

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 the 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 hierarchy:
 - 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

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

$$E \longrightarrow a$$

$$E \longrightarrow b$$

$$S \longrightarrow \epsilon$$

$$S \longrightarrow ES$$

Notational shorthand:

$$E \longrightarrow a \mid b$$

$$S \longrightarrow \epsilon \mid ES$$

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

$$S \longrightarrow \epsilon$$

$$S \longrightarrow aSb$$

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

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 non-terminal 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

$$E \longrightarrow E + E$$

$$E \longrightarrow E - E$$

$$E \longrightarrow E * E$$

$$E \longrightarrow E / E$$

$$E \longrightarrow (E)$$

$$E \longrightarrow \text{id}$$

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

Context-free Grammars

Production: rule with *non-terminal* symbol on left hand side, and a (possibly empty) sequence of terminal or non-terminal 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:** α, β, γ
- **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 wsw
- 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 wsw^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 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 the 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

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

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

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

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 $\text{id} + \text{id} \in \mathcal{L}(E)$?

$$\text{id} + \text{id} \longleftarrow E + \text{id}$$

$$\longleftarrow E + E$$

$$\longleftarrow E$$

Derivations

Grammar: $E \rightarrow E + E$
 $E \rightarrow id$

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 α derives β 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 α such that $S \xRightarrow{*} \alpha$

Derivations

- Rightmost derivation: Rightmost non-terminal is replaced first:

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

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

- Leftmost derivation: Leftmost non-terminal is replaced first:

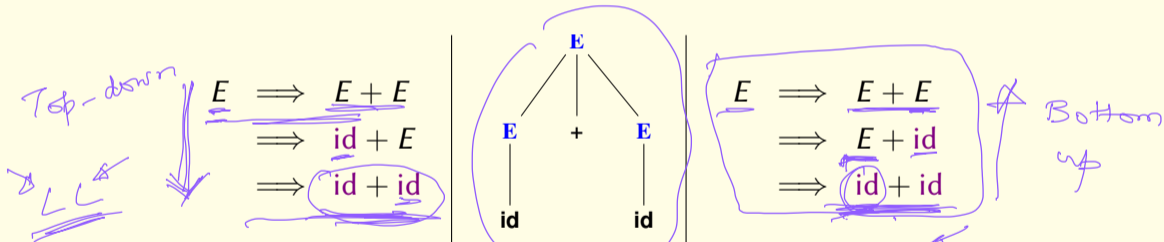
$$\begin{aligned} E &\Longrightarrow E + E \\ &\Longrightarrow \text{id} + E \\ &\Longrightarrow \text{id} + \text{id} \end{aligned}$$

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

Parse Trees

$$E \rightarrow \underline{id}$$

Graphical Representation of Derivations



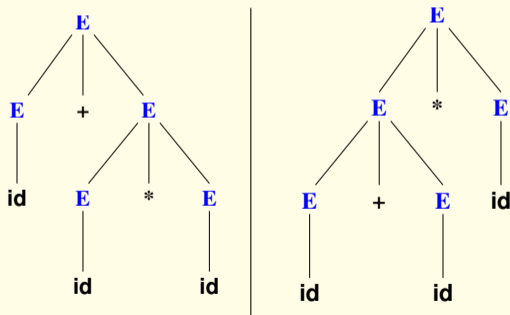
A **Parse Tree** succinctly captures the structure of a sentence.

Recursive descent \rightarrow recursive, backtracks
Predictive parsing LL Antlr

Ambiguity

A Grammar is *ambiguous* if there are multiple parse trees for the same sentence.

Example: $\text{id} + \text{id} * \text{id}$



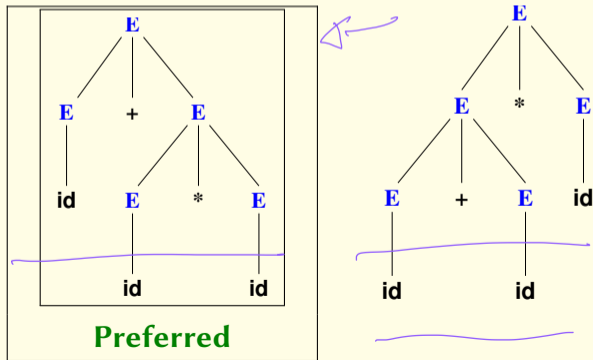
Disambiguation

Express Preference for one parse tree over others.

Example: $id + id * id$

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

$5ab$



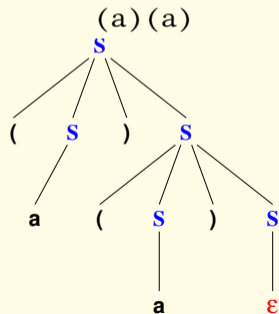
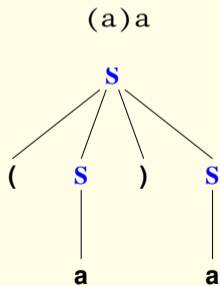
Parsing

Construct a parse tree for a given string.

$$S \rightarrow (S)S$$

$$S \rightarrow a$$

$$S \rightarrow \epsilon$$



A Procedure for Parsing

Grammar: $S \rightarrow a$

→ Top-down

→ Bottom-up

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

Predictive Parsing

Grammar: $S \rightarrow a$
 $S \rightarrow \epsilon$

```
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.) ↩

Grammar:

$S \rightarrow (S)S$
$S \rightarrow a$
$S \rightarrow \epsilon$

$S \rightarrow (A) | (B)$

$\text{parse_S}()$

$\text{parse_A}()$

$\text{parse_B}()$

```
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;  
  }
```

$LL(1)$

$LL(k)$

$LR(1)$

Predictive Parsing (contd.)

Grammar:

$$S \longrightarrow (S)S$$
$$S \longrightarrow a$$
$$S \longrightarrow \epsilon$$

```
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

$$A \longrightarrow A a$$

$$A \longrightarrow b$$

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

$$A \longrightarrow bA'$$

$$A' \longrightarrow aA'$$

$$A' \longrightarrow \epsilon$$

Removing Left Recursion

More generally,

$$A \longrightarrow A\alpha_1 | \cdots | A\alpha_m$$

$$A \longrightarrow \beta_1 | \cdots | \beta_n$$

Can be transformed into

$$A \longrightarrow \beta_1 A' | \cdots | \beta_n A'$$

$$A' \longrightarrow \alpha_1 A' | \cdots | \alpha_m A' | \epsilon$$

Removing Left Recursion: An Example

$$E \longrightarrow E + E$$

$$E \longrightarrow \text{id}$$

⇓

$$E \longrightarrow \text{id } E'$$

$$E' \longrightarrow + E E'$$

$$E' \longrightarrow \epsilon$$

Predictive Parsing: Restrictions

May not be able to choose a *unique* production

$$S \longrightarrow a B d$$

$$B \longrightarrow b$$

$$B \longrightarrow bc$$

Left-factoring can help:

$$S \longrightarrow a B d$$

$$B \longrightarrow bC$$

$$C \longrightarrow c|\epsilon$$

Predictive Parsing: Restrictions

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

Recursive Descent Parsing

$$S \longrightarrow a B d$$

$$B \longrightarrow b$$

$$B \longrightarrow bc$$

Recursive Descent Parsing

Grammar:

$$S \longrightarrow a B d$$
$$B \longrightarrow b$$
$$B \longrightarrow bc$$

```
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;  
  }  
}
```

Non-recursive 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

Grammar:

$$\begin{array}{ll} A \longrightarrow a & S \longrightarrow A S B \\ B \longrightarrow b & S \longrightarrow \epsilon \end{array}$$

Parsing Table:

NONTERMINAL	INPUT SYMBOL		
	a	b	EOF
<i>S</i>	$S \longrightarrow A S B$	$S \longrightarrow \epsilon$	$S \longrightarrow \epsilon$
<i>A</i>	$A \longrightarrow a$		
<i>B</i>		$B \longrightarrow b$	

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 ( $M[X, input\_token] = X \longrightarrow Y_1, Y_2, \dots, Y_k$ ) {  
            pop(stack);  
            for i = k downto 1 do  
                push(stack, Yi);  
        }  
    else /* Syntax Error */  
}
```

FIRST and FOLLOW

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

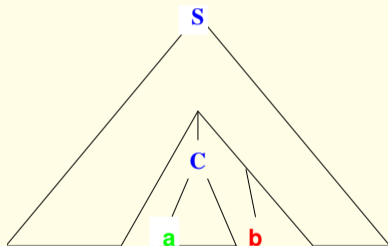
- **FIRST**(X) = First character of any string that can be derived from X

$$\text{FIRST}(S) = \{ (, a, \epsilon \}.$$

- **FOLLOW**(A) = First character that, in any derivation of a string in the language, appears immediately after A .

$$\text{FOLLOW}(S) = \{), \text{EOF} \}$$

FIRST and FOLLOW (contd.)



$a \in \text{FIRST}(C)$
 $b \in \text{FOLLOW}(C)$

FIRST and FOLLOW

FIRST(X): First terminal in some α such that $X \xRightarrow{*} \alpha$.

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

Grammar:

$A \rightarrow a$	$S \rightarrow A S B$
$B \rightarrow b$	$S \rightarrow \epsilon$

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

$First(A) = \{ a \}$ $Follow(A) = \{ a, b \}$

$First(B) = \{ b \}$ $Follow(B) = \{ b, EOF \}$

Definition of FIRST

Grammar:

$$\begin{array}{ll} A \longrightarrow a & S \longrightarrow A S B \\ B \longrightarrow b & S \longrightarrow \epsilon \end{array}$$

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

$\alpha =$	Property of $FIRST(\alpha)$
a , a terminal	$a \in FIRST(\alpha)$
A , a nonterminal	$A \longrightarrow \epsilon \in G \implies \epsilon \in FIRST(\alpha)$ $A \longrightarrow \beta \in G, \beta \neq \epsilon \implies FIRST(\beta) \subseteq FIRST(\alpha)$
$X_1 X_2 \cdots X_k$, a string of terminals and non-terminals	$FIRST(X_1) - \{\epsilon\} \subseteq FIRST(\alpha)$ $FIRST(X_i) \subseteq FIRST(\alpha)$ 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

Grammar: $A \rightarrow a$ $S \rightarrow A S B$
 $B \rightarrow b$ $S \rightarrow \epsilon$

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

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

A Procedure to Construct Parsing Tables

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


LL(1) Grammars

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

$$\begin{array}{l} E \longrightarrow \text{id } E' \\ E' \longrightarrow + E E' \\ E' \longrightarrow \epsilon \end{array}$$

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	$E \rightarrow \text{id } E'$		
E'		$E' \rightarrow + E E'$	$E' \rightarrow \epsilon$

$\$E$	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 $A \rightarrow \alpha \mid \beta \in G$

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

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

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

$E \Rightarrow E+T$	$E \Rightarrow E+T$
$\Rightarrow T+T$	$\Rightarrow E+\text{id}$
$\Rightarrow \text{id}+T$	$\Rightarrow T+\text{id}$
$\Rightarrow \text{id}+\text{id}$	$\Rightarrow \text{id}+\text{id}$
LEFTMOST	RIGHTMOST

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

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

Reduction \equiv Derivation⁻¹.


Handles

The image contains two handwritten diagrams in purple ink. The left diagram shows the string $A + A + id + id$. A box is drawn around the substring $A + id$. A larger box is drawn around the entire string $A + A + id + id$. The right diagram shows the production $A \rightarrow A + id$. A box is drawn around the right-hand side $A + id$.

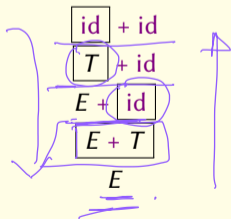
Informally, a “handle” of a sentential form is a substring that matches the right side of a production, and whose reduction to the non-terminal on the left hand side of the production represents one step along the reverse rightmost derivation.

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: **id + id**:



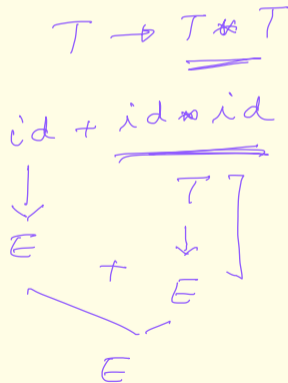
Handles

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

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

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.

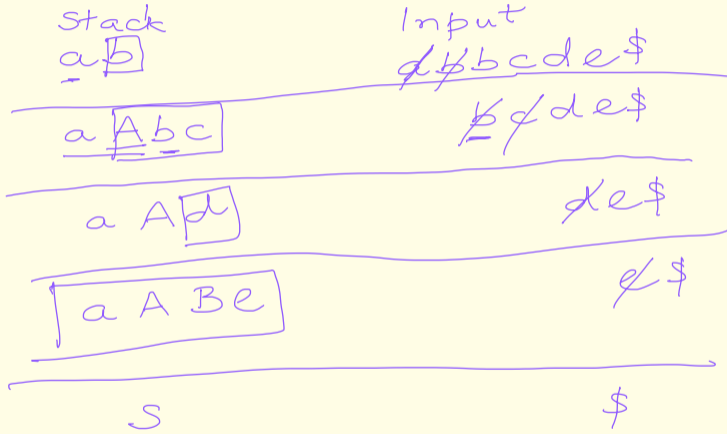
and nonterminals

terminals

Shift-Reduce Parsing: An Example

$S \rightarrow aAbe$
 $A \rightarrow Abc|b$
 $B \rightarrow d$

To parse: a b b c d e



Shift-Reduce Parsing: An Example

E	\rightarrow	$E+T$
E	\rightarrow	T
T	\rightarrow	id

$A \rightarrow \underline{b}$
 $B \rightarrow b$

STACK	INPUT STREAM	ACTION
\$	<u>id</u> + id \$	<u>shift</u>
\$ 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	\$	ACCEPT

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 α .

$\boxed{id + id} + id$

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 β at top), and w is unread input.

Identification of Handles and Relationship to Conflicts

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

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

Case 2: With $\alpha\beta_1\beta_2$ on the 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 $\alpha\underline{\beta_1}$ where $\beta = \underline{\beta_1}\beta_2$

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

$A \rightarrow \boxed{aBC}$

$C \rightarrow \boxed{de}$

Knutz

r states
 r non terminals

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

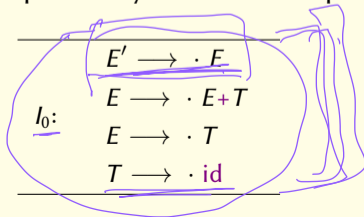
How to Construct LR Parsing Table?

Key idea: Construct an FSA to recognize RHS of productions

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

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

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



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

$E' \rightarrow E$
 $S' \rightarrow \cdot S$

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 $A \rightarrow \alpha \cdot B\beta$, and there is a production $B \rightarrow \gamma$, then s should include $B \rightarrow \cdot \gamma$
- goto operation: For a set I of items, $\text{goto}[I, X]$ is the closure of all items $A \rightarrow \alpha X \cdot \beta$ for each $A \rightarrow \alpha \cdot X\beta$ 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) → states

Step 2: Construct a DFA for recognizing items

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

Item Set Construction

1. Augment the grammar with a rule $S' \rightarrow S$, 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

$E' \rightarrow E$

$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

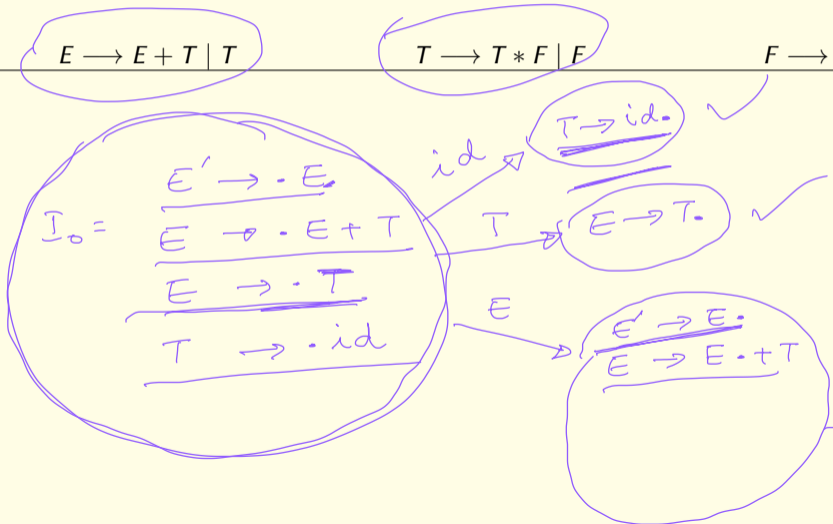
$F \rightarrow (E) \mid id$

$I_0: E' \rightarrow \cdot E$

$I_1: E' \rightarrow E \cdot$

$I_2: E \rightarrow T \cdot$

$I_3: T \rightarrow F \cdot$



Item Set Construction (Contd.)

$E' \rightarrow E$

$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

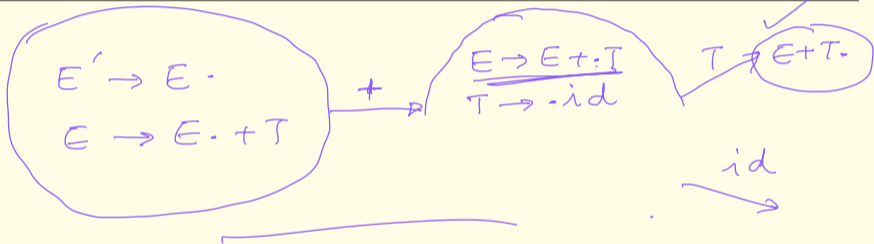
$F \rightarrow (E) \mid id$

$I_4 : F \rightarrow (\cdot E)$

$I_5 : F \rightarrow id \cdot$

$I_6 : E \rightarrow E + \cdot T$

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



Item Set Construction (Contd.)

$$\frac{E' \longrightarrow E \qquad E \longrightarrow E + T \mid T \qquad T \longrightarrow T * F \mid F \qquad F \longrightarrow (E) \mid id}{}$$

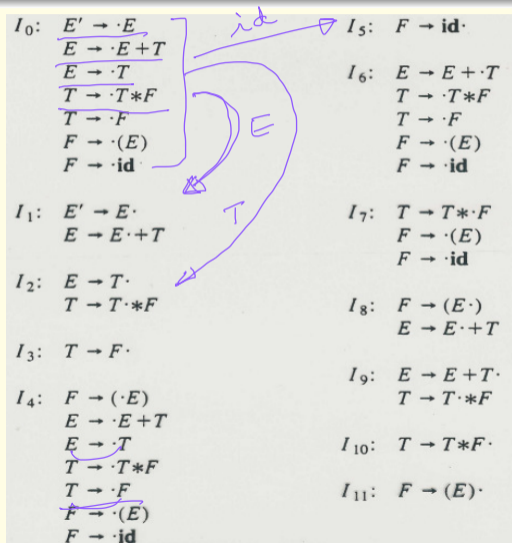
$$I_8 : F \longrightarrow (E \cdot)$$

$$I_9 : E \longrightarrow E + T \cdot$$

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

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

Item Sets for the Example



SLR(1) Parse Table for the Example Grammar

STATE	<u>action</u>					<u>goto</u>			
	<u>id</u>	+	*	()	\$	<i>E</i>	<i>T</i>	<i>F</i>
<u>0</u>	<u>s5</u>			s4			<u>1</u>	2	<u>3</u>
1		s6				acc			
2		r2	s7		r2	r2			
3		<u>r4</u>	r4		r4	r4			
4	s5			s4			8	2	3
<u>5</u>		<u>r6</u>	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			

id + id * id

~~id~~
|
0 F 3

(F)

Defining *action* and *goto* tables

- Let I_0, I_1, \dots, I_n be the item sets constructed before



- Define *action* as follows

LR(0)

- 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

$action[i, a] = \text{"shift } j\text{"}$

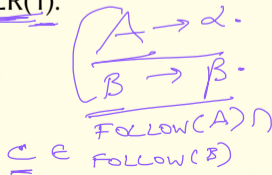
- If $A \rightarrow \alpha \cdot$ is in I_i then $action[i, a] = \text{"reduce } A \rightarrow \alpha\text{"}$ for every $a \in FOLLOW(A)$

- If $S' \rightarrow S \cdot$ is in I_i then $action[i, \$] = \text{"accept"}$

- 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"



Defining *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 $action[i, \$] = \text{"accept"}$

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"

$\hookrightarrow a \cdot B$

$\left[\begin{array}{l} A \rightarrow \alpha \cdot \\ A \rightarrow \alpha \cdot a B \end{array} \right]$

LR(0)

b

$\left[\begin{array}{l} A \rightarrow \alpha \cdot \\ B \rightarrow \beta \cdot \end{array} \right]$

$FOLLOW(A) \cap$

$\subseteq FOLLOW(B)$

Deficiencies of SLR Parsing

SLR = LR(0) item sets + $LR(1)$
1 - look ahead for reduction

SLR(1) treats all occurrences of a RHS on stack as identical.

Only a few of these reductions may lead to a successful parse.

$S \rightarrow ab \mid ba$

Example:

$S \rightarrow AaAb$	$A \rightarrow \epsilon$
$S \rightarrow BbBa$	$B \rightarrow \epsilon$

$FOLLOW(A) = \{a, b\}$

$FOLLOW(B) = \{a, b\}$

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

Since $FOLLOW(A) = FOLLOW(B)$, we have reduce/reduce conflict in state 0.

LR(1) Item Sets

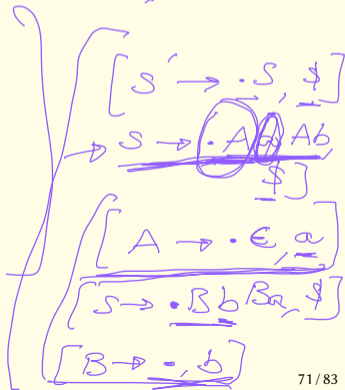
Construct LR(1) items of the form $A \rightarrow \alpha \cdot \beta, a$, which means:

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

$S \rightarrow AaAb$	$A \rightarrow \epsilon$
$S \rightarrow BbBa$	$B \rightarrow \epsilon$

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

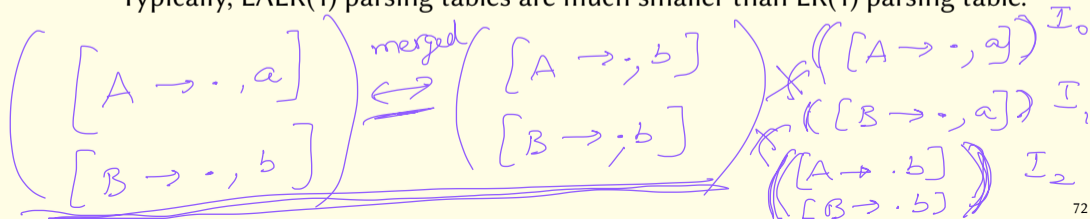
LR(1) parsing: Parse tables built using LR(1) item sets.

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.

LR(1) item sets
that are identical
except for the look
ahead
↳ merged



YACC

Yet Another Compiler Compiler:

LALR(1) parser generator.

- Grammar rules are 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 are 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()`.

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

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

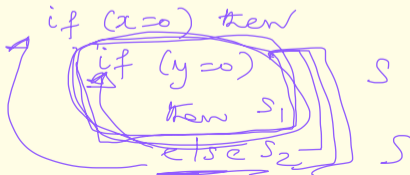
Using Yacc

```
%{  
    ... 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

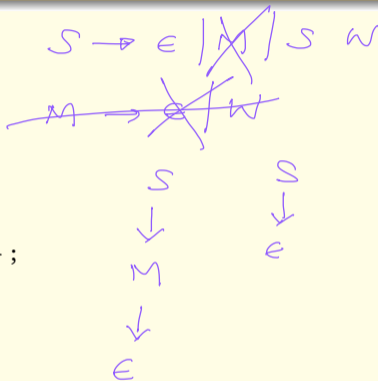
if Stmt \rightarrow if Expr then S
if Expr then S else S

- 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 ylex (void);  
void yyerror (char const *);  
%}  
/* Bison Declarations */  
%token NUM  
%left '-', '+',  
%left '*', '/'  
%left NEG /* negation--unary minus */  
%right '^' /* exponentiation */
```

$$5 + 5 * 3$$
$$\frac{(5 + 5) * 3}{5 + (5 * 3)}$$
$$\frac{55 + 3 *}{553 * +}$$

$$\frac{(-5) * 5}{5 - (5 * 5)}$$

lower precedence

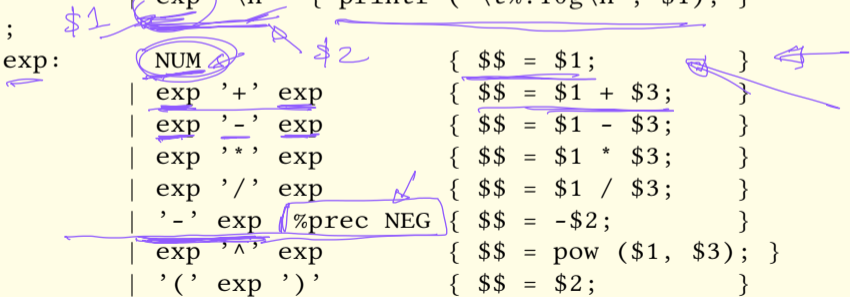
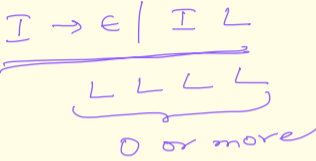
higher

Infix Calculator (Continued)

```

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

```



Error Recovery

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

5 + (x) 5 \n
\$
~~exp~~ error

- Pop stack contents to expose a state where an 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

$\text{expseq1: exp | expseq1 ', ' exp;}$

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

$\text{expseq1: exp | exp ', ' expseq1;}$

is a right-recursive definition

LL
5, 5, 5, 5, 5

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