

CS 460 - Compilers

Arian Izadi

Spring 2024

1 Languages

Syntax is the rules for what a syntactically correct program looks like. Semantics is the meaning of a program.

When does it matter the order of evaluation (right to left vs left to right)? When the code has side effects, an example of this is postfix vs prefix increment (a++ vs ++a).

Compilers for a language L, move from front end \rightarrow intermediate representation \rightarrow back end.

- Front end: Lexical Analysis, Syntax Analysis, and Semantic Analysis
- Intermediate: Intermediate Code
- Back end: Optimizer and Code Generation

1.1 Lexical Analysis & Scanning

Lexical analysis, a scanner, is the process of converting a stream of characters into a stream of tokens.

1. Find all terminals in the grammar.
2. Write the Scanner.
 - (a) Do we use a DFA, NFA, or PDA?
 - (b) Look at token types. All tokens can be expressed by a regular expression.
 - i. Symbols: Semicolon, commas, etc.
 - ii. Keywords: for, while, etc.
 - iii. Variables: x, y, etc.
 - iv. Numbers: 1, 3.14, 0x64, etc.

Chomsky Language Hierarchy

- Type 0: Unrestricted (Turing Machines)
- Type 1: Context Sensitive
- Type 2: Context Free (PDA)
- Type 3: Regular Expressions (NFA, DFA)

Both RE and CFG have 1 non-terminal on the left of any combination of terminals and non-terminals on the right.

Example 1:

$S \rightarrow X$ $X \rightarrow aXb|d$ not regular: $a^n db^n$

Example 2:

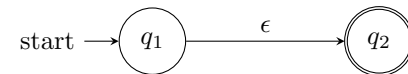
$S \rightarrow X$ $X \rightarrow aX|b$ regular: a^*b

Example 3:

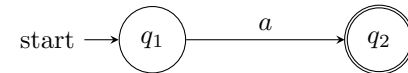
$S \rightarrow X$ $X \rightarrow aY|\epsilon$ $Y \rightarrow bX$ regular: $(ab)^*$

An NFA for recognizing tokens, construct NFA for each construct of RE.

\mathcal{E} :



$a\epsilon\Sigma$:



Any RE can be turned into an NFA using these rules. If all the tokens of a language are represented by RE's, r_1, \dots, r_n . Create an NFA for each RE.

2 Parsing

2.1 Types of Parsers

- $LL(k)$ – leftmost derivation
 - (Hard Explain) Always develop the leftmost non terminal in a sentential form.
 - Above means: Start at the start symbol and work towards the input.
- $LR(k)$ – rightmost derivation
 - Reverse rightmost derivation.
 - Above means: Start at input and work backwards to the start symbol.

Types of $LL(k)$ parsers:

- $LL(0)$ – No look ahead and no left recursion.
- $LL(1)$ – One look ahead and no left recursion.
- $LL(k)$ – k look ahead and no left recursion.

Example 1: a^*b^* $\delta(a^*b^*) = \epsilon, a, b, aa, ab, bb, \dots$

$S \rightarrow A \quad A \rightarrow aA|B \quad B \rightarrow bB|\epsilon$

Try $w = aabb$:

$S \xrightarrow{LM} A \xrightarrow{LM} aA \xrightarrow{LM} aaA \xrightarrow{LM} aaB \xrightarrow{LM} aabB \xrightarrow{LM} aabbB \xrightarrow{LM} aabb$

What if B could generate strings starting with an a ?

$S \rightarrow A \quad A \rightarrow aA|B \quad B \rightarrow bB|a|\epsilon$

Try $w = aabb$:

It cannot be parsed by an $LL(1)$ parser.

A is not $LL(1)$ when

- $\exists w_1 \in \lambda(\alpha_1)$ or $\exists w_2 \in \lambda(\alpha_2)$
- $w_1 = aw_1$ or $w_2 = aw_2$

Example 2: A set representing tokens that a sentential form can start with: Call it $FIRST(\alpha)$

1. $\epsilon : FIRST(\epsilon) = \{\epsilon\}$
2. $a \in \Sigma : FIRST(a) = \{a\}$
3. a mixed string: $Y = Y_1Y_2 \dots Y_n$
 - (a) $FIRST(Y_1 \dots Y_n) \supseteq FIRST(Y_1) \cup \{\epsilon\}$
 - (b) if $\epsilon \in FIRST(Y_1)$:
 - i. $Y_1, Y_2, \dots, Y_n \rightarrow \epsilon$
 - ii. $Y_1, Y_2, \dots, Y_n \rightarrow Y_2, \dots, Y_n$

Example 3: $X \rightarrow X_1 | \dots | X_m$

$FIRST(X) = FIRST(X_1) \cup FIRST(X_2) \cup \dots \cup FIRST(X_m)$

Example 4: General Case

For $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_n$

if $\exists i, j$ such that $i \neq j$ and $FIRST(\alpha_i) \cap FIRST(\alpha_j) \neq \emptyset$ then A is not $LL(1)$.

MUST BE DONE PAIRWISE, only needs to exist one pair to not be $LL(1)$.

$\forall i, j : i \neq j (1 \leq i, j \leq n) : FIRST(\alpha_i) \cap FIRST(\alpha_j) = \emptyset$ then A is probably (necessary, but not sufficient) $LL(1)$.

Example 5: General Case

For $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_n$

\vdots

$A_m \rightarrow \alpha_{m1} | \alpha_{m2} | \dots | \alpha_{mn_m}$

Example 6:

$$A \rightarrow B|C$$

$$\vdots$$

$$B \rightarrow \dots$$

$$C \rightarrow \dots$$

$$\text{FIRST}(B) = \{a, b\} \text{ FIRST}(C) = \{c, f\}$$

$$w = a\beta \text{ where } a \in \Sigma \text{ and } \beta \in \Sigma^*$$

* a is the lookahead token.

$$A \rightarrow_{LM} B \rightarrow_{LM} aX$$

* aX matches $a\beta$, so we can derive β from X .

2.2 Grammar Transformation**2.2.1 Left vs Right Recursion**

$$A \rightarrow Aa|b$$

$$\delta(A) = \{b, ba, baa, \dots\} = ba^*$$

$$\text{FIRST}(A) = \{b\} \text{ FIRST}(Aa) = \{b\}$$

* Both have $b \rightarrow$ and is not $LL(1)$.

* Just make it right recursive.

$$A \rightarrow bA'$$

$$A' \rightarrow aA'|\epsilon$$
In General:

$$A \rightarrow A\alpha_1|A\alpha_2|\dots|A\alpha_n|\beta_1|\beta_2|\dots|\beta_m$$

$$\delta(A) = (B_1|\dots|B_n)(\alpha_1|\dots|\alpha_n)$$

$$A \rightarrow \beta_1A'|\dots|\beta_mA'$$

$$A' \rightarrow \alpha_1A'|\dots|\alpha_nA'|\epsilon$$
Example 4.6:

$$E \rightarrow E + b|b \text{ similar to } A \rightarrow A\alpha|\beta$$

$$E \rightarrow bE'$$

$$E' \rightarrow +bE'|\epsilon$$
Example 4.7:

$$E \rightarrow E + b|E - c|b|C \text{ similar to } A \rightarrow A\alpha_1|A\alpha_2|\beta_1|\beta_2$$

$$E \rightarrow bE'|CE'$$

$$E' \rightarrow +bE'|-cE'|\epsilon$$
2.2.2 Common Subexpressions**Example 4.8:**

$$A \rightarrow abC|abD$$

* is this $LL(1)$?

No, because $\text{FIRST}(abC) \cap \text{FIRST}(abD) \neq \emptyset$.

* if we change it to $A \rightarrow abA'$ and $A' \rightarrow c|D$, then it is $LL(1)$.

* is this $LL(1)$?

Only if $\text{FIRST}(c) \cap \text{FIRST}(D) = \emptyset$.

For the following:

$$A \rightarrow A_1|\dots|A_n$$

Let α be the longest common prefix of A_1, \dots, A_n such that $A_i = \alpha A'_i$. Then $A \rightarrow \alpha A'$ and $A' \rightarrow A_1|\dots|A_n$.

2.2.3 Grammar Ambiguity

Definition: A grammar is ambiguous if there exists a string w such that there are two or more different parse trees for the same string.

$$A \rightarrow aA|bA|c \text{ with } w = ababc.$$

$$A \rightarrow aA \rightarrow abA \rightarrow abaA \rightarrow ababA \rightarrow ababc$$
Simple Example:

$$A \rightarrow aA|a|\epsilon$$
Example 4.9:

$$E \rightarrow E + E|E - E|id \text{ is this ambiguous?}$$

Yes, because $id + id - id$ can be parsed as $(id + id) - id$ or $id + (id - id)$.

PHASE 2 ADVICE: AST is never the answer to the left hand side, use a subclass of it.

What if we want +, -, *, /?

$E \rightarrow E + E | E - E | E * E | E / E | id$

Consider $id + id * id$.

* is higher precedence than +, so $id + (id * id)$.

When dealing with precedence, the lower precedence should be higher in the parse tree.

$E \rightarrow E + T | E - T | T$

$T \rightarrow T * F | T / F | F$

$F \rightarrow id | (E)$

* The bracket is the highest precedence.

* This is also bad because it is left recursive.

Eliminate Left Recursion:

$E \rightarrow TE'$

$E' \rightarrow +TE' | -TE' | \epsilon$

$T \rightarrow FT'$

$T' \rightarrow *FT' | /FT' | \epsilon$

$F \rightarrow id | (E)$

```

parse_A() {
    if (lookahead == 'b') {
        parse_B();
    } else if (lookahead == 'c') {
        parse_C();
    } else {
        // will not check error vs epsilon
        // this is a problem
        // if none of alpha i are epsilon
    }
}

```

2.3 Recursive Descent Parsing

$A \rightarrow B | C$

$B \rightarrow b | \epsilon$

$C \rightarrow c$

* lookahead = \$

$A \rightarrow B \rightarrow \epsilon$

```

parse_x();
parse_y();
match('a');
parse_z();

```