (₃ L ₅) ₆		reduce $S \rightarrow (L)$
S	₁ S ₄	\$	done

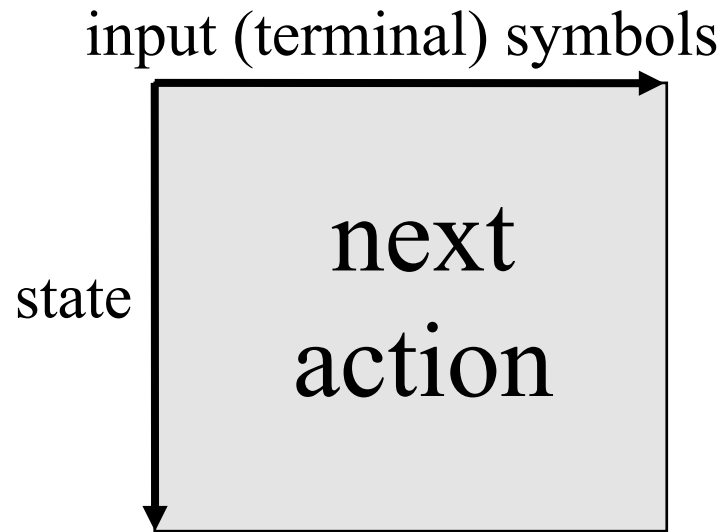
Optimization

- Don't need to rerun DFA from beginning on every reduction
- On reducing $X \rightarrow \gamma$ with stack $\alpha\gamma$:
 - pop γ off stack, revealing prefix α *and state*
 - take single step in DFA from top state
 - push X onto stack with new DFA state

$((L)$, y) state = 6

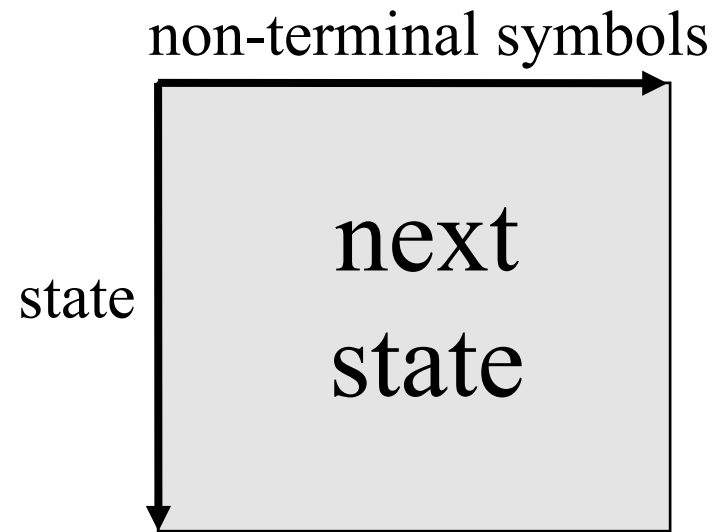
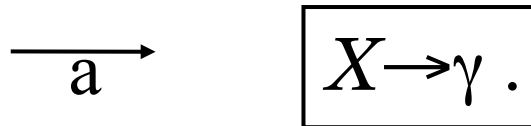
$(S$, y) state = ?

Implementation: LR parsing table



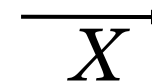
Action table

Used at every step to
decide whether to shift
or reduce



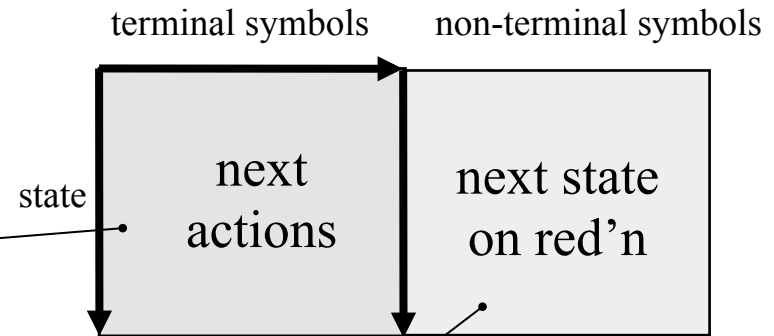
Goto table

Used only when reducing,
to determine next state



Shift-reduce parsing table

action table



1. shift and goto state n
 2. reduce using $X \rightarrow \gamma$
 - pop symbols γ off stack
 - using state label of top (end) of stack, look up X in **goto table** and go to that state
- DFA + stack = push-down automaton (PDA)

List grammar parsing table

	()	<i>id</i>	,	\$	<i>S</i>	<i>L</i>
1	s3		s2			g4	
2	$S \rightarrow id$	$S \rightarrow id$	$S \rightarrow id$	$S \rightarrow id$	$S \rightarrow id$		
3	s3		s2			g7	g5
4					accept		
5		s6		s8			
6	$S \rightarrow (L)$	$S \rightarrow (L)$	$S \rightarrow (L)$	$S \rightarrow (L)$	$S \rightarrow (L)$		
7	$L \rightarrow S$	$L \rightarrow S$	$L \rightarrow S$	$L \rightarrow S$	$L \rightarrow S$		
8	s3		s2			g9	
9	$L \rightarrow L,S$	$L \rightarrow L,S$	$L \rightarrow L,S$	$L \rightarrow L,S$	$L \rightarrow L,S$		

Shift-reduce parsing

- Grammars can be parsed bottom-up using a DFA + stack
 - DFA processes stack σ to decide what reductions might be possible given
 - *shift-reduce parser* or *push-down automaton (PDA)*
 - Compactly represented as *LR parsing table*
- State construction converts grammar into states that decide action to take