

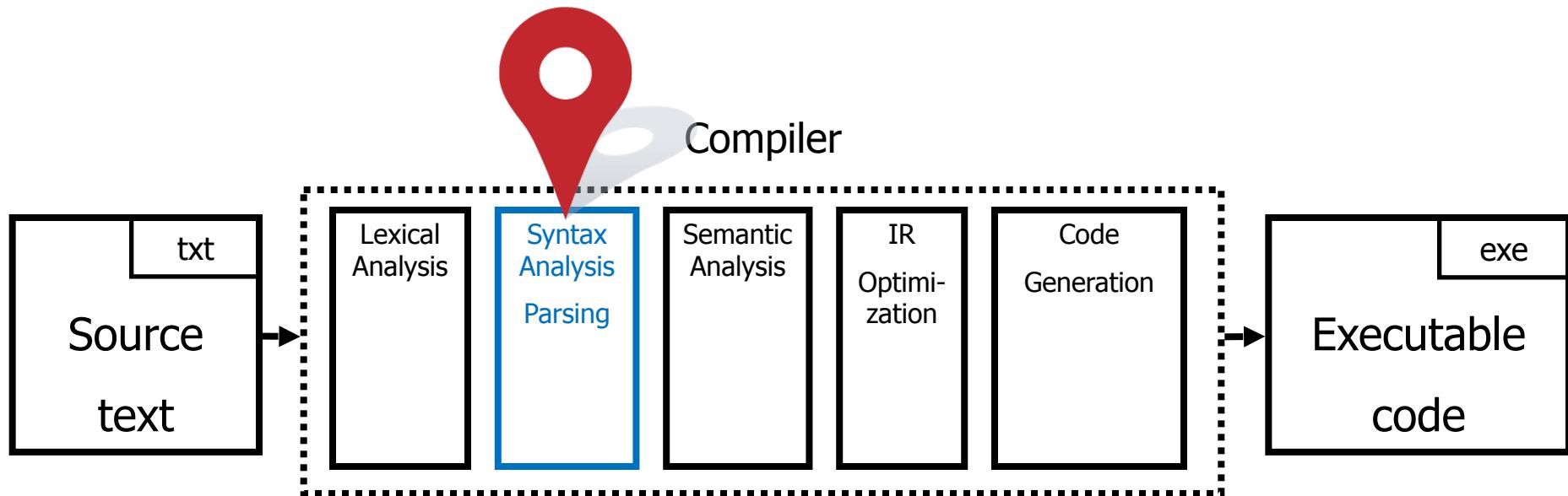
# THEORY OF COMPIRATION

## LECTURE 02

SYNTAX  
ANALYSIS

TOP-DOWN PARSING

# You are here



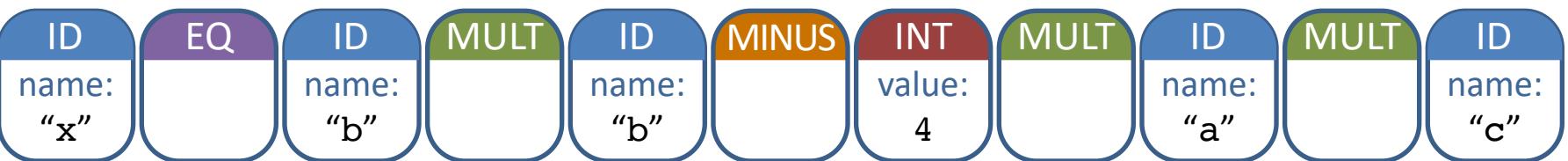
# Previously: from characters to tokens

txt

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

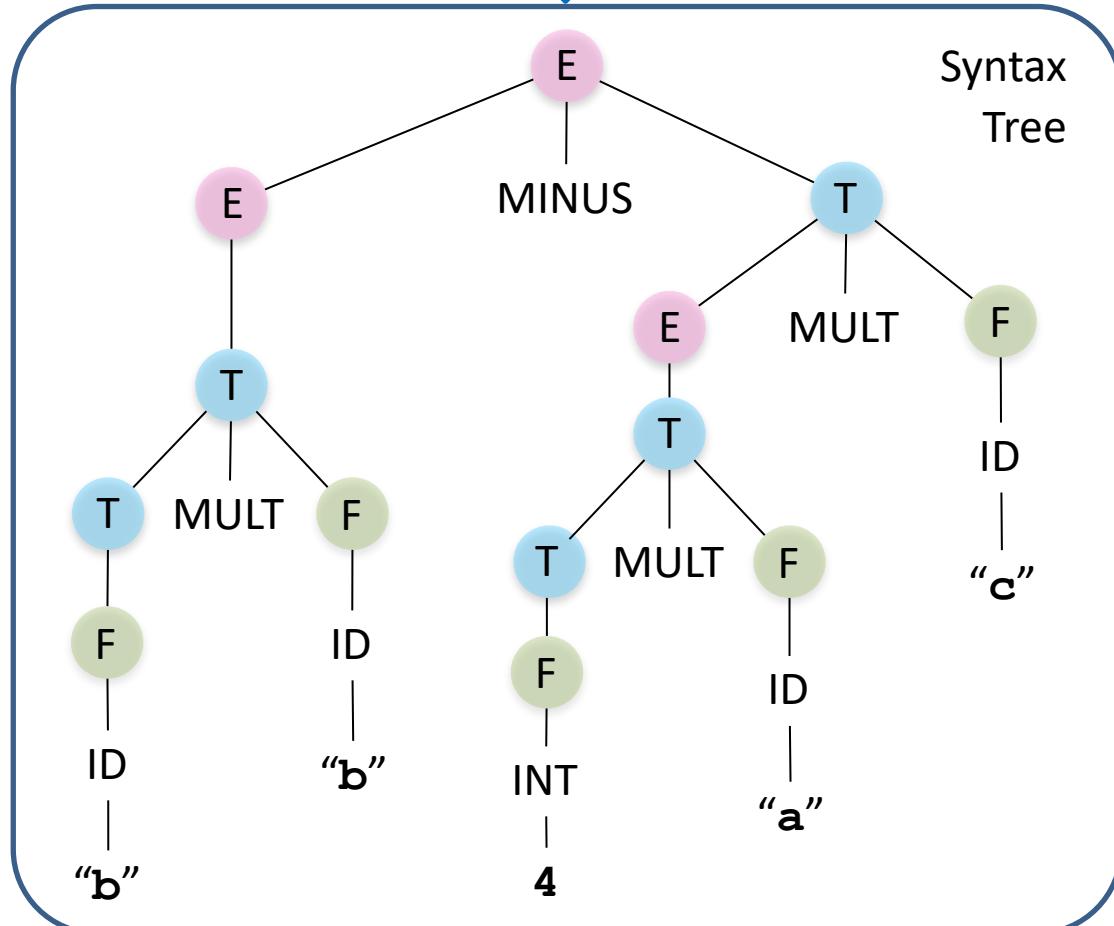


Token Stream



# Today: from tokens to syntax tree

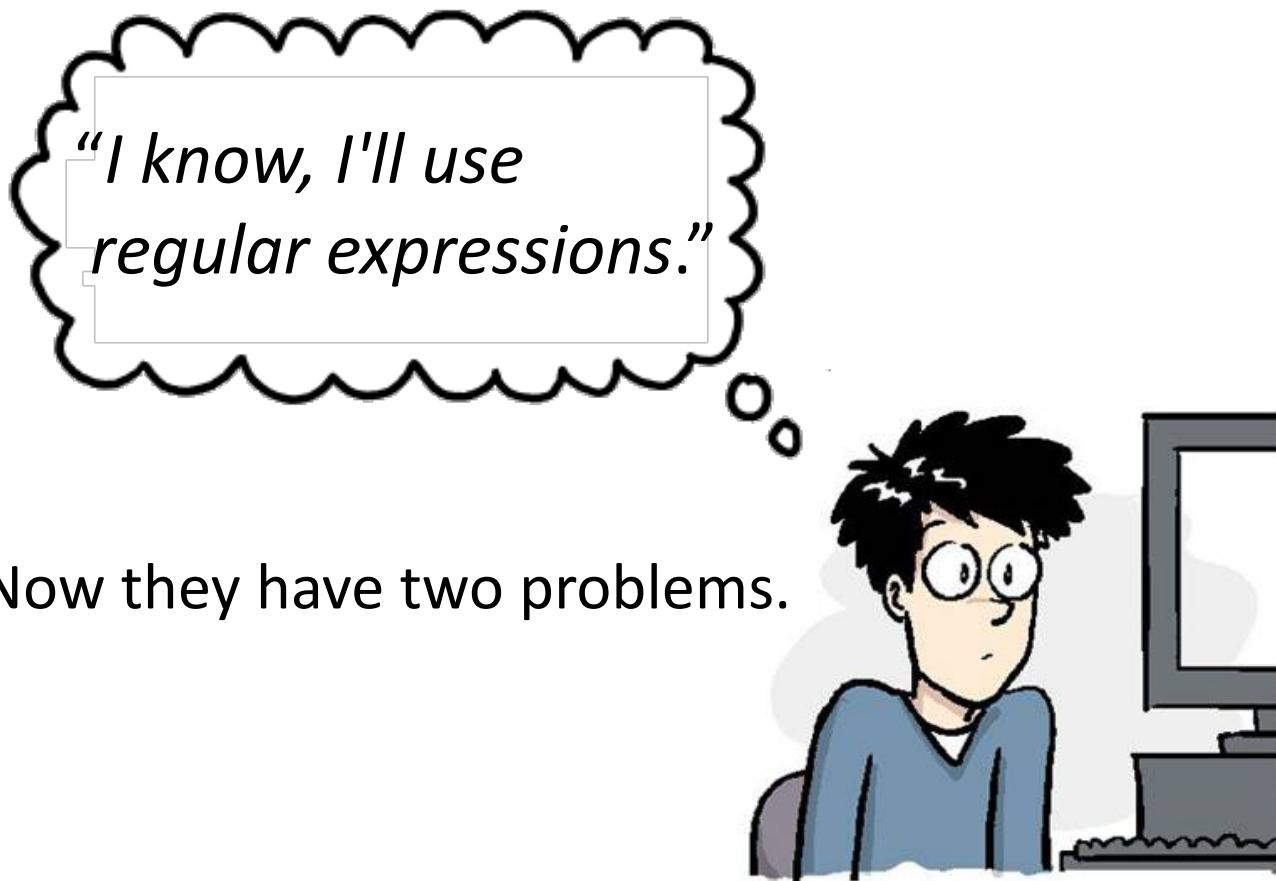
`<ID,"b"> <MULT> <ID,"b"> <MINUS> <INT,4> <MULT> <ID,"a"> <MULT> <ID,"c">`



# Parsing

- Goals
  - ▶ Decide whether the sequence of tokens a valid program in the source language
  - ▶ Construct a structured representation of the input
  - ▶ Error detection and reporting
- Challenges
  - ▶ How do you describe the programming language?
  - ▶ How do you check validity of an input?
  - ▶ Where do you report an error?

Some people, when confronted with a problem, think —



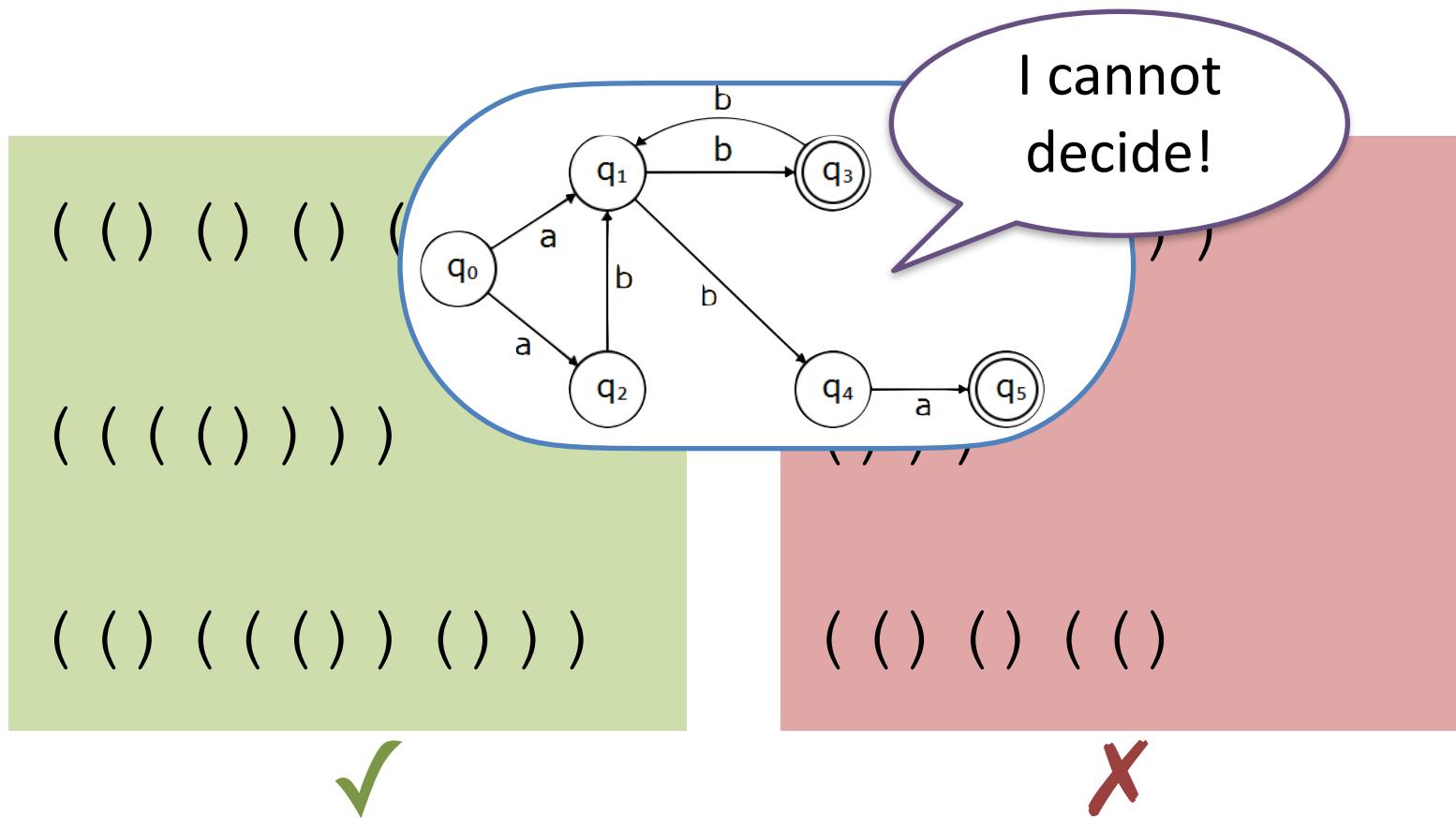
Now they have two problems.

The plural of *regex* is *regrets*.



# Regular Expressions: Computability Bound

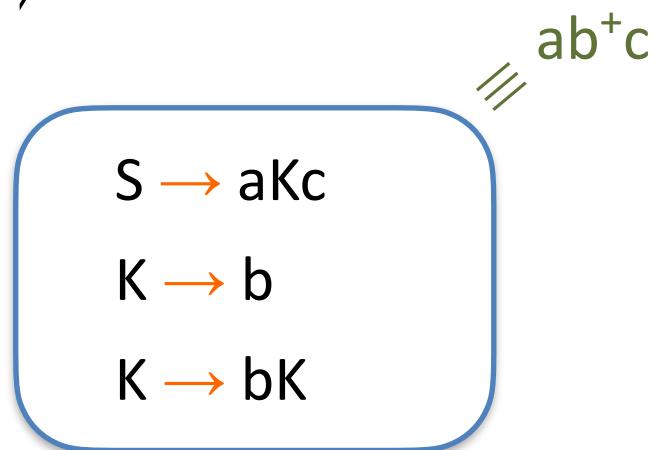
$L^()$  = language of balanced parentheses over  $\Sigma = \{‘(‘, ‘)’\}$



# Context Free Grammars

$$G = \langle V, T, P, S \rangle$$

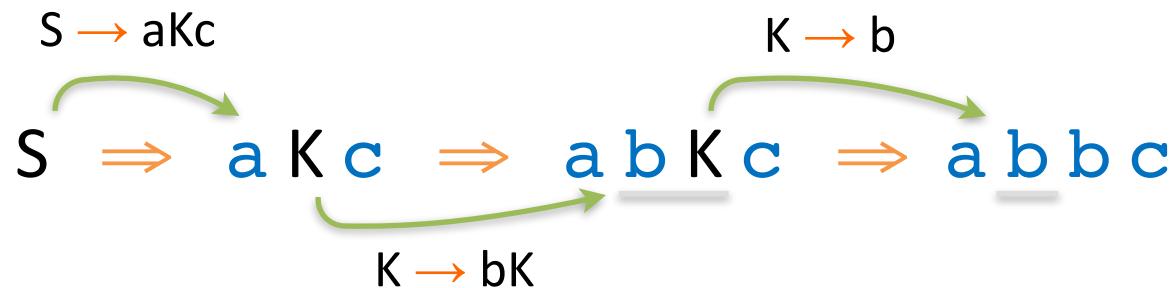
- $V$  – non-terminals  
(syntactic variables)
- $T$  – terminals (tokens)
- $P$  – derivation rules
  - ▶ Each rule of the form  
 $V \rightarrow (T \cup V)^*$
- $S$  – start symbol



$\begin{aligned} S &\Rightarrow aKc \\ &\Rightarrow abKc \\ &\Rightarrow abbKc \\ &\Rightarrow abbbc \end{aligned}$

# Terminology

- **Derivation** — a sequence of replacements of non-terminals using the derivation rules



- **Language** — the set of strings of terminals (“words”) derivable from the start symbol
- **Sentential form** — the result of a partial derivation, in which there may be non-terminals

# Context Free Grammars > Regular Expressions

$L^0$  from before is accepted  
by the following CFG:

$$S \rightarrow \epsilon$$

$$S \rightarrow ( S )$$

$$S \rightarrow SS$$

$$S \Rightarrow ( S )$$

$$\Rightarrow ( SS )$$

$$\Rightarrow ( ( S ) S )$$

$$\Rightarrow ( ( ) S )$$

$$\Rightarrow ( ( ) ( S ) )$$

$$\Rightarrow ( ( ) ( SS ) )$$

$$\Rightarrow ( ( ) ( ( S ) S ) )$$

$$\Rightarrow ( ( ) ( ( S ) ( S ) ) )$$

$$\Rightarrow ( ( ) ( ( ( S ) ) ( S ) ) )$$

$$\Rightarrow ( ( ) ( ( ( ) ) ( S ) ) )$$

$$\Rightarrow ( ( ) ( ( ( ) ) ( ) ) )$$

# Example

$S \rightarrow S ; S$

$S \rightarrow id := E$

$E \rightarrow id \mid num$

$\quad \mid E + E \mid E * E \mid ( E )$

$V = \{ S, E \}$

$T = \{ id, num, +, *, (, ), :=, ; \}$

# Derivation

$\text{id} := \text{id} ; \text{id} := \text{id} + \text{id}$

input

lex

$x := z ;$   
 $y := x + z$

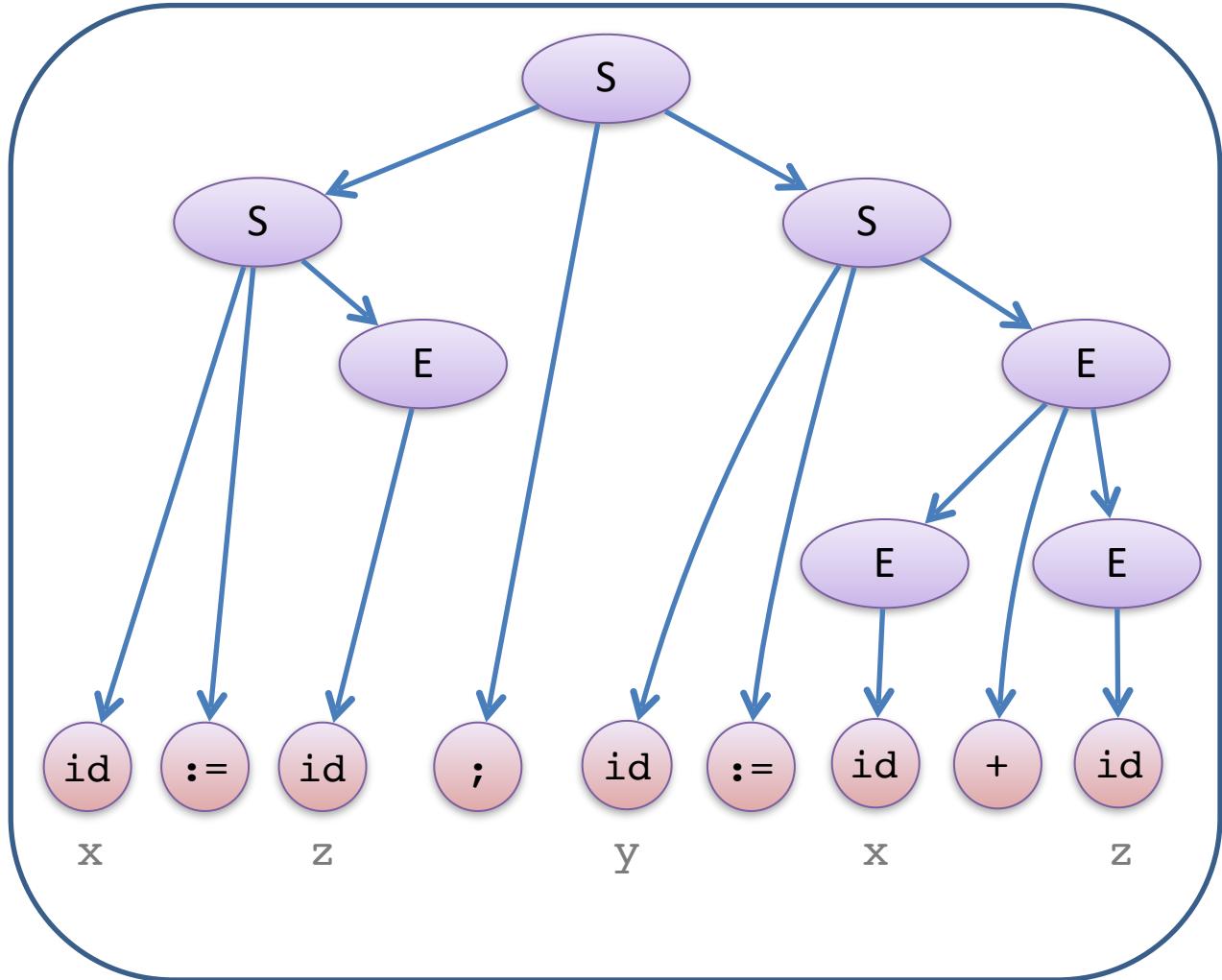
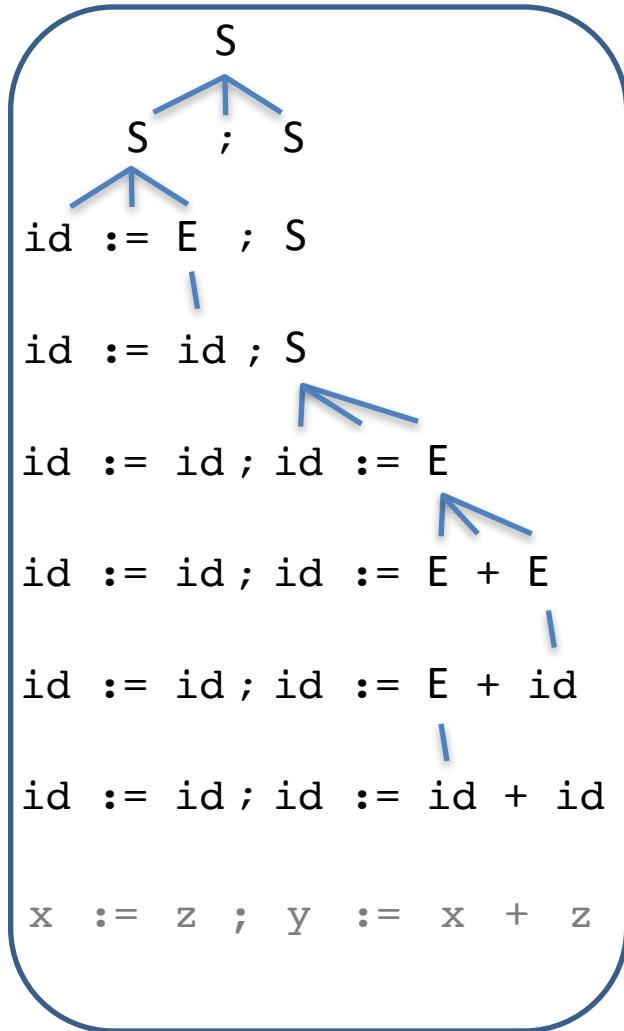
grammar

$S \rightarrow S ; S$   
 $S \rightarrow \text{id} := E$   
 $E \rightarrow \text{id} | E + E | E * E | ( E )$

$S$   
 $S ; S$   
 $\text{id} := E ; S$   
 $\text{id} := E ; \text{id} := E$   
 $\text{id} := \text{id} ; \text{id} := E$   
 $\text{id} := \text{id} ; \text{id} := E + E$   
 $\text{id} := \text{id} ; \text{id} := E + \text{id}$   
 $\text{id} := \text{id} ; \text{id} := \text{id} + \text{id}$   
  
 $x := y ; y := x + z$

# Parse Tree

also: Syntax Tree, Derivation Tree



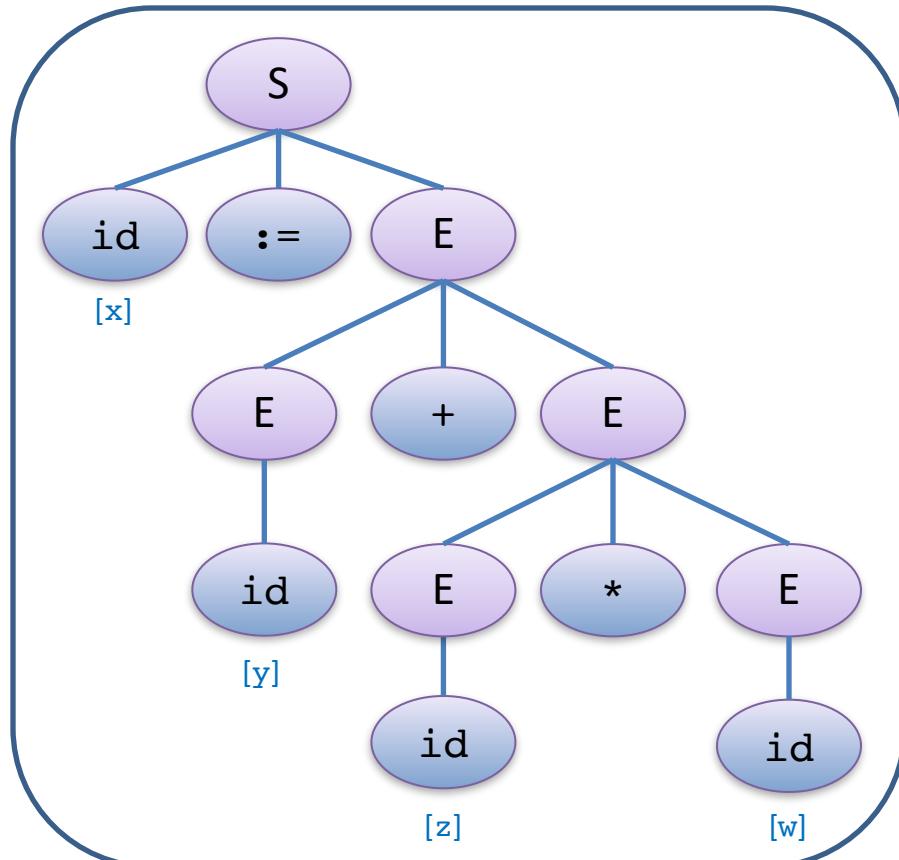
# Questions

- How did we know which rule to apply on every step?
- Does it matter?
  - ▶ Can we get “stuck” if we make a wrong choice?
  - ▶ Would we always get the same parse tree?

# Ambiguity

$x := y + z * w$

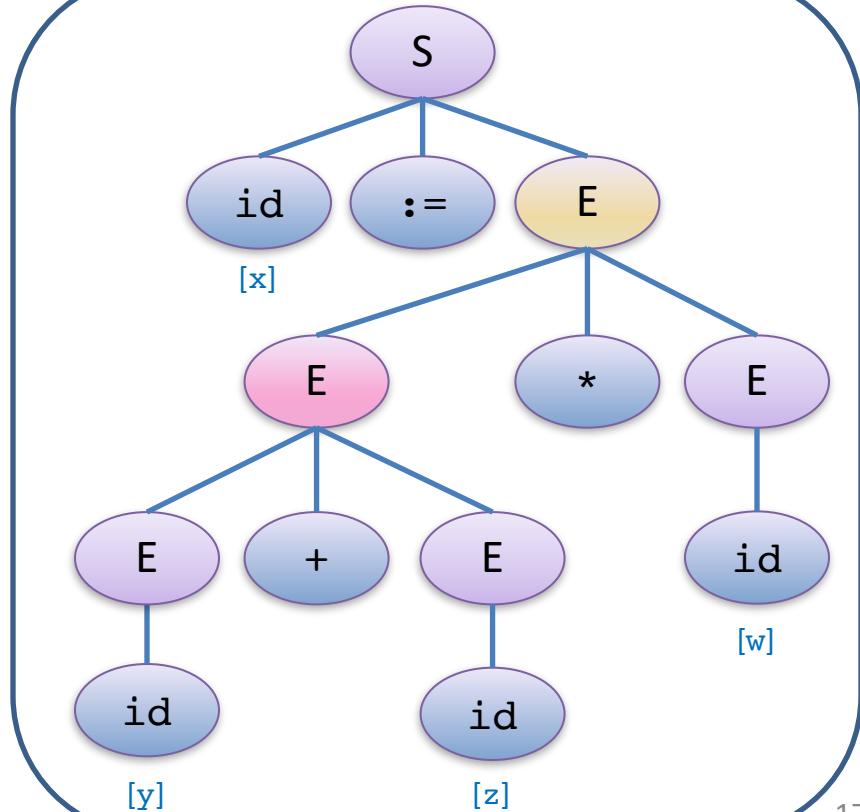
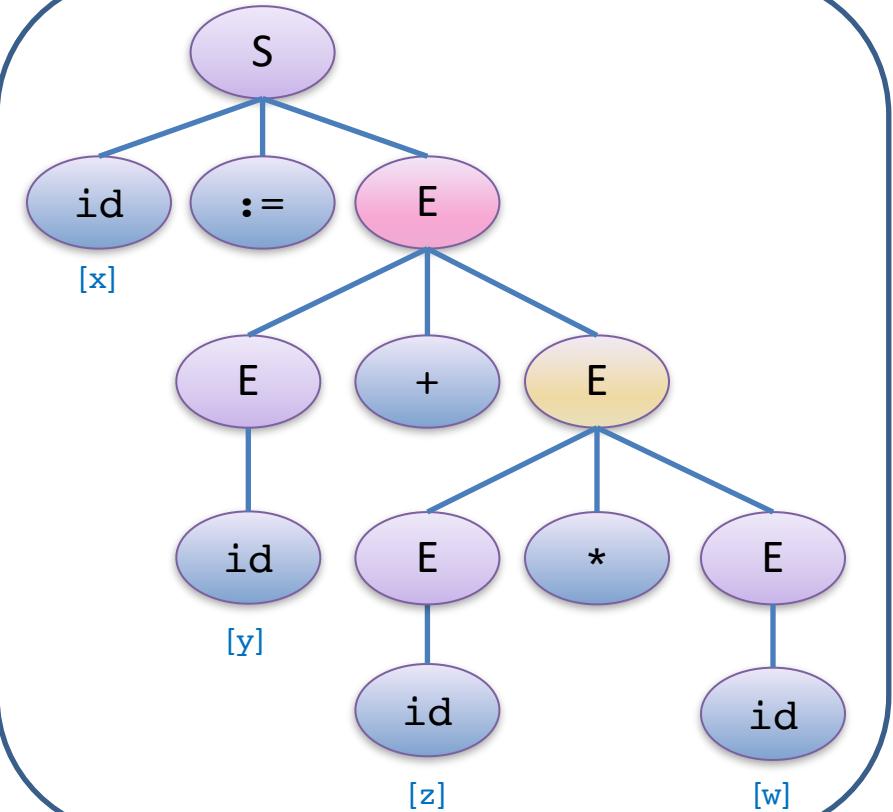
$S \rightarrow S ; S$   
 $S \rightarrow id := E$   
 $E \rightarrow id \mid E + E \mid E * E \mid ( E )$



# Ambiguity

$x := y + z * w$

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

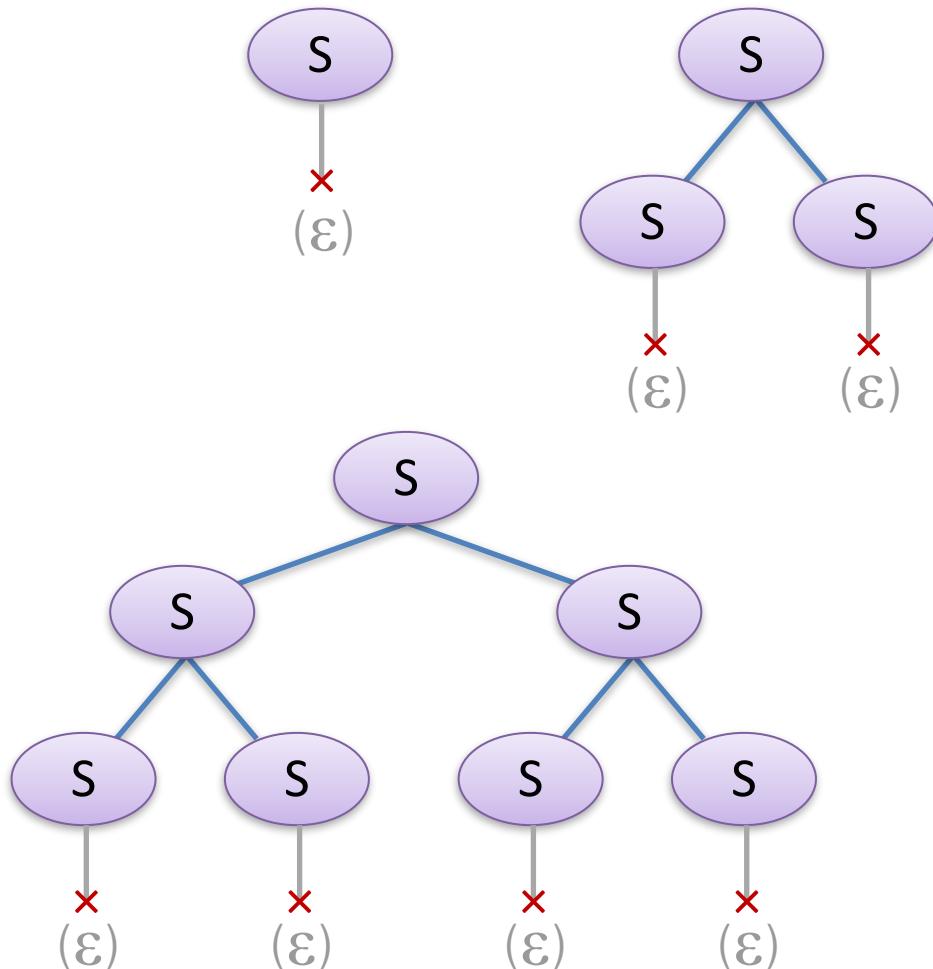


# Ambiguity

$S \rightarrow \epsilon$

$S \rightarrow ( S )$

$S \rightarrow SS$



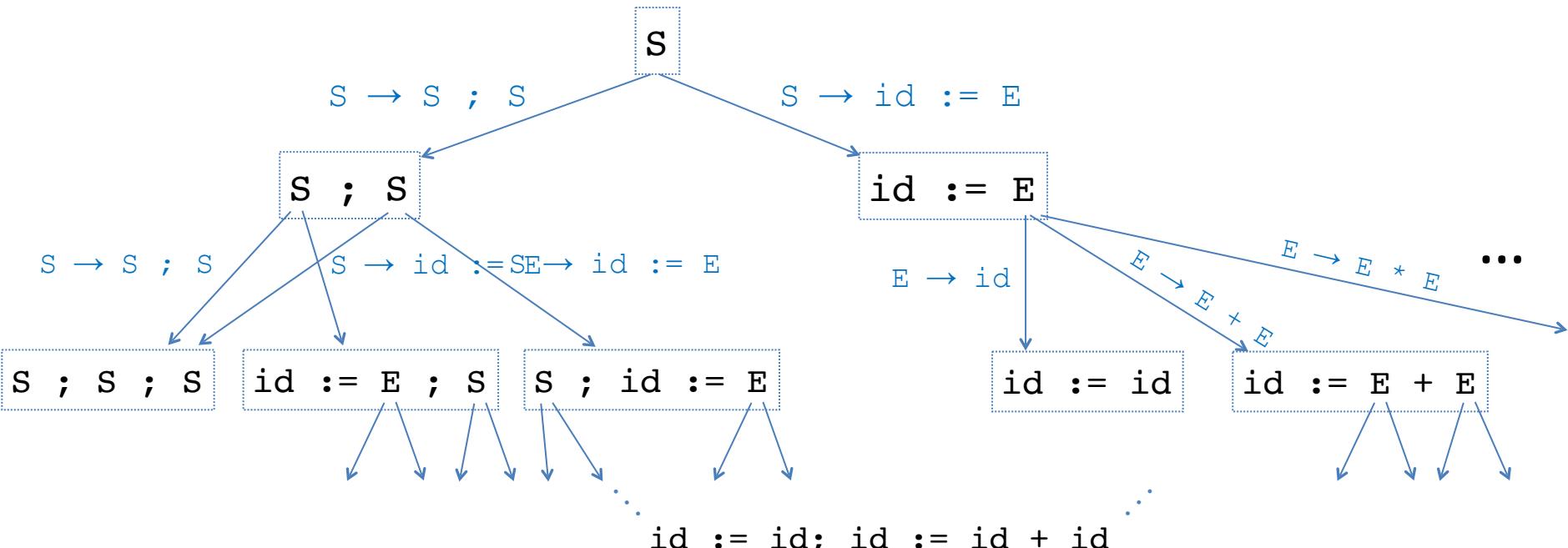
# Complexity of Parsing

- Parsing can be seen as a search problem
  - ▶ Can you find a derivation from the start symbol to the input word?
  - ▶ Easy (but very expensive) to solve with backtracking

# “Brute-force” Parsing

```
x := z;  
y := x + z
```

```
S → S ; S  
S → id := E  
E → id | E + E | E * E | ( E )
```



(not a parse tree... a **search** for a parse by exhaustively applying all rules)

# Chomsky Hierarchy

Turing machine

Linear-bounded  
non-deterministic  
Turing machine

Recursively enumerable

Context sensitive

Context free

Regular

Non-deterministic  
pushdown  
automaton

Finite-state  
automaton

Hm, that looks  
about right.



# Complexity of Parsing

- We already know: a regular language can be recognized by a deterministic finite automaton.
- A context free language can be recognized by a **non-deterministic pushdown automaton**.

## Non-deterministic pushdown automaton?!

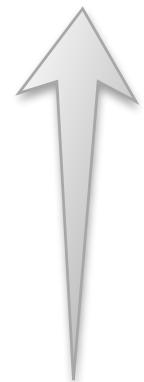
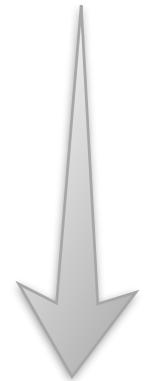
- We want ***efficient*** parsers that we can ***execute***
  - ▶ **Deterministic** computation model preferred

# Parsing Algorithms

- *Earley* and *CYK* are parsing algorithms that can be used to parse any context-free language — with a (worst case) time complexity of  $O(n^3)$
  - We want *more* efficient parsers
    - ▶ Linear in input size
- ⇒ **We will sacrifice generality for efficiency**

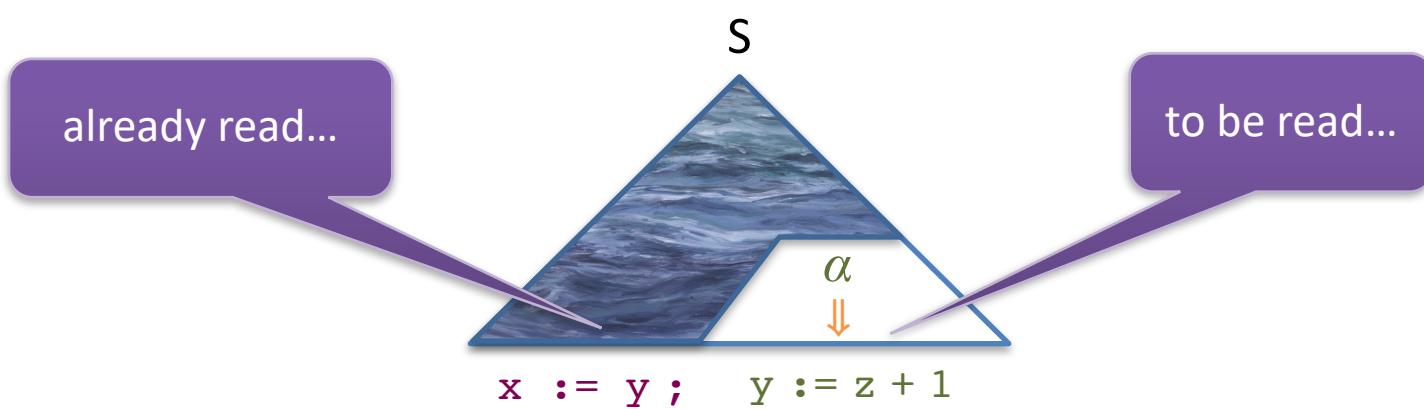
# Efficient Parsers

- Top-down (predictive)
  - ▶ Construct the leftmost derivation
  - ▶ Apply rules “from left to right”
  - ▶ Predict what rule to apply based on nonterminal and token
  
- Bottom up (shift-reduce)
  - ▶ Construct the rightmost derivation
  - ▶ Apply rules “from right to left”
  - ▶ Reduce a right-hand side of a production to its non-terminal

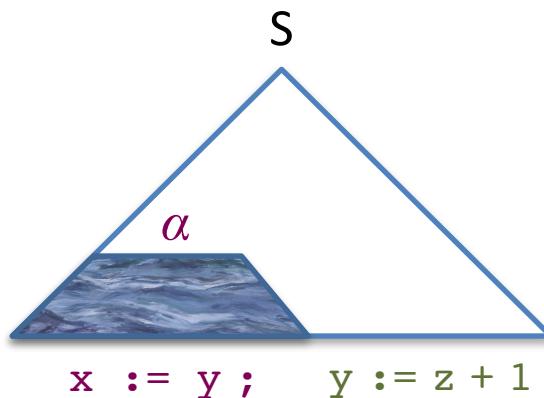


# Efficient Parsers

- Top-down (predictive)

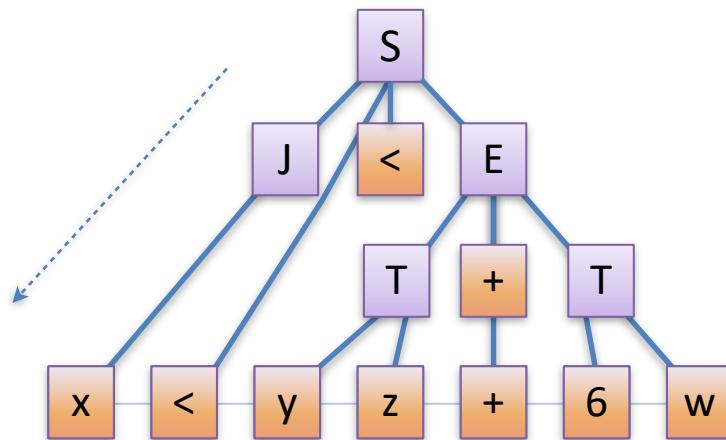


- Bottom-up (shift-reduce)



# Parser Classes – Illustrated

- Top-down (predictive)



$S \rightarrow J < E$   
 $J \rightarrow ID$   
 $E \rightarrow T + T$   
 $T \rightarrow ID\ ID \mid NUM\ ID$

# Top-down Parsing

- Given a grammar  $G = \langle V, T, P, S \rangle$  and a word  $w$
- Goal: derive  $w$  using  $G$ 
  - ▶  $S \xrightarrow{*} w$
- Idea
  - ▶ Apply production to leftmost nonterminal
  - ▶ Pick production rule based on next input token



# Leftmost Derivation

input

```
x := z;  
y := x + z
```

grammar

```
S → S ; S  
S → id := E  
E → id | E + E | E * E | ( E )
```

► always expand  
leftmost  
non-terminal

```
S  
S ; S  
id := E ; S  
id := id ; S  
id := id ; id := E  
id := id ; id := E + E  
id := id ; id := id + E  
id := id ; id := id + id
```

S → S ; S  
S → id := E  
E → id  
S → id := E  
E → E + E  
E → id  
E → id

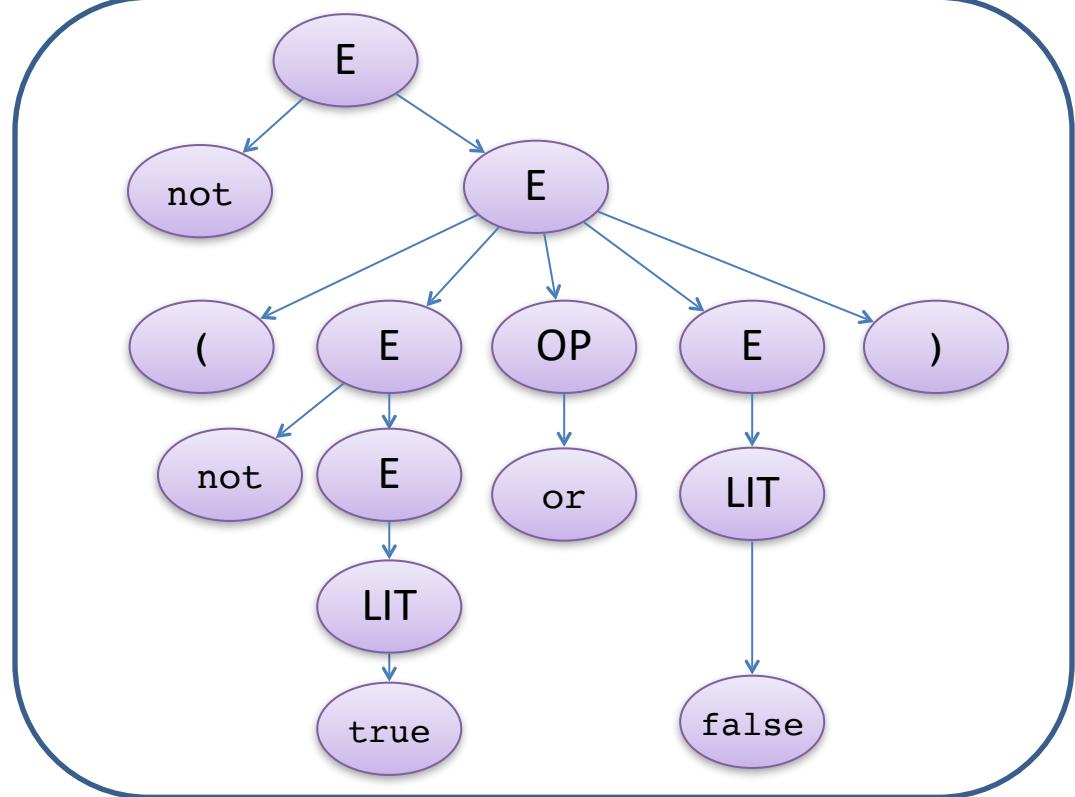
**x := z ; y := x + z**

# Boolean Expressions Example

not ( not true  
or false )

$E \rightarrow \text{LIT} \mid (\ E \text{ OP } E ) \mid \text{not } E$   
 $\text{LIT} \rightarrow \text{true} \mid \text{false}$   
 $\text{OP} \rightarrow \text{and} \mid \text{or} \mid \text{xor}$

$E \Rightarrow$   
 $\text{not } E \Rightarrow$   
 $\text{not} ( E \text{ OP } E ) \Rightarrow$   
 $\text{not} ( \text{not } E \text{ OP } E ) \Rightarrow$   
 $\text{not} ( \text{not } \text{LIT} \text{ OP } E ) \Rightarrow$   
 $\text{not} ( \text{not } \text{true} \text{ OP } E ) \Rightarrow$   
 $\text{not} ( \text{not } \text{true} \text{ or } E ) \Rightarrow$   
 $\text{not} ( \text{not } \text{true} \text{ or } \text{LIT} ) \Rightarrow$   
 $\text{not} ( \text{not } \text{true} \text{ or } \text{false} )$



Production to apply is known from **next input token**

# Recursive Descent Parsing

- Define a subroutine for every nonterminal
  - Every subroutine works as follows:
    - ▶ Select applicable production rule
    - ▶ Iterate over symbols in rhs
      - ▷ For terminal symbols —  
check that it matches the next input token
      - ▷ For nonterminal symbols —  
calls (recursively) their corresponding functions
- Use Lookahead*
- Left-to-right, depth-first:  
yields Leftmost derivation*

# Matching tokens

```
void match(token t) {  
    if (current == t)  
        current = next_token();  
    else  
        error;  
}
```

Consume token from input  
Set Lookahead to next token

- Variable **current** holds the *Lookahead*

# Functions for Nonterminals

```
E → LIT | ( E OP E ) | not E
```

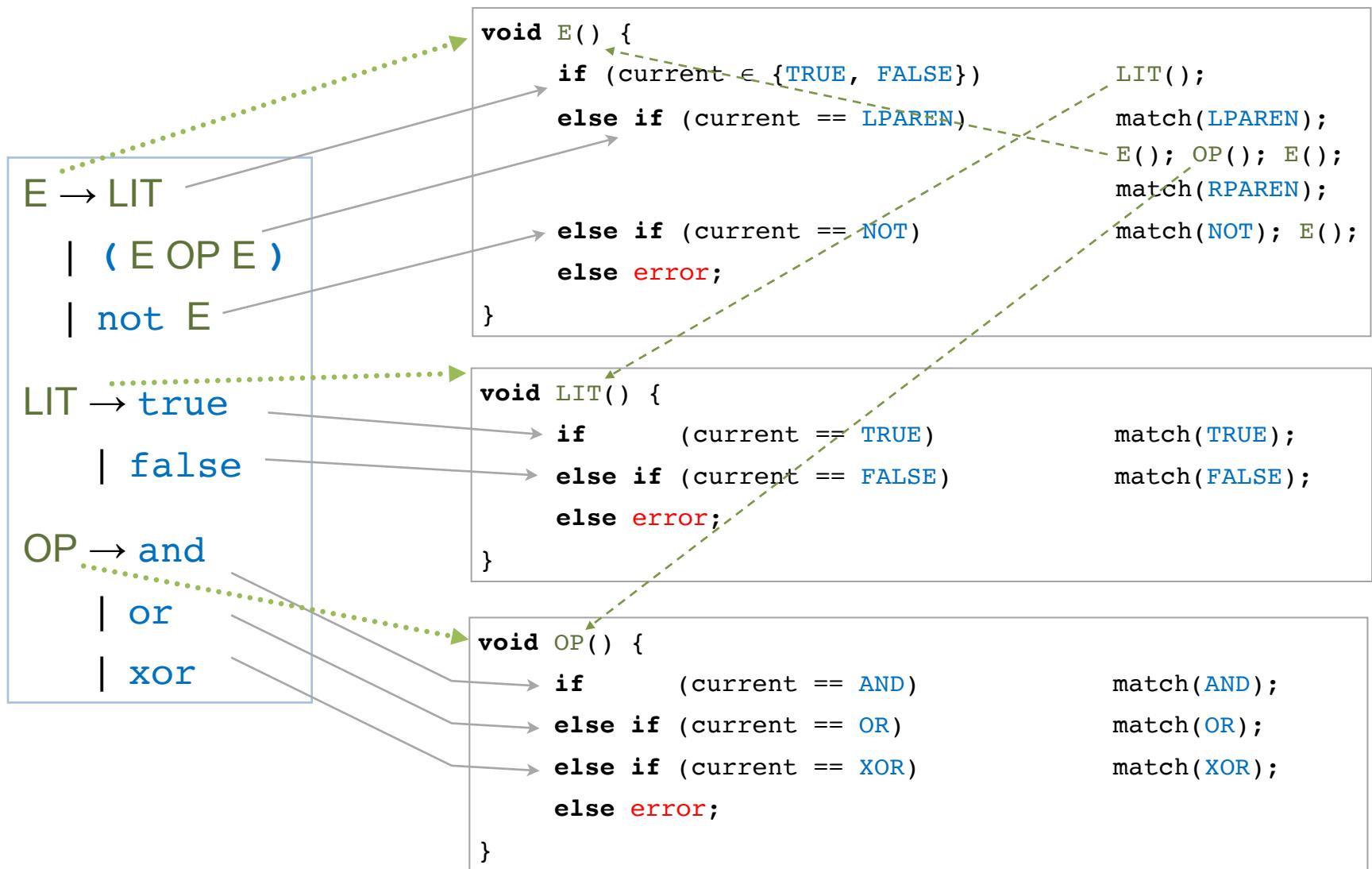
```
LIT → true | false
```

```
OP → and | or | xor
```

```
void E() {
    if (current ∈ {TRUE, FALSE}) // E → LIT
        LIT();
    else if (current == LPAREN) // E → ( E OP E )
        match(LPAREN); E(); OP(); E(); match(RPAREN);
    else if (current == NOT) // E → not E
        match(NOT); E();
    else
        error;
}

void LIT() {
    if (current == TRUE) match(TRUE);
    else if (current == FALSE) match(FALSE);
    else error;
}
```

# Functions for Nonterminals



# Recursive Descent – General Structure

```
void A() {  
    choose an A-production,  $A \rightarrow X_1 X_2 \cdots X_k$  ;  
    for (i = 1; i ≤ k; i++) {  
        if ( $X_i \in T$ )           match( $X_i$ ) ;  
        else /*  $X_i \in V$  */   call function  $X_i()$  ;  
    }  
}
```

The diagram illustrates the general structure of a recursive descent parser. It shows a pseudocode implementation of the A() function. A red oval highlights the 'choose' step, which selects an A-production. A red arrow points from this oval to the opening brace of the for-loop, indicating that this choice leads to the processing of each symbol in the production. A green dashed arrow points from the production rule  $A \rightarrow X_1 X_2 \cdots X_k$  to the variable  $X_i$ , representing the current symbol being matched or generated.

- How do you pick the right A-production?
  - ▶ Can we generalize the rationale behind all those **if** statements we saw in the example?

# LL( $k$ ) Grammars

- A grammar is in the class LL( $k$ ) when it can be parsed via:
  - ▶ Top-down analysis
  - ▶ Scanning the input from **left to right** (L)
  - ▶ Producing the **leftmost derivation** (L)
  - ▶ With **lookahead** of  $k$  tokens ( $k$ )
- A language is said to be LL( $k$ ) if it has an LL( $k$ ) grammar
- The simplest case is LL(1), which we discuss next

# FIRST Sets

- To formalize the property (of a grammar) that we can determine a rule using a single lookahead token, we define the FIRST sets.
- For a **sentential form**  $\alpha \in (V \cup T)^*$ 
  - ▶  $\text{FIRST}(\alpha)$  = all terminals that  $\alpha$  **can start with**
  - ▶ i.e.,  $\sigma \in \text{FIRST}(\alpha) \Leftrightarrow$  there exists some derivation sequence  $\alpha \Rightarrow^* \sigma w$
  - ▶ Special case:  $\varepsilon \in \text{FIRST}(\alpha) \Leftrightarrow \alpha \Rightarrow^* \varepsilon$  ( $\alpha$  is *nullable*)

!

(not really  
a terminal)

# FIRST Sets

- For every terminal  $\sigma$ , nonterminal  $N$ :
  - ▶  $\text{FIRST}(\varepsilon) = \{\varepsilon\}$
  - ▶  $\text{FIRST}(\sigma) = \{\sigma\}$   $\sigma$  terminal
  - ▶  $\text{FIRST}(N) = \bigcup \{\text{FIRST}(\alpha) \mid (N \rightarrow \alpha) \in P\}$   $N$  nonterminal

$$\text{FIRST}(x\alpha) = \begin{cases} \text{FIRST}(x) & \varepsilon \notin \text{FIRST}(x) \\ \text{FIRST}(x) \setminus \{\varepsilon\} & \varepsilon \in \text{FIRST}(x) \\ \cup \text{FIRST}(\alpha) \end{cases}$$

# FIRST Sets

$$E \rightarrow LIT \mid ( E \text{ OP } E ) \mid \text{not } E$$
$$LIT \rightarrow \text{true} \mid \text{false}$$
$$OP \rightarrow \text{and} \mid \text{or} \mid \text{xor}$$

- In our Boolean expressions example:

►  $\text{FIRST}( LIT ) = \{ \text{true}, \text{false} \}$

►  $\text{FIRST}( ( E \text{ OP } E ) ) = \{ '(' \}$

►  $\text{FIRST}( \text{not } E ) = \{ \text{not} \}$

These are all the alternatives for  $E \rightarrow \alpha$

# Recursive Descent (revised)

```
void A() {  
    find an A-production,  $A \rightarrow \overbrace{X_1 X_2 \dots X_k}^{\alpha}$ ,  
    such that current  $\in \text{FIRST}(\alpha)$   
  
    for (i = 1; i ≤ k; i++) {  
        if ( $X_i \in T$ )           match( $X_i$ );  
        else /*  $X_i \in V$  */     call function  $X_i()$ ;  
    }  
}
```

# FIRST Sets

- No intersection between FIRST sets  $\Rightarrow$  can *always*\* pick a single rule
- If the FIRST sets intersect, may need longer lookahead
  - ▶  $LL(1) \subset LL(2) \subset LL(3) \subset \dots$
  - ▶ But the size of FIRST sets grow exponentially with  $k$

\* except when null productions exist — see next slide

# FOLLOW Sets

- What do we do with nullable alternatives?

- Example:

$$S \rightarrow AB \mid c \quad A \rightarrow a \mid \epsilon \quad B \rightarrow b$$

- Need to know what may come after A in the language to select the right production
- For *every* nonterminal A
  - ▶ FOLLOW(A) = set of tokens that can immediately follow A (in some derivation)

# FOLLOW Sets

- (Fixpoint computation, similar to FIRST)
    - ▶  $\text{FOLLOW}(S) = \{\$\}$       $S$  = initial symbol,  $\$$  = end-of-input marker
    - ▶  $\text{FOLLOW}(N) =$   $N$  nonterminal  
 $(N \neq S)$ 

$$\bigcup \{\text{FIRST}(\omega) \setminus \varepsilon \mid (A \rightarrow \alpha N \omega) \in P\} \cup$$

$$\bigcup \{\text{FOLLOW}(A) \mid (A \rightarrow \alpha N \omega) \in P, \varepsilon \in \text{FIRST}(\omega)\}$$

# Recursive Descent (revised again)

```
void A() {  
    find an A-production,  $A \rightarrow \overbrace{X_1 X_2 \dots X_k}^{\alpha}$ ,  
    such that current  $\in \text{FIRST}(\alpha)$   
    or  $\epsilon \in \text{FIRST}(\alpha)$  and current  $\in \text{FOLLOW}(A)$   
  
    for (i = 1; i <= k; i++) {  
        if ( $X_i \in T$ ) match( $X_i$ );  
        else /*  $X_i \in V$  */ call function  $X_i()$ ;  
    }  
}
```

This is the basis for an LL(1) parser.

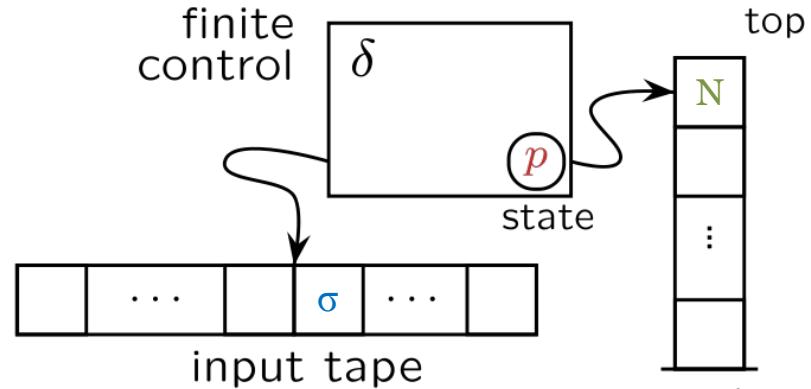
(but it is still recursive — stay tuned)

# LL( $k$ ) Parsers

- **Recursive Descent**
  - ▶ Manual construction
  - ▶ Uses recursion
- **Wanted**
  - ▶ A parser that can be generated automatically
  - ▶ Iterative — does not use recursion
  - ▶ Uses a deterministic computational model

# LL(1) Parsing with Pushdown Automata

- A PDA uses
  - ▶ Prediction stack
  - ▶ Input stream
  - ▶ Transition table



- $\delta : \text{nonterminals} \times \text{tokens} \rightarrow \text{production alternative}$   
(right-hand side)
- Entry  $\delta[N, \sigma]$  for nonterminal  $N$  and terminal  $\sigma$  contains the alternative of  $N$  that would be predicted when current input starts with  $\sigma$

# Example Transition Table

- |                                   |                                    |                                 |
|-----------------------------------|------------------------------------|---------------------------------|
| (1) $E \rightarrow LIT$           | (4) $LIT \rightarrow \text{true}$  | (6) $OP \rightarrow \text{and}$ |
| (2) $E \rightarrow ( E OP E )$    | (5) $LIT \rightarrow \text{false}$ | (7) $OP \rightarrow \text{or}$  |
| (3) $E \rightarrow \text{not } E$ |                                    | (8) $OP \rightarrow \text{xor}$ |

which rule should  
be used

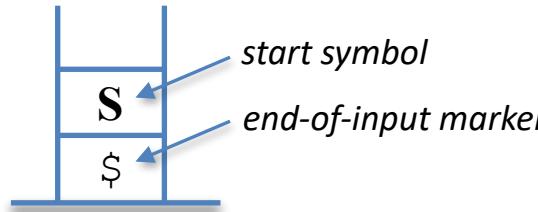
		next input token								
		(	)	not	true	false	and	or	xor	\$
current nonterminal	$\delta$	2		3	1	1				
	LIT				4	5				
	OP						6	7	8	

# The Transition Table

- Constructing the transition table is easy
  - It relies on FIRST and FOLLOW
  - Based on the concept of our recursive descent earlier
- You will construct FIRST, FOLLOW, and the table in the tutorials

# LL( $k$ ) parsing with pushdown automata

- Initial state —



- Two possible moves  $(\sigma = \text{next token})$

- ▶ **Predict**

When top of stack is nonterminal  $N$ , pop  $N$ , lookup  $\delta[N, \sigma]$ .

If  $\delta[N, \sigma]$  is defined, push  $\delta[N, \sigma]$  on prediction stack.

Otherwise – **syntax error**.

- ▶ **Match**

When top of stack is a terminal  $t$ , it must be equal to next input token.

If  $t = \sigma$ , pop  $t$  and consume  $\sigma$ .

If  $t \neq \sigma$  – **syntax error**.

- Parsing terminates when prediction stack is empty.

- ▶ At this point the input *must be finished*  $\Rightarrow$  **success**.

# Simple Example

aacbb\$

$A \rightarrow aAb \mid c$

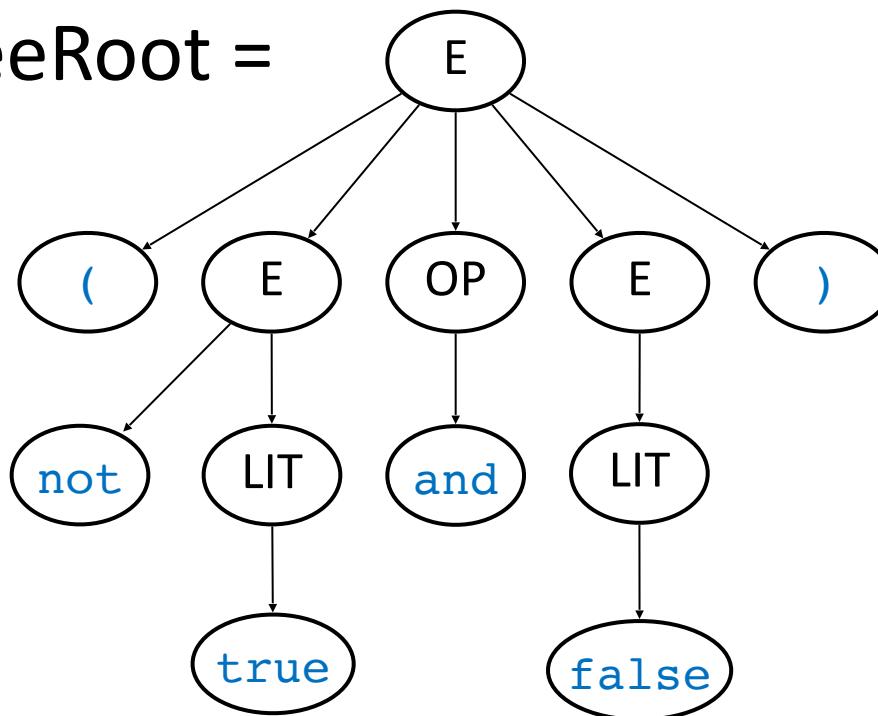
Stack content (top is left)	Remaining input	Action
► A \$	a a c b b \$	predict(A,a) = aAb
aAb \$	a a c b b \$	match(a,a)
Ab \$	a c b b \$	predict(A,a) = aAb
aAbb \$	a c b b \$	match(a,a)
Abb \$	c b b \$	predict(A,c) = c
c b b \$	c b b \$	match(c,c)
b b \$	b b \$	match(b,b)
b \$	b \$	match(b,b)
\$	\$	match(\$,\$) accept
	a	b
► A	aAb	c

# Building the Syntax Tree

- Input = “( not true and false )”;



- Node treeRoot =



# Building the Syntax Tree

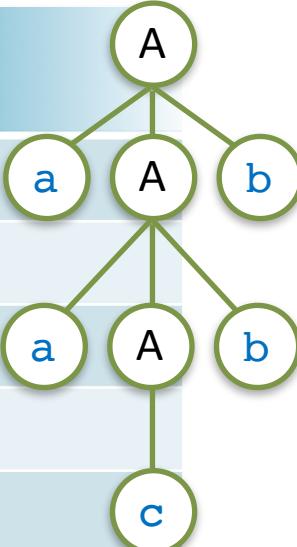
- Can add an action to perform on each production — in the LL( $k$ ) case, at **predict**
- Every symbol on the parser stack (except **\$**) represents a node in the syntax tree.
  - ▶ The initial entry for  $S$  is the root.
  - ▶ When the top of the stack is a nonterminal **N** and parser performs **predict(N,  $X_1X_2\dots X_k$ )**:
    - ▶ Create nodes for  $X_i$ 's as child nodes of **N**

# Simple Example

aacbb\$

$A \rightarrow aAb \mid c$

Stack content (top is left)	Remaining input	Action
A \$	a a c b b \$	$\text{predict}(A,a) = aAb$
aAb \$	a a c b b \$	$\text{match}(a,a)$
Ab \$	a c b b \$	$\text{predict}(A,a) = aAb$
aAbb \$	a c b b \$	$\text{match}(a,a)$
Abb \$	c b b \$	$\text{predict}(A,c) = c$
c b b \$	c b b \$	$\text{match}(c,c)$
b b \$	b b \$	$\text{match}(b,b)$
b \$	b \$	$\text{match}(b,b)$
\$	\$	$\text{match}($,$)$ accept



	a	b	c
A	aAb		c

# LL(1) Conflicts

term → ID | indexed

indexed → ID [ expr ]

- FIRST( ID ) = { ID }
- FIRST( indexed ) = { ID }

} Both are alternatives for term →  $\alpha$

FIRST/FIRST conflict

⇒ This grammar is not in LL(1). Can we “fix” it?

# Left Factoring

- Rewrite the grammar to be in LL(1)

```
term → ID | indexed
```

```
indexed → ID [ expr ]
```



```
term → ID after_ID
```

```
after_ID → [ expr ] | ε
```

Intuition: just like factoring  $x \cdot y + x \cdot z$  into  $x \cdot (y + z)$

# Left Factoring

## another example

```
S → if E then { S } else { S }
| if E then { S }
| ...
```



```
S → if E then { S } S'
| ...
S' → else { S } | ε
```

# Are We Done?

$$S \rightarrow A \ a \ b$$
$$A \rightarrow a \mid \epsilon$$

- Select a rule for  $A$  with  $a$  in the look-ahead:
  - ▶ Should we pick (1)  $A \rightarrow a$  or (2)  $A \rightarrow \epsilon$  ?

- ▶  $\text{FIRST}(a) = \{ 'a' \}$
  - ▶  $\text{FIRST}(\epsilon) = \{ \epsilon \}$        $\text{FOLLOW}(A) = \{ 'a' \}$
- } alternatives for  
 $A \rightarrow a$

FIRST/FOLLOW conflict

⇒ The grammar is not in LL(1). Can we fix *that*?

# Grammatical Substitution

$$S \rightarrow A \ a \ b$$
$$A \rightarrow a \mid \epsilon$$


Substitute A in S

$$S \rightarrow a \ a \ b \mid a \ b$$


Left factoring

$$S \rightarrow a \text{ after\_a}$$
$$\text{after\_a} \rightarrow a \ b \mid b$$

# So Far

- Can determine if a grammar is in LL(1)
  - ▶ The FIRST and FOLLOW sets
  - ▶ Algorithms for computing these – in tutorials
- Have some techniques for modifying a grammar to find an equivalent grammar in LL(1)
  - ▶ Left factoring
  - ▶ Assignment
- Now let's look at a third example and present one more such technique

# Left Recursion

$$E \rightarrow E + \text{term} \mid \text{term}$$

- Left recursion cannot be handled with a bounded lookahead
- What can we do?
- Theorem: Any grammar with left recursion has an equivalent grammar without left recursion

# Left Recursion Removal

$$N \rightarrow N\alpha \mid \beta$$

$G_1$



$$\begin{aligned} N &\rightarrow \beta N' \\ N' &\rightarrow \alpha N' \mid \epsilon \end{aligned}$$

$G_2$

$$L(G_1) = \{\beta, \beta\alpha, \beta\alpha\alpha, \beta\alpha\alpha\alpha, \dots\}$$

$$N' \Rightarrow^* \epsilon, \alpha, \alpha\alpha, \alpha\alpha\alpha, \dots$$

$$L(G_2) = \{\beta, \beta\alpha, \beta\alpha\alpha, \beta\alpha\alpha\alpha, \dots\}$$

- ▶ For our example:

$$\begin{aligned} E &\rightarrow E + \text{term} \\ &\mid \text{term} \end{aligned}$$



$$\begin{aligned} E &\rightarrow \text{term } E' \\ E' &\rightarrow + \text{ term } E' \mid \epsilon \end{aligned}$$

# Error Handling

- Types of errors
  - ▶ Lexical errors
  - ▶ Syntax errors
  - ▶ Semantic errors (e.g., type mismatch)
  - ▶ Runtime errors (e.g., division by zero)
- ▶ Requirements
  - Report the error clearly
  - Recover and continue so more errors can be discovered
  - Be efficient

# Simple Example, Bad Word

abcbb\$

$A \rightarrow aAb \mid c$

Stack content	Remaining input	Action
A\$	abcbb\$	$\text{predict}(A,a) = aAb$
aAb\$	abcbb\$	$\text{match}(a,a)$
Ab\$	bcbbs	$\text{predict}(A,b) =$ error

	a	b	c
A	aAb		c

# Simple Example, Bad Word

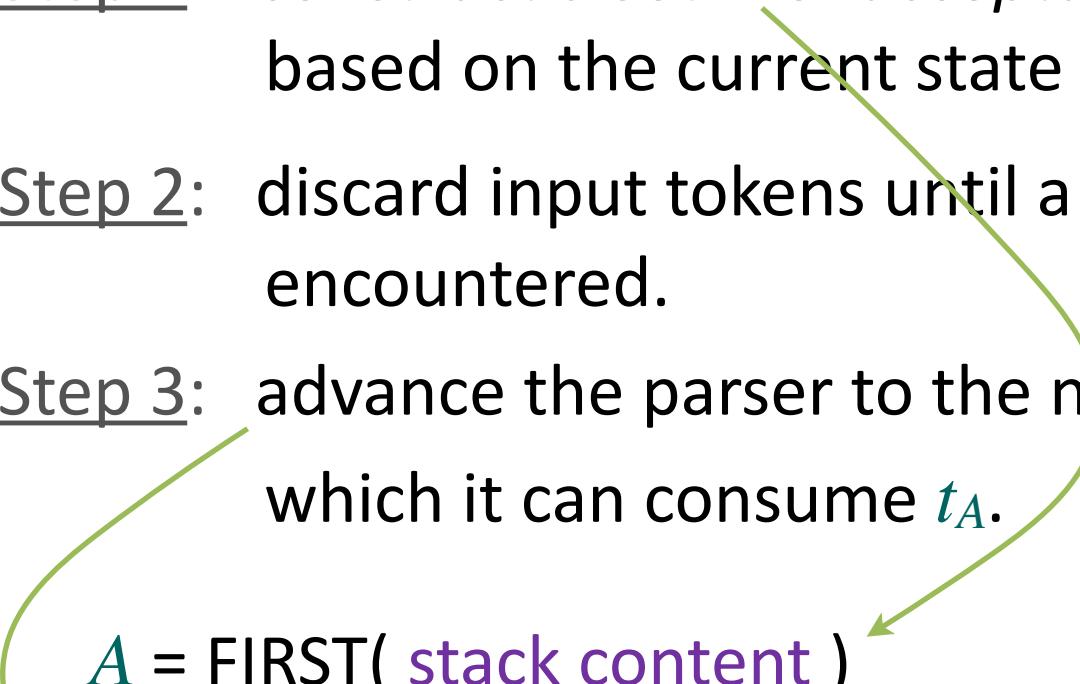
a~~bc~~bb\$

A → aAb | c

Stack content	Remaining input	Action
A\$	a <del>bc</del> <ins>b</ins> b\$	predict(A,a) = aAb
aAb\$	a <del>bc</del> <ins>b</ins> b\$	match(a,a)
Ab\$	<del>b</del> cbb\$	predict(A,b) = X skip
Ab\$	cbb\$	predict(A,a) = A → c
cb\$	cbb\$	match(c,c)
b\$	bb\$	match(b,b)
\$	b\$	match(\$,b) = X skip
\$	\$	match(\$,\$)

	a	b	c
A	aAb		c

# Error Handling in LL Parsers

- The **acceptable-set method**
    - ▶ 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  $t_A$ .
- $A = \text{FIRST}(\text{ stack content })$
- don't need to do anything
- panic mode
- 
- 

# Error Handling in LL Parsers

- The **acceptable-set method**
    - ▶ 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  $t_A$ .
- $A = \text{FOLLOW}(\text{ stack top })$
- pop the topmost nonterminal
- FOLLOW-set  
error recovery**

# Error Handling

- Different compilers adopt different approaches.
- Some examples
  - ▶ Panic mode (or another acceptable-set method): drop tokens until reaching a synchronizing token, like a semicolon, a right parenthesis, end of file, etc.
  - ▶ Phrase-level recovery: attempting local changes: replace “,” with “;”, eliminate or add a “;”, etc.
  - ▶ Error production: anticipate errors and automatically handle them by adding special rules to the grammar
  - ▶ Global correction: find the minimum modification to the program that will make it derivable in the grammar
    - ▷ Not a practical solution, except with very small grammars

# A Note About ANTLR



- ANTLR = ANother Tool for Language Recognition
- LL(\*) algorithm
  - ▶ Like LL( $k$ ) on steroids
  - ▶ Notable extensions:
    - ▶ Repeat operators — like in regex
    - ▶ Lookahead predicates — allow for unbounded lookahead at the cost of backtracking
- \* There's a nice demo at [lab.antlr.org](http://lab.antlr.org)

\* LL(\*): *The Foundation of the ANTLR Parser Generator*, Parr and Fischer, PLDI 2011

# Summary

- Parsing
  - ✓ Top-down or bottom-up
- Top-down parsing
  - ✓ Recursive descent
  - ✓ LL(k) grammars
  - ✓ LL(k) parsing with pushdown automata
- LL(k) parsers limitations
  - ✓ Cannot deal with common prefixes and left recursion
  - ✓ Left-recursion removal might result in complicated grammar

# Coming Up

