

# LECTURE 11

## INTRODUCTION TO SYNTAX ANALYSIS

# SUBJECTS

**Context free grammars**

**Derivations**

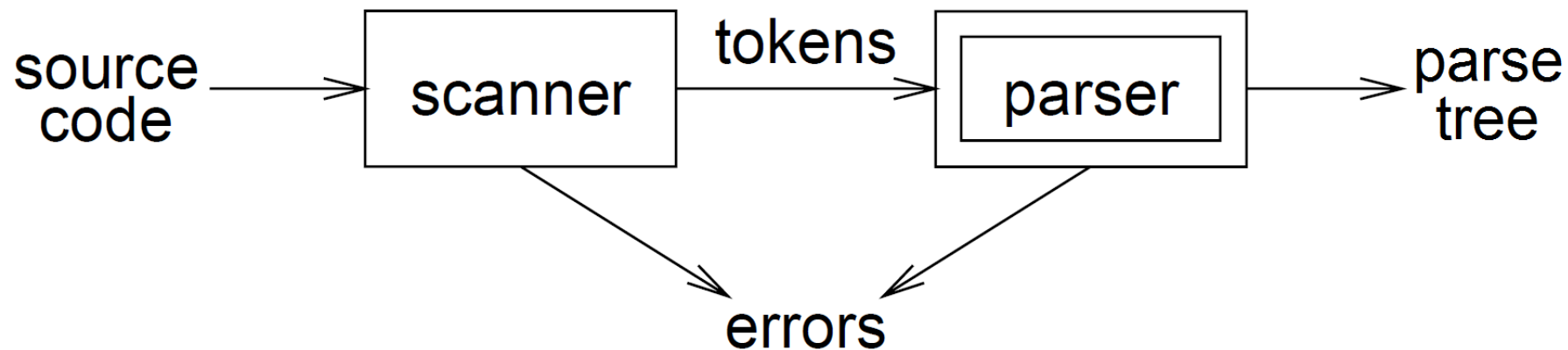
**Parse Trees**

**Ambiguity**

**Top-down parsing**

**Left recursion**

# THE ROLE OF PARSER



# CONTEXT FREE GRAMMARS

**A Context Free Grammar (CFG) consists of**

- Terminals
- Nonterminals
- Start symbol
- Productions

**A language that can be generated by a grammar is said to be a **context-free language****

# CONTEXT FREE GRAMMARS

**Terminals:** are the basic symbols from which strings are formed

- These are the tokens that were produced by the Lexical Analyser

**Nonterminals:** are syntactic variables that denote sets of strings

- One nonterminal is distinguished as the **start symbol**

The **productions** of a grammar specify the manner in which the terminal and nonterminals can be combined to form strings

# EXAMPLE OF GRAMMAR

The grammar with the following productions defines simple arithmetic expressions

$$\begin{aligned}\langle \text{expr} \rangle &::= \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \\ \langle \text{expr} \rangle &::= \text{id} \\ \langle \text{expr} \rangle &::= \text{num} \\ \langle \text{op} \rangle &::= + \\ \langle \text{op} \rangle &::= - \\ \langle \text{op} \rangle &::= * \\ \langle \text{op} \rangle &::= /\end{aligned}$$

In this grammar, the terminal symbols are  
*num, id + - \* /*

The nonterminal symbols are  $\langle \text{expr} \rangle$  and  $\langle \text{op} \rangle$ , and  $\langle \text{expr} \rangle$  is the start symbol

# DERIVATIONS

$\langle \text{expr} \rangle \Rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle$

*is read "expr derives expr op expr"*

$\langle \text{expr} \rangle \Rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle$

$\Rightarrow \text{id} \langle \text{op} \rangle \langle \text{expr} \rangle$

$\Rightarrow \text{id} * \langle \text{expr} \rangle$

$\Rightarrow \text{id} * \text{id}$

is called a **derivation** of  $\text{id} * \text{id}$  from  $\text{expr}$ .

# DERIVATIONS

If  $A ::= \gamma$  is a production and  $\alpha$  and  $\beta$  are arbitrary strings of grammar symbols, we can say:

$$\alpha A \beta \Rightarrow \alpha \gamma \beta$$

If  $\alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n$ , we say  $\alpha_1$  **derives**  $\alpha_n$ .



# DERIVATIONS

$\Rightarrow$  means “derives in one step.”

$\Rightarrow^*$  means “derives in zero or more steps.”

• *if  $\alpha \Rightarrow^* \beta$  and  $\beta \Rightarrow \gamma$  then  $\alpha \Rightarrow^* \gamma$*

$\Rightarrow^+$  means “derives in one or more steps.”

If  $S \Rightarrow^* \alpha$ , where  $\alpha$  may contain nonterminals, then we say that  $\alpha$  is a sentential form

- If  $\alpha$  does not contain any nonterminals, we say that  $\alpha$  is a sentence

# DERIVATIONS

**G:** grammar

**S:** start symbol

**L(G):** the language generated by G

Strings in **L(G)** may contain only terminal symbols of **G**

A string of terminals **w** is said to be in **L(G)** if and only if  $S \xRightarrow{+} w$

- As we said, the string **w** is called a sentence of **G**

A language that can be generated by a grammar is said to be a context-free language

- If two grammars generate the same language, the grammars are said to be equivalent

# DERIVATIONS

We have already seen the following production rules:

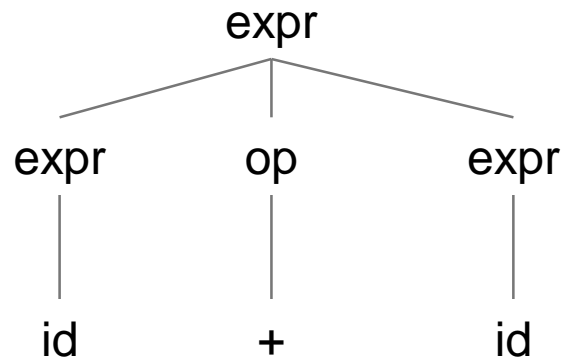
$$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \mid \text{id} \mid \text{num}$$
$$\langle \text{op} \rangle ::= + \mid - \mid * \mid /$$

The string `id+id` is a sentence of the above grammar because

$$\begin{aligned} \langle \text{expr} \rangle &\Rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \\ &\Rightarrow \text{id} \langle \text{op} \rangle \langle \text{expr} \rangle \\ &\Rightarrow \text{id} + \langle \text{expr} \rangle \\ &\Rightarrow \text{id} + \text{id} \end{aligned}$$

We write  $\langle \text{expr} \rangle \xRightarrow{*} \text{id} + \text{id}$

# PARSE TREE



This is called:  
**Leftmost derivation**

# TWO PARSE TREES

Let us again consider the arithmetic expression grammar.

For the line of code:

**x+2\*y**

(we are not considering the semi colon for now)



**Grammar:**

$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \mid \text{id} \mid \text{num}$

$\langle \text{op} \rangle ::= + \mid - \mid * \mid /$

# TWO PARSE TREES

Let us again consider the arithmetic expression grammar.

The sentence `id + id * id` has two distinct **leftmost** derivations:

```

<expr>  ⇒ <expr> <op> <expr>
        ⇒ id <op> <expr>
        ⇒ id + <expr>
        ⇒ id + <expr> <op> <expr>
        ⇒ id + id <op> <expr>
        ⇒ id + id * <expr>
        ⇒ id + id * id
    
```

```

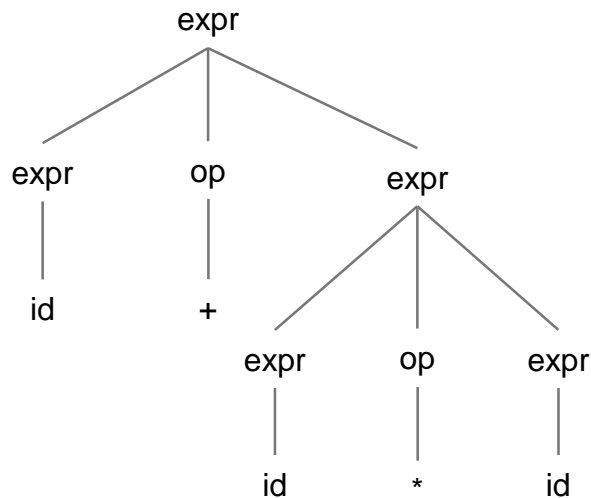
<expr>  ⇒ <expr> <op> <expr>
        ⇒ <expr> <op> <expr> <op> <expr>
        ⇒ id <op> <expr> <op> <expr>
        ⇒ id + <expr> <op> <expr>
        ⇒ id + id <op> <expr>
        ⇒ id + id * <expr>
        ⇒ id + id * id
    
```

**Grammar:**

```

<expr> ::= <expr> <op> <expr> | id | num
<op>  ::= + | - | * | /
    
```

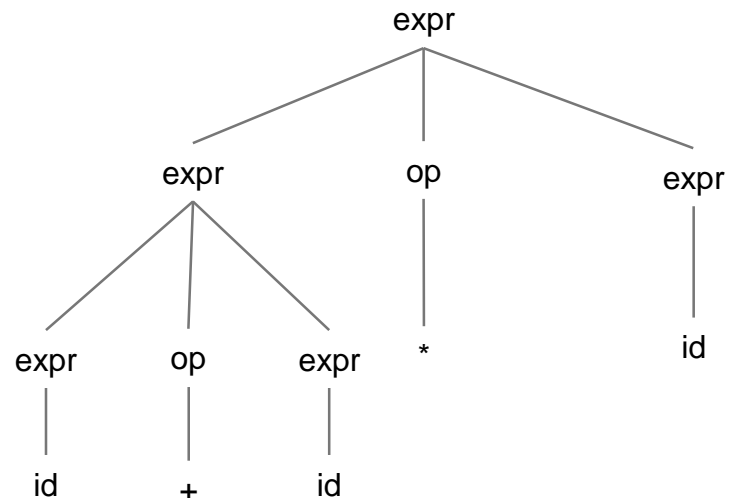
# TWO PARSE TREES



Equivalent to:  
id+(id\*id)

**Grammar:**

$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \mid \text{id} \mid \text{num}$   
 $\langle \text{op} \rangle ::= + \mid - \mid * \mid /$



Equivalent to:  
(id+id)\*id **✗**

# PRECEDENCE

**The previous example highlights a problem in the grammar:**

- It does not enforce precedence
- It has not implied an order of evaluation

**We can expand the production rules to add precedence**

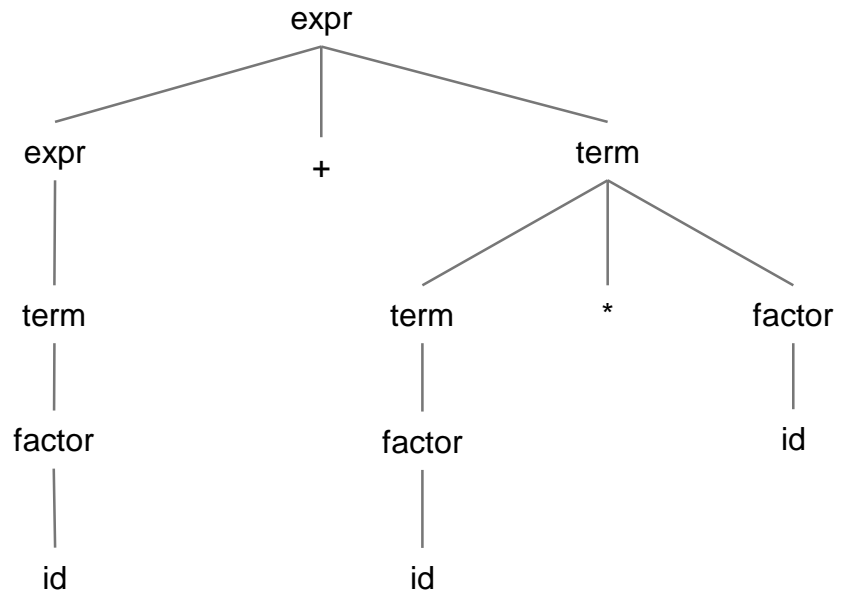
$$\begin{aligned}\langle \text{expr} \rangle &::= \langle \text{expr} \rangle + \langle \text{term} \rangle \\ &\quad | \langle \text{expr} \rangle - \langle \text{term} \rangle \\ &\quad | \langle \text{term} \rangle \\ \langle \text{term} \rangle &::= \langle \text{term} \rangle * \langle \text{factor} \rangle \\ &\quad | \langle \text{term} \rangle / \langle \text{factor} \rangle \\ &\quad | \langle \text{factor} \rangle \\ \langle \text{factor} \rangle &::= \text{num} \\ &\quad | \text{id}\end{aligned}$$



# APPLYING PRECEDENCE UPDATE

The sentence `id + id * id` has only one **leftmost** derivation now:

```
<expr>  => <expr> + <term>
        => <term> + <term>
        => <factor> + <term>
        => id + <term>
        => id + <term> * <factor>
        => id + <factor> * <factor>
        => id + id * <factor>
        => id + id * id
```



**Grammar:**

```
<expr>    ::= <expr> + <term> | <expr> - <term> | <term>
<term>    ::= <term> * <factor> | <term> / <factor> | <factor>
<factor>  ::= num | id
```

# AMBIGUITY

A grammar that produces more than one parse tree for some sentence is said to be *ambiguous*.

**Example:**

$$\begin{array}{lcl} \langle \text{stmt} \rangle & ::= & \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \\ & & | \quad \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle \\ & & | \quad \text{other stmts} \end{array}$$

**Consider the following statement:**

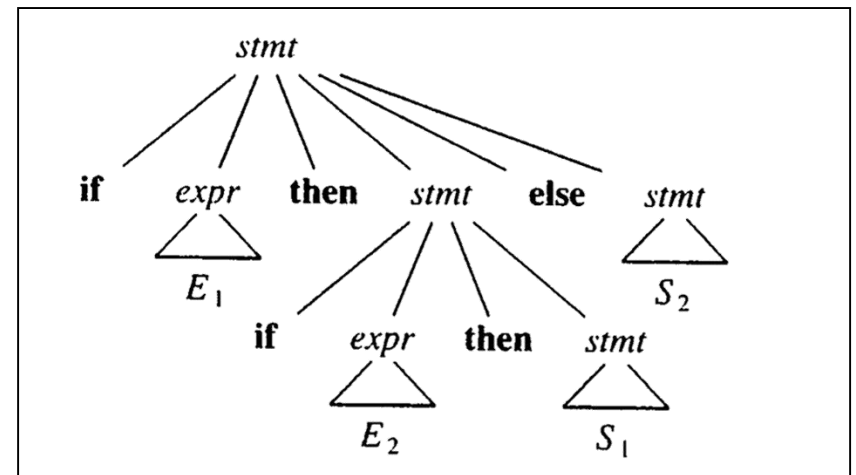
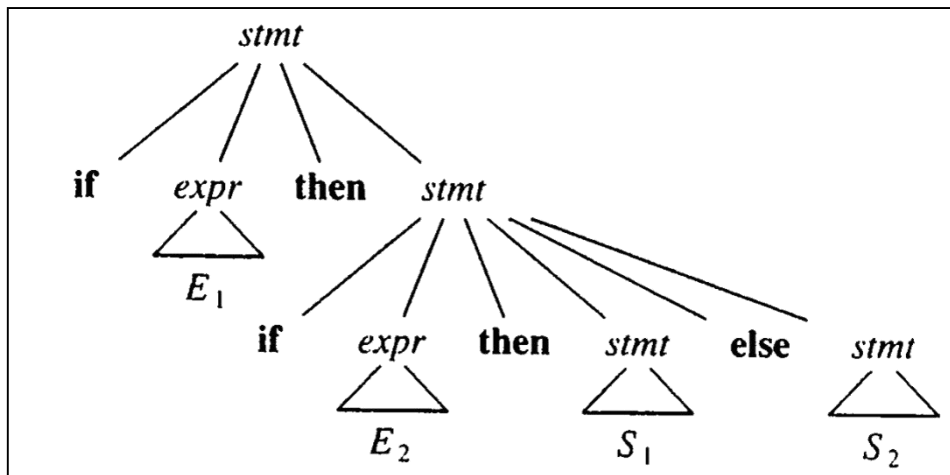
if  $E_1$  then if  $E_2$  then  $S_1$  else  $S_2$

**It has two derivations**

- It is a context free ambiguity

# AMBIGUITY

A grammar that produces more than one parse tree for some sentence is said to be *ambiguous*



# ELIMINATING AMBIGUITY

**Sometimes an ambiguous grammar can be rewritten to eliminate the ambiguity.**

- E.g. “match each else with the closest unmatched then”
- This is most likely the intention of the programmer

```
⟨stmt⟩      ::=  ⟨matched⟩  
              |  ⟨unmatched⟩  
⟨matched⟩   ::=  if ⟨expr⟩ then ⟨matched⟩ else ⟨matched⟩  
              |  other stmts  
⟨unmatched⟩ ::=  if ⟨expr⟩ then ⟨stmt⟩  
              |  if ⟨expr⟩ then ⟨matched⟩ else ⟨unmatched⟩
```

# MAPPING THIS TO A JAVA EXAMPLE

In Java, the grammar rules are slightly different than the previous example

Below is a (very simplified) version of these rules

```
<stmt>          ::= <matched>
                  | <unmatched>

<matched>        ::= if ( <expr> ) <matched> else <matched>
                  | other stmts

<unmatched>      ::= if ( <expr> ) < stmt >
                  | if ( <expr> ) <matched> else <unmatched>
```

# MAPPING THIS TO A JAVA EXAMPLE

For the following piece of code

```
if (x==0)
    if (y==0)
        z = 0;
    else
        z = 1;
```

After running the lexical analyser, we get the following list of tokens:

```
if ( id == num ) if (id == num) id = num ; else id = num ;
```

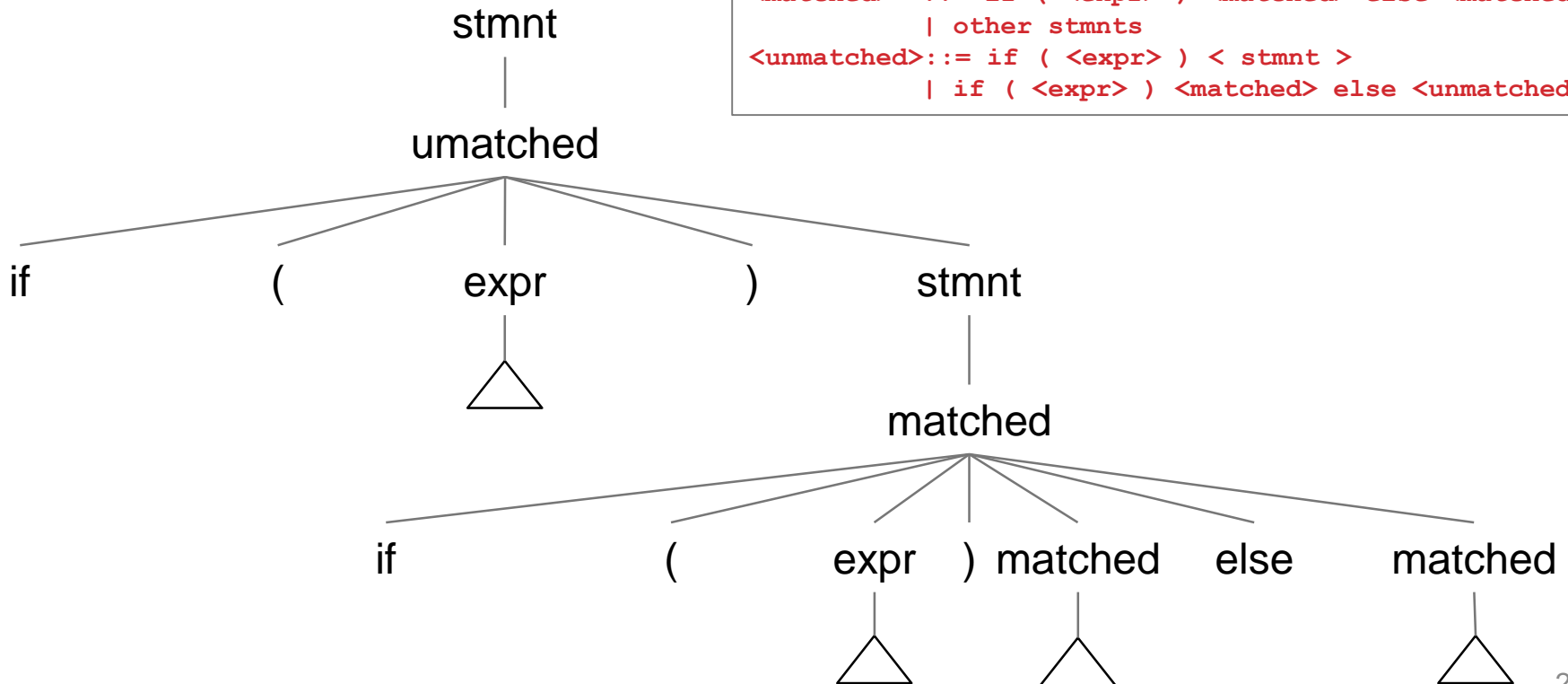
# MAPPING THIS TO A JAVA EXAMPLE

## Token input string

```
if ( id == num ) if (id == num) id = num ; else id = num ;
```

## Grammar

```
<stmt>      ::= <matched>  
              | <unmatched>  
<matched>   ::= if ( <expr> ) <matched> else <matched>  
              | other stmts  
<unmatched> ::= if ( <expr> ) <stmt>  
              | if ( <expr> ) <matched> else <unmatched>
```



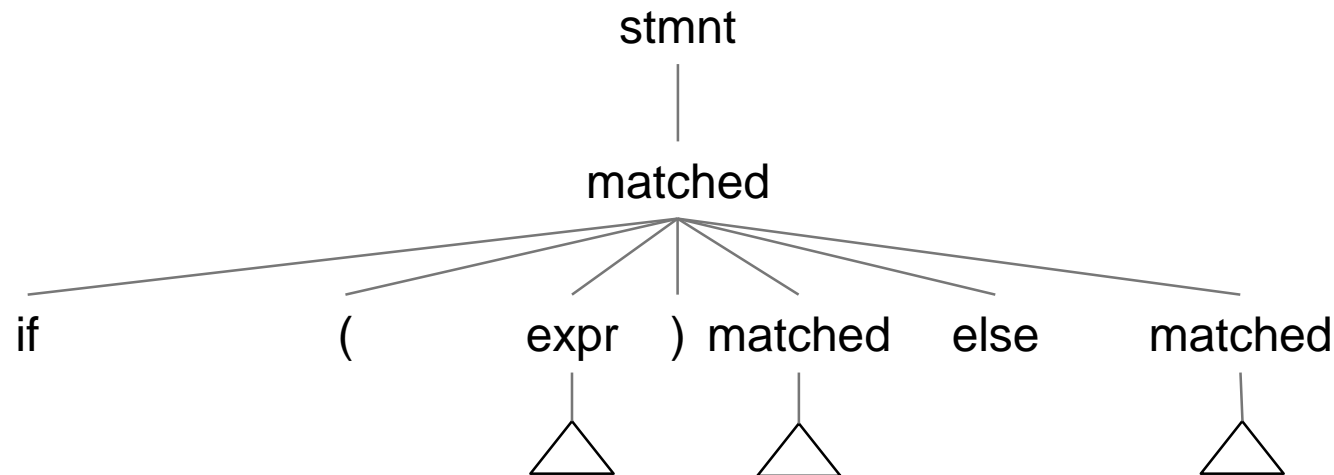
# MAPPING THIS TO A JAVA (ANOTHER) EXAMPLE

## Token input string

```
if ( id == num ) else id = num ;
```

## Grammar

```
<stmt>      ::= <matched>  
              | <unmatched>  
<matched>   ::= if ( <expr> ) <matched> else <matched>  
              | other stmts  
<unmatched> ::= if ( <expr> ) <stmt>  
              | if ( <expr> ) <matched> else <unmatched>
```





# TOP DOWN PARSING

A top-down parser starts with the root of the parse tree, labelled with the start or goal symbol of the grammar

To build a parse tree, we repeat the following steps until the leaves of the parse tree match the input string

1. At a node labelled  $A$ , select a production  $A ::= \alpha$  and construct the appropriate child for each symbol of  $\alpha$
2. When a terminal is added to the parse tree that does not match the input string, backtrack
3. Find the next nonterminal to be expanded

# TOP DOWN PARSING

Top-down parsing can be viewed as an attempt to find a leftmost derivation for an input string

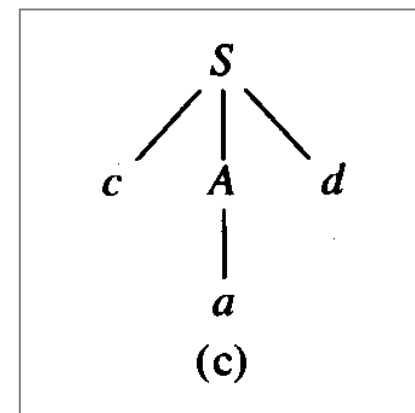
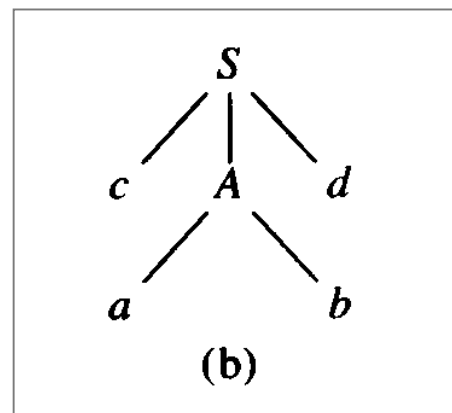
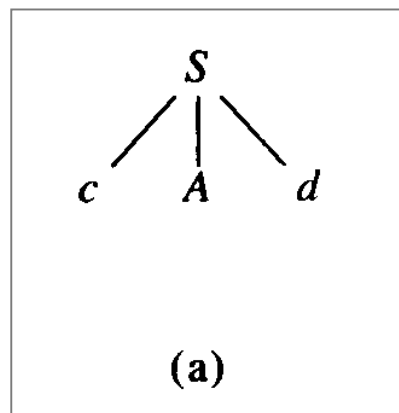
Example:

**Grammar:**

$\langle S \rangle ::= c\langle A \rangle d$   
 $\langle A \rangle ::= ab \mid a$

**Input string**

cad



*We need to backtrack!*

# EXPRESSION GRAMMAR

Recall our grammar for simple expressions:

$$\begin{array}{lcl} \langle \text{expr} \rangle & ::= & \langle \text{expr} \rangle + \langle \text{term} \rangle \\ & & | \langle \text{expr} \rangle - \langle \text{term} \rangle \\ & & | \langle \text{term} \rangle \\ \langle \text{term} \rangle & ::= & \langle \text{term} \rangle * \langle \text{factor} \rangle \\ & & | \langle \text{term} \rangle / \langle \text{factor} \rangle \\ & & | \langle \text{factor} \rangle \\ \langle \text{factor} \rangle & ::= & \text{num} \\ & & | \text{id} \end{array}$$

Consider the input string:

`id - num * id`

# EXAMPLE



Prod'n	Sentential form	Input
1	$\langle \text{expr} \rangle$	$\uparrow \text{id} \quad - \quad \text{num} \quad * \quad \text{id}$

## Reference Grammar

$\langle \text{expr} \rangle$	$::=$	$\langle \text{expr} \rangle + \langle \text{term} \rangle$
		$\langle \text{expr} \rangle - \langle \text{term} \rangle$
		$\langle \text{term} \rangle$
$\langle \text{term} \rangle$	$::=$	$\langle \text{term} \rangle * \langle \text{factor} \rangle$
		$\langle \text{term} \rangle / \langle \text{factor} \rangle$
		$\langle \text{factor} \rangle$
$\langle \text{factor} \rangle$	$::=$	<code>num</code>
		<code>id</code>

# EXAMPLE



## Reference Grammar

```

<expr> ::= <expr> + <term>
        | <expr> - <term>
        | <term>
<term>  ::= <term> * <factor>
        | <term> / <factor>
        | <factor>
<factor> ::= num
         | id
    
```

Prod'n	Sentential form	Input
1	<expr>	↑ id - num * id
2	<expr> + <term>	↑ id - num * id
4	<term> + <term>	↑ id - num * id
7	<factor> + <term>	↑ id - num * id
9	id + <term>	↑ id - num * id
-	id + <term>	id ↑ - num * id
-	<expr>	↑ id - num * id
3	<expr> - <term>	↑ id - num * id
4	<term> - <term>	↑ id - num * id
7	<factor> - <term>	↑ id - num * id
9	id - <term>	↑ id - num * id
-	id - <term>	id ↑ - num * id
-	id - <term>	id - ↑ num * id
7	id - <factor>	id - ↑ num * id
8	id - num	id - ↑ num * id
-	id - num	id - num ↑ * id
-	id - <term>	id - ↑ num * id
5	id - <term> * <factor>	id - ↑ num * id
7	id - <factor> * <factor>	id - ↑ num * id
8	id - num * <factor>	id - ↑ num * id
-	id - num * <factor>	id - num ↑ * id
-	id - num * <factor>	id - num * ↑ id
9	id - num * id	id - num * ↑ id
-	id - num * id	id - num * id ↑

# EXAMPLE

## Another possible parse for id – num \* id

Prod'n	Sentential form	Input
1	$\langle \text{expr} \rangle$	$\uparrow \text{id} - \text{num} * \text{id}$
2	$\langle \text{expr} \rangle + \langle \text{term} \rangle$	$\uparrow \text{id} - \text{num} * \text{id}$
2	$\langle \text{expr} \rangle + \langle \text{term} \rangle + \langle \text{term} \rangle$	$\uparrow \text{id} - \text{num} * \text{id}$
2	$\langle \text{expr} \rangle + \langle \text{term} \rangle + \dots$	$\uparrow \text{id} - \text{num} * \text{id}$
2	$\langle \text{expr} \rangle + \langle \text{term} \rangle + \dots$	$\uparrow \text{id} - \text{num} * \text{id}$
2	$\dots$	$\uparrow \text{id} - \text{num} * \text{id}$

### Reference Grammar

$\langle \text{expr} \rangle$	$::=$	$\langle \text{expr} \rangle + \langle \text{term} \rangle$
		$\langle \text{expr} \rangle - \langle \text{term} \rangle$
		$\langle \text{term} \rangle$
$\langle \text{term} \rangle$	$::=$	$\langle \text{term} \rangle * \langle \text{factor} \rangle$
		$\langle \text{term} \rangle / \langle \text{factor} \rangle$
		$\langle \text{factor} \rangle$
$\langle \text{factor} \rangle$	$::=$	$\text{num}$
		$\text{id}$

**If the parser makes the wrong choices, expansion does not terminate**

- This is not a good property for a parser to have
- Parsers should terminate, eventually...

# LEFT RECURSION

A grammar is left recursive if:

“It has a nonterminal **A** such that there is a derivation  $A \xRightarrow{+} A\alpha$  for some string  $\alpha$ ”

*Top down parsers with left derivation  
cannot handle left-recursion in a grammar*

# ELIMINATING LEFT RECURSION

Consider the grammar fragment:

$$\begin{array}{lcl} \langle \text{foo} \rangle & ::= & \langle \text{foo} \rangle \alpha \\ & | & \beta \end{array}$$

Where  $\alpha$  and  $\beta$  do not start with  $\langle \text{foo} \rangle$

We can re-write this as:

$$\begin{array}{lcl} \langle \text{foo} \rangle & ::= & \beta \langle \text{bar} \rangle \\ \langle \text{bar} \rangle & ::= & \alpha \langle \text{bar} \rangle \\ & | & \epsilon \end{array}$$

Where  $\langle \text{bar} \rangle$  is a new non-terminal

This Fragment contains no left recursion



# ELIMINATING LEFT RECURSION

Similarly,

$$\begin{aligned} \langle \text{foo} \rangle &::= \langle \text{foo} \rangle \alpha_1 \\ &| \langle \text{foo} \rangle \alpha_2 \\ &\vdots \\ &| \langle \text{foo} \rangle \alpha_n \\ &| \beta_1 \\ &| \beta_2 \\ &\vdots \\ &| \beta_n \end{aligned}$$

Can be re-written as:

$$\begin{aligned} \langle \text{foo} \rangle &::= \beta_1 \langle \text{bar} \rangle \\ &| \beta_2 \langle \text{bar} \rangle \\ &\vdots \\ &| \beta_n \langle \text{bar} \rangle \\ \langle \text{bar} \rangle &::= \alpha_1 \langle \text{bar} \rangle \\ &| \alpha_2 \langle \text{bar} \rangle \\ &\vdots \\ &| \alpha_n \langle \text{bar} \rangle \\ &| \epsilon \end{aligned}$$

# EXAMPLE

Our expression grammar contains two cases of left-recursion

$$\begin{aligned} \langle \text{expr} \rangle &::= \langle \text{expr} \rangle + \langle \text{term} \rangle \\ &\quad | \langle \text{expr} \rangle - \langle \text{term} \rangle \\ &\quad | \langle \text{term} \rangle \\ \langle \text{term} \rangle &::= \langle \text{term} \rangle * \langle \text{factor} \rangle \\ &\quad | \langle \text{term} \rangle / \langle \text{factor} \rangle \\ &\quad | \langle \text{factor} \rangle \end{aligned}$$

Applying the transformation gives

$$\begin{aligned} \langle \text{expr} \rangle &::= \langle \text{term} \rangle \langle \text{expr}' \rangle \\ \langle \text{expr}' \rangle &::= + \langle \text{term} \rangle \langle \text{expr}' \rangle \\ &\quad | - \langle \text{term} \rangle \langle \text{expr}' \rangle \\ &\quad | \epsilon \\ \langle \text{term} \rangle &::= \langle \text{factor} \rangle \langle \text{term}' \rangle \\ \langle \text{term}' \rangle &::= * \langle \text{factor} \rangle \langle \text{term}' \rangle \\ &\quad | / \langle \text{factor} \rangle \langle \text{term}' \rangle \\ &\quad | \epsilon \end{aligned}$$

With this grammar, a top-down parser will

- Terminate (for sure)
- Backtrack on some inputs

# LOOK AHEAD...

**We saw that top-down parsers may need to backtrack when they select the wrong production**

**Therefore, we might need to look ahead in order to avoid backtracking**

**This is where predictive parsers come in handy**

- LL(1): left to right scan, left-most derivation, 1-token look ahead
- LR(1): left to right scan, right most derivation, 1-token look ahead

# **THANK YOU!**

## **QUESTIONS?**