

Compilation

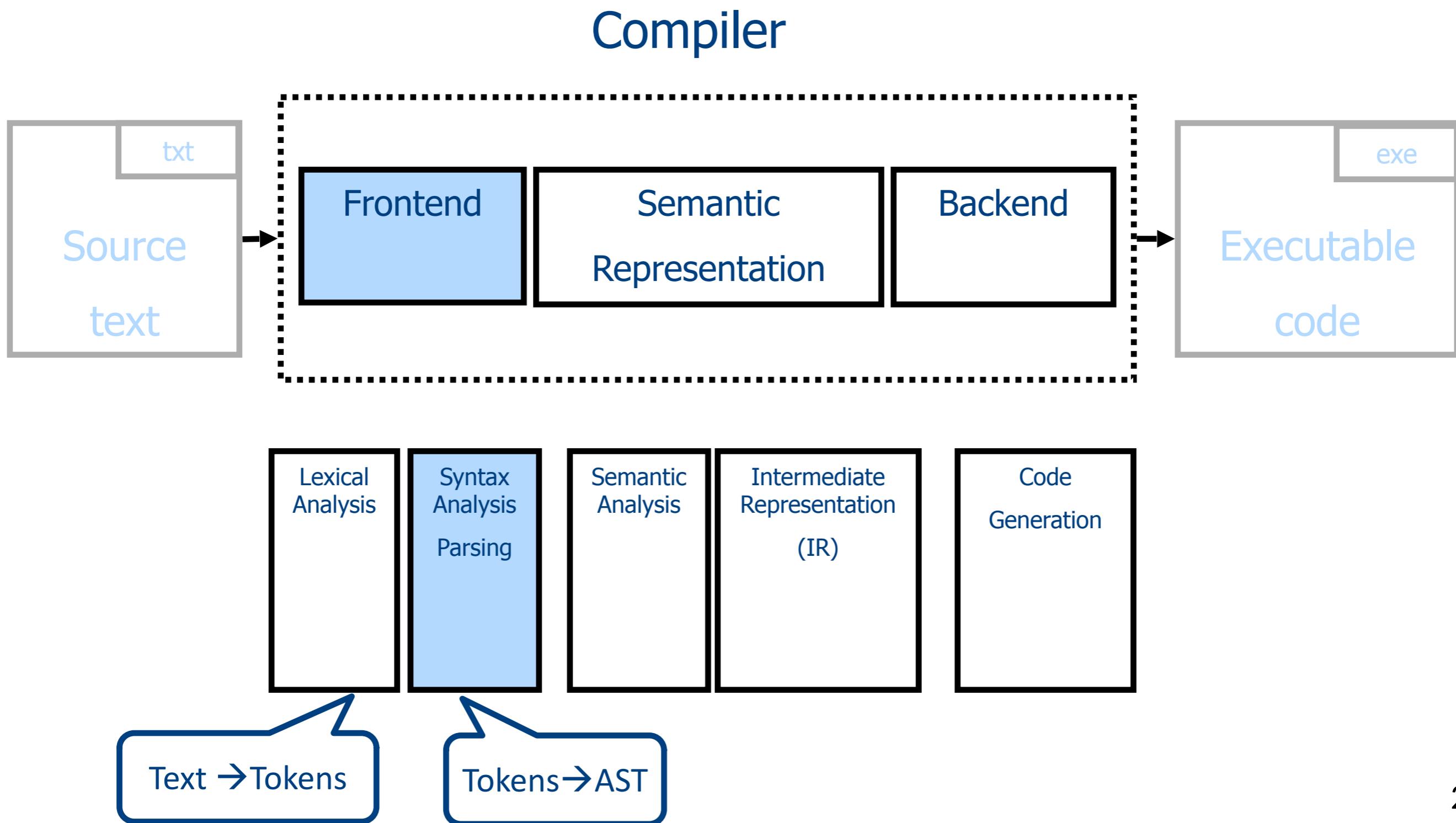
0368-3133

Lecture 4:

Syntax Analysis:

Bottom-Up Parsing (Part II)

Conceptual Structure of a Compiler



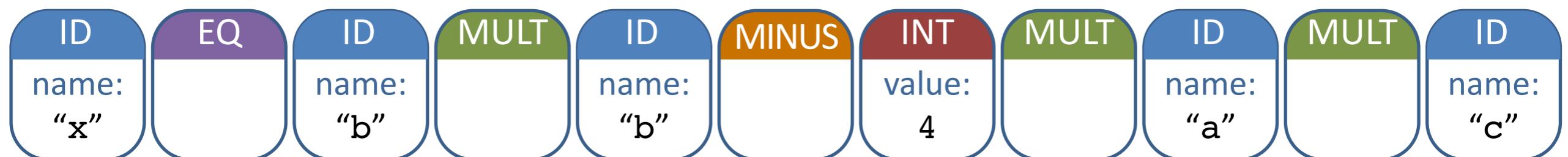
Lexing: from Characters to Tokens

txt

```
x = b*b - 4*a*c
```

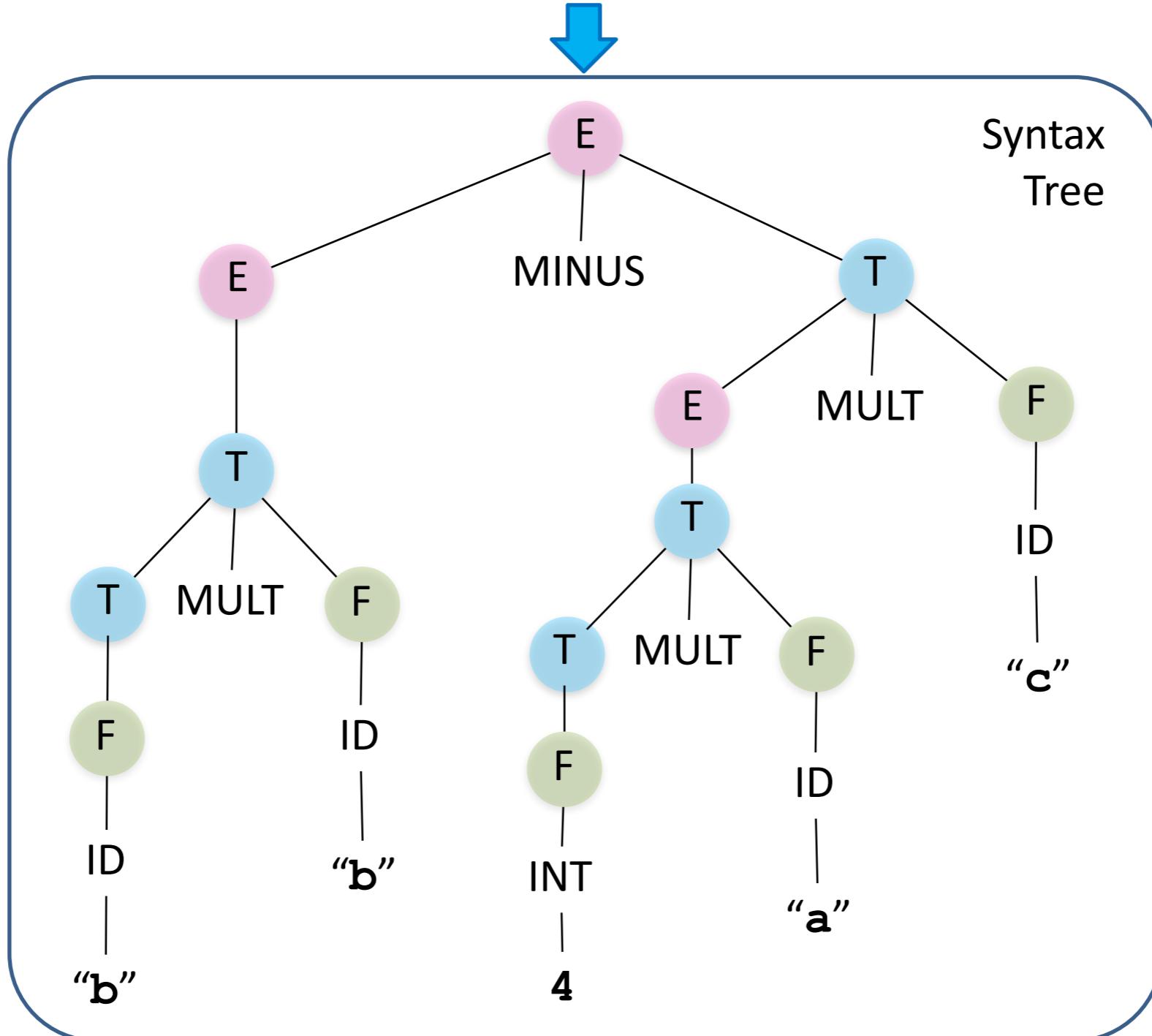


Token Stream



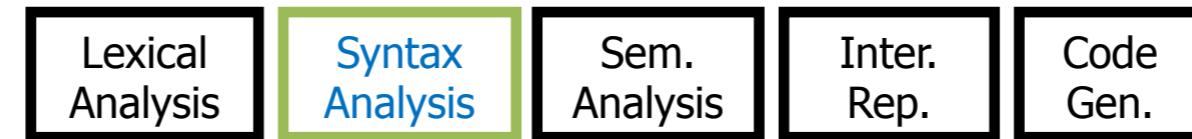
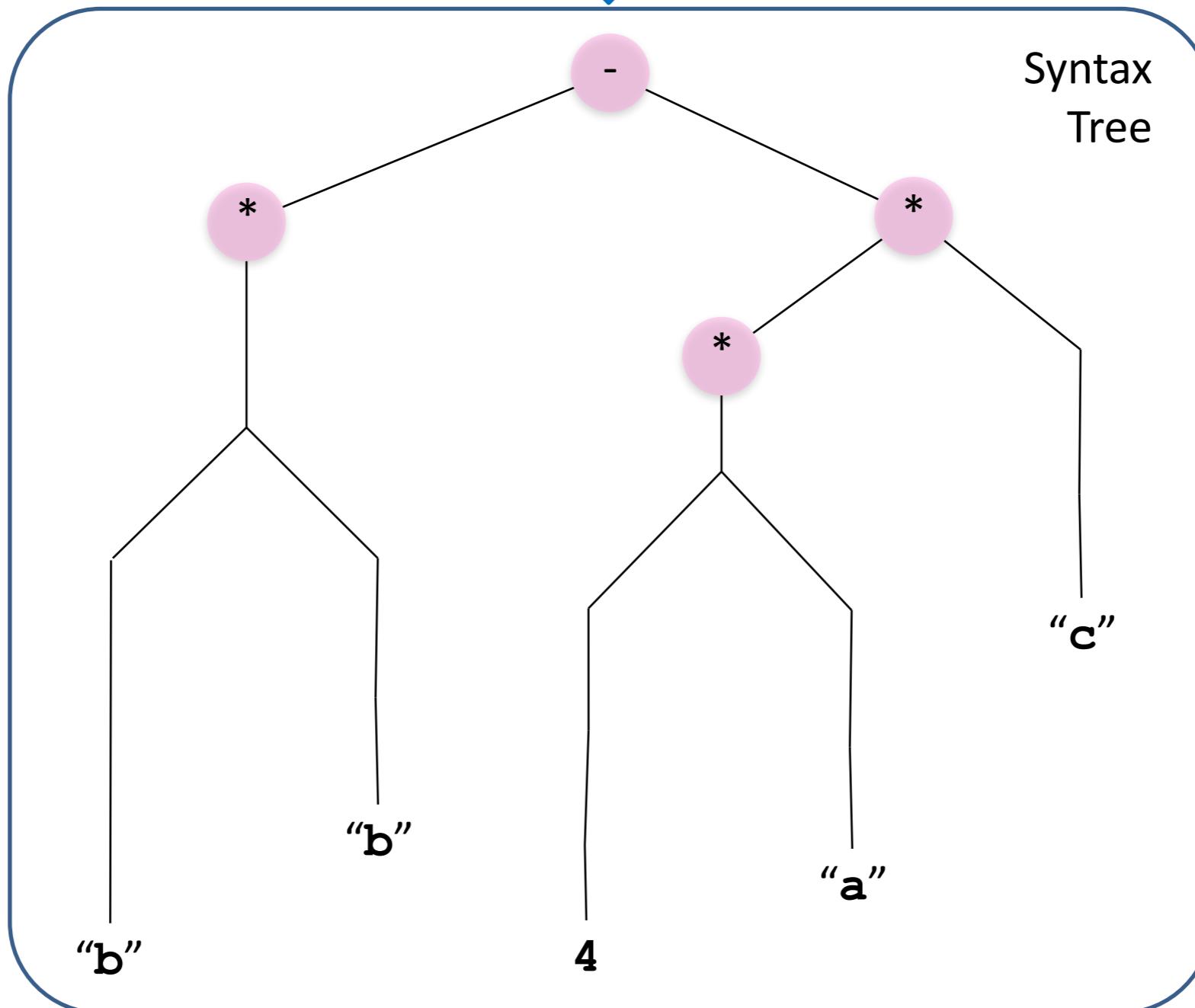
Parsing: from Tokens to Syntax Tree

`<ID,"b"> <MULT> <ID,"b"> <MINUS> <INT,4> <MULT> <ID,"a"> <MULT> <ID,"c">`



Parsing: from Tokens to AST

$\langle \text{ID}, "b" \rangle \langle \text{MULT} \rangle \langle \text{ID}, "b" \rangle \langle \text{MINUS} \rangle \langle \text{INT}, 4 \rangle \langle \text{MULT} \rangle \langle \text{ID}, "a" \rangle \langle \text{MULT} \rangle \langle \text{ID}, "c" \rangle$

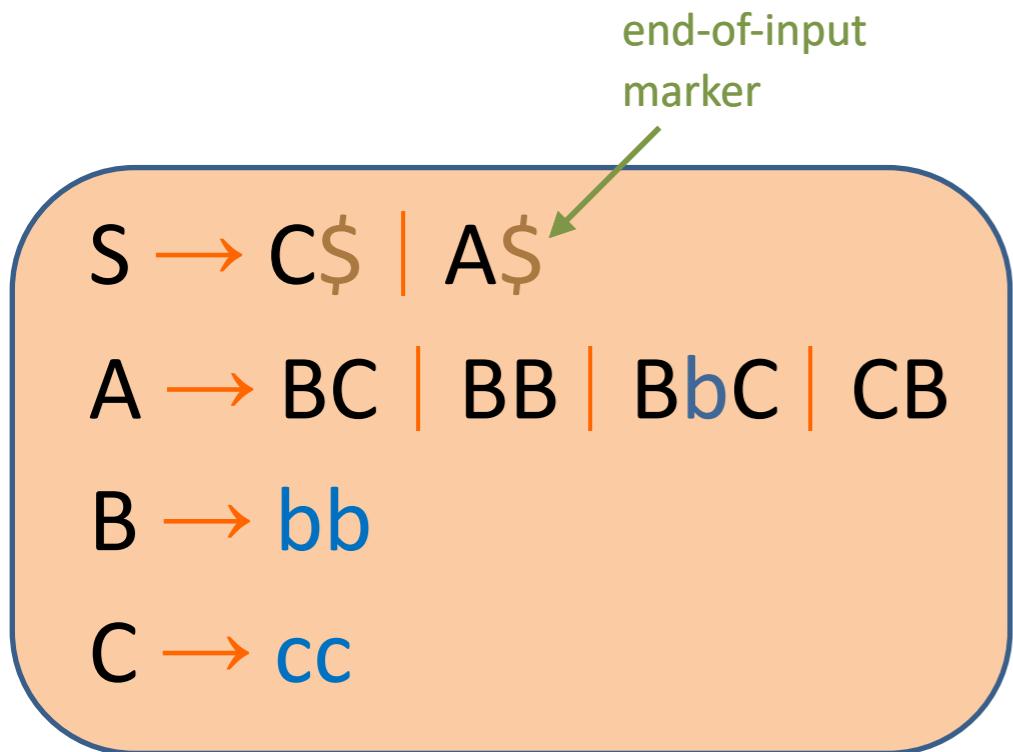


Bottom-Up Parsing

- Bottom-Up parsers
 - Scan input from **left** to right
 - Construct the parse tree in a **bottom-up** manner
 - Add a node when
 - All its leaves (tokens) have been read
 - All its children have been added
 - Reduce* the input to the start symbol
 - Reduction in reverse order of the **rightmost** derivation
 - Using a lookahead of ***k*** tokens
- A grammar that can be parsed like that is called an ***LR(k)*** grammar

* Reduction (step): $\beta \alpha \mu \rightarrow \beta \text{A} \mu$ if $A \rightarrow \alpha$ be a production rule

Example

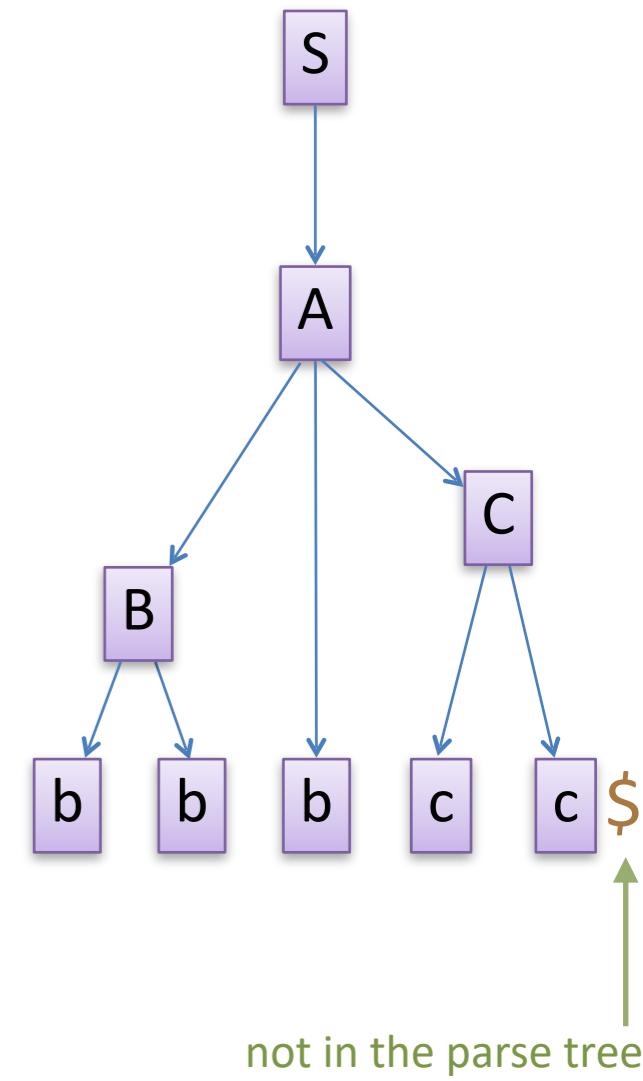


Rightmost derivation

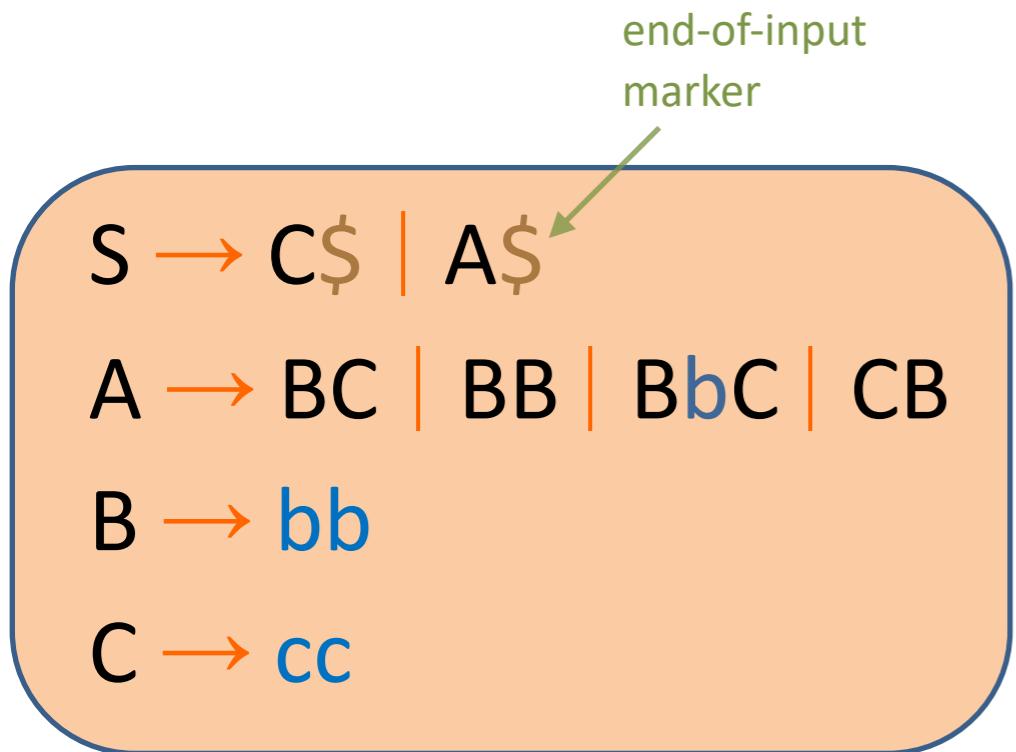
$w = \text{bbbcc}\$$

$S \Rightarrow$
 $A\$ \Rightarrow$
 $BbC\$ \Rightarrow$
 $Bbcc\$ \Rightarrow$
 $\text{bbbcc}\$$

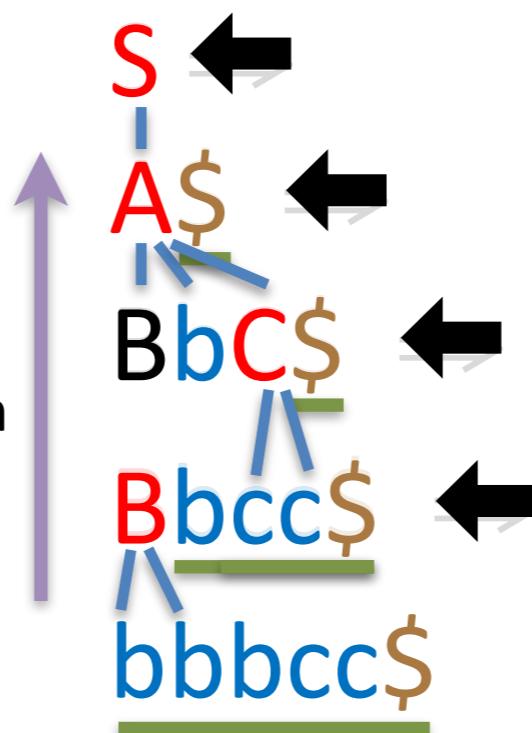
Parse tree



Example



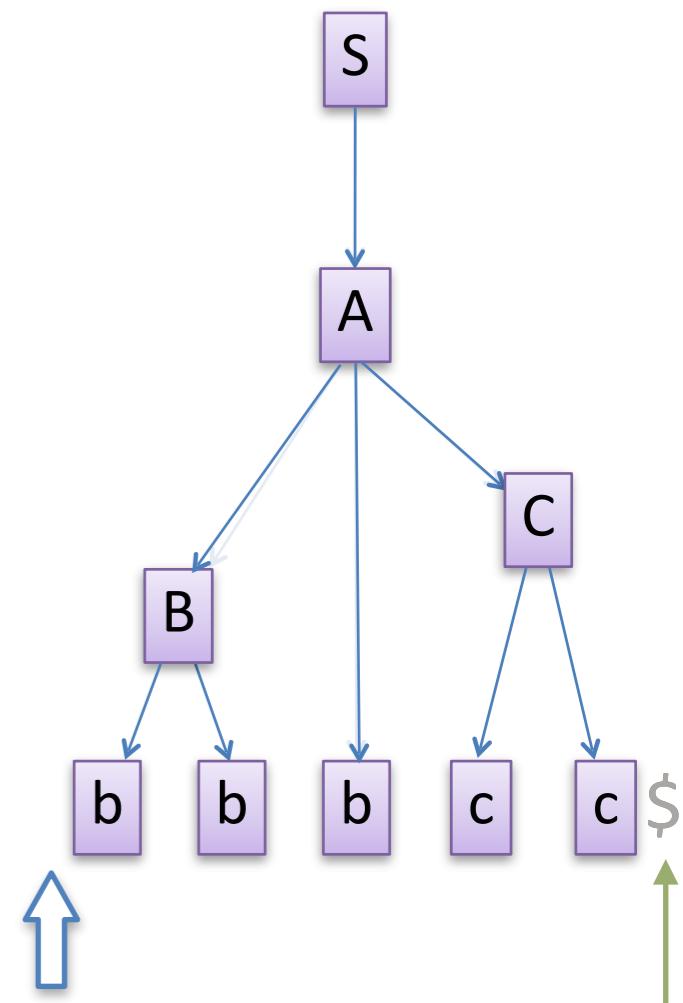
desired reduction



Rightmost derivation

$w = \text{bbbcc\$}$

Parse tree



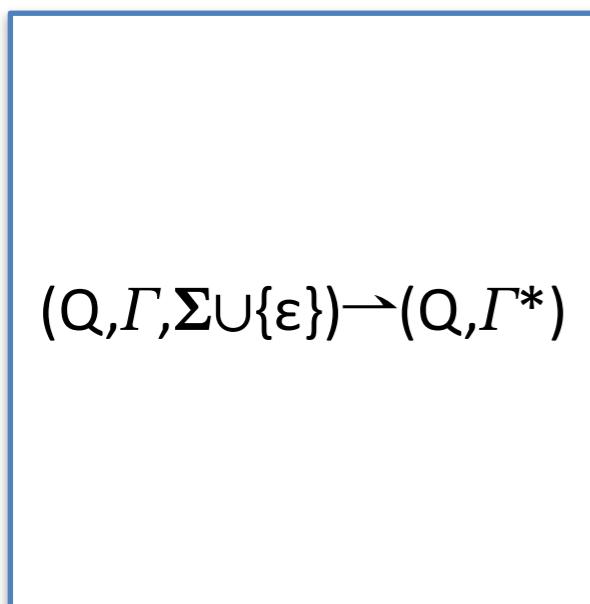
not in the parse tree

How does the Parser Work?

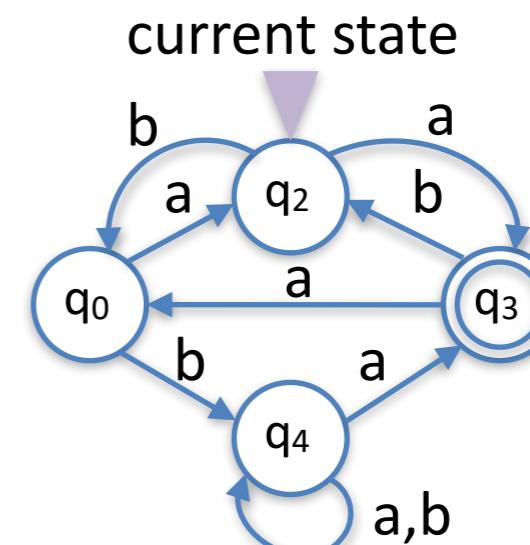
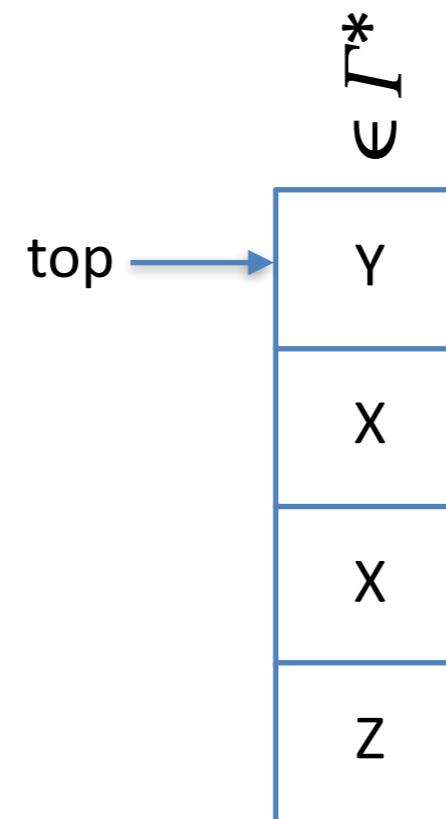
Shift & Reduce (LR) Parsers

deterministic

The parser uses a **pushdown automaton (DPDA)**



Transition function



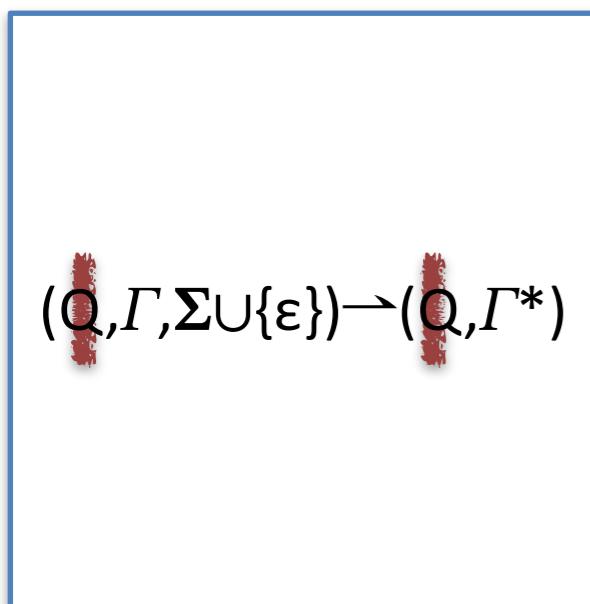
$w = \boxed{b \ a \ a \ b \ a \ a \ b \ a \ a \ b \ b \ b \ a \ b \ a \ b \ b \ b \ a} \in \Sigma^*$

current position

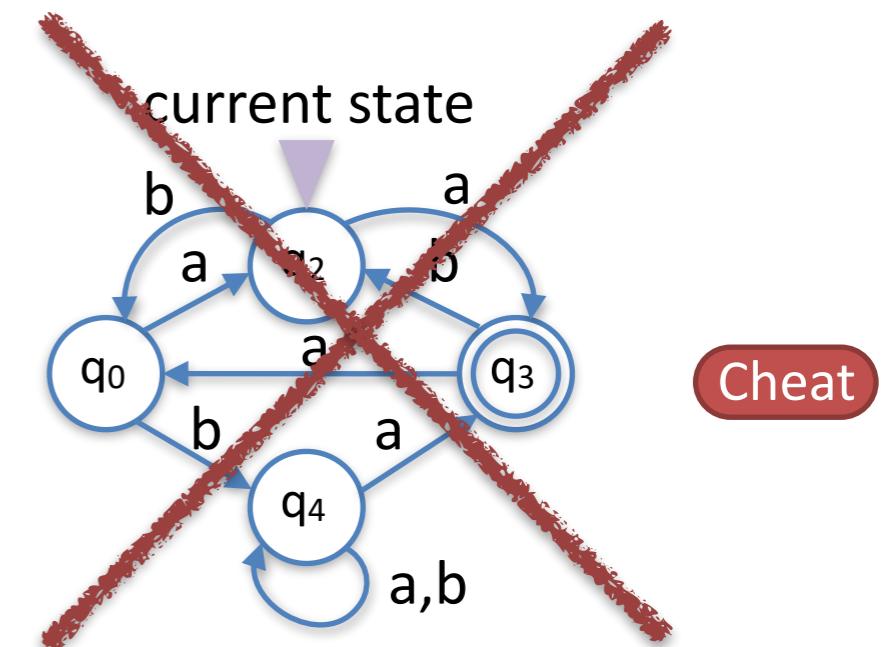
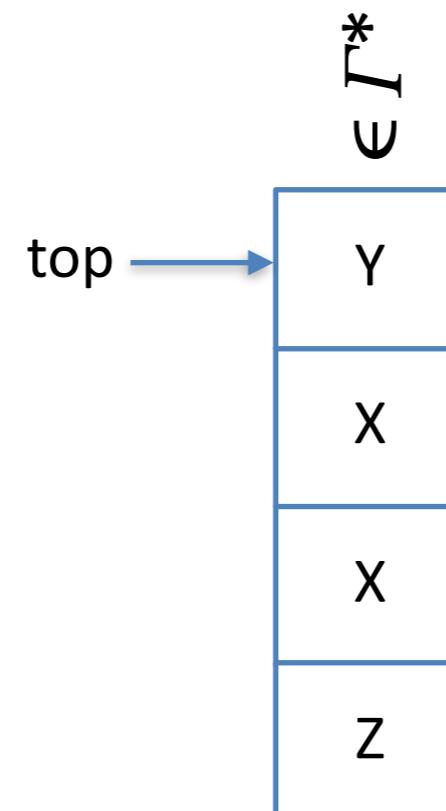
Shift & Reduce (LR) Parsers

deterministic

The parser uses a **pushdown automaton (DPDA)**



Transition function



$w = \boxed{b \ a \ a \ b \ a \ a \ b \ a \ a \ b \ b \ b \ a \ b \ a \ b \ b \ b \ a} \in \Sigma^*$

current position

Shift & Reduce (LR) Parsers

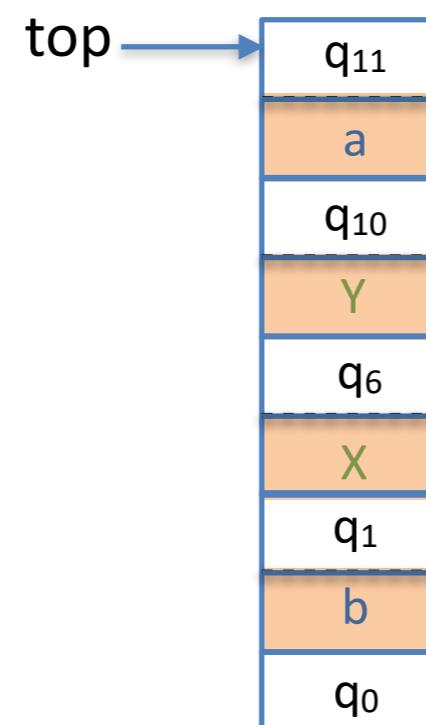
deterministic

The parser uses a **pushdown automaton** (DPDA)

Encodes a (confusingly named)
"LR Automaton"



Transition function table

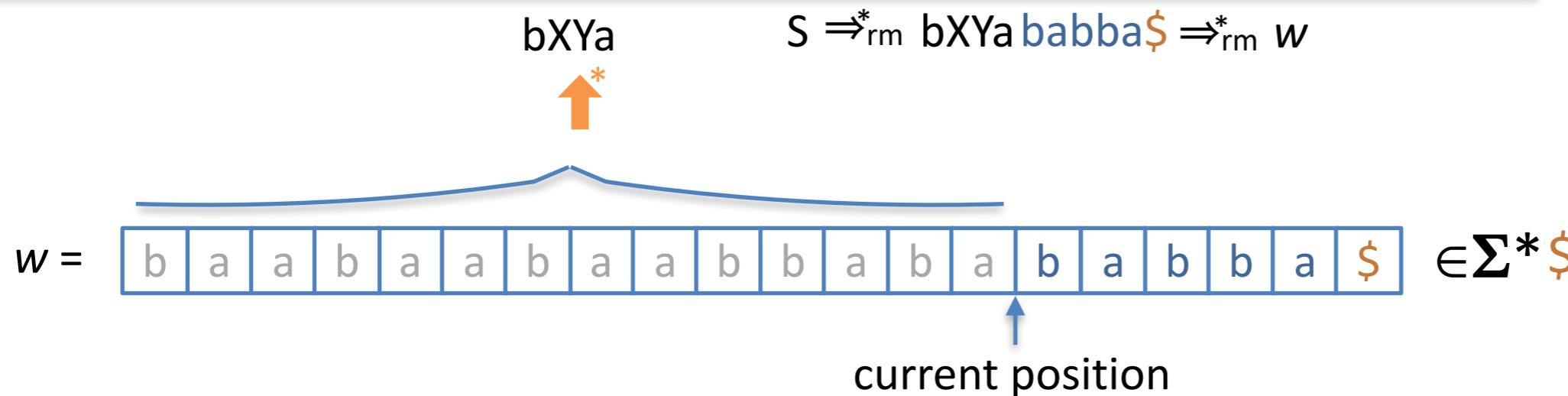


Stack

confusingly named "state"
confusingly named "initial state"

$S \rightarrow bXYba\$$
 $X \rightarrow aabX \mid b$
 $Y \rightarrow Yab \mid ab$

CFL



Shift & Reduce (LR) Parsers

deterministic

The parser uses a **pushdown automaton (DPDA)**

The PDA maintains a **stack** of grammar symbols and states

In each step, we either **shift** or **reduce**

- shift**
 - Moving a symbol from the input to the stack
 - (computing & pushing a new state)

- reduce**
 - Popping the symbols of the right-hand side of the reduced rule (and their paired states)
 - Pushing the non-terminal in the left-hand side of the reduced rule
 - (computing & pushing a new state)

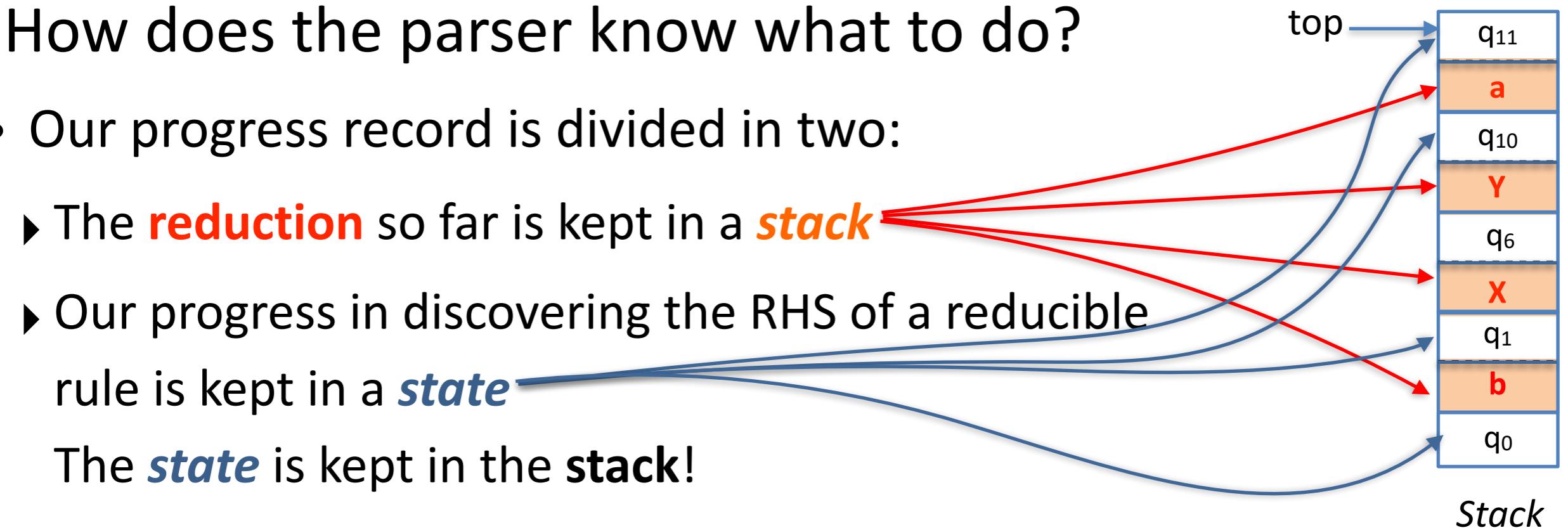
Shift & Reduce (LR) Parsers

deterministic

The parser uses a **pushdown automaton (DPDA)**

How does the parser know what to do?

- Our progress record is divided in two:
 - ▶ The **reduction** so far is kept in a **stack**
 - ▶ Our progress in discovering the RHS of a reducible rule is kept in a **state**
- The **state** is kept in the **stack**!
 - ▶ (*We also keep the position in the input*)
 - ▶ A **transition table** will tell the iterative algorithm “what to do” (**shift?** **reduce?**) based on the current (top) state and the next token(s)



Important Bottom-Up LR-Parsers

- LR(0) – simplest, explains basic ideas
- SLR(1) – simple, explains lookahead
- LR(1) – complicated, very powerful, expensive
- LALR(1) – complicated, powerful enough, used by automatic tools

LR(0) Items

For a production rule

$$N \rightarrow \alpha \beta$$

in the grammar,

$$A \rightarrow \bullet \alpha \beta$$

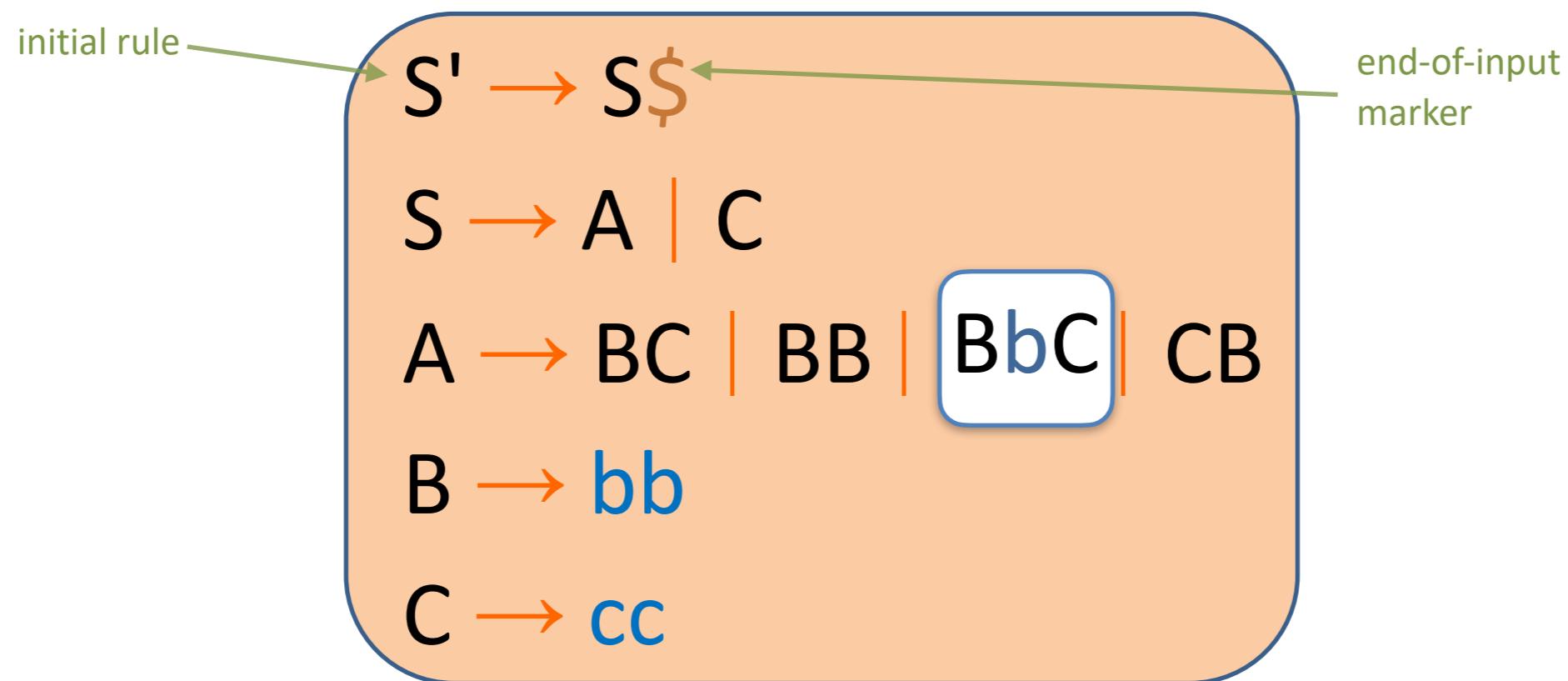
$$A \rightarrow \alpha \bullet \beta \quad (\beta \neq \epsilon)$$

Shift Item

$$A \rightarrow \alpha \beta \bullet$$

Reduce Item

LR(0) Items



$A \rightarrow \bullet BbC$ $A \rightarrow B \bullet bC$ $A \rightarrow Bb \bullet C$ Shift Items

$A \rightarrow BbC \bullet$ $B \rightarrow bb \bullet$

Reduce Items

LR(0) States

- LR(0) state = set of LR(0) items
- A shift state = a set of shift items
 - e.g., $q_3 = \{B \rightarrow \bullet bb, B \rightarrow b\bullet b\}$
 - There are some rules regulating which sets can be shift states.
We'll discuss this in a minute
- A reduce state = a set containing a single reduce item
 - e.g., $q_6 = \{A \rightarrow BC\bullet\}$
 - We'll see later on why it must be a singleton state

S' \rightarrow S\$
S \rightarrow A | C
A \rightarrow BC | BB | BbC | CB
B \rightarrow bb
C \rightarrow cc

LR(0) States - Closure Rules

- A **shift state** containing an item $A \rightarrow \alpha \bullet X \beta$, where X is a non-terminal, must include all items of the form $X \rightarrow \bullet \delta$ for every production rule in the grammar
- For a set of items C , $\text{Clos}(C)$ denotes the closure of C according to the above rule

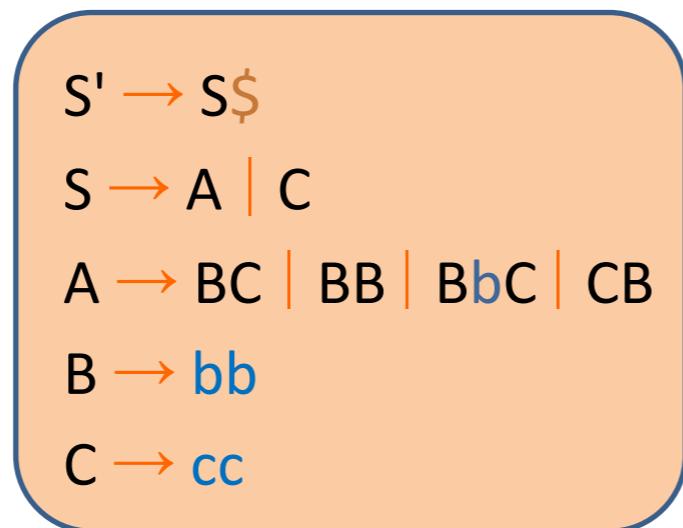
```
S' → S$  
S → A | C  
A → BC | BB | BbC | CB  
B → bb  
C → cc
```

e.g., $\{A \rightarrow B \bullet C, C \rightarrow \bullet cc\}$ is a shift state, but $\{A \rightarrow \bullet BC\}$ isn't

- The **initial state** is the closure of the initial rule
 - e.g., $\text{Clos}(\{S' \rightarrow \bullet S\$}\}) = \{S' \rightarrow \bullet S\$, S \rightarrow \bullet A, S \rightarrow \bullet C, A \rightarrow \bullet BC,$
 $A \rightarrow \bullet BB, A \rightarrow \bullet BbC, A \rightarrow \bullet CB,$
 $B \rightarrow \bullet bb, C \rightarrow \bullet cc\}$

LR(0) Transitions - nextState

- For every symbol (terminal or variable) K, and every (closed) set of items K
- $\text{nextSet}(K, X) = \text{Clos}(\{N \rightarrow \alpha X \bullet \beta \mid (N \rightarrow \alpha \bullet X \beta) \in K\})$



e.g., $K = \{A \rightarrow \bullet BC, C \rightarrow \bullet cc\}$

$\text{nextSet}(K, c) = \{ C \rightarrow c \bullet c \}$

$\text{nextSet}(K, B) = \{ A \rightarrow B \bullet C, C \rightarrow \bullet cc \}$

- For any possible state q, $\text{nextState}(q, X) = \text{nextSet}(q, X)$

LR(0) Actions

- The LR(0) state at the top of the stack determines the action of the PDA
- If the top state is a **shift state**, the PDA performs a **shift**
 - **shift state** = contains only shift items
- If the top state is a **reduce state**, the PDA performs a **reduce**
 - **reduce state** = contains only a (single) reduce item

LR(0) Actions

- **shift** moves (for token t)
 - Fire when the top state q contains a shift item(s) $X \rightarrow \alpha \bullet t \beta$ and the next input token is t
 - Determines the next top state q' based on the current top state q and the shifted token t : compute $\text{nextState}(q, t)$
 - Push t and q' to the stack
- **reduce** moves (for rule $N \rightarrow \alpha$)
 - Fire when the top state q contains a reduce item $N \rightarrow \alpha \bullet$
 - Pop α (and its paired states), resulting in a new top state q_r
 - Determines the next top state q' based on the discovered state q_r and N : compute $\text{nextState}(q_r, N)$
 - Push N and q' to the stack

LR(0) Actions

- **shift** moves (for token t)
 - Fire when the top state q contains a shift item(s) $X \rightarrow \alpha \bullet t \beta$ } Need the stack
and the next input token is t } Looks at input
 - Determines the next top state q' based on the current top state q and the shifted token t : compute $\text{nextState}(q, t)$ } q' can be precomputed (for a given q, t)
 - Push t and q' to the stack } Need the stack
- **reduce** moves (for rule $N \rightarrow \alpha$)
 - Fire when the top state q contains a reduce item $N \rightarrow \alpha \bullet$ } Need the stack
Does **not** Look at input
 - Pop α (and its paired states), resulting in a new top state q_r } Need the stack
 - Determines the next top state q' based on the discovered state q_r and N : compute $\text{nextState}(q_r, N)$ } q' can be precomputed (for a given q_r, t)
 - Push N and q' to the stack } Need the stack

LR(0) Automaton

- We can pre-compute the set of possible states stored in the stack and the potential state transitions ahead of parsing time
- The pre-computation results in an **LR(0) automaton**

Simplified fixpint construction of the LR(0) automaton for $G = (V, T, P, S)$

States' := States := {InitialState}

Transitions = {}

do

States' := States $\cup \{nextState(q, X) \mid q \in States, X \in V \cup T\}$

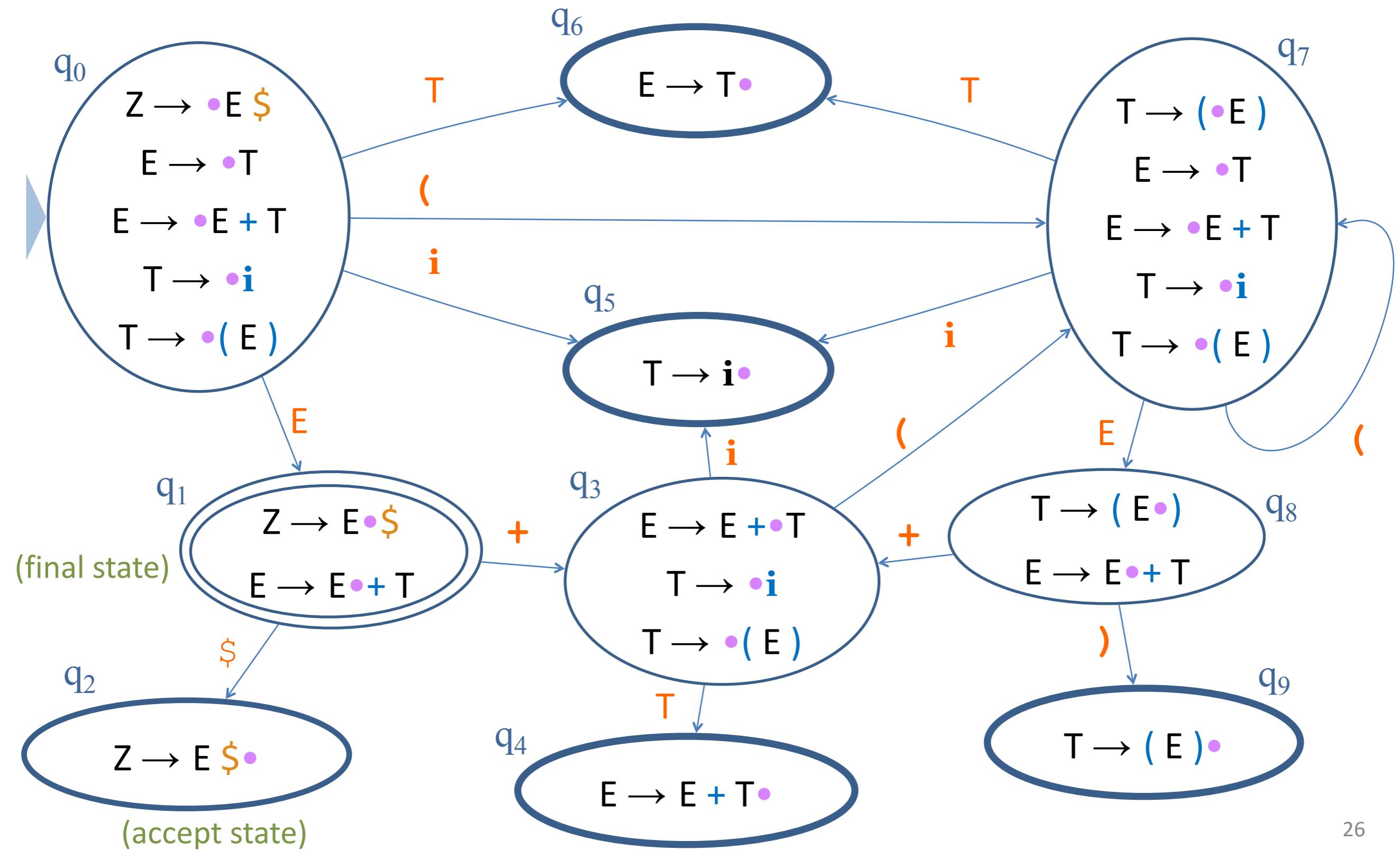
Transitions := Transitions $\cup \{q \xrightarrow{X} q' \mid q, q' \in States, X \in V \cup T, q' = nextState(q)\}$

until States' = States

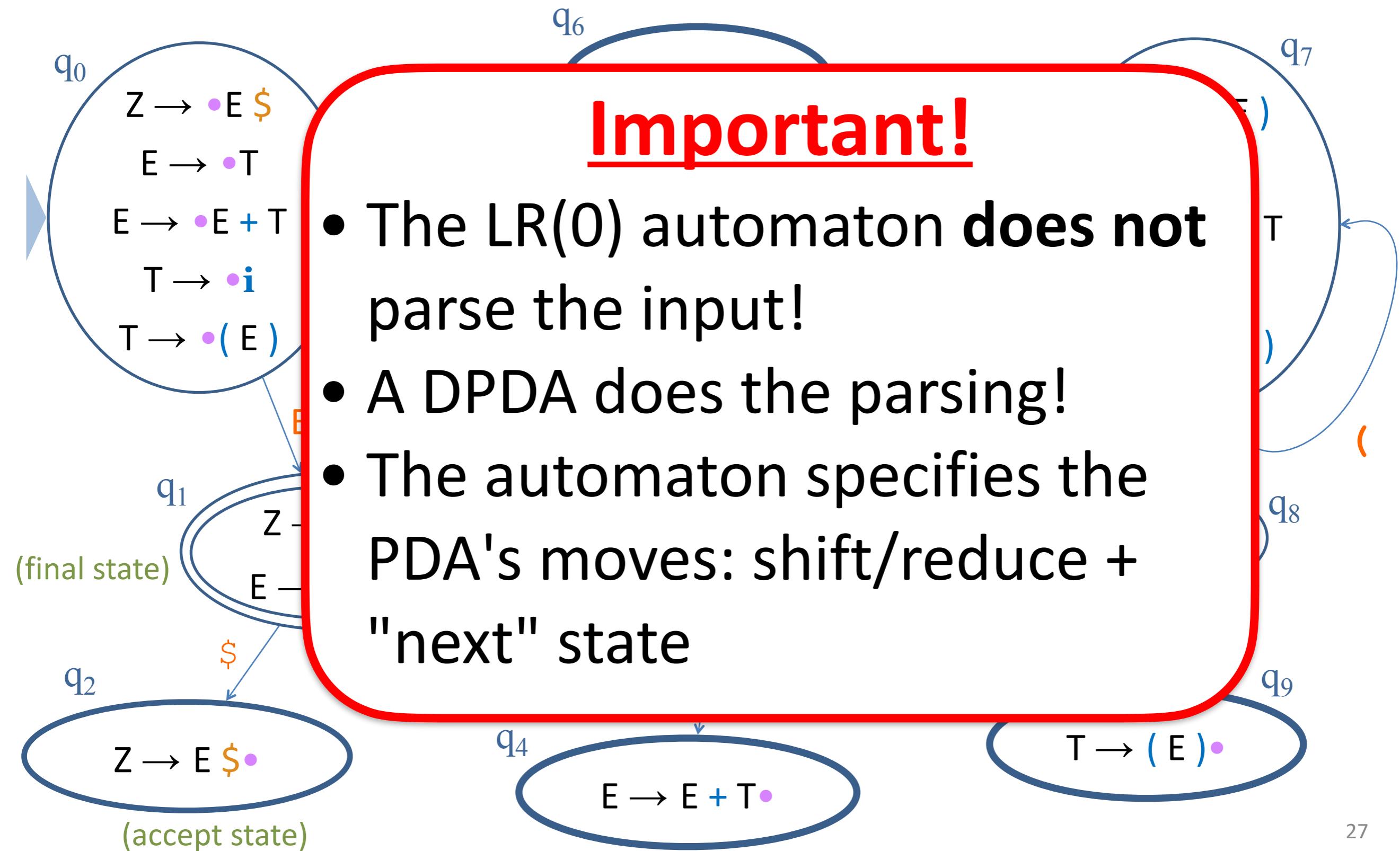
Example Grammar

$$Z \rightarrow E \$$$
$$E \rightarrow T \mid E + T$$
$$T \rightarrow i \mid (E)$$

LR(0) Automaton (Precomputed):

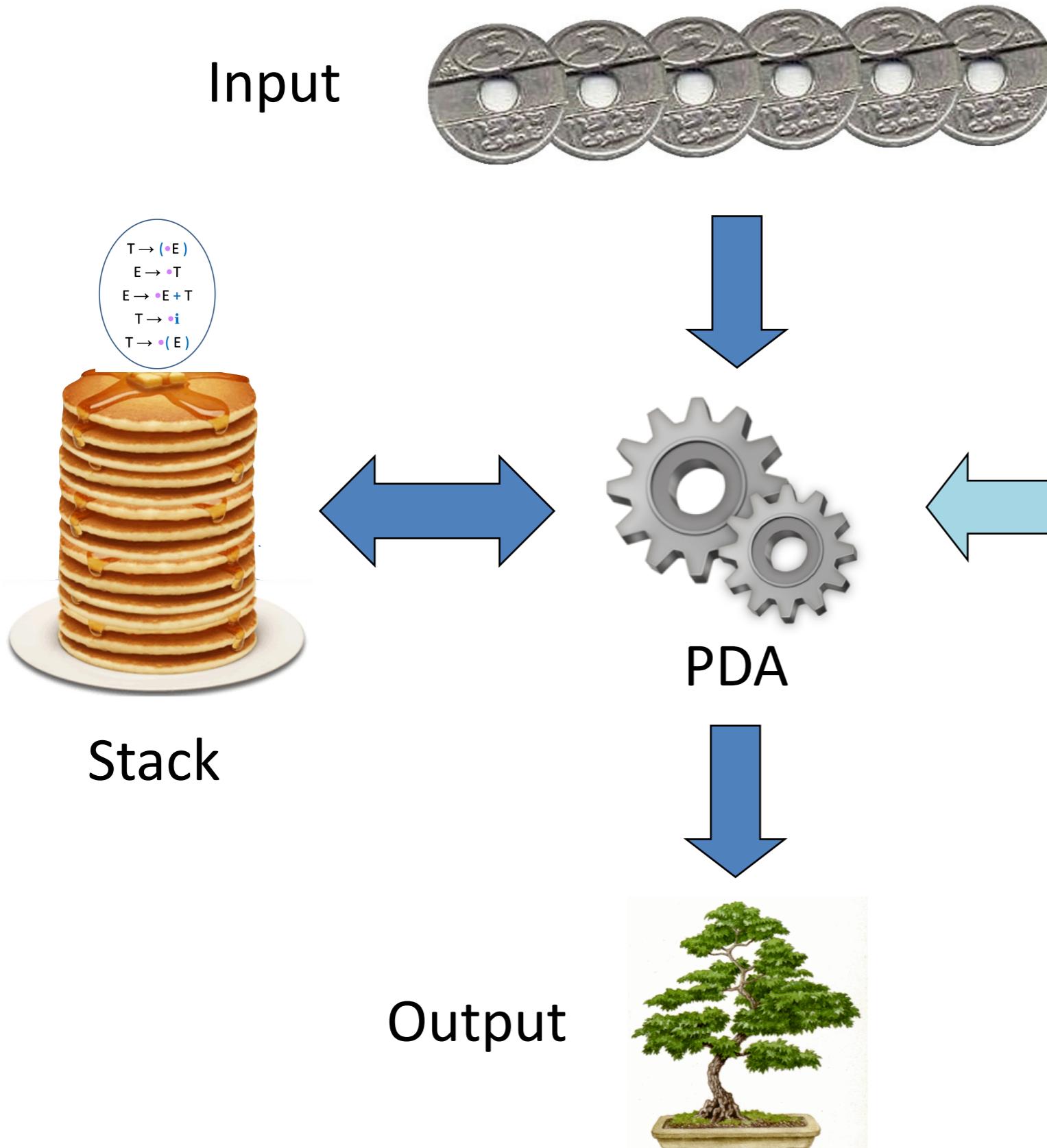


LR(0) Automaton (Precomputed):

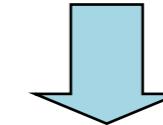


How do parsers represent and use
the LR(0) automaton?

LR(0) Parsing



LR(0)-Automaton

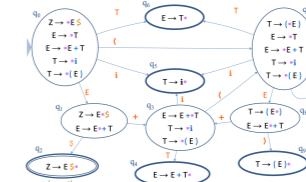


ACTION Table

GOTO Table

Transition table

G



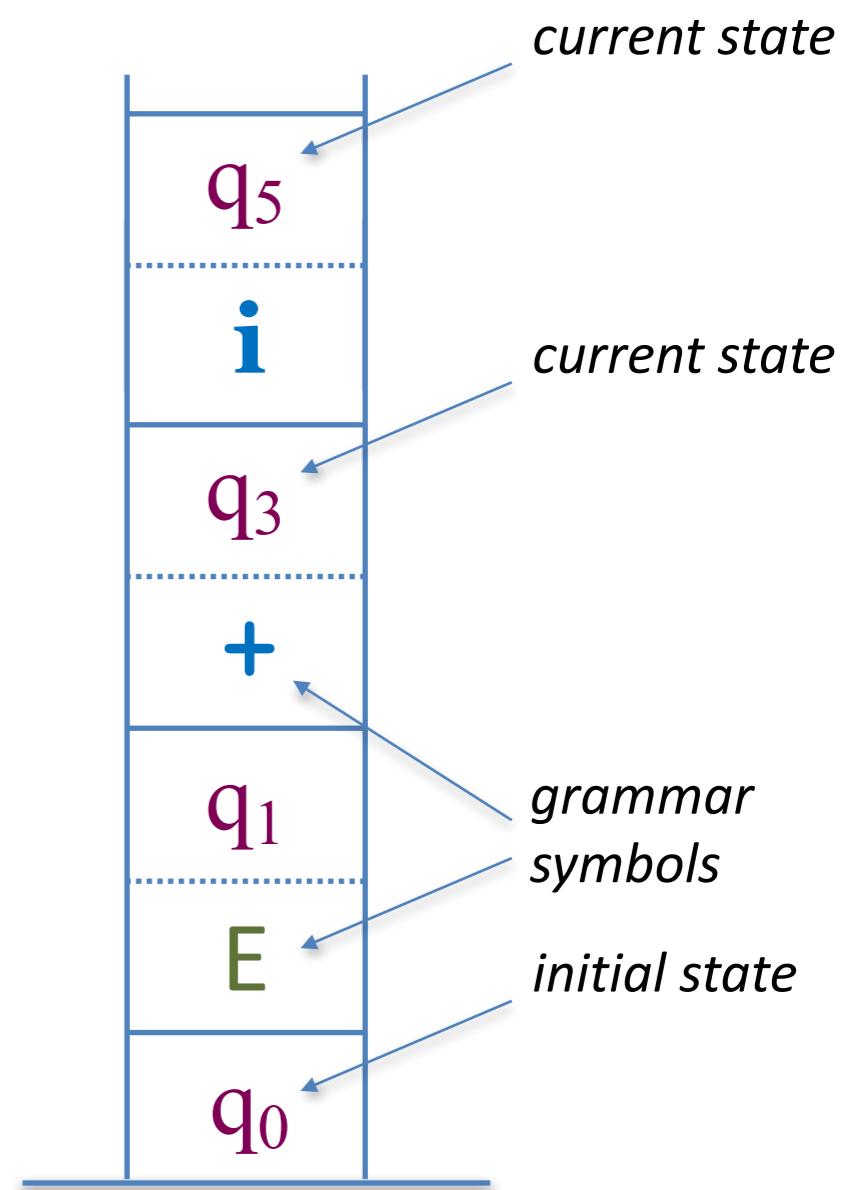
The Stack



- The stack contains states and grammar symbols
- The initial stack contains q_0 only
- The rest of the stack contains pairs of
 - (state, token)
 - (state, nonterminal)

Aside: in a situation where we just want to check the syntax and return true/false, we could do without the symbols

Aside II: sometimes we will look at the stack as an alternating sequence of states and symbols

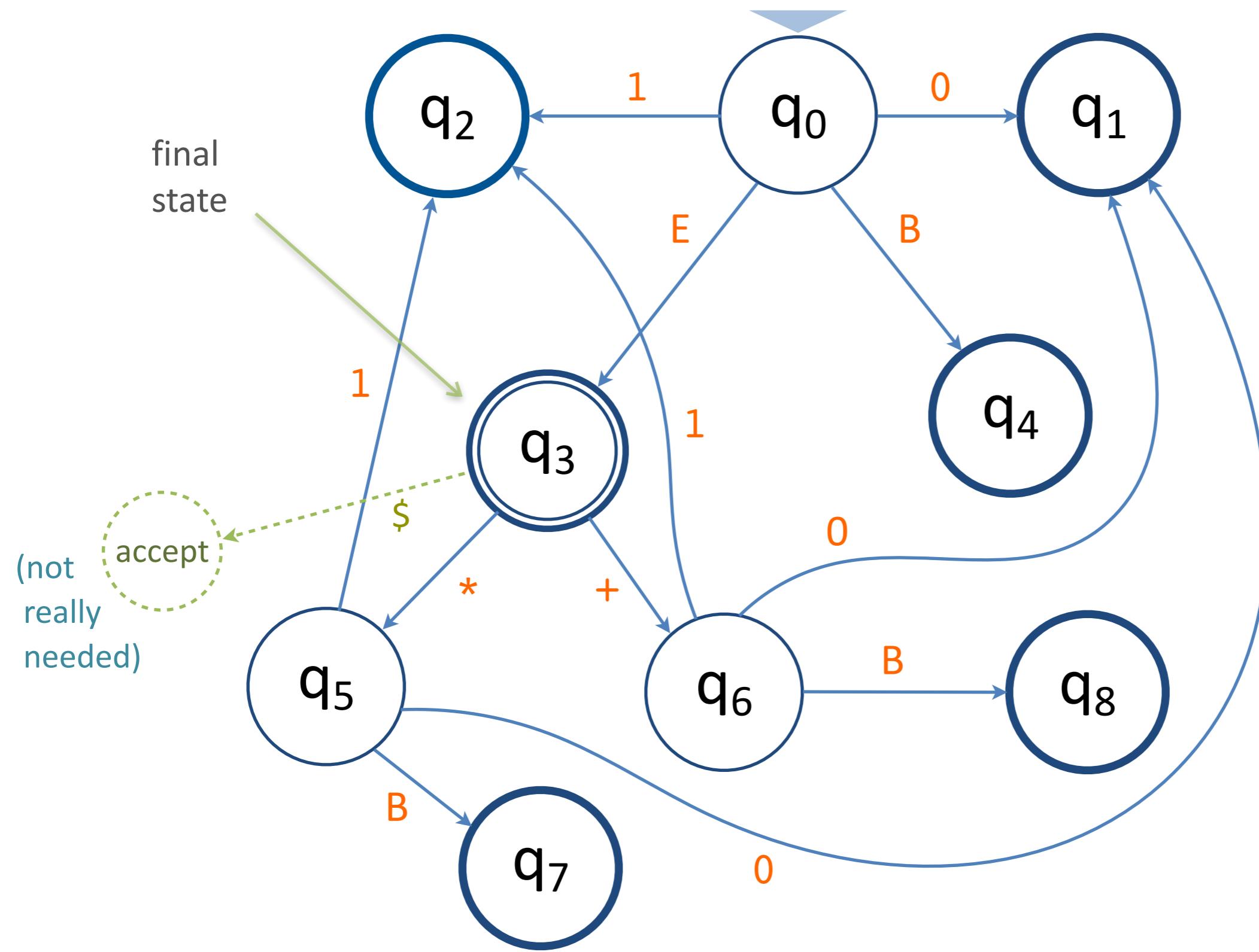


Automaton \rightsquigarrow Tables

$$S \rightarrow E \$$$

$$E \rightarrow E * B \quad | \quad E + B \quad | \quad B$$

$$B \rightarrow 0 \quad | \quad 1$$



q_0

$S \rightarrow \bullet E \$$

$E \rightarrow \bullet E * B$

$E \rightarrow \bullet E + B$

$E \rightarrow \bullet B$

$B \rightarrow \bullet 0 \quad B \rightarrow \bullet 1$

$E \rightarrow B\bullet$

q_3

$S \rightarrow E^* \bullet B \$$

$B \rightarrow \bullet 0 \quad B \rightarrow \bullet 1$

1

The ACTION Table

- At each step we need to decide whether to **shift** the next token to the stack (and move to the appropriate state) or **reduce** a production rule from the grammar
- **ACTION[q, t]** table tells us what to do based on current state **q** and next token **t**:

shift *i* : shift and move to q_i

reduce *j*: reduce according to production rule (*j*)

also: **accept** and **error** conditions

↓
(empty cells)

LR(0) Action Table

	ACTION				
	*	+	0	1	\$
q_0			s_1	s_2	
q_1	r_4	r_4	r_4	r_4	r_4
q_2	r_5	r_5	r_5	r_5	r_5
q_3	s_5	s_6			acc
q_4	r_3	r_3	r_3	r_3	r_3
q_5			s_1	s_2	
q_6			s_1	s_2	
q_7	r_1	r_1	r_1	r_1	r_1
q_8	r_2	r_2	r_2	r_2	r_2

shift “1” and
go to state q_2

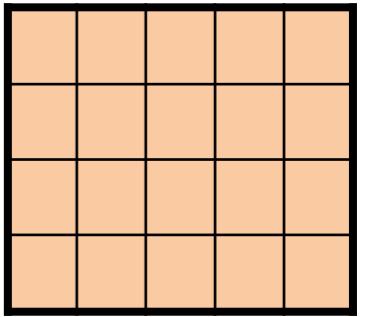
reduce by
rule (4) $B \rightarrow 0$

accept

- (0) $S \rightarrow E \$$
- (1) $E \rightarrow E * B$
- (2) $E \rightarrow E + B$
- (3) $E \rightarrow B$
- (4) $B \rightarrow 0$
- (5) $B \rightarrow 1$

s_i = shift to state q_i

r_j = reduce using rule (j)



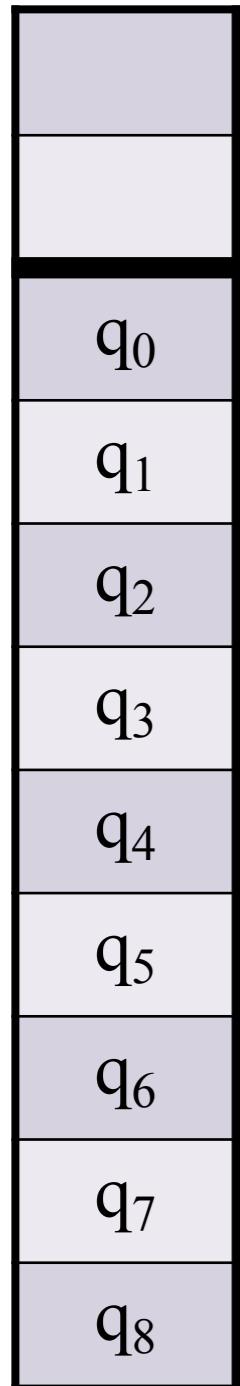
The GOTO Table

- Defines what to do on **reduce** actions
 - ▶ After reducing a right-hand side to the deriving non-terminal, we need to decide what the next state is
- This is determined by the previous state (which is on the stack) and the variable we got
 - ▶ Suppose we **reduce** according to $N \rightarrow \beta$;
 - ▶ We remove β from the stack, and look at the state q that is now at the top. **GOTO**[q, N] specifies the next state.

Note – this can be a little confusing:

- * q is the state **after** popping β
- * N is the left-hand side of the rule just used in **reduce**

LR(0) Goto Table



GOTO	
E	B
g3	g4
	g7
	g8

if after popping the stack due to reduction to B, the top most state is q_0 then **push q_4** to the stack

- (0) $S \rightarrow E \$$
- (1) $E \rightarrow E * B$
- (2) $E \rightarrow E + B$
- (3) $E \rightarrow B$
- (4) $B \rightarrow 0$
- (5) $B \rightarrow 1$

LR(0) Transition Table

	ACTION					GOTO	
	*	+	0	1	\$	E	B
q ₀			s ₁	s ₂		g ₃	g ₄
q ₁	r ₄						
q ₂	r ₅						
q ₃	s ₅	s ₆			acc		
q ₄	r ₃						
q ₅			s ₁	s ₂			g ₇
q ₆			s ₁	s ₂			g ₈
q ₇	r ₁						
q ₈	r ₂						

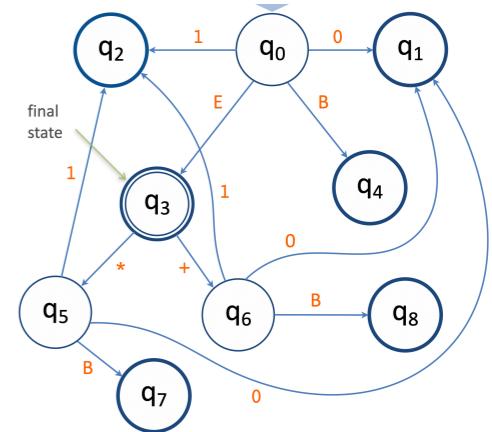
if after popping the stack due to reduction to B, the top most state is q₀ then **push** q₄ to the stack

- (0) S → E \$
- (1) E → E * B
- (2) E → E + B
- (3) E → B
- (4) B → 0
- (5) B → 1

s i = **shift** to state q_i

r j = **reduce** using rule (j)

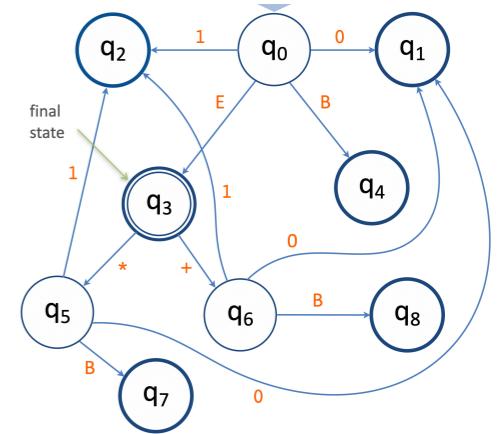
Building the Tables



- A row for each state.
 - If there is a transition from q_i to q_j upon seeing X , then in row q_i and column x we write j .
(except $X = \$$;
see next slide)

State	ACTION					GOTO	
	*	+	0	1	\$	E	B
q ₀			1	2		3	4
q ₁							
q ₂							
q ₃	5	6					
q ₄							
q ₅			1	2			7
q ₆			1	2			8
q ₇							
q ₈							

Building the tables: accept



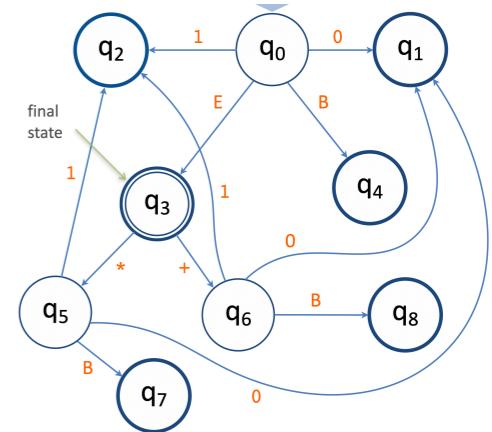
State	ACTION					GOTO	
	*	+	0	1	\$	E	B
q0			1	2		3	4
q1							
q2							
q3	5	6			acc		
q4							
q5			1	2			7
q6			1	2			8
q7							
q8							

- Add **accept** in column **\$** for each state that has the item $S \rightarrow E \bullet \$$ (so-called *final states*).

q3

$S \rightarrow E \bullet \$$
 $E \rightarrow E \bullet ^* B$
 $E \rightarrow E \bullet + B$

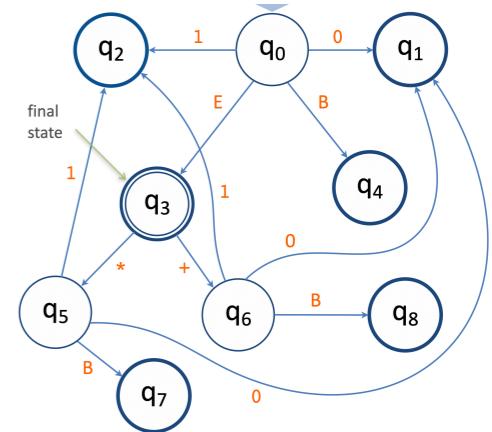
Building the Tables: GOTO



State	ACTION					GOTO	
	*	+	0	1	\$	E	B
q0			1	2		g3	g4
q1							
q2							
q3	5	6			acc		
q4							
q5			1	2			g7
q6			1	2			g8
q7							
q8							

- Any number n in the **GOTO** table becomes **goto n**.

Building the Tables: Shift



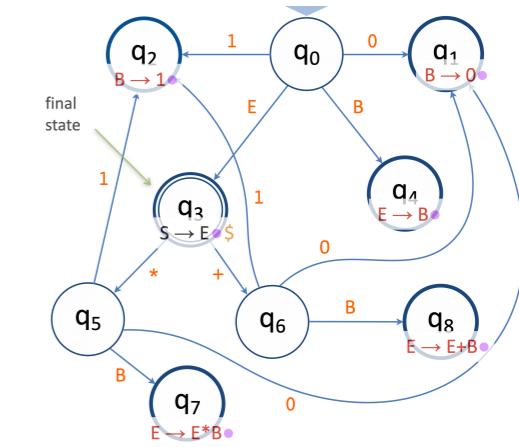
State	ACTION					GOTO	
	*	+	0	1	\$	E	B
q_0			s_1	s_2		g_3	g_4
q_1							
q_2							
q_3	s_5	s_6			acc		
q_4							
q_5			s_1	s_2			g_7
q_6			s_1	s_2			g_8
q_7							
q_8							

- Any number n in the ACTION table becomes shift n.

- (0) $S \rightarrow E \$$
- (1) $E \rightarrow E * B$
- (2) $E \rightarrow E + B$
- (3) $E \rightarrow B$
- (4) $B \rightarrow 0$
- (5) $B \rightarrow 1$

Using the Tables: Reduce

	ACTION				GOTO	
	+	0	1	\$	E	B
q_0		s_1	s_2		g_3	g_4
q_1	r4	r4	r4	r4	r4	
q_2	r5	r5	r5	r5	r5	
q_3	s_5	s_6			acc	
q_4	r3	r3	r3	r3	r3	
q_5			s_1	s_2		g_7
q_6			s_1	s_2		g_8
q_7	r1	r1	r1	r1	r1	
q_8	r2	r2	r2	r2	r2	



For any state which includes an item $A \rightarrow \alpha \bullet$, such that $A \rightarrow \alpha$ is production rule (m):

Fill *all columns* of that state in the ACTION table with reduce m.

It means that – when a reduce is possible, we execute it *without* checking the next token.

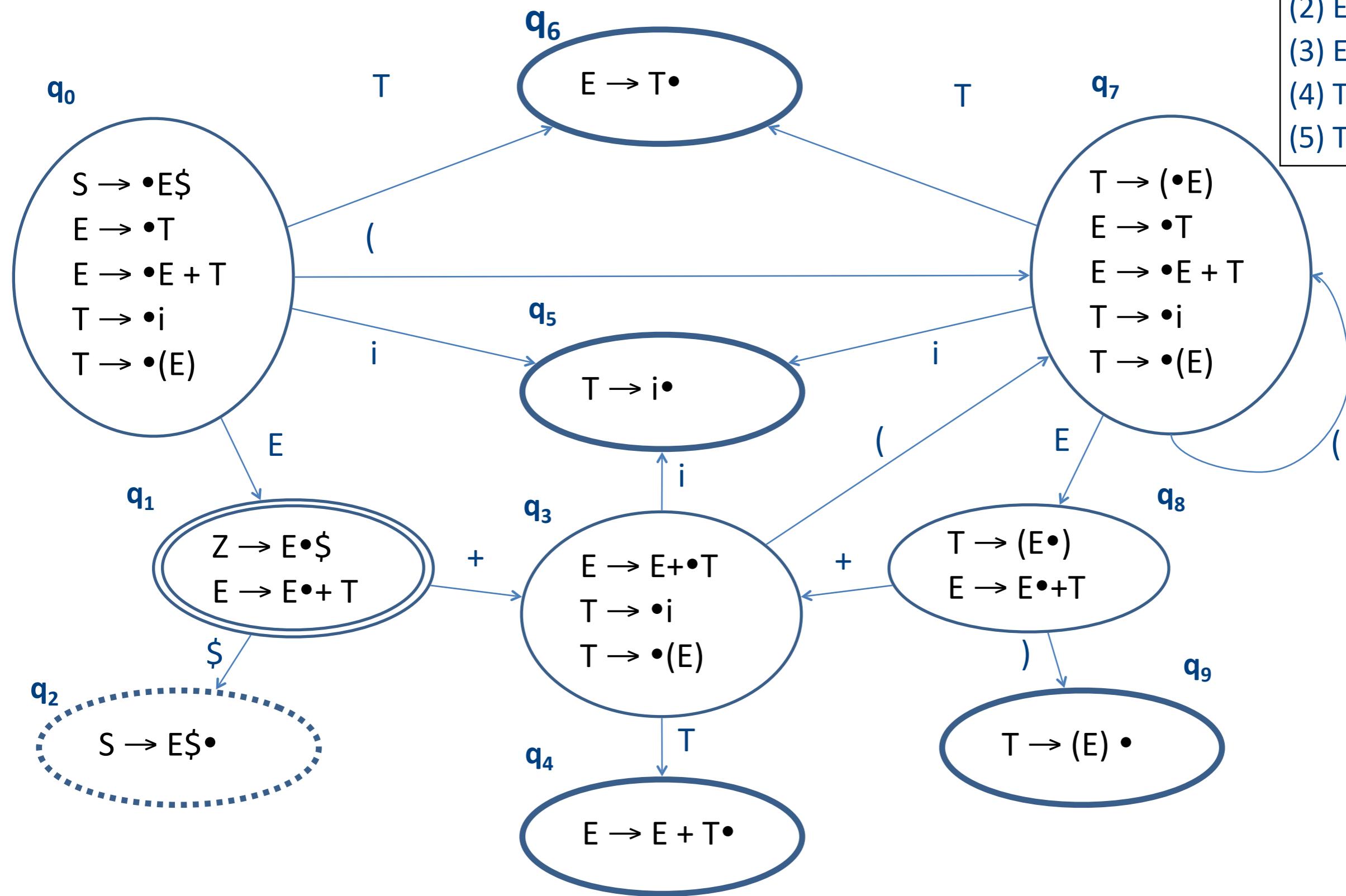
$k = 0$

The Parsing Algorithm, Formally

- Initialize the stack to q_0
- Repeat until halting:
 - ▶ Let q = top of stack, t = next token;
Consider ACTION[q, t] –
 - “shift i ”:
 - Remove t from the input; push t and q_i on the stack.
 - “reduce j ”, where rule (j) is $N \rightarrow \beta$:
 - Pop $|\beta|$ pairs from the stack;
let q' = the state at the top of the stack now.
 - Push N and the state GOTO[q', N] on the stack.
 - “accept”: halt successfully.
 - empty cell: halt with an error.

Example Run

- (1) $S \rightarrow E \$$
- (2) $E \rightarrow T$
- (3) $E \rightarrow E + T$
- (4) $T \rightarrow id$
- (5) $T \rightarrow (E)$



GOTO/ACTION Table

State	ACTION						GOTO	
	i	+	()	\$	E	T	
q ₀	s ₅		s ₇			g ₁	g ₆	
q ₁		s ₃			acc			
q ₃	s ₅		s ₇				g ₄	
q ₄	r ₃							
q ₅	r ₄							
q ₆	r ₂							
q ₇	s ₅		s ₇			g ₈	g ₆	
q ₈		s ₃		s ₉				
q ₉	r ₅							

- (1) Z → E \$
- (2) E → T
- (3) E → E + T
- (4) T → i
- (5) T → (E)

sn = shift to state n

rm = reduce using rule number (m)

Parsing id+id\$

Stack grows this way →

Stack	Input	Action
0	id + id \$	s5

Initialize with state 0

S	action						goto	
	id	+)	(\$	E	T	
0	s5		s7			g1	g6	
1		s3			acc			
2								
3	s5		s7				g4	
4	r3	r3	r3	r3	r3			
5	r4	r4	r4	r4	r4			
6	r2	r2	r2	r2	r2			
7	s5		s7			g8	g6	
8		s3		s9				
9	r5	r5	r5	r5	r5			



rn = reduce using rule number n
sm = shift to state m

Parsing id+id\$

Stack grows this way →

Stack	Input	Action
0	id + id \$	s5

Initialize with state 0

S	action						goto	
	id	+)	(\$	E	T	
0	s5			s7			g1	g6
1			s3			acc		
2								
3	s5			s7				g4
4	r3	r3	r3	r3	r3			
5	r4	r4	r4	r4	r4			
6	r2	r2	r2	r2	r2			
7	s5			s7			g8	g6
8			s3		s9			
9	r5	r5	r5	r5	r5			

rn = reduce using rule number n
sm = shift to state m

- | |
|---------------------------|
| (1) $S \rightarrow E \$$ |
| (2) $E \rightarrow T$ |
| (3) $E \rightarrow E + T$ |
| (4) $T \rightarrow id$ |
| (5) $T \rightarrow (E)$ |

Parsing $id + id \$$

Stack grows this way →

Stack	Input	Action
0	$id + id \$$	s5
0 id 5	$+ id \$$	r4

S	action						goto	
	id	+)	(\$	E	T	
0	s5		s7			g1	g6	
1		s3			acc			
2								
3	s5		s7				g4	
4	r3	r3	r3	r3	r3			
5	r4	r4	r4	r4	r4			
6	r2	r2	r2	r2	r2			
7	s5		s7			g8	g6	
8		s3		s9				
9	r5	r5	r5	r5	r5			

rn = reduce using rule number n
sm = shift to state m

Parsing id+id\$

Stack grows this way →

- (1) $S \rightarrow E \$$
- (2) $E \rightarrow T$
- (3) $E \rightarrow E + T$
- (4) $T \rightarrow id$
- (5) $T \rightarrow (E)$

Stack	Input	Action
0	id + id \$	s5
0 id 5	+ id \$	r4

S	action						goto	
	id	+)	(\$	E	T	
0	s5		s7			g1	g6	
1		s3			acc			
2								
3	s5		s7				g4	
4	r3	r3	r3	r3	r3			
5	r4	r4	r4	r4	r4			
6	r2	r2	r2	r2	r2			
7	s5		s7			g8	g6	
8		s3		s9				
9	r5	r5	r5	r5	r5			

pop id 5

rn = reduce using rule number n
sm = shift to state m

- (1) $S \rightarrow E \$$
- (2) $E \rightarrow T$
- (3) $E \rightarrow E + T$
- (4) $T \rightarrow id$
- (5) $T \rightarrow (E)$

Parsing $id + id \$$

Stack grows this way →

Stack	Input	Action
0	$id + id \$$	s5
0 id 5	$+ id \$$	r4

push T 6

S	action						goto	
	id	+)	(\$	E	T	
0	s5		s7			g1	g6	
1		s3				acc		
2								
3	s5		s7				g4	
4	r3	r3	r3	r3	r3			
5	r4	r4	r4	r4	r4			
6	r2	r2	r2	r2	r2			
7	s5		s7			g8	g6	
8		s3		s9				
9	r5	r5	r5	r5	r5			

rn = reduce using rule number n
sm = shift to state m

Parsing id+id\$

Stack grows this way →

Stack	Input	Action
0	id + id \$	s5
0 id 5	+ id \$	r4
0 T 6	+ id \$	r2

S	action						goto	
	id	+)	(\$	E	T	
0	s5		s7			g1	g6	
1		s3			acc			
2								
3	s5		s7				g4	
4	r3	r3	r3	r3	r3			
5	r4	r4	r4	r4	r4			
6	r2	r2	r2	r2	r2			
7	s5		s7			g8	g6	
8		s3		s9				
9	r5	r5	r5	r5	r5			

rn = reduce using rule number n
sm = shift to state m

- | |
|---------------------------|
| (1) $S \rightarrow E \$$ |
| (2) $E \rightarrow T$ |
| (3) $E \rightarrow E + T$ |
| (4) $T \rightarrow id$ |
| (5) $T \rightarrow (E)$ |

Parsing $id + id \$$

Stack grows this way →

Stack	Input	Action
0	$id + id \$$	s5
0 id 5	$+ id \$$	r4
0 T 6	$+ id \$$	r2
0 E 1	$+ id \$$	s3

S	action						goto	
	id	+)	(\$	E	T	
0	s5		s7			g1	g6	
1		s3			acc			
2								
3	s5		s7				g4	
4	r3	r3	r3	r3	r3			
5	r4	r4	r4	r4	r4			
6	r2	r2	r2	r2	r2			
7	s5		s7			g8	g6	
8		s3		s9				
9	r5	r5	r5	r5	r5			

rn = reduce using rule number n
sm = shift to state m

- | |
|---------------------------|
| (1) $S \rightarrow E \$$ |
| (2) $E \rightarrow T$ |
| (3) $E \rightarrow E + T$ |
| (4) $T \rightarrow id$ |
| (5) $T \rightarrow (E)$ |

Parsing $id + id \$$

Stack grows this way →

Stack	Input	Action
0	$id + id \$$	s5
0 id 5	$+ id \$$	r4
0 T 6	$+ id \$$	r2
0 E 1	$+ id \$$	s3
0 E 1 + 3	$id \$$	s5

S	action						goto	
	id	+)	(\$	E	T	
0	s5		s7			g1	g6	
1		s3			acc			
2								
3	s5		s7				g4	
4	r3	r3	r3	r3	r3			
5	r4	r4	r4	r4	r4			
6	r2	r2	r2	r2	r2			
7	s5		s7			g8	g6	
8		s3		s9				
9	r5	r5	r5	r5	r5			

rn = reduce using rule number n
sm = shift to state m

- | |
|---------------------------|
| (1) $S \rightarrow E \$$ |
| (2) $E \rightarrow T$ |
| (3) $E \rightarrow E + T$ |
| (4) $T \rightarrow id$ |
| (5) $T \rightarrow (E)$ |

Parsing $id + id \$$

Stack grows this way →

Stack	Input	Action
0	$id + id \$$	s5
0 id 5	$+ id \$$	r4
0 T 6	$+ id \$$	r2
0 E 1	$+ id \$$	s3
0 E 1 + 3	$id \$$	s5
0 E 1 + 3 id 5	$\$$	r4

S	action						goto	
	id	+)	(\$	E	T	
0	s5		s7			g1	g6	
1		s3			acc			
2								
3	s5		s7				g4	
4	r3	r3	r3	r3	r3			
5	r4	r4	r4	r4	r4			
6	r2	r2	r2	r2	r2			
7	s5		s7			g8	g6	
8		s3		s9				
9	r5	r5	r5	r5	r5			

rn = reduce using rule number n
sm = shift to state m

- (1) $S \rightarrow E \$$
- (2) $E \rightarrow T$
- (3) $E \rightarrow E + T$
- (4) $T \rightarrow id$
- (5) $T \rightarrow (E)$

Parsing $id + id \$$

Stack grows this way →

Stack	Input	Action
0	$id + id \$$	s5
0 id 5	$+ id \$$	r4
0 T 6	$+ id \$$	r2
0 E 1	$+ id \$$	s3
0 E 1 + 3	$id \$$	s5
0 E 1 + 3 id 5	$\$$	r4
0 E 1 + 3 T 4	$\$$	r3

S	action						goto	
	id	+)	(\$	E	T	
0	s5		s7			g1	g6	
1		s3			acc			
2								
3	s5		s7				g4	
4	r3	r3	r3	r3	r3	r3		
5	r4	r4	r4	r4	r4			
6	r2	r2	r2	r2	r2	r2		
7	s5		s7			g8	g6	
8		s3		s9				
9	r5	r5	r5	r5	r5	r5		

rn = reduce using rule number n
sm = shift to state m

Parsing id+id\$

Stack grows this way →

- (1) $S \rightarrow E \$$
- (2) $E \rightarrow T$
- (3) $E \rightarrow E + T$
- (4) $T \rightarrow id$
- (5) $T \rightarrow (E)$

Stack	Input	Action
0	id + id \$	s5
0 id 5	+ id \$	r4
0 T 6	+ id \$	r2
0 E 1	+ id \$	s3
0 E 1 + 3	id \$	s5
0 E 1 + 3 id 5	\$	r4
0 E 1 + 3 T 4	\$	r3
0 E 1	\$	acc

S	action						goto	
	id	+)	(\$	E	T	
0	s5		s7			g1	g6	
1		s3			acc			
2								
3	s5		s7				g4	
4	r3	r3	r3	r3	r3			
5	r4	r4	r4	r4	r4			
6	r2	r2	r2	r2	r2			
7	s5		s7			g8	g6	
8		s3		s9				
9	r5	r5	r5	r5	r5			

rn = reduce using rule number n
sm = shift to state m

Exercise: Parsing: $0^* 0 + 1$

	ACTION					GOTO	
	*	+	0	1	\$	E	B
q ₀			s ₁	s ₂		3	4
q ₁	r ₄						
q ₂	r ₅						
q ₃	s ₅	s ₆			acc		
q ₄	r ₃						
q ₅			s ₁	s ₂			7
q ₆			s ₁	s ₂			8
q ₇	r ₁						
q ₈	r ₂						

(0) $S \rightarrow E \$$

(1) $E \rightarrow E * B$

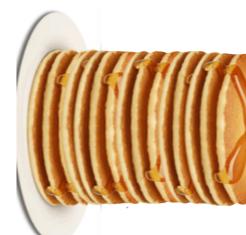
(2) $E \rightarrow E + B$

(3) $E \rightarrow B$

(4) $B \rightarrow 0$

(5) $B \rightarrow 1$

- Initialize the stack to q_0
- Repeat until halting:
 - ▶ Let $q = \text{top of stack}$, $t = \text{next token}$;
Consider $\text{ACTION}[q, t]$ —
 - “shift i ”:
 - Remove t from the input; push t and q_i on the stack.
 - “reduce j ”, where rule (j) is $N \rightarrow \beta$:
 - Pop $|\beta|$ pairs from the stack; let q' = the state at the top of the stack now.
 - Push N and the state $\text{GOTO}[q', N]$ on the stack.
 - “accept”: halt successfully.
 - empty cell: halt with an error.



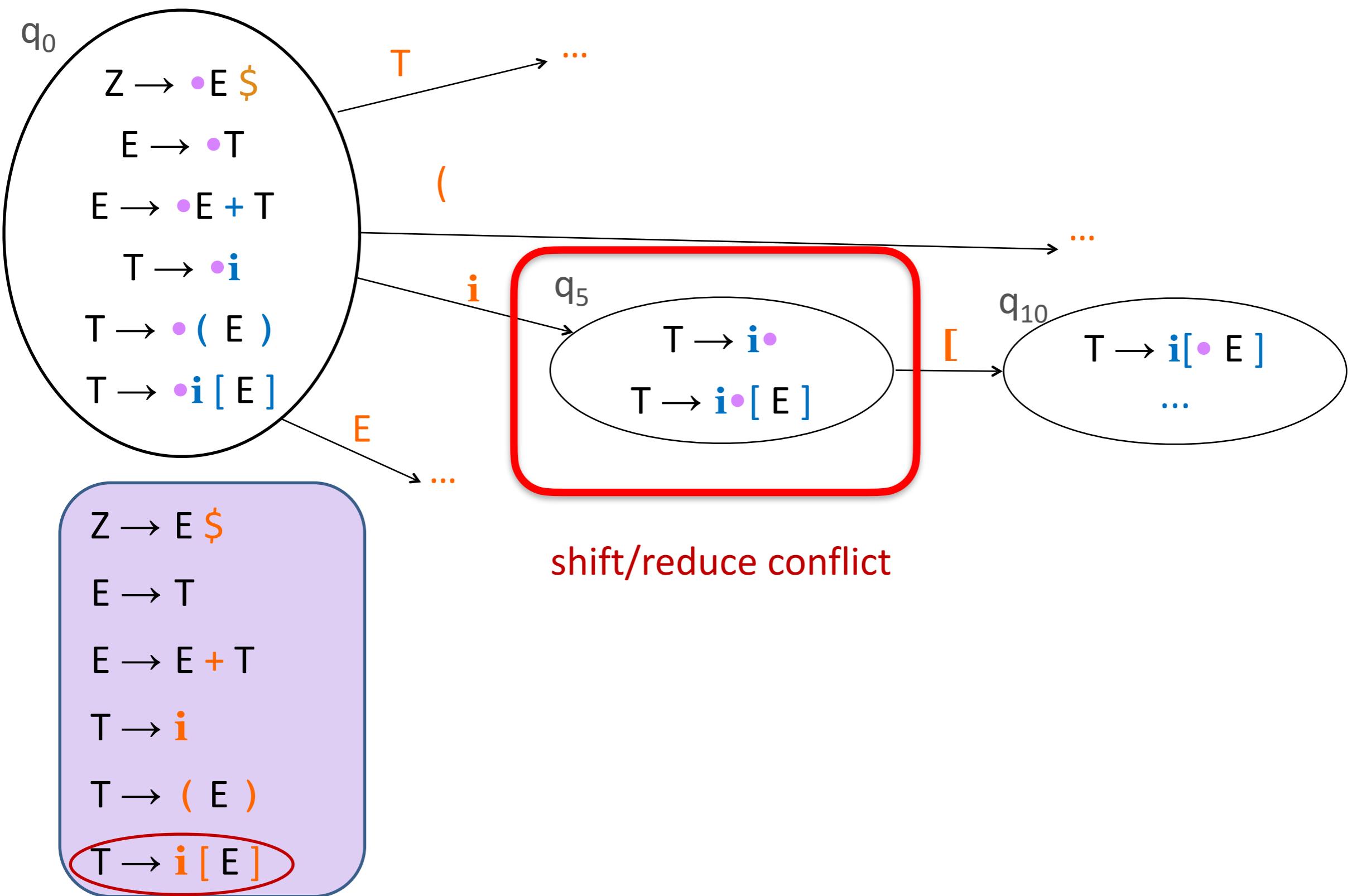
q₀

0 * 0 + 1 \$

Are we done?

- Can make a transition diagram for any grammar
- Can make a **GOTO** table for every grammar
- ...but the states are not always clear on what to do
⇒ **Cannot** make a *deterministic* ACTION table
for every grammar

LR(0) Conflicts

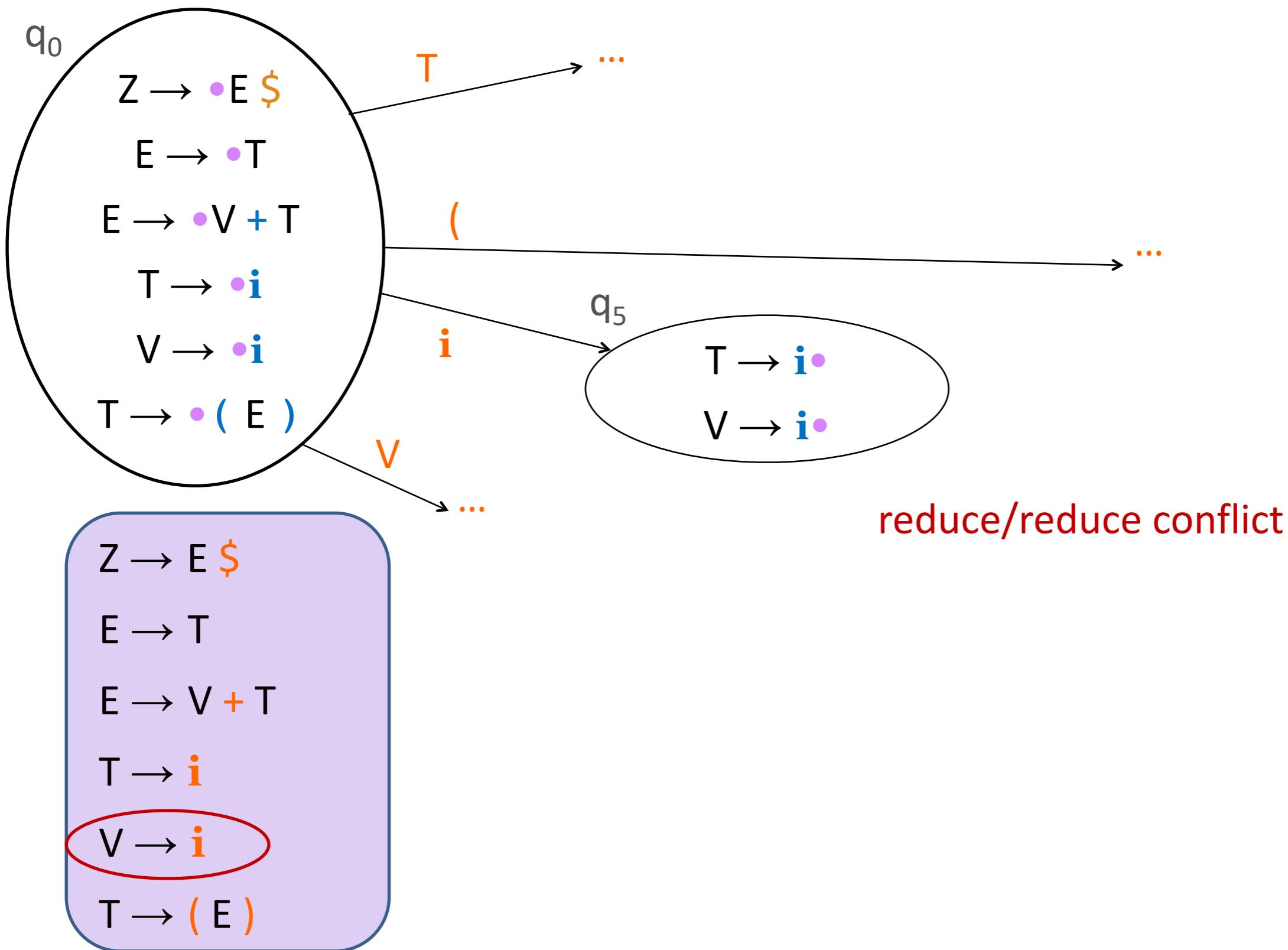


View in Action/Goto Table

- shift/reduce conflict...

	i	+	()	\$	[E	T
q0	s5		s7				1	6
q1		s3			acc			
q2	r1	r1	r1	r1	r1	r1		
q3	s5		s7					4
q4	r3	r3	r3	r3	r3	r3		
q5	r4	r4	r4	r4	r4	r4	r4/s10	
q6	r2	r2	r2	r2	r2	r2		
q7	s5		s7				8	6
q8		s3		s9				
q9	r5	r5	r5	r5	r5	r5		

LR(0) Conflicts



View in Action/Goto Table

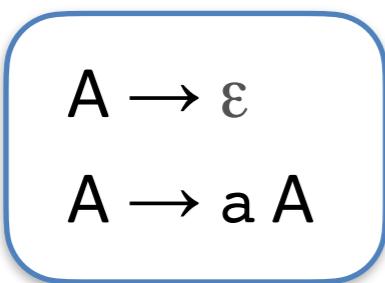
- reduce/reduce conflict...

	i	+	()	\$	E	T
q0	s5		s7			1	6
q1		s3			acc		
q2	r1	r1	r1	r1	r1		
q3	s5		s7				4
q4	r3	r3	r3	r3	r3		
q5	r4/r5	r4/r5	r4/r5	r4/r5	r4/r5		
q6	r2	r2	r2	r2	r2		
q7	s5		s7			8	6
q8		s3		s9			
q9	r5	r5	r5	r5	r5		

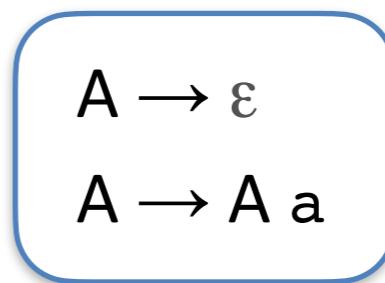
Can there be a shift/shift conflict?

LR(0) vs. ϵ -Rules

- Whenever a nonterminal has an ϵ production, it will be reduced as soon as it is reached in the grammar (remember, *without looking at the next token*).
- If the variable has another production with a terminal prefix, there is an inherent **shift/reduce** conflict



✗ Not good



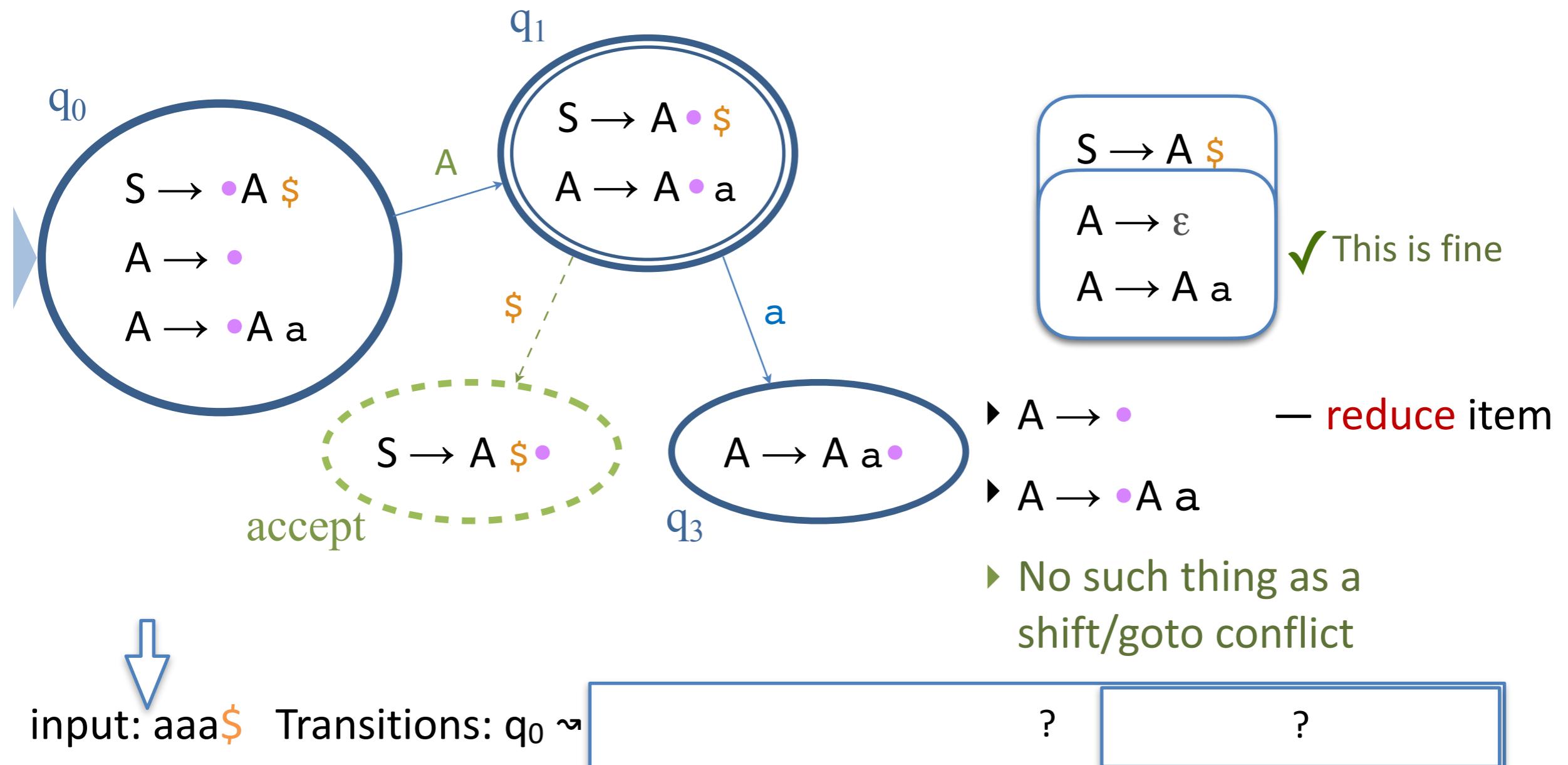
✓ This is fine

- ▶ $A \rightarrow \bullet$ — **reduce item**
- ▶ $A \rightarrow \bullet a A$ — **shift item**
- ▶ both are in the closure of any item of the form $\{P \rightarrow \alpha \bullet A \beta\}$

- ▶ $A \rightarrow \bullet$ — **reduce item**
- ▶ $A \rightarrow \bullet A a$
- ▶ No such thing as a **shift/goto conflict**

Good News: Left Recursion

- Left recursion is perfectly fine for LR(0) parsing and does not cause an infinite loop.



Can we do better?

Can we handle:

- Shift/reduce conflicts ?
- Reduce/reduce conflicts ?
- ϵ -rules: $A \rightarrow \epsilon \quad A \rightarrow aA$?

Important Bottom-Up LR-Parsers

- LR(0) – simplest, explains basic ideas
- SLR(1) – simple, explains lookahead
- LR(1) – complicated, very powerful, expensive
- LALR(1) – complicated, powerful enough, used by automatic tools

LL(k) grammars: k = number of tokens we lookahead when we do a reduce operation

Back to Action/Goto Table

- Remember? Reductions ignore the input...

	i	+	()	\$	E	T
q0	s4		s6			g1	g5
q1		s2			acc		
q2	s4		s6				g3
q3	r3	r3	r3	r3	r3		
q4	r4	r4	r4	r4	r4		
q5	r2	r2	r2	r2	r2		
q6	s4		s6			g7	g5
q7		s2		s8			
q8	r5	r5	r5	r5	r5		

SLR(1) Grammars

- aka Simple LR(1) or SLR
- A string should only be reduced to a nonterminal N if the look-ahead is a token that can follow N
 - ▶ A reduce item $N \rightarrow \alpha^\bullet$ is applicable only when the lookahead is in $\text{FOLLOW}(N)$

$$\text{FOLLOW}(N) = \{t \in T \mid S \Rightarrow^* \beta N t \gamma\} \cup \{\$\mid S \Rightarrow^* \beta N\}$$

- Differs from LR(0) only on the original **reduce** rows
- Allows us to sometimes not reduce, instead **shift** (or do nothing – detect error sooner)

SLR Grammars and FOLLOW

$$\begin{array}{l} E \rightarrow E * B \mid E + B \mid B \\ \text{(1)} \quad \text{(2)} \quad \text{(3)} \\ B \rightarrow 0 \mid 1 \\ \text{(4)} \quad \text{(5)} \end{array}$$

Derivation (in reverse):

0 + 0 * 1

B + 0 * 1

'+' ∈ FOLLOW(B)

E + 0 * 1

E + B * 1

'*' ∈ FOLLOW(B)

E * 1

E * B

'*' ∈ FOLLOW(E)

E

\$ ∈ FOLLOW(B)

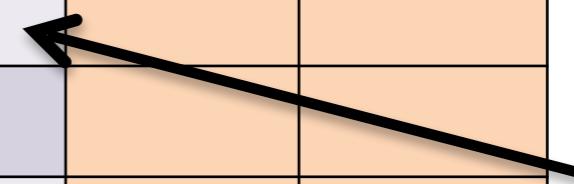
When reducing a sub-string to some nonterminal **N**, while the lookahead is a token **a** → “**Na**” is going to appear in the **sentential form**

GOTO/ACTION Table

	i	+	()	\$	E	T
q0	s4		s6			g1	g5
q1		s2			acc		
q2	s4		s6				g3
q3		r3		r3	r3		
q4	r4	r4	r4	r4	r4		
q5	r2	r2	r2	r2	r2		
q6	s4		s6			g7	g5
q7		s2			s8		
q8	r5	r5	r5	r5	r5		

- (1) Z → E \$
- (2) E → T
- (3) E → E + T
- (4) T → i
- (5) T → (E)

The tokens that can follow E are '+' ')' and '\$'.



$\text{FOLLOW}(E) = \{+,), \$\}$

GOTO/ACTION Table

	i	+	()	\$	E	T
q0	s4		s6			g1	g5
q1		s2			acc		
q2	s4		s6				g3
q3		r3		r3	r3		
q4	r4	r4	r4	r4	r4		
q5		r2		r2	r2		
q6	s4		s6			g7	g5
q7		s2		s8			
q8	r5	r5	r5	r5	r5		

- (1) Z → E \$
- (2) E → T
- (3) E → E + T
- (4) T → i
- (5) T → (E)

The tokens that can follow E are '+' ')' and '\$'.

$\text{FOLLOW}(E) = \{+,), \$\}$

GOTO/ACTION Table

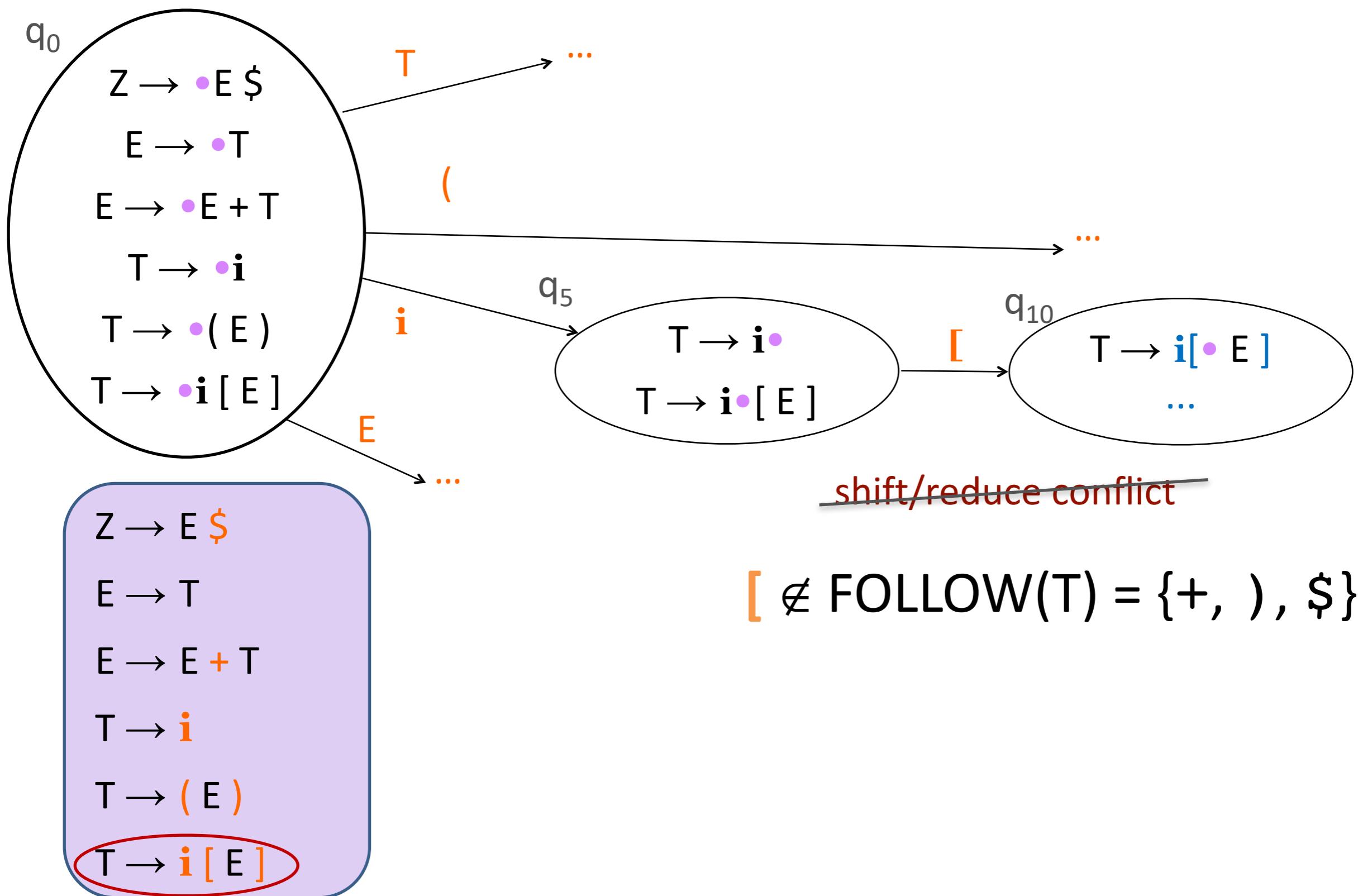
	i	+	()	\$	E	T
q0	s4		s6			g1	g5
q1		s2			acc		
q2	s4		s6				g3
q3		r3		r3	r3		
q4		r4		r4	r4		
q5		r2		r2	r2		
q6	s4		s6			g7	g5
q7		s2		s8			
q8		r5		r5	r5		

- (1) $Z \rightarrow E \$$
- (2) $E \rightarrow T$
- (3) $E \rightarrow E + T$
- (4) $T \rightarrow i$
- (5) $T \rightarrow (E)$

Same for T

$\text{FOLLOW}(T) = \{+,), \$\}$

Now let's add “ $T \rightarrow i [E]$ ”



Now let's add “T → i [E]”

	i	+	()	[]	\$	E	T
q ₀	s ₄		s ₆					g ₁	g ₅
q ₁		s ₂					acc		
q ₂	s ₄		s ₆						g ₃
q ₃		r ₃		r ₃			r ₃		
q ₄		r ₄		r ₄	s ₉	yay 😊	r ₄		
q ₅		r ₂		r ₂			r ₂		
q ₆	s ₄		s ₆					g ₇	g ₅
q ₇		s ₂		s ₈					
q ₈		r ₅		r ₅			r ₅		
q ₉	s ₄		s ₆					g ₁₀	
:									

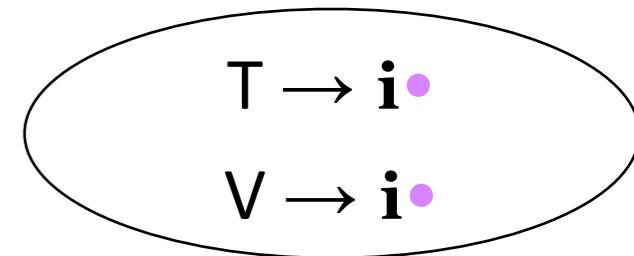
- (1) Z → E \$
- (2) E → T
- (3) E → E + T
- (4) T → i
- (5) T → (E)
- (6) T → i [E]

FOLLOW(T) = {+,), \$}

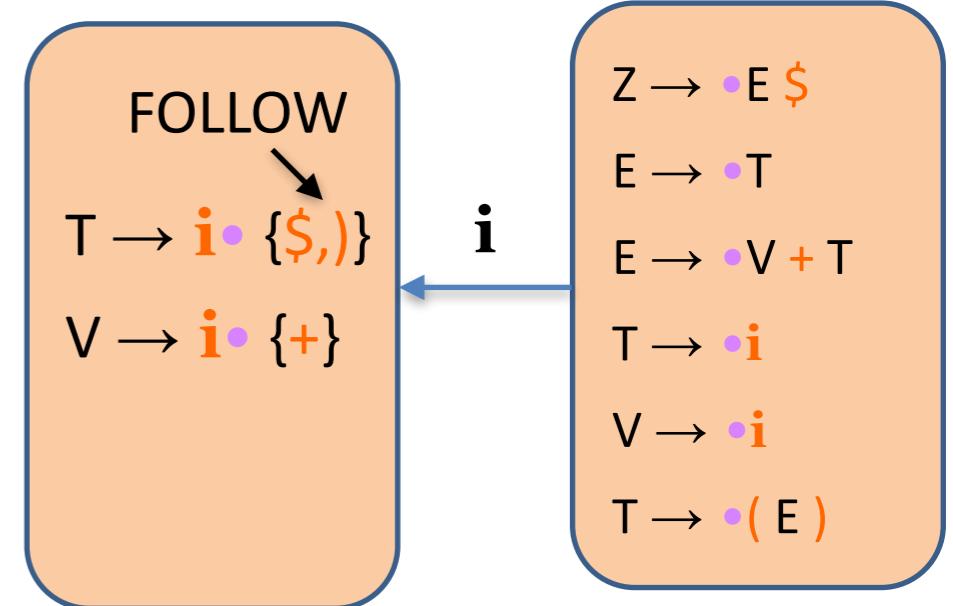
SLR(1): check next token when reducing

- Previous example demonstrates the elimination of a **shift/reduce** conflict
- Can eliminate **reduce/reduce** conflicts when conflicting rules' left-hand sides satisfy:

$$\text{FOLLOW}(T) \cap \text{FOLLOW}(V) = \emptyset$$



- But cannot resolve all shift/reduce and reduce/reduce conflicts!

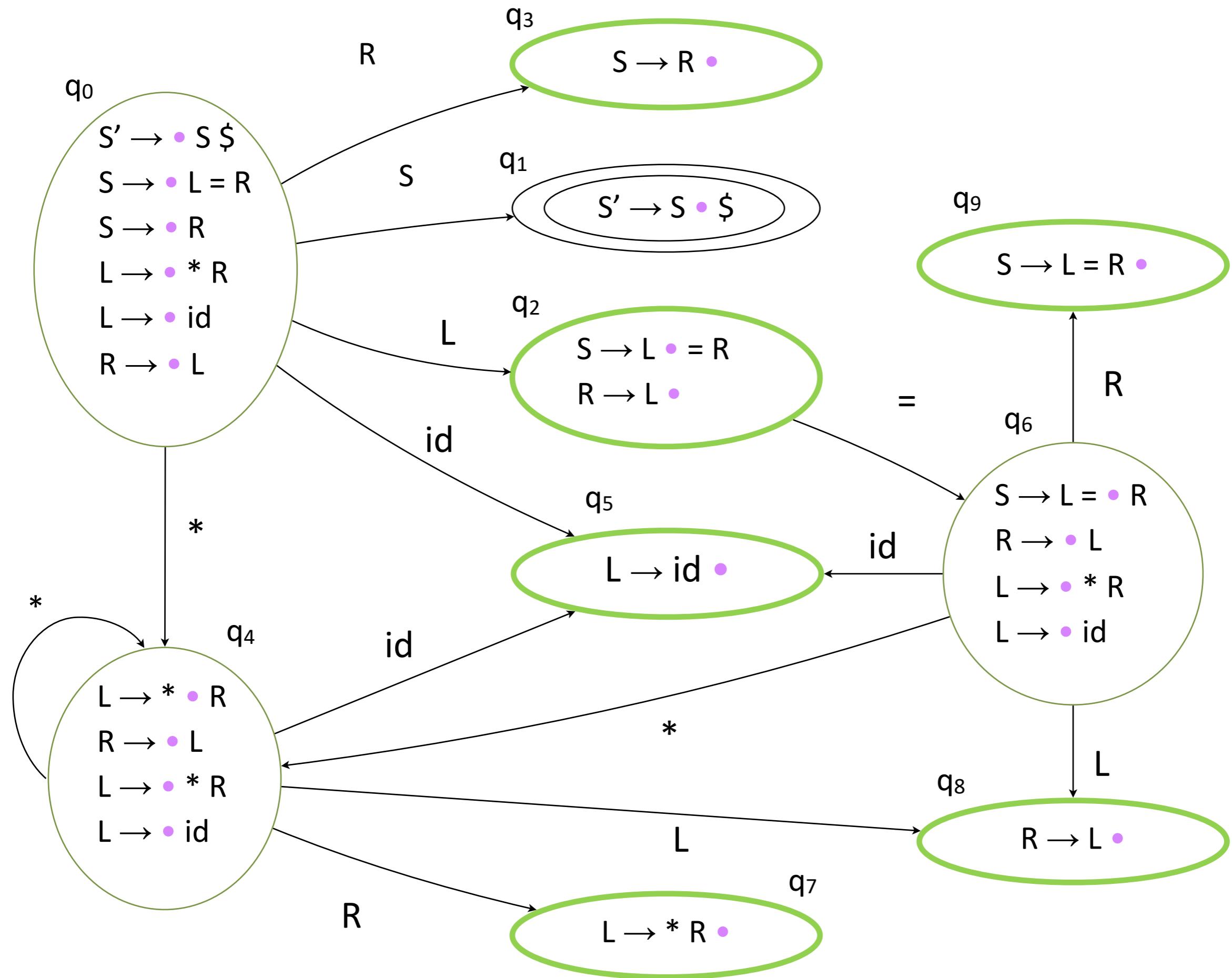


Consider this non-LR(0) grammar

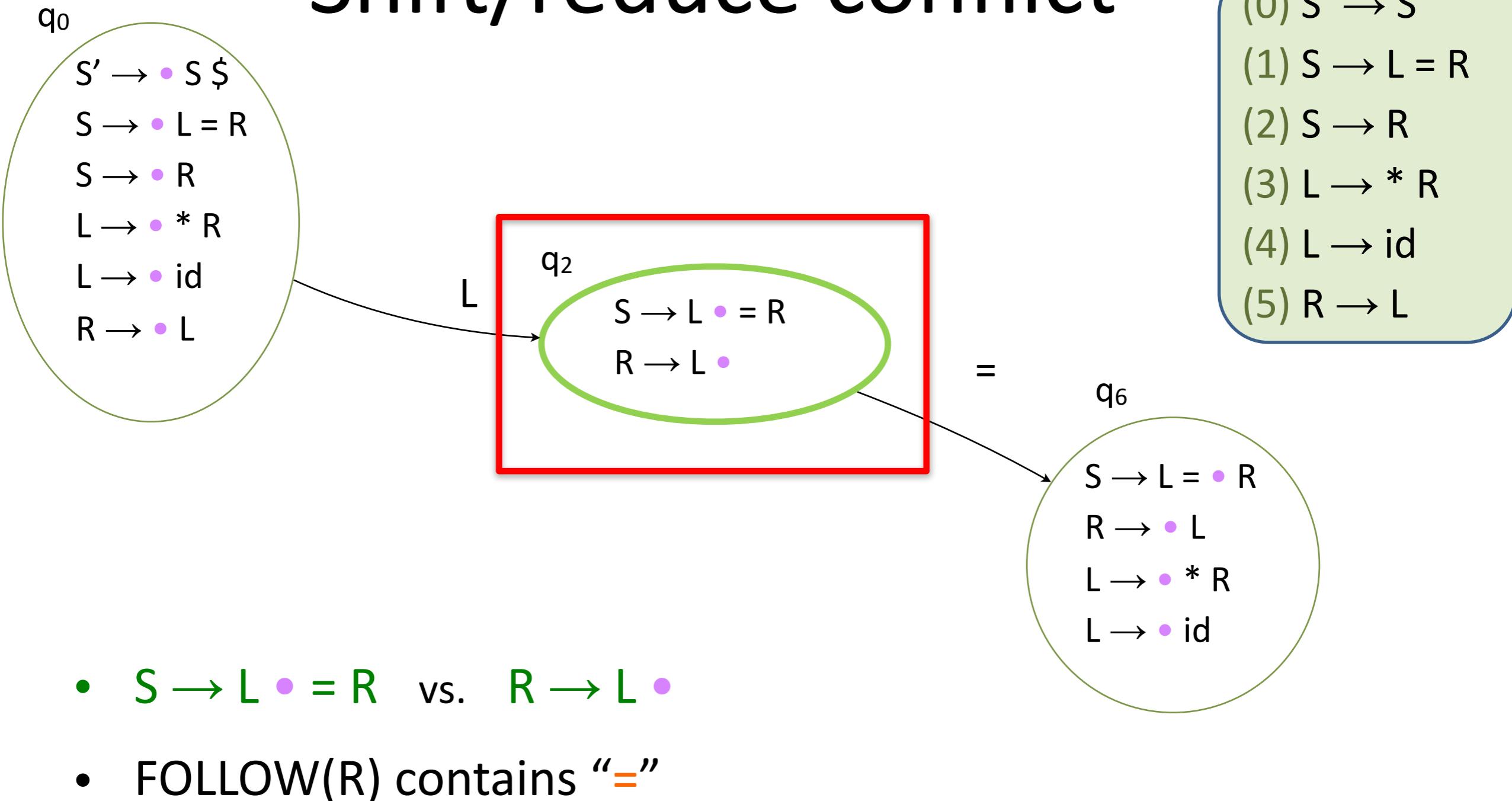
- (0) $S' \rightarrow S \$$
- (1) $S \rightarrow L = R$
- (2) $S \rightarrow R$
- (3) $L \rightarrow * R$
- (4) $L \rightarrow id$
- (5) $R \rightarrow L$

id
***id
id = id
*id = id
id = *id
***id = **id

...



Shift/reduce conflict



- $S \rightarrow L \cdot = R$ vs. $R \rightarrow L \cdot$
- FOLLOW(R) contains “=”

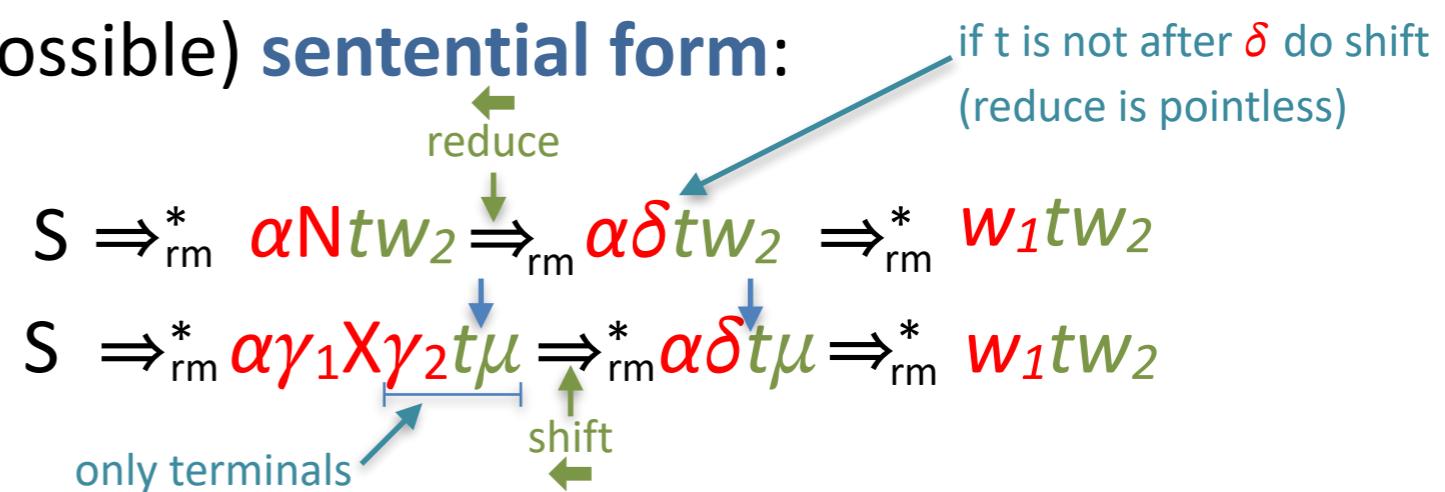
$$S \Rightarrow L = R \Rightarrow * \overbrace{R = R}^{\text{FOLLOW}(R)}$$

⇒ SLR **cannot resolve** the conflict in this case

Resolving the Conflict

- In SLR(1): a reduce item $N \rightarrow \delta \bullet$ is applicable when the lookahead is in $\text{FOLLOW}(N)$
- But when we reach a **reduce state**, there is a whole sequence of **reduction steps** that generate a prefix of (possible) **sentential form**:

input = $w_1 t w_2 \xrightarrow{} \alpha \delta \xrightarrow{?} \alpha N$*



- Given αN , we can ask what the next token of N may be in any sentential form whose prefix is αN :

$$\text{FOLLOW}(\alpha N) = \{ t \mid S \Rightarrow_{\text{rm}}^* \alpha N t \beta \}$$

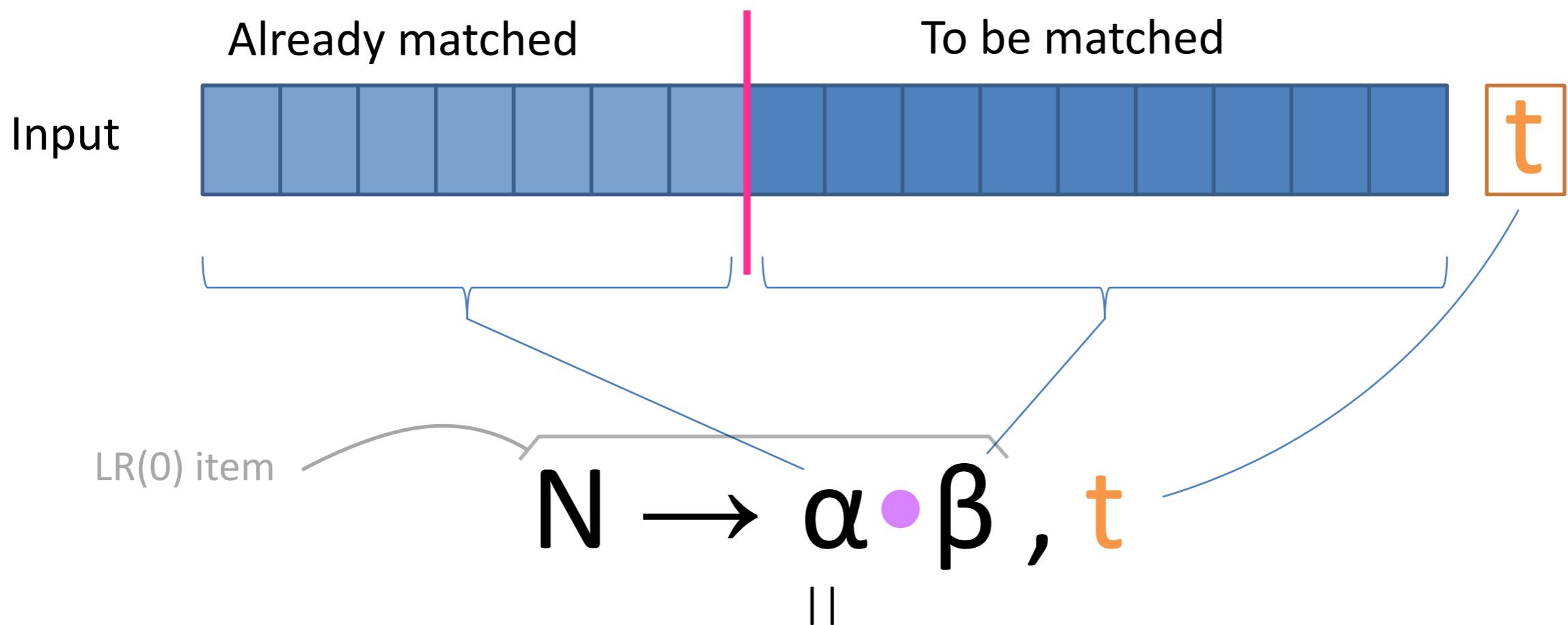
- Looking at $\text{FOLLOW}(\alpha N)$ is more restrictive than looking at the FOLLOW of the last variable N .

$$\text{FOLLOW}(\alpha N) \subseteq \text{FOLLOW}(N) = \{ t \in \text{FIRST}(\beta) \mid S \Rightarrow_{\text{rm}}^* \gamma N \beta \}$$

- ▶ In a way, $\text{FOLLOW}(N)$ merges look-ahead for all possible occurrences of N in a sentential form

LR(1) Item

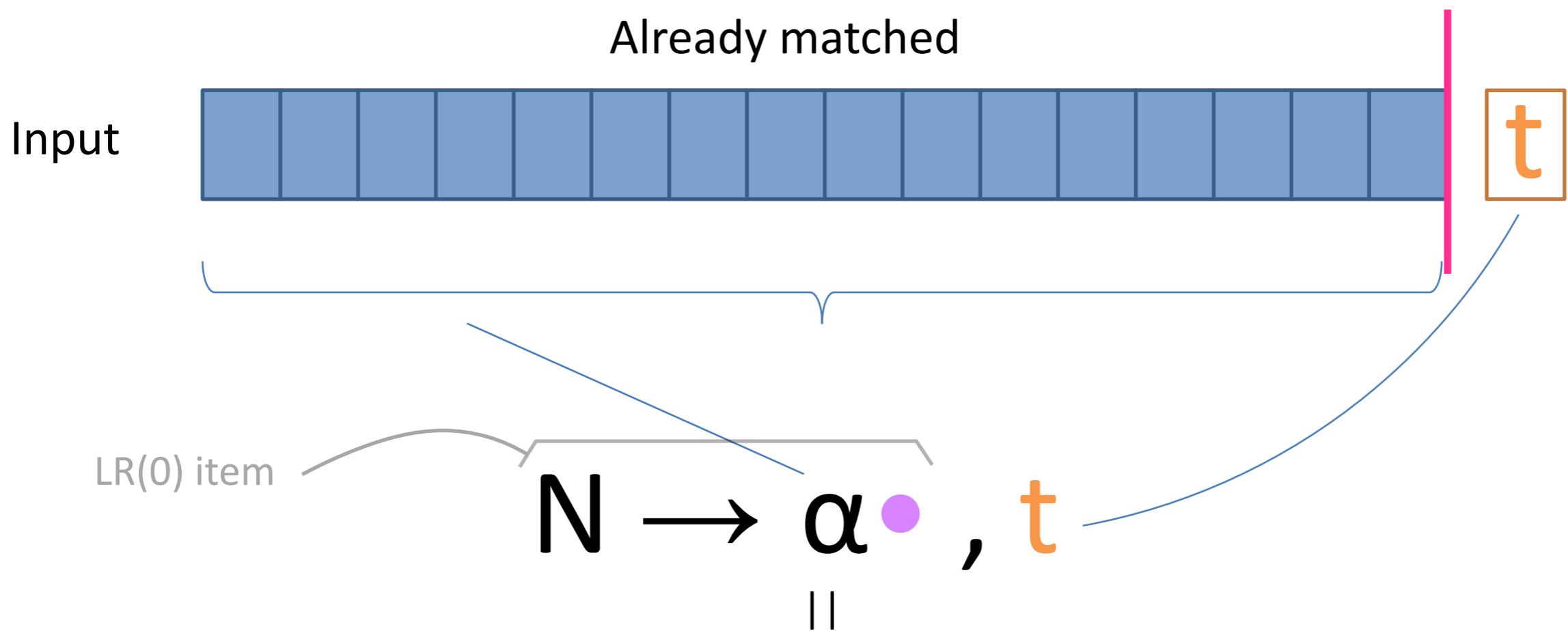
LR(1) tracks FOLLOW(αN) by keeping a look-ahead with each LR(1) item



So far we've matched α , expecting to see β ,
followed by the lookahead t

LR(1) Item

LR(1) tracks FOLLOW(αN) by keeping a look-ahead with each LR(1) item



So far we've matched α , if the next token is the lookahead t ,
we should reduce α to N

LR(1) Item

- Example: the production $L \rightarrow id$ yields the following LR(1) items

(0) $S' \rightarrow S$
(1) $S \rightarrow L = R$
(2) $S \rightarrow R$
(3) $L \rightarrow * R$
(4) $L \rightarrow id$
(5) $R \rightarrow L$

$L \rightarrow \bullet id$
 $L \rightarrow id \bullet$

$L \rightarrow \bullet id, *$
 $L \rightarrow \bullet id, =$
 $L \rightarrow \bullet id, id$
 $L \rightarrow \bullet id, \$$
 $L \rightarrow id \bullet, *$
 $L \rightarrow id \bullet, =$
 $L \rightarrow id \bullet, id$
 $L \rightarrow id \bullet, \$$

will never be generated, since $*, id \notin FOLLOW(L)$
(by construction)

SLR
Reduce only when next token is in $FOLLOW(L) = \{ =, \$ \}$

LR(1)
Refines FOLLOW by using the RHS of the rules

Creating the states for LR(1)

- We start with the initial state:
 - ▶ q_0 will be the closure of: $\{S' \rightarrow \bullet S, \$\}$

- Closure set for LR(1):

If the set contains an item of the form

$$A \rightarrow \alpha \bullet B \beta, c$$

then it must *also* contain an item

$$B \rightarrow \bullet \delta, d$$

for every production $B \rightarrow \delta$ and every token $d \in \text{FIRST}(\beta c)$

$$\begin{aligned}\text{FIRST}(\beta c) &= \{t \in T \mid \beta c \Rightarrow^* t\gamma\} \\ &= \text{FIRST}(\beta) \setminus \{\epsilon\} \cup \{c \mid \beta \Rightarrow^* \epsilon\}\end{aligned}$$

Closure of $\{S' \rightarrow \bullet S, \$\}$

No need to have
\$ in the grammar
explicitly

- We would like to add rules that start with S, but keep track of possible lookahead.

- ▶ $S' \rightarrow \bullet S, \$$
- ▶ $S \rightarrow \bullet L = R, \$$
- ▶ $S \rightarrow \bullet R, \$$
- ▶ $L \rightarrow \bullet * R, =$
- ▶ $L \rightarrow \bullet id, =$
- ▶ $R \rightarrow \bullet L, \$$
- ▶ $L \rightarrow \bullet * R, \$$
- ▶ $L \rightarrow \bullet id, \$$

– Rules for S

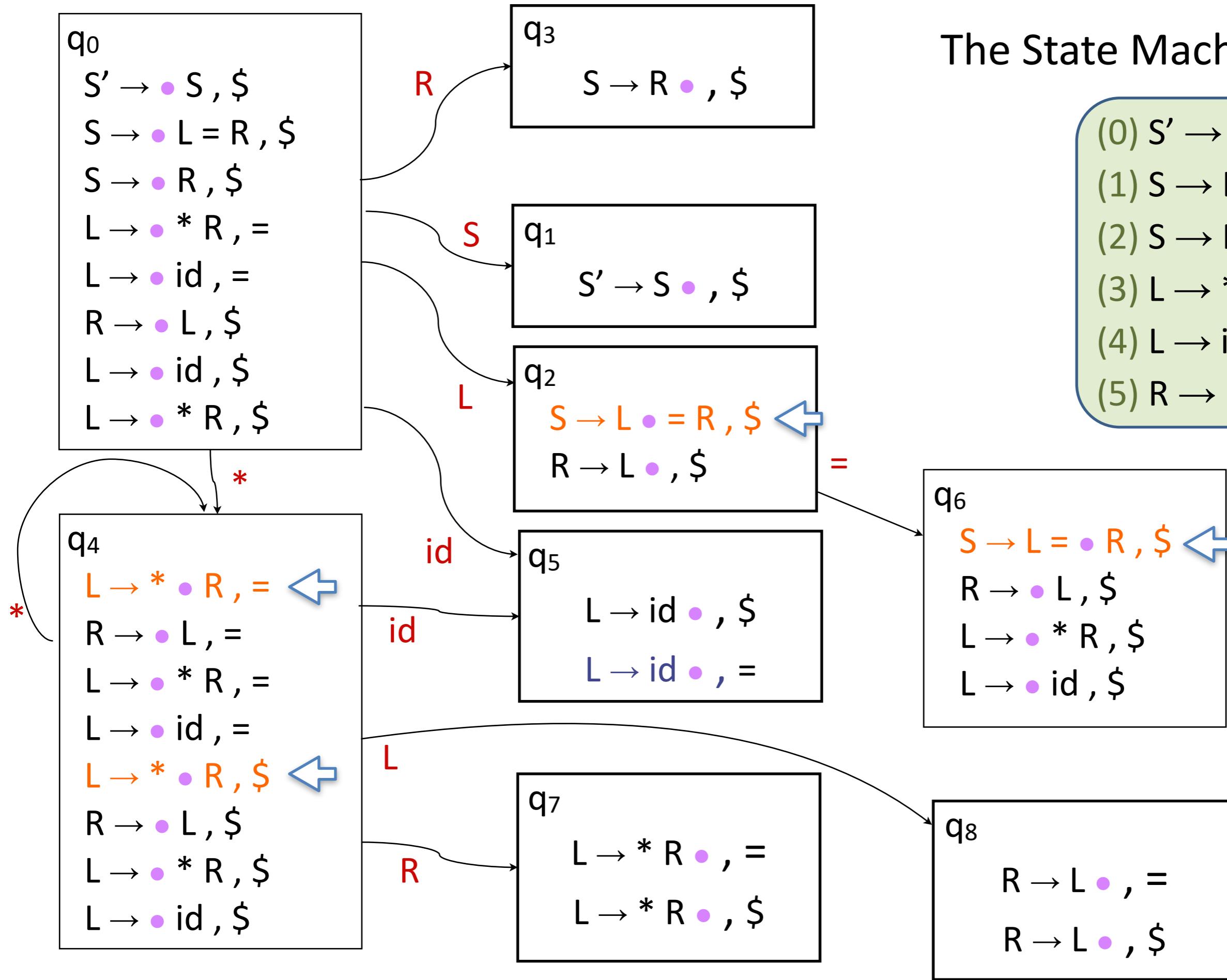
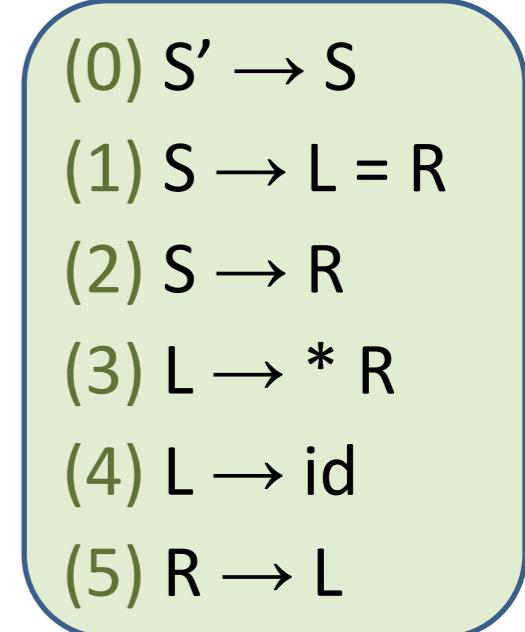
– Rules for L

– Rules for R

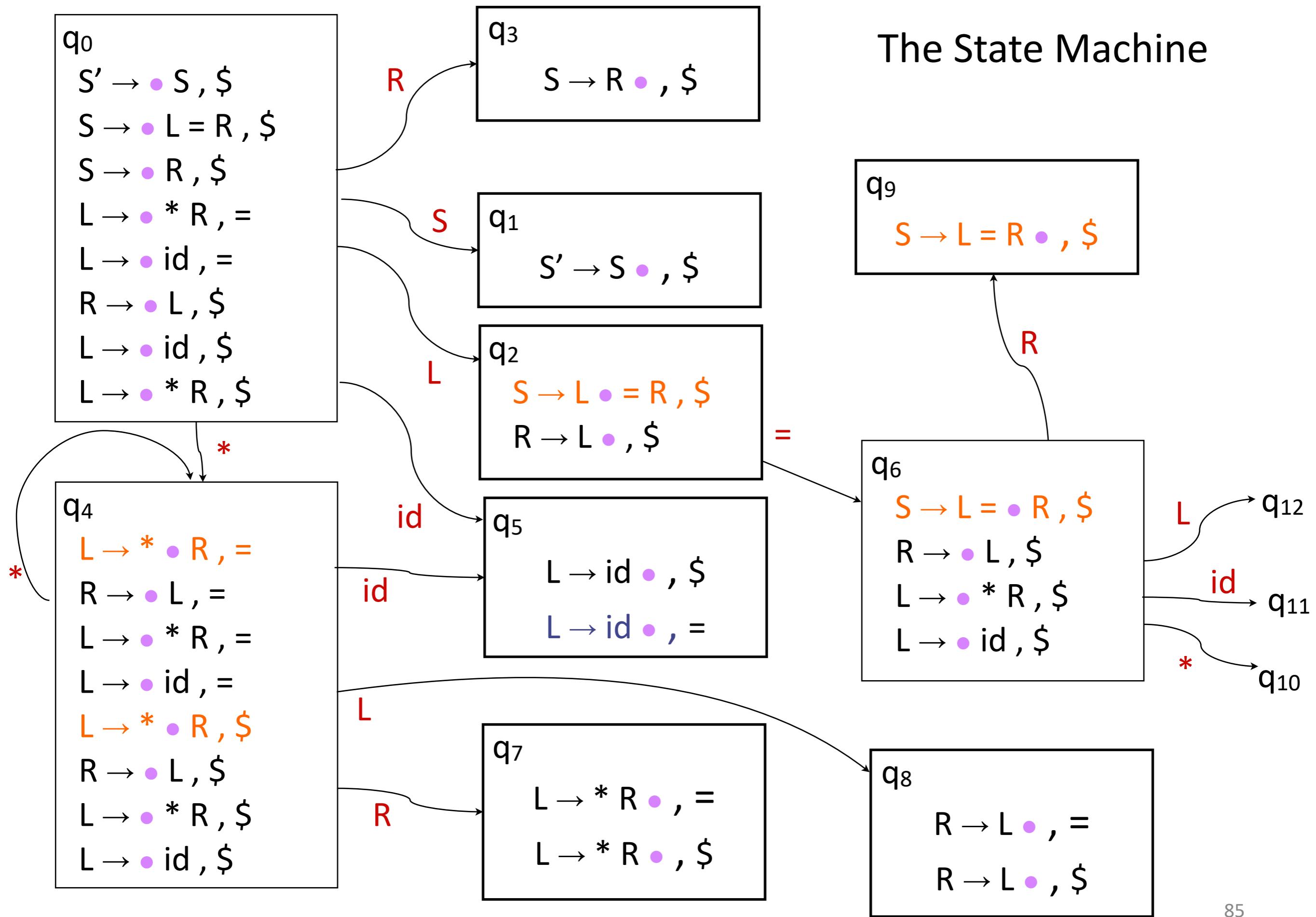
– More rules for L

- (0) $S' \rightarrow S$
- (1) $S \rightarrow L = R$
- (2) $S \rightarrow R$
- (3) $L \rightarrow * R$
- (4) $L \rightarrow id$
- (5) $R \rightarrow L$

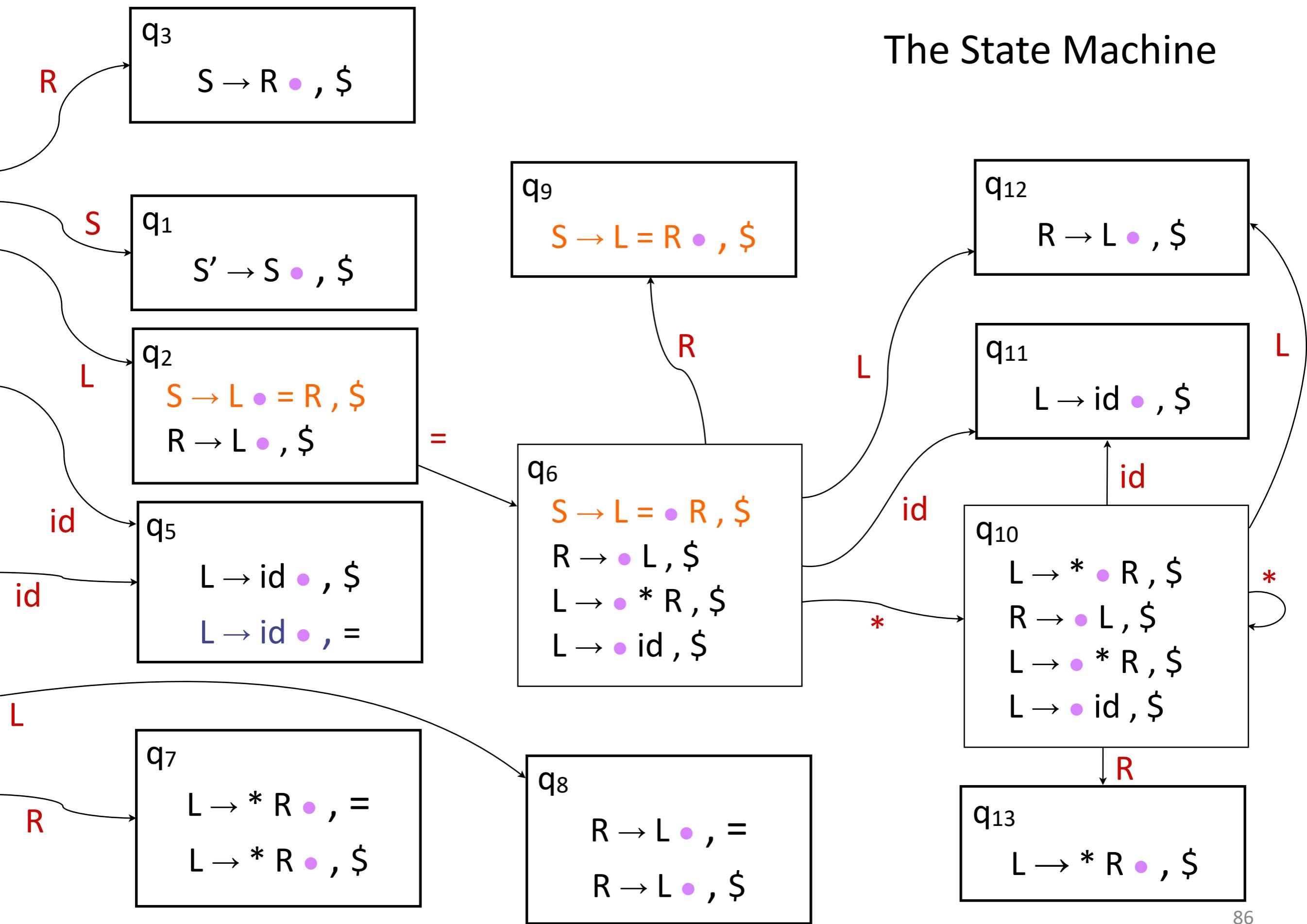
The State Machine



The State Machine



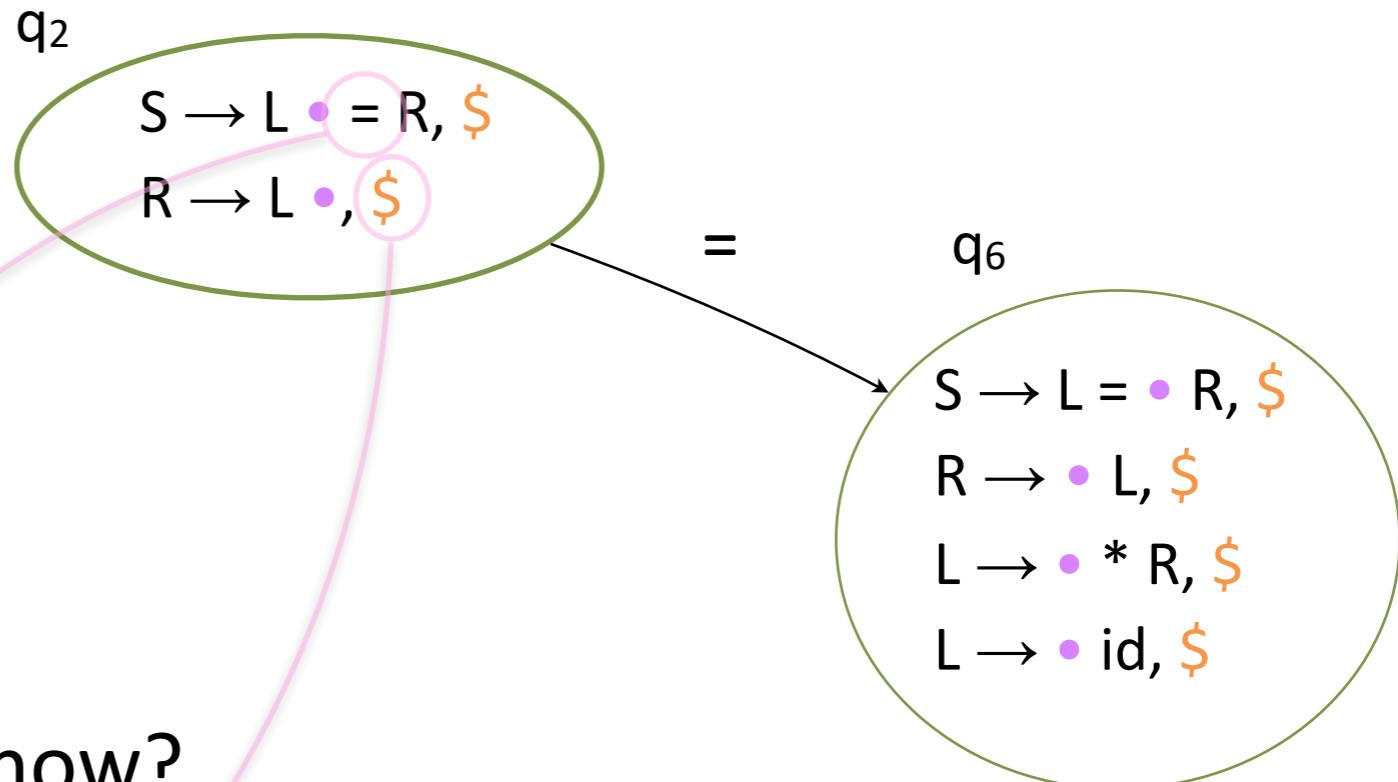
The State Machine



Back to the conflict

(0) $S' \rightarrow S$
 (1) $S \rightarrow L = R$
 (2) $S \rightarrow R$
 (3) $L \rightarrow * R$
 (4) $L \rightarrow id$
 (5) $R \rightarrow L$

FOLLOW(R) contains “=”



- Is there a conflict now?

State	ACTION				GOTO	
	*	=	id	\$	E	T
q_2		s_6		r5		

Building the Tables

- Similarly to LR(0) and SLR, we start with the automaton.
- Turn each transition in a token column to a shift.
- The variables' columns form the GOTO section.
- The “acc” is put in the \$ column, for any state that contains $(S' \rightarrow S^\bullet, \$)$.
- For any state that contains an item of the form $(A \rightarrow \beta^\bullet, a)$, where $A \rightarrow \beta$ is rule number (m), use “**reduce m** ” for the row of this state and the column of token a .

Building the LR(1) Table

	ACTION				GOTO		
	id	*	=	\$	S	R	L
q ₀	s ₅	s ₄			g ₁	g ₃	g ₂
q ₁				acc			
q ₂			s ₆	r ₅			
q ₃				r ₂			
q ₄	s ₅	s ₄				g ₇	g ₈
q ₅			r ₄	r ₄			
q ₆	s ₁₁	s ₁₀				g ₉	g ₁₂
q ₇			r ₃	r ₃			
q ₈			r ₅	r ₅			
q ₉				r ₁			
q ₁₀	s ₁₁	s ₁₀				g ₁₃	g ₁₂
q ₁₁				r ₄			
q ₁₂				r ₅			
q ₁₃				r ₃			

- (0) $S' \rightarrow S$
- (1) $S \rightarrow L = R$
- (2) $S \rightarrow R$
- (3) $L \rightarrow * R$
- (4) $L \rightarrow id$
- (5) $R \rightarrow L$

q₂
 $S \rightarrow L \bullet = R, \$$
 $R \rightarrow L \bullet, \$$

FOLLOW(R) =
{=, \$}

Bottom Up Parsing

*LR(k)

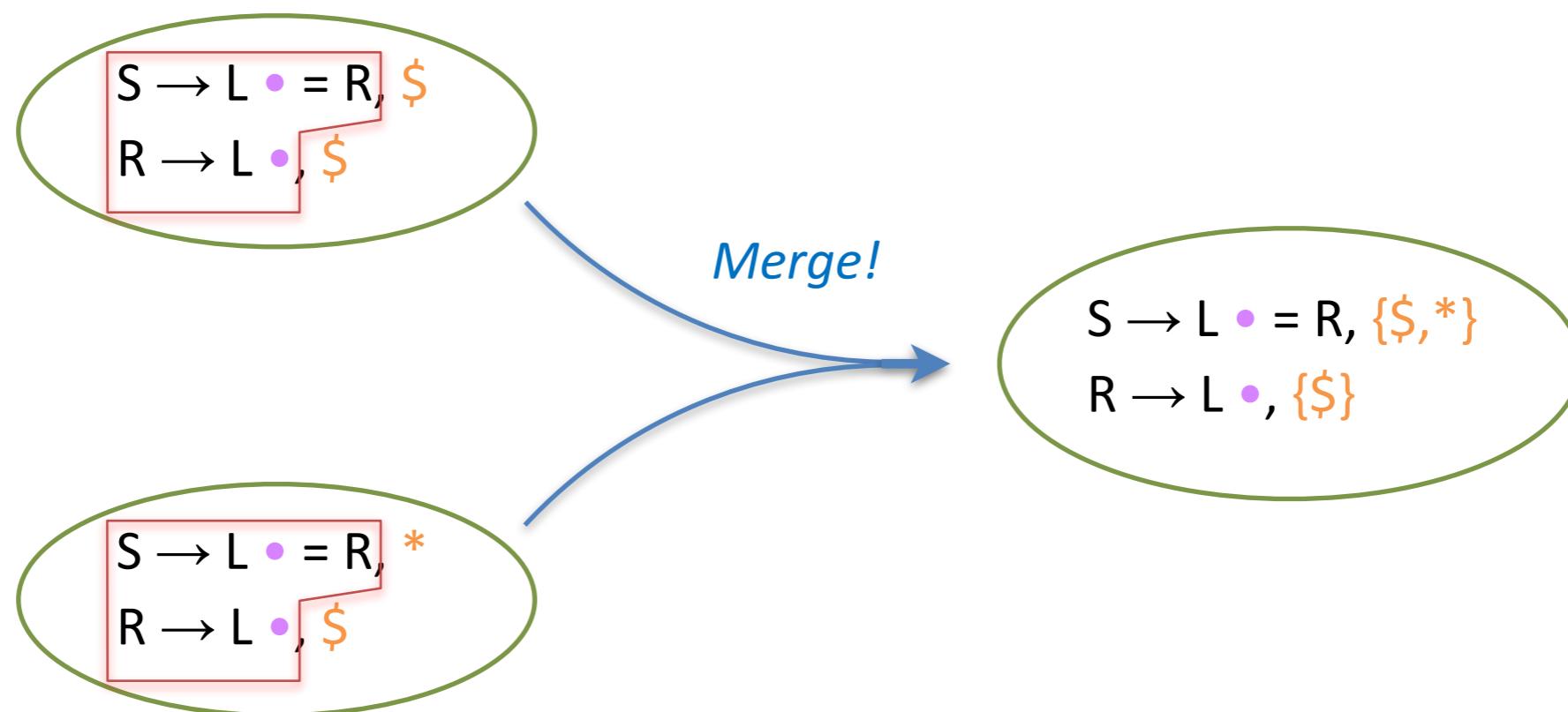
*SLR

- All follow the same table+stack algorithm
- Differ on type of “LR Items”, translation to table

LR(0)	SLR(1)	LR(1)
$N \rightarrow \alpha \bullet \beta$	$N \rightarrow \alpha \bullet \beta$	$N \rightarrow \alpha \bullet \beta, \sigma$
	$N \rightarrow \alpha \bullet \{ FOLLOW(N) \}$	

LALR

- Goal. LR(1) creates a lot of states that look similar. Can we reduce the number of states?
- Idea. Merge states that have exactly the same set of LR(0) items (**core**), representing the lookaheads as *sets*.



* Of course, we may get *more* conflicts as a result.

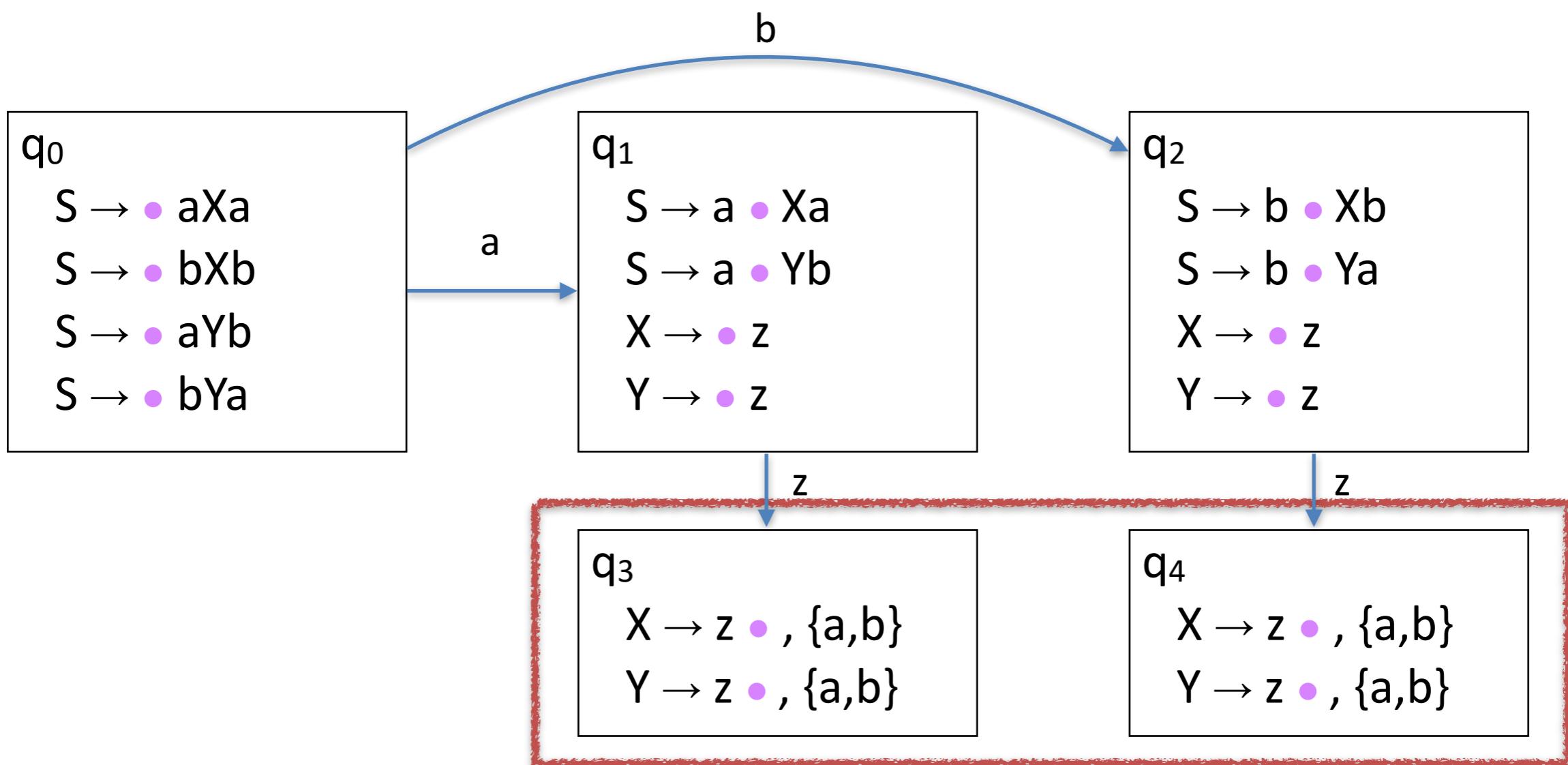
Example: SLR(1) Parsing

$$S \rightarrow aXa \mid bXb \mid aYb \mid bYa$$

$$X \rightarrow z$$

$$Y \rightarrow z$$

SLR(1) Automaton (partial)



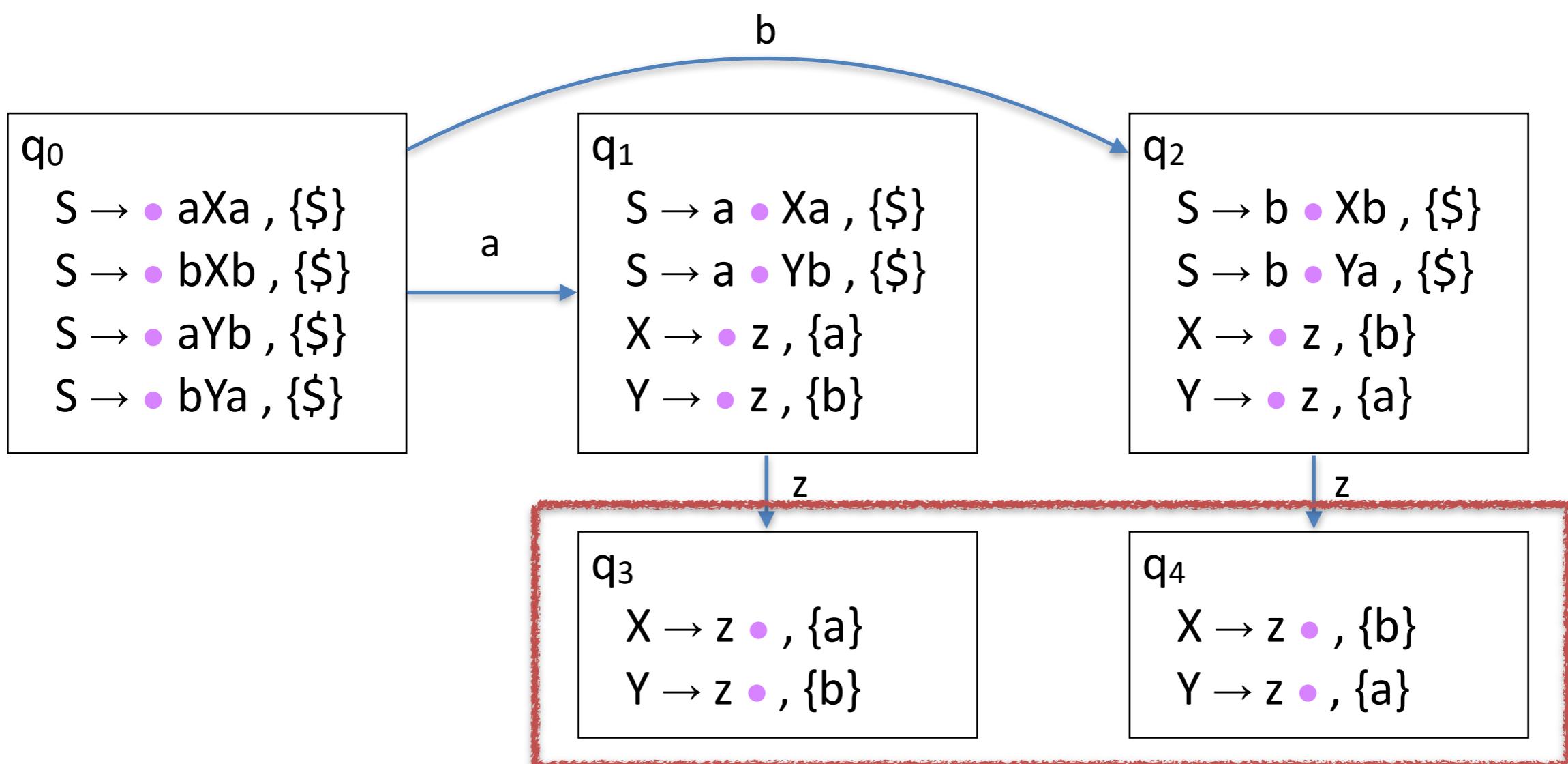
Example: LR(1) Parsing

$$S \rightarrow aXa \mid bXb \mid aYb \mid bYa$$

$$X \rightarrow z$$

$$Y \rightarrow z$$

LR(1) Automaton (partial)



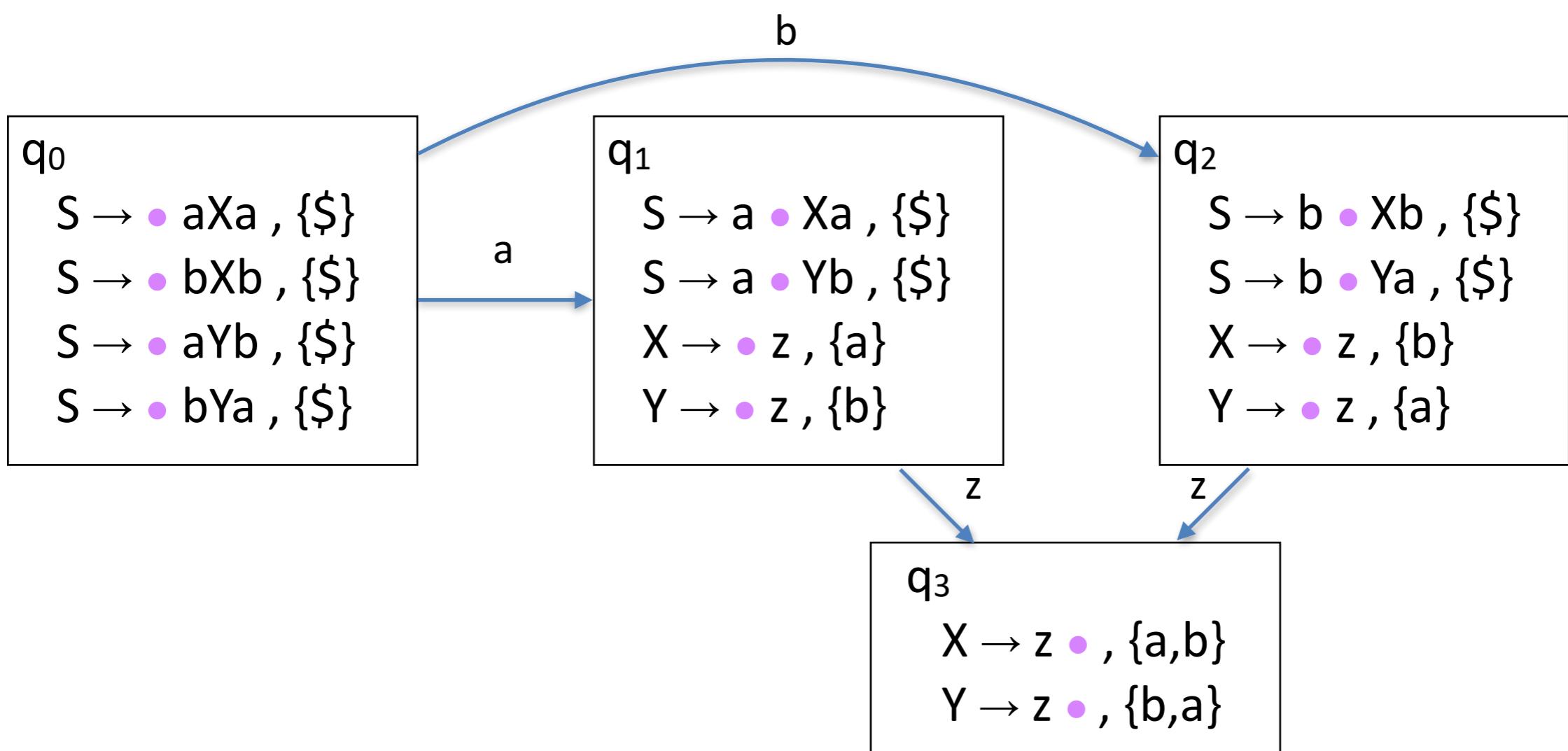
LALR Parsing

$$S \rightarrow aXa \mid bXb \mid aYb \mid bYa$$

$$X \rightarrow z$$

$$Y \rightarrow z$$

LR(1) Automaton (partial)



Example: SLR(1) Parsing

$$S \rightarrow A \mid xb$$

$$A \rightarrow aAb \mid B$$

$$B \rightarrow x$$

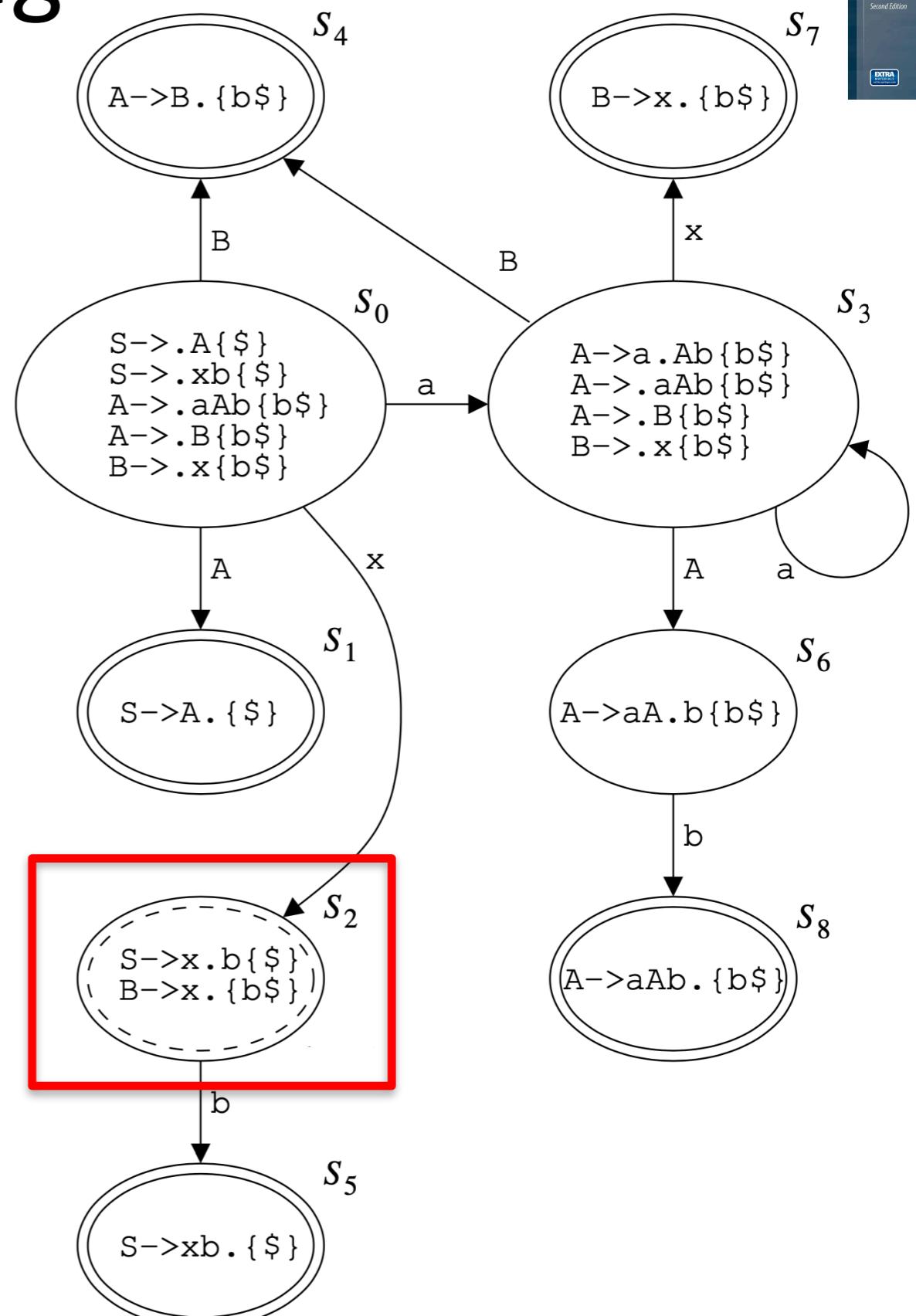
Remark: not LL(1) either — why?

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

$$\text{FOLLOW}(A) = \{ b\$ \}$$

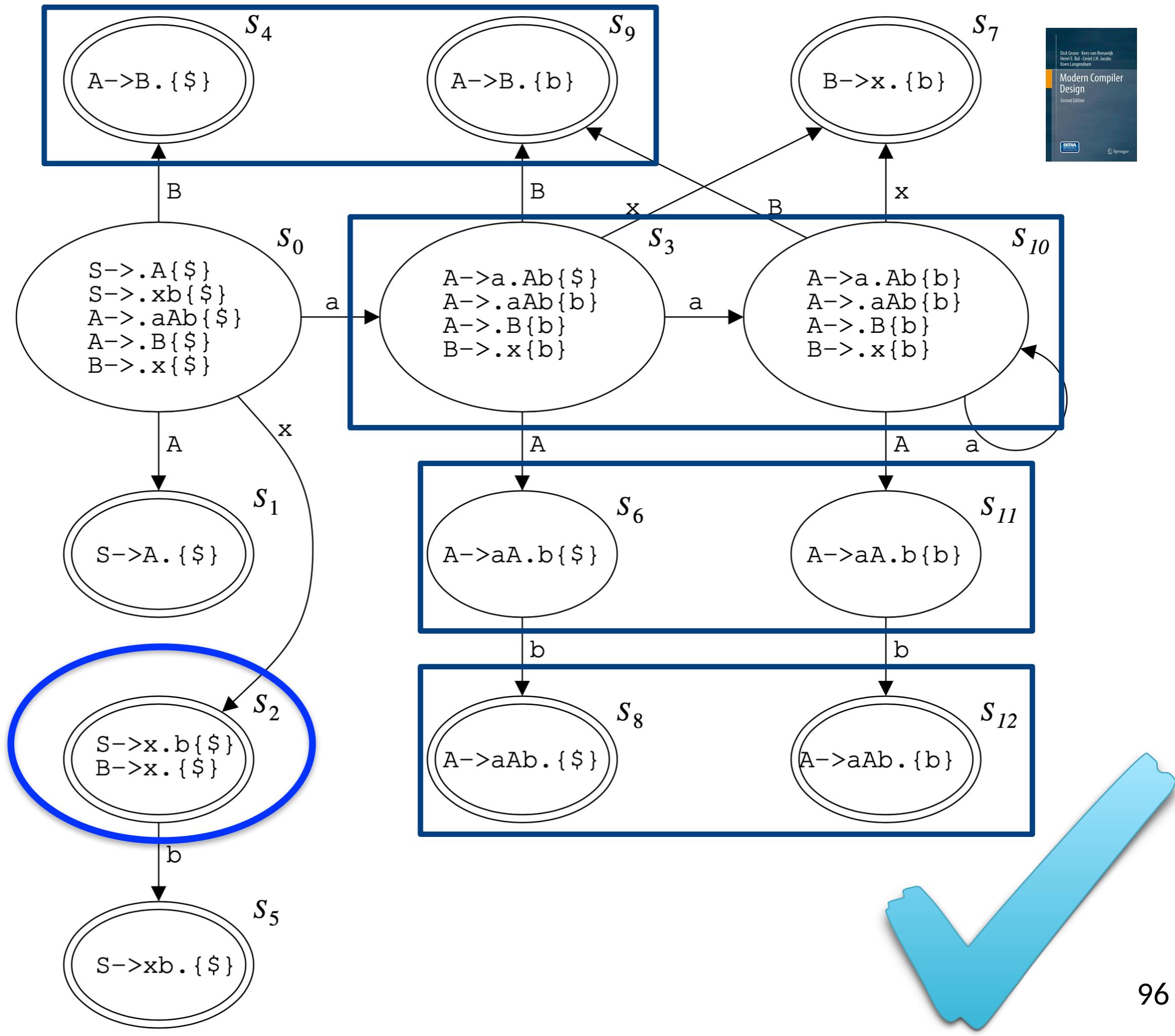
$$\text{FOLLOW}(B) = \{ b\$ \}$$

Follow set is written in automaton inside {}

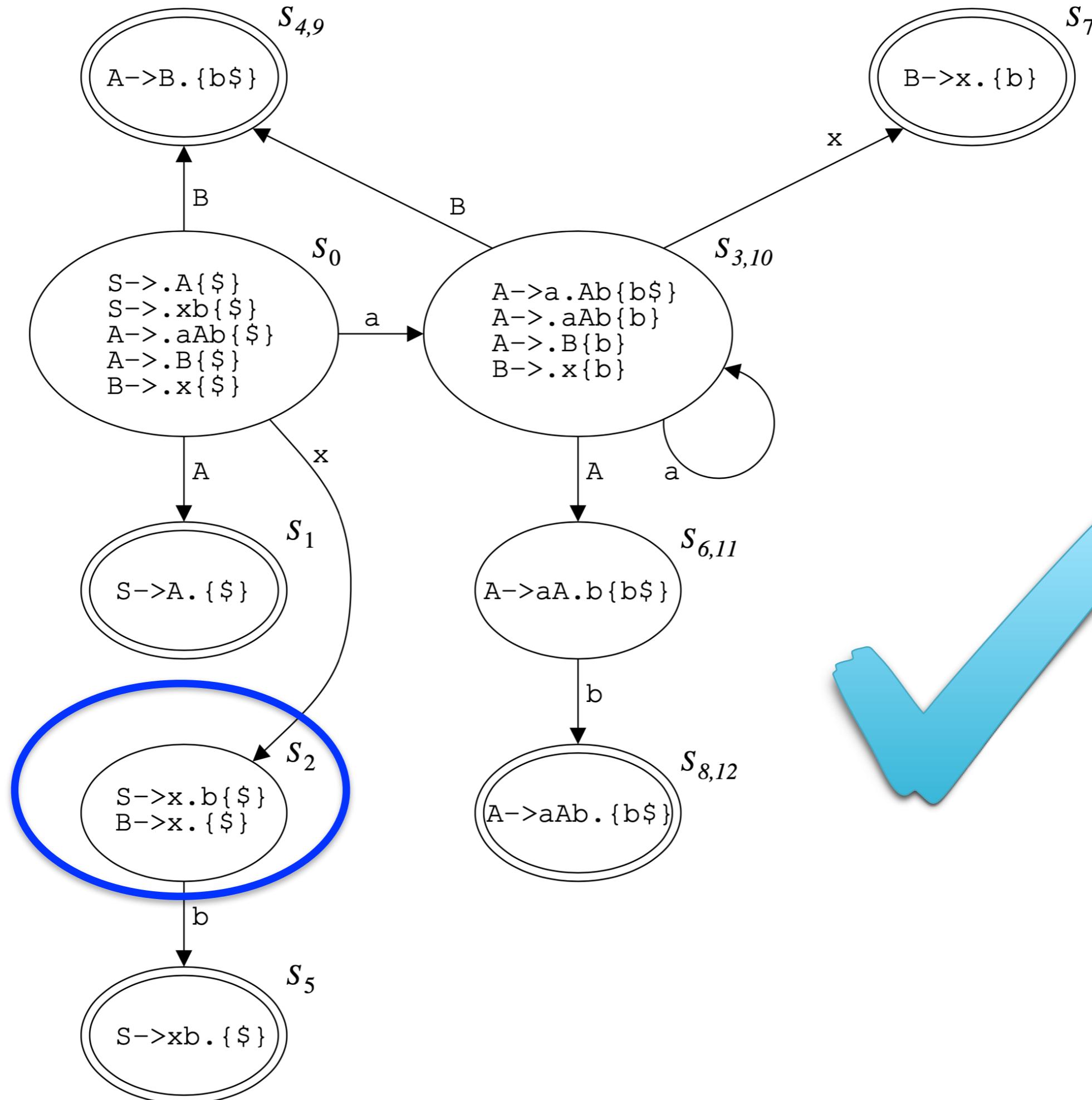


SLR(1) Automaton

Example: LR(1) Automaton

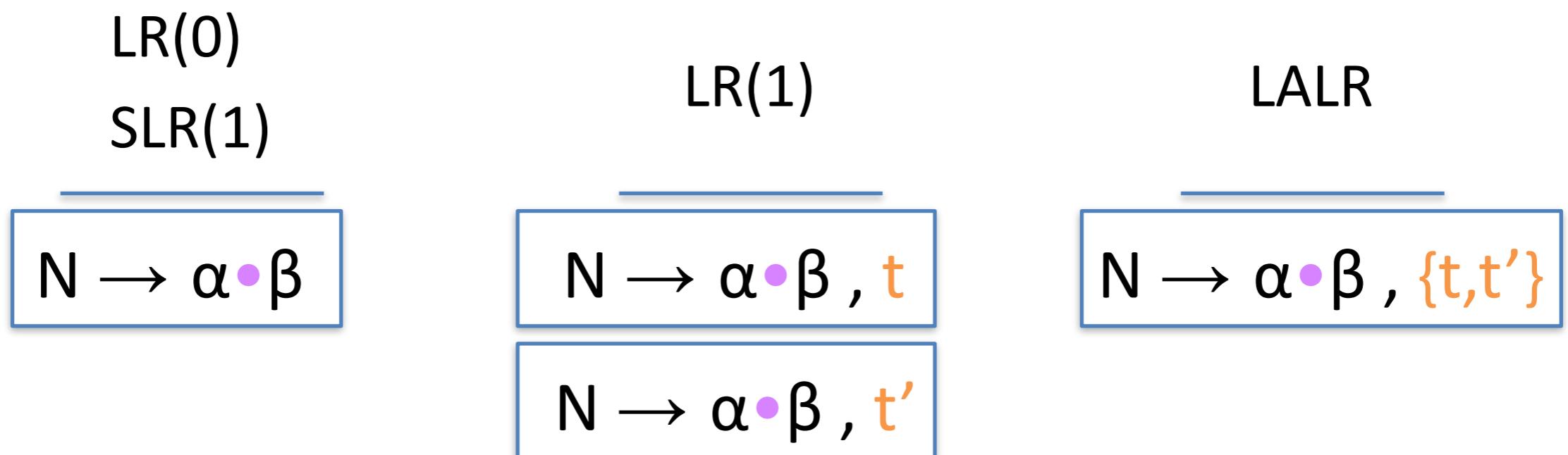


Example: LALR Automaton



Bottom-up Parsing

- * LR(k)
 - * SLR
 - * LALR
-
- All follow the same pushdown-based algorithm
 - Differ on granularity of lookahead in “LR Items”



Building the Parse Tree





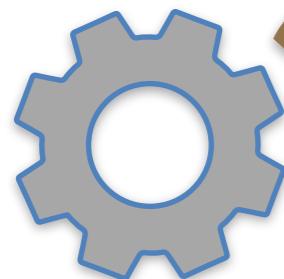
Building the Parse Tree

- Done at the time of **reduce**.

State	ACTION					GOTO	
	*	+	0	1	\$	E	T
q ₈	r ₂	r ₂			r ₂		

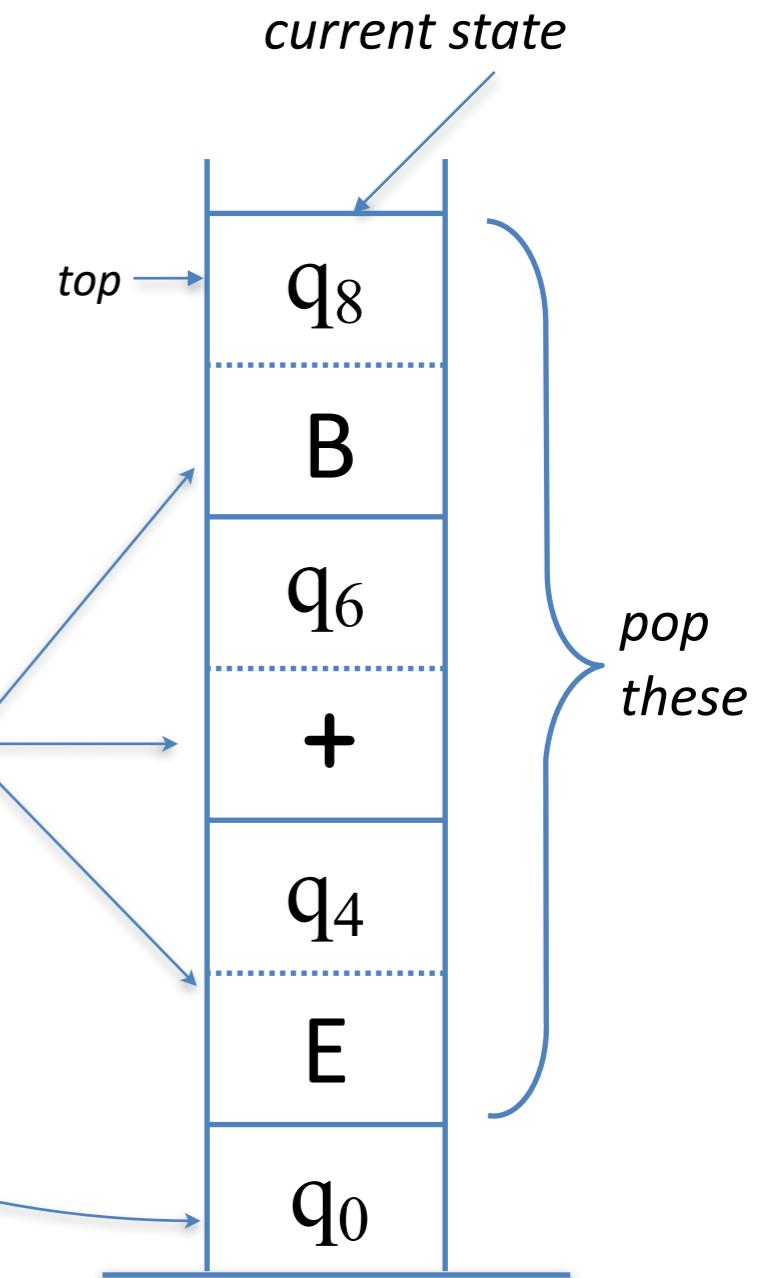
{ :

- (0) S → E
- (1) F → E * B
- (2) E → E + B
- (3) E → B
- (4) B → 0
- (5) B → 1



```

RESULT = new Node("E");
RESULT.addChild(stack[top-6]);
RESULT.addChild(stack[top-4]);
RESULT.addChild(stack[top-2]);
pop(6);
next = GOTO[stack[top], "E"];
push(RESULT);
push(next);
    
```



Building the Parse Tree (“CUP”)



- Done at the time of **reduce**.

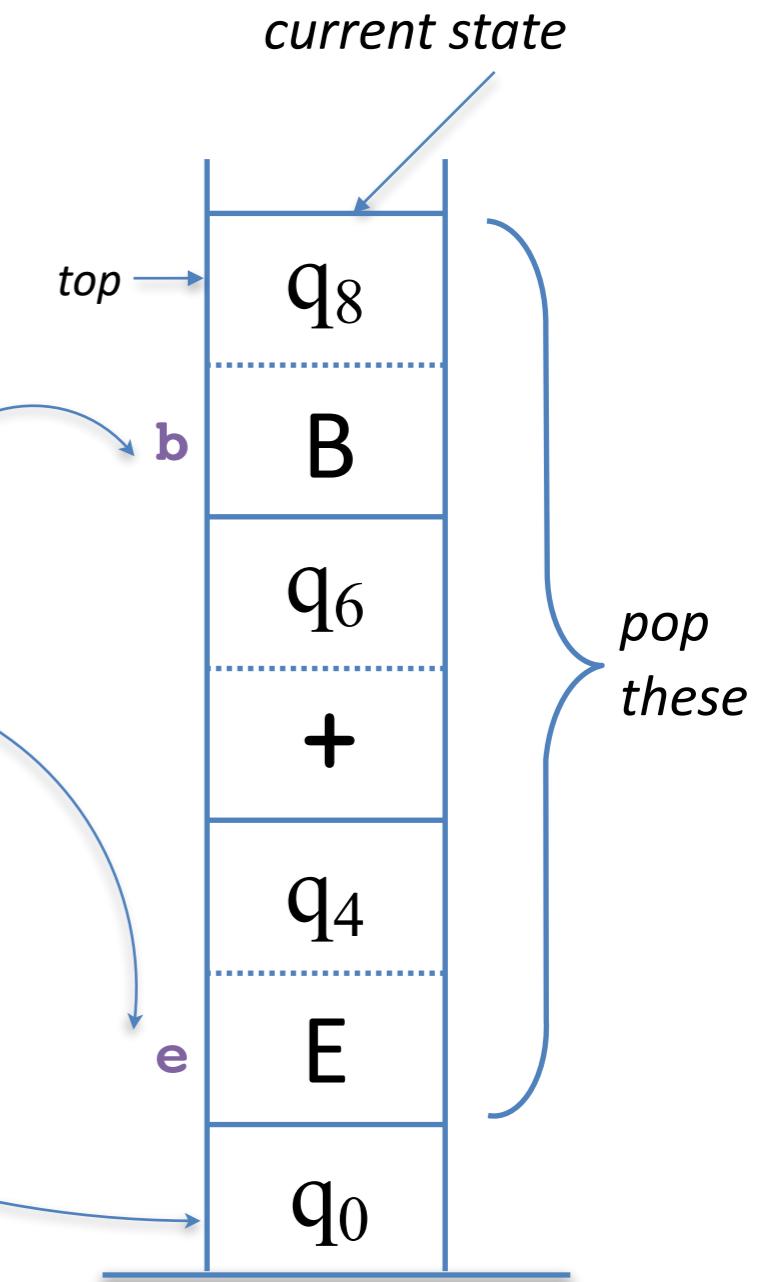
State	ACTION					GOTO	
	*	+	0	1	\$	E	T
q8	r2	r2			r2		

E ::= E:e PLUS B:b {:
 RESULT = new Node("E");
 RESULT.addChild(e);
 RESULT.addChild(new Node("+"));
 RESULT.addChild(b);
 pop(6);
 next = GOTO[stack[top], "E"];
 push(result);
 push(next);

- (0) S → E
- (1) F → E * B
- (2) E → E + B
- (3) E → B
- (4) B → 0
- (5) B → 1



: } :



Building the Parse Tree (“Bison”)



- Done at the time of **reduce**.

State	ACTION					GOTO	
	*	+	0	1	\$	E	T
q ₈	r ₂	r ₂			r ₂		

{

- (0) S → E
- (1) F → E * B
- (2) E → E + B
- (3) E → B
- (4) B → 0
- (5) B → 1

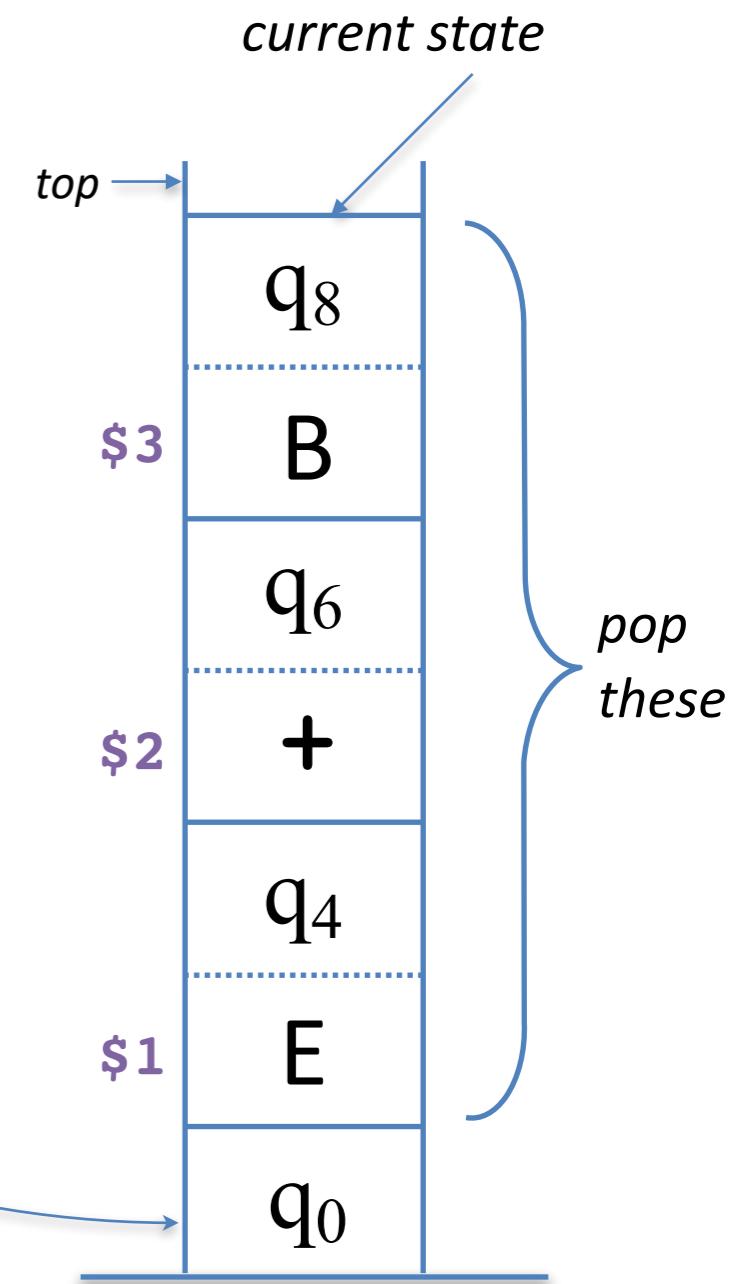


```

$$_ = new Node("E");
$$.addChild($1);
$$.addChild($2);
$$.addChild($3);

pop(6);

next = GOTO[stack[top-1], "E"];
push($$);
push(next);
    
```



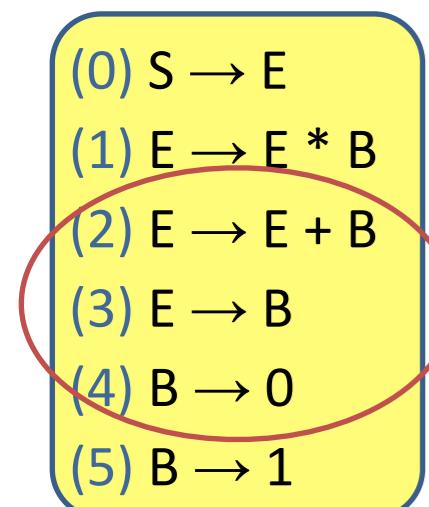
Building an AST

- More useful representation of syntax tree
 - Less clutter
 - Actual level of detail depends on your design
- Basis for semantic analysis
- Later annotated with various information
 - Type information
 - Computed values
- Technically – a class hierarchy of abstract syntax tree nodes

Building the Abstract Syntax Tree



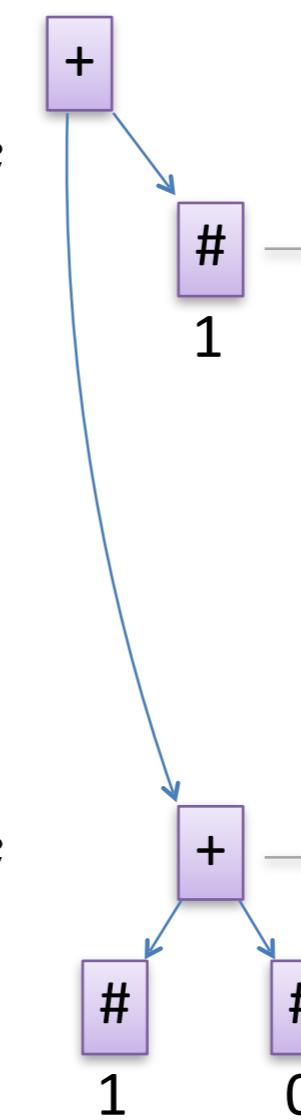
- Generally – just “skip over” the creation of some internal nodes and you get an AST



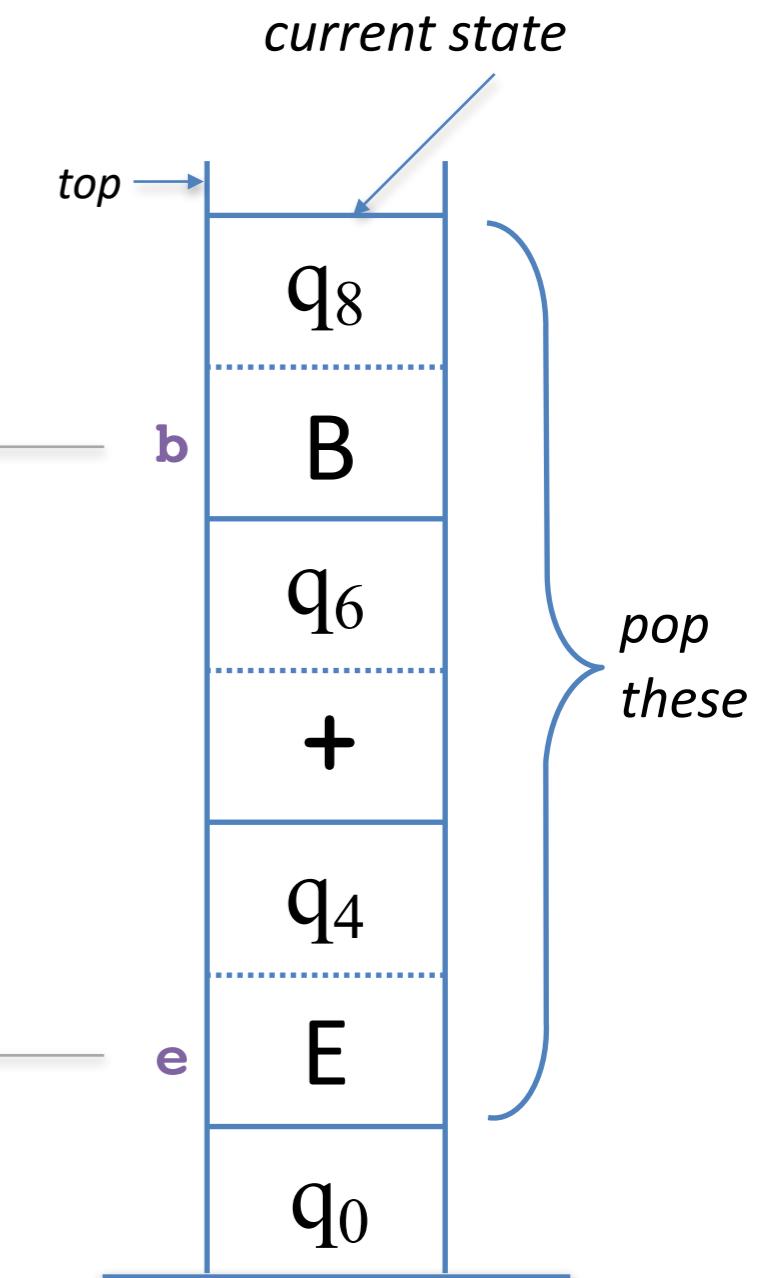
```

E ::= E:e PLUS B:b { :
    RESULT = new Node("+");
    RESULT.addChild(e);
    RESULT.addChild(b);
} |
B:b { :
    RESULT = b;
}
B ::= 0 { :
    RESULT = new Node("#");
    RESULT.value = 0;
}

```



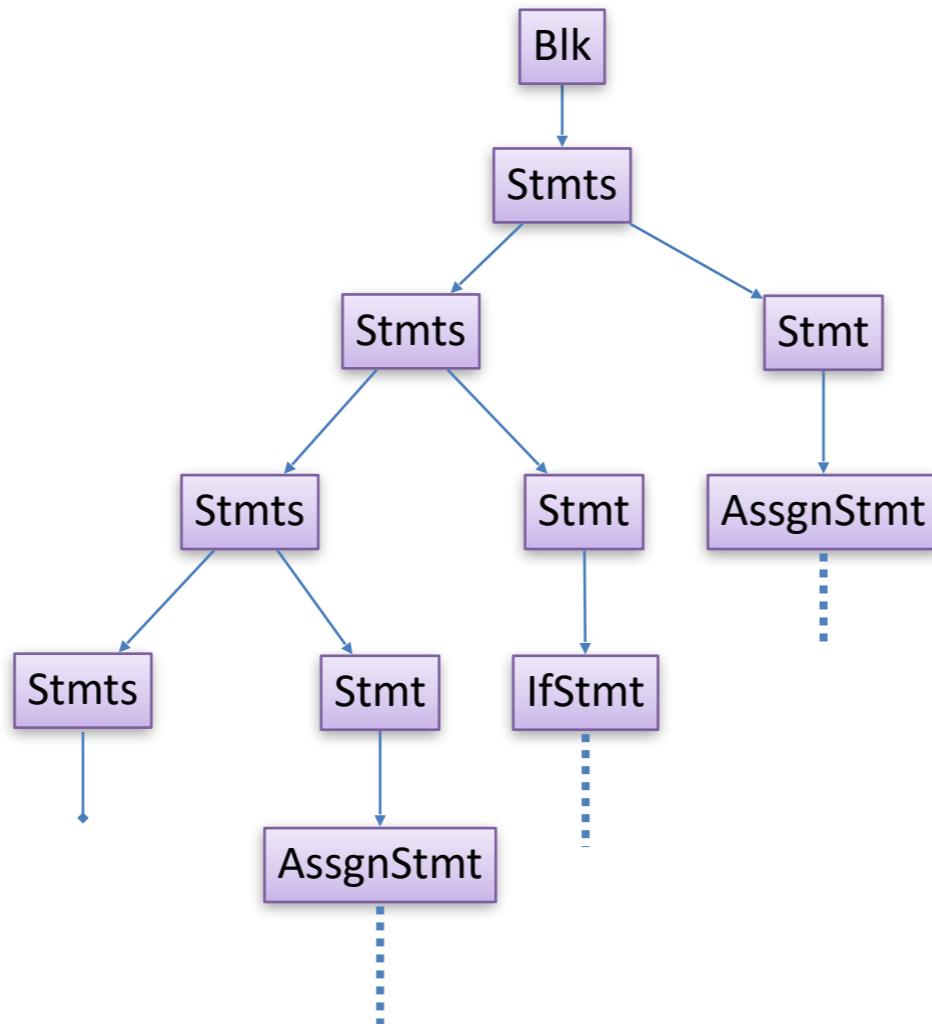
input = 1 + 0 + 1





Building the **Abstract Syntax Tree**

- There is no "The AST" for a grammar
- However, clearly the parse tree is full of artifacts of the grammar we don't care about

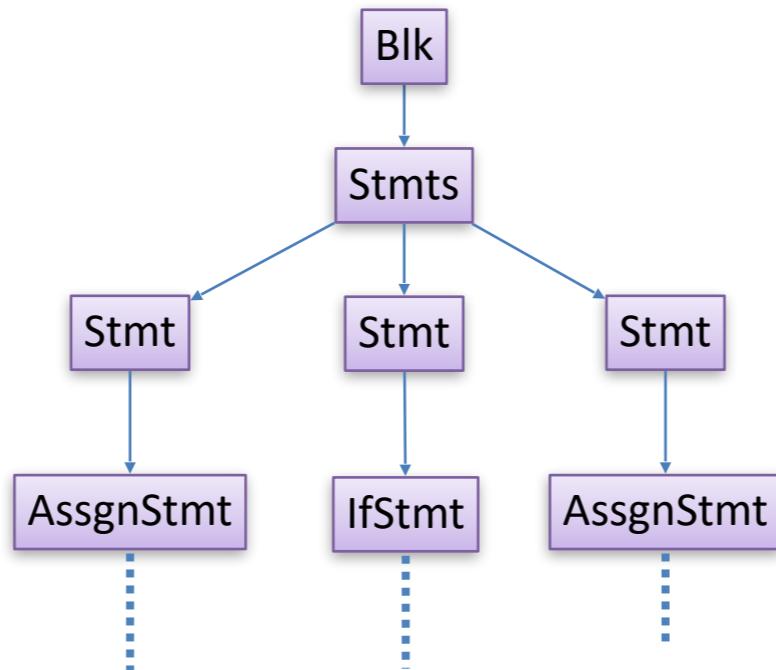


- (0) $\text{Blk} \rightarrow \text{Stmts}$
- (1) $\text{Stmts} \rightarrow \epsilon$
- (2) $\text{Stmts} \rightarrow \text{Stmts Stmt}$
- (3) $\text{Stmt} \rightarrow \text{IfStmt}$
- (4) $\text{Stmt} \rightarrow \text{AssgnStmt}$
- ⋮



Building the **Abstract Syntax Tree**

- This is nicer:

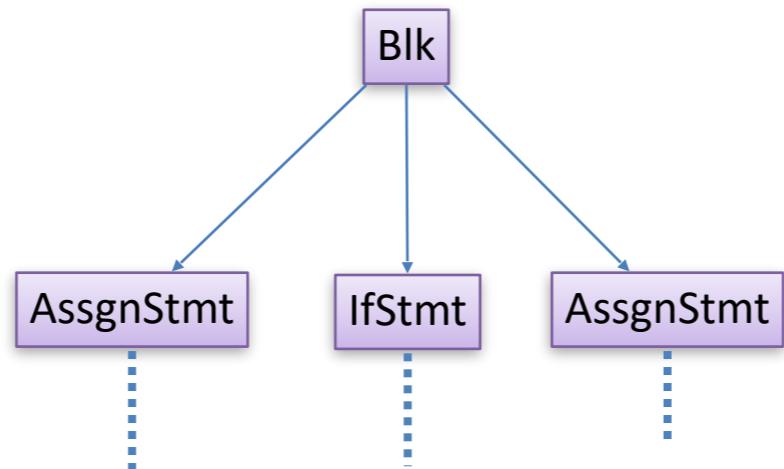


(0) $\text{Blk} \rightarrow \text{Stmts}$
(1) $\text{Stmts} \rightarrow \epsilon$
(2) $\text{Stmts} \rightarrow \text{Stmts Stmt}$
(3) $\text{Stmt} \rightarrow \text{IfStmt}$
(4) $\text{Stmt} \rightarrow \text{AssgnStmt}$
⋮



Building the **Abstract Syntax Tree**

- Even nicer:



(0) $\text{Blk} \rightarrow \text{Stmts}$
(1) $\text{Stmts} \rightarrow \epsilon$
(2) $\text{Stmts} \rightarrow \text{Stmts Stmt}$
(3) $\text{Stmt} \rightarrow \text{IfStmt}$
(4) $\text{Stmt} \rightarrow \text{AssgnStmt}$
⋮



Building the **Abstract Syntax Tree**

- We use the parser actions to skip over and merge nodes as we go

```
Blk → Stmtss:ss {  
    RESULT = new Block();  
    RESULT.statements = s.statements;  
}  
  
Stmts → ε {  
    RESULT = new StatementList();  
}  
  
Stmts → Stmtss:ss Stmt:s {  
    RESULT = ss;  
    RESULT.statements.push(s);  
}  
  
Stmt → IfStmt:ifs {  
    RESULT = ifs;  
};  
}
```

- (0) Blk → Stmtss
- (1) Stmtss → ε
- (2) Stmtss → Stmtss Stmt
- (3) Stmt → IfStmt
- (4) Stmt → AssgnStmt
- ⋮

```
class Block {  
public:  
    vector<Statement> statements;  
}  
  
class StatementList{  
public:  
    vector<Statement> statements;  
}  
  
class Statement;  
  
class IfStatement: public Statement {  
...  
}
```

Error Handling

Error Handling



MCD

S3.2

It has been suggested that with today's fast interactive systems, there is no point in continuing program processing after the first error has been detected, since the user can easily correct the error and then recompile in less time than it would take to read the next error message. But users like to have some idea of how many syntax errors there are left in their program; recompiling several times, each time expecting it to be the last time, is demoralizing. We therefore like to continue the parsing and give as many error messages as there are syntax errors. This means that we have to do error recovery.

Error Handling in LR Parsers



Recovery without modifying the stack

- The **acceptable-set method**

MCD
S3.5.10

- ▶ Step 1: construct a set A of *acceptable* terminals based on the current state of the parser.
- ▶ Step 2: discard input tokens until a token $t_A \in A$ is encountered.
- ▶ Step 3: advance the parser to the next state in which it can consume t_A .

Not practical:
throws away
important tokens
and produce bad
results

A = Correct tokens(**top state**)

don't need to do anything

panic mode



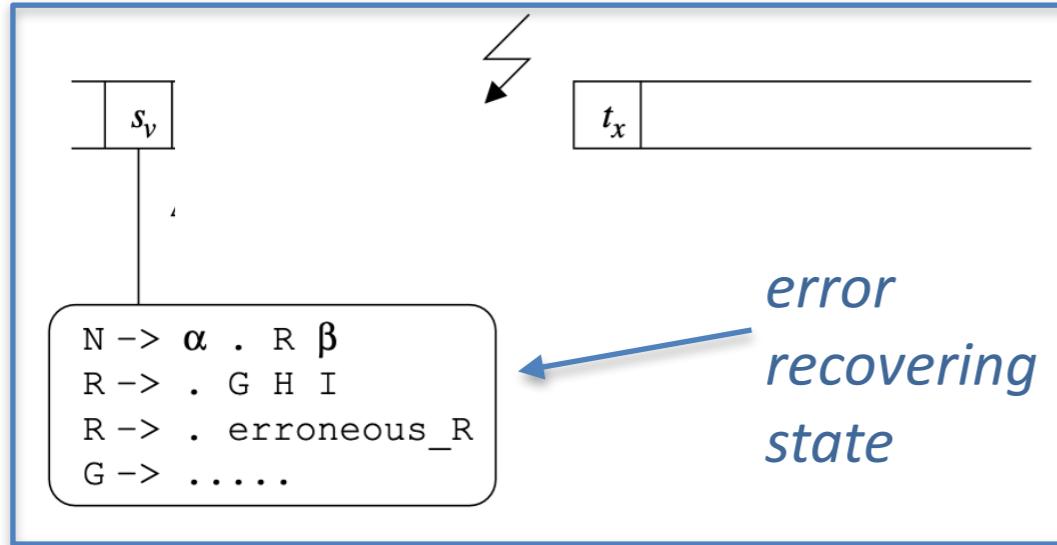
Error Handling in LR Parsers



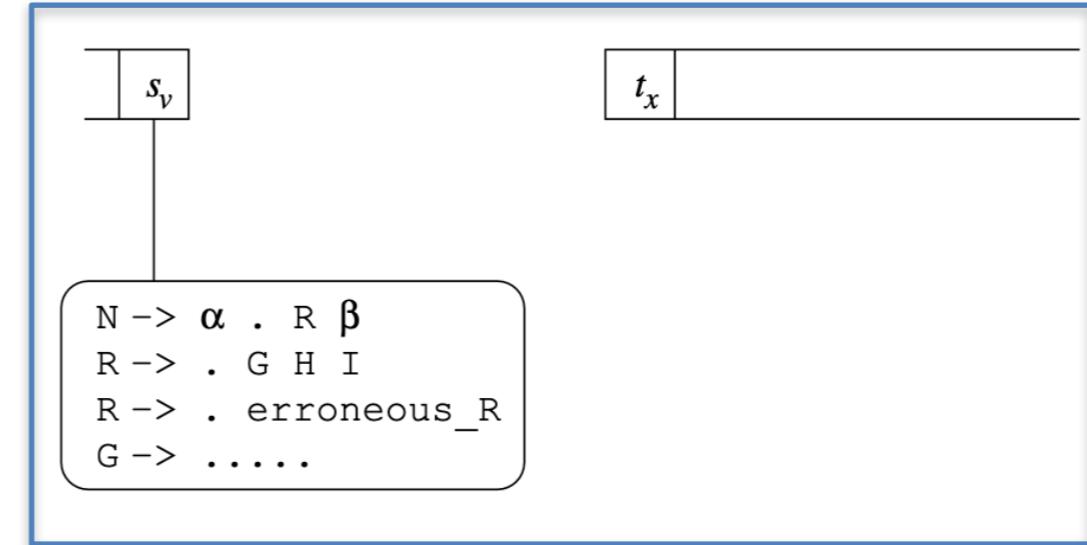
Recovery with stack modification

- **Error-recovering non-terminals**

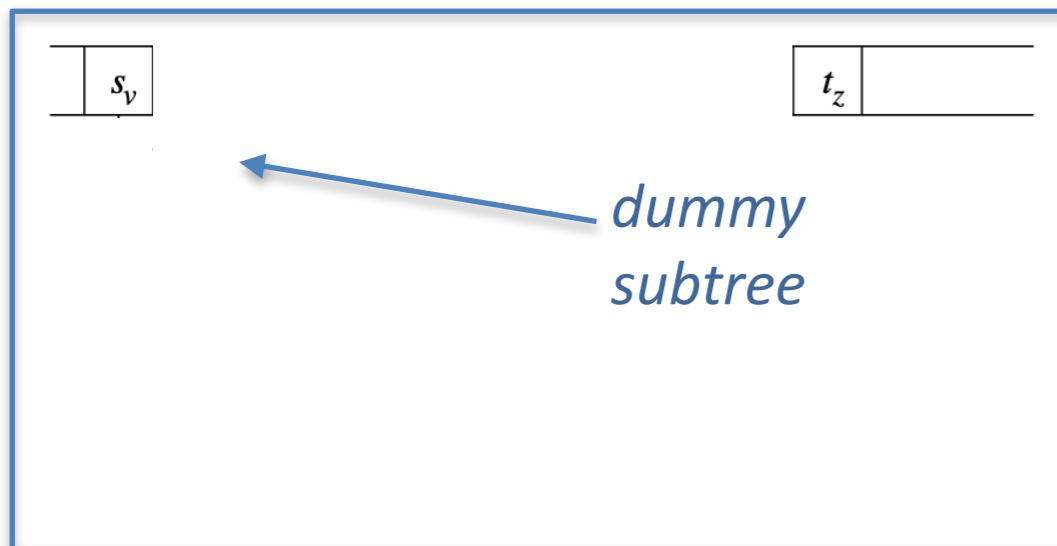
- ▶ Step 1: Define some non-terminal as error-recovering
 - “big names”, e.g., declarations, statements
 - $R \rightarrow G\ H\ I \mid \text{errorneous_}R$
 - Any state containing an item $N \rightarrow \alpha \bullet R\beta$ is an **error recovering state**
- ▶ Step 2: When an error occurs, given up on that non-terminal
 - Pop stack until finding an (error recovering) state with $N \rightarrow \alpha \bullet R\beta$
 - Drop tokens until finding one that "ends" R



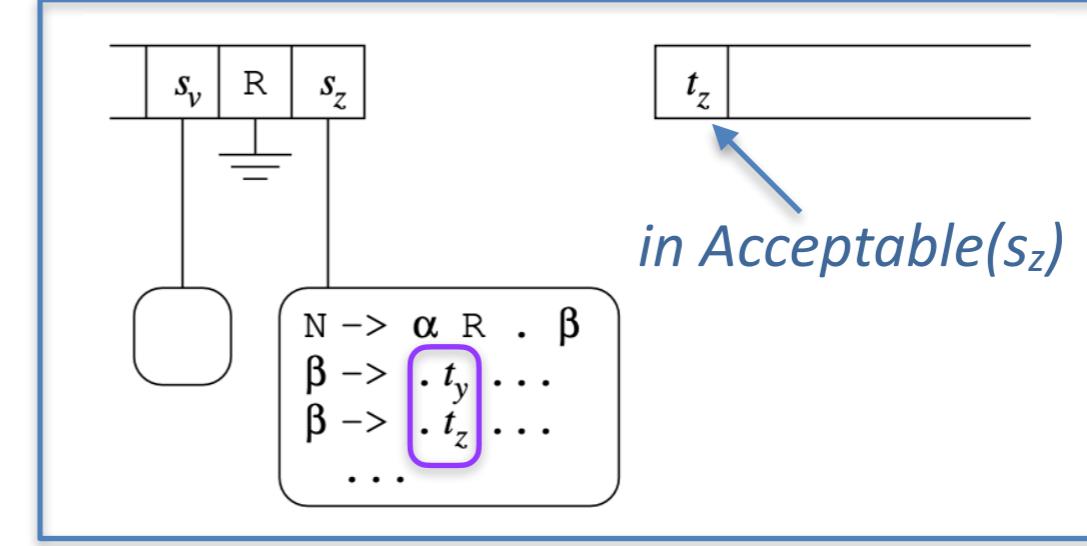
1. Detecting an error



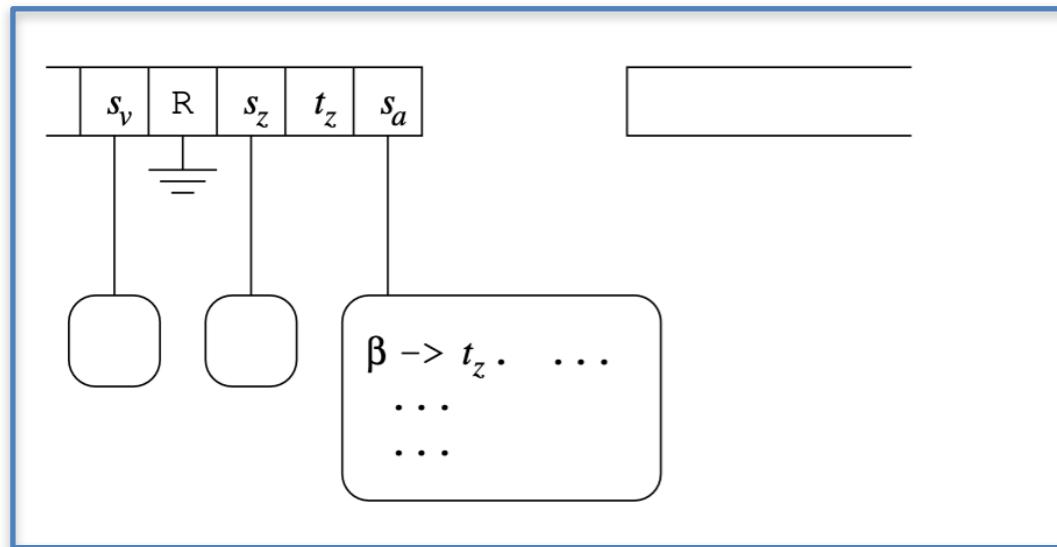
2. Popping to the latest R's error recovery state



3. Repairing the stack



4. Repairing the input



4. Restarting the parser

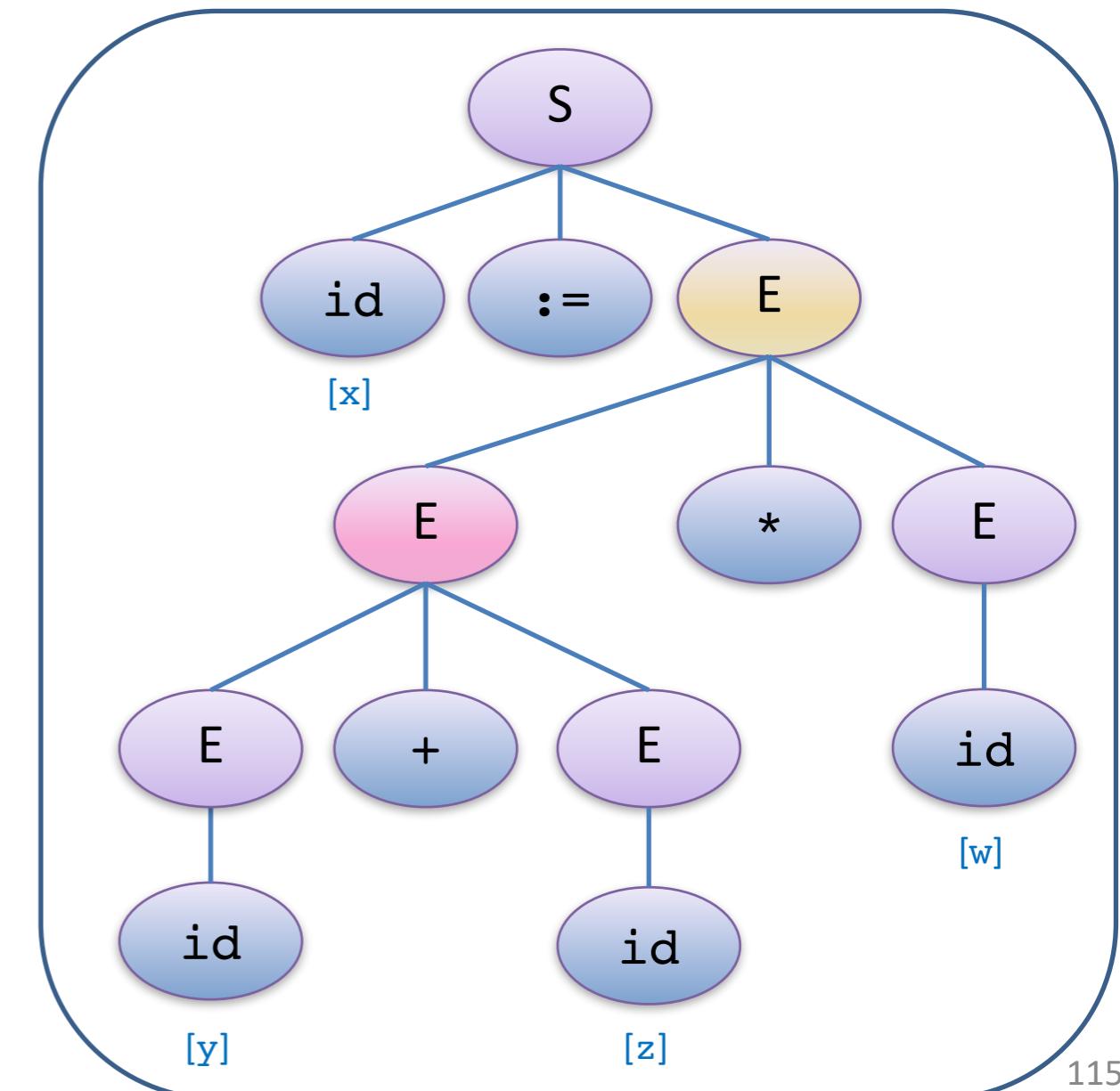
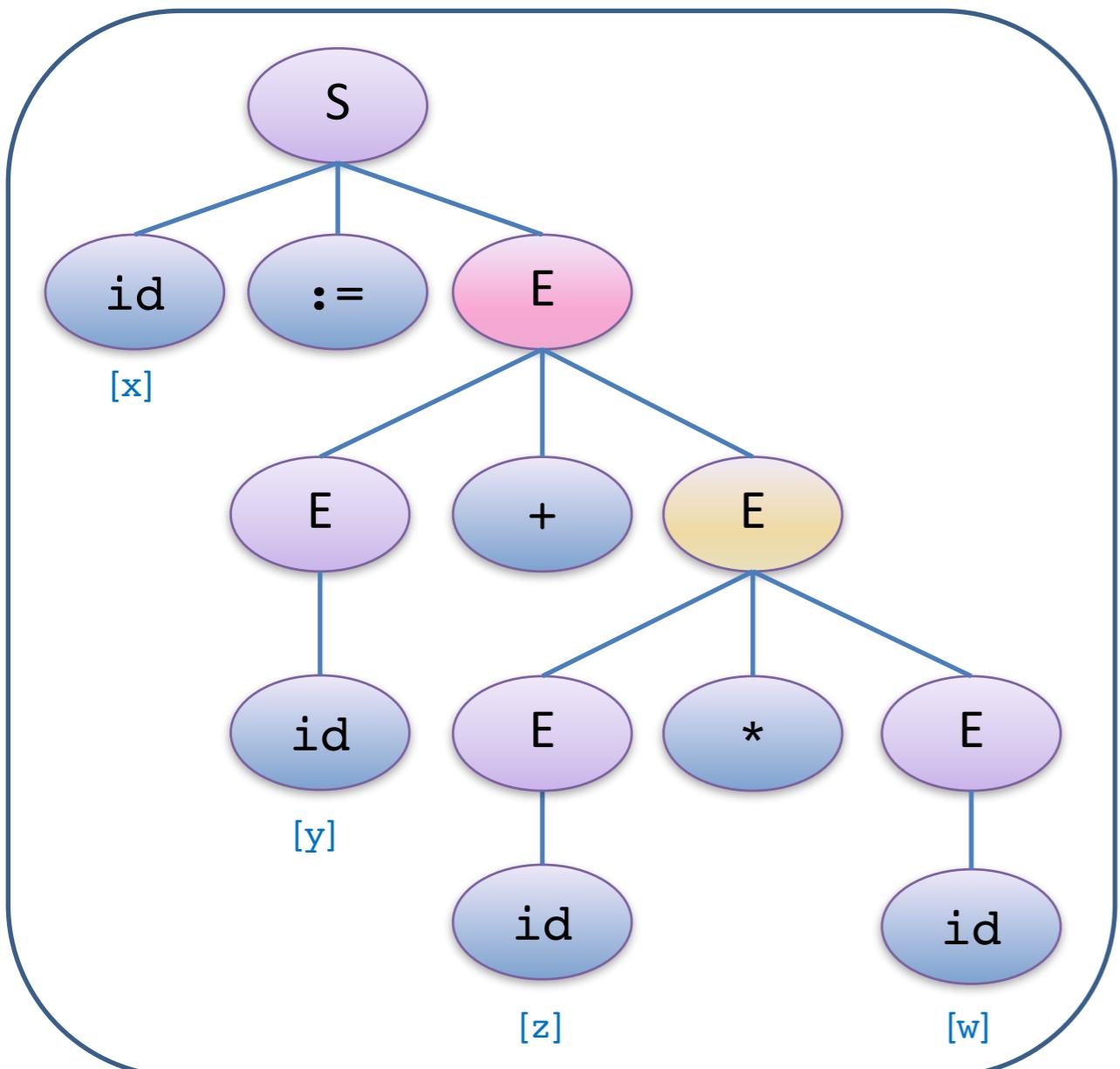
- ✓ Cannot loop
- ✓ Restricts damage to parse tree to a known place
 - ✗ Though discards correctly parsed subtrees
- ✓ Reasonable chance to successfully continue
 - Future phases must deal with dummy trees
 - Must not generate code

Ad Hock Treatment of Ambiguity

Ambiguity

$x := y + z * w$

$S \rightarrow S ; S$
 $S \rightarrow id := E$
 $E \rightarrow id \mid E + E \mid E * E \mid (E)$



Precedence Ambiguity

$x := y + z * w$

$S \rightarrow S ; S$
 $S \rightarrow id := E$
 $E \rightarrow id \mid E + E \mid E * E \mid (E)$

$S \rightarrow id := E \rightarrow id := E + E \rightarrow id := E + E * E \rightarrow id := E + E * id \rightarrow id := E + id * id \rightarrow id := id + id * id$

$id := id + id * id \rightarrow id := E + id * id \rightarrow id := E + E * id \rightarrow id := E + E * E \rightarrow id := E + E \rightarrow id := E \rightarrow S$

$S \rightarrow id := \bullet E \{ ;, \$ \}$
 $E \rightarrow \bullet E + E \{ +, *, ;, \$ \}$
 $E \rightarrow \bullet E * E \{ +, *, ;, \$ \}$
 $E \rightarrow \bullet id \{ +, *, ;, \$ \}$

E

$S \rightarrow id := E \bullet \{ ;, \$ \}$
 $E \rightarrow E \bullet + E \{ +, *, ;, \$ \}$
 $E \rightarrow E \bullet * E \{ +, *, ;, \$ \}$

$+ E$

$E \rightarrow E + E \bullet \{ +, *, ;, \$ \}$
 $E \rightarrow E \bullet + E \{ +, *, ;, \$ \}$
 $E \rightarrow E \bullet * E \{ +, *, ;, \$ \}$

id

$E \rightarrow id \bullet \{ +, *, ;, \$ \}$

I would like to say:
* precedes +, so please, dear parser,
prefer shift * over reduce $E + E$

$+ < *$

Associativity Ambiguity

$x := y + z + w$

$S \rightarrow S ; S$
 $S \rightarrow id := E$
 $E \rightarrow id \mid E + E \mid E * E \mid (E)$

$S \rightarrow id := E \rightarrow id := E + E \rightarrow id := E + E + E \rightarrow id := E + E + id \rightarrow id := E + id + id \rightarrow id := id + id + id$

$id := id + id + id \rightarrow id := E + id + id \rightarrow id := E + E + id \rightarrow id := E + E + E \rightarrow id := E + E \rightarrow id := E \rightarrow S$

$S \rightarrow id := \bullet E \{ ;, \$ \}$
 $E \rightarrow \bullet E + E \{ +, *, ;, \$ \}$
 $E \rightarrow \bullet E * E \{ +, *, ;, \$ \}$
 $E \rightarrow \bullet id \{ +, *, ;, \$ \}$

E

$S \rightarrow id := E \bullet \{ ;, \$ \}$
 $E \rightarrow E \bullet + E \{ +, *, ;, \$ \}$
 $E \rightarrow E \bullet * E \{ +, *, ;, \$ \}$

$+ E$

$E \rightarrow E + E \bullet \{ +, *, ;, \$ \}$
 $E \rightarrow E \bullet + E \{ +, *, ;, \$ \}$
 $E \rightarrow E \bullet * E \{ +, *, ;, \$ \}$

id

$E \rightarrow id \bullet \{ +, *, ;, \$ \}$

I would like to say:
 $+$ is left associative, so please, dear parser,
prefer reduce $E + E$ over shift



left
associative

“Automatically” Resolving Conflicts in LR Parser

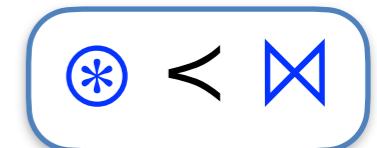
- **Shift-reduce conflicts**

- Highest (user-specified) precedence wins:

- $P \rightarrow \alpha \bullet \otimes \beta \{ \dots \}$



- $Q \rightarrow \mu \bowtie R \bullet \{ \dots \}$ or $Q \rightarrow \mu \bowtie \bullet \{ \dots \}$



- (User-specified) associativity wins:

- $P \rightarrow \alpha \bullet \otimes \beta \{ \dots \}$



- $Q \rightarrow \mu \otimes R \bullet \{ \dots \}$ or $Q \rightarrow \mu \otimes \bullet \{ \dots \}$

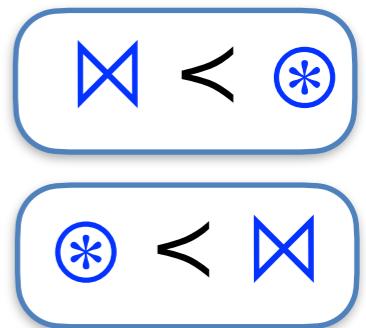


- Default: use shift

a bit more general: allow
for different terminals with
equal precedence

“Automatically” Resolving Conflicts in LR Parser

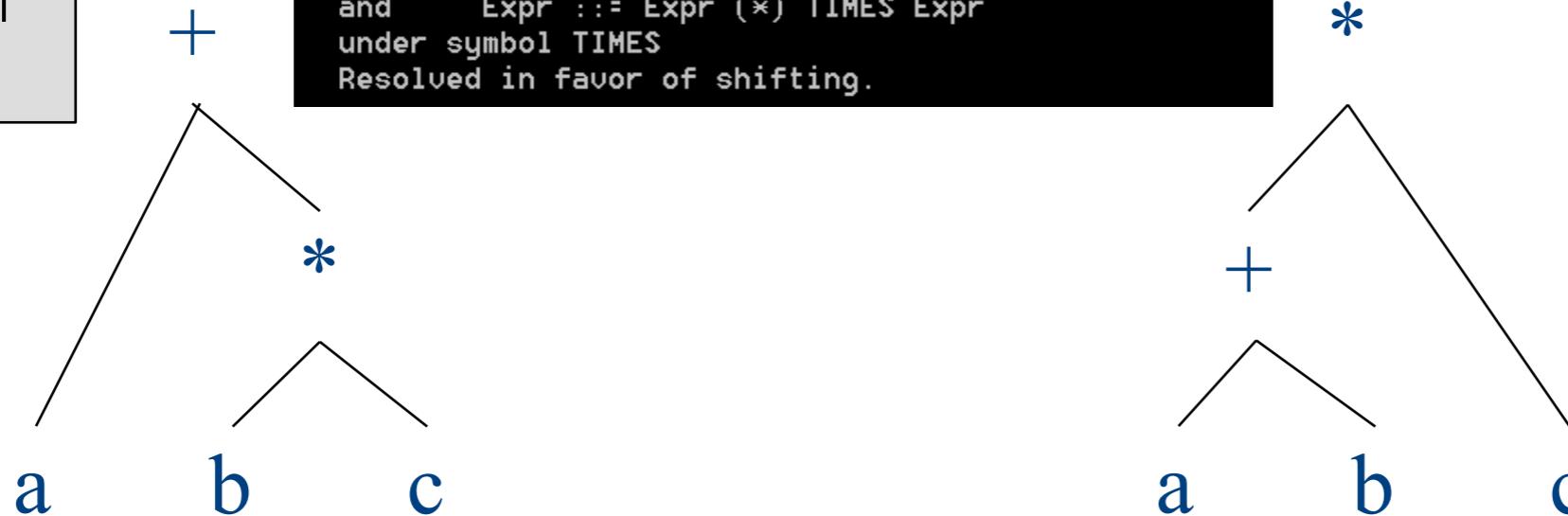
- **Shift-reduce conflicts**
 - Highest (user-specified) precedence win:
 - $P \rightarrow \alpha \cdot \circledast \beta \{ \dots \}$
 - $Q \rightarrow \mu \bowtie R \cdot \{ \dots \circledast \dots \}$ or $Q \rightarrow \mu \bowtie \cdot \{ \dots \circledast \dots \}$
 - Solves if/else ambiguity
 - $\text{if_S} \rightarrow \text{if } (\text{BoolExpr}) S \cdot \{ \dots \text{else} \dots \}$
 - $\text{if_S} \rightarrow \text{if } (\text{BoolExpr}) S \cdot \text{else } S \{ \dots \text{else} \dots \}$
- **Reduce-reduce conflicts**
 - First rule wins



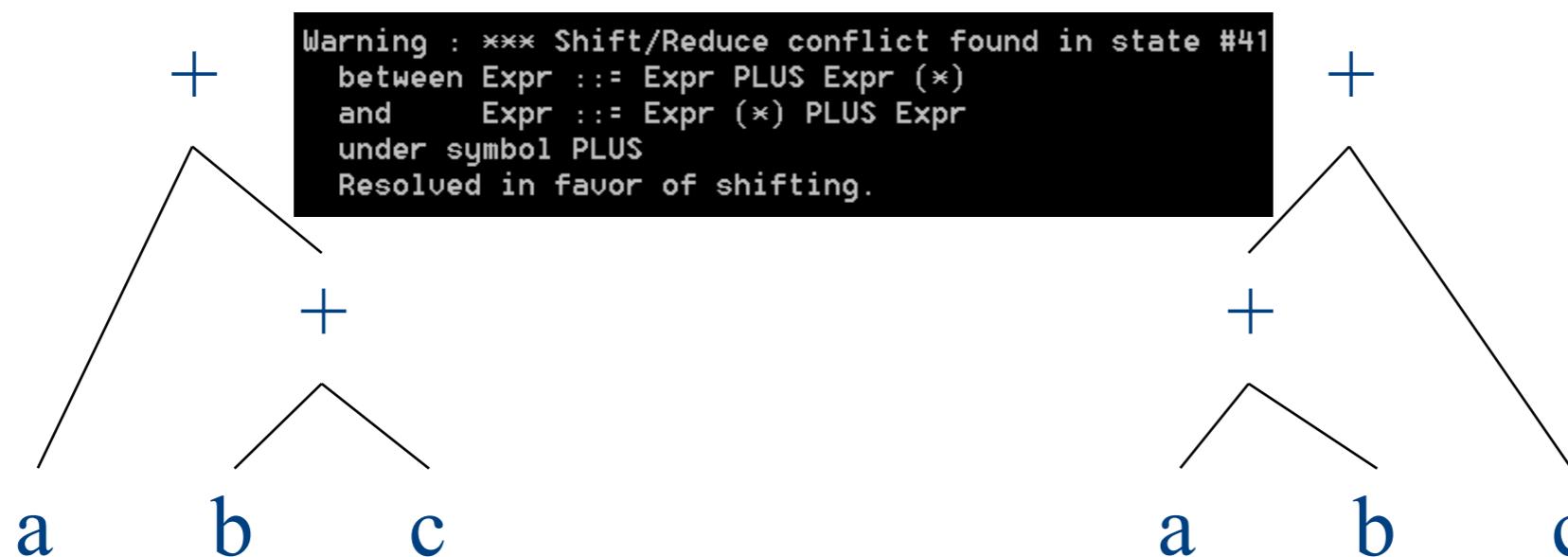
Ambiguities as conflicts in CUP

```
E → id |  
      E + E |  
      E * E
```

Warning : *** Shift/Reduce conflict found in state #41
between Expr ::= Expr PLUS Expr (*)
and Expr ::= Expr (*) TIMES Expr
under symbol TIMES
Resolved in favor of shifting.

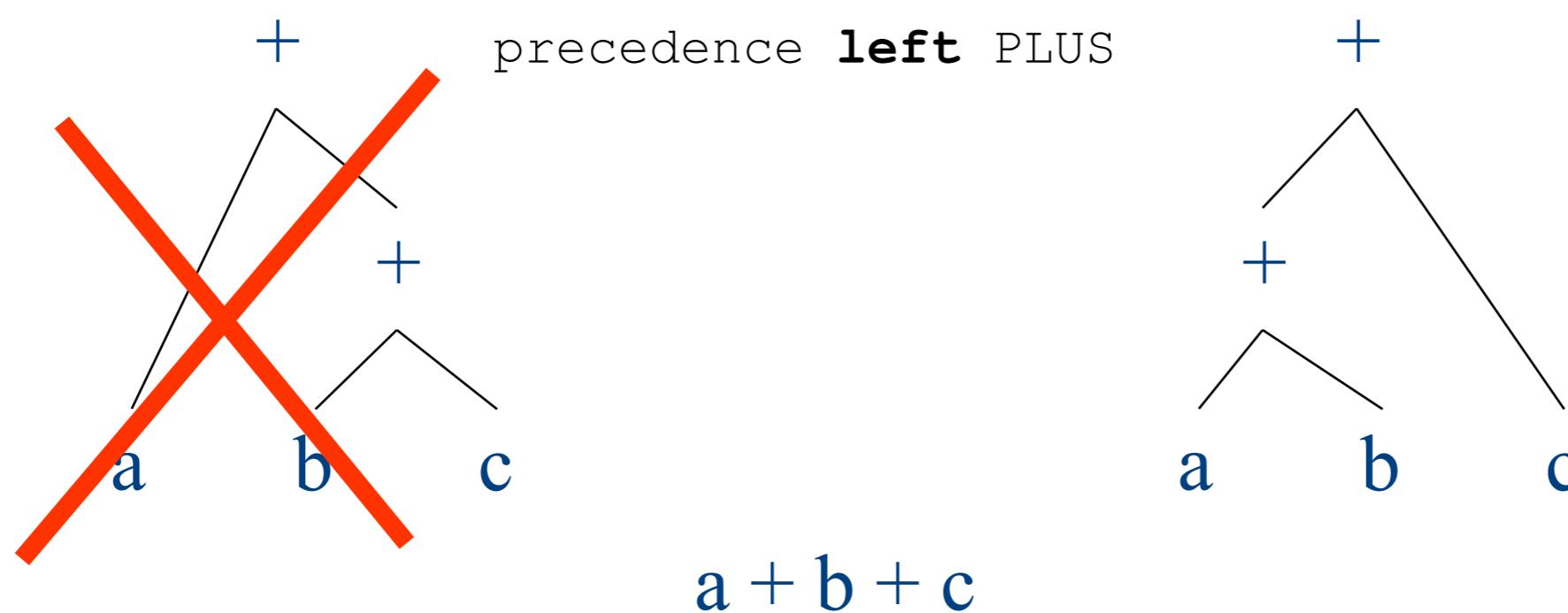
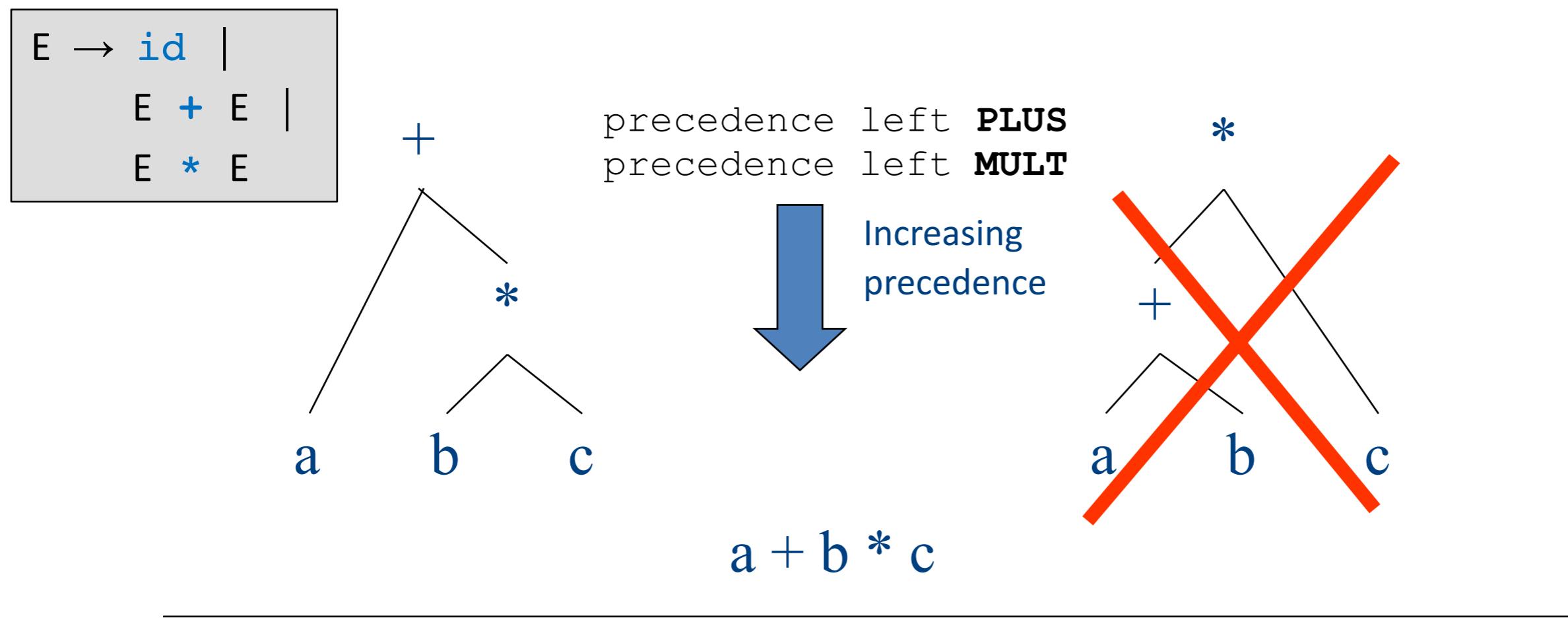


$a + b * c$



$a + b + c$

Ambiguities as conflicts in CUP



Automated Parser Generation (e.g., via CUP)

Will be discussed more
thoroughly in the recitation