

# **Parsing (2)**

**15-411/15-611 Compiler Design**

Ben L. Titzer and Seth Copen Goldstein

Feb 13, 2025

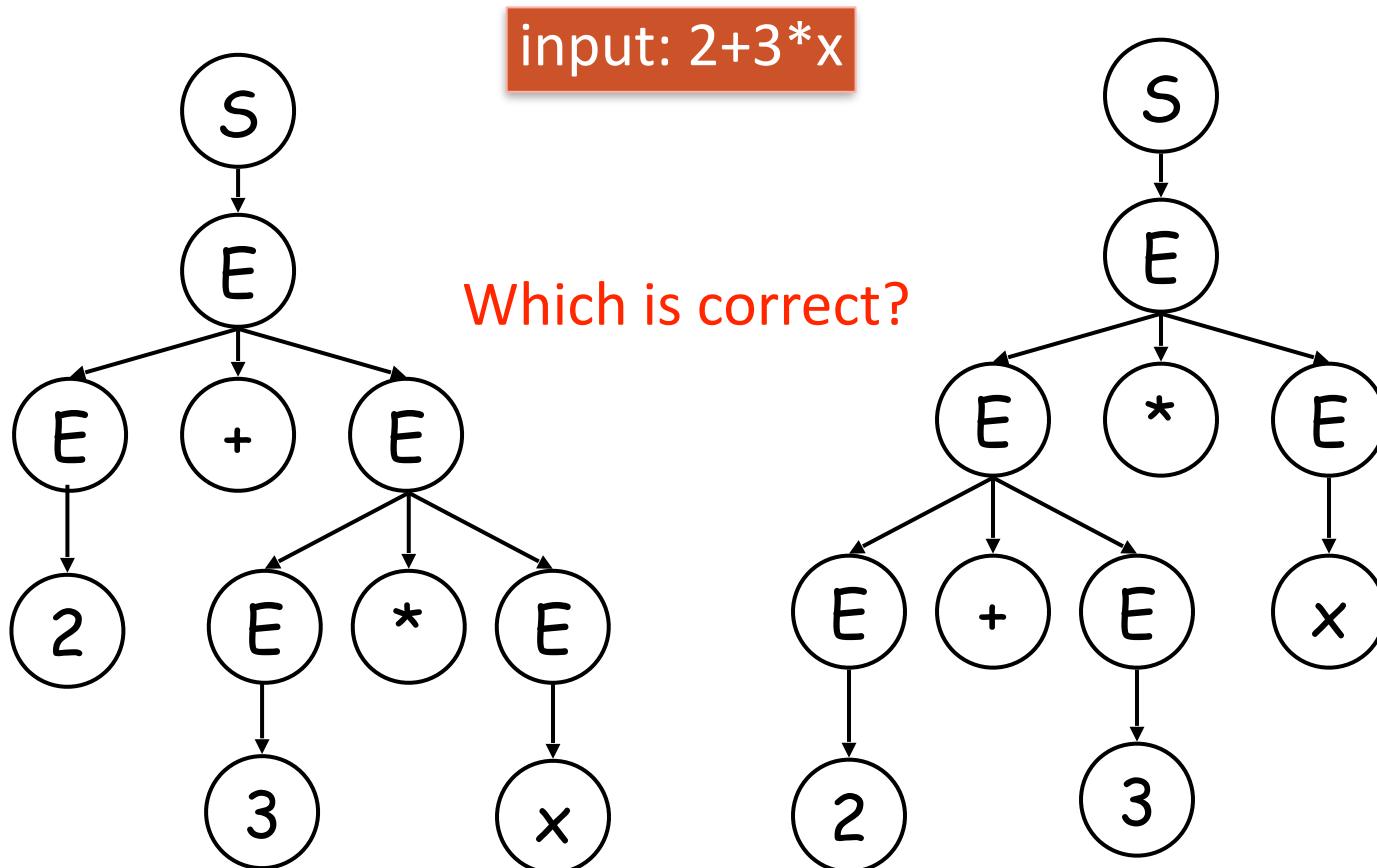
# Today

## Parsing

- Bottom-up parsers
  - constructing state machine
  - LR0
  - SLR
  - LR( $k$ ) & LALR
  - Handling Ambiguity
- Typechecking

# Ambiguity in Grammars

- Some grammars have more than one way to parse a given input, i.e. are ambiguous



# Resolving Ambiguity

- Ambiguity is a problem with the grammar
- One possible fix: Add precedence with more non-terminals
- In this example, one for each level of precedence:
  - $(+, -)$  exp
  - $(*, /)$  term
  - **(id, int)** factor
  - Make sure parse derives sentences that respect the precedence
  - Make sure that extra levels of precedence can be bypassed, i.e., “x” is still legal

# A Better Exp Grammar

1	S	:= Exp	S
2	Exp	:= Exp + Term	by 1 $\Rightarrow$ Exp
3	Exp	:= Exp - Term	by 2 $\Rightarrow$ Exp + Term
4	Exp	:= Term	by 4 $\Rightarrow$ Term + Term
5	Term	:= Term * Factor	by 7 $\Rightarrow$ Factor + Term
6	Term	:= Term / Factor	by 9 $\Rightarrow$ int <sub>2</sub> + Term
7	Term	:= Factor	by 5 $\Rightarrow$ int <sub>2</sub> + Term * Factor
8	Factor	:= id	by 7 $\Rightarrow$ int <sub>2</sub> + Factor * Factor
9	Factor	:= int	by 9 $\Rightarrow$ int <sub>2</sub> + int <sub>3</sub> * Factor
			by 8 $\Rightarrow$ int <sub>2</sub> + int <sub>3</sub> * id..

What is the parse tree?

# Parsing a Grammar

- Top-Down
  - start at root of parse-tree
  - pick a production and expand to match input
  - may require backtracking
  - if no backtracking required, predictive
- Bottom-up
  - start at leaves of tree
  - recognize valid prefixes of productions
  - consume input and change state to match
  - use stack to track state

# Computing FIRST( $\alpha$ )

- Given  $X := A \ B \ C$ ,  $\text{FIRST}(X) = \text{FIRST}(A \ B \ C)$
- Can we ignore B or C?
- Consider:

$A := a$

|

$B := b$

| A

$C := c$

# Computing FIRST( $\alpha$ )

- Given  $X := A B C$ ,  $\text{FIRST}(X) = \text{FIRST}(A B C)$
- Can we ignore B or C?
- Consider:

A := a

A :=

B := b

B := A

C := c

- $\text{FIRST}(X)$  must also include  $\text{FIRST}(C)$
- IOW:
  - Must keep track of NTs that are **nullable**
  - For **nullable** NTs, determine **FOLLOWs(NT)**

# nullable(A)

- **nullable(A)** is
  - true if A can derive the empty string
  - false otherwise
- For example:

$B := X Y b$

$X := x$

|  
 $Y Y$

$Y :=$

In this case, **nullable(X) = nullable(Y) = true**

**nullable(B) = false**

# FOLLOW(A)

- FOLLOW(A) is the set of terminals that can immediately follow A in a sentential form.
- I.e.,  
 $a \in \text{FOLLOW}(A)$  iff  $S \Rightarrow^* \alpha A a \beta$  for some  $\alpha$  and  $\beta$

# Bottom-up parsers

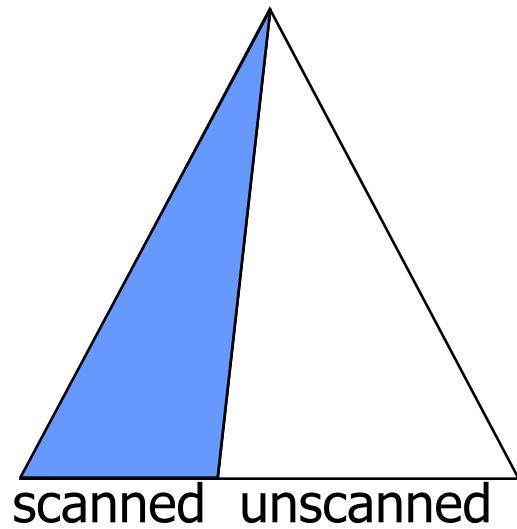
- What is the inherent restriction of top-down parsing, e.g., with LL( $k$ ) grammars?

# Bottom-up parsers

- What is the inherent restriction of top-down parsing, e.g., with LL( $k$ ) grammars?
- Bottom-up parsers use the entire right-hand side of the production
- LR( $k$ ):
  - Left-to-right parse,
  - Rightmost derivation (in reverse),
  - $k$  look ahead tokens

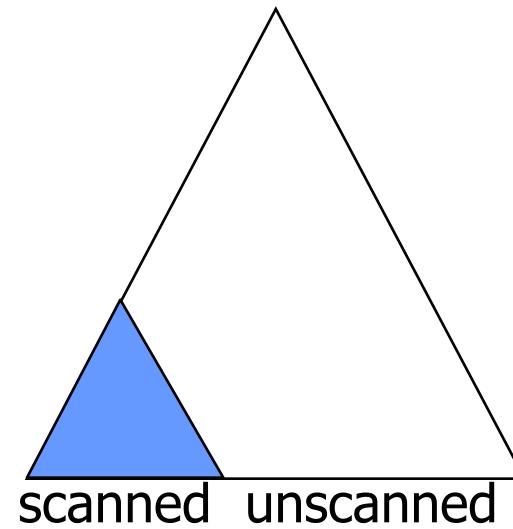
# Top-down vs. Bottom-up

LL( $k$ ), recursive descent



Top-down

LR( $k$ ), shift-reduce



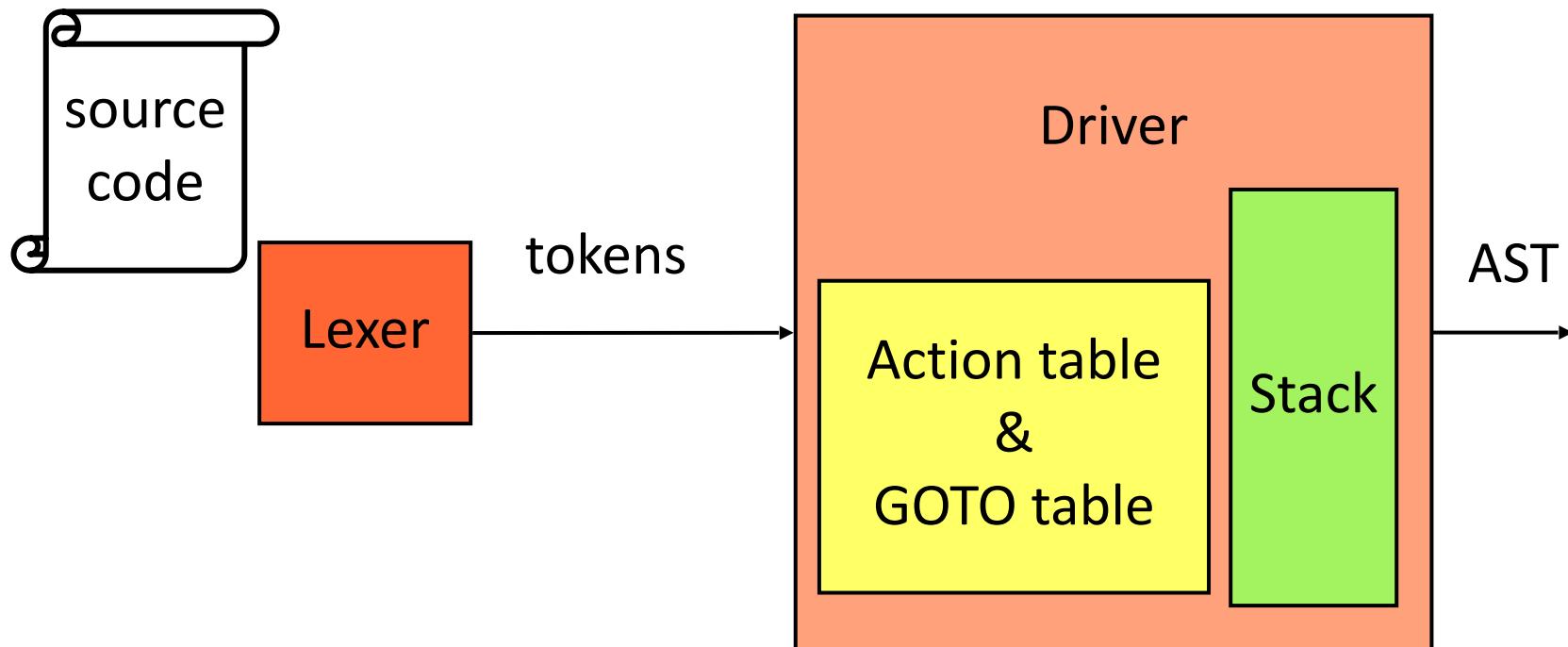
Bottom-up

# A Shift-Reduce Parser

- Implement as a FSM with a stack
- Stack holds sequences of symbols
- Input stream holds remaining source
- Four actions:
  - shift: push token from input stream onto stack
  - reduce: right-end of a handle ( $\beta$  of  $A \rightarrow \beta$ ) is at top of stack, pop handle ( $\beta$ ), push A
  - accept: success
  - error: syntax error discovered

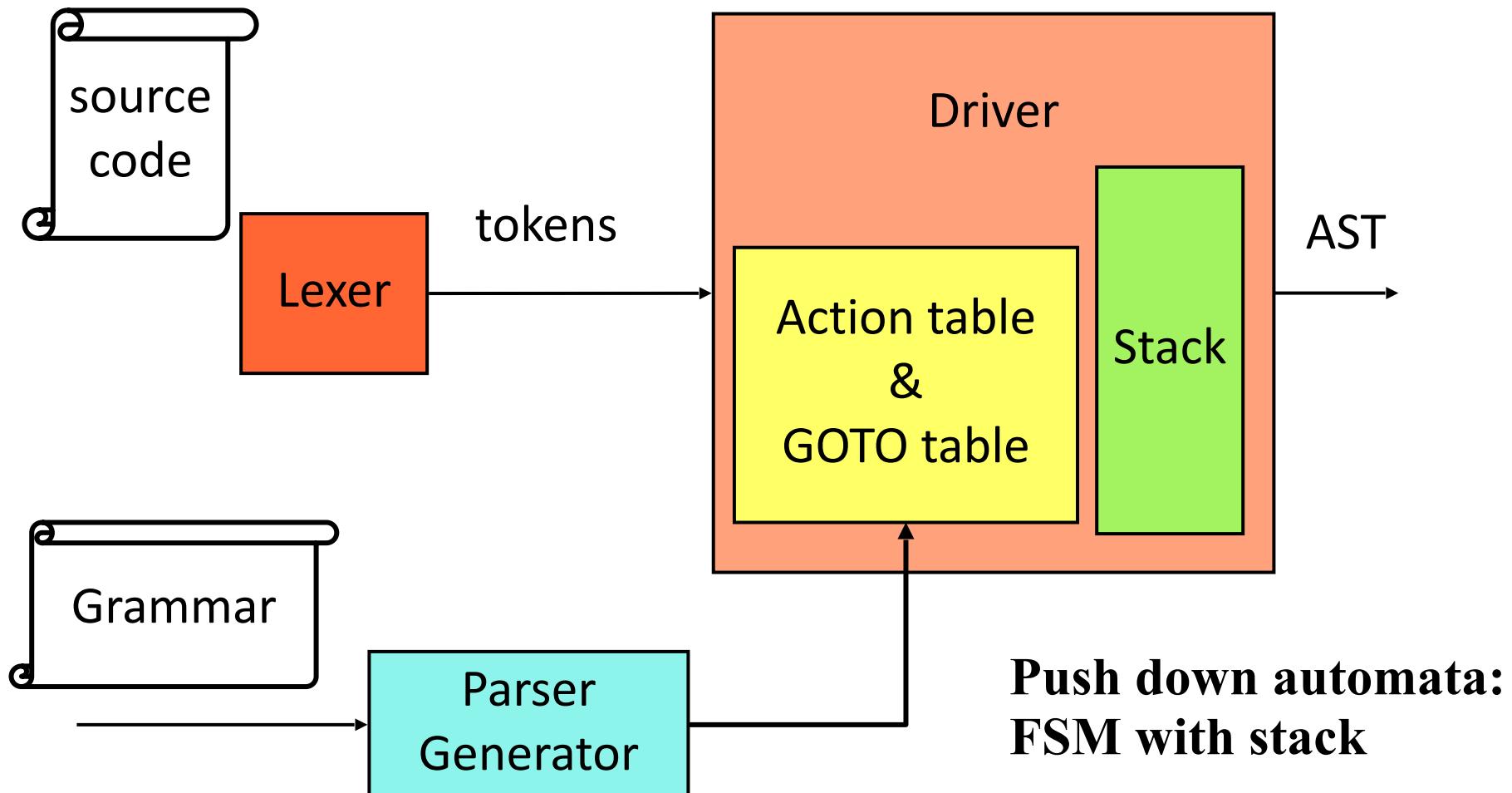
Key is recognizing handles efficiently

# Table-driven LR(k) parsers



**Push down automata:  
FSM with stack**

# Table-driven LR(k) parsers



# Parser Loop

Driver

- Same code regardless of grammar
  - only tables change
- (Very) General Algorithm:
  - Based on table contents, top of stack, and current input token either
    - **shift**: push token onto stack and read next token
    - **reduce**: replace part of stack with the correct rule (NT) that derived it
    - **accept**: successfully parsed entire input
    - **error**: input not in language

# Stack

- Represents the input parsed so far
- Contents?
  - Symbols: terminals (and non-terminals)
  - Must also store previously seen *states*
    - the context of the current position
  - In fact, nonterminals unnecessary
    - include for readability

Stack

x + y $\bullet$  + z

T  
+  
T

# Parser Tables

Action table  
&  
GOTO table

## Action table

- given state  $s$  and **terminal  $a$**  tells parser loop what action (shift, reduce, accept, reject) to perform

## Goto table

- used when performing reduction; given a state  $s$  and **nonterminal  $X$**  says what state to transition to

# Parser Tables

Action table  
&  
GOTO table

**sN** push state  $N$  onto stack

**rR** reduce by rule  $R$

**gN** goto state  $N$

**a** accept

**error**

0	$S \rightarrow E\$$
1	$E \rightarrow T + E$
2	$E \rightarrow T$
3	$T \rightarrow \text{identifier}$

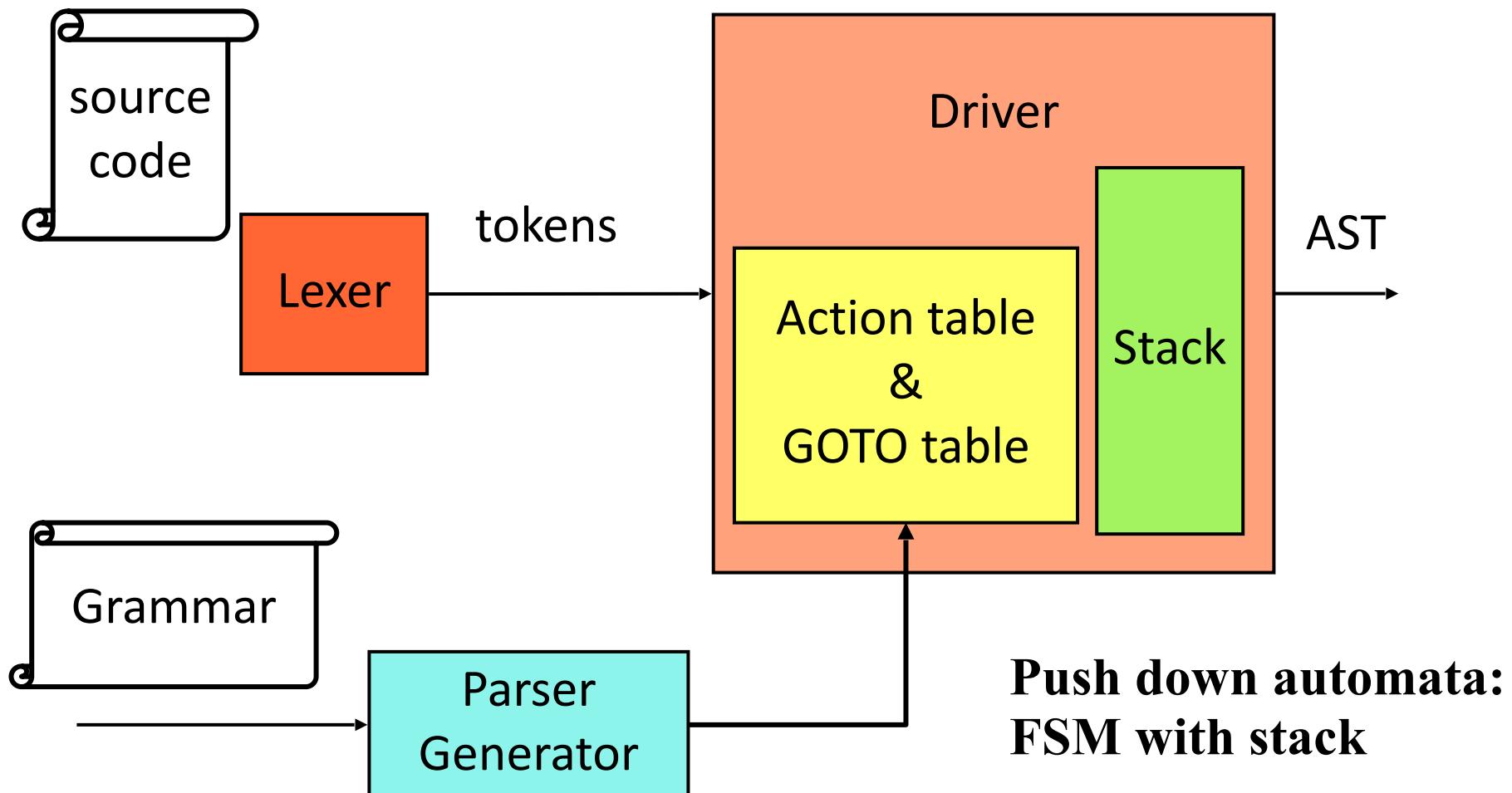
state	action			goto	
	<i>ident</i>	+	\$	E	T
0	<b>s3</b>			<b>g1</b>	<b>g2</b>
1			<b>a</b>		
2		<b>s4</b>	<b>r2</b>		
3			<b>r3</b>	<b>r3</b>	
4	<b>s3</b>			<b>g5</b>	<b>g2</b>
5			<b>r1</b>		

# Parser Loop Revisited

Driver

```
while(true)
    s = state on top of stack
    a = current input token
    if(action[s][a] == sN)                      shift
        push N
        read next input token
    else if(action[s][a] == rR)                  reduce
        pop rhs of rule R from stack
        x = lhs of rule R
        N = state on top of stack
        push goto[N][x]
    else if(action[s][a] == a)                  accept
        return success
    else                                         error
        return failure
```

# Table-driven LR(k) parsers



# The parser generator

Parser  
Generator

- Creates the **action** and **GOTO** tables.
- Creates the states
  - Each state indicates how much of a handle we have seen
  - each state is a set of *items*

# Items

- Items are used to identify handles.
- LR( $k$ ) items have the form:  
[ production-with-dot, lookahead]
- For example,  $A \rightarrow a X b$  has 4 LR(0) items
  - $[A \rightarrow \bullet a X b]$
  - $[A \rightarrow a \bullet X b]$
  - $[A \rightarrow a X \bullet b]$
  - $[A \rightarrow a X b \bullet]$

The  $\bullet$  indicates how much of the handle we have recognized.

# What LR(0) Items Mean

- $[X \rightarrow \bullet \alpha \beta \gamma]$   
input is consistent with  $X \rightarrow \alpha \beta \gamma$
- $[X \rightarrow \alpha \bullet \beta \gamma]$   
input is consistent with  $X \rightarrow \alpha \beta \gamma$  and we have already recognized  $\alpha$
- $[X \rightarrow \alpha \beta \bullet \gamma]$   
input is consistent with  $X \rightarrow \alpha \beta \gamma$  and we have already recognized  $\alpha \beta$
- $[X \rightarrow \alpha \beta \gamma \bullet]$   
input is consistent with  $X \rightarrow \alpha \beta \gamma$  and we can reduce to  $X$

# Generating the States

- Start with start production.
- In this case, “ $S \rightarrow E\$$ ”

$S \rightarrow \bullet E\$$

- 0  $S \rightarrow E\$$
- 1  $E \rightarrow T + E$
- 2  $E \rightarrow T$
- 3  $T \rightarrow \text{identifier}$

- Each state is consistent with what we have already shifted from the input and what is possible to reduce. So, what other items should be in this state?

# Completing a state

- For each item in a state, add in all other consistent items.

$S \rightarrow \bullet E\$$

$E \rightarrow \bullet T + E$

$E \rightarrow \bullet T$

$T \rightarrow \bullet identifier$

0  $S \rightarrow E\$$

1  $E \rightarrow T + E$

2  $E \rightarrow T$

3  $T \rightarrow identifier$

- This is called, taking the closure of the state.

# Closure\*

```
closure(state)
repeat
    foreach item  $A \rightarrow a \bullet X b$  in state
        foreach production  $X \rightarrow w$ 
            state.add( $X \rightarrow \bullet w$ )
until state does not change
return state
```

*Intuitively:*

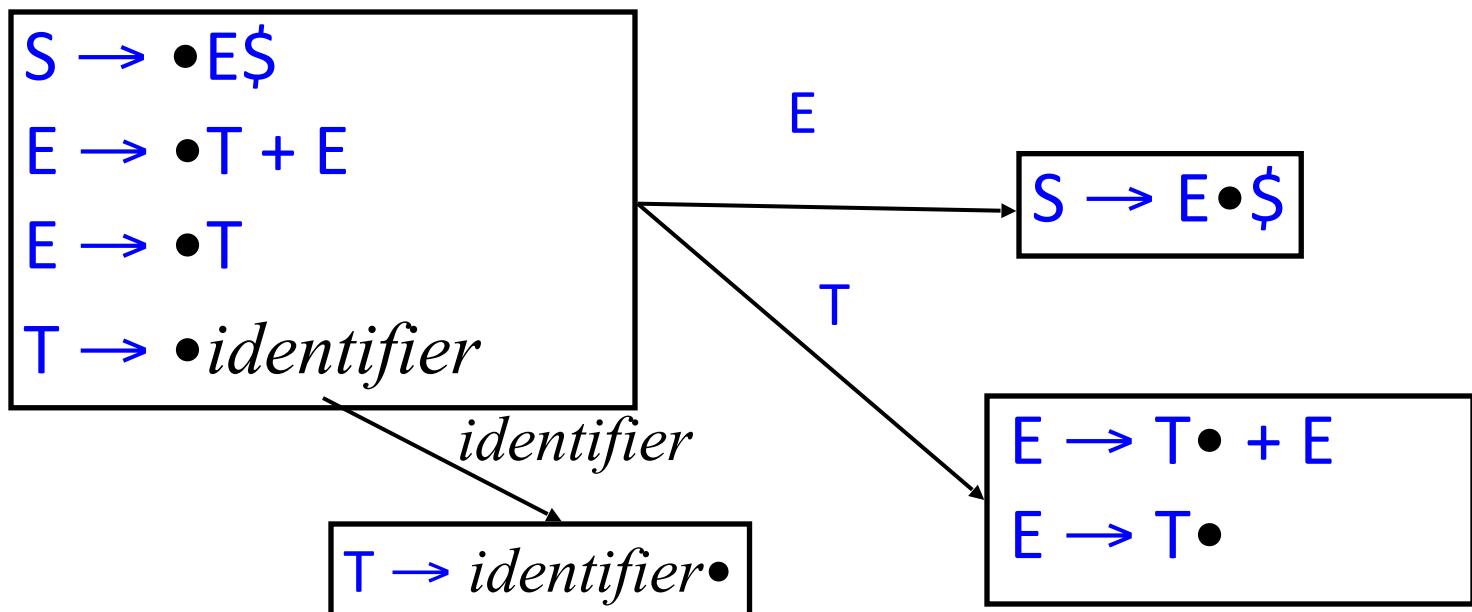
*Given a set of items, add all production rules that could produce the nonterminal(s) at the current position in each item*

\*: for LR(0) items

# What about the other states?

- How do we decide what the other states are?
- How do we decide what the transitions between states are?

0  $S \rightarrow E\$$   
1  $E \rightarrow T + E$   
2  $E \rightarrow T$   
3  $T \rightarrow \text{identifier}$



# Next(state, sym)

- Next function determines what state to goto based on current state and symbol being recognized.
- For Non-terminal, this is used to determine the GOTO table.
- For terminal, this is used to determine the shift action.

# Constructing states

```
initial_state = closure({start production})  
state_set.add(initial_state)  
state_queue.push(initial_state)  
  
while (!state_queue.empty())  
    s = state_queue.pop()  
    foreach item A → a•Xb in s  
        n = closure(next(s, X))  
        if (!state_set.contains(n))  
            state_set.add(n)  
            state_queue.push(n)
```

*A state is a set of LR(0) items*

*get “next” state*

# Closure\*

$\text{closure}(\{S \rightarrow \bullet E\$ \}) =$

$S \rightarrow \bullet E\$$

- 0  $S \rightarrow E\$$
- 1  $E \rightarrow T + E$
- 2  $E \rightarrow T$
- 3  $T \rightarrow \text{identifier}$

\*: for LR(0) items

# Closure\*

$\text{closure}(\{S \rightarrow \bullet E\$ \}) =$

$S \rightarrow \bullet E\$$

$E \rightarrow \bullet T + E$

$E \rightarrow \bullet T$

$T \rightarrow \bullet identifier$

**0**  $S \rightarrow E\$$

**1**  $E \rightarrow T + E$

**2**  $E \rightarrow T$

**3**  $T \rightarrow identifier$

\*: for LR(0) items

# Next

```
next(state, X)
    ret = empty
    foreach item A → a•xb in state
        ret.add(A → aX•b)
    return ret
```

- 0  $S \rightarrow E\$$
- 1  $E \rightarrow T + E$
- 2  $E \rightarrow T$

3  $T \rightarrow \text{identifier}$

initial:

$S \rightarrow \bullet E\$$

$E \rightarrow \bullet T + E$

$E \rightarrow \bullet T$

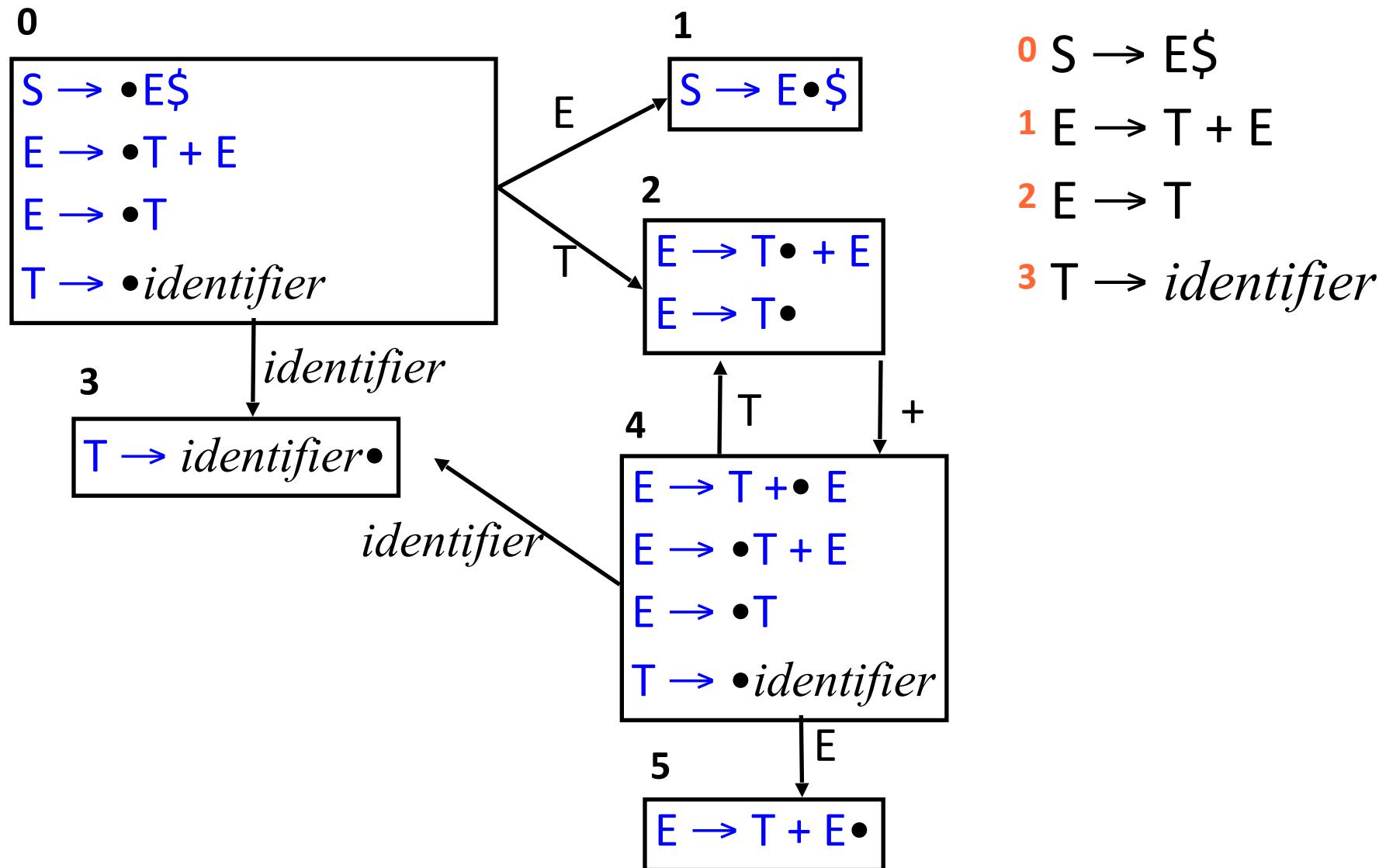
$T \rightarrow \bullet \text{identifier}$

next(initial, E)

next(initial, T)

next(initial, identifier)

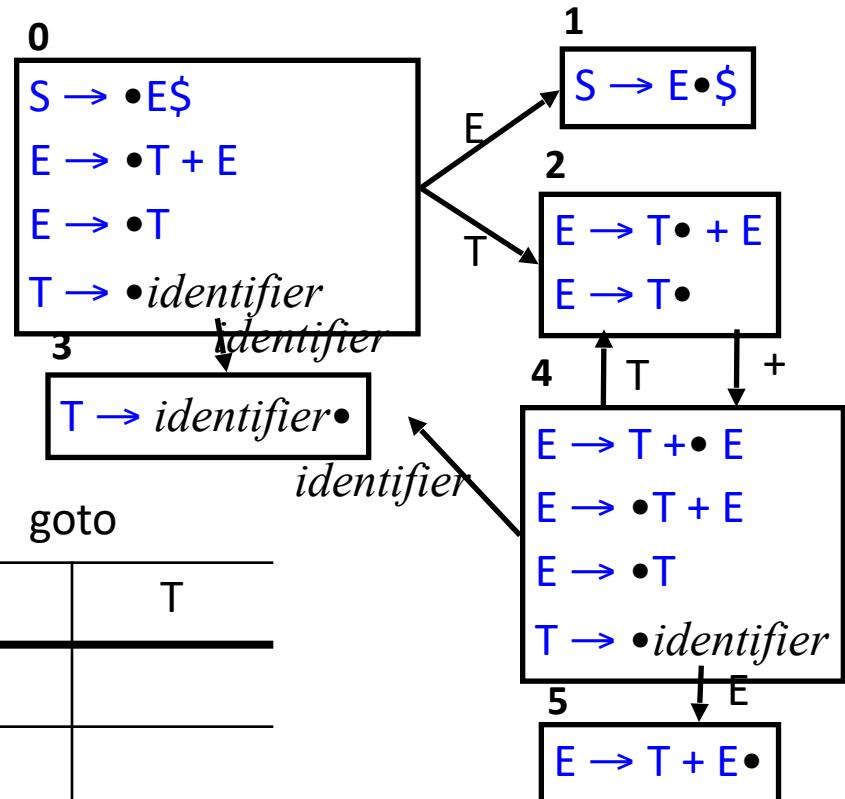
# Example



# Parse Tables for LR(0) parser

What can we fill out?

state	action			goto	
	<i>ident</i>	+	\$	E	T
0					
1					
2					
3					
4					
5					



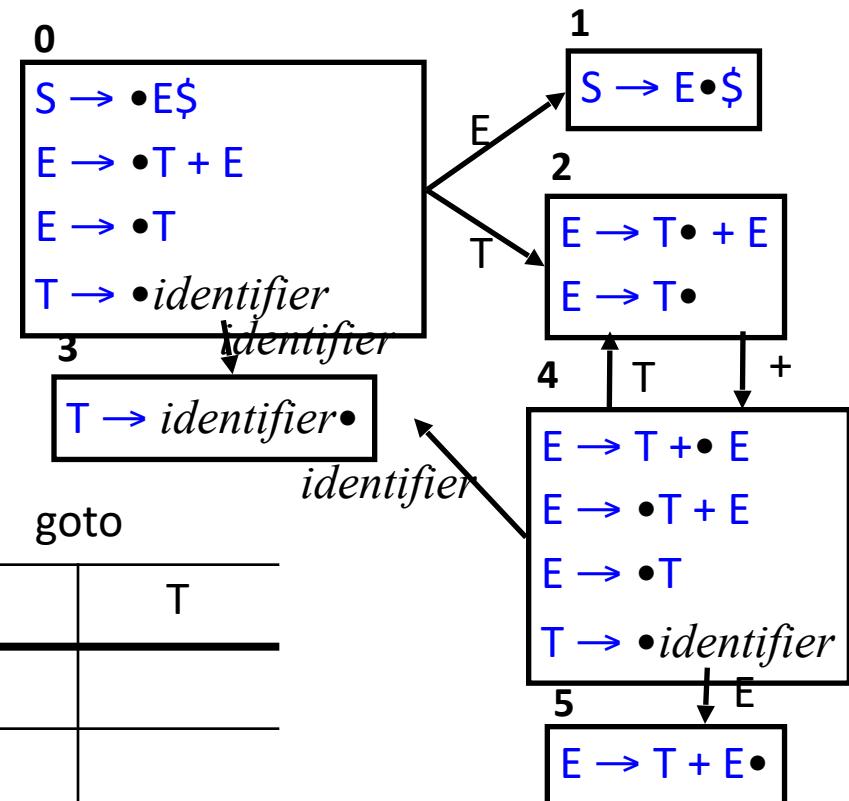
- 0  $S \rightarrow E\$$
- 1  $E \rightarrow T + E$
- 2  $E \rightarrow T$
- 3  $T \rightarrow \text{identifier}$

# Parse Tables for LR(0) parser

shift

transition on terminal

state	action			goto	
	<i>ident</i>	+	\$	E	T
0	s3				
1					
2		s4			
3					
4	s3				
5					



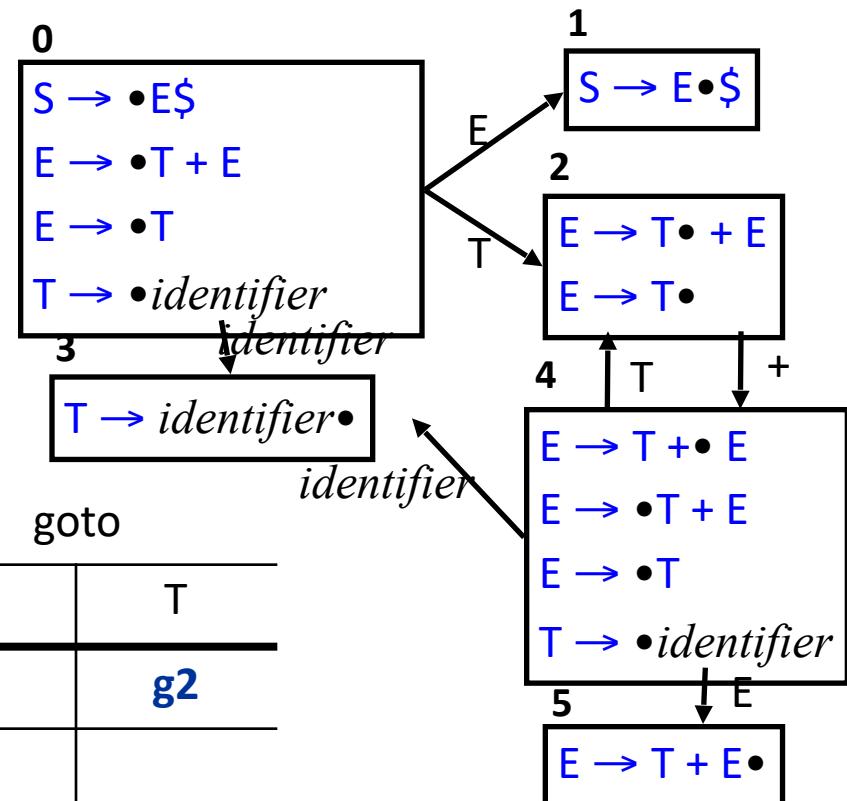
- 0  $S \rightarrow E\$$
- 1  $E \rightarrow T + E$
- 2  $E \rightarrow T$
- 3  $T \rightarrow \text{identifier}$

# Parse Tables for LR(0) parser

goto

transition on nonterminal

state	action			goto	
	<i>ident</i>	+	\$	E	T
0	s3			g1	g2
1					
2		s4			
3					
4	s3			g5	g2
5					

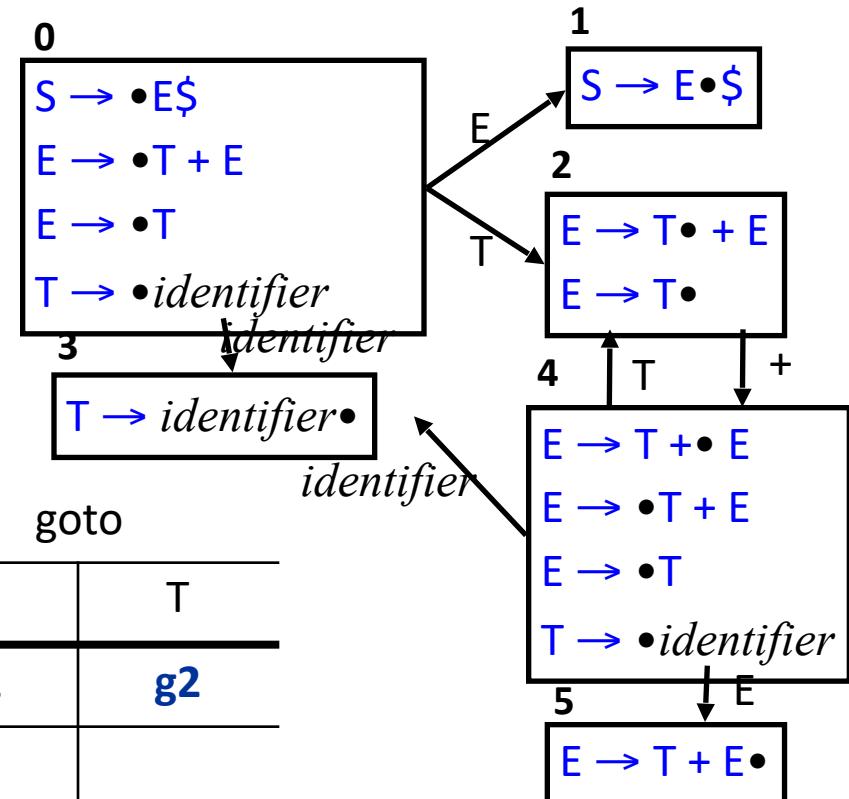


- 0  $S \rightarrow E \$$
- 1  $E \rightarrow T + E$
- 2  $E \rightarrow T$
- 3  $T \rightarrow \text{identifier}$

# Parse Tables for LR(0) parser

accept  
about to shift \$

state	action			goto	
	<i>ident</i>	+	\$	E	T
0	s3			g1	g2
1			a		
2		s4			
3					
4	s3			g5	g2
5					



- 0  $S \rightarrow E\$$
- 1  $E \rightarrow T + E$
- 2  $E \rightarrow T$
- 3  $T \rightarrow identifier$

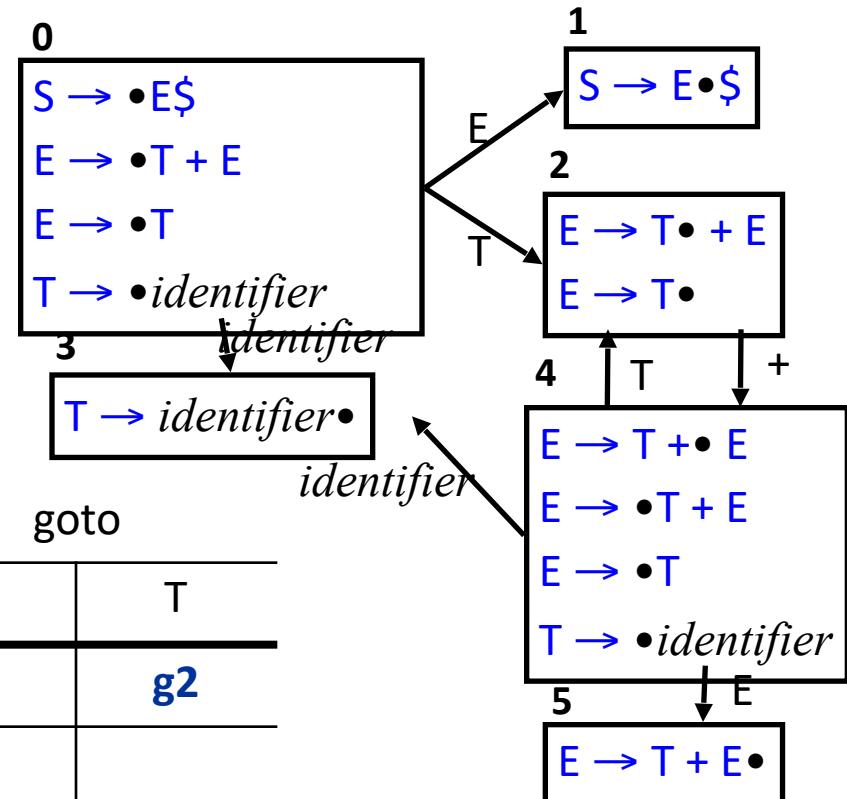
# Parse Tables for LR(0) parser

**reduce**

item has dot at end

$A \rightarrow W\bullet$

state	action			goto	
	<i>ident</i>	+	\$	E	T
0	s3			g1	g2
1			a		
2		s4			
3					
4	s3			g5	g2
5					



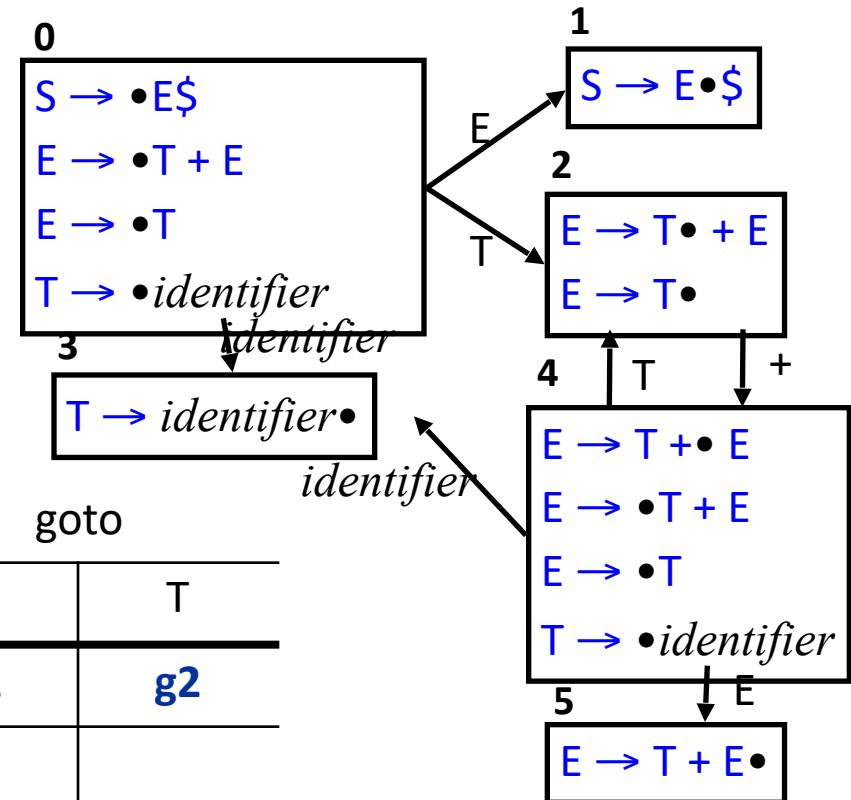
- 0  $S \rightarrow E \$$
- 1  $E \rightarrow T + E$
- 2  $E \rightarrow T$
- 3  $T \rightarrow \text{identifier}$

# LR(0)

No lookahead

reduce state for *all* nonterminals

state	action			goto	
	<i>ident</i>	+	\$	E	T
0	s3			g1	g2
1			a		
2	r2	r2/s4	r2		
3	r3	r3	r3		
4	s3			g5	g2
5	r1	r1	r1		



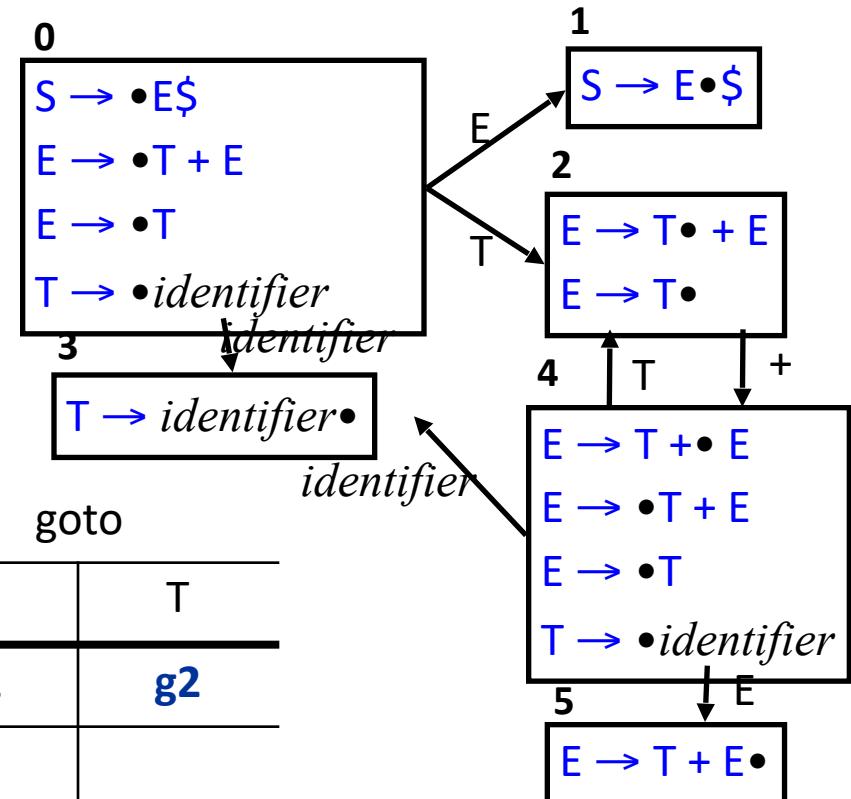
- 0  $S \rightarrow E\$$
- 1  $E \rightarrow T + E$
- 2  $E \rightarrow T$
- 3  $T \rightarrow identifier$

# LR(0)

## shift/reduce conflict

need to be pickier about  
when we reduce

state	action			goto	
	<i>ident</i>	+	\$	E	T
0	s3			g1	g2
1			a		
2	r2	r2/s4	r2		
3	r3	r3	r3		
4	s3			g5	g2
5	r1	r1	r1		

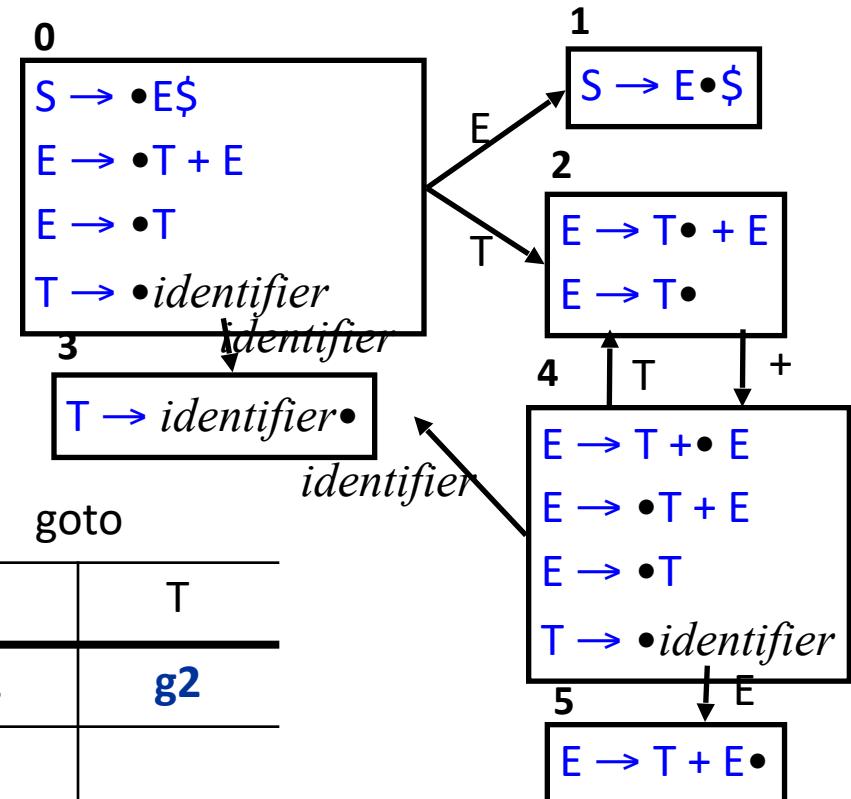


- 0  $S \rightarrow E \$$
- 1  $E \rightarrow T + E$
- 2  $E \rightarrow T$
- 3  $T \rightarrow \text{identifier}$

# SLR - Simple LR

Only reduce in position  $(s, a)$   
by rule  $R: A \rightarrow w$  if  $a$  is in the  
*follow set* of  $A$

state	action			goto	
	<i>ident</i>	+	\$	E	T
0	s3			g1	g2
1			a		
2		s4			
3					
4	s3			g5	g2
5					



- 0  $S \rightarrow E \$$
- 1  $E \rightarrow T + E$
- 2  $E \rightarrow T$
- 3  $T \rightarrow \text{identifier}$

# Reminder: Follow sets

## **FOLLOW(X)**

set of terminals that can appear immediately after the nonterminal X in some sentential form

I.e.,  $t \in \text{FOLLOW}(X)$  iff  $S \Rightarrow^* \alpha X t \beta$  for some  $\alpha$  and  $\beta$

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

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

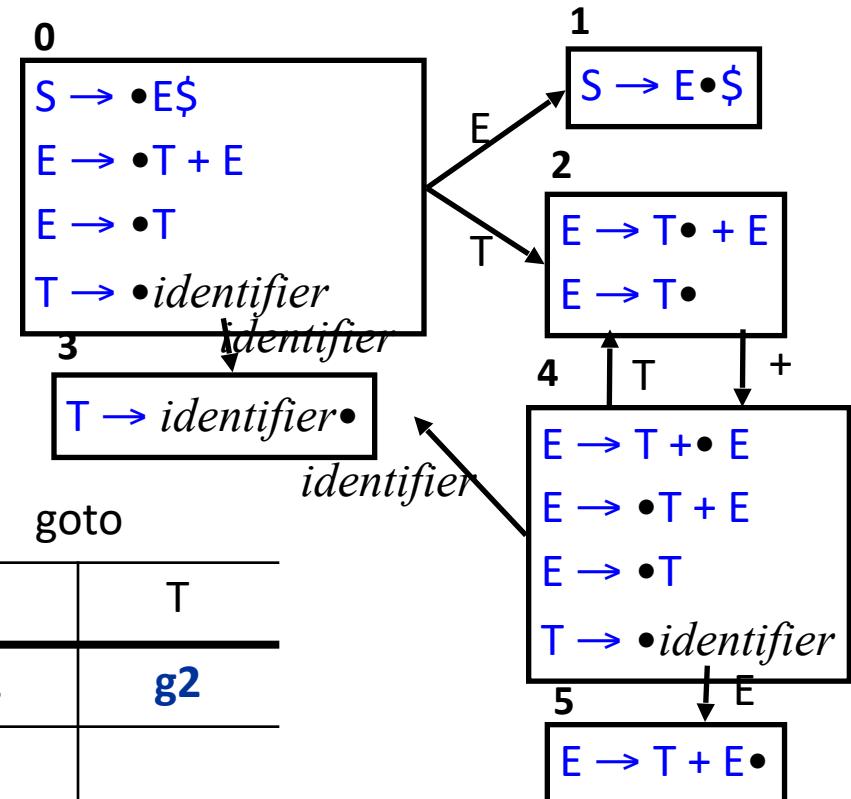
- 0  $S \rightarrow E \$$
- 1  $E \rightarrow T + E$
- 2  $E \rightarrow T$
- 3  $T \rightarrow \text{identifier}$

# SLR - Reduce using follow sets

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

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

state	action			goto	
	<i>ident</i>	+	\$	E	T
0	s3			g1	g2
1			a		
2		s4	r2		
3		r3	r3		
4	s3			g5	g2
5			r1		



- 0  $S \rightarrow E \$$
- 1  $E \rightarrow T + E$
- 2  $E \rightarrow T$
- 3  $T \rightarrow \text{identifier}$

# SLR Limitations

- SLR uses LR(0) item sets
- Can remove some (but not all) shift/reduce conflicts using follow set
- Doesn't fix all problems

# Problem with SLR

- Reduce on ALL terminals in FOLLOW set

S	→	L = R
		R
L	→	* R
		id
R	→	L

2

S → L • = R
R → L •

- $\text{FOLLOW}(R) = \text{FOLLOW}(L)$
- But, we should never reduce  $R \rightarrow L$  on '='  
i.e.,  $R=...$  is not a viable prefix for a right sentential form
- Thus, there should be no reduction in state 2
- How can we solve this?

# LR(1) Items

- An LR(1) item is an LR(0) item combined with a single terminal (*the lookahead*)
- $[X \rightarrow \alpha \bullet \beta, a]$  Means
  - $\alpha$  is at top of stack
  - Input string is derivable from  $\beta a$
- In other words, when we reduce  $X \rightarrow \alpha\beta$ ,  $a$  must be the lookahead symbol.
- Or, Only put ‘reduce by  $X \rightarrow \alpha\beta$ ’ in **action [s,a]**
- Can construct states as before, but have to modify closure

# What LR(1) Items Mean

- $[X \rightarrow \bullet \alpha \beta \gamma, a]$   
input is consistent with  $X \rightarrow \alpha \beta \gamma$
- $[X \rightarrow \alpha \bullet \beta \gamma, a]$   
input is consistent with  $X \rightarrow \alpha \beta \gamma$  and we have already recognized  $\alpha$
- $[X \rightarrow \alpha \beta \bullet \gamma, a]$   
input is consistent with  $X \rightarrow \alpha \beta \gamma$  and we have already recognized  $\alpha \beta$
- $[X \rightarrow \alpha \beta \gamma \bullet, a]$   
input is consistent with  $X \rightarrow \alpha \beta \gamma$  and if lookahead symbol is  $a$ , then we can reduce to  $X$

# LR(1) Closure

```
closure(state)
    repeat
        foreach item [ $A \rightarrow \alpha \bullet X\beta$ , t] in state
            foreach production  $X \rightarrow w$ 
                and each terminal  $t'$  in FIRST( $\beta t$ )
                    state.add([ $X \rightarrow \bullet w$ ,  $t'$ ])
    until state does not change
    return state
```

# Closure

$\text{closure}(\{S \rightarrow \bullet E\$, ?\}) =$

$S \rightarrow \bullet E\$, \quad ?$

**0**  $S \rightarrow E\$$

**1**  $E \rightarrow L = R$

**2**  $E \rightarrow R$

**3**  $L \rightarrow id$

**4**  $L \rightarrow *R$

**5**  $R \rightarrow L$

$\text{FIRST}( \$ ) = \{ \$ \}$

# Closure

$\text{closure}(\{S \rightarrow \bullet E\$, ?\}) =$

$S \rightarrow \bullet E\$, \quad ?$   
 $E \rightarrow \bullet L = R, \quad \$$   
 $E \rightarrow \bullet R, \quad \$$

**0**  $S \rightarrow E\$$   
**1**  $E \rightarrow L = R$   
**2**  $E \rightarrow R$   
**3**  $L \rightarrow id$   
**4**  $L \rightarrow *R$   
**5**  $R \rightarrow L$

FIRST(\$)  
FIRST(=R\$)

{\$}  
{=}

# Closure

$\text{closure}(\{S \rightarrow \bullet E\$, ?\}) =$

$S \rightarrow \bullet E\$,$       ?

$E \rightarrow \bullet L = R,$       \\$

$E \rightarrow \bullet R,$       \\$

$L \rightarrow \bullet id,$       =

$L \rightarrow \bullet *R,$       =

$0 S \rightarrow E\$$

$1 E \rightarrow L = R$

$2 E \rightarrow R$

$3 L \rightarrow id$

$4 L \rightarrow *R$

$5 R \rightarrow L$

FIRST(\\$)      {\\$}

FIRST(=R\\$)      {=}

# Closure

$\text{closure}(\{S \rightarrow \bullet E\$, ?\}) =$

$S \rightarrow \bullet E\$,$       ?

$E \rightarrow \bullet L = R,$       \\$

$E \rightarrow \bullet R,$       \\$

$L \rightarrow \bullet id,$       =

$L \rightarrow \bullet *R,$       =

$R \rightarrow \bullet L,$       \\$

$0 S \rightarrow E\$$

$1 E \rightarrow L = R$

$2 E \rightarrow R$

$3 L \rightarrow id$

$4 L \rightarrow *R$

$5 R \rightarrow L$

FIRST(\\$)      {\\$}

FIRST(=R\\$)      {=}

# Closure

$\text{closure}(\{S \rightarrow \bullet E\$, ?\}) =$

$S \rightarrow \bullet E\$,$       ?

$E \rightarrow \bullet L = R,$       \\$

$E \rightarrow \bullet R,$       \\$

$L \rightarrow \bullet id,$       =

$L \rightarrow \bullet *R,$       =

$R \rightarrow \bullet L,$       \\$

$L \rightarrow \bullet id,$       \\$

$L \rightarrow \bullet *R,$       \\$

$0 S \rightarrow E\$$

$1 E \rightarrow L = R$

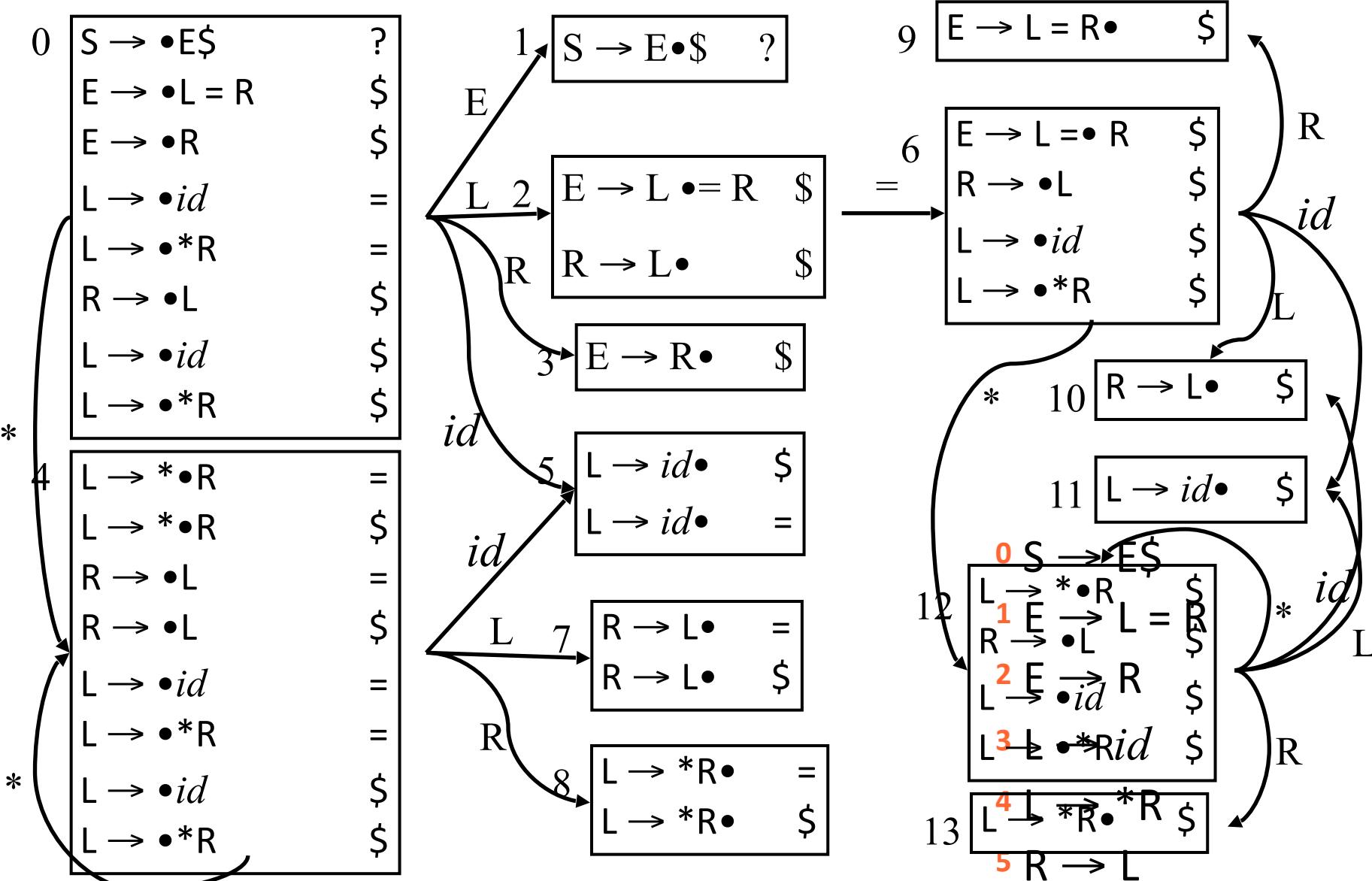
$2 E \rightarrow R$

$3 L \rightarrow id$

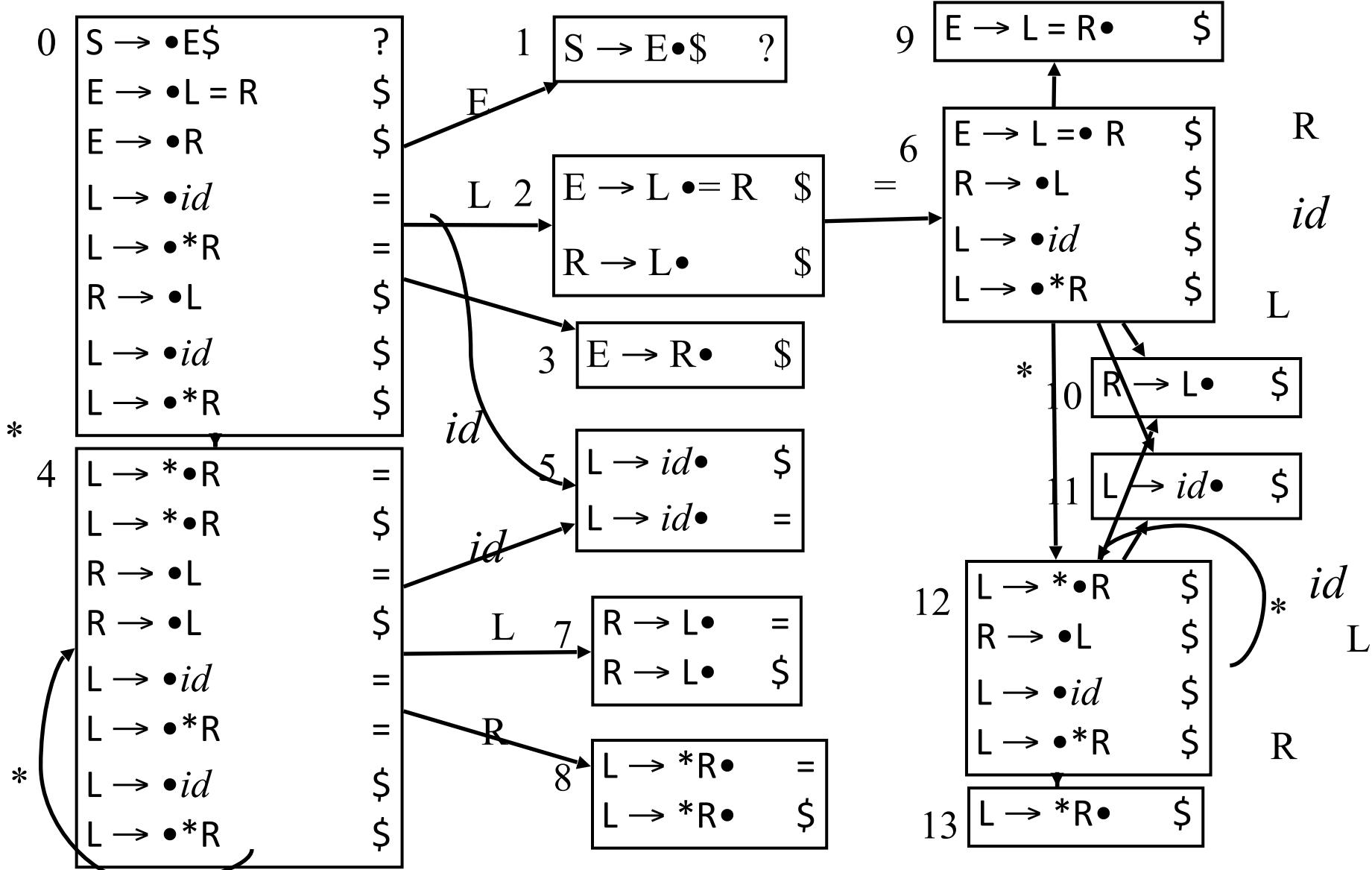
$4 L \rightarrow *R$

$5 R \rightarrow L$

# LR(1) Example



# LR(1) Example



# Parsing Table for LR(1)

- 14 states versus 10 LR(0) states
- In general, the number of states (and therefore size of the parsing table) is much larger with LR(1) items

# LALR: Lookahead LR

- More powerful than SLR
- Given LR(1) states, merge states that are identical except for lookaheads
- End up with ~ same size table as SLR
- Can this introduce conflicts?



# LALR

- Can generate parse table without constructing LR(1) item sets
  - construct LR(0) item sets
  - compute *lookahead* sets
    - more precise than follow sets
- LALR is used by most parser generators (e.g., bison, lalrpop, ocamllyacc,...)

# Recap

- LR(0)      not very useful
- SLR          uses follow sets to reduce
- LALR        uses lookahead sets
- LR(1)        uses full lookahead context

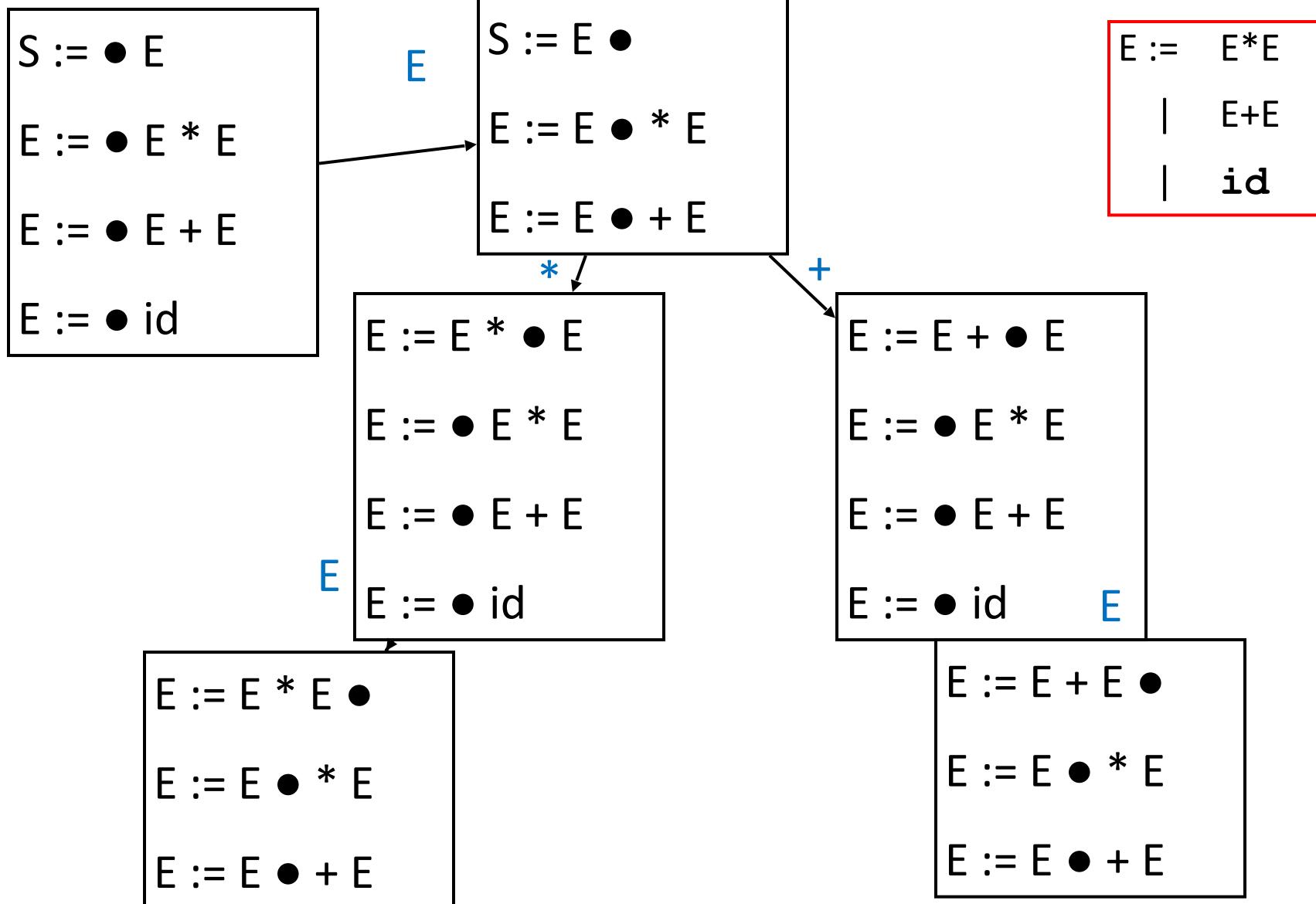
# Power of shift-reduce parsers

- There are unambiguous grammars which cannot be parsed with shift-reduce parsers.
- Such grammars can have
  - shift/reduce conflicts
  - reduce/reduce conflicts
- There grammars are not LR( $k$ )
- But, we can often choose shift or reduce to recognize what want.

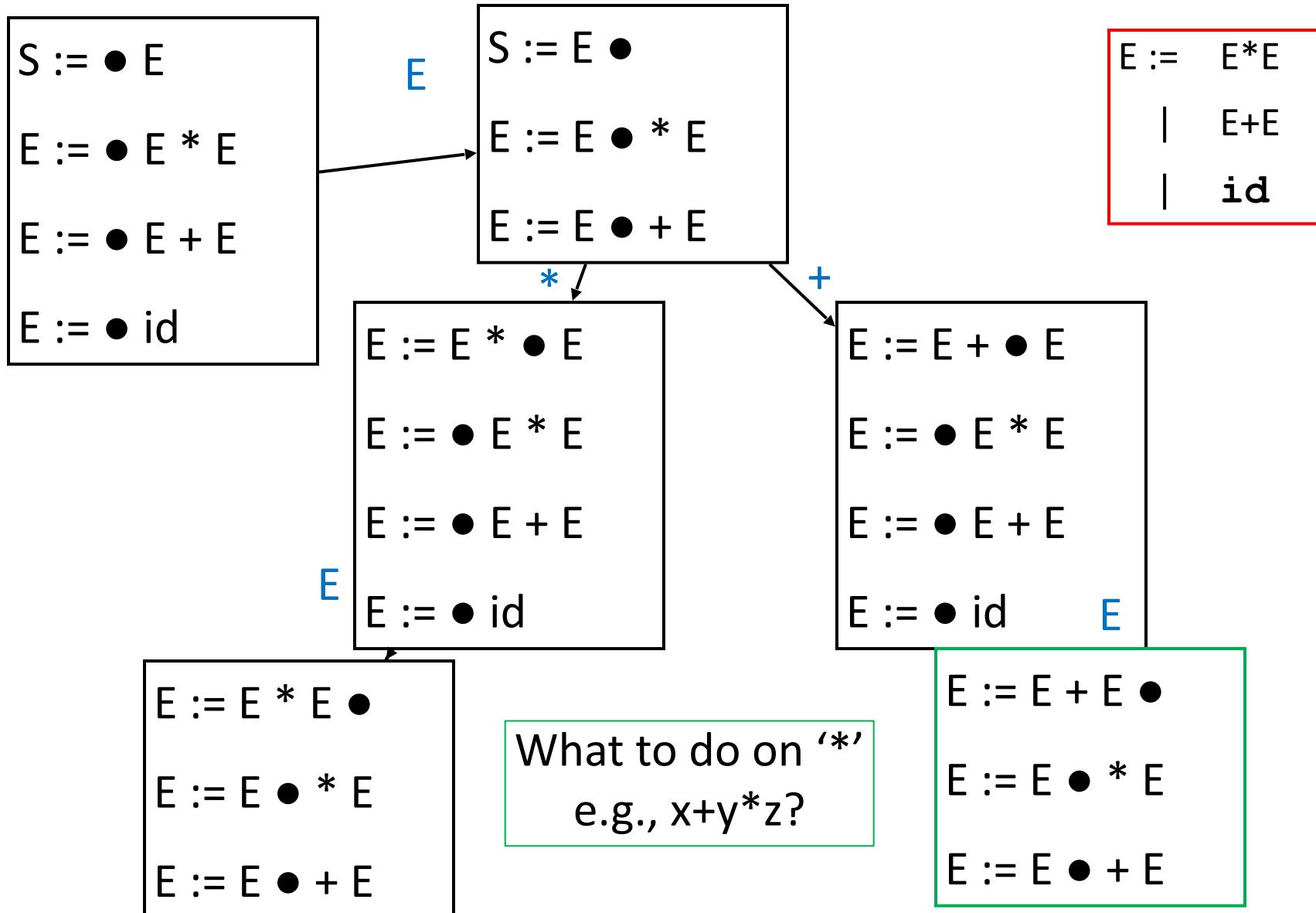
# Handling precedence and associativity

- Straightforwardly written grammars for infix will have shift/reduce conflicts
- Associativity:  $a + b \bullet + c$ 
  - prefer reduce => left-associative  $(a + b) + c$
  - prefer shift => right-associative  $a + (b + c)$
- Precedence:  $a * b \bullet + c$ 
  - prefer reduce => higher prec  $(a * b) + c$
  - prefer shift => lower prec  $a * (b + c)$

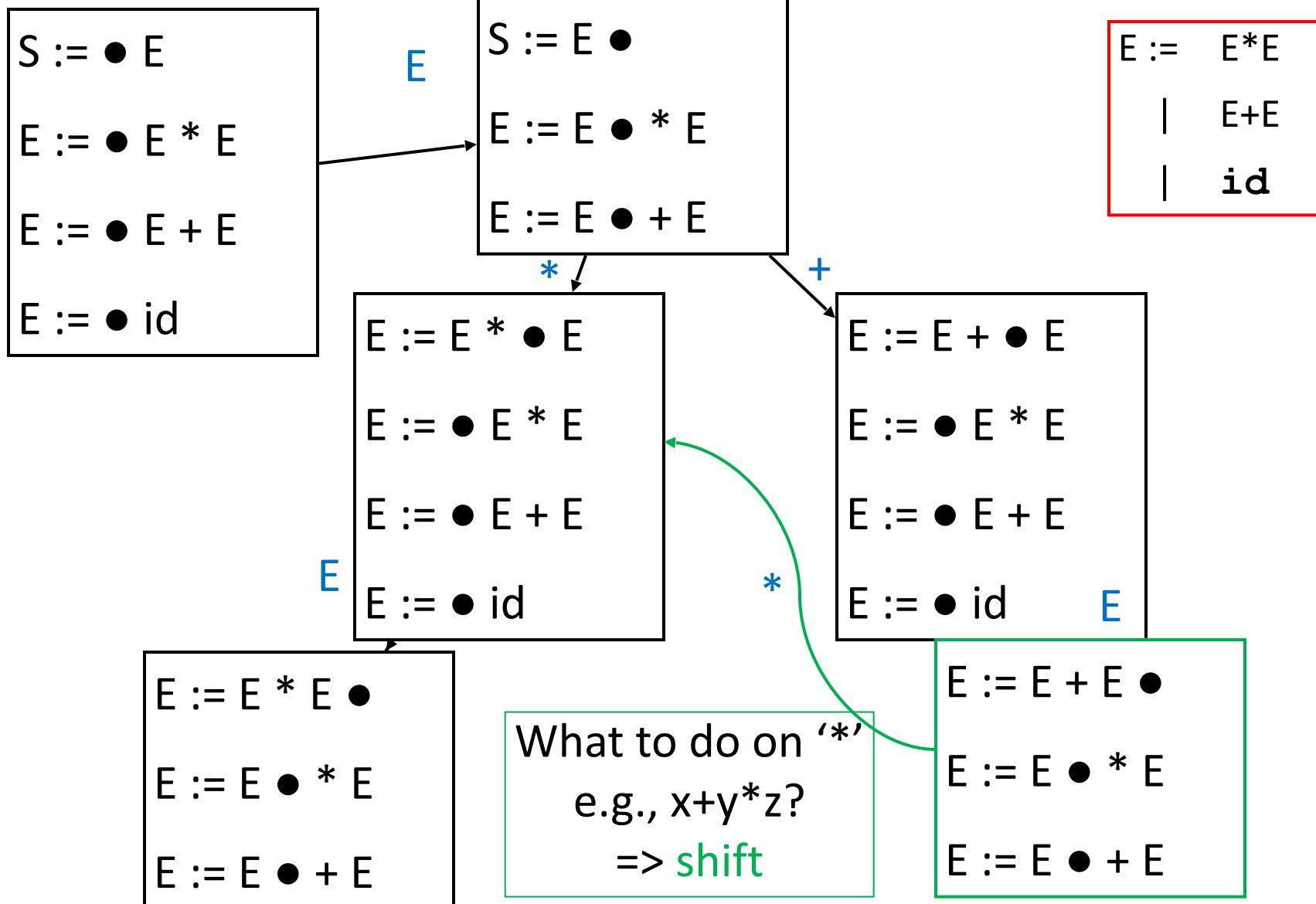
# Expression Grammars & Precedence



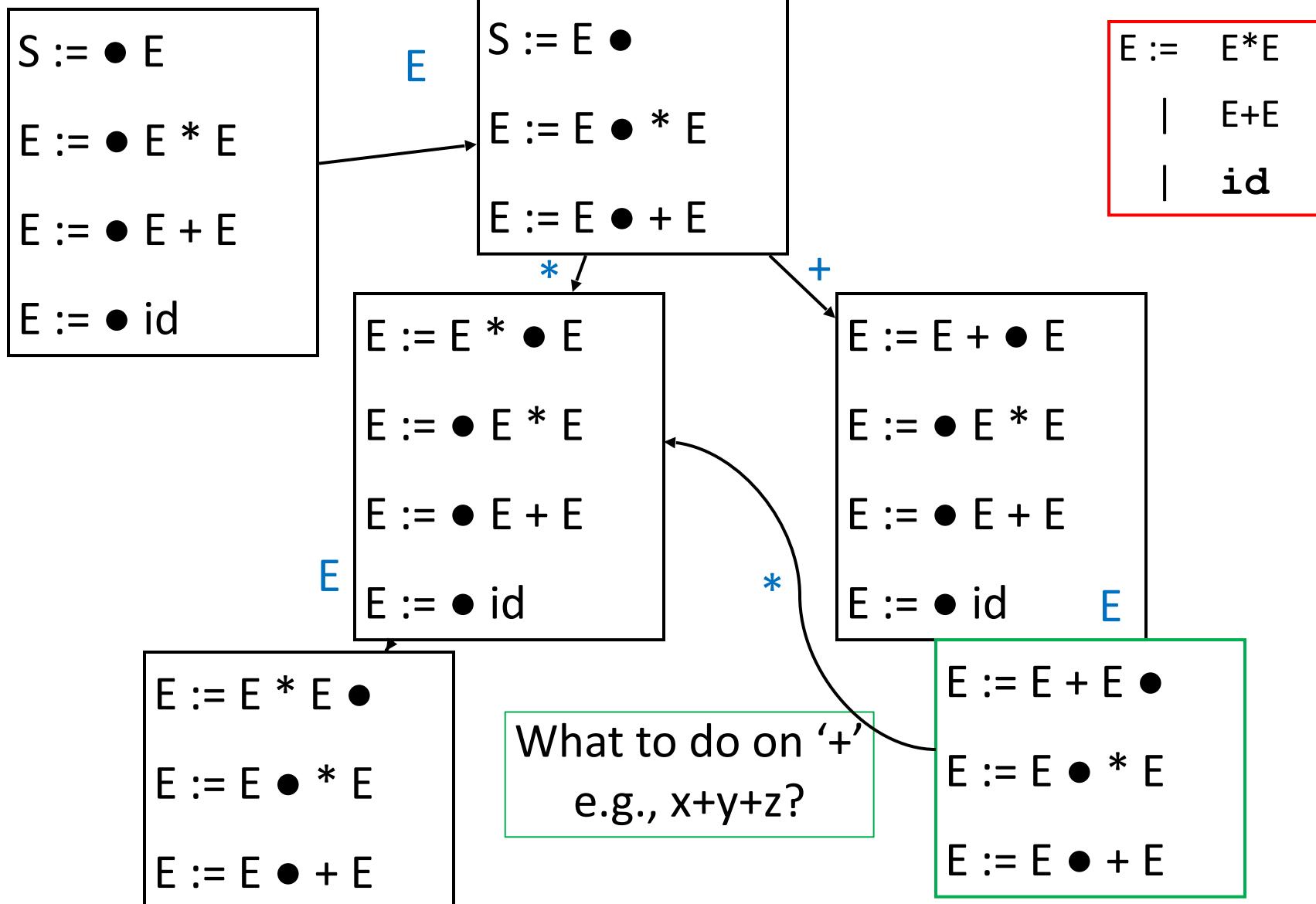
# Expression Grammars & Precedence



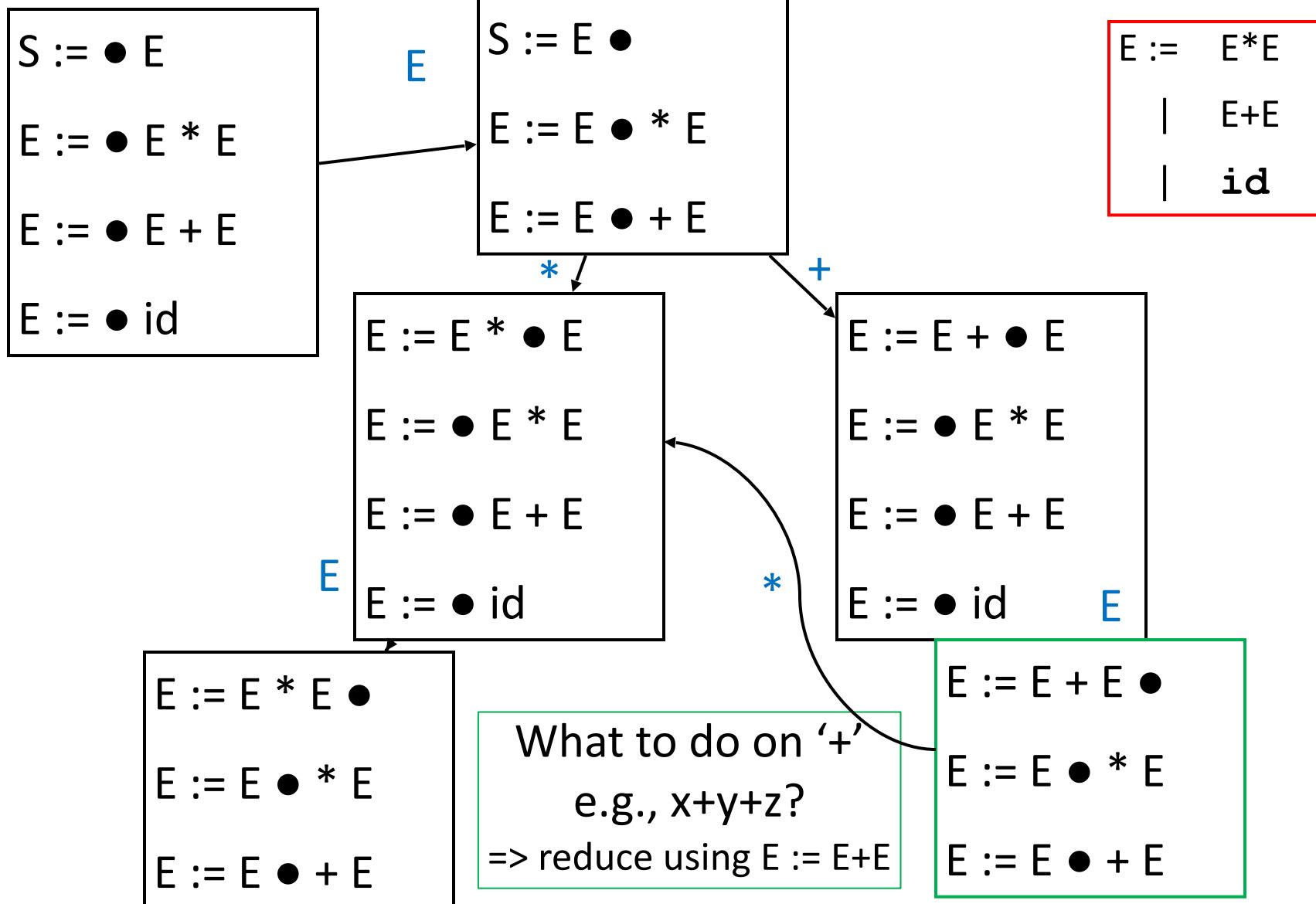
# Expression Grammars & Precedence



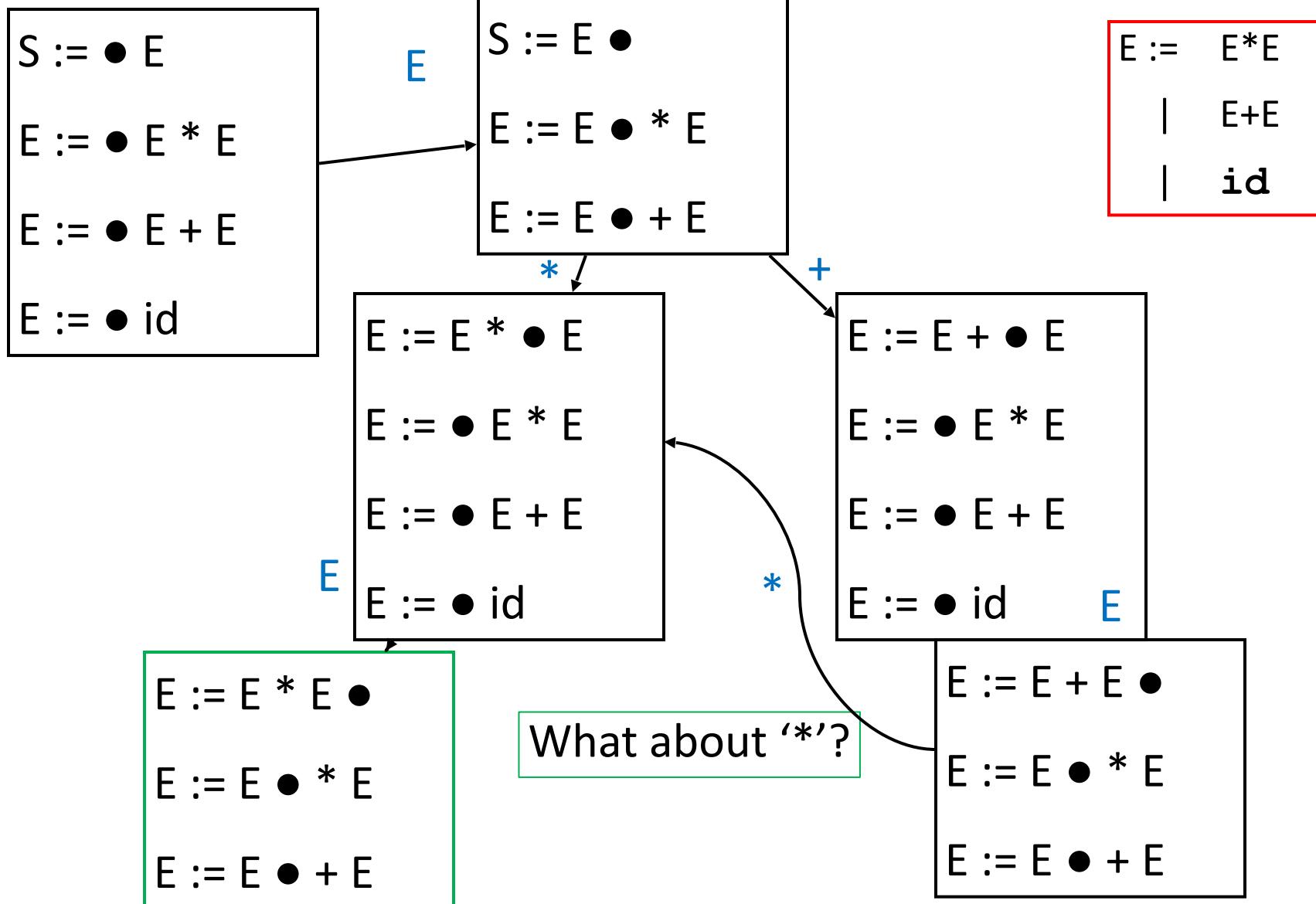
# Expression Grammars & Precedence



# Expression Grammars & Precedence



# Expression Grammars & Precedence



# Bison

- Precedence and Associativity declarations
- Precedence derived from order of directives: from lowest to highest
- Associativity from %left, %right, %nonassoc
- Can be attached to rules as well (This can solve the dangling if-else problem)
- Use output of generator showing items and transitions to debug s/r and r/r errors

# Dangling Else

$S := \text{if } E \text{ then } S$

|  $\text{if } E \text{ then } S \text{ else } S$

| **other**

- We can be in the following state:

$\dots \text{if } E \text{ then } S \quad \text{else} \dots \$$

- What do we do?
  - shift the **else** (hoping to reduce by second rule), or
  - reduce by first rule