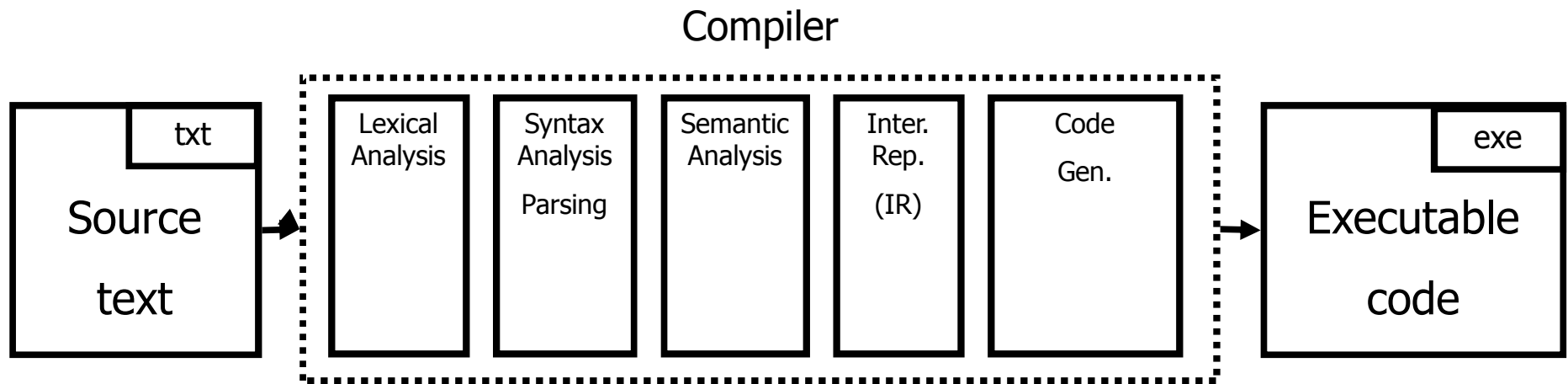


Lecture 03 – Syntax analysis: top-down parsing

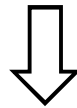
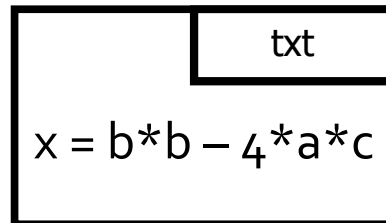
THEORY OF COMPILATION

Eran Yahav

You are here



Last Week: from characters to tokens



Token
Stream

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

Last Week: Regular Expressions

Basic Patterns	Matching
x	The character x
.	Any character, usually except a new line
[xyz]	Any of the characters x,y,z
Repetition Operators	
R?	An R or nothing (=optionally an R)
R*	Zero or more occurrences of R
R+	One or more occurrences of R
Composition Operators	
R ₁ R ₂	An R ₁ followed by R ₂
R ₁ R ₂	Either an R ₁ or R ₂
Grouping	
(R)	R itself

Parsing

- Goals
 - Is a sequence of tokens a valid program in the language?
 - Construct a structured representation of the input text
 - 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?

Context free grammars

$$G = (V, T, P, S)$$

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

Why do we need context free grammars?

$S \rightarrow SS$

$S \rightarrow (S)$

$S \rightarrow ()$

Example

$$S \rightarrow S;S$$
$$S \rightarrow id := E$$
$$E \rightarrow id \mid E + E \mid E * E \mid (E)$$
$$V = \{S, E\}$$
$$T = \{id, '+', '*', '(', ')'\}$$

Derivation

input

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

grammar

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

S	-----	$S \rightarrow S;S$
$S \ ; \ S$	-----	$S \rightarrow id := E$
$id := E \ ; \ S$	-----	$E \rightarrow id$
$id := id \ ; \ S$	-----	$S \rightarrow id := E$
$id := id \ ; \ id := E$	-----	$E \rightarrow E + E$
$id := id \ ; \ id := E + E$	-----	$E \rightarrow id$
$id := id \ ; \ id := E + id$	-----	$E \rightarrow id$
$id := id \ ; \ id := id + id$	-----	
$x := z \ ; \ y := x + z$		

Parse Tree

S

S ; S

id := E ; S

id := id ; S

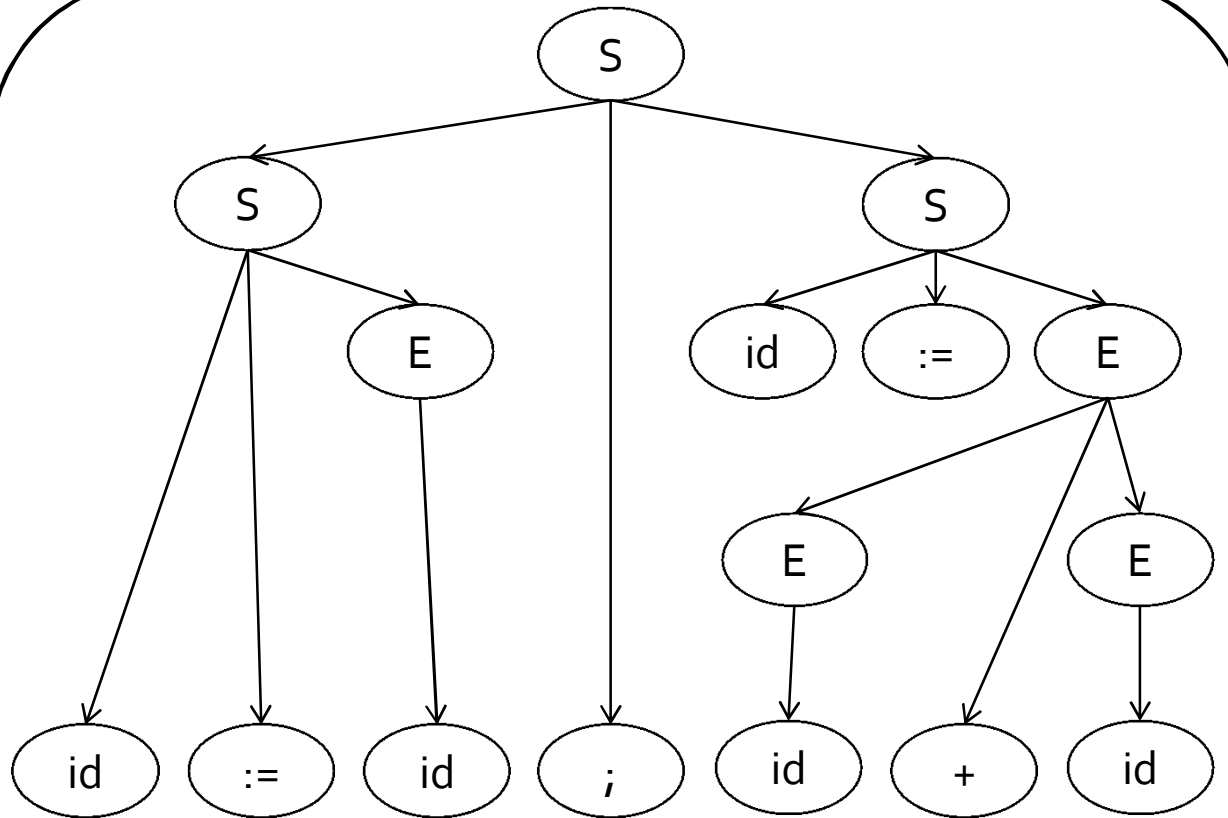
id := id ; id := E

id := id ; id := E + E

id := id ; id := E + id

id := id ; id := id + id

x := z ; y := x + z



Questions

- How did we know which rule to apply on every step?
- Does it matter?
- Would we always get the same result?

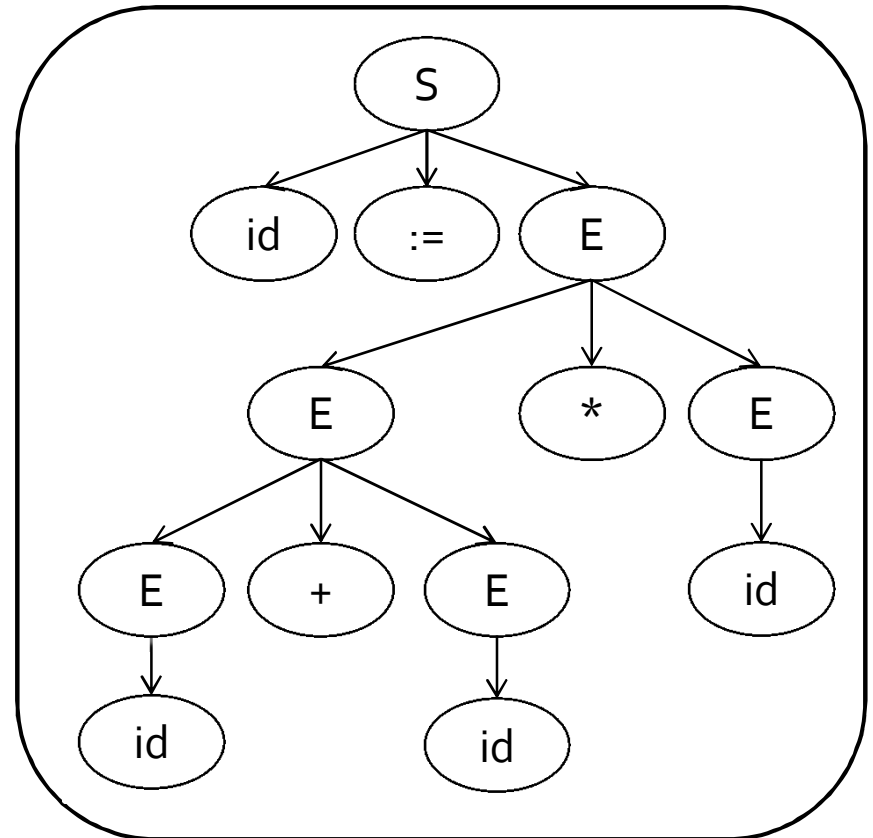
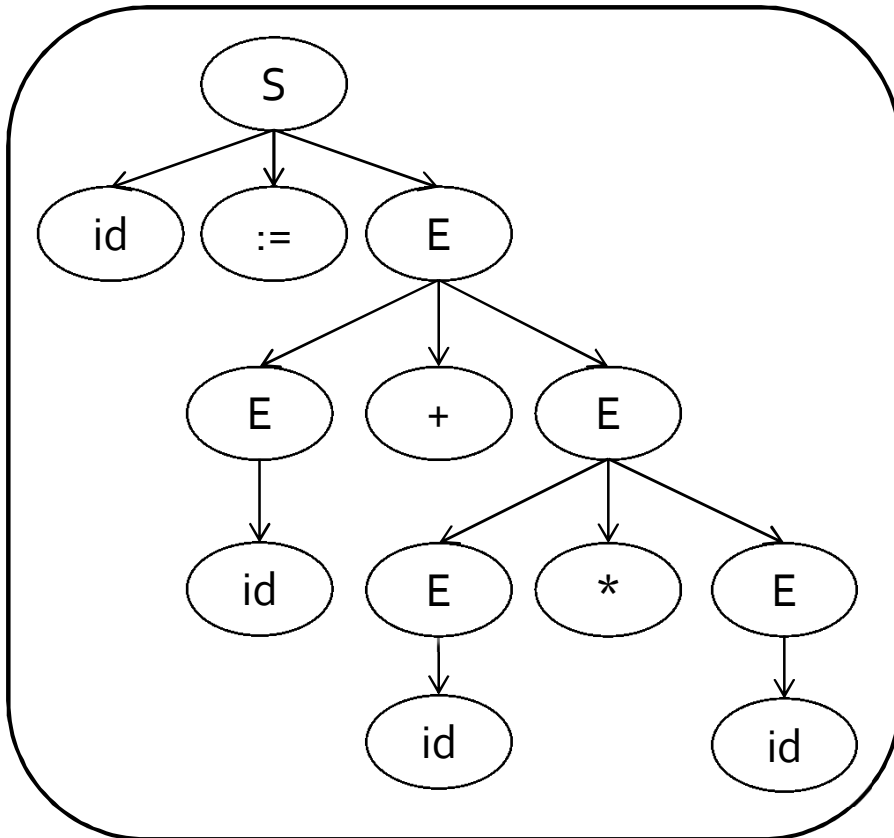
Ambiguity

$x := y + z * w$

$S \rightarrow S ; S$

$S \rightarrow id := E$

$E \rightarrow id \mid E + E \mid E * E \mid (E)$



Leftmost/rightmost Derivation

- Leftmost derivation
 - always expand leftmost non-terminal
- Rightmost derivation
 - Always expand rightmost non-terminal
- Allows us to describe derivation by listing the sequence of rules
 - always know what a rule is applied to
- Orders of derivation applied in our parsers (coming soon)

Leftmost Derivation

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

$S \rightarrow S;S$

$S \rightarrow id := E$

$E \rightarrow id \mid E + E \mid E * E \mid (E)$

S	-----	$S \rightarrow S;S$
$S \ ; \ S$	-----	$S \rightarrow id := E$
$id := E \ ; \ S$	-----	$E \rightarrow id$
$id := id \ ; \ S$	-----	$S \rightarrow id := E$
$id := id \ ; \ id := E$	-----	$E \rightarrow E + E$
$id := id \ ; \ id := E + E$	-----	$E \rightarrow id$
$id := id \ ; \ id := id + E$	-----	$E \rightarrow id$
$id := id \ ; \ id := id + id$	-----	
$x := z \ ; \ y := x + z$		

Rightmost Derivation

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

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

S		$S \rightarrow S;S$
$S ; S$		
$S ; id := E$		$S \rightarrow id := E$
$S ; id := E + E$		$E \rightarrow E + E$
$S ; id := E + id$		$E \rightarrow id$
$S ; id := id + id$		$E \rightarrow id$
$id := E ; id := id + id$		$S \rightarrow id := E$
$id := id ; id := id + id$		$E \rightarrow id$
$x := z ; y := x + z$		

Bottom-up Example

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

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

$id := id ; id := id + id$	
-----	$E \rightarrow id$
$id := E ; id := id + id$	
-----	$S \rightarrow id := E$
$S ; id := id + id$	
-----	$E \rightarrow id$
$S ; id := E + id$	
-----	$E \rightarrow id$
$S ; id := E + E$	
-----	$E \rightarrow E + E$
$S ; id := E$	
-----	$S \rightarrow id := E$
$S ; S$	
-----	$S \rightarrow S;S$
S	

Bottom-up picking left alternative on every step \rightarrow Rightmost derivation when going top-down

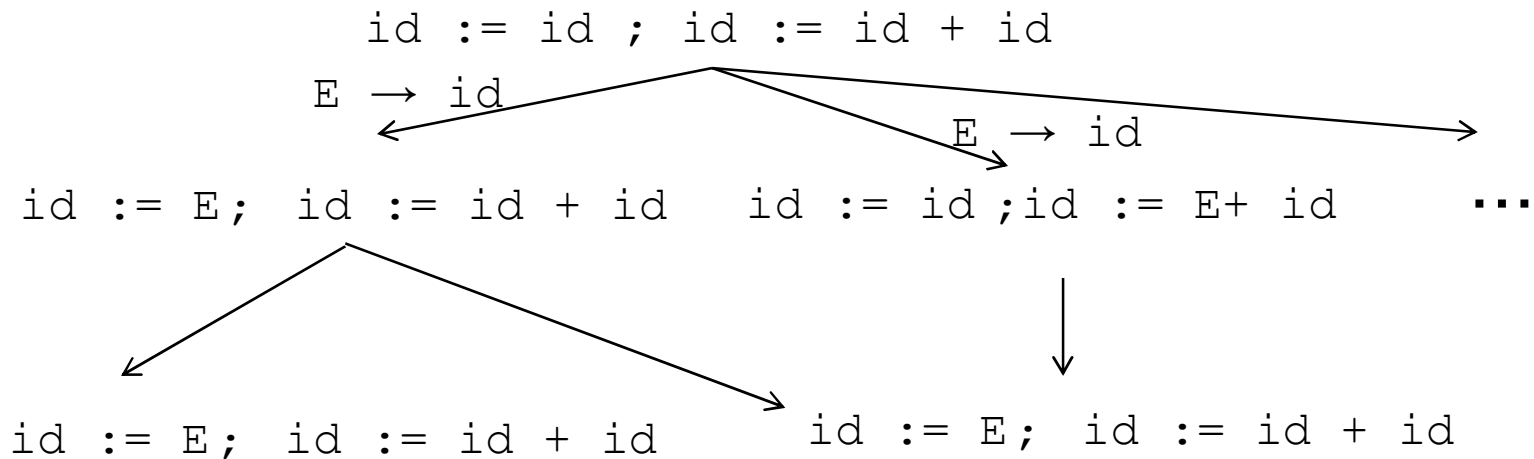
Parsing

- A context free language can be recognized by a non-deterministic pushdown automaton
- 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
- CYK parser can be used to parse any context-free language but has complexity $O(n^3)$
- We want efficient parsers
 - Linear in input size
 - Deterministic pushdown automata
 - We will sacrifice generality for efficiency

“Brute-force” Parsing

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

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



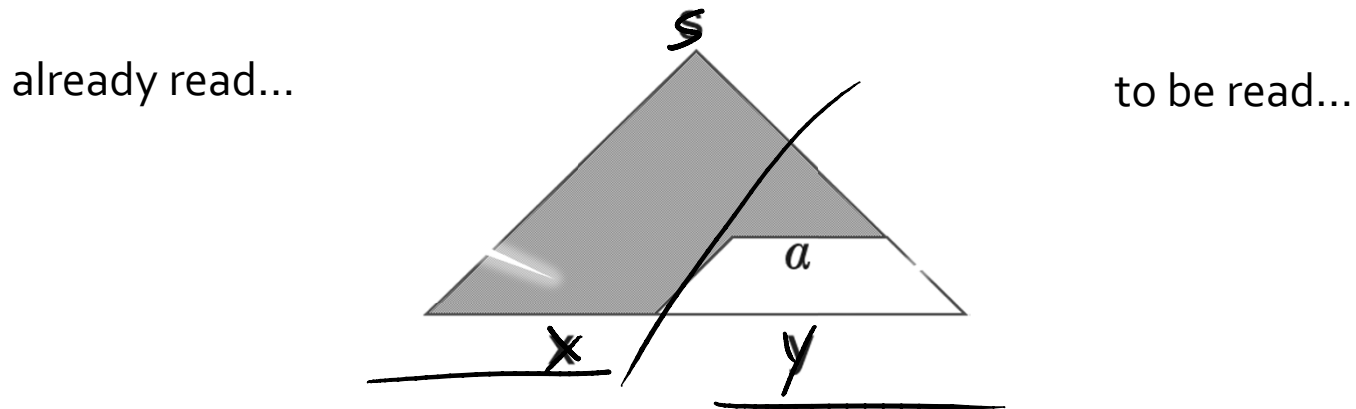
(not a parse tree... a search for the parse tree by exhaustively applying all rules)

Efficient Parsers

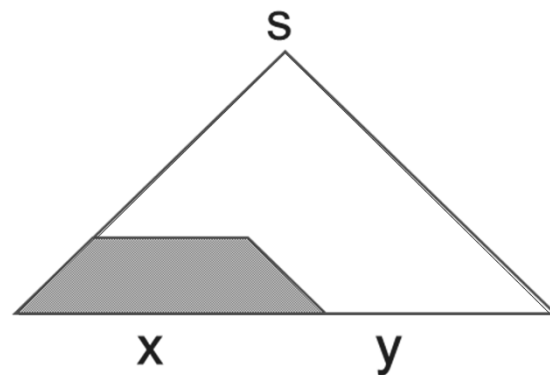
- Top-down (predictive)
 - Construct the leftmost derivation
 - Apply rules “from left to right”
 - Predict what rule to apply based on nonterminal and token
- Bottom up (shift reduce)
 - Construct the rightmost derivation
 - Apply rules “from right to left”
 - Reduce a right-hand side of a production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

Efficient Parsers

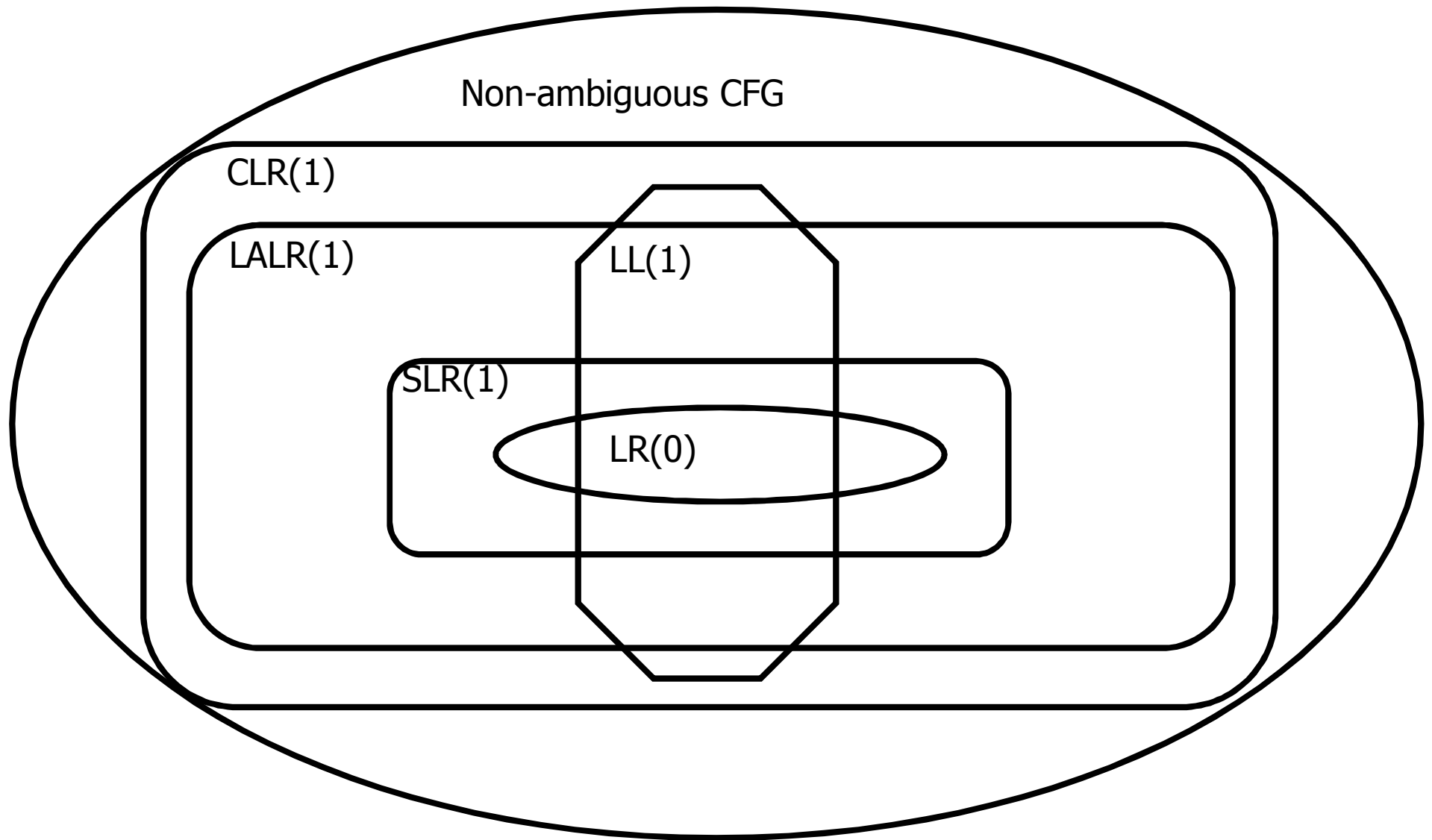
- Top-down (predictive parsing)



- Bottom-up (shift reduce)



Grammar Hierarchy



Top-down Parsing

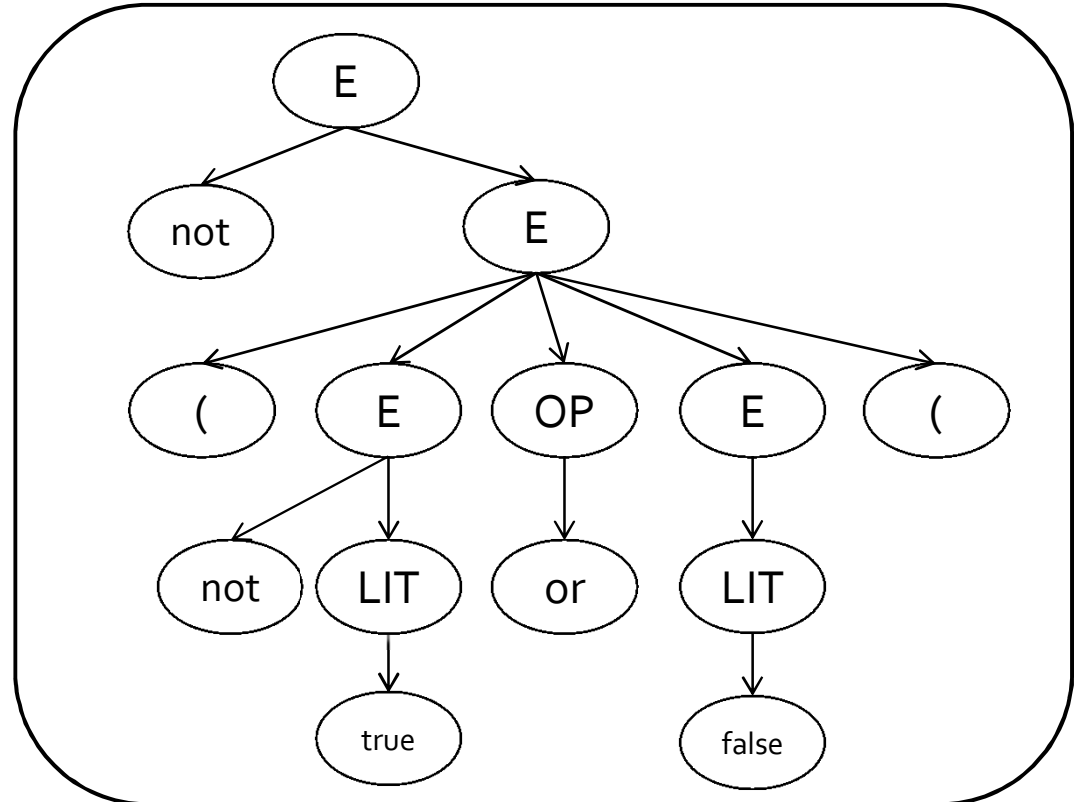
- Given a grammar $G=(V,T,P,S)$ and a word w
- Goal: derive w using G
- Idea
 - Apply production to leftmost nonterminal
 - Pick production rule based on next input token
- General grammar
 - More than one option for choosing the next production based on a token
- Restricted grammars (LL)
 - Know exactly which single rule to apply
 - May require some lookahead to decide

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 function for every nonterminal
- Every function work as follows
 - Find applicable production rule
 - Terminal function checks match with next input token
 - Nonterminal function calls (recursively) other functions
- If there are several applicable productions for a nonterminal, use lookahead

Matching tokens

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

- Variable `current` holds the current input token

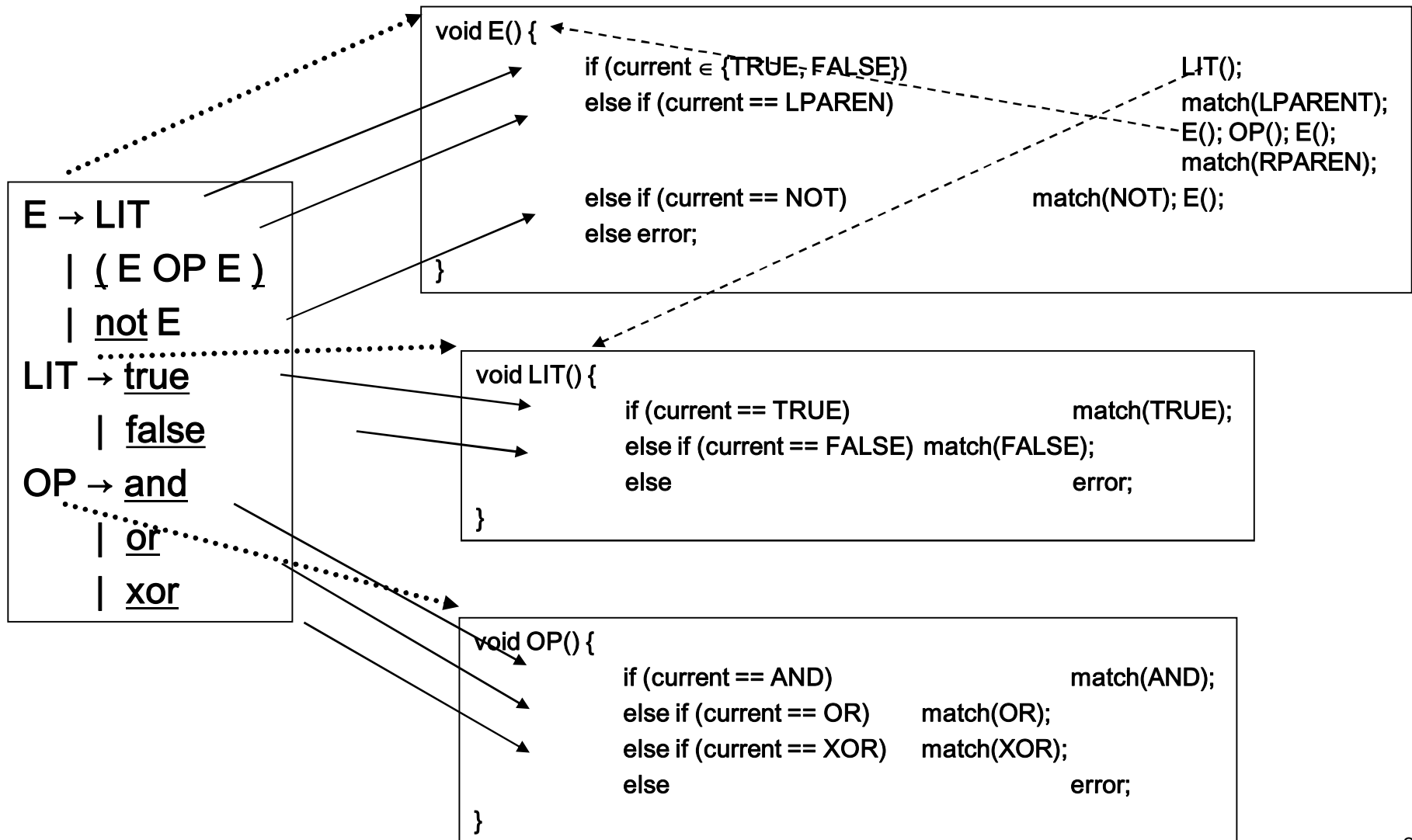
functions for nonterminals

```
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



Adding semantic actions

- Can add an action to perform on each production rule
- Can build the parse tree
 - Every function returns an object of type Node
 - Every Node maintains a list of children
 - Function calls can add new children

Building the parse tree

```
Node E() {
    result = new Node();
    result.name = "E";
    if (current ∈ {TRUE, FALSE}) // E → LIT
        result.addChild(LIT());
    else if (current == LPAREN) // E → ( E OP E )
        result.addChild(match(LPAREN));
        result.addChild(E());
        result.addChild(OP());
        result.addChild(E());
        result.addChild(match(RPAREN));
    else if (current == NOT) // E → not E
        result.addChild(match(NOT));
        result.addChild(E());
    else error;
    return result;
}
```

Recursive Descent

```
void A() {
    choose an A-production,  $A \rightarrow X_1X_2\dots X_k$ ;
    for (i=1; i ≤ k; i++) {
        if ( $X_i$  is a nonterminal)
            call procedure  $X_i()$ ;
        elseif ( $X_i == \text{current}$ )
            advance input;
        else
            report error;
    }
}
```

- How do you pick the right A-production?
- Generally – try them all and use backtracking
- In our case – use lookahead

Recursive descent: are we done?

```
term → ID | indexed_elem  
indexed_elem → ID [ expr ]
```

- The function for `indexed_elem` will never be tried...
 - What happens for input of the form
 - `ID [expr]`

Recursive descent: are we done?

```
S → A a b  
A → a | ε
```

```
int S() {  
    return A() && match(token('a')) && match(token('b'));  
}  
int A() {  
    return match(token('a')) || 1;  
}
```

- What happens for input "ab" ?
- What happens if you flip order of alternatives and try "aab"?

Recursive descent: are we done?

$E \rightarrow E - \text{term}$

```
int E() {  
    return E() && match(token('-')) && term();  
}
```

- What happens with this procedure?
- Recursive descent parsers cannot handle left-recursive grammars

Figuring out when it works...

①

```
term → ID | indexed_elem  
indexed_elem → ID [ expr ]
```

②

```
S → A a b  
A → a | ε
```

③

```
E → E - term
```

3 examples where we got into trouble with our recursive descent approach

FIRST sets

- For every production rule $A \rightarrow \alpha$
 - $\text{FIRST}(\alpha)$ = all terminals that α can start with
 - i.e., every token that can appear as first in α under some derivation for α
- In our Boolean expressions example
 - $\text{FIRST}(\text{LIT}) = \{ \text{true}, \text{false} \}$
 - $\text{FIRST}((E \text{ OP } E)) = \{ '(' \}$
 - $\text{FIRST}(\text{not } E) = \{ \text{not} \}$
- No intersection between FIRST sets => can always pick a single rule
- If the FIRST sets intersect, may need longer lookahead
 - $\text{LL}(k)$ = class of grammars in which production rule can be determined using a lookahead of k tokens
 - $\text{LL}(1)$ is an important and useful class

FOLLOW Sets

- What do we do with nullable alternatives?
 - Use what comes afterwards to predict the right production
- For every production rule $A \rightarrow \alpha$
 - $\text{FOLLOW}(A)$ = set of tokens that can immediately follow A
- Can predict the alternative A_k for a non-terminal N when the lookahead token is in the set
 - $\text{FIRST}(A_k) \cup (\text{if } A_k \text{ is nullable then } \text{FOLLOW}(N))$

LL(k) Grammars

- A grammar is in the class LL(K) when it can be derived via:
 - Top down derivation
 - 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) when it has an LL(k) grammar

Back to our 1st example

```
term → ID | indexed_elem  
indexed_elem → ID [ expr ]
```

- $\text{FIRST}(\text{ID}) = \{ \text{ID} \}$
- $\text{FIRST}(\text{indexed_elem}) = \{ \text{ID} \}$
- FIRST/FIRST conflict

Left factoring

- Rewrite the grammar to be in LL(1)

```
term → ID | indexed_elem  
indexed_elem → ID [ expr ]
```



```
term → ID after_ID  
after_ID → [ expr ] | ε
```

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

Left factoring - another example

```
S → if E then S else S  
   | if E then S  
   | T
```



```
S → if E then S S'  
   | T  
S' → else S | ε
```

Back to our 2nd example

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

- $\text{FIRST}(S) = \{ 'a' \}$, $\text{FOLLOW}(S) = \{ \}$
- $\text{FIRST}(A) = \{ 'a' \varepsilon \}$, $\text{FOLLOW}(A) = \{ 'a' \}$
- FIRST/FOLLOW conflict

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

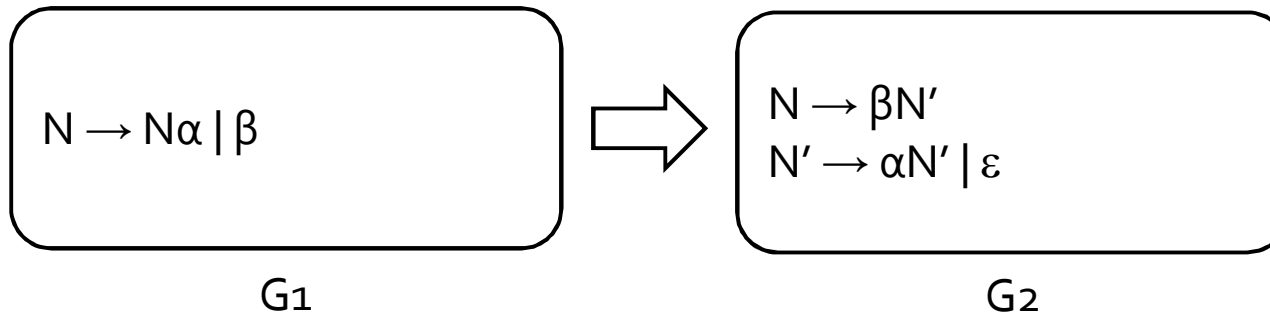
$S \rightarrow a \text{ after_A}$
 $\text{after_A} \rightarrow a b \mid b$

Back to our 3rd example

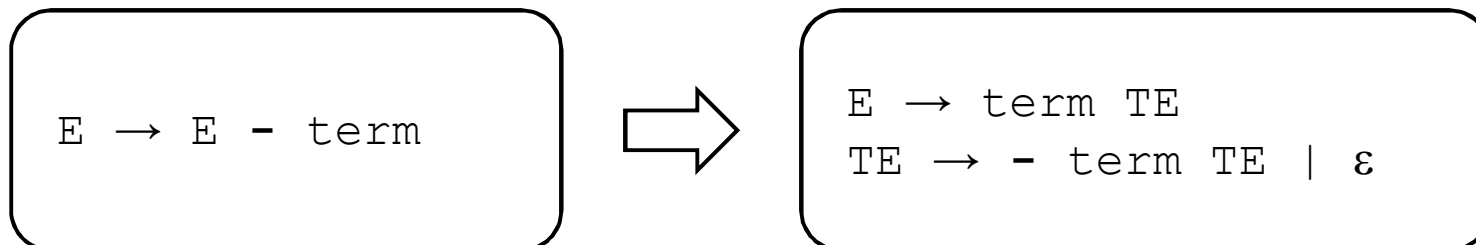
$E \rightarrow E - \text{term}$

- Left recursion cannot be handled with a bounded lookahead
- What can we do?

Left recursion removal



- $L(G_1) = \beta, \beta\alpha, \beta\alpha\alpha, \beta\alpha\alpha\alpha, \dots$
- $L(G_2) = \text{same}$
- For our 3rd example:



LL(k) Parsers

- Recursive Descent
 - Manual construction
 - Uses recursion
- Wanted
 - A parser that can be generated automatically
 - Does not use recursion

LL(k) parsing with pushdown automata

- Pushdown automaton uses
 - Prediction stack
 - Input stream
 - Transition table
 - nonterminals x tokens -> production alternative
 - Entry indexed by nonterminal N and token t contains the alternative of N that must be predicated when current input starts with t

LL(k) parsing with pushdown automata

- Two possible moves
 - Prediction
 - When top of stack is nonterminal N , pop N , lookup $\text{table}[N,t]$. If $\text{table}[N,t]$ is not empty, push $\text{table}[N,t]$ on prediction stack, otherwise – syntax error
 - Match
 - When top of prediction stack is a terminal T , must be equal to next input token t . If $(t == T)$, pop T and consume t . If $(t \neq T)$ syntax error

- Parsing terminates when prediction stack is empty. If input is empty at that point, success. Otherwise, syntax error

Example transition table

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

Input tokens

	()	not	true	false	and	or	xor	\$
E	2		3	1	1				
LIT				4	5				
OP						6	7	8	

Nonterminals

Which rule should be used

Simple Example

aacbb\$

$A \rightarrow aAb \mid c$

Input suffix	Stack content	Move
aacbb\$	A\$	predict(A,a) = $A \rightarrow aAb$
aacbb\$	aAb\$	match(a,a)
acbb\$	Ab\$	predict(A,a) = $A \rightarrow aAb$
acbb\$	aAbb\$	match(a,a)
cbb\$	Abb\$	predict(A,c) = $A \rightarrow c$
cbb\$	cbb\$	match(c,c)
bb\$	bb\$	match(b,b)
b\$	b\$	match(b,b)
\$	\$	match(\$,\$) – success

	a	b	c
A	$A \rightarrow aAb$		$A \rightarrow c$

Simple Example

abcbb\$

$A \rightarrow aAb \mid c$

Input suffix	Stack content	Move
abcbb\$	A\$	predict(A,a) = $A \rightarrow aAb$
abcbb\$	aAb\$	match(a,a)
bcbbs\$	Ab\$	predict(A,b) = ERROR

	a	b	c
A	$A \rightarrow aAb$		$A \rightarrow c$

Error Handling

- Mentioned last time
 - Lexical errors
 - Syntax errors
 - Semantic errors (e.g., type mismatch)

Error Handling and Recovery

$$x = a * (p+q * (-b * (r-s));$$

- Where should we report the error?
- The valid prefix property
- Recovery is tricky
 - Heuristics for dropping tokens, skipping to semicolon, etc.

Error Handling in LL Parsers

c\$

$S \rightarrow a c \mid b S$

Input suffix	Stack content	Move
c\$	S\$	predict(S,c) = ERROR

- Now what?
 - Predict bS anyway “missing token b inserted in line XXX”

	a	b	c
S	$S \rightarrow a c$	$S \rightarrow b S$	

Error Handling in LL Parsers

c\$

$S \rightarrow a c \mid b S$

Input suffix	Stack content	Move
bc\$	S\$	predict(b,c) = $S \rightarrow bS$
bc\$	bS\$	match(b,b)
c\$	S\$	Looks familiar?

- Result: infinite loop

	a	b	c
S	$S \rightarrow a c$	$S \rightarrow b S$	

Error Handling

- Requires more systematic treatment
- Enrichment
 - Acceptable-set method
 - Not part of course material

Summary

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

Coming up next time

- More syntax analysis