



Programming Languages

Lecture05 – Top-Down Parsing



남 범 석

bnam@skku.edu

Parsing Techniques

❖ *Top-down parsers*

- *LL = Left-to-right input scan, Leftmost derivation*
- Start at the root of the parse tree and grow toward leaves
- Pick a production & try to match the input
- Bad “pick” \Rightarrow may need to backtrack
- Some grammars are backtrack-free (*predictive parsing*)

❖ *Bottom-up parsers*

- *LR = Left-to-right input scan, Rightmost derivation*
- Start at the leaves and grow toward root
- As input is consumed, encode possibilities in an internal state
- Start in a state valid for legal first tokens
- Bottom-up parsers handle a large class of grammars



Top-down Parser

- Problems in Top-down parser
 - Backtrack \Rightarrow predictive parser
 - Left-recursion \Rightarrow may result in infinite loop
- Predictive Top-down parser
 - LL(1)
 - 1 means “predict based on k tokens of lookahead”
 - Left factoring transforms some non-LL(1) to LL(1)

CFG with precedence

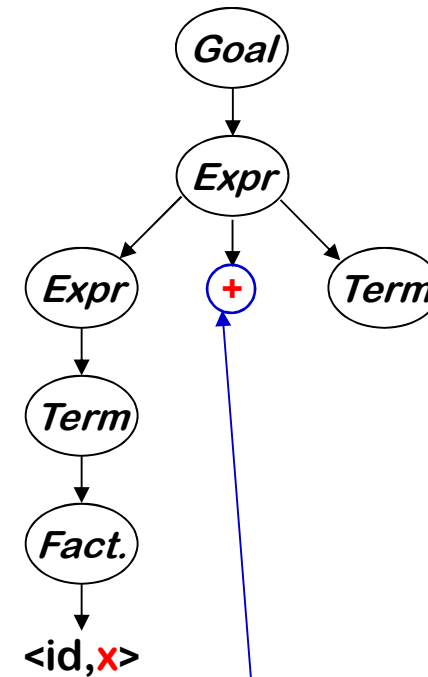
1	$Goal \rightarrow Expr$
2	$Expr \rightarrow Expr + Term$
3	$Expr - Term$
4	$Term$
5	$Term \rightarrow Term * Factor$
6	$Term / Factor$
7	$Factor$
8	$Factor \rightarrow \underline{number}$
9	\underline{id}

And the input $x - 2 * y$

Top-down Parser – backtrack (1)

- Let's try $x - 2 * y$: *Leftmost derivation, choose productions in an order that exposes problems*

Rule	Sentential Form	Input
—	Goal	$\uparrow x - 2 * y$
1	Expr	$\uparrow x - 2 * y$
2	Expr + Term	$\uparrow x - 2 * y$
4	Term + Term	$\uparrow x - 2 * y$
7	Factor + Term	$\uparrow x - 2 * y$
9	$\langle id, x \rangle + Term$	$\uparrow x - 2 * y$
—	$\langle id, x \rangle + Term$	$x \uparrow - 2 * y$

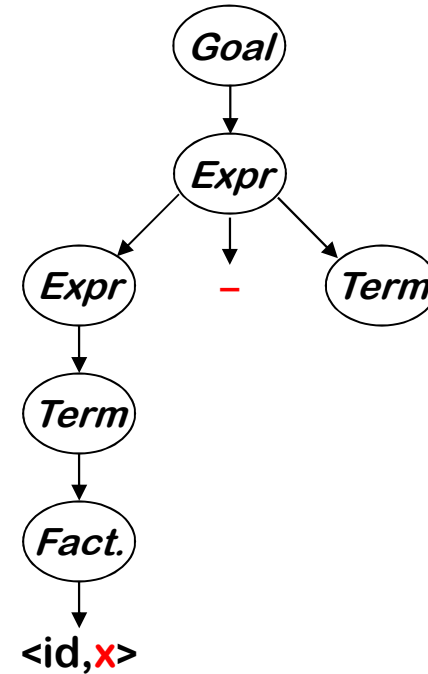


- This worked well, except that “-” doesn't match “+”
- The parser must backtrack to here

Top-down Parser – backtrack (2)

- Continuing with $x - 2 * y$:

Rule	Sentential Form	Input
—	<i>Goal</i>	$\uparrow x - 2 * y$
1	<i>Expr</i>	$\uparrow x - 2 * y$
3	<i>Expr - Term</i>	$\uparrow x - 2 * y$
4	<i>Term - Term</i>	$\uparrow x - 2 * y$
7	<i>Factor - Term</i>	$\uparrow x - 2 * y$
9	$\langle id, x \rangle - Term$	$\uparrow x - 2 * y$
9	$\langle id, x \rangle - Term$	$x \uparrow - 2 * y$
—	$\langle id, x \rangle - Term$	$x - \uparrow 2 * y$



This time, “-” and “-” matched

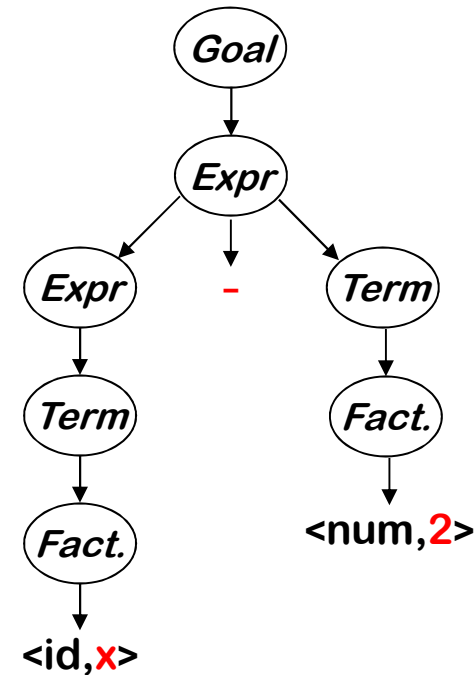
We can advance past “-” to look at “2”

⇒ Now, we need to expand *Term* - the last *NT* on the fringe

Top-down Parser – backtrack (3)

- Trying to match the "2" in $x - 2 * y$:

Rule	Sentential Form	Input
—	$\langle id, x \rangle - Term$	$x - \uparrow 2 * y$
7	$\langle id, x \rangle - Factor$	$x - \uparrow 2 * y$
9	$\langle id, x \rangle - \langle num, 2 \rangle$	$x - \uparrow 2 * y$
—	$\langle id, x \rangle - \langle num, 2 \rangle$	$x - 2 \uparrow * y$

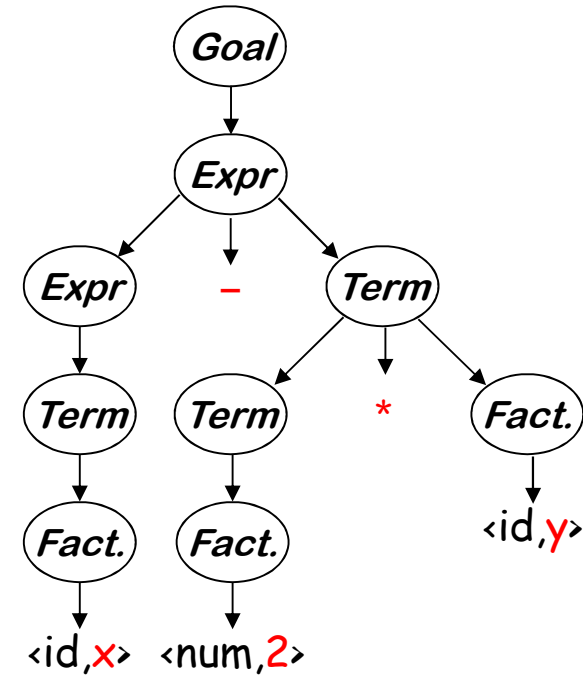


- Where are we?
 - "2" matches "2"
 - We have more input, but no *NTs* left to expand
 - The expansion terminated too soon
 ⇒ Need to backtrack

Top-down Parser – backtrack (4)

- Trying again with "2" in $x - 2 * y$:

Rule	Sentential Form	Input
—	$\langle id, x \rangle - Term$	$\underline{x} - \uparrow \underline{2} * y$
5	$\langle id, x \rangle - Term * Factor$	$\underline{x} - \uparrow \underline{2} * y$
7	$\langle id, x \rangle - Factor * Factor$	$\underline{x} - \uparrow \underline{2} * y$
8	$\langle id, x \rangle - \langle num, 2 \rangle * Factor$	$\underline{x} - \uparrow \underline{2} * y$
—	$\langle id, x \rangle - \langle num, 2 \rangle * Factor$	$\underline{x} - \underline{2} \uparrow * y$
—	$\langle id, x \rangle - \langle num, 2 \rangle * Factor$	$\underline{x} - \underline{2} * \uparrow y$
9	$\langle id, x \rangle - \langle num, 2 \rangle * \langle id, y \rangle$	$\underline{x} - \underline{2} * \uparrow y$
—	$\langle id, x \rangle - \langle num, 2 \rangle * \langle id, y \rangle$	$\underline{x} - \underline{2} * y \uparrow$



- This time, we matched & consumed all the input
 \Rightarrow Success!



Left Recursion

- Top-down parsers cannot handle left-recursive grammars
 - $S \rightarrow S\alpha \mid \beta$
 - Formally,
A grammar is left recursive if $\exists A \in NT$ such that
 \exists a derivation $A \Rightarrow^+ A\alpha$, for some string $\alpha \in (NT \cup T)^+$
- Our expression grammar is left recursive
 - This can lead to non-termination in a top-down parser
 - We would like to convert the left recursion to right recursion

Eliminating Left Recursion

- Remove left recursion

- Original grammar

$$F \rightarrow F \alpha$$
$$| \beta$$

where neither α nor β starts with F

- Rewrite the above as

$$F \rightarrow \beta P$$
$$P \rightarrow \alpha P$$
$$| \varepsilon$$

where P is a new non-terminal

- *Accepts the same language, but uses only right recursion*

Eliminating Left Recursion

- The expression grammar contains two cases of left recursion

$$\begin{array}{l} \text{Expr} \rightarrow \text{Expr} + \text{Term} \\ \quad | \text{Expr} - \text{Term} \\ \quad | \text{Term} \end{array} \qquad \begin{array}{l} \text{Term} \rightarrow \text{Term} * \text{Factor} \\ \quad | \text{Term} / \text{Factor} \\ \quad | \text{Factor} \end{array}$$

- Applying the transformation yields

$$\begin{array}{l} \text{Expr} \rightarrow \text{Term Expr}' \\ \text{Expr}' \rightarrow + \text{Term Expr}' \\ \quad | - \text{Term Expr}' \\ \quad | \varepsilon \end{array} \qquad \begin{array}{l} \text{Term} \rightarrow \text{Factor Term}' \\ \text{Term}' \rightarrow * \text{Factor Term}' \\ \quad | / \text{Factor Term}' \\ \quad | \varepsilon \end{array}$$

- These fragments use only right recursion
- They retain the original left associativity (evaluate left to right)



Predictive Top-down Parsing

- Basic idea
 - Given $A \rightarrow \alpha \mid \beta$,
the parser must be able to choose between α or β
 - First and Follow sets
 - Simple recursive descent parsers
 - Table-driven LL(I) parsers

- LL(I) Parser
 - L = scan input left to right
 - L = Leftmost derivation
 - 1 = lookahead is enough to pick right production rule to use
 - Given production rules
 $S \rightarrow \alpha \mid \beta$
the parser should be able to choose between α or β using one
lookahead
 - No Backtracking
 - No Left Recursion

First Sets

- For some rhs $\alpha \in G$,
FIRST(α) is the set of terminals that appear as the first symbol in some string that derives from α
 - $x \in \text{FIRST}(\alpha)$ iff $\alpha \Rightarrow^* x\gamma$, for some γ
 - Some number of derivations gets us x at the beginning
- Goal \rightarrow SheepNoise
- SheepNoise \rightarrow SheepNoise baa | baa

For SheepNoise:

FIRST(Goal) = { baa }

FIRST(SheepNoise) = { baa }

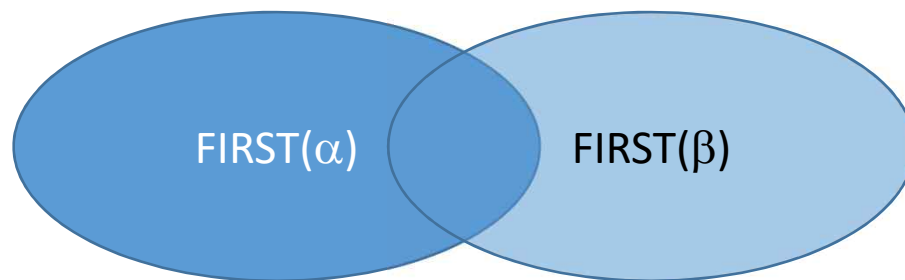
FIRST(baa) = { baa }

The LL(1) Property (first version)

- If $A \rightarrow \alpha$ and $A \rightarrow \beta$ both appear in the grammar, we would like

$$\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \emptyset$$

- This would allow the parser to make a correct choice with a lookahead of exactly one symbol !



Does not have LL(1) Property

What about ε -productions?

- If $A \rightarrow \alpha$ and $A \rightarrow \beta$ and $\varepsilon \in \text{FIRST}(\alpha)$, then we need to ensure

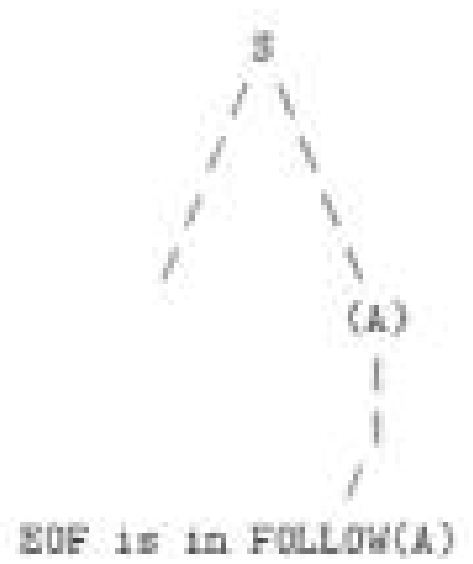
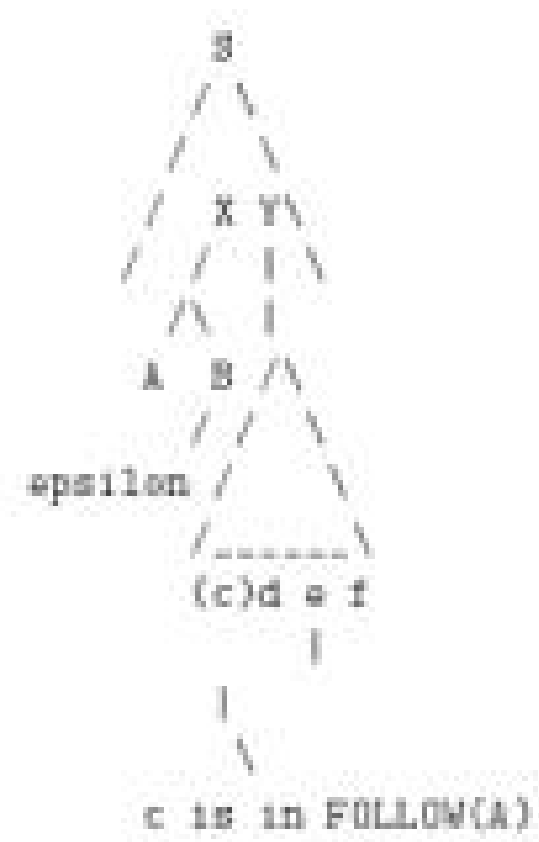
$$\text{FIRST}(\beta) \cap \text{FOLLOW}(A) = \emptyset$$

where,

$\text{FOLLOW}(A)$ = the set of terminal symbols that can immediately follow A in a sentential form

- Formally,
 $\text{FOLLOW}(A) = \{t \mid (t \text{ is a terminal and } G \Rightarrow^* \alpha A t \beta)$
or $(t \text{ is eof and } G \Rightarrow^* \alpha A)\}$

Follow Sets Intuition



FIRST⁺ sets

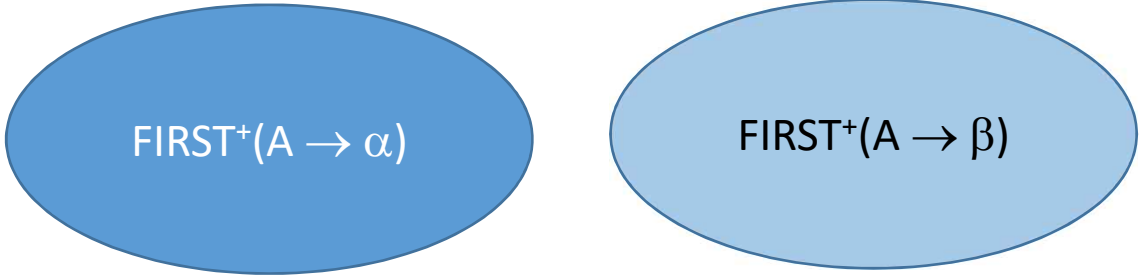
- Definition of FIRST⁺(A → α)
if $\varepsilon \in \text{FIRST}(\alpha)$

$$\text{FIRST}^+(A \rightarrow \alpha) = \text{FIRST}(\alpha) \cup \text{FOLLOW}(A)$$

else

$$\text{FIRST}^+(A \rightarrow \alpha) = \text{FIRST}(\alpha)$$

- Grammar is LL(1) iff $A \rightarrow \alpha$ and $A \rightarrow \beta$ implies
 $\text{FIRST}^+(A \rightarrow \alpha) \cap \text{FIRST}^+(A \rightarrow \beta) = \emptyset$



FIRST⁺(A → α)

FIRST⁺(A → β)



Example:

Grammar

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \varepsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \varepsilon$$

$$F \rightarrow (E) \mid \text{id}$$

First and Follow

Symbol	FIRST	FOLLOW
E	(,id	\$,)
E'	+, ϵ	\$,)
T	(,id	+, \$,)
T'	*, ϵ	+, \$,)
F	(,id	*, +, \$,)



What if my grammar is not LL(1)?

- Can we transform a non-LL(1) grammar into an LL(1) grammar?
- In general, the answer is no
- In some cases, however, the answer is yes
- Perform:
 - Eliminate left-recursion (Slide #10)
 - Perform left factoring

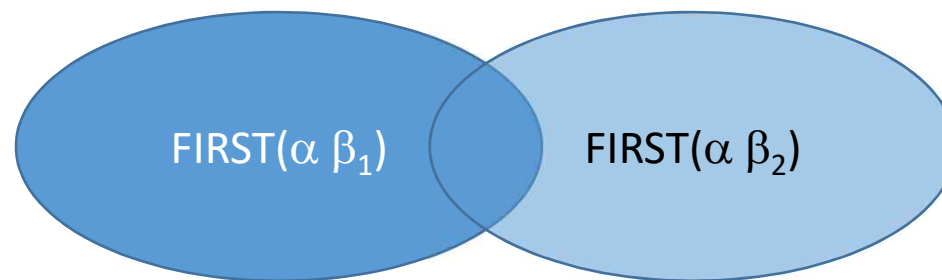
What if my grammar is not LL(1)?

- Given grammar G with productions

$$A \rightarrow \alpha \beta_1$$

$$A \rightarrow \alpha \beta_2$$

if α derives anything other than epsilon and
 $\text{FIRST}^+(A \rightarrow \alpha \beta_1) \cap \text{FIRST}^+(A \rightarrow \alpha \beta_2) \neq \emptyset$



This grammar is not LL(1)

Left Factoring

- If we pull the common prefix, a , into a separate production, we may make the grammar LL(1).

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1$$

$$| \beta_2$$

Now, if $\text{FIRST}^+(A' \rightarrow \beta_1) \cap \text{FIRST}^+(A' \rightarrow \beta_2) = \emptyset$,
G may be LL(1)

Left Factoring

$\forall A \in NT$,
find the longest prefix α that occurs in two
or more right-hand sides of A

if $\alpha \neq \epsilon$ then replace all of the A productions,

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n \mid \gamma,$$

with

$$A \rightarrow \alpha Z \mid \gamma$$

$$Z \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

where Z is a new element of NT

Repeat until no common prefixes remain

$$\begin{array}{l} A \rightarrow \alpha\beta_1 \\ \mid \alpha\beta_2 \\ \mid \alpha\beta_3 \end{array}$$



$$\begin{array}{l} A \rightarrow \alpha Z \\ Z \rightarrow \beta_1 \\ \mid \beta_2 \\ \mid \beta_3 \end{array}$$

Left Factoring

(An example)

- Consider the following fragment of the expression grammar

$Factor \rightarrow \underline{Identifier}$
 $\quad \quad | \underline{Identifier} [ExprList]$
 $\quad \quad | \underline{Identifier} (ExprList)$

$FIRST^+(rhs_1) = \{ \underline{Identifier} \}$
 $FIRST^+(rhs_2) = \{ \underline{Identifier} \}$
 $FIRST^+(rhs_3) = \{ \underline{Identifier} \}$

- After left factoring, it becomes

$Factor \rightarrow \underline{Identifier} Arguments$
 $Arguments \rightarrow [ExprList]$
 $\quad \quad | (ExprList)$
 $\quad \quad | \varepsilon$

$FIRST^+(rhs_1) = \{ []$
 $FIRST^+(rhs_2) = \{ ()$
 $FIRST^+(rhs_3) = \{ \varepsilon \} \cup FOLLOW(Factor)$
 \Rightarrow It has the $LL(1)$ property

- This form has the same syntax, with the $LL(1)$ property

Left Recursion & Left Factoring (Generality)

Question

By *eliminating left recursion* and *left factoring*, can we transform an arbitrary CFG to a form where it meets the $LL(1)$ condition? (and can be parsed predictively with a single token lookahead?)

Answer

Given a CFG that doesn't meet the $LL(1)$ condition, it is undecidable whether or not an equivalent $LL(1)$ grammar exists.

Language that Cannot Be LL(1)

Example that has no LL(1) grammar

$$\{a^n 0 b^n \mid n \geq 1\} \cup \{a^n 1 b^{2n} \mid n \geq 1\}$$

$$G \rightarrow \underline{a} A \underline{b}$$
$$| \underline{a} B \underline{bb}$$

$$A \rightarrow \underline{a} A \underline{b}$$
$$| \underline{0}$$

$$B \rightarrow \underline{a} B \underline{bb}$$
$$| \underline{1}$$

Problem: need an unbounded number of a characters before you can determine whether you are in the A group or the B group.

Automate Predictive Parsing

- Given a grammar that has the $LL(1)$ property
 - Can write a simple routine to recognize each lhs
 - Code is both simple & fast
- Consider $A \rightarrow \beta_1 \mid \beta_2 \mid \beta_3$, with
 - $FIRST^+(\beta_1) \cap FIRST^+(\beta_2) = \emptyset$
 - $FIRST^+(\beta_2) \cap FIRST^+(\beta_3) = \emptyset$
 - $FIRST^+(\beta_1) \cap FIRST^+(\beta_3) = \emptyset$

```
/* find an A */  
if (current_token  $\in$   $FIRST^+(\beta_1)$ )  
    find a  $\beta_1$  and return true  
else if (current_token  $\in$   $FIRST^+(\beta_2)$ )  
    find a  $\beta_2$  and return true  
else if (current_token  $\in$   $FIRST^+(\beta_3)$ )  
    find a  $\beta_3$  and return true  
else  
    report an error and return false
```

Grammars with the $LL(1)$ property are called predictive grammars because the parser can “predict” the correct expansion at each point in the parse.

Parsers that capitalize on the $LL(1)$ property are called predictive parsers.

One kind of predictive parser is the recursive descent parser.

Predictive Parsing Example

■ Expression grammar, after transformation

1	<i>Goal</i>	→	<i>Expr</i>
2	<i>Expr</i>	→	<i>Term Expr'</i>
3	<i>Expr'</i>	→	+ <i>Term Expr'</i>
4			- <i>Term Expr'</i>
5			ϵ
6	<i>Term</i>	→	<i>Factor Term'</i>
7	<i>Term'</i>	→	* <i>Factor Term'</i>
8			/ <i>Factor Term'</i>
9			ϵ
10	<i>Factor</i>	→	<u>id</u>
11			<u>number</u>

This produces a parser with six mutually recursive routines:

- *Goal*
- *Expr*
- *EPrime*
- *Term*
- *TPrime*
- *Factor*

Each recognizes one *NT* or *T*

The term descent refers to the direction in which the parse tree is built.

Recursive Descent Parser

- A couple of routines from the expression parser

Goal()

```
token ← next_token();  
if (Expr() = true & token = EOF)  
  then next compilation step;  
  else  
    report syntax error;  
    return false;
```

Expr()

```
if (Term() = false)  
  then return false;  
  else return Eprime();
```

looking for EOF,
found other
token

Factor()

```
if (token = Number) then  
  token ← next_token();  
  return true;  
else if (token = Identifier) then  
  token ← next_token();  
  return true;  
else  
  report syntax error;  
  return false;
```

global variable

EPrime, *Term*, & *TPRime* follow the
same basic lines

looking for Number or Identifier,
found other token instead

Parse Tree - Recursive Descent Parser

- To build a parse tree:
 - Augment parsing routines to build nodes
 - Pass nodes between routines using a stack
 - Node for each symbol on *rhs*
 - Action is to pop *rhs* nodes, make them children of *lhs* node, and push this subtree
- To build an Abstract Syntax Tree (AST)
 - Non-terminals are not present in the AST
 - Build fewer nodes
 - Put them together in a different order

$Expr \rightarrow Term Expr'$

```
Expr()  
  result ← true;  
  if (Term() = false)  
    then return false;  
    else if (EPrime() = false)  
      then result ← false;  
      else  
        build an Expr node  
        pop EPrime node  
        pop Term node  
        make EPrime & Term  
          children of Expr  
        push Expr node  
  return result;
```

Success ⇒ build a piece of the parse tree

Building Table-driven Parser

■ Strategy

- Encode knowledge in a table
- Use a standard “skeleton” parser to interpret the table

■ Example

- The non-terminal Factor has two expansions
 - Identifier or Number
- Need a row for every NT & a column for every T
- Table might look like:

Terminal Symbols

	+	-	*	/	Id.	Num	EOF
<u>Factor</u>	⊖	-	-	-	10	11	-

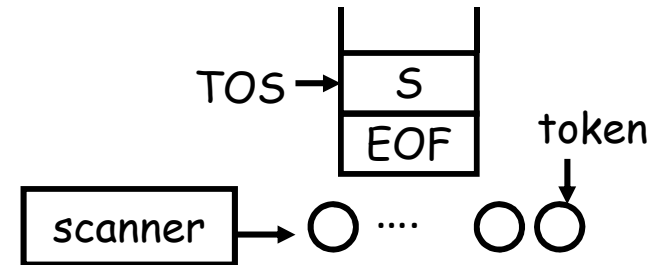
Non-terminal Symbols

Error on '+'

Reduce by rule 10 on 'x'

LL(1) Skeleton Parser

```
token ← next_token()
push EOF onto Stack
push the start symbol, S, onto Stack
TOS ← top of Stack
loop forever
  if TOS = EOF and token = EOF then
    break & report success
  else if TOS is a terminal then
    if TOS matches token then
      pop Stack
      token ← next_token()
    else report error looking for TOS
  else
    if TABLE[TOS,token] is  $A \rightarrow B_1 B_2 \dots B_k$  then
      pop Stack
      push  $B_k, B_{k-1}, \dots, B_1$ 
    else report error expanding TOS
TOS ← top of Stack
```



exit on success

// recognized TOS

// TOS is a non-terminal

// get rid of A
// in that order



Building LL(1) table

- Building the complete table for LL(1)
 - Need a row for every NT & a column for every T
 - Need an algorithm to build the table
- Filling in $TABLE[X,y]$, $X \in NT, y \in T \cup \{EOF\}$
 - entry is the rule $X \rightarrow \beta$, if $y \in FIRST^+(\beta)$
 - entry is error, otherwise
- If any entry is defined multiple times, G is not LL(1)

Example:

Grammar

$E \rightarrow TE'$

$E' \rightarrow +TE' \mid \varepsilon$

$T \rightarrow FT'$

$T' \rightarrow *FT' \mid \varepsilon$

$F \rightarrow (E) \mid id$

Symbol	FIRST	FOLLOW
E	(,id	\$,)
E'	+, ε	\$,)
T	(,id	+, \$,)
T'	*, ε	+, \$,)
F	(,id	*, +, \$,)

Building the table

	Id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$			$F \rightarrow (E)$		

Input=id+id*id

Stack	Input buffer
$\$E$	$id+id*id\$$
$\$E'T$	$Id+id*id\$$
$\$E'T'F$	$Id+id*id\$$
$\$E'T'id$	$Id+id*id\$$
$\$E'T'$	$+id*id\$$
$\$E'$	$+id*id\$$
$\$E'T+$	$+id*id\$$
$\$E'T$	$id*id\$$

Input=id+id*id (Cont'd)

Stack	Input Buffer
$\$E'T'F$	$id*id\$$
$\$E'T'id$	$id*id\$$
$\$E'T'$	$*id\$$
$\$E'T'F*$	$*id\$$
$\$E'T'F$	$id\$$
$\$E'T'id$	$id\$$
$\$E'T'$	$\$$
$\$E'$	$\$$
$\$$	Accepted

Grammar that requires LL(k) parsing

Grammar:

$S \rightarrow iEtSS' \mid a$

$S' \rightarrow eS \mid \varepsilon$

$E \rightarrow b$

	FIRST	FOLLOW
S	a,i	e, \$
S'	e, ε	e, \$
E	b	t

Note that this is

If then else statement

Parse Table

	a	b	e	i	t	\$
S	$S \rightarrow a$			$S \rightarrow iEtSS'$		
S'			$S' \rightarrow \epsilon$ $S' \rightarrow eS$			$S' \rightarrow \epsilon$
E		$E \rightarrow b$				

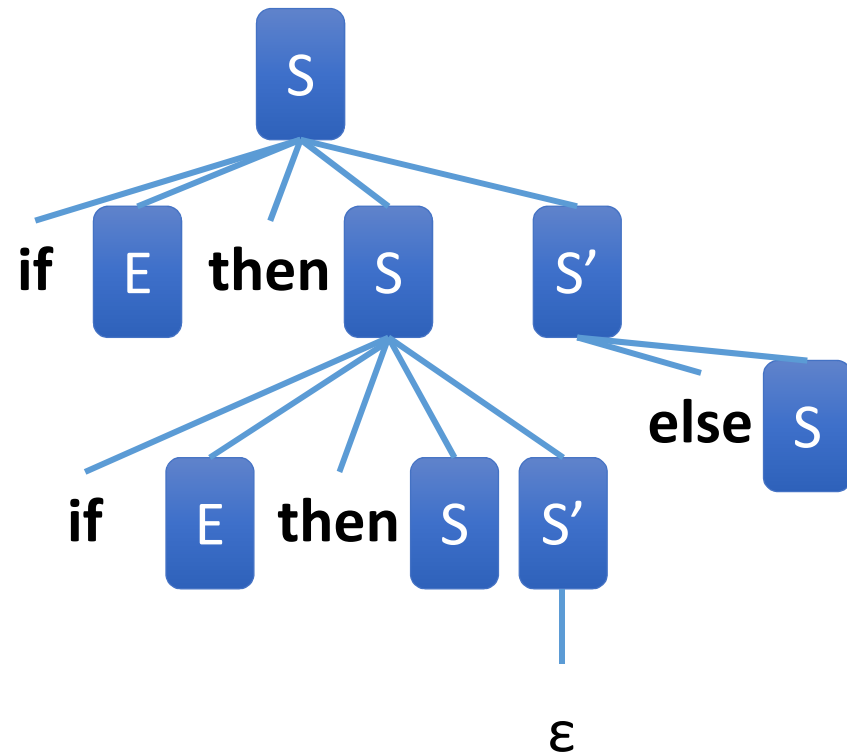
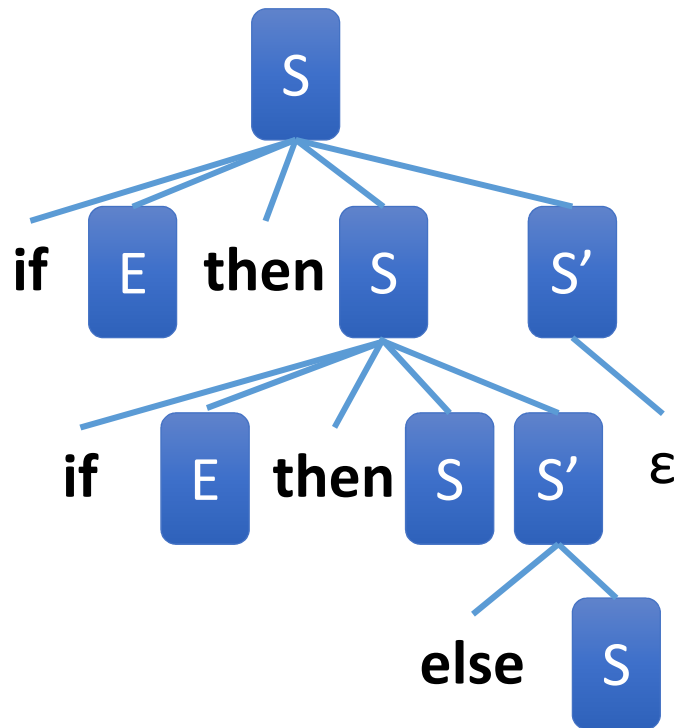
Ambiguity

Ambiguous grammar

$S \rightarrow iEtSS' \mid a$

$S' \rightarrow eS \mid \varepsilon$

$E \rightarrow b$



Ambiguous grammar and LL(2) Parsing

- The grammar is ambiguous and it is evident by the fact that we have two entries corresponding to $M[S',e]$ containing $S \rightarrow \epsilon$ and $S' \rightarrow eS$. This ambiguity can be resolved if we choose $S' \rightarrow eS$ i.e. associating the else's with the closest previous "then".
- Note that the ambiguity will be solved if we use LL(2) parser, i.e. always see for the two input symbols. How?
 - When input is 'e' then it looks at next input. Depending on the next input we choose the appropriate rule.

$S \rightarrow iES' \mid a$

$S' \rightarrow tS \mid tSeS$

$E \rightarrow b$



Summary

- Top-down parser
 - Use leftmost derivation
 - Bad pick of rewrite rule results in backtrack
- Left recursion removal
 - Avoid non-terminating top-down parser
- Predictive parsing
 - LL(1) property ensures only one production rule is chosen by looking ahead one terminal symbol.
- Left factoring
 - Transform some non-LL(1) to LL(1)
- Automatic top-down parser generation
 - Recursive decent parser
 - Building LL(1) table: $f(X,y) \rightarrow P$ (where $X \in NT, y \in T$)