

## Parsing

A.k.a. **Syntax Analysis**

- Recognize *sentences* in a language.
- Discover the structure of a document/program.
- Construct (implicitly or explicitly) a tree (called as a **parse tree**) to represent the structure.
- The above tree is used later to guide translation.

## Grammars

The syntactic structure of a language is defined using *grammars*.

- Grammars (like regular expressions) specify a set of strings over an alphabet.
- Efficient *recognizers* (like DFA) can be constructed to efficiently determine whether a string is in the language.
- Language heirarchy:
  - Finite Languages (FL)  
Enumeration
  - Regular Languages (RL  $\supset$  FL)  
Regular Expressions
  - Context-free Languages (CFL  $\supset$  RL)  
Context-free Grammars

## Regular Languages

---

Languages represented by regular expressions  $\equiv$  Languages recognized by finite automata

---

Examples:

- ✓  $\{a, b, c\}$
- ✓  $\{\epsilon, a, b, aa, ab, ba, bb, \dots\}$
- ✓  $\{(ab)^n \mid n \geq 0\}$
- ×  $\{a^n b^n \mid n \geq 0\}$

## Grammars

Notation where recursion is explicit.

Examples:

- $\{\epsilon, a, b, aa, ab, ba, bb, \dots\}$ :

$$\begin{aligned} E &\longrightarrow a \\ E &\longrightarrow b \\ S &\longrightarrow \epsilon \\ S &\longrightarrow ES \end{aligned}$$

Notational shorthand:

$$\begin{aligned} E &\longrightarrow a \mid b \\ S &\longrightarrow \epsilon \mid ES \end{aligned}$$

- $\{a^n b^n \mid n \geq 0\}$  :

$$\begin{aligned} S &\longrightarrow \epsilon \\ S &\longrightarrow aSb \end{aligned}$$

- $\{w \mid \text{no. of } a\text{'s in } w = \text{no. of } b\text{'s in } w\}$

## Context-free Grammars

- **Terminal Symbols:** Tokens
- **Nonterminal Symbols:** set of strings made up of tokens
- **Productions:** Rules for constructing the set of strings associated with nonterminal symbols.  
Example:  $Stmt \longrightarrow \text{while } Expr \text{ do } Stmt$

**Start symbol:** nonterminal symbol that represents the set of all strings in the language.

### Example

$$\begin{aligned} E &\longrightarrow E + E \\ E &\longrightarrow E - E \\ E &\longrightarrow E * E \\ E &\longrightarrow E / E \\ E &\longrightarrow ( E ) \\ E &\longrightarrow \text{id} \end{aligned}$$

$$\mathcal{L}(E) = \{\text{id}, \text{id} + \text{id}, \text{id} - \text{id}, \dots, \text{id} + (\text{id} * \text{id}) - \text{id}, \dots\}$$

## Context-free Grammars

Production: rule with *nonterminal* symbol on left hand side, and a (possibly empty) sequence of terminal or nonterminal symbols on the right hand side.

Notations:

- **Terminals:** lower case letters, digits, punctuation
- **Nonterminals:** Upper case letters
- **Arbitrary Terminals/Nonterminals:**  $X, Y, Z$
- **Strings of Terminals:**  $u, v, w$
- **Strings of Terminals/Nonterminals:**  $\alpha, \beta, \gamma$
- **Start Symbol:**  $S$

## Context-Free Vs Other Types of Grammars

- Context-free grammar (CFG): Productions of the form  $NT \longrightarrow [NT|T]^*$
- Context-sensitive grammar (CSG): Productions of the form  $[t|NT]^* NT[t|NT]^* \longrightarrow [t|NT]^*$
- Unrestricted grammar: Productions of the form  $[t|NT]^* \longrightarrow [t|NT]^*$

## Examples of Non-Context-Free Languages

- Checking that variables are declared before use. If we simplify and abstract the problem, we see that it amounts to recognizing strings of the form  $ws w$
- Checking whether the number of actual and formal parameters match. Abstracts to recognizing strings of the form  $a^n b^m c^n d^m$
- In both cases, the rules are not enforced in grammar, but deferred to type-checking phase
- Note: Strings of the form  $ws w^R$  and  $a^n b^n c^m d^m$  can be described by a CFG

## What types of Grammars Describe These Languages?

- Strings of 0's and 1's of the form  $xx$
- Strings of 0's and 1's in which 011 doesn't occur
- Strings of 0's and 1's in which each 0 is immediately followed by a 1
- Strings of 0's and 1's with equal number of 0's and 1's.

## Language Generated by Grammars, Equivalence of Grammars

- How to show that a grammar  $G$  generates a language  $\mathcal{M}$ ? Show that
  - $\forall s \in \mathcal{M}$ , show that  $s \in \mathcal{L}(G)$
  - $\forall s \in \mathcal{L}(G)$ , show that  $s \in \mathcal{M}$
- How to establish that two grammars  $G_1$  and  $G_2$  are equivalent?  
Show that  $\mathcal{L}(G_1) = \mathcal{L}(G_2)$

### Grammar Examples

$$S \longrightarrow 0S1S \mid 1S0S \mid \epsilon$$

What is the language generated by this grammar?

### Grammar Examples

$$S \longrightarrow 0A \mid 1B \mid \epsilon$$

$$A \longrightarrow 0AA \mid 1S$$

$$B \longrightarrow 1BB \mid 0S$$

What is the language generated by this grammar?

## The Two Sides of Grammars

**Specify** a set of strings in a language.

**Recognize** strings in a given language:

- Is a given string  $x$  in the language?  
Yes, if we can construct a *derivation* for  $x$
- Example: Is  $\text{id} + \text{id} \in \mathcal{L}(E)$ ?

$$\begin{aligned} \text{id} + \text{id} &\longleftarrow E + \text{id} \\ &\longleftarrow E + E \\ &\longleftarrow E \end{aligned}$$

### Derivations

<b>Grammar:</b>	$E \rightarrow E + E$
	$E \rightarrow id$

$E$  derives  $id + id$ :  $E \Rightarrow E + E$   
 $\Rightarrow E + id$   
 $\Rightarrow id + id$

- $\alpha A \beta \Rightarrow \alpha \gamma \beta$  iff  $A \rightarrow \gamma$  is a production in the grammar.
- $\alpha \xRightarrow{*} \beta$  if  $\alpha$  derives  $\beta$  in zero or more steps.  
Example:  $E \xRightarrow{*} id + id$
- **Sentence:** A sequence of terminal symbols  $w$  such that  $S \xRightarrow{+} w$  (where  $S$  is the start symbol)
- **Sentential Form:** A sequence of terminal/nonterminal symbols  $\alpha$  such that  $S \xRightarrow{*} \alpha$

### Derivations

- **Rightmost derivation:** Rightmost nonterminal is replaced first:

$$\begin{aligned}
 E &\Rightarrow E + E \\
 &\Rightarrow E + id \\
 &\Rightarrow id + id
 \end{aligned}$$

Written as  $E \xRightarrow{*}_{rm} id + id$

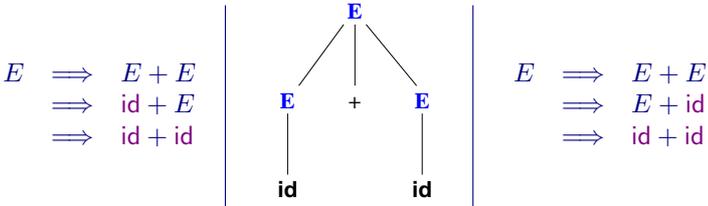
- **Leftmost derivation:** Leftmost nonterminal is replaced first:

$$\begin{aligned}
 E &\Rightarrow E + E \\
 &\Rightarrow id + E \\
 &\Rightarrow id + id
 \end{aligned}$$

Written as  $E \xRightarrow{*}_{lm} id + id$

### Parse Trees

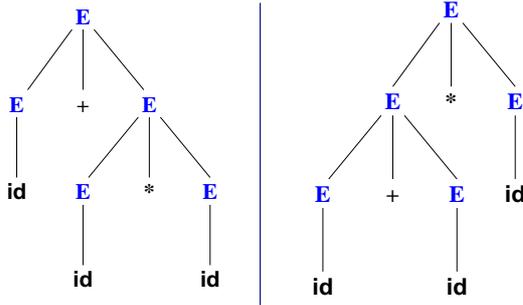
Graphical Representation of Derivations



A **Parse Tree** succinctly captures the *structure* of a sentence.

### Ambiguity

A Grammar is *ambiguous* if there are multiple parse trees for the same sentence.  
Example:  $id + id * id$

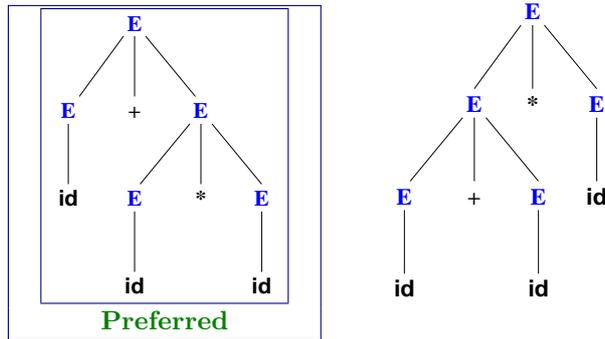


## Disambiguation

**Express Preference for one parse tree over others.**

Example:  $id + id * id$

The usual precedence of  $*$  over  $+$  means:



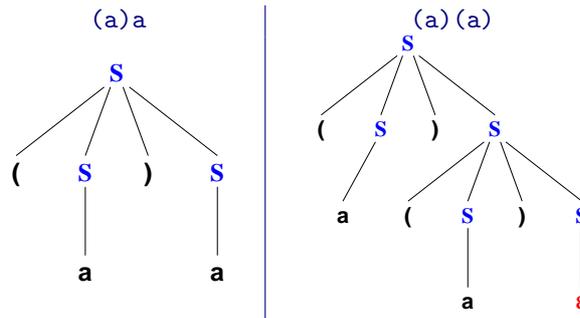
## Parsing

Construct a parse tree for a given string.

$S \rightarrow (S)S$

$S \rightarrow a$

$S \rightarrow \epsilon$



## A Procedure for Parsing

**Grammar:**  $S \rightarrow a$

```

procedure parse_S() {
  switch (input_token) {
    case TOKEN_a:
      consume(TOKEN_a);
      return;
    default:
      /* Parse Error */
  }
}

```

## Predictive Parsing

<b>Grammar:</b>	$S \rightarrow a$
	$S \rightarrow \epsilon$

---

```
procedure parse_S() {  
  switch (input_token) {  
    case TOKEN_a: /* Production 1 */  
      consume(TOKEN_a);  
      return;  
    case TOKEN_EOF: /* Production 2 */  
      return;  
    default:  
      /* Parse Error */  
  }  
}
```

---

### Predictive Parsing (Contd.)

<b>Grammar:</b>	$S \rightarrow (S)S$
	$S \rightarrow a$
	$S \rightarrow \epsilon$

---

```
procedure parse_S() {  
  switch (input_token) {  
    case TOKEN_OPEN_PAREN: /* Production 1 */  
      consume(TOKEN_OPEN_PAREN);  
      parse_S();  
      consume(TOKEN_CLOSE_PAREN);  
      parse_S();  
      return;  
  }
```

### Predictive Parsing (contd.)

<b>Grammar:</b>	$S \rightarrow (S)S$
	$S \rightarrow a$
	$S \rightarrow \epsilon$

```
    case TOKEN_a: /* Production 2 */  
      consume(TOKEN_a);  
      return;  
    case TOKEN_CLOSE_PAREN:  
    case TOKEN_EOF: /* Production 3 */  
      return;  
    default:  
      /* Parse Error */
```

---

### Predictive Parsing: Restrictions

## Grammar cannot be left-recursive

Example:  $E \rightarrow E + E \mid a$

```
procedure parse_E() {  
  switch (input_token) {  
    case TOKEN_a: /* Production 1 */  
      parse_E();  
      consume(TOKEN_PLUS);  
      parse_E();  
      return;  
    case TOKEN_a: /* Production 2 */  
      consume(TOKEN_a);  
      return;  
  }  
}
```

## Removing Left Recursion

$$\begin{aligned} A &\rightarrow A a \\ A &\rightarrow b \end{aligned}$$

$$\mathcal{L}(A) = \{b, ba, baa, baaa, baaaa, \dots\}$$

$$\begin{aligned} A &\rightarrow bA' \\ A' &\rightarrow aA' \\ A' &\rightarrow \epsilon \end{aligned}$$

## Removing Left Recursion

More generally,

$$\begin{aligned} A &\rightarrow A\alpha_1 \mid \dots \mid A\alpha_m \\ A &\rightarrow \beta_1 \mid \dots \mid \beta_n \end{aligned}$$

Can be transformed into

$$\begin{aligned} A &\rightarrow \beta_1 A' \mid \dots \mid \beta_n A' \\ A' &\rightarrow \alpha_1 A' \mid \dots \mid \alpha_m A' \mid \epsilon \end{aligned}$$

## Removing Left Recursion: An Example

$$\begin{aligned} E &\rightarrow E + E \\ E &\rightarrow \text{id} \end{aligned}$$

↓

$$\begin{aligned} E &\rightarrow \text{id } E' \\ E' &\rightarrow + E E' \\ E' &\rightarrow \epsilon \end{aligned}$$

## Predictive Parsing: Restrictions

May not be able to choose a *unique* production

$$\begin{aligned} S &\longrightarrow a B d \\ B &\longrightarrow b \\ B &\longrightarrow bc \end{aligned}$$

Left-factoring can help:

$$\begin{aligned} S &\longrightarrow a B d \\ B &\longrightarrow bC \\ C &\longrightarrow c|\epsilon \end{aligned}$$

## Predictive Parsing: Restrictions

In general, though, we may need a backtracking parser:

Recursive Descent Parsing

$$\begin{aligned} S &\longrightarrow a B d \\ B &\longrightarrow b \\ B &\longrightarrow bc \end{aligned}$$

## Recursive Descent Parsing

<b>Grammar:</b>	$S \longrightarrow a B d$
	$B \longrightarrow b$
	$B \longrightarrow bc$

```
procedure parse_B() {
  switch (input_token) {
    case TOKEN_b: /* Production 2 */
      consume(TOKEN_b);
      return;
    case TOKEN_b: /* Production 3 */
      consume(TOKEN_b);
      consume(TOKEN_c);
      return;
  }
}
```

## Nonrecursive Parsing

Instead of recursion,

use an explicit *stack* along with the parsing table.

Data objects:

- **Parsing Table:**  $M(A, a)$ , a two-dimensional array, dimensions indexed by nonterminal symbols ( $A$ ) and terminal symbols ( $a$ ).
- A **Stack** of terminal/nonterminal symbols
- **Input stream** of tokens

The above data structures manipulated using a *table-driven parsing program*.

## Table-driven Parsing

<b>Grammar:</b>	$A \rightarrow a$	$S \rightarrow A S B$
	$B \rightarrow b$	$S \rightarrow \epsilon$

Parsing Table:

NONTERMINAL	INPUT SYMBOL		
	a	b	EOF
$S$	$S \rightarrow A S B$	$S \rightarrow \epsilon$	$S \rightarrow \epsilon$
$A$	$A \rightarrow a$		
$B$		$B \rightarrow b$	

## Table-driven Parsing Algorithm

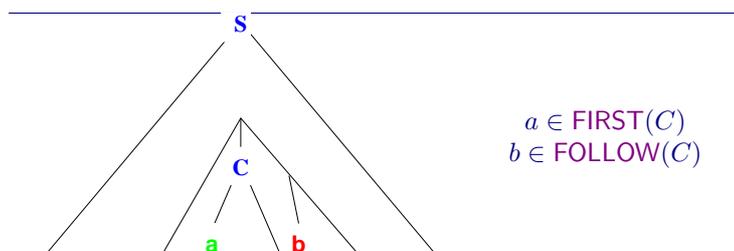
```

stack initialized to EOF.
while (stack is not empty) {
    X = top(stack);
    if (X is a terminal symbol)
        consume(X);
    else /* X is a nonterminal */
        if (M[X, input_token] = X → Y1, Y2, ..., Yk) {
            pop(stack);
            for i = k downto 1 do
                push(stack, Yi);
        }
    else /* Syntax Error */
}
    
```

## FIRST and FOLLOW

**Grammar:**  $S \rightarrow (S)S \mid a \mid \epsilon$

- $FIRST(X)$  = First character of any string that can be derived from  $X$   
 $FIRST(S) = \{ (, a, \epsilon \}$ .
- $FOLLOW(A)$  = First character that, in any derivation of a string in the language, appears immediately after  $A$ .  
 $FOLLOW(S) = \{ ), EOF \}$



## FIRST and FOLLOW

$FIRST(X)$ : First terminal in some  $\alpha$  such that  $X \xRightarrow{*} \alpha$ .  
 $FOLLOW(A)$ : First terminal in some  $\beta$  such that  $S \xRightarrow{*} \alpha A \beta$ .

<b>Grammar:</b>	$A \longrightarrow a$	$S \longrightarrow A S B$
	$B \longrightarrow b$	$S \longrightarrow \epsilon$

$First(S) = \{ a, \epsilon \}$        $Follow(S) = \{ b, EOF \}$   
 $First(A) = \{ a \}$        $Follow(A) = \{ a, b \}$   
 $First(B) = \{ b \}$        $Follow(B) = \{ b, EOF \}$

### Definition of FIRST

<b>Grammar:</b>	$A \longrightarrow a$	$S \longrightarrow A S B$
	$B \longrightarrow b$	$S \longrightarrow \epsilon$

$FIRST(\alpha)$  is the smallest set such that

$\alpha =$	Property of $FIRST(\alpha)$
$a$ , a terminal	$a \in FIRST(\alpha)$
$A$ , a nonterminal	$A \longrightarrow \epsilon \in G \implies \epsilon \in FIRST(\alpha)$ $A \longrightarrow \beta \in G, \beta \neq \epsilon \implies FIRST(\beta) \subseteq FIRST(\alpha)$
$X_1 X_2 \dots X_k$ , a string of terminals and nonterminals	$FIRST(X_1) - \{\epsilon\} \subseteq FIRST(\alpha)$ $FIRST(X_i) \subseteq FIRST(\alpha)$ if $\forall j < i \quad \epsilon \in FIRST(X_j)$ $\epsilon \in FIRST(\alpha)$ if $\forall j < k \quad \epsilon \in FIRST(X_j)$

### Definition of FOLLOW

<b>Grammar:</b>	$A \longrightarrow a$	$S \longrightarrow A S B$
	$B \longrightarrow b$	$S \longrightarrow \epsilon$

$FOLLOW(A)$  is the smallest set such that

$A$	Property of $FOLLOW(A)$
$= S$ , the start symbol	$EOF \in FOLLOW(S)$ Book notation: $\$ \in FOLLOW(S)$
$B \longrightarrow \alpha A \beta \in G$	$FIRST(\beta) - \{\epsilon\} \subseteq FOLLOW(A)$
$B \longrightarrow \alpha A$ , or $B \longrightarrow \alpha A \beta, \epsilon \in FIRST(\beta)$	$FOLLOW(B) \subseteq FOLLOW(A)$

### A Procedure to Construct Parsing Tables

```

procedure table_construct(G) {
  for each  $A \longrightarrow \alpha \in G$  {
    for each  $a \in FIRST(\alpha)$  such that  $a \neq \epsilon$ 
      add  $A \longrightarrow \alpha$  to  $M[A, a]$ ;
    if  $\epsilon \in FIRST(\alpha)$ 
      for each  $b \in FOLLOW(A)$ 
        add  $A \longrightarrow \alpha$  to  $M[A, b]$ ;
  }
}

```

## LL(1) Grammars

Grammars for which the parsing table constructed earlier has no multiple entries.

$$\begin{aligned} E &\longrightarrow \text{id } E' \\ E' &\longrightarrow + E E' \\ E' &\longrightarrow \epsilon \end{aligned}$$

NONTERMINAL	INPUT SYMBOL		
	id	+	EOF
$E$	$E \longrightarrow \text{id } E'$		
$E'$		$E' \longrightarrow + E E'$	$E' \longrightarrow \epsilon$

## Parsing with LL(1) Grammars

NONTERMINAL	INPUT SYMBOL		
	id	+	EOF
$E$	$E \longrightarrow \text{id } E'$		
$E'$		$E' \longrightarrow + E E'$	$E' \longrightarrow \epsilon$

$\$E$	id + id\$	$E \Rightarrow \text{id}E'$
$\$E'\text{id}$	id + id\$	
$\$E'$	+ id\$	$\Rightarrow \text{id}+EE'$
$\$E'E+$	+ id\$	
$\$E'E$	id\$	$\Rightarrow \text{id}+\text{id}E'E'$
$\$E'E'\text{id}$	id\$	
$\$E'E'$	\$	$\Rightarrow \text{id}+\text{id}E'$
$\$E'$	\$	$\Rightarrow \text{id}+\text{id}$
$\$$	\$	

## LL(1) Derivations

Left to Right Scan of input

Leftmost Derivation

(1) look ahead 1 token at each step

Alternative characterization of LL(1) Grammars:

Whenever  $A \longrightarrow \alpha \mid \beta \in G$

1.  $FIRST(\alpha) \cap FIRST(\beta) = \{ \}$ , and
2. if  $\alpha \xRightarrow{*} \epsilon$  then  $FIRST(\beta) \cap FOLLOW(A) = \{ \}$ .

**Corollary:** No Ambiguous Grammar is LL(1).

## Leftmost and Rightmost Derivations

$$\begin{aligned} E &\longrightarrow E+T \\ E &\longrightarrow T \\ T &\longrightarrow \text{id} \end{aligned}$$

Derivations for  $\text{id} + \text{id}$ :

$\begin{aligned} E &\Rightarrow E+T \\ &\Rightarrow T+T \\ &\Rightarrow \text{id}+T \\ &\Rightarrow \text{id}+\text{id} \end{aligned}$ <p style="text-align: center; color: red;">LEFTMOST</p>		$\begin{aligned} E &\Rightarrow E+T \\ &\Rightarrow E+\text{id} \\ &\Rightarrow T+\text{id} \\ &\Rightarrow \text{id}+\text{id} \end{aligned}$ <p style="text-align: center; color: red;">RIGHTMOST</p>
------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	--	---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

## Bottom-up Parsing

Given a stream of tokens  $w$ , *reduce* it to the start symbol.

$$\begin{array}{l} \hline E \longrightarrow E+T \\ E \longrightarrow T \\ T \longrightarrow \text{id} \\ \hline \end{array}$$

Parse input stream:  $\text{id} + \text{id}$ :

$$\begin{array}{l} \text{id} + \text{id} \\ T + \text{id} \\ E + \text{id} \\ E + T \\ E \end{array}$$

**Reduction**  $\equiv$  **Derivation**<sup>-1</sup>.

## Handles

Informally, a “handle” of a string is a substring that matches the right side of a production, and whose reduction to the nonterminal on the left hand side of the production represents one step along the reverse rightmost derivation.

## Handles

A structure that furnishes a means to perform reductions.

$$\begin{array}{l} \hline E \longrightarrow E+T \\ E \longrightarrow T \\ T \longrightarrow \text{id} \\ \hline \end{array}$$

Parse input stream:  $\text{id} + \text{id}$ :

$$\begin{array}{l} \boxed{\text{id}} + \text{id} \\ \boxed{T} + \text{id} \\ E + \boxed{\text{id}} \\ \boxed{E + T} \\ E \end{array}$$

## Handles

Handles are substrings of sentential forms:

1. A substring that matches the right hand side of a production
2. Reduction using that rule can lead to the start symbol

$$\begin{array}{l} E \implies \boxed{E + T} \\ \implies E + \boxed{\text{id}} \\ \implies \boxed{T} + \text{id} \\ \implies \boxed{\text{id}} + \text{id} \end{array}$$

**Handle Pruning:** replace handle by corresponding LHS.

## Shift-Reduce Parsing

Bottom-up parsing.

- **Shift:** Construct leftmost handle on top of stack

- **Reduce:** Identify handle and replace by corresponding RHS
- **Accept:** Continue until string is reduced to start symbol and input token stream is empty
- **Error:** Signal parse error if no handle is found.

## Implementing Shift-Reduce Parsers

- **Stack** to hold grammar symbols (corresponding to tokens seen thus far).
- **Input stream** of yet-to-be-seen tokens.
- **Handles** appear on top of stack.
- Stack is initially empty (denoted by \$).
- Parse is successful if stack contains only the start symbol when the input stream ends.

## Shift-Reduce Parsing: An Example

$$\begin{array}{l}
 \hline
 S \rightarrow aABe \\
 A \rightarrow Abc|b \\
 B \rightarrow d \\
 \hline
 \end{array}
 \quad \text{To parse: } a b b c d e$$

## Shift-Reduce Parsing: An Example

$$\begin{array}{l}
 \hline
 E \rightarrow E+T \\
 E \rightarrow T \\
 T \rightarrow id \\
 \hline
 \end{array}$$

STACK	INPUT STREAM	ACTION
\$	id + id \$	shift
\$ id	+ id \$	reduce by $T \rightarrow id$
\$ T	+ id \$	reduce by $E \rightarrow T$
\$ E	+ id \$	shift
\$ E +	id \$	shift
\$ E + id	\$	reduce by $T \rightarrow id$
\$ E + T	\$	reduce by $E \rightarrow E+T$
\$ E	\$	<b>ACCEPT</b>

## More on Handles

**Handle:** Let  $S \xRightarrow{*}_{rm} \alpha Aw \xRightarrow{rm} \alpha \beta w$ .

Then  $A \rightarrow \beta$  is a handle for  $\alpha \beta w$  at the position immediately following  $\alpha$ .

**Notes:**

- For unambiguous grammars, every right-sentential form has a unique handle.
- In shift-reduce parsing, handles always appear on top of stack, i.e.,  $\alpha \beta$  is in the stack (with  $\beta$  at top), and  $w$  is unread input.

## Identification of Handles and Relationship to Conflicts

**Case 1:** With  $\alpha \beta$  on stack, don't know if we have a handle on top of stack, or we need to shift some more input to get  $\beta x$  which is a handle.

- Shift-reduce conflict
- Example: if-then-else

**Case 2:** With  $\alpha \beta_1 \beta_2$  on stack, don't know if  $A \rightarrow \beta_2$  is the handle, or  $B \rightarrow \beta_1 \beta_2$  is the handle

- Reduce-reduce conflict
- Example:  $E \rightarrow E - E \mid - E \mid id$

## Viable Prefix

Prefix of a right-sentential form that does not continue beyond the rightmost handle.

With  $\alpha \beta w$  example of the previous slides, a viable prefix is something of the form  $\alpha \beta_1$  where  $\beta = \beta_1 \beta_2$

## LR Parsing

- Stack contents as  $s_0 X_1 s_1 X_2 \cdots X_m s_m$
- Its actions are driven by two tables, *action* and *goto*

Parser Configuration:  $(\underbrace{s_0 X_1 s_1 X_2 \cdots X_m s_m}_{\text{stack}}, \underbrace{a_i a_{i+1} \cdots a_n \$}_{\text{unconsumed input}})$

$action[s_m, a_i]$  can be:

- shift  $s$ : new config is  $(s_0 X_1 s_1 X_2 \cdots X_m s_m a_i s, a_{i+1} \cdots a_n \$)$
- reduce  $A \rightarrow \beta$ : Let  $|\beta| = r$ ,  $goto[s_{m-r}, A] = s$ : new config is  $(s_0 X_1 s_1 X_2 \cdots X_{m-r} s_{m-r} A s, a_i a_{i+1} \cdots a_n \$)$
- error: perform recovery actions
- accept: Done parsing

## LR Parsing

- *action* and *goto* depend only on the state at the top of the stack, not on all of the stack contents
  - The  $s_i$  states compactly summarize the “relevant” stack content that is at the top of the stack.
- You can think of *goto* as the action taken by the parser on “consuming” (and shifting) nonterminals
  - similar to the shift action in the *action* table, except that the transition is on a nonterminal rather than a terminal
- The *action* and *goto* tables define the transitions of an FSA that accepts RHS of productions!

## Example of LR Parsing Table and its Use

- See Text book Algorithm 4.7: (follows directly from description of LR parsing actions 2 slides earlier)
- See expression grammar (Example 4.33), its associated parsing table in Fig 4.31, and the use of the table to parse  $id * id + id$  (Fig 4.32)

## LR Versus LL Parsing

Intuitively:

- LL parser needs to guess the production based on the first symbol (or first few symbols) on the RHS of a production
- LR parser needs to guess the production *after* seeing all of the RHS

Both types of parsers can use next  $k$  input symbols as look-ahead symbols (LL( $k$ ) and LR( $k$ ) parsers)

- Implication:  $LL(k) \subset LR(k)$

## How to Construct LR Parsing Table?

**Key idea:** Construct an FSA to recognize RHS of productions

- States of FSA remember which parts of RHS have been seen already.
- We use “ $\cdot$ ” to separate seen and unseen parts of RHS

**LR(0) item:** A production with “ $\cdot$ ” somewhere on the RHS. Intuitively,

- ▷ grammar symbols before the “ $\cdot$ ” are on stack;
- ▷ grammar symbols after the “ $\cdot$ ” represent symbols in the input stream.

$$\begin{array}{l}
 \hline
 I_0: \quad E' \longrightarrow \cdot E \\
 \quad \quad E \longrightarrow \cdot E + T \\
 \quad \quad E \longrightarrow \cdot T \\
 \quad \quad T \longrightarrow \cdot id \\
 \hline
 \end{array}$$

## How to Construct LR Parsing Table?

- If there is no way to distinguish between two different productions at some point during parsing, then the same state should represent both.
  - *Closure* operation: If a state  $s$  includes LR(0) item  $A \longrightarrow \alpha \cdot B\beta$ , and there is a production  $B \longrightarrow \gamma$ , then  $s$  should include  $B \longrightarrow \cdot \gamma$
  - *goto* operation: For a set  $I$  of items,  $goto[I, X]$  is the closure of all items  $A \longrightarrow \alpha X \cdot \beta$  for each  $A \longrightarrow \alpha \cdot X\beta$  in  $I$

**Item set:** A set of items that is closed under the *closure* operation, corresponds to a state of the parser.

## Constructing Simple LR (SLR) Parsing Tables

**Step 1:** Construct LR(0) items (Item set construction)

**Step 2:** Construct a DFA for recognizing items

**Step 3:** Define *action* and *goto* based from the DFA

## Item Set Construction

1. Augment the grammar with a rule  $S' \longrightarrow S$ , and make  $S'$  the new start symbol

2. Start with initial set  $I_0$  corresponding to the item  $S' \rightarrow \cdot S$
3. apply *closure* operation on  $I_0$ .
4. For each item set  $I$  and grammar symbol  $X$ , add  $goto[I, X]$  to the set of items
5. Repeat previous step until no new item sets are generated.

### Item Set Construction

$$\frac{E' \rightarrow E}{I_0 : E' \rightarrow \cdot E} \quad \frac{E \rightarrow E + T \mid T}{I_3 : T \rightarrow F \cdot} \quad \frac{T \rightarrow T * F \mid F}{I_4 : F \rightarrow (\cdot E)} \quad \frac{F \rightarrow (E) \mid id}{I_5 : F \rightarrow id \cdot}$$

$$I_1 : E' \rightarrow E \cdot$$

$$I_2 : E \rightarrow T \cdot$$

### Item Set Construction (Continued)

$$\frac{E' \rightarrow E}{I_6 : E' \rightarrow E + \cdot T} \quad \frac{E \rightarrow E + T \mid T}{I_8 : F \rightarrow (E \cdot)} \quad \frac{T \rightarrow T * F \mid F}{I_9 : E \rightarrow E + T \cdot} \quad \frac{F \rightarrow (E) \mid id}{I_{10} : T \rightarrow T * F \cdot}$$

$$I_7 : T \rightarrow T * \cdot F$$

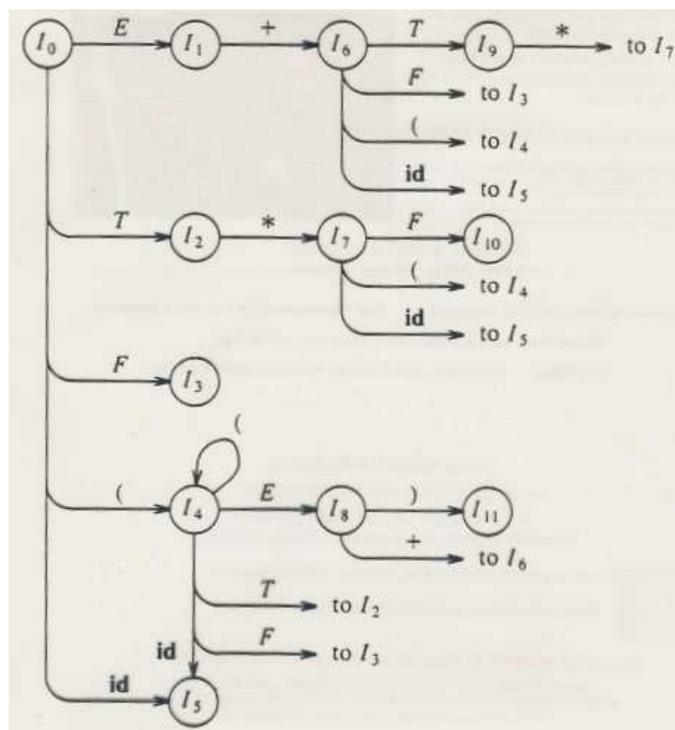
$$I_{10} : T \rightarrow T * F \cdot$$

$$I_{11} : F \rightarrow (E) \cdot$$

### Item Sets for the Example

$I_0: E' \rightarrow \cdot E$ $E \rightarrow \cdot E + T$ $E \rightarrow \cdot T$ $T \rightarrow \cdot T * F$ $T \rightarrow \cdot F$ $F \rightarrow \cdot (E)$ $F \rightarrow \cdot id$	$I_5: F \rightarrow id \cdot$
$I_1: E' \rightarrow E \cdot$ $E \rightarrow E \cdot + T$	$I_6: E \rightarrow E + \cdot T$ $T \rightarrow \cdot T * F$ $T \rightarrow \cdot F$ $F \rightarrow \cdot (E)$ $F \rightarrow \cdot id$
$I_2: E \rightarrow T \cdot$ $T \rightarrow T \cdot * F$	$I_7: T \rightarrow T * \cdot F$ $F \rightarrow \cdot (E)$ $F \rightarrow \cdot id$
$I_3: T \rightarrow F \cdot$	$I_8: F \rightarrow (E \cdot)$ $E \rightarrow E \cdot + T$
$I_4: F \rightarrow (\cdot E)$ $E \rightarrow \cdot E + T$ $E \rightarrow \cdot T$ $T \rightarrow \cdot T * F$ $T \rightarrow \cdot F$ $F \rightarrow \cdot (E)$ $F \rightarrow \cdot id$	$I_9: E \rightarrow E + T \cdot$ $T \rightarrow T \cdot * F$
	$I_{10}: T \rightarrow T * F \cdot$
	$I_{11}: F \rightarrow (E) \cdot$

### Constructing DFA to Recognize Viable Prefixes



### SLR(1) Parse Table for the Example Grammar

STATE	action					goto			
	id	+	*	(	)	\$	E	T	F
0	s5			s4			1	2	3
1		s6				acc			
2		r2	s7		r2	r2			
3		r4	r4		r4	r4			
4	s5			s4			8	2	3
5		r6	r6		r6	r6			
6	s5			s4				9	3
7	s5			s4					10
8		s6			s11				
9		r1	s7		r1	r1			
10		r3	r3		r3	r3			
11		r5	r5		r5	r5			

### Define action and goto tables

- Let  $I_0, I_1, \dots, I_n$  be the item sets constructed before
- Define *action* as follows
  - If  $A \rightarrow \alpha \cdot a\beta$  is in  $I_i$  and there is a DFA transition to  $I_j$  from  $I_i$  on symbol  $a$  then  $action[i, a] = \text{“shift } j\text{”}$
  - If  $A \rightarrow \alpha \cdot$  is in  $I_i$  then  $action[i, a] = \text{“reduce } A \rightarrow \alpha\text{”}$  for every  $a \in FOLLOW(A)$
  - If  $S' \rightarrow S \cdot$  is in  $I_i$  then  $action[I_i, \$] = \text{“accept”}$
- If any conflicts arise in the above procedure, then the grammar is *not* SLR(1).
- *goto* transition for LR parsing defined directly from the DFA transitions.
- All undefined entries in the table are filled with “error”

### Deficiencies of SLR Parsing

SLR(1) treats all occurrences of a RHS on stack as identical.  
 Only a few of these reductions may lead to a successful parse.  
 Example:

$$\begin{array}{l} \hline S \rightarrow AaAb \quad A \rightarrow \epsilon \\ S \rightarrow BbBa \quad B \rightarrow \epsilon \\ \hline \end{array}$$

$I_0 = \{[S' \rightarrow \cdot S], [S \rightarrow \cdot AaAb], [S \rightarrow \cdot BbBa], [A \rightarrow \cdot], [B \rightarrow \cdot]\}$ .  
 Since  $FOLLOW(A) = FOLLOW(B)$ , we have reduce/reduce conflict in state 0.

### LR(1) Item Sets

Construct LR(1) items of the form  $A \rightarrow \alpha \cdot \beta, a$ , which means:

The production  $A \rightarrow \alpha\beta$  can be applied when the next token on input stream is  $a$ .

$$\begin{array}{l} \hline S \rightarrow AaAb \quad A \rightarrow \epsilon \\ S \rightarrow BbBa \quad B \rightarrow \epsilon \\ \hline \end{array}$$

An example LR(1) item set:

$I_0 = \{[S' \rightarrow \cdot S, \$], [S \rightarrow \cdot AaAb, \$], [S \rightarrow \cdot BbBa, \$], [A \rightarrow \cdot, a], [B \rightarrow \cdot, b]\}$ .

## LR(1) and LALR(1) Parsing

**LR(1) parsing:** Parse tables built using LR(1) item sets.

**LALR(1) parsing:** *Look Ahead* LR(1)

Merge LR(1) item sets; then build parsing table.

Typically, LALR(1) parsing tables are much smaller than LR(1) parsing table.

## YACC

Yet Another Compiler Compiler:

LALR(1) parser generator.

- Grammar rules written in a specification (.y) file, analogous to the regular definitions in a lex specification file.
- Yacc translates the specifications into a parsing function `yyparse()`.

`spec.y`  $\xrightarrow{\text{yacc}}$  `spec.tab.c`

- `yyparse()` calls `yylex()` whenever input tokens need to be consumed.
- `bison`: GNU variant of `yacc`.

## Using Yacc

```
%{
... C headers (#include)
%}

... Yacc declarations:
    %token ...
    %union{...}
    precedences
%%
... Grammar rules with actions:

Expr:  Expr TOK_PLUS Expr
      |  Expr TOK_MINUS Expr
      ;
%%
... C support functions
```

## YACC

Yet Another Compiler Compiler:

LALR(1) parser generator.

- Grammar rules written in a specification (.y) file, analogous to the regular definitions in a lex specification file.
- Yacc translates the specifications into a parsing function `yyparse()`.

`spec.y`  $\xrightarrow{\text{yacc}}$  `spec.tab.c`

- `yyparse()` calls `yylex()` whenever input tokens need to be consumed.
- `bison`: GNU variant of `yacc`.

## Using Yacc

```

%{
    ... C headers (#include)
}%

... Yacc declarations:
    %token ...
    %union{...}
    precedences
%%
... Grammar rules with actions:

Expr:  Expr TOK_PLUS Expr
      | Expr TOK_MINUS Expr
      ;
%%
... C support functions

```

## Conflicts and Resolution

- Operator precedence works well for resolving conflicts that involve operators
  - But use it with care – only when they make sense, not for the sole purpose of removing conflict reports
- Shift-reduce conflicts: Bison favors shift
  - Except for the dangling-else problem, this strategy does not ever seem to work, so don't rely on it.

## Reduce-Reduce Conflicts

```

sequence: /* empty */
         { printf ("empty sequence\n"); }
         | maybeward
         | sequence word
         { printf ("added word %s\n", $2); };

maybeward: /* empty */
          { printf ("empty maybeward\n"); }
          | word
          { printf ("single word %s\n", $1); };

```

In general, grammar needs to be rewritten to eliminate conflicts.

## Sample Bison File: Postfix Calculator

```

input:    /* empty */
         | input line
;
line:    '\n'
         | exp '\n'    { printf ("\t%.10g\n", $1); }
;
exp:     NUM          { $$ = $1;          }
         | exp exp '+' { $$ = $1 + $2;    }
         | exp exp '-' { $$ = $1 - $2;    }
         | exp exp '*' { $$ = $1 * $2;    }
         | exp exp '/' { $$ = $1 / $2;    }
         /* Exponentiation */
         | exp exp '^' { $$ = pow ($1, $2); }
         /* Unary minus */
         | exp 'n'    { $$ = -$1;        };
%%

```

## Infix Calculator

```

%{
#define YYSTYPE double
#include <math.h>
#include <stdio.h>
int yylex (void);
void yyerror (char const *);
%}
/* Bison Declarations */
%token NUM
%left '-' '+'
%left '*' '/'
%left NEG    /* negation--unary minus */
%right '^'   /* exponentiation */

```

## Infix Calculator (Continued)

```

%% /* The grammar follows. */
input:    /* empty */
         | input line
;
line:    '\n'
         | exp '\n' { printf ("\t%.10g\n", $1); }
;
exp:     NUM          { $$ = $1;          }
         | exp '+' exp { $$ = $1 + $3;    }
         | exp '-' exp { $$ = $1 - $3;    }
         | exp '*' exp { $$ = $1 * $3;    }
         | exp '/' exp { $$ = $1 / $3;    }
         | '-' exp %prec NEG { $$ = -$2;    }
         | exp '^' exp { $$ = pow ($1, $3); }
         | '(' exp ')' { $$ = $2;        }
;
%%

```

## Error Recovery

```
line:      '\n'  
         | exp '\n'  { printf ("\t%.10g\n", $1); }  
         | error '\n' { yyerror;                };
```

- Pop stack contents to expose a state where `error` token is acceptable
- Shift error token onto the stack
- Discard input until reaching a token that can follow this error token

Error recovery strategies are never perfect — some times they lead to cascading errors, unless carefully designed.

## Left Versus Right Recursion

```
expseq1:  exp | expseq1 ',' exp;
```

is a left-recursive definition of a sequence of `exp`'s, whereas

```
expseq1:  exp | exp ',' expseq1;
```

is a right-recursive definition

- Left-recursive definitions are no-no for LL parsing, but yes-yes for LR parsing
- Right-recursive definition is bad for LR parsing as it needs to shift entire list on stack before any reduction — increases stack usage