

# Top-Down Parsing

---

## Parser Classes

- **Top-down** LL (read Left to right, follow Leftmost derivation)
- **Bottom-up** LR (read Left to right, follow Rightmost derivation)
- Named after how the parse tree is built: top to bottom or bottom to top
- Also differ in power

## Top-Down Parser

### Power and Implementation

- Less powerful
  - There are grammars that can be parsed with bottom-up parsers but not with top-down parsers with the same lookaheads
- More complicated for back-end operations unless some limitations are in place
- Easier to implement, especially *recursive descent*

### Operations and Properties

- Builds the parse tree from the root down
- Follows leftmost derivation, therefore, it is called LL(k)
- Uses stack memory
  - Systems stack in recursive descent
  - Expands the topmost symbol on the stack (leftmost in derivation)
- Tokens from program being parsed are **never** on the stack
  - the stack contains tokens from the grammar that are **expected** in the program

### Push-Down Algorithm

- Utilize stack memory
- Start with pushing initial the nonterminal S
- At any moment
  - if a **terminal** is on the top of the stack
    - if it matches the incoming token
      - the terminal is popped  
and the token is consumed (next token is requested)
    - error otherwise
  - if a <nonterminal> is on the top of the stack
    - the nonterminal is replaced with one of its RHS productions
    - if more than one production, must *predict* the correct one
    - errors possible in a predictive parser if none predicted

# Top-Down Parsing

---

- Successful termination
  - Stack is empty and **EOF** is the remaining token

## Example

Use the unambiguous expression grammar to top-down parse the following program: `id+id*id`

## Problems in Top-Down Parsing

1. **Left recursive productions** (direct or indirect)
  - infinite stack growth
  - can always be handled by a removal procedure that does not alter the grammar just writes it differently
2. **Non-deterministic productions** (more than one production for a <nonterminal>)
  - determinism is verified by *First* and *Follow* sets
  - sets must be pairwise disjoint for the same <nonterminal>
  - may be handled (make deterministic), but not guaranteed
    - *left-factorization* is the general technique for that

## Preparing Grammar for LL(1)

- Not every grammar can be converted to LL(k=1). The objective is to avoid the above problems. The steps are:
  1. Remove left recursion (guaranteed)
  2. Verify k=1 using First and Follow sets. If k≤1, done. Otherwise
    - Left factorize and verify again
    - This step is not guaranteed to succeed

## Left Recursion

- Top-down parsers cannot handle left-recursion, direct or indirect
- Any direct left-recursion can be removed with equivalent grammar
- Indirect left-recursions can be replaced by direct, and subsequently removed
- This process is guaranteed to succeed but care must be taken if removing multiple recursions
- Removing direct left recursion
  - separate all left recursive from the other productions for each nonterminal

$$A \rightarrow A\alpha \mid A\beta \mid \gamma \mid \delta$$
  - introduce a new nonterminal X
  - change nonrecursive productions to
$$A \rightarrow \gamma X \mid \delta X$$
  - replace recursive productions by
$$X \rightarrow \epsilon \mid \alpha X \mid \beta X$$

## Example

Remove left recursion from  $\{E \rightarrow E+T \mid T\}$ . Result is  $\{E \rightarrow TA', A' \rightarrow +TA' \mid \epsilon\}$

# Top-Down Parsing

---

## Nondeterminism

- Grammar requires the number of lookaheads equal to the lookahead needed for the worst nonterminal (max k)
  - Thus computing k (the number of lookaheads) needed for a grammar comes down to computing k for each nonterminal and then taking the max value
- Nonterminal is deterministic when there is only one choice of a production
  - No lookahead needed to make decision thus  $k=0$
- Nonterminal is non-deterministic when there is more than one production thus  $k>0$
- The objective is to ensure that 1 lookahead ( $k=1$ ) will be enough to make the right decision for each nonterminal

There are two parts here

- *First* and *Follow* sets verify whether  $k=1$  or  $k>1$
- *Left factorization* potentially reduces k
- There are LL( $k>1$ ) grammar that cannot be expressed as LL( $k-1$ ) and thus converting grammar to LL(1) is not always possible

## Example

**if then [else]** is an example of not-LL(1) construct and cannot be reduced to LL(1), but it may be solved in LL(1)-parser using other techniques such as by ordering productions (not covered here).

## Grammar Verification for LL(1)

1. Remove all left recursions first
2. Nonterminals without choice have  $k=0$   
 $N \rightarrow \alpha$
3. Nonterminals with multiple productions must be verified whether  $k=1$  or  $k>1$   
 $N \rightarrow \alpha \mid \beta$   
 $N \rightarrow \alpha \mid \epsilon$   
 $N \rightarrow \alpha \mid \beta \mid \epsilon$
4. The entire grammar is LL(max k over all nonterminals).  
Thus, each nonterminal with multiple productions must be LL(1) for the entire grammar to be LL(1)

## Nondeterministic Nonterminal Verification for LL(1)

1. Compute *First* sets for each RHS choice
2. Compute *Follow* set of LHS if there is the empty production
3. All the above sets must be pairwise disjoint (excluding the empty symbol) for the nonterminal to be LL(1)
  - otherwise the nonterminal is LL( $k>1$ ) and thus the entire grammar is LL( $k>1$ )

# Top-Down Parsing

---

## Example

Verify LL(1) or LL(k>1)

$S \rightarrow \alpha$   $k=0$

$N \rightarrow \beta$   $k=0$

$M \rightarrow \alpha \mid \beta$   $k=1$  (suppose)

Then the entire grammar is LL(1).

## Left Factorization

- Combines alternative productions starting with the same prefixes
  - this delays decisions about predictions until new tokens are seen
  - this is a form of extending the lookahead by utilizing the stack
  - bottom-up parsers extend this idea even further
  - from implementation perspective, it may not be necessary to perform left-factorization but rather to delay decision in the code

## Example

Assume grammar

$S \rightarrow \mathbf{ee} \mid \mathbf{bAc} \mid \mathbf{bAe}$

$A \rightarrow \mathbf{d} \mid \mathbf{cA}$

When the parsed program is **bcde**, it is impossible to decide production on S while looking only at the incoming 1 token (**b** here)

Rewrite by combining common prefix **bA**

$S \rightarrow \mathbf{ee} \mid \mathbf{bA(c \mid e)} \rightarrow \mathbf{ee} \mid \mathbf{bAX}$

$X \rightarrow \mathbf{c} \mid \mathbf{e}$

$A \rightarrow \mathbf{d} \mid \mathbf{cA}$

## First Set

The name *First* comes from the fact that the set contains token which may show up first on the stack when the production is used and before a token is consumed.

We compute *First* of right hand sides of non-deterministic productions (more than one choice) and not individual nonterminals unless needed in the algorithm.

$N \rightarrow \alpha \mid \beta$

Must compute  $\text{First}(\alpha)$  and  $\text{First}(\beta)$

**FIRST( $\alpha$ ) algorithm:**

$\alpha$  is a single element (terminal/token or nonterminal) or  $\epsilon$

if  $\alpha = \text{terminal } y$  then  $\text{FIRST}(\alpha) = \{y\}$

if  $\alpha = \epsilon$  then  $\text{FIRST}(\alpha) = \{\epsilon\}$

---

## Top-Down Parsing

---

if  $\alpha$  is nonterminal and we have productions  $\alpha \rightarrow \beta_1 \mid \beta_2 \mid \dots$  then  
 $\text{FIRST}(\alpha) = \cup \text{FIRST}(\beta_i)$

$$\alpha = X_1 X_2 \dots X_n$$

set  $\text{FIRST}(\alpha) = \{\}$

for  $j=1..n$  include  $\text{FIRST}(X_j)$  in  $\text{FIRST}(\alpha)$ , but STOP when  $X_j$  is not *nullable*

$X$  is nullable if  $X \rightarrow \epsilon$  directly or indirectly

$X = \epsilon$  or  $\epsilon$  is a member of  $\text{FIRST}(X)$

terminal is of course never nullable

if  $X_n$  was reached and is also nullable then include  $\epsilon$  in  $\text{FIRST}(\alpha)$

### Example

$S \rightarrow Ab \mid Bc$        $\text{First}(Ab) = \{\mathbf{h}, \mathbf{i}, \mathbf{c}, \mathbf{d}\}$     $\text{First}(Bc) = \{e, \mathbf{g}\}$  thus  $k=1$   
 $A \rightarrow Df \mid CA$   
 $B \rightarrow \mathbf{g}A \mid \mathbf{e}$   
 $C \rightarrow \mathbf{d}C \mid \mathbf{c}$   
 $D \rightarrow \mathbf{h} \mid \mathbf{i}$

### Follow Set

The name Follow comes from the fact that the set contains tokens that physically follow the given nonterminal on the stack (are below).

$A \rightarrow \alpha \mid \epsilon$

Note that Follow set never contains  $\epsilon$ .

**FOLLOW(A) algorithm:**

If  $A=S$  then put end marker (e.g., EOFtk) into  $\text{FOLLOW}(A)$  and continue

Find all productions with  $A$  on rhs:  $Q \rightarrow \alpha A \beta$

- if  $\beta$  begins with a terminal  $q$  then  $q$  is in  $\text{FOLLOW}(A)$
- if  $\beta$  begins with a nonterminal then  $\text{FOLLOW}(A)$  includes  $\text{FIRST}(\beta) - \{\epsilon\}$
- if  $\beta=\epsilon$  or when  $\beta$  is nullable then include  $\text{FOLLOW}(Q)$  in  $\text{FOLLOW}(A)$

### Grammar is LL( $k \leq 1$ ) When

- No left recursion
- For all nonterminals with multiple productions, the *FIRST* sets of the RHSs are disjoint (excluding  $\epsilon$ )

# Top-Down Parsing

---

- In addition, if any RHS is  $\epsilon$  or nullable, the *FOLLOW* of that nonterminal must be disjoint from all *FIRST* sets for the nonterminal

## Example

Rewrite the last expression grammar until proven LL(1).

Removing left recursion gives

$E \rightarrow TQ$		$k=0$
$Q \rightarrow +TQ \mid -TQ \mid \epsilon$	$FIRST(+TQ)=\{+\}$ $FIRST(-TQ)=\{-\}$ $FOLLOW(Q)=\{EOFtk\}$	$k=1$
$T \rightarrow FR$		$k=0$
$R \rightarrow *FR \mid /FR \mid \epsilon$	$FIRST(*FR)=\{*\}$ $FIRST(/FR)=\{/ \}$ $FOLLOW(R)=\{+ - \} EOFtk\}$	$k=1$
$F \rightarrow (E) \mid id$	$FIRST(id)=\{id\}$ $FIRST((E))=\{(\}$	$k=1$

Entire grammar is LL(1)

## Predictions in LL(1)

- Ensure grammar is LL( $k \leq 1$ )
- In the top-down algorithm
  - When terminal is on top of the stack, just match as before
    - remove from the stack and get next token from the scanner
  - When nonterminal is on top of the stack
    - if there is single production for this nonterminal, just apply it
    - otherwise look at the next token from the scanner and
      - if it is in one of the *FIRST* sets for the RHSs, apply that RHS
      - otherwise if there is empty production, apply the empty production
      - otherwise there is an error