

# **Compilers**

2018A7PS0193P



## CHAPTER 1

### Process of Translation

The process of compilers converting from source code (the high level language) to target code (the machine language) is known as translation. It consists of many steps, which are described below.

#### 1. Lexical Analysis

Lexical analysis follows the following steps:

- (1) Identify the valid set of characters in the language
- (2) Break the sequence of characters into appropriate words or tokens (keywords, numbers, operators, etc.)
- (3) Find out whether these tokens are valid or not.

The key goal of the lexical analyzer is to break a sentence into a series of words/tokens. These breaks are generally done via certain separators. The tokens are recognized via some rules encoded into a Finite State Machine.

During the lexical phase, we may experience the following lexical errors:

- Occurrence of illegal characters
- Exceeding the length of identifier

The output of the lexical analyzer will be a sequence of tokens and their “type”, which signifies whether it is a identifier, operator, etc.

#### 2. Syntax Analysis

The syntax analysis takes the sequence of tokens as an input, and generates a parse tree. In case the syntax is not correct according to the grammar rules, it flags a syntactical error. This is modelled using Context Free Grammars that will be recognized using PDAs or Table Driven processes.

### 3. Semantic Analysis

Semantic analysis takes the parse tree as an input, and outputs a disambiguated parse tree. It performs the following:

- Check Semantics
- Error reporting (types, etc.)
- Disambiguate overloaded operators (meaning of operators depends on operands)
- Type coercion (type casting)
- Uniqueness checking (redeclaration of variables)

As such, the disambiguated parse tree gives us an unambiguous representation of the parse tree.

The phases mentioned till now comprise the **front end** of the compiler, where the source code is handled. After this, the compiler works on generating the target code.

### 4. Code optimization

This is an optional phase that modifies the programs to run faster and consume less resources like memory, registers, etc. However, it will not change the representation of the program.

Some examples of machine independent code optimization done is:

- **Common sub-expression elimination:** The compiler searches for instances of identical expressions and analyzes whether it is worthwhile to replace them with a single variable holding the computer value.
- **Copy Propagation:** The compiler replaces the occurrences of targets of direct assignments with their value. For instance, if we had the code `y=x; z = 3+y;`, copy propagation would yield `z=3+x;`
- **Dead code elimination:** The compiler removes code which does not affect the program results, including unreachable code and unused variables.
- **Code Motion:** The compiler moves statements and expressions out of the body of a loop if they are loop-invariant, i.e., doing so does not affect program semantics.
- **Strength Reduction:** The compiler replaces expensive operations with equivalent but less expensive operations.
- **Constant Folding:** The compiler recognizes and evaluates constant expressions at compile time instead of runtime.

## 5. Code Generation

This is the process of mapping from source level abstractions (identifiers, values, etc) to target machine abstractions (registers, memory, etc.). This is a two step process - initially, intermediate code gets generated from the disambiguated parse tree, which is used to generate the final machine code.

During code generation, we have to do the following:

- Map identifiers to locations (memory or registers)
- Map source code operators to opcodes or sequences of opcodes.
- Transform conditionals and iterations to a test/jump or compare instructions
- We use layout parameter passing protocols - the locations for parameters, return values, etc.

## 6. Post Translation Optimizations

Unlike in the code optimization phase where we perform machine independent code optimizations, this does machine-dependent code optimizations. This is an optional phase as well, where we may remove unneeded operations or rearrange to prevent hazards. It is a flexible phase, and may occur at any time in the back-end of the compiler.

## 7. Symbol Table

The symbol table contains information required about the source program identifiers during compilation, including:

- Category of variable
- Data type
- Quantity stored in structure
- Scope information
- Address in Memory

The symbol table must be present in every phase of the compiler, and is used in all the phases to get information about the identifiers.

## 8. Advantages and Disadvantages of Compilers

The advantages of compilers are:

- Highly modular in nature

- It is retargetable. This means that if there is a single language and multiple machines, then we can use the same front end.
- Source code and machine independent optimizations are possible.

The limitations of the compiler are:

- Design of programming languages has a huge effect on the performance of compilers.
- Lots of work is repeatable. For  $S$  languages and  $M$  machines,  $S \cdot M$  compilers are needed. This is known as the  $S * M$  problem of compilers.

The  $S * M$  problem is generally solved by introducing some common intermediate language, called the **Universal Intermediate Language Generator**. Some common machine independent intermediate code generation techniques are:

- Postfix Notation
- Three Address code
- Syntax tree
- Directed Acyclic Graph

## CHAPTER 2

# Lexical Analysis

### 1. Functions of the Lexical Analyzer

The lexical analyzer performs the following functions:

- Take high level language as input and output a sequence of tokens
- It generally cleans the code, by stripping off blanks, tabs newlines and comments.
- Keeps track of the line numbers for associated error messages

The lexical analyzer is modelled using regular expressions. As such, its implementation is done with a DFA. An example of one rule is  $L \cdot (L + D)^*$ , where  $L$  refers to a letter and  $D$  refers to a digit.

DEFINITION 1.1. A token is a string of characters which logically belong together, e.g. keywords, number, identifiers, etc.

DEFINITION 1.2. A pattern is the set of strings for which the same token is produced.

DEFINITION 1.3. A lexeme is a sequence of characters matched by a pattern to form the corresponding token.

Now that we understand the definitions, we can see what the lexical analyzer actually does - it transforms strings to the token and passes the lexeme as its corresponding attribute. For instance, the integer 43 would become `<num,43>`.

### 2. Working of the Lexical Analyzer

The lexical analyzer reads the character one by one from the source code into the lexeme. When it reaches a separator, it assigns a token to the lexeme based on certain rules, and continues to read the characters once more.

However, reading the lexemes character by character is slow, and involves many IO operations. This is done from a buffer instead of directly from the file. Moreover, the prefix of a lexeme is often not enough to determine the token - think of the lexemes `=` and `==`. We

instead use a lookahead pointer to determine the appropriate token for a lexeme, and then push back the characters that we do not need in the current lexeme.

### 3. Symbol Table and the Lexical Analyzer

The lexical analyzer also interfaces with the symbol table. When the lexical analyzer discovers a lexeme constituting an identifier, it enters that lexeme to the symbol table. Sometimes, information regarding the token of a particular lexeme may also be store in the symbol table. As such, the symbol table must implement the following operations:

- (1) `insert(s,t)` : Save lexeme `s` and token `t` and return pointer.
- (2) `lookup(s)` : return the index of entry for lexeme `s` or '0' if `s` is not found.

To make the symbol table space efficient, we save lexemes in some separate memory, and instead store pointers to the lexemes in the symbol table.

The rule for identifying an identifier and a keyword is generally the same. To be able to tokenize the identifiers and keywords separately, we initialize the symbol table with the list of keywords, say, by calling `insert("if",keyword)`.

### 4. Challenges in Development of Lexical Analyzer

- **Free vs Fixed Lexemes** : A language could specify that lexemes must be in a free or a fixed format. For instance, in a free format, code could look like this.

```
flag = flag
* 6;
```

But in the case of fixed format, this must be entirely in one line. An example of a fixed format language is Python, while a free format language is C.

- **Whitespaces** : How do we deal with whitespaces? Some languages ignore whitespaces until a separator is reached (or interpret contextually), while some languages consider the whitespaces as separators themselves. The former is much more complicated to implement than the latter.
- **Maximal Munch** : The principle of maximal munch directs the lexical analyzer to consume as much available input as possible while creating a construct. This allows us to deal with lexemes like `iff`, and correctly assign it as a identifier rather than the keyword `if`.



### 5. Techniques for specifying tokens

DEFINITION 5.1. Consider  $R_i$  is a regular expression and  $N_i$  is a unique name, then a regular definition is a series of definitions of the following form

$$N_1 \rightarrow R_1$$

$$N_2 \rightarrow R_2$$

...

$$N_n \rightarrow R_n$$

where each  $R_i$  is a regular expression over  $\sum \cup \{N_1, N_2, \dots, N_n\}$ .

Hence, by assigning a special name  $N_i$  to the regular expression  $R_i$ , we are in effect defining macros, that remove redundancy in later parts.

The following is an example regular definition for identifiers:

$$\text{Alphabet} \rightarrow A|B|C|\dots|Z|a|b|c|\dots|z$$

$$\text{Digit} \rightarrow 0|1|2|\dots|9$$

$$\text{Identifier} \rightarrow \text{Alphabet}(\text{Alphabet}|\text{Digit})^*$$

This too comes with its own challenges. Regular expressions often fail when identifying the appropriate token, and may pass the invalid tokens to the subsequent translation phases of the compiler (how?). They are only language specifications. Tokenization is a implementation problem.

Tokenization can be done via the following steps:

- (1) Construct regular expressions for lexemes of each token
- (2) Construct  $R$  matching all lexemes of tokens, so  $R = R_1 + R_2 + \dots$ , in some well defined precedence order.
- (3) Consider the input stream to be  $S = s_1s_2\dots s_n$ . For  $i \in [1, n]$ , verify whether  $s_1\dots s_i \in L(R)$ .
- (4) If  $s_1\dots s_i \in L(R) \implies s_1\dots s_i \in L(R_x)$  for some  $x$ . We choose the smallest  $x$  to be the class of  $s_1\dots s_i$ .
- (5) Discard the tokenized input and go back to step 3.

The procedure gives preference to tokens specified earlier using regular expressions. If  $s_1..s_i \in L(R)$  and  $s_1..s_j \in L(R)$ , we choose the longest prefix, in accordance with the principle of maximal munch.

To implement our regular definitions and recognize tokens, we use **transition diagrams**. They are shown diagrammatically in the same way as Finite Automata. Transitions can be labelled with a symbol, a group of symbols, or regular definitions. A few states may be **retracting states** that indicates that the lexeme does not include the symbol that can bring us to the accepting state.

Sometimes, we may want to push back extra characters into the token stream (think of  $>$  and  $>=$ , we may want to push back the extra character read if the token is  $>$ ). We mark those states with a  $*$ , to show that we must push back extra characters.

Let us consider the example of hexadecimal and octal constant. The regular definition would be:

$$hex \rightarrow 0|1|2|\dots|9|A|B|C|\dots F$$

$$oct \rightarrow 0|1|2|\dots|7$$

$$Qualifier \rightarrow u|U|l|L$$

$$OctalConstant \rightarrow 0oct^+(Qualifier|\epsilon)$$

$$HexadecimalConstant \rightarrow 0(x|X)hex^+(Qualifier|\epsilon)$$

The transition diagram for the given regular definition is given in Fig 1.

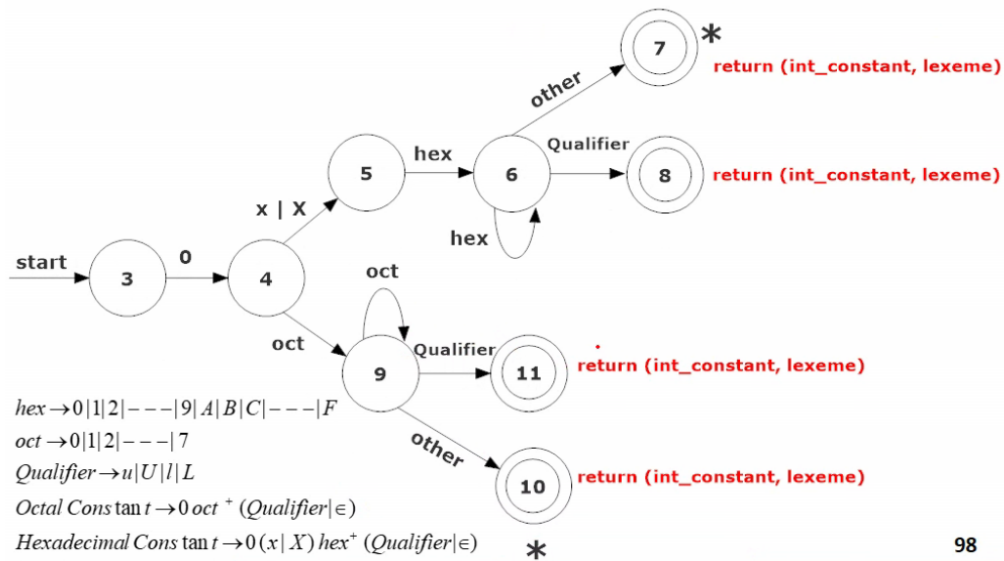


FIGURE 1. Transition diagram for the Hex and Octal constants

The retracting states are not given. If we get a “bad character”, we should report it as a lexical error.

Let us consider another example of a generalized expression for unsigned numbers. The regular definition is:

$$\textit{Digit} \rightarrow 0|1|2|\dots|9$$

$$\textit{Digits} \rightarrow \textit{Digit}^+$$

$$\textit{Fraction} \rightarrow \textit{.Digits}|\epsilon$$

$$\textit{Exponent} \rightarrow (E(+|-|\epsilon)\textit{Digits})|\epsilon$$

$$\textit{Number} \rightarrow \textit{Digits} \cdot \textit{Fraction} \cdot \textit{Exponent}$$

The resulting transition diagram from this is in Fig 2.

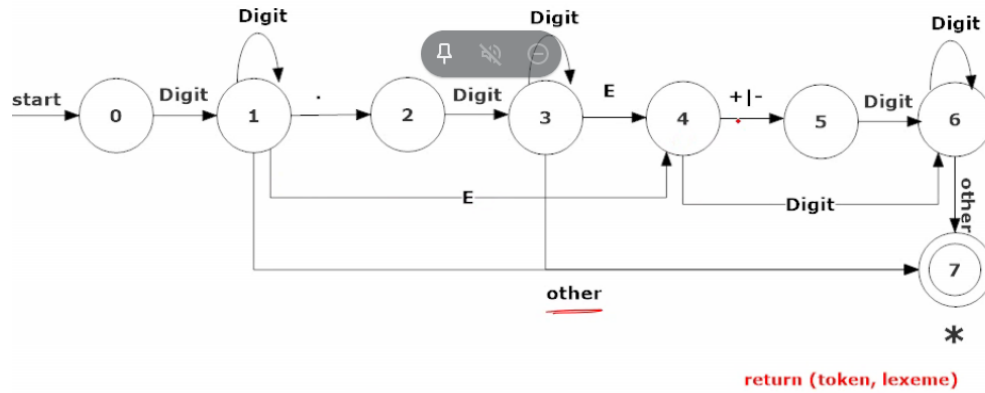


FIGURE 2. Transition Diagram for Unsigned Numbers

The issue here is this often introduces a lot of complexity. Instead, we may add split this into multiple transition diagrams, to improve implementation complexity. This also is seen empirically to give better speeds. These multiple transition diagrams can be appropriately combined to generate a lexical analyzer.

The matching process should always start with some transition diagram. If failure occurs in one transition diagram, we retract the forward pointer to the start state, and activate the next transition diagram. If there is failure in all the transition diagrams, we throw a lexical error.



## CHAPTER 3

# Syntax Analysis

### 1. Grammar Rules

Regular definitions are not enough for syntax analysis, since languages are generally not regular. Instead, we use Context Free Grammars to model our language, given by  $G = \{T, N, P, S\}$ , where

- $T$  is the set of tokens (or terminals)
- $N$  is the set of non-terminals
- $P$  is the set of production rules
- $S$  is the start symbol.

However, just creating the CFG is not the end of our problems - we must choose a method of derivation (left vs right), choose the right non-terminal to expand, choose the right production rule and maintain some precedence order.

When we replace the leftmost non-terminal in any sentential form of the grammar, it is called **leftmost derivation**. Choosing the rightmost non-terminal makes it **rightmost derivation**. If a grammar generates more than one leftmost/rightmost derivation of a string, it is said to be **ambiguous**.

The end goal of our syntax analysis is to create the **parse tree** - where the root is a start symbol, internal nodes are non-terminals and leaf nodes are labelled by tokens. An ambiguous grammar can have more than one parse tree.

Unfortunately, no well defined algorithm exists to remove ambiguity in a grammar. However, we can manually convert any ambiguous grammar into an unambiguous grammar. We can do this by enforcing associativity and precedence rules. This is pretty well covered in PPL and TOC.

Remember that right associativity is created via right recursion ( $N \rightarrow \alpha N$ ) and left associativity is created via left recursion ( $N \rightarrow N\alpha$ )

In programming languages, ambiguity can come in many forms. For instance, consider the sentence `if A then if B then C1 else C2`. This could be:

```

if A then
{
    if B then
        C1
    else
        C2
}

```

or:

```

if A then
{
    if B then
        C1
}
else
    C2

```

Since there is ambiguity in its meaning, we must resolve it. One way to do this is to relate an **if** with the closest unrelated **else** (the second interpretation). Another way is to match each **else** with the closest previous **then** (the first interpretation). The unambiguous grammar would then be:

$$stmt \rightarrow matchedstmt | unmatchedstmt$$

$$matchedstmt \rightarrow \text{if } exp \text{ then } matchedstmt \text{ else } matchedstmt | others$$

$$unmatchedstmt \rightarrow \text{if } exp \text{ then } stmt | \text{if } exp \text{ then } matchedstmt \text{ else } unmatchedstmt$$

## 2. Parsing

Resolving ambiguity is the problem of the language specification, not of the syntax analysis. Parsing is the job of implementation. We can do this in two ways:

- **Top down parsing:** Construction of the parse tree starting from the root node (start symbol) and proceeds towards the leaves (terminals).
- **Bottom up parsing:** Construction of the parse tree starts from the terminal nodes and builds up towards the root node (start symbol).

### 3. Top Down Parsing

In top down parsing, we repeat the following two steps:

- (1) At a node labelled with non-terminal  $A$ , select one of the productions of  $A$  and construct the children nodes.
- (2) Find the next node at which the subtree is to be constructed.

To choose the production rule, as in step 1, backtracking must be done. We try every production rule that we can apply and see which one works. The parser must be intelligent enough to choose the production rule to apply based on the token being pointed to by the pointer.

### 4. First Set

If we could determine the first character of the string produced when a production rule is applied, we could compare it to the current token and choose the production rule correctly. For this, we use a concept called the First Set.

The **First set**, denoted by  $\text{FIRST}(X)$  for a grammar symbol  $X$  is the set of tokens that begin the strings derivable from  $X$ . If there is a production rule:

$$A \rightarrow \alpha$$

then  $\text{First}(A)$  is the set of tokens that appear as the first token in strings generated from  $\alpha$ .

To compute  $\text{FIRST}(X)$ , we use the following rules:

- (1) If  $X$  is a terminal, then  $\text{FIRST}(X) = \{X\}$
- (2) If  $X$  is a non terminal and  $X \rightarrow Y_1Y_2...Y_k$  is a production for  $k \geq 1$ , then place  $a$  in  $\text{First}(X)$  if for some  $i$ ,  $a$  is in  $\text{First}(Y_i)$  and  $\epsilon$  is in all of  $\text{First}(Y_j)$  for  $j < i$ . If  $\epsilon$  is in all  $Y_i$ , then add  $\epsilon$  to  $\text{First}(X)$ .
- (3) If  $X \rightarrow \epsilon$ , then add  $\epsilon$  to  $\text{First}(X)$ .

Consider the following example:

$$S \rightarrow ABCDE$$

$$A \rightarrow a|\epsilon$$

$$B \rightarrow b|\epsilon$$

$$C \rightarrow c$$

$$D \rightarrow d|\epsilon$$

$$E \rightarrow e|\epsilon$$

Then, the first sets are as following:

- $\text{First}(S) = \{a, b, c\}$
- $\text{First}(A) = \{a, \epsilon\}$
- $\text{First}(B) = \{b, \epsilon\}$
- $\text{First}(C) = \{c\}$
- $\text{First}(D) = \{d, \epsilon\}$
- $\text{First}(E) = \{e, \epsilon\}$

### 5. Recursive Descent Parser

Let us look at our first parser - the **recursive descent parser**. RDP is a top down method of syntax analysis in which a set of recursive procedures are executed to parse the stream of tokens. A procedure is associated with each non-terminal of the grammar. The procedure generally implements one of the production rules of the grammar.

In each procedure, we perform a match operation on hitting any token on the RHS of the grammar with the current token in the input that needs to be parsed. For example, if we have the string *aba* and the rule  $S \rightarrow abB$ , it will match the tokens *a* in both the RHS and the string, and move on to the next token. If there is no match, we throw a syntax error.

This approach, of course, has its own limitations. Consider a grammar with two productions:

$$X \rightarrow \gamma_1$$

$$X \rightarrow \gamma_2$$

Suppose  $\text{First}(\gamma_1) \cap \text{First}(\gamma_2) \neq \emptyset$ , and that *a* is the common terminal symbol. In cases like this, we need to perform backtracking to choose the right production rule. However, RDP does not support backtracking right now. To support it, all productions should be tried in some order. Failure for some productions implies we need to try remaining productions. We would report an error only when there are no other rules.

Moreover, a recursive descent parser may loop forever on productions of the form:

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

This is known as left recursion. We can remove it from the grammar by rewriting the rule as:

$$A \rightarrow \beta_1 A' | \beta_2 A' | \dots | \beta_m A'$$

$$A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \epsilon$$



Left recursion can be hidden as well, where expanding production rules in some order could lead to left recursion. This is called hidden left recursion, and we must handle it in the same way as we did left recursion.

Left factoring is the process of removing the common left factor that appears in two productions of the same non-terminal. An example of a rule to remove is:

$$A \rightarrow \alpha\beta_1|\alpha\beta_2|\dots|\alpha\beta_n$$

We can remove the left factoring and get:

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1|\beta_2|\dots|\beta_n$$

In the initial form, the parser would be confused about which production rule to try. With the second case, we can kick the can down the road, and use a (hopefully) single lookahead pointer to decide what to do with  $A'$ .

## 6. Follow Set

Consider the following grammar:

$$A \rightarrow aBb$$

$$B \rightarrow c|\epsilon$$

And suppose we want to parse an input string “ab”. How would the parser know to use  $\epsilon$  for  $B$ ? To do this we use the **follow set**.

Follow( $X$ ) for a non terminal  $X$  is the set of symbols that might follow the derivation of  $X$  in an input stream. The follow of the start symbol  $S$  is  $\{\$$ . The follow set can never be  $\epsilon$ , since if it ends with an epsilon, the follow would be the ending symbol of the string - the  $\$$ . The steps to compute the follow set is:

- Always include  $\$$  in Follow( $S$ ).
- If there is a production rule  $A \rightarrow \alpha B \beta$ , where  $B$  is a non terminal, then Follow( $B$ ) = First( $\beta$ ).
- If there is a production rule  $A \rightarrow \alpha B \beta$ , where  $B$  is a non terminal, and First( $\beta$ ) contains  $\epsilon$ , then everything in Follow( $A$ ) is in Follow( $B$ ).
- If there is a production rule  $A \rightarrow \alpha B$ , then Follow( $B$ ) = Follow( $A$ ).

## 7. Predictive Parsing

Now that we have talked about all the issues that can face RDP, let us look at a new parser - **predictive parsing**. This is a non recursive top down parsing method, which recognizes LL(1) languages. LL languages are those that can be parsed by an LL parser, which parses the input from left to right, and constructs a leftmost derivation. The 1 means that we use one input symbol of lookahead to make parsing action decisions. The predictive parser makes use of a parse table and a stack to process the input token stream.

The parse table is a two dimensional array  $M[X, a]$  where  $X$  is a non terminal and  $a$  is a terminal. This parse table tells the parser which production rule to use when we have a non-terminal  $X$  on top of the stack and the lookahead pointer points to  $a$ . This parse table is generated by using the First and Follow set of every non-terminal in the grammar.

To construct the parse table, we do the following steps for each production rule  $A \rightarrow \alpha$ :

- (1) For each terminal  $a$  in  $\text{First}(\alpha)$ ,  $M[A, a] = A \rightarrow \alpha$
- (2) If  $\epsilon$  is in  $\text{First}(\alpha)$ ,  $M[A, b] = A \rightarrow \alpha$  for each terminal  $b$  in  $\text{Follow}(A)$ .
- (3) If  $\epsilon$  is in  $\text{First}(\alpha)$ , and  $\$$  is in  $\text{Follow}(A)$ ,  $M[A, \$] = A \rightarrow \alpha$ .

Now we can finally see the algorithm for predictive parsing. Consider that  $\$$  is a special token at the bottom of the stack and also terminates the input string. Assume that initially, the predictive parser has  $X$  symbol on top of the stack and  $a$  is the current input symbol. Then, the predictive parser does the following:

- If  $X = a = \$$ , then stop.
- If  $X = a \neq \$$ , then pop  $X$  and increment the lookahead pointer.
- If  $X$  is a non terminal, then
  - If  $M[X, a] = X \rightarrow PQR$ , then pop  $X$  and push  $R, Q, P$ .
  - Else, throw an error

Predictive parsing needs a LL(1) Grammar. A grammar is LL(1) if the constructed parse table has no multiple entries. If it does have multiple entries of in the same cell, it cannot be LL(1), and cannot be parsed using predictive parsing.

Another way to tell if a grammar is LL(1) is to consider rules of the type:

$$A \rightarrow \alpha_1 | \alpha_2 | \cdots | \alpha_n$$

All rules of this type must be such that  $\bigcap_{i=1}^n First(\alpha_i) = \phi$  and rules of the type:

$$A \rightarrow \alpha | \epsilon$$

must be such that  $First(\alpha) \cap Follow(A) = \phi$ .

## 8. Error Recovery

The compiler must recover from errors and identify as many errors as possible. For this, the most frequently used error technique is **panic mode**.

In panic mode, when an error is encountered anywhere in the statement, the rest of the statement is ignored by not processing the input from erroneous input to delimiters. This mode prevents the parser from developing infinite loops and is considered as the easiest way for recovery of the errors.

The error is detected when an entry in the parse table is empty. We then skip over symbols in the input until a token in a selected set of **synchronizing tokens** appears. If the error occurs for the parse table entry  $M[A, a]$ , we place the symbols in  $Follow(A)$  in a synchronizing set in the parse table. So, we skip tokens until an element in the synchronizing set appears, then pop  $A$  and continue parsing from that token.

## 9. Bottom Up Parsing and Shift-Reduce Parsing

In bottom up parsing, we design a parse tree for an input string starting from the leaf nodes and going towards the root. This is equivalent to finding a series of steps to reduce a string  $w$  of input to the start symbol of the grammar by tracing out the rightmost derivations of  $w$  in reverse.

Bottom up parsing is done via **shift reduce parsing**. Shift reduce parsing splits the string into two parts, separated by a special character ‘.’. The left part is a string of terminals and non terminals (on the stack), and the right part is a string of terminals (rest of the input string). Initially, we assume the input to be  $.w$ . The parser does one of two operations:

- **Shift:** It moves terminal symbol from the right part of the string to the left part of the string. If string before the shift is  $\alpha.pqr$ , then the string after shift is  $\alpha p.qr$ . This is equivalent to pushing a terminal onto the stack.
- **Reduce:** It occurs immediately on the left of ‘.’ and identifies a string name as RHS of a production and replaces it by the LHS. If string before reduce action is  $\alpha\beta.pqr$  and  $A \rightarrow \beta$  is a production then string after reduction is  $\alpha A.pqr$ . This is equivalent to popping the RHS of the production, and pushing the LHS onto the stack.

Bottom up parsing is capable of handling left recursive grammars.

A string that matches the RHS of a production and whose replacement gives a step in the reverse of right most derivation is called a **handle**. Our goal is to always reduce the handle and not just any RHS. Hence the shift reduce parser must be capable of detecting these handles.

A shift reduce parser could also find a **conflict**. A shift-reduce conflict is when both a shift and a reduce operation are valid. A reduce-reduce conflict is when reduction can be done by more than one production rule.

## 10. The LR(0) Parser

Now let us look at our first bottom up parser, the LR(0) parser. LR parsers are those that read from left to right, producing a rightmost derivation in reverse. LR(0) uses 0 tokens of lookahead to make it's predictions. Before we discuss the exact algorithm we must understand some concepts.

**10.1. Augmentation of Grammar.** If  $G$  is a grammar with a start symbol  $S$ , the augmented grammar  $G'$  which has a new start symbol  $S'$  and a additional production:

$$S' \rightarrow S$$

When the parser reduces by this new rule, it will stop immediately in the accept state.

**10.2. LR(0) Items.** An LR(0) item of a grammar is a production of  $G$  with a special symbol '.' at some position on the RHS. So, the production  $A \rightarrow XYZ$  gives four LR(0) items:

$$A \rightarrow .XYZ$$

$$A \rightarrow X.YZ$$

$$A \rightarrow XY.Z$$

$$A \rightarrow XYZ.$$

Each item indicates how much of a production has been seen at a point in the process of parsing. For instance, the second rule indicates that we have seen an input string derivable from  $X$  and hope to see a string derivable from  $YZ$  on the next input.

**10.3. Closure.** Let  $I$  be a set of items for a grammar  $G$ . Then the  $\text{Closure}(I)$  is a set constructed as follows:

- Every item  $I$  is in  $\text{closure}(I)$
- If  $A \rightarrow \alpha.B\beta$  is in  $\text{closure}(I)$  and  $B \rightarrow \gamma$  is a production then  $B \rightarrow .\gamma$  is in  $\text{closure}(I)$ .

**10.4. Goto.** The Goto operation, defined by  $\text{Goto}(I, X)$ , where  $I$  is a set of items and  $X$  is a grammar symbol, is closure of set of items  $A \rightarrow \alpha X \beta$  such that  $A \rightarrow \alpha X \beta$  is in  $I$ . If  $I$  is a set of items for some valid prefix  $\alpha$ , then  $\text{Goto}(I, X)$  is set of valid items for prefix  $\alpha X$ .

**10.5. LR(0) Automaton.** The goal of the LR(0) parser is to create a LR(0) automaton, also called the Goto graph. The states of the automaton are the sets of items of the LR(0) collection, and the transitions are given by the Goto function. The start state of this automaton is  $\text{Closure}(\{S' \rightarrow .S\})$ . All the states are accepting states.

This automaton helps us with shift-reduce decisions. Suppose that the string  $\gamma$  of grammar symbols takes the LR(0) automaton from the start state 0 to some state  $j$ . Then, shift on the next input symbol  $a$  if state  $j$  has a transition on  $a$ . Otherwise, we reduce; the items in state  $j$  tell us which production to use.

**10.6. Parsing Table.** This LR(0) automaton is encoded in the **parsing table**. It consists of two parts - a parsing action function “Action” and a goto function “Goto”.

The Action function takes as arguments a state  $i$  and a terminal  $a$  (or \$, the input end marker). The value of  $\text{Action}(i, a)$  can have one of four forms:

- (1) Shift  $j$ , where  $j$  is a state. The action taken by the parser effectively shifts input  $a$  onto the stack, but uses state  $j$  to represent  $a$ . This is denoted by  $sj$ .
- (2) Reduce  $A \rightarrow \beta$ . The action of the parser effectively reduces  $\beta$  on the top of the stack to head  $A$ . This is denoted by  $rj$ , where  $j$  is the number of the production.
- (3) Accept the input and finish parsing. This is denoted by  $\text{acc}$ .
- (4) Error. This is denoted by a blank cell.

The Goto function maps a state  $i$  and a non terminal  $A$  to state  $j$  if  $\text{Goto}(I_i, A) = I_j$

**10.7. Is this Grammar LR(0)?** A grammar is LR(0) if its LR(0) parsing table does not contain multiple defined entries. Another way is to check if it has any shift-reduce/reduce-reduce conflicts.

**10.8. Contents of LR Parser.** The configuration of a LR parser is defined by the 2-tuple (Stack Contents, Remaining Input). Hence, the initial configuration would be  $(S_0, a_0 a_1 a_2 \dots a_n \$)$

**10.9. Example.** Let us consider the grammar

$$S \rightarrow AA$$

$$A \rightarrow aA|b$$

And parse the input *aabb*. The grammar will be augmented to become:

$$0 : S' \rightarrow S$$

$$1 : S \rightarrow AA$$

$$2 : A \rightarrow aA$$

$$3 : A \rightarrow b$$

Now we can generate the Goto graph. This is done by first starting with the start state  $S' \rightarrow .S$ . This state is populated with other LR(0) items which exist in its closure. Then, from each state, we make transitions over terminals (whichever are possible), where each transition is a goto (or a shift).

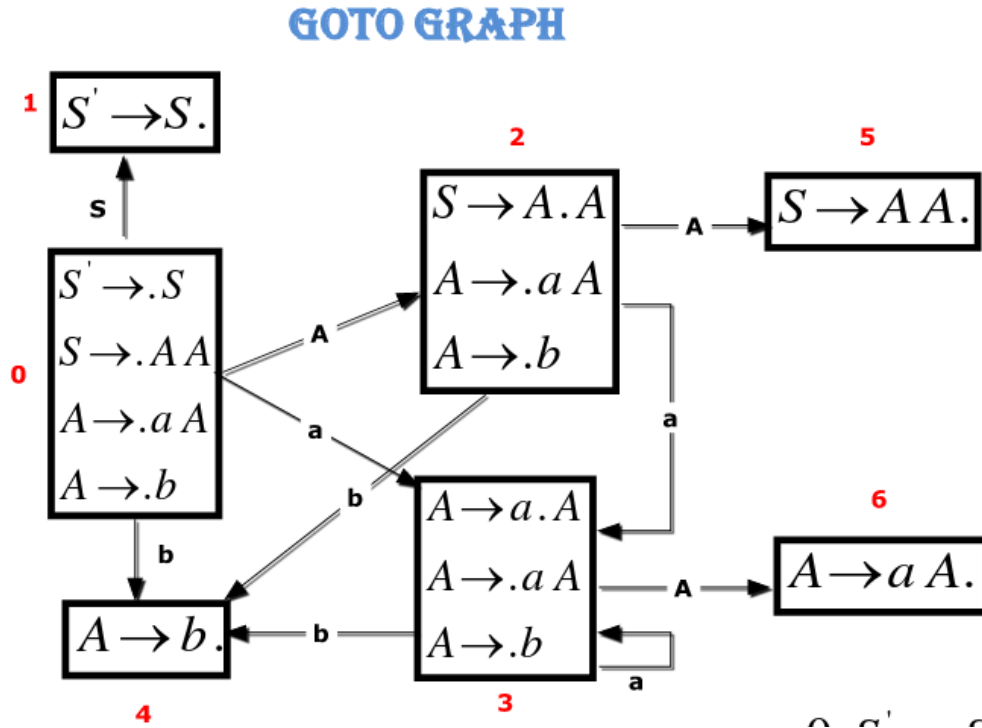


FIGURE 1. Goto graph

Each edge in the automaton is the goto, and every state is the closure. Using this graph, we can encode it as a parsing table (Fig 2). The goto operations over terminals are equivalent to shift operations. If a state has no outgoing goto operations, we must reduce by the item in that state, over every terminal.

**LR(0) PARSING TABLE**

	Action			Goto	
	a	b	\$	S	A
<b>0</b>	S3	S4		1	2
<b>1</b>			accept		
<b>2</b>	S3	S4			5
<b>3</b>	S3	S4			6
<b>4</b>	R3	R3	R3		
<b>5</b>	R1	R1	R1		
<b>6</b>	R2	R2	R2		

FIGURE 2. Parsing Table

The parsing operation of  $aabb$  can be seen in Table 1. In the 4th step, we see we pop out  $b4$  and push  $A6$ . How did we decide to move to state 6? This is because of the Goto operation. We checked the preceding state  $a3$  and check it's goto over non terminal  $A$ . This was a transition to state 6.

**10.10. Limitations of the LR(0) Parser.** The primary limitation of the LR(0) parser is that it does not take the lookaheads into account at all - it reduces indiscriminately, regardless of the input token. As such, it only works when states have a single reduce action!

**11. The SLR(1) Parser**

The SLR(1) parser is much like the LR(0) parser, except that the “reduce” operations are only written for terminals which are in Follow of the variable whose production is reduced. This way, the SLR parser uses 1 token of lookahead to make it's predictions, which is the reason that we include the 1.

**11.1. Constructing the SLR(1) Parsing Table.** The algorithm for constructing the parsing table is as follows:

- (1) Construct  $C = \{I_0, \dots, I_n\}$  the collection of sets of LR(0) items
- (2) If  $A \rightarrow \alpha.a\beta$  is in  $I_i$  and  $\text{goto}(I_i, a) = I_j$ , then  $\text{action}[i, a] = \text{shift } j$ .
- (3) If  $A \rightarrow \alpha.$  is in  $I_i$  then  $\text{action}[i, a] = \text{reduce } A \rightarrow \alpha$  for all  $a$  in  $\text{Follow}(A)$ .
- (4) If  $S' \rightarrow S.$  is in  $I_i$ , then  $\text{action}[i, \$] = \text{accept}$

TABLE 1. Parsing Process

Stack	Input	Action
0	<i>aabb</i> \$	Shift : Push <i>a</i> and state 3 onto the stack
0a3	<i>abb</i> \$	Shift : Push <i>a</i> and state 3 onto the stack
0a3a3	<i>bb</i> \$	Shift Push <i>b</i> and state 4 onto the stack
0a3a3b4	<i>b</i> \$	Reduce by $A \rightarrow b$ . Pop 2 symbols from the stack and push <i>A</i> and state 6 onto the stack
0a3a3A6	<i>b</i> \$	Reduce by $A \rightarrow aA$ . Pop 4 symbols from the stack and push <i>A</i> and state 6 onto the stack
0a3A6	<i>b</i> \$	Reduce by $A \rightarrow aA$ . Pop 4 symbols from the stack and push <i>A</i> and state 2 onto the stack.
0A2	<i>b</i> \$	Shift : Push <i>b</i> and state 4 on stack.
0A2b4	\$	Reduce by $A \rightarrow b$ . Pop 2 symbols from the stack and push <i>A</i> and state 5 onto the stack
0A2A5	\$	Reduce by $S \rightarrow AA$ . Pop 4 symbols from the stack and push state 1 onto the stack.
0S1	\$	Accept.

(5) If  $\text{goto}(I_i, A) = I_j$  then  $\text{goto}[i, A] = j$  for all non terminals *A*

(6) All undefined entries are errors

As we can see, the only difference between this and the LR(0) parser is that the parser will not only reduce in a given state with  $A \rightarrow \alpha..$  Instead, it will make smarter decisions based on the follow set.

**11.2. Is it SLR(1)?** As long as there are no conflicts in the parsing table, the grammar is *SLR*(1).

**11.3. Limitations of the SLR(1) Parser.** The SLR(1) parser is better, but also does not take into account all the information available. Assume we have a configuration such as  $X \rightarrow u..$ , we know that this corresponds to having *u* as a handle at the top of the stack that we can reduce. We do this whenever the input symbol is in  $\text{Follow}(X)$ . However this may not always be desirable, since elements below *u* in the stack may preclude *u* from being a handle for reduction. In other words, SLR only considers the top elements of the stack! We may need to divide the states into parts so that we could consider the different ways in which certain reductions can be carried out.



## 12. The Canonical LR(1) Parser

**12.1. The LR(1) Item.** An LR(1) item has the form  $[I, t]$  where  $I$  is in LR(0) item and  $t$  is a lookahead token. As the dot moves through the right hand side of  $I$ , token  $t$  remains attached to it. When  $I$  is of the form  $A \rightarrow \alpha.\beta$  and  $\beta \neq \epsilon$ , the lookahead token has no meaning. However, the LR(1) item  $[A \rightarrow \alpha., t]$  calls for a reduce action if the lookahead is  $t$ .

**12.2. The Closure Operation.** Closure is done in a different way from LR(0) items. The algorithm is as follows:

---

**Algorithm 1: CLOSURE(I)**


---

```

while No more items are added to I do
  for each item  $[A \rightarrow \alpha.B\beta, a]$  in  $I$  do
    for each production  $B \rightarrow \gamma$  in  $G'$  do
      for each terminal  $b$  in  $First(\beta a)$  do
        | add  $[B \rightarrow .\gamma, b]$  to set  $I$ ;
      end
    end
  end
end
return  $I$ ;

```

---

**12.3. The Goto Operation.** For LR(1) items, the algorithm to find the Goto is:

---

**Algorithm 2: GOTO(I)**


---

```

initialize  $J$  to be the empty set;
for each item  $[A \rightarrow \alpha.X\beta, a]$  in  $I$  do
  | add item  $[A \rightarrow \alpha.X.\beta, a]$  to set  $J$ ;
end
return Closure( $J$ );

```

---

This is essentially the same as with LR(0) items, but explicitly stating to keep the lookahead token when performing the goto.

**12.4. The CLR(1) Automaton.** Just like before, we can create an automaton representing the operation of the parser. The states are the sets of items and the transitions are the goto over terminals and non-terminals. To get the LR(1) automaton, we can use the

algorithm 3. Notice it does exactly the same thing as in an LR(0) parser, but stated in an algorithmic form.

---

**Algorithm 3:** items( $G'$ )

---

Let  $C = CLOSURE(\{[S' \rightarrow .S, \$]\})$

```

while no new sets of items are added to  $C$  do
  for Each set of items  $I$  in  $C$  do
    for Each grammar symbol  $X$  do
      if  $GOTO(I, X)$  is not empty and not in  $C$  then
        | add  $GOTO(I, X)$  to  $C$ ;
      end
    end
  end
end

```

---

**12.5. The CLR(1) Parsing Table.** Just like before, we have to generate a parsing table to encode the automaton. It is generated as follows:

- (1) Construct  $C' = \{I_0, \dots, I_n\}$ , the collection of sets of LR(1) items for  $G'$ .
- (2) State  $i$  of the parser is constructed from  $I_i$ . The parsing action for state  $i$  is determined as follows:
  - (a) If  $[A \rightarrow \alpha.a\beta, b]$  is in  $I_i$  and  $GOTO(I_i, a) = I_j$ , then set  $Action[i, a]$  to shift  $j$ . Here  $a$  is a terminal.
  - (b) If  $[A \rightarrow \alpha., a]$  is in  $I_i$ ,  $A \neq S'$ , then set  $Action[i, a]$  to reduce  $A \rightarrow \alpha$ .
  - (c) If  $[S' \rightarrow S., \$]$  is in  $I_i$ , then set  $Action[i, \$]$  to accept.
- (3) The goto transitions for state  $i$  are constructed for all nonterminals  $A$  using the rule: If  $GOTO(I_i, A) = I_j$ , then  $GOTO[i, A] = j$ .
- (4) All other entries are error.
- (5) The initial state of the parser is the one constructed from the set of items containing  $[S' \rightarrow .S, \$]$

**12.6. Is a Language CLR(1)?** A language is CLR(1) if and only if the parsing table does not have repeated entries.

### 13. The LALR parser

**13.1. Why Do We Need It?** Generally, LALR parsers are used over CLR parsers. This is because CLR parsers would have thousands of states for a grammar that an LALR

parser could parse with hundreds of states. It is hence more economical to use these in practice. LALR reduces the number of states by merging similar sets. For instance, if we have two sets:

$$\begin{aligned} [A \rightarrow b., a/b] \\ [A \rightarrow b., \$] \end{aligned}$$

We could merge them to create

$$[A \rightarrow b., a/b/\$]$$

Those two sets are said to share a common **core**.

**13.2. Parse Table Construction.** The construction of a parse table for LALR works as follows:

- (1) Construct  $C = \{I_0, I_1, \dots, I_n\}$ , the collection of sets of LR(1) items
- (2) Find all the sets sharing the same cores, and merge them into their union.
- (3) Let  $C' = \{J_0, J_1, \dots, J_m\}$  be the resulting sets. The parsing actions for each of these states are the same as in CLR(1).
- (4) The Goto table is constructed as follows. If  $J$  is the union of one or more sets of LR(1) items, that is  $J = I_1 \cup I_2 \cup \dots \cup I_k$ , then the cores of  $\text{GOTO}(I_1, X)$ ,  $\text{GOTO}(I_2, X)$ , ...,  $\text{GOTO}(I_k, X)$  are the same, since  $I_1, I_2, \dots, I_k$  all have the same core. Let  $K$  be the union of all sets of items having the same core as  $\text{GOTO}(I_1, X)$ . Then  $\text{GOTO}(J, X) = K$ .

The meaning of (4) is that we can just directly merge the rows of the CLR(1) parsing tables of  $I_1, \dots, I_k$  to get the LALR parsing table. In case there is an overlap, that overlap would be a merged set present in  $C'$ .

**13.3. Relative Power.** The relative power of LALR(1) is actually less than CLR(1), since it merges some states that could have earlier contained some contextual information. When talking about parsers,

$$LR(0) \subset SLR \subset LALR(1) \subset CLR(1)$$

The LL(1) parser has no direct relation to them - however,  $LL(1) \subset LR(1)$ . This relation is in fact true for all  $k$  in LL( $k$ ) and LR( $k$ ).