

Front End: Syntax Analysis

Bottom-Up Parsing

Parsers

- Top-down
Construct **leftmost** derivations starting from the start symbol.
- Bottom-up
Construct (reverse) **rightmost** derivations starting from the input string by **reducing** it to the start symbol.

Both parsers are guided by the input string in the search of a derivation.

Bottom-up Parsing

Intuition: construct the parse tree from the leaves to the root.

Grammar:

$S \rightarrow AB$

$A \rightarrow x \mid Y$

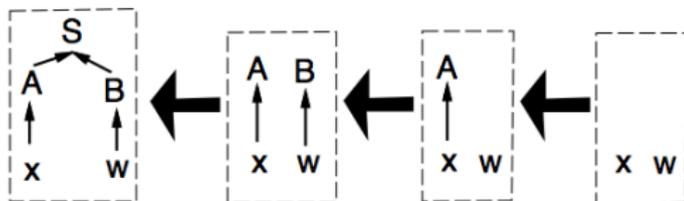
$B \rightarrow w \mid Z$

$Y \rightarrow xb$

$Z \rightarrow wp$

Example:

Input xw .



Bottom-up parsing

Constructing a parse tree for an input string starting from the leaves towards the root.

Example II

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

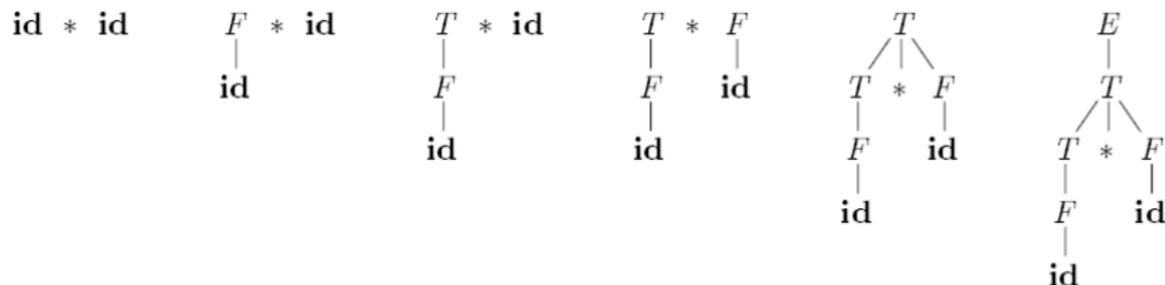


Figure 4.25: A bottom-up parse for `id * id`

Shift-Reduce Parsing

A **shift-reduce parser** is a form of bottom-up parser whose primary operations are

- **Shift**: shift the next input symbol
- **Reduce**: identify the *handle* and replace it with the head of the appropriate production.

A *reduction step* is the reverse of a derivation step (= a non-terminal is replaced by the body of one of its productions). Thus, reducing corresponds to constructing a derivation in reverse.

Example: The parse in Fig. 4.25 corresponds to the rightmost derivation

$$E \Rightarrow T \Rightarrow T * F \Rightarrow T * \mathbf{id} \Rightarrow F * \mathbf{id} \Rightarrow \mathbf{id} * \mathbf{id}$$

Handle

A **handle** is a substring that matches the body of a production.

RIGHT SENTENTIAL FORM	HANDLE	REDUCING PRODUCTION
$\mathbf{id_1 * id_2}$	$\mathbf{id_1}$	$F \rightarrow \mathbf{id}$
$F * \mathbf{id_2}$	F	$T \rightarrow F$
$T * \mathbf{id_2}$	$\mathbf{id_2}$	$F \rightarrow \mathbf{id}$
$T * F$	$T * F$	$T \rightarrow T * F$
T	T	$E \rightarrow T$

Figure 4.26: Handles during a parse of $\mathbf{id_1 * id_2}$

Handle: Formal Definition

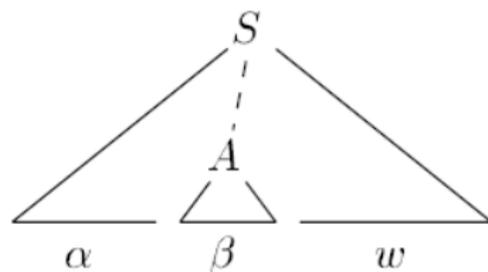


Figure 4.27: A handle $A \rightarrow \beta$ in the parse tree for $\alpha\beta w$

Definition

A **handle** for $\gamma = \alpha\beta w$ s.t. $S \Rightarrow_{rm}^* \alpha A w \Rightarrow_{rm} \gamma$ is

- 1 a **production rule** $A \rightarrow \beta$, and
- 2 a **position** p in γ where β can be located.

Stack Implementation

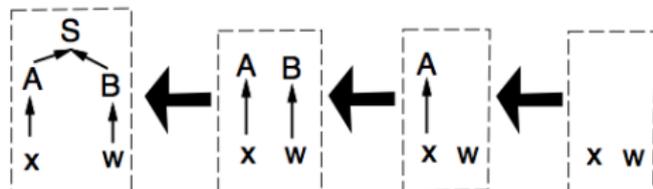
A **stack** holds grammar symbols and an **input buffer** holds the rest of the string to be parsed.

STACK	INPUT	ACTION
\$	id₁ * id₂ \$	shift
\$ id₁	* id₂ \$	reduce by $F \rightarrow \mathbf{id}$
\$ F	* id₂ \$	reduce by $T \rightarrow F$
\$ T	* id₂ \$	shift
\$ T *	id₂ \$	shift
\$ T * id₂	\$	reduce by $F \rightarrow \mathbf{id}$
\$ T * F	\$	reduce by $T \rightarrow T * F$
\$ T	\$	reduce by $E \rightarrow T$
\$ E	\$	accept

Figure 4.28: Configurations of a shift-reduce parser on input **id₁*id₂**

Example

STACK	INPUT	ACTION
\$	$xw\$$	shift
$\$x$	$w\$$	reduce by $A \rightarrow x$
$\$A$	$w\$$	shift
$\$Aw$	$\$$	reduce by $B \rightarrow w$
$\$AB$	$\$$	reduce by $S \rightarrow AB$
$\$S$	$\$$	accept



$$S \xRightarrow{rm} AB \xRightarrow{rm} Aw \xRightarrow{rm} xw.$$

Conflict I: Shift-Reduce

There are CF grammars for which the shift-reduce parsing does not work.

$stmt \rightarrow$	if <i>expr</i> then <i>stmt</i>
	if <i>expr</i> then <i>stmt</i> else <i>stmt</i>
	other

With the following configuration we cannot decide whether to shift or to reduce.

STACK	INPUT
\$... if <i>expr</i> then <i>stmt</i>	else ...

Conflict II: Reduce-Reduce

Consider a language where procedures and arrays share the same syntax.

- (1) $stmt \rightarrow id (parameter_list)$
- (2) $stmt \rightarrow expr := expr$
- (3) $parameter_list \rightarrow parameter_list , parameter$
- (4) $parameter_list \rightarrow parameter$
- (5) $parameter \rightarrow id$
- (6) $expr \rightarrow id (expr_list)$
- (7) $expr \rightarrow id$
- (8) $expr_list \rightarrow expr_list , expr$
- (9) $expr_list \rightarrow expr$

Figure 4.30: Productions involving procedure calls and array references

Which production should we choose with configuration

STACK	INPUT
\$... id (id	id, id)...

LR Parsing

LR(k) parsing, introduced by D. Knuth in 1965, is today the most prevalent type of bottom-up parsing.

L is for *Left-to-right* scanning of the input,

R is for reverse *Rightmost* derivation,

k is the number of lookahead tokens.

Different types:

- Simple LR or **SLR**, the easiest method for constructing shift-reduce parsers,
- **Canonical LR**,
- **LALR**.

The last two types are used in the majority of LR parsers.

The LR-Parsing Model

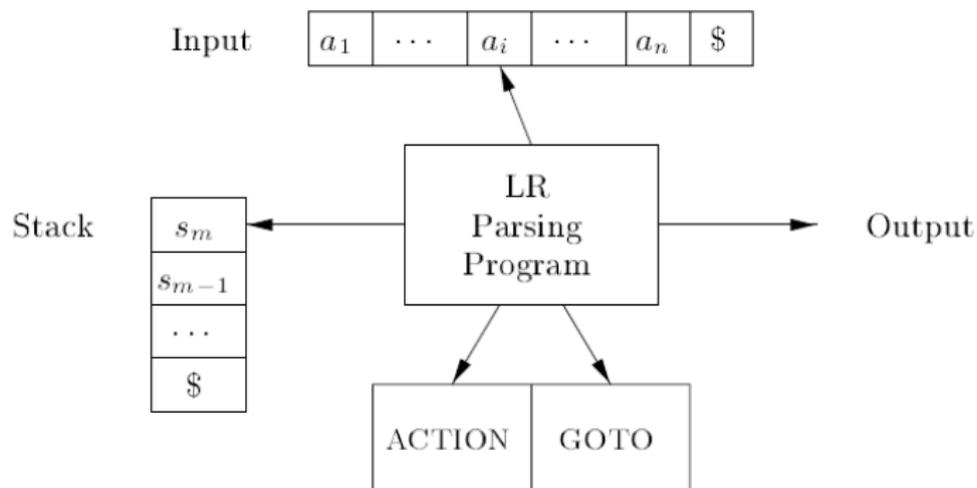


Figure 4.35: Model of an LR parser

The **parsing table**, consisting of the ACTION and GOTO functions, is the only variable part. The stack content is a sequence of **states**, corresponding each to a grammar symbol.

The LR-Parsing Algorithm

All LR-parsers behave as summarised below: the only difference is the info held by the parsing table.

```
let  $a$  be the first symbol of  $w\$$ ;  
while(1) { /* repeat forever */  
    let  $s$  be the state on top of the stack;  
    if ( ACTION[ $s, a$ ] = shift  $t$  ) {  
        push  $t$  onto the stack;  
        let  $a$  be the next input symbol;  
    } else if ( ACTION[ $s, a$ ] = reduce  $A \rightarrow \beta$  ) {  
        pop  $|\beta|$  symbols off the stack;  
        let state  $t$  now be on top of the stack;  
        push GOTO[ $t, A$ ] onto the stack;  
        output the production  $A \rightarrow \beta$ ;  
    } else if ( ACTION[ $s, a$ ] = accept ) break; /* parsing is done */  
    else call error-recovery routine;  
}
```

Figure 4.36: LR-parsing program

Conflict Resolution

STACK	INPUT	ACTION
\$	$\mathbf{id}_1 * \mathbf{id}_2$ \$	shift
\$ \mathbf{id}_1	* \mathbf{id}_2 \$	reduce by $F \rightarrow \mathbf{id}$
\$ F	* \mathbf{id}_2 \$	reduce by $T \rightarrow F$
\$ T	* \mathbf{id}_2 \$	shift
\$ $T *$	\mathbf{id}_2 \$	shift
\$ $T * \mathbf{id}_2$	\$	reduce by $F \rightarrow \mathbf{id}$
\$ $T * F$	\$	reduce by $T \rightarrow T * F$
\$ T	\$	reduce by $E \rightarrow T$
\$ E	\$	accept

Figure 4.28: Configurations of a shift-reduce parser on input $\mathbf{id}_1 * \mathbf{id}_2$

How does a shift-reduce parser know when to shift and when to reduce?

Example:

In Fig. 4.28, how does the parser know that T is not yet a handle and that the appropriate action is a *shift*?

Constructing LR-Parsing Table

LR parsers are table-driven, similarly to the non-recursive LL parsers.

In order to recognise the right-hand side of a production, an LR parser must be able to recognise handles of right sentential forms when they appear on top of the stack.

Idea: maintaining states to keep track of where we are in a parse can help an LR parser to decide when to shift and when to reduce. Construct a **Finite Automaton**.

The **SLR** method constructs a parsing table on the base of **LR(0) items** and **LR(0) automata**.

Items

Definition

An **LR(0) item** (or simply **item**) of a grammar G is a production of G with a dot at some position of the body.

Example:

For the production $A \rightarrow XYZ$ we get the items

$$A \rightarrow \bullet XYZ$$
$$A \rightarrow X \bullet YZ$$
$$A \rightarrow XY \bullet Z$$
$$A \rightarrow XYZ \bullet$$

A **state** in our FA is a set of items.

Closure

Given a set of items I , the **closure** of I is computed as follows:

```
SetOfItems CLOSURE( $I$ ) {  
     $J = I$ ;  
    repeat  
        for ( each item  $A \rightarrow \alpha \cdot B \beta$  in  $J$  )  
            for ( each production  $B \rightarrow \gamma$  of  $G$  )  
                if (  $B \rightarrow \cdot \gamma$  is not in  $J$  )  
                    add  $B \rightarrow \cdot \gamma$  to  $J$ ;  
    until no more items are added to  $J$  on one round;  
    return  $J$ ;  
}
```

Figure 4.32: Computation of CLOSURE

Intuition: if item $A \rightarrow \alpha \bullet B \beta$ is in $\text{CLOSURE}(I)$, then at some point the parser might see a substring derivable from $B \beta$ as input.

Example

$E' \rightarrow E$

$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

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

If $I = \{ E' \rightarrow \bullet E \}$,
then $\text{CLOSURE}(I)$ is



Example

$E' \rightarrow E$
 $E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid \mathbf{id}$

If $I = \{ E' \rightarrow \cdot E \}$,
then $\text{CLOSURE}(I)$ is



$E' \rightarrow \cdot E$
 $E \rightarrow \cdot E + T$
 $E \rightarrow \cdot T$
 $T \rightarrow \cdot T * F$
 $T \rightarrow \cdot F$
 $F \rightarrow \cdot (E)$
 $F \rightarrow \cdot \mathbf{id}$

The GOTO Function

If $A \rightarrow \alpha \bullet X\beta$ is in I , $\text{GOTO}(I, X)$ contains $\text{CLOSURE}(A \rightarrow \alpha \bullet X\beta)$.

```
void items( $G'$ ) {  
     $C = \{\text{CLOSURE}(\{[S' \rightarrow \cdot S]\})\}$ ;  
    repeat  
        for ( each set of items  $I$  in  $C$  )  
            for ( each grammar symbol  $X$  )  
                if (  $\text{GOTO}(I, X)$  is not empty and not in  $C$  )  
                    add  $\text{GOTO}(I, X)$  to  $C$ ;  
    until no new sets of items are added to  $C$  on a round;  
}
```

Figure 4.33: Computation of the canonical collection of sets of LR(0) items

G =

$E' \rightarrow E$
 $E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid \text{id}$

I =

$E' \rightarrow E \cdot$
 $E \rightarrow E \cdot + T$

GOTO(I,+) =



G =

$E' \rightarrow E$
 $E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid id$

I =

$E' \rightarrow E \cdot$
 $E \rightarrow E \cdot + T$

If $[A \rightarrow \alpha \cdot X \beta] \in I$,
GOTO(I, X) contains
CLOSURE($A \rightarrow \alpha X \cdot \beta$)

GOTO(I, +) =

$E \rightarrow E + \cdot T$
 $T \rightarrow \cdot T * F$
 $T \rightarrow \cdot F$
 $F \rightarrow \cdot (E)$
 $F \rightarrow \cdot id$

The LR(0) Automaton

G' : augmented grammar

LR(0) automaton for G'



$\langle \mathbf{Q}, \mathbf{q}_0, \mathbf{GOTO}: \mathbf{Q} \times (\mathbf{T}_{G'} \cup \mathbf{N}_{G'}) \rightarrow \mathbf{Q}, \mathbf{F} \rangle$

where:

$\mathbf{Q} = \mathbf{F} = \text{items}(G')$,

$\mathbf{q}_0 = \text{CLOSURE}(\{S' \rightarrow \bullet S\})$

Example

Construction of the LR(0) automaton for the augmented grammar:

$E' \rightarrow E$
 $E \rightarrow E + T \mid T$
 $T \rightarrow T^* F \mid F$
 $F \rightarrow (E) \mid \mathbf{id}$

$E' \rightarrow E$
 $E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid id$

