

---

# CMPSC 160

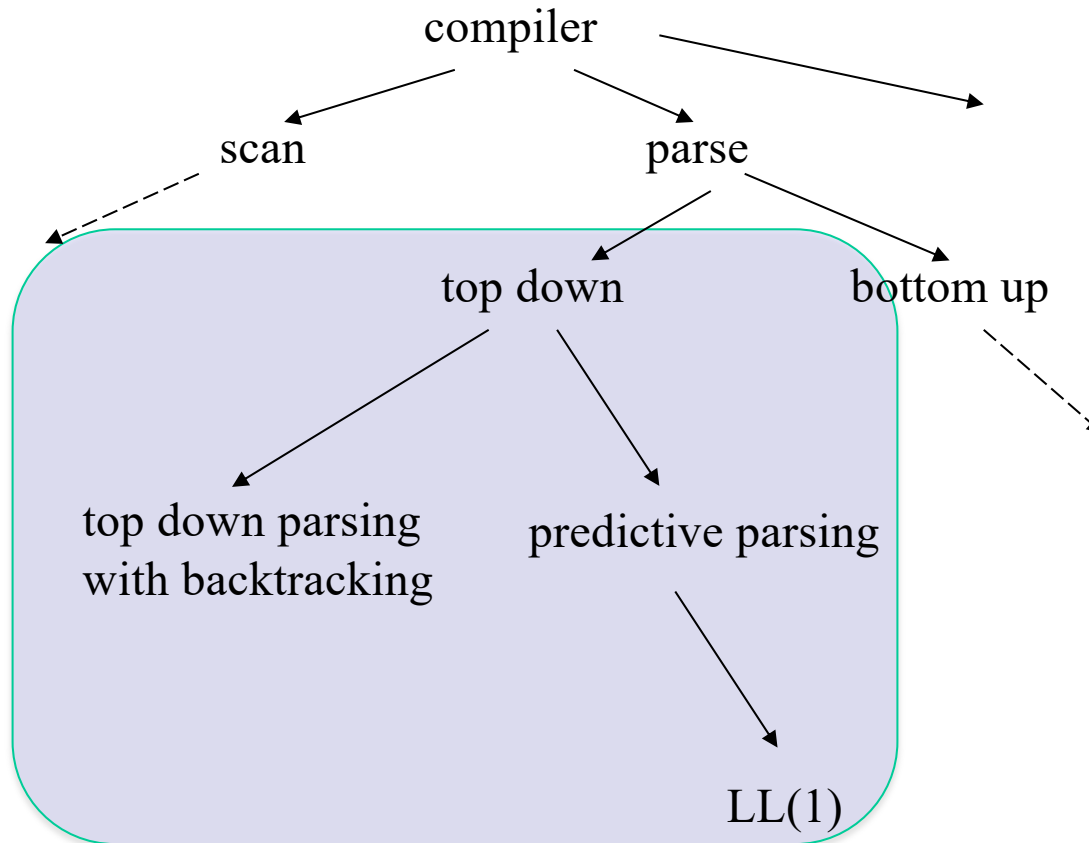
## Translation of Programming Languages

Lecture 5: Top-down parsing +  
Introduction to Bottom-up Parsing

---

# Where are we in the process?

---



# Predictive Parsing with LL(1) Grammars

---

## Basic idea

Given  $A \rightarrow \alpha \mid \beta$ , the parser should be able to choose between  $\alpha$  and  $\beta$  based on peeking at the next token in the stream

**LL(1)**: Left-to-right scan of the input, Leftmost derivation, 1-token look-ahead

A grammar  $G$  is LL(1) if for each set of its productions

$A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ :

FIRST( $\alpha_1$ ), FIRST( $\alpha_2$ ), ..., FIRST( $\alpha_n$ ), are all pair-wise disjoint

If  $\alpha_j \Rightarrow^* \varepsilon$ , then FIRST( $\alpha_j$ )  $\cap$  FOLLOW( $A$ ) =  $\emptyset$  for all  $1 \leq i \leq n, i \neq j$

In other words, LL(1) grammars

- during leftmost derivation, productions are uniquely predictable with a **one** token lookahead
-

# General Predictive Parsing

---

How much look-ahead is needed?

- In general, an arbitrarily large amount
- Use the Cocke-Younger, Kasami, or Earley's algorithm
  - Complexity is  $O(|x|^3)$  where  $x$  is the input string

Fortunately,

- Large subclasses of context free grammars can be parsed efficiently with limited look-ahead
  - Linear complexity,  $O(|x|)$  where  $x$  is the input string
- Most programming language constructs fall in those subclasses

Among the interesting subclasses is the *LL(1)* grammars.

---

# Recursive Descent Parsing

1	$S$	$\rightarrow$	if $E$ then $S$ else $S$
2			begin $S L$
3			print $E$
4	$L$	$\rightarrow$	end
5			; $S L$
6	$E$	$\rightarrow$	num = num

```
void S() {
    switch(lookahead) {
        case IF: match(IF); E(); match(THEN); S();
                 match(ELSE); S(); break;
        case BEGIN: match(BEGIN); S(); L(); break;
        case PRINT: match(PRINT); E(); break;
        default: error();
    }
}

void E() { match(NUM); match(EQ); match(NUM); }
```

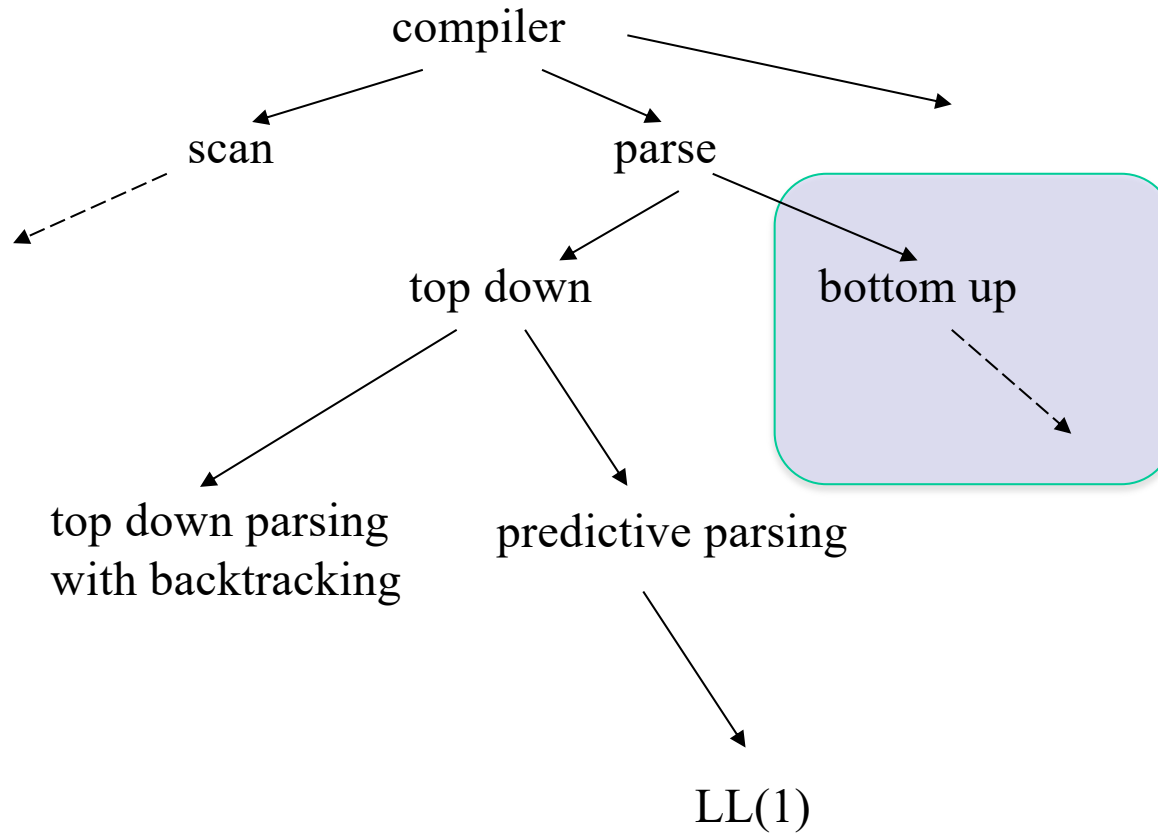
```
void match(int token) {
    if (lookahead==token)
        lookahead=getNextToken();
    else
        error();
}
```

```
void L() {
    switch(lookahead) {
        case END: match(END); break;
        case SEMI: match(SEMI); S();
                  L(); break;
        default: error();
    }
}
```

```
void main() {
    lookahead=getNextToken();
    S();
    match EOF;
}
```

# Where are we in the process?

---

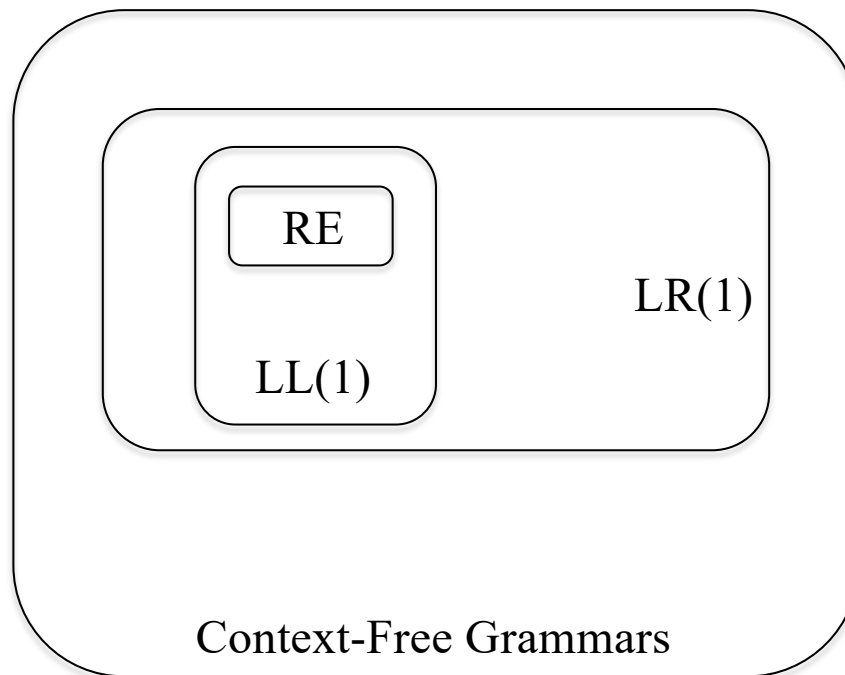


# Parsing Techniques

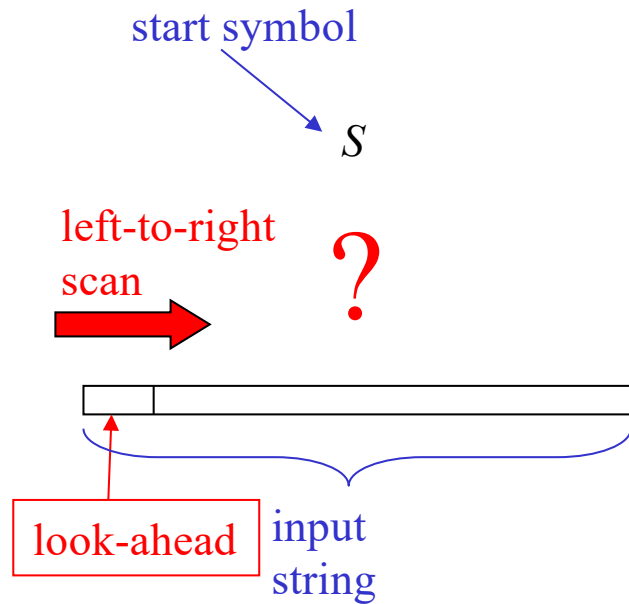
---

*Top-down parsers*    (*LL(1), recursive descent parsers*)

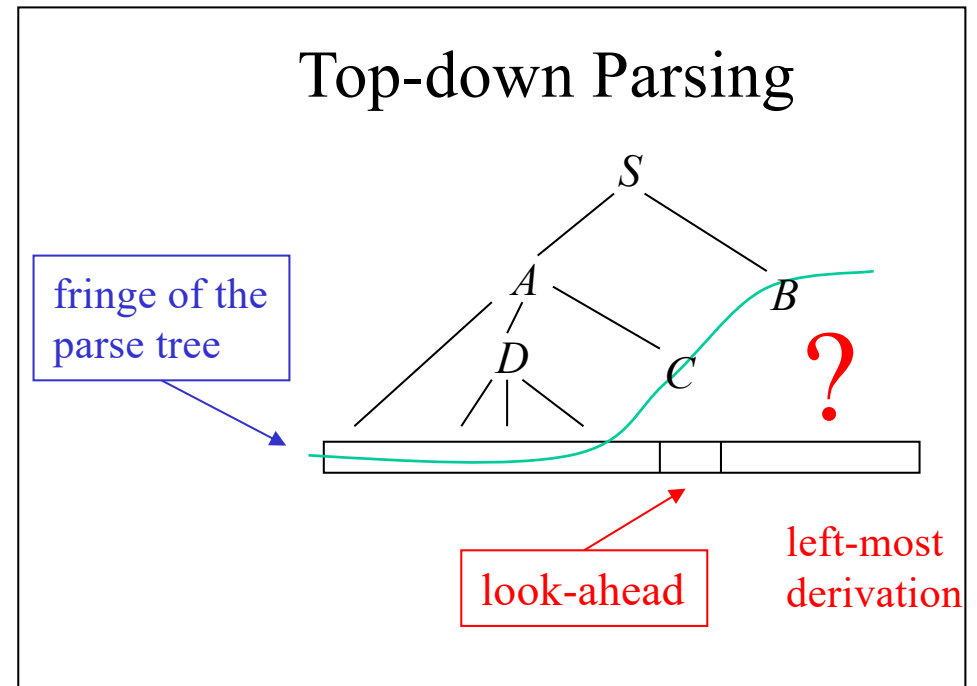
*Bottom-up parsers*    (*LR(1), shift-reduce parsers*)



# Top-Down Parsing v.s. Bottom-Up Parsing



When we start with a predictive top-down parser, the look-ahead symbol we read from our input string **MUST** fully specify the parse tree from  $S$  to the input symbol. In the example, **we have to know that  $S \rightarrow AB$  before we even see any of  $B$**

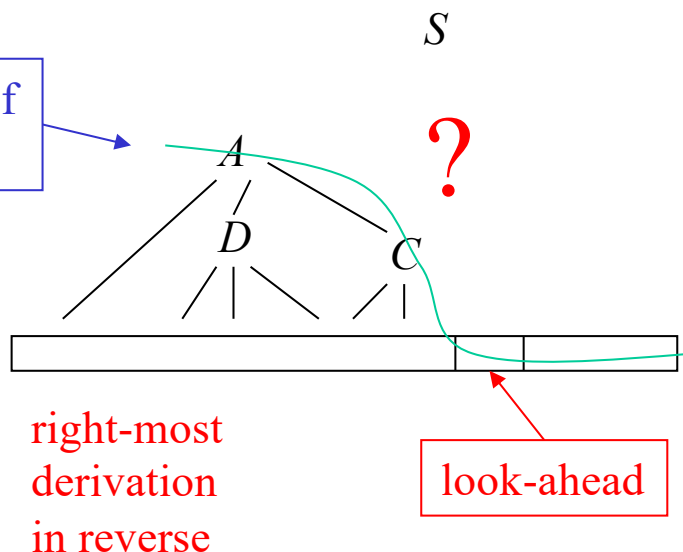




# Top-Down Parsing v.s. Bottom-Up Parsing

## Bottom-up Parsing

upper fringe of the parse tree



In a bottom-up parser, we can **delay this decision** because we only need to build the tree up above the part of the input string we have examined so far.

In the graphical example on the left, you can see that even though we are at the same point in the input string, the production  $S \rightarrow AB$  has not been specified yet. **This delayed decision allows us to parse more grammars than predictive top-down parsing (LL).**

As a nice side effect, bottom-up parsing allows us to handle **left-recursive** grammars without modification

# Bottom-Up Parsers

---

- Given a stream of tokens  $w$ , reduce it to the start symbol.

$$\begin{array}{l} \hline E \longrightarrow E+T \\ E \longrightarrow T \\ T \longrightarrow id \\ \hline \end{array}$$

Parsing for **id + id**:

id + id  
T + id  
E + id  
E + T  
E

**Reduction  $\equiv$  Derivation in Reverse (Rightmost Derivation)**

**LR: Left-to-right scanning, Rightmost Derivation**

---

# Shift-Reduce Parsing: An Example

$$\begin{array}{l}
 \hline
 E \longrightarrow E+T \\
 E \longrightarrow T \\
 T \longrightarrow id \\
 \hline
 \end{array}$$

STACK	INPUT STREAM	ACTION
\$	id + id \$	<b>shift</b>
\$ id	+ id \$	<b>reduce by <math>T \longrightarrow id</math></b>
\$ T	+ id \$	reduce by $E \longrightarrow T$
\$ E	+ id \$	shift
\$ E +	id \$	shift
\$ E + id	\$	reduce by $T \longrightarrow id$
\$ E + T	\$	reduce by $E \longrightarrow E+T$
\$ E	\$	<b>ACCEPT</b>

Stack: sentential form + rightmost derivation in reverse

Input stream: from left to right

Two key actions: shift or reduce

Key question: when to reduce? which production rule to choose?

# Handles

---

- **Handle:** A structure that furnishes a means to perform reductions

$$\begin{array}{c} \hline E \longrightarrow E+T \\ E \longrightarrow T \\ T \longrightarrow \text{id} \\ \hline \end{array}$$

- Handles are substrings of sentential forms:
    - A substring that matches the right-hand side of a production
    - Reduction using that rule can lead to the start symbol
  - **Handle Pruning:** The process of discovering a handle & reducing it to the appropriate left-hand side is called handle pruning.
-

# Shift-Reduce Parsing

---

## Key operations (based on the stack and input stream):

- **Shift:** Construct leftmost handle on top of stack
- **Reduce:** Identify handle and replace by corresponding RHS
- **Accept:** Continue until string is reduced to start symbol and input token stream is empty
- **Error:** Signal parse error if no handle is found.

## Implementations:

- Stack to hold grammar symbols (corresponding to tokens seen thus far).
  - Stack is initially empty (denoted by \$).
  - Input stream of yet-to-be-seen tokens.
  - Start with the key operations above.
  - **Handles appear on top of stack.**
  - Parse is successful if stack contains only the start symbol when the input stream ends.
-

# Handle-pruning, Bottom-up Parsers

One implementation of a simple *shift-reduce parser*

```
push $
lookahead = get_next_token( )
repeat until (top of stack == start symbol and lookahead == $)
  if the top of the stack is a handle  $\alpha \rightarrow \beta$ 
    then /* reduce  $\beta$  to  $\alpha$  */
      pop  $|\beta|$  symbols off the stack
      push  $\alpha$  onto the stack
  else if (token  $\neq$  $)
    then /* shift */
      push lookahead
      lookahead = get_next_token( )
```

How do errors show up?

- failure to find a handle
- hitting \$ and needing to shift (final else clause)

Either generates an error

“pop  $|\beta|$  symbols off the stack” means that if we have some production  $A \rightarrow BCx$  then we need to pop  $B$ ,  $C$ , and  $x$  off the stack (  $|b| = 3$  in this case ) before we push  $A$

# Keys for Shift-Reduce Parsing

---

- Identify a handle in string.
  - Top of stack is the rightmost end of the handle.
- What is the leftmost end?
  - If there are **multiple productions** with the handle on the RHS, which one to choose?

# Existence of a Unique Handler for Rightmost Derivation

---

Insight

*If  $G$  is unambiguous, then every right-sentential form has a unique handle.*

**If we can find those handles, we can build a derivation!**

Sketch of Proof:

- 1  $G$  is unambiguous  $\Rightarrow$  rightmost derivation is unique (through a list of sentential forms  $\gamma_1, \dots, \gamma_n$ )
- 2  $\Rightarrow$  a unique production  $\alpha \rightarrow \beta$  applied to take sentential form  $\gamma_{i-1}$  to  $\gamma_i$
- 3  $\Rightarrow$  a unique position  $k$  at which  $\alpha \rightarrow \beta$  is applied
- 4  $\Rightarrow$  a unique handle  $\langle \alpha \rightarrow \beta, k \rangle$ : replacing  $\beta$  at  $k$  with  $\alpha$ .

This all follows from the definitions

---



# LR Parsers

---

- LR(k) parsers could parse a set of CFGs, for which the **handles can be deterministically indented with k lookahead of the input stream.**
    - The key part is recognizing the handles
  - **Wait a minute... recognizing... we studied something that recognizes strings: State machines (aka automata)!**
  - The big picture is that LR parsers use a state machine (to recognize handles) in coordination with a stack (to handle the recursive nature of grammars) to parse an input.
-

# A Simple Example of LR Parsing

- For input stream: aa

$S$	$\rightarrow$	$BC$
$B$	$\rightarrow$	$a$
$C$	$\rightarrow$	$a$

STACK	INPUT STREAM	ACTION
\$	a a \$	shift
\$ a	a \$	reduce by $B \rightarrow a$
\$ B	a \$	shift
\$ B a	\$	reduce by $C \rightarrow a$
\$ B C	\$	reduce by $S \rightarrow BC$
\$ S	\$	<b>ACCEPT</b>

Any insights about why we should select these handles respectively???

Can you turn these insights into general principles???

# A Simple Example of LR(0) Parsing

1	$S' \rightarrow S$
2	$S \rightarrow BC$
3	$B \rightarrow a$
4	$C \rightarrow a$

STACK	INPUT	STATE	ACTION
\$	a a \$	$S' \rightarrow \bullet S$ $S \rightarrow \bullet BC$ $B \rightarrow \bullet a$	shift
\$ a	a \$	$B \rightarrow a \bullet$	reduce by 3
\$ B	a \$	$S \rightarrow B \bullet C$ $C \rightarrow \bullet a$	shift
\$ B a	\$	$C \rightarrow a \bullet$	reduce by 4
\$ B C	\$	$S \rightarrow BC \bullet$	reduce by 2
\$ S	\$	$S' \rightarrow S \bullet$	<b>ACCEPT</b>

- **States:** Productions with “•” somewhere on the RHS.
  - Grammar symbols before the “•” are on stack
  - Grammar symbols after the “•” represent symbols in the input stream.
- **Closure:** What other items are “equivalent” to the given item?

Given an item  $A \rightarrow \alpha \bullet B \beta$ ,  $\text{closure}(A \rightarrow \alpha \bullet B \beta)$  is the smallest set that contains the item

$A \rightarrow \alpha \bullet B \beta$ , and every item in  $\text{closure}(B \rightarrow \bullet \gamma)$  for every production  $B \rightarrow \gamma \in \text{CFG}$

# Example: States of an LR(0) parser

$E' \rightarrow E$	$E \rightarrow T$
$E \rightarrow E+T$	$T \rightarrow id$

$l_0$	$= closure(\{E' \rightarrow \bullet E\})$	$E' \rightarrow \bullet E$ $E \rightarrow \bullet E+T$ $E \rightarrow \bullet T$ $T \rightarrow \bullet id$
$l_1$	$= goto(l_0, E)$	$E' \rightarrow E \bullet$ $E \rightarrow E \bullet +T$
$l_2$	$= goto(l_0, T)$	$E \rightarrow T \bullet$
$l_3$	$= goto(l_0, id)$	$T \rightarrow id \bullet$
$l_4$	$= goto(l_1, +)$	$E \rightarrow E+ \bullet T$ $T \rightarrow \bullet id$
$l_5$	$= goto(l_4, T)$	$E \rightarrow E+T \bullet$

- Instead of generating the states along with the inputs to parse, we could do a preprocessing and store all possible states (along with their transition) in a table before the parsing process.
- The parsing process becomes largely a table lookup process (just like the scanner part).

# States to Table

$l_0$	$= \text{closure}(\{E' \rightarrow \bullet E\})$	$E' \rightarrow \bullet E$ $E \rightarrow \bullet E+T$ $E \rightarrow \bullet T$ $T \rightarrow \bullet \text{id}$
$l_1$	$= \text{goto}(l_0, E)$	$E' \rightarrow E \bullet$ $E \rightarrow E \bullet +T$
$l_2$	$= \text{goto}(l_0, T)$	$E \rightarrow T \bullet$
$l_3$	$= \text{goto}(l_0, \text{id})$	$T \rightarrow \text{id} \bullet$
$l_4$	$= \text{goto}(l_1, +)$	$E \rightarrow E+ \bullet T$ $T \rightarrow \bullet \text{id}$
$l_5$	$= \text{goto}(l_4, T)$	$E \rightarrow E+T \bullet$

Action Table				Goto table	
	id	+	\$	E	T
0	S, 3			1	2
1		S, 4	A		
2	R3	R3	R3		
3	R4	R4	R4		
4	S, 3				5
5	R2	R2	R2		

There are in fact two different tables

The **ACTION table** tells you if you should shift, reduce, accept, or throw an error.

It additionally tells you how to update the state on a shift

The **GOTO table** tells you how to update the state on a reduce action

# Quiz: States of an LR(0) parser

---

example

0  $S' \rightarrow S \$$

1  $S \rightarrow V = E$

2  $S \rightarrow E$

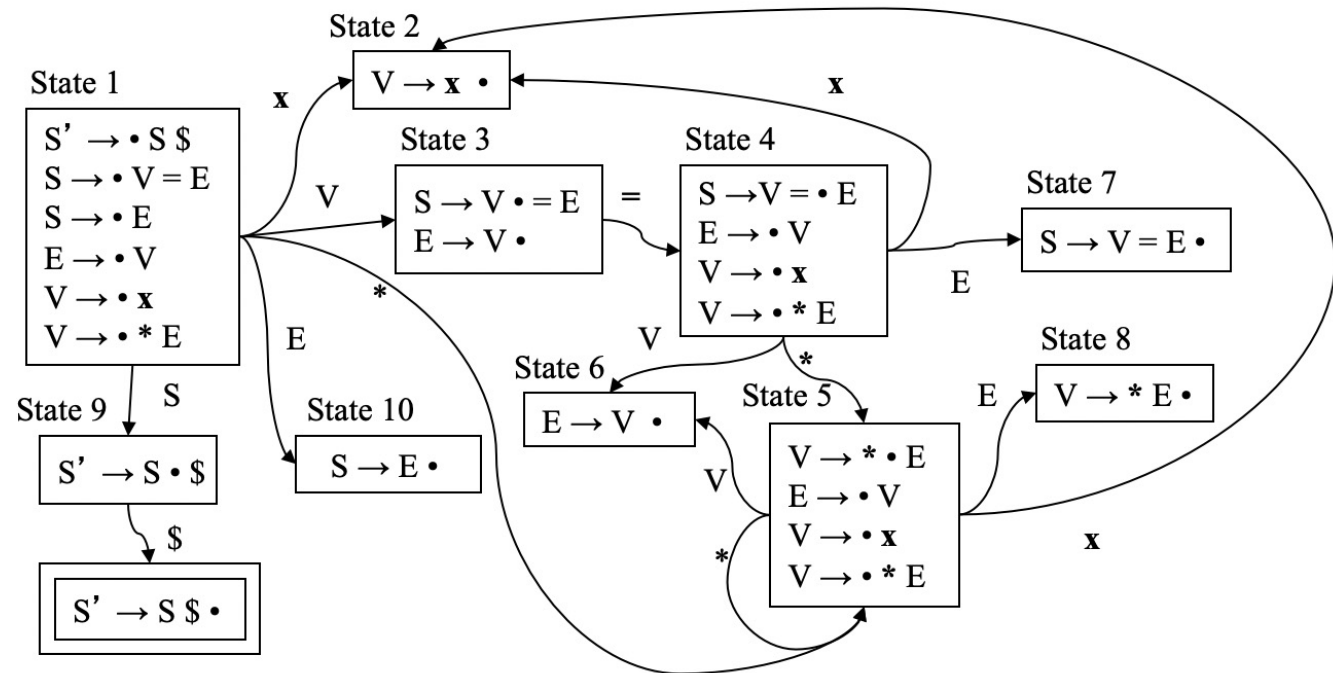
3  $E \rightarrow V$

4  $V \rightarrow x$

5  $V \rightarrow * E$

# Example: States of an LR(0) parser

example	
0	$S' \rightarrow S \$$
1	$S \rightarrow V = E$
2	$S \rightarrow E$
3	$E \rightarrow V$
4	$V \rightarrow x$
5	$V \rightarrow * E$



**Closure:** What other items are “equivalent” to the given item?

Given an item  $A \rightarrow \alpha \cdot B \beta$ ,  $\text{closure}(A \rightarrow \alpha \cdot B \beta)$  is the smallest set that contains the item  $A \rightarrow \alpha \cdot B \beta$ , and every item in  $\text{closure}(B \rightarrow \cdot \gamma)$  for every production  $B \rightarrow \gamma \in \text{CFG}$

# Things behind the automation: NFA to DFA

For example, to match  $[S \rightarrow \cdot V = E]$  we need to first match  $V$ . Let us connect the **NFA states together with  $\epsilon$ -transitions** whenever one state needs to make a “subroutine” call to another state.

example	
0	$S' \rightarrow S \$$
1	$S \rightarrow V = E$
2	$S \rightarrow E$
3	$E \rightarrow V$
4	$V \rightarrow x$
5	$V \rightarrow * E$

