

# CS445: Compiler Design

<b>Compiler Block Diagram</b>	<b>1</b>
<b>Lex, Flex, Yacc, and Bison</b>	<b>1</b>
<b>Definition of Grammar</b>	<b>1</b>
<b>Ambiguity in Grammars</b>	<b>2</b>
<b>Top-down Parsing</b>	<b>6</b>
<b>Recursive Descent Parsers</b>	<b>6</b>
<b>LL(1) Grammars and Their Parsers</b>	<b>9</b>
<b>Parsing Sentential Forms Using Parse Tables</b>	<b>12</b>
<b>Left Recursion Removal in Grammars</b>	<b>14</b>
<b>Left Factoring in Grammars</b>	<b>15</b>
<b>Computing First and Follow Sets of Grammars</b>	<b>16</b>
<b>Constructing Parse Tables for Grammars</b>	<b>20</b>
<b>LR(1) and Bottom-up Parsers</b>	<b>22</b>
<b>LR(1) Parsing</b>	<b>29</b>
<b>LALR(1) Parsing</b>	<b>31</b>

## Compiler Block Diagram

Source code → scanner → tokens → parser → parse tree → semantic analyzer → annotated tree → optimizer → intermediate representation (IR) → code generation → target language

## Lex, Flex, Yacc, and Bison

In yacc, what is the difference between a shift/reduce conflict and reduce/reduce conflict?

A reduce/reduce conflict occurs when we have a set of things that can be reduced to different nonterminals (for example  $x$  or  $y$ ). We don't know which terminal to reduce to. In a shift/reduce, we can choose to reduce everything on the RHS to the LHS now, or we can shift something else onto the symbol stack and move on to a later reduction. Dangling else is a shift/reduce conflict.

## Definition of Grammar

A given grammar generates the language (a programming language in our case). The language is a set of programs  $L = \{P_0, P_1, \dots, P_n\}$  that are generated by the grammar. A recognizer has to recognize each program as being a member of the language.

## Ambiguity in Grammars

Given a grammar where number contains all integers:

$\text{exp} \rightarrow \text{exp op exp} \mid (\text{exp}) \mid \text{number}$   
 $\text{op} \rightarrow + \mid - \mid *$

Show that a legal derivation exists for the sentential form “ $(34 - 3) * 42$ ” using the productions in the grammar. A sentential form is any valid member of the language generated by the grammar. Notice the grammar is ambiguous as we can substitute for either  $\text{exp}$  in the RHS  $\text{exp} \rightarrow \text{exp op exp}$ .

The legal derivation:

$\text{exp} \rightarrow \text{exp op exp}$   
 $\text{exp} \rightarrow \text{exp op number}$   
 $\text{exp} \rightarrow \text{exp} * \text{number}$   
 $\text{exp} \rightarrow (\text{exp}) * \text{number}$   
 $\text{exp} \rightarrow (\text{exp op exp}) * \text{number}$   
 $\text{exp} \rightarrow (\text{exp op number}) * \text{number}$   
 $\text{exp} \rightarrow (\text{exp} - \text{number}) * \text{number}$   
 $\text{exp} \rightarrow (\text{number} - \text{number}) * \text{number}$   
 $\text{exp} \rightarrow (\text{number} - \text{number}) * \text{number}$

The derivation shows “ $(34 - 3) * 42$ ” is a legal member of the language generated by the grammar. In this example, we wrote the rightmost derivation as we always choose to substitute the RHS of  $\text{exp}$ . We call this an LALR(1) derivation.

Now, show the leftmost derivation for “ $(34 - 3) * 42$ .”

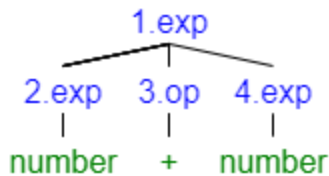
$\text{exp} \rightarrow \text{exp op exp}$   
 $\text{exp} \rightarrow (\text{exp}) \text{ op exp}$   
 $\text{exp} \rightarrow (\text{exp op exp}) \text{ op exp}$   
 $\text{exp} \rightarrow (\text{number op exp}) \text{ op exp}$   
 $\text{exp} \rightarrow (\text{number} - \text{exp}) \text{ op exp}$   
 $\text{exp} \rightarrow (\text{number} - \text{number}) \text{ op exp}$   
 $\text{exp} \rightarrow (\text{number} - \text{number}) * \text{exp}$   
 $\text{exp} \rightarrow (\text{number} - \text{number}) * \text{number}$

This leftmost derivation shows “ $(34 - 3) * 42$ ” is a legal member of the language generated by the grammar.

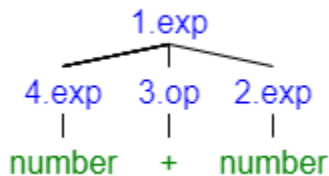
We can write a leftmost derivation for “ $4 + 7$ ” using the same grammar:

1.  $\text{exp} \rightarrow \text{exp op exp}$
2.  $\text{exp} \rightarrow \text{number op exp}$
3.  $\text{exp} \rightarrow \text{number} + \text{exp}$
4.  $\text{exp} \rightarrow \text{number} + \text{number}$

The corresponding parse tree for this derivation is:



We traverse the above tree in the order of the leftmost derivation, called preorder traversal. We can traverse the same tree in reverse postorder, as shown below. This corresponds to the rightmost derivation.



In compilers, we use abstract syntax trees (ASTs) to determine if the sentential form is a member of the grammar. Abstract syntax trees remove unnecessary terminals and improve efficiency. A parse tree has “everything” meaning all nonterminals that appear in the grammar must appear in the parse tree and the leaves in the tree must be terminals. In an AST, we remove the unnecessary nonterminals—those which don’t provide additional information. This saves memory and decreases tree traversal time.

We want to codify precedence in our op symbols, following the PEMDAS concept from math. The original grammar was:

$$\text{exp} \rightarrow \text{exp op exp} \mid (\text{exp}) \mid \text{number}$$

$$\text{op} \rightarrow + \mid - \mid *$$

We can convert it to give multiplication precedence over addition/subtraction:

$$\text{exp} \rightarrow \text{exp addop exp} \mid \text{term}$$

$$\text{addop} \rightarrow + \mid -$$

$$\text{term} \rightarrow \text{term mulop term} \mid \text{factor}$$

$$\text{mulop} \rightarrow *$$

$$\text{factor} \rightarrow (\text{exp}) \mid \text{number}$$

In this grammar, mulop is lower than addop and will appear lower in the parse tree. This is called a precedence cascade. We can further add associativity to the language. In this case, it is left associative because we write exp on the left in the first production. This grows the tree down to the right making it left recursive.

$$\text{exp} \rightarrow \text{exp addop term} \mid \text{term}$$

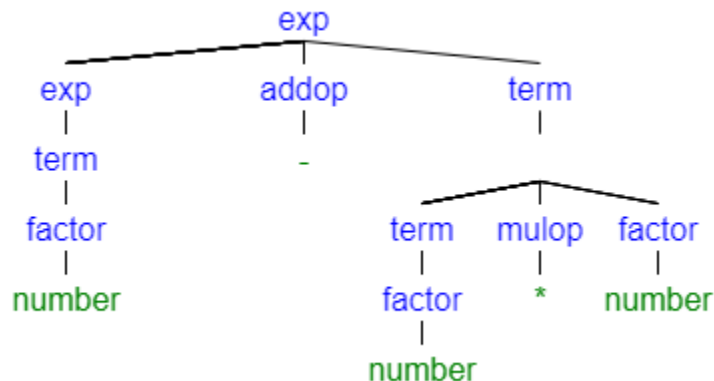
$$\text{addop} \rightarrow + \mid -$$

$$\text{term} \rightarrow \text{term mulop factor} \mid \text{factor}$$

$$\text{mulop} \rightarrow *$$

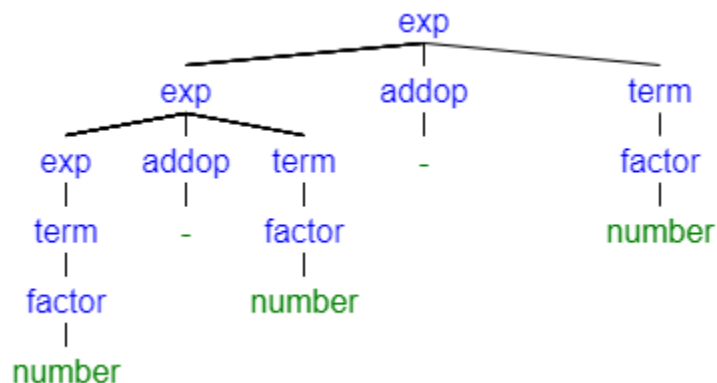
$$\text{factor} \rightarrow (\text{exp}) \mid \text{number}$$

We can construct a parse tree for the sentential form “34 - 3 \* 42.” Is it a valid sentential form?



Clearly, the sentential form is generated by the grammar. In order to resolve the second level of the tree (exp addop term), we must first go into the higher precedence (appearing lower in the grammar and parse tree) mulop.

Now, do the same for the sentential form “34 - 3 - 42.”



It is clear there is left associativity, meaning the tree “grows” to the left, meaning we need to resolve the left side before the right.

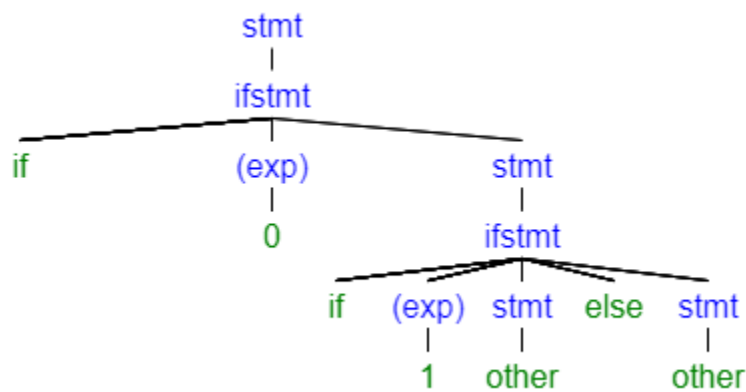
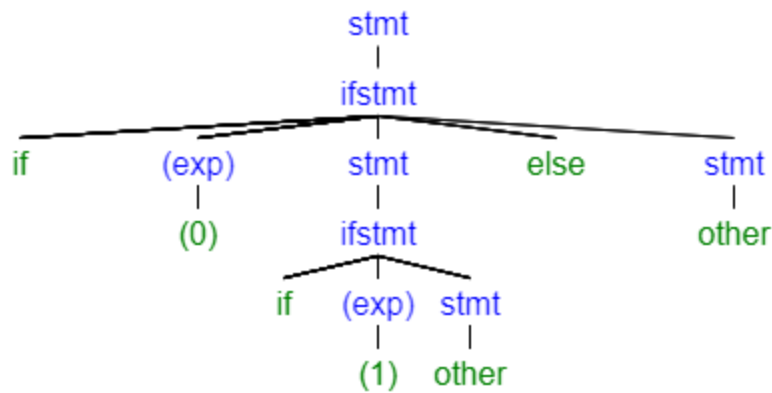
Here we have a grammar for simple if statements:

stmt  $\rightarrow$  ifstmt | other

ifstmt  $\rightarrow$  if (exp) stmt | if (exp) stmt else stmt

exp  $\rightarrow$  0 | 1

If we have the sentential form “if (0) if (1) other else other” it is ambiguous as we can derive two different parse trees that show the sentential form is a member of the language generated by the grammar. The two valid parse trees are:



We want to enforce a rule where the else is matched to the closest unmatched if. This solves the dangling else ambiguity.

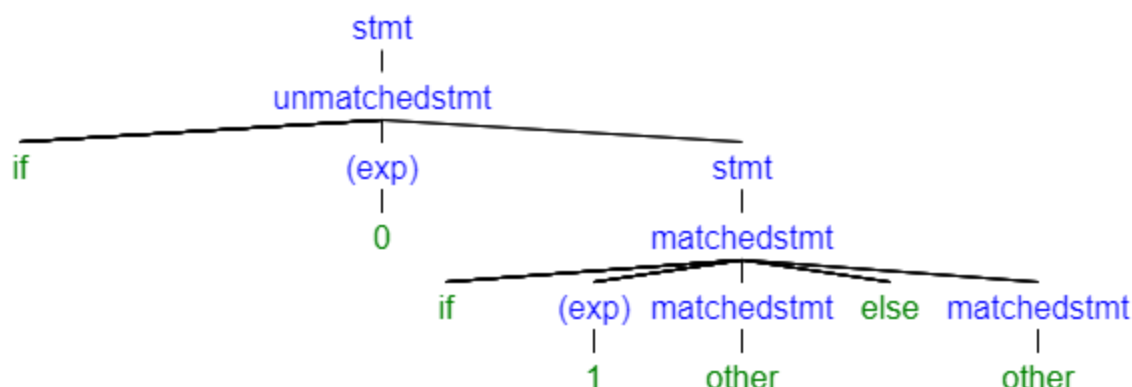
$\text{stmt} \rightarrow \text{matchedstmt} \mid \text{unmatchedstmt}$

$\text{matchedstmt} \rightarrow \text{if}(\text{exp}) \text{matchedstmt} \text{else} \text{matchedstmt} \mid \text{other}$

$\text{unmatchedstmt} \rightarrow \text{if}(\text{exp}) \text{stmt} \mid \text{if}(\text{exp}) \text{matchedstmt} \text{else} \text{unmatchedstmt}$

$\text{exp} \rightarrow 0 \mid 1$

We can prove "if (0) if (1) other else other" is a sentential form for the language generated by the new grammar.



There is no other valid parse tree that can be constructed for the sentential form. The grammar is not ambiguous. In programming languages, we often use `{}` to disambiguate the dangling else.

An unambiguous grammar for statements separated by semicolons is:

$\text{stmtseq} \rightarrow \text{stmt} ; \text{stmtseq} \mid \text{stmt}$

$\text{stmt} \rightarrow s$

The grammar is made ambiguous with the change:

$\text{stmtseq} \rightarrow \text{stmtseq} ; \text{stmtseq} \mid \text{stmt}$

$\text{stmt} \rightarrow s$

However, this is an unessential ambiguity as it does not have a meaningful impact on the language. The two grammars still generate the same things even if the second grammar is ambiguous.

## Top-down Parsing

In top-down parsing, we traverse from the root of the tree downward. We can divide top-down parsers into two classes: predictive and backtracking. In predictive top-down parsers, we consider the current state and input to predict the construct being parsed at that time. In backtracking, we are able to backtrack to a saved state if errors are encountered. We save states in a stack. Top-down parsers have an implicit preorder traversal. We will study two top-down parsers in the predictive class: recursive descent and LL(1).

Note that the first letter in LL(1) indicates how the input is considered. In LL(1), it is left-to-right. The second letter indicates whether we are doing a leftmost or rightmost derivation (called reduction). In LL(1), we do leftmost derivation. The number 1 indicates how many tokens we consider at most before making our prediction.

## Recursive Descent Parsers

Consider the below grammar.

$\text{exp} \rightarrow \text{exp addop term} \mid \text{term}$   
 $\text{addop} \rightarrow + \mid -$   
 $\text{term} \rightarrow \text{term mulop factor} \mid \text{factor}$   
 $\text{mulop} \rightarrow *$   
 $\text{factor} \rightarrow (\text{exp}) \mid \text{number}$

In a recursive descent parser, we consider a production as being a rule (or programmatic function) that specifies how to “write the code” to recognize the left hand side. For a factor, we may write the pseudocode:

```

procedure factor();
begin
    case token of
        '(':
            match('(');
            exp();
            match(')');
        number:
            match(number);
        else error;
    end case;
end factor;

```

If it is a factor, the token must start with '(' as in (exp) or number. We write pseudocode for match as follows:

```

procedure match(expTok);
begin[
    if (token = expTok) then
        getToken();
    else
        error();
    endif
end match;

```

In match, we have an error if the next input is not the correct token. Otherwise, it is correct, and we consume the token.

If we were to write a procedure for exp, it would call term which would call factor which would call exp (recursive).

$\text{ifstmt} \rightarrow \text{if (exp) stmt} \mid \text{if (exp) stmt else stmt}$

In the above simple grammar, both the productions have the same “prefix,” that is “if (exp) stmt.” If we see an “else,” we can assume we still need to parse the “if (exp) stmt else stmt” RHS. Otherwise, we assume we have successfully parsed to the first production. We can write pseudocode for this:

```

procedure ifstmt()
begin
    match(if);
    match('(');
    exp();
    match(')');
    stmt();
    if (token = else) then
        match(else);
        stmt();
    endif;
end ifstmt;

```

In the previously discussed grammar,

```

exp → exp addop term | term
addop → + | -
term → term mulop factor | factor
mulop → *
factor → (exp) | number

```

we encounter an issue when creating a recursive descent parser. We must define an `exp()` function to parse the `exp` variable on the RHS of the first production. However, the LHS of the production also contains the `exp` variable. Therefore, `exp` is left recursive in the production. Partially converting the BNF form to EBNF allows us to solve this “chicken or the egg” problem.

The RHS of the first production in the new EBNF production represents a term followed by 0 or more `addop` coupled with a term.

```
exp → term {addop term}
```

We then write the pseudocode for `exp()`.

```

procedure exp()
begin
    term();
    while (token = + or token = '-') do
        match(token);
        term();
    end while
end exp;

procedure term()

```



```

begin
    factor();
    while (token = '*') do
        match(token);
    end while
end term

```

How do we preserve the left associativity of + and -? We can rewrite exp().

```

procedure exp() : return integer
begin
    var tmp : integer;
    tmp := term();
    while(token = '+' or token = '-') do
        case token of
            '+':
                match('+');
                tmp := tmp + term();
            '-':
                match('-');
                tmp := tmp - term();
        end case;
    end while;
    return tmp;
end exp;

```

Let's say we have a grammar with many finite terminals or variables.

$$A \rightarrow a \mid b \mid \dots \mid z$$

We need to unambiguously determine which RHS our LHS is to be reduced to at a given point in parsing. This requires computations of "first sets." We must examine all possible terminals that may match the string in question: First(a), First(b), ..., First(z) where none are subsets of each other.

If these notes on first sets don't make sense to you it's because Wilder spent all of 3 minutes explaining the concept with the bulk of that being this quote:

"This requires computation of what's called the first sets of each token or of each nonterminal the first set of each nonterminal that we have to be able to be looking at to determine which right hand side which nonterminal thing it is that we're reducing the left hand to."

# LL(1) Grammars and Their Parsers

Consider the grammar:

$$S \rightarrow (S)S \mid \epsilon$$

It is of the form Dyck as the parentheses are matched. Recursive descent LL(1) parsing requires a stack. We use the stack instead of recursive calls to procedures. In our stack the “\$” character represents the bottom of the stack or end of input.

Is the sentential form “()” a member of the language generated by the grammar?

In LL(1) parsing, we either consume input or perform a reduction. When we consume input we call it a match like in the previously written pseudocode for a match() procedure.

Stack	Input	Action	Notes
\$S	()\$	$S \rightarrow (S)S$	Reduction; input doesn't change
\$S)S(	()\$	match “(“	We “lookahead” to the next symbol in input which is “(“ and match it.
\$S)S	)\$	$S \rightarrow \epsilon$	Reduce the S to $\epsilon$ as we can't match “)”
\$S)	)\$	match “)”	
\$S	\$	$S \rightarrow \epsilon$	We still have S in the stack; we can reduce it to $\epsilon$
\$	\$		We accept the sentential form through LL(1) parsing

Generalized, a top-down LL(1) parser looks like:

\$startsymbol    inputstring\$

.

. sequence of reductions and/or matches

.

\$                      \$                      string accepted

At any point when parsing, we can do one of two things:

1. Replace the nonterminal at the top of the symbol stack with a string (nonterminals or terminals) using the relevant production. This is called a reduce (also known as a generate).
2. Match a token on the top of the symbol stack with the next token of the input string.

Perform a leftmost reduction to prove the sentential form “()” is a member of the language generated by  $S \rightarrow (S)S \mid \epsilon$ .

$$\begin{array}{ll}
 S \rightarrow (S)S & [S \rightarrow (S)S] \\
 \rightarrow ()S & [S \rightarrow \epsilon] \\
 \rightarrow () & [S \rightarrow \epsilon]
 \end{array}$$



You can see the reductions in square brackets are the same as those in the “Action” column of the LL(1) table (ignoring match actions). This shows we performed a leftmost reduction in the LL(1) table.

We must make a choice when we have a nonterminal at the top of the parsing stack: we can choose any production. However, we must attempt a match if the top of the stack is a terminal. There is an error if we can’t do a match.

We can create a more formal lookup table for LL(1) parsing called the  $M[N, T]$  table. Note: M = machine, N = nonterminals, T = terminals. The table shows us which production to use for every nonterminal coupled with a terminal.

There are rules for the table:

- If  $A \rightarrow \alpha$  is a production in the grammar and there exists a derivation that goes from  $\alpha$  through some sequence of derivations to  $a\beta$  where  $a$  is a token (not a nonterminal), we add  $A \rightarrow \alpha$  to the table at entry  $M[A, a]$ .
- If  $A \rightarrow \alpha$  is a production in the grammar and there exists derivations where  $\alpha \rightarrow \epsilon$  and  $S\$ \rightarrow BA\alpha$  (see below note) where  $S$  is the start symbol of the grammar and  $a$  is a token (or  $\$$ ), we add  $A \rightarrow \alpha$  to the table at entry  $M[A, a]$ .

Note: he wrote this symbol here  he wrote another symbol here . I think it's a gamma but was nowhere else in the lecture, so I don't know what it represents.

In the LL(1) lookup table (called a LUT), each row represents a nonterminal and each column a terminal. Each entry specifies the production that should be used to perform a reduction for the nonterminal in the  $M[N, T]$  column and terminal in the left side of the input string.

For the previous example for the grammar  $S \rightarrow (S)S \mid \epsilon$  and sentential form “()”, the table appears as follows:

$M[N, T]$	(	)	\$
S	$S \rightarrow (S)S$	$S \rightarrow \epsilon$	$S \rightarrow \epsilon$

The definition of an LL(1) grammar from the textbook:

A grammar is  $LL(1)$  if the associated parse table has at most one production in every entry in the matrix. If we have two or more, it is not LL(1) as it would be ambiguous and the parsing cannot deterministically predict the correct production to use and we cannot backtrack.

We can write a pseudocode procedure for LL(1) parsing.

```

push start symbol onto parsing stack
while top of stack != $ do
    if top of stack is a terminal a and next token in input = a then
        pop a off parsing stack (a match has occurred)
        advance input to next token
    else if top of stack is a nonterminal A and next input is a token a or $ and parsing table
        entry M[A, a] contains production  $A \rightarrow x_1 x_2 \dots x_n$  then
        Pop parsing stack for  $i := n$  down to 1 do
            push  $x_i$  onto parsing stack
    else
        error
    endif

    if top of stack = $ and input = $
        accept the input
    else
        error
    endif
endwhile

```

## Parsing Sentential Forms Using Parse Tables

Another example of constructing an LL(1) parse table:

```

stmt  $\rightarrow$  ifstmt | other
ifstmt  $\rightarrow$  if (exp) stmt elsepart
elsepart  $\rightarrow$  else stmt |  $\epsilon$ 
exp  $\rightarrow$  0 | 1

```

M[N, T]	if	other	else	0	1	\$
stmt	stmt $\rightarrow$ ifstmt	stmt $\rightarrow$ other	error	error	error	error
ifstmt	ifstmt $\rightarrow$ if (exp) stmt elsepart	error	error	error	error	error
elsepart	error	error	elsepart $\rightarrow$ else stmt elsepart $\rightarrow$ $\epsilon$	error	error	elsepart $\rightarrow$ $\epsilon$

exp	error	error	error	exp $\rightarrow$ 0	exp $\rightarrow$ 1	error
-----	-------	-------	-------	---------------------	---------------------	-------

We have an ambiguity as there are two productions in an entry (the dangling else). So, we do not have a valid LL(1) parse tree.

An example of parsing the dangling else is shown in the sentential form “if (0) if (1) other else other.”

stack	input	action
\$stmt	if (0) if (1) other else other\$	stmt $\rightarrow$ ifstmt
\$ifstmt	if (0) if (1) other else other\$	ifstmt $\rightarrow$ if (exp) stmt elsepart
\$elsepart stmt )exp( if	if (0) if (1) other else other\$	match “if”
\$elsepart stmt )exp(	(0) if (1) other else other\$	match “(”
\$elsepart stmt )exp	0) if (1) other else other\$	exp $\rightarrow$ 0
\$elsepart stmt )0	0) if (1) other else other\$	match “0”
\$elsepart stmt )	) if (1) other else other\$	match “)”
\$elsepart stmt	if (1) other else other\$	stmt $\rightarrow$ ifstmt
\$elsepart ifstmt	if (1) other else other\$	ifstmt $\rightarrow$ if (exp) stmt elsepart
\$elsepart elsepart stmt )exp( if	if (1) other else other\$	match “if”
\$elsepart elsepart stmt )exp(	(1) other else other\$	match “(”
\$elsepart elsepart stmt )exp	1) other else other\$	exp $\rightarrow$ 1
\$elsepart elsepart stmt )1	1) other else other\$	match “1”
\$elsepart elsepart stmt )	) other else other\$	match “)”
\$elsepart elsepart stmt	other else other\$	stmt $\rightarrow$ other
\$elsepart elsepart other	other else other\$	match “other”
\$elsepart elsepart	else other\$	elsepart $\rightarrow$ else stmt (ambiguous; we could choose elsepart $\rightarrow \epsilon$ )

\$elsepart stmt else	else other\$	match “else”
\$elsepart stmt	other\$	stmt $\rightarrow$ other
\$elsepart other	other\$	match “other”
\$elsepart	\$	elsepart $\rightarrow \epsilon$
\$	\$	accept

Clearly, “if (0) if (1) other else other” is a member of the language generated by the grammar.

Not all languages that we want to use work with LL(1) grammars. However, we can modify some grammars to work with LL(1) parsing. One method is left recursion removal.

## Left Recursion Removal in Grammars

We want to remove left recursion to remove ambiguity and conform to LL(1) parsing requirements.

$\text{exp} \rightarrow \text{exp addop term} \mid \text{term}$

The above grammar is left recursive meaning the parse tree will grow to the left for a given sentential form. It is also left associative and the left recursion is immediate. If we switched “exp” with “term” in the RHS it would be right recursive and right associative.

We can remove the left recursion and maintain left associativity to allow LL(1) parsing. This gets more complicated when left recursion is not immediate.

For example, the below grammar requires a substitution before left recursion is apparent. It has general immediate left recursion.

$A \rightarrow Bb \mid \dots$

$B \rightarrow Aa \mid \dots$

Consider the following grammar with simple immediate left recursion.

$\text{exp} \rightarrow \text{exp addop term} \mid \text{term}$

It is of the form “ $A \rightarrow A\alpha \mid \beta$ ” and we can rewrite it as:

$A \rightarrow \beta A'$

$A' \rightarrow \alpha A' \mid \epsilon$

This form no longer has left recursion.

Here is what this looks like in the exp grammar:

$$\begin{aligned}\text{exp} &\rightarrow \text{term exp}' \\ \text{exp}' &\rightarrow \text{addop term exp}' \mid \varepsilon\end{aligned}$$

Below is a grammar with general immediate left recursion.

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_n \mid B_1 \mid B_2 \mid \dots \mid B_n$$

We can reduce it to the below grammar.

$$\begin{aligned}A &\rightarrow B_1A' \mid B_2A' \mid \dots \mid B_nA' \\ A' &\rightarrow \alpha_1A' \mid \alpha_2A' \mid \dots \mid \alpha_nA'\end{aligned}$$

We can do the same with the below grammar.

$$\begin{aligned}\text{exp} &\rightarrow \text{exp} + \text{term} \mid \text{exp term} \mid \text{term} \\ \text{exp}' &\rightarrow + \text{term exp}' \mid - \text{term exp}' \mid \varepsilon\end{aligned}$$

This change does not modify the language, just the grammar. Left associativity is maintained. We can now use LL(1) parsing for the grammar.

## Left Factoring in Grammars

Left factoring is required when we have two or more grammar choices that begin with the same prefix. Consider the general form of a grammar below.

$$A \rightarrow \alpha\beta \mid \alpha\rho$$

There is a common prefix “ $\alpha$ .” We need to correctly choose which production to use during parsing. We can fix this by factoring the grammar to:

$$\begin{aligned}A &\rightarrow \alpha A' \\ A' &\rightarrow \beta \mid \rho\end{aligned}$$

Now, the prefix is only present in one production. We must capture all of the commonality on the left for this to work.

A more complex example of a grammar with this issue is:

$$\begin{aligned}\text{stmtseq} &\rightarrow \text{stmt}; \text{stmtseq} \mid \text{stmt} \\ \text{stmt} &\rightarrow s\end{aligned}$$

Each production has the “s” prefix. We factor the grammar to remove the ambiguity.

$$\begin{aligned}\text{stmtseq} &\rightarrow \text{stmt stmtseq}' \\ \text{stmtseq}' &\rightarrow ; \text{stmtseq} \mid \epsilon\end{aligned}$$

Let's now attempt to factor the below grammar.

$$\text{ifstmt} \rightarrow \text{if (exp) stmt} \mid \text{if (exp) stmt else stmt}$$

Here we have the prefix “if (exp) stmt.” In a recursive descent parser each production is a “rule” or procedure. We need to perform left factoring to ensure the correct choice is made during LL(1) parsing. The left factored version of the grammar is:

$$\begin{aligned}\text{ifstmt} &\rightarrow \text{if (exp) stmt elsepart} \\ \text{elsepart} &\rightarrow \text{else stmt} \mid \epsilon\end{aligned}$$

Below is another case where the grammar does not comply with the requirements of LL(1) parsing.

$$\begin{aligned}\text{stmt} &\rightarrow \text{assignstmt} \mid \text{callstmt} \mid \text{other} \\ \text{assignstmt} &\rightarrow \text{identifier} := \text{exp} \\ \text{callstmt} &\rightarrow \text{identifier}(\text{explist})\end{aligned}$$

Here, identifier is a common prefix for an assign and call statement. Before left factoring, we first bring assignstmt and callstmt into a single production.

$$\text{stmt} \rightarrow \text{identifier} := \text{exp} \mid \text{identifier}(\text{explist}) \mid \text{other}$$

Then, we left factor.

$$\begin{aligned}\text{stmt} &\rightarrow \text{identifier stmt}' \mid \text{other} \\ \text{stmt}' &\rightarrow := \text{exp} \mid (\text{explist})\end{aligned}$$

The grammar now satisfies the requirements of LL(1) parsers.

## Computing First and Follow Sets of Grammars

LL(1) parsing falls under the category of predictive parsing. When doing predictive parsing, we cannot backtrack. This means we cannot make mistakes when determining which production is used for parsing at a given time. Given at most one token in the input string we must determine which production to use for reduction or matching. Computing the first sets for a grammar is helpful when doing so. The First(x) function is what we use to create the parse table for LL(1). Generally, we perform left factoring and left recursion removal before calculating the first sets.



First(x) is computed as follows:

1. If x is a terminal or  $\epsilon$  then  $\text{First}(x) = \{x\}$ .
2. If x is a nonterminal then for each production  $x \rightarrow x_1 x_2 \dots x_n$  (any combination of terminals and nonterminals),  $\text{First}(x)$  contains  $\text{First}(x_1) - \{\epsilon\}$ . If for some  $i < n$  all sets  $\text{First}(x_1) \dots \text{First}(x_i)$  contains  $\epsilon$ , then  $\text{First}(x)$  contains  $\text{First}(x_{i+1}) - \{\epsilon\}$ . If all sets  $\text{First}(x_1) \dots \text{First}(x_n)$  contains  $\epsilon$  then  $\text{First}(x)$  also contains  $\epsilon$ .

Define F(x) for a sequence of terminals and nonterminals  $\alpha = x_1 x_2 \dots x_n$  as:

$\text{First}(\alpha)$  contains  $\text{First}(x_1) - \{\epsilon\}$ . For each  $i = 2 \dots n$ , if  $\text{First}(x_k)$  contains  $\epsilon$  for all  $k = 1 \dots i - 1$ , then  $\text{First}(\alpha)$  contains  $\text{First}(x_i) - \{\epsilon\}$  and if for all  $i = 1 \dots n$ ,  $\text{First}(x_i)$  contains  $\epsilon$ , then  $\text{First}(\alpha)$  contains  $\epsilon$ .

We can write pseudocode for computing First(x).

```

for all nonterminals A do
    First(A) := {}
while any changes to any First(A) do
    for each production  $A \rightarrow x_1 x_2 \dots x_n$  do
        k := 1
        continue := true
        while (continue = true) and k <= n do
            add  $\text{First}(x_k) - \{\epsilon\}$  to First(A)
            if  $\epsilon$  is not a member of First(A) then
                continue := true
            k := k + 1
        if (continue = true) then
            add  $\{\epsilon\}$  to First(A)

```

Consider the below grammar.

```

exp → exp addop term | term
addop → + | -
term → term mulop factor | factor
mulop → *
factor → (exp) | number

```

When computing the first sets, we first break up the productions into individual lines:

```

exp → exp addop term
exp → term
addop → +
addop → -
term → term mulop factor
term → factor

```

$\text{mulop} \rightarrow *$   
 $\text{factor} \rightarrow (\text{exp})$   
 $\text{factor} \rightarrow \text{number}$

The goal is to have a set of possible RHS terminals and nonterminals for a given LHS. We have an error if there is an empty first set for the LHS provided.

The steps for computing the first sets for the above grammar are shown below. We begin with empty sets and continue to loop as long as there are changes to any first set for any nonterminal in the grammar. We cannot calculate the first set when there is recursion in the production or the first part of the RHS references an empty first set. When we reach the last column, no more changes to the first sets occur, and we are done computing the first sets.

$\text{First}(\text{exp}) = \{\}$	$\text{First}(\text{exp}) = \{\}$	$\text{First}(\text{exp}) = \{\}$	$\text{First}(\text{exp}) = \{ (, \text{number} \}$
$\text{First}(\text{addop}) = \{\}$	$\text{First}(\text{addop}) = \{ +, - \}$	$\text{First}(\text{addop}) = \{ +, - \}$	$\text{First}(\text{addop}) = \{ +, - \}$
$\text{First}(\text{term}) = \{\}$	$\text{First}(\text{term}) = \{\}$	$\text{First}(\text{term}) = \{ (, \text{number} \}$	$\text{First}(\text{term}) = \{ (, \text{number} \}$
$\text{First}(\text{mulop}) = \{\}$	$\text{First}(\text{mulop}) = \{ * \}$	$\text{First}(\text{mulop}) = \{ * \}$	$\text{First}(\text{mulop}) = \{ * \}$
$\text{First}(\text{factor}) = \{\}$	$\text{First}(\text{factor}) = \{ (, \text{number} \}$	$\text{First}(\text{factor}) = \{ (, \text{number} \}$	$\text{First}(\text{factor}) = \{ (, \text{number} \}$

A nullable nonterminal can be reduced to  $\epsilon$  through a series of replacements. In the below grammar, “elsepart” is nullable.

$\text{stmt} \rightarrow \text{ifstmt} \mid \text{other}$   
 $\text{ifstmt} \rightarrow \text{if}(\text{exp}) \text{stmt} \text{elsepart}$   
 $\text{elsepart} \rightarrow \text{else stmt} \mid \epsilon$   
 $\text{exp} \rightarrow 0 \mid 1$

We can transfer the above grammar into individual lines:

$\text{stmt} \rightarrow \text{ifstmt}$   
 $\text{stmt} \rightarrow \text{other}$   
 $\text{ifstmt} \rightarrow \text{if}(\text{exp}) \text{stmt} \text{elsepart}$   
 $\text{elsepart} \rightarrow \text{else stmt}$   
 $\text{elsepart} \rightarrow \epsilon$   
 $\text{exp} \rightarrow 0$   
 $\text{exp} \rightarrow 1$

$\text{First}(\text{stmt}) = \{\}$	$\text{First}(\text{stmt}) = \{\text{other}\}$	$\text{First}(\text{stmt}) = \{\text{other, if}\}$
$\text{First}(\text{ifstmt}) = \{\}$	$\text{First}(\text{ifstmt}) = \{\text{if}\}$	$\text{First}(\text{ifstmt}) = \{\text{if}\}$
$\text{First}(\text{elsepart}) = \{\}$	$\text{First}(\text{elsepart}) = \{\text{else, } \epsilon\}$	$\text{First}(\text{elsepart}) = \{\text{else, } \epsilon\}$
$\text{First}(\text{exp}) = \{\}$	$\text{First}(\text{exp}) = \{0, 1\}$	$\text{First}(\text{exp}) = \{0, 1\}$

Another grammar is:

$\text{stmtseq} \rightarrow \text{stmt stmtseq}'$   
 $\text{stmtseq}' \rightarrow ; \text{stmtseq} \mid \epsilon$   
 $\text{stmt} \rightarrow s$

First, we write each production in its own line

$\text{stmtseq} \rightarrow \text{stmt stmtseq}'$   
 $\text{stmtseq}' \rightarrow ; \text{stmtseq}$   
 $\text{stmtseq}' \rightarrow \epsilon$   
 $\text{stmt} \rightarrow s$

$\text{First}(\text{stmtseq}) = \{\}$	$\text{First}(\text{stmtseq}) = \{\}$	$\text{First}(\text{stmtseq}) = \{s\}$
$\text{First}(\text{stmtseq}') = \{\}$	$\text{First}(\text{stmtseq}') = \{;, \epsilon\}$	$\text{First}(\text{stmtseq}') = \{;, \epsilon\}$
$\text{First}(\text{stmt}) = \{\}$	$\text{First}(\text{stmt}) = \{s\}$	$\text{First}(\text{stmt}) = \{s\}$

Given a nonterminal  $A$ , the set  $\text{Follow}(A)$  is composed of terminals (and possibly  $\$$ ) and is defined as follows:

1. If  $A$  is the start symbol  $\$ \in \text{Follow}(A)$
2. If there exists a production  $B \rightarrow \alpha A \gamma$  then  $\text{First}(\gamma) - \{\epsilon\} \in \text{Follow}(A)$
3. If there exists a production  $B \rightarrow \alpha A \gamma$  such that  $\epsilon$  is in  $\text{First}(\gamma)$  then  $\text{Follow}(B) \in \text{Follow}(A)$

Note:  $\epsilon$  is never a member of a follow set, only first sets.

We can write pseudocode for computing the follow sets:

```

Follow(startsymbol) := {$}
for all nonterminals A != startsymbol do
    Follow(A) := {}
while any changes to any follow sets do
    for each production  $A \rightarrow x_1 x_2 \dots x_n$  do
        for each  $x_i$  that is a nonterminal do
            add  $\text{First}(x_{i+1} x_{i+2} \dots x_n) - \{\epsilon\}$  to  $\text{Follow}(x_i)$ 
            if  $\epsilon$  is in  $\text{First}(x_{i+1} x_{i+2} \dots x_n)$  then

```

add Follow(A) to Follow( $x_i$ )

Calculate the follow sets for the previously used example grammar and its first sets.

$\text{stmt} \rightarrow \text{ifstmt}$   
 $\text{stmt} \rightarrow \text{other}$   
 $\text{ifstmt} \rightarrow \text{if (exp) stmt elsepart}$   
 $\text{elsepart} \rightarrow \text{else stmt}$   
 $\text{elsepart} \rightarrow \epsilon$   
 $\text{exp} \rightarrow 0$   
 $\text{exp} \rightarrow 1$

$\text{First}(\text{stmt}) = \{\text{other}, \text{if}\}$   
 $\text{First}(\text{ifstmt}) = \{\text{if}\}$   
 $\text{First}(\text{elsepart}) = \{\text{else}, \epsilon\}$   
 $\text{First}(\text{exp}) = \{0, 1\}$

$\text{Follow}(\text{stmt}) = \{\$, \text{else}\}$	$\text{Follow}(\text{stmt}) = \{\$, \text{else}\}$
$\text{Follow}(\text{ifstmt}) = \{\$\}$	$\text{Follow}(\text{ifstmt}) = \{\$, \text{else}\}$
$\text{Follow}(\text{elsepart}) = \{\$\}$	$\text{Follow}(\text{elsepart}) = \{\$, \text{else}\}$
$\text{Follow}(\text{exp}) = \{\}$	$\text{Follow}(\text{exp}) = \{\}$

If the nonterminal is the last one in the RHS, use the Follow() of the LHS nonterminal.

Let's compute the follow sets for another grammar we previously looked at.

$\text{stmtseq} \rightarrow \text{stmt stmtseq}'$   
 $\text{stmtseq}' \rightarrow ; \text{stmtseq}$   
 $\text{stmtseq}' \rightarrow \epsilon$   
 $\text{stmt} \rightarrow s$

$\text{First}(\text{stmtseq}) = \{s\}$   
 $\text{First}(\text{stmtseq}') = \{;, \epsilon\}$   
 $\text{First}(\text{stmt}) = \{s\}$

$\text{Follow}(\text{stmtseq}) = \{\$\}$	$\text{Follow}(\text{stmtseq}) = \{\$\}$
$\text{Follow}(\text{stmtseq}') = \{\$\}$	$\text{Follow}(\text{stmtseq}') = \{\$\}$
$\text{Follow}(\text{stmt}) = \{;\}$	$\text{Follow}(\text{stmt}) = \{;, \$\}$

# Constructing Parse Tables for Grammars

Now that we can compute the first and follow sets, we can begin to construct parse tables. Use the instructions below to construct an LL(1) parse table for a given grammar.

For each nonterminal  $A$  and production  $A \rightarrow \alpha$  do:

1. For each token  $a$  in  $\text{First}(\alpha)$  add  $A \rightarrow \alpha$  to entry  $M[A, a]$ .
2. If  $\epsilon$  is in  $\text{First}(\alpha)$ , for each element  $a$  of  $\text{Follow}(A)$  (it must be a token or  $\$$ ; not  $\epsilon$ ) add  $A \rightarrow \alpha$  to  $M[A, a]$ .

A grammar in BNF is LL(1) if and only if:

1. For every production  $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ ,  $\text{First}(\alpha_i) \cap \text{First}(\alpha_j)$  is empty for all  $i$  and  $j$  such that  $1 \leq i$  and  $j \leq n$  and  $i \neq j$ .
2. For every nonterminal  $A$  such that  $\text{First}(A)$  contains  $\epsilon$ ,  $\text{First}(A) \cap \text{Follow}(A)$  is empty.

Construct a parse table for the grammar.

$\text{exp} \rightarrow \text{term exp}'$ $\text{exp}' \rightarrow \text{addop term exp}' \mid \epsilon$ $\text{addop} \rightarrow + \mid -$ $\text{term} \rightarrow \text{factor term}'$ $\text{term}' \rightarrow \text{mulop factor term}' \mid \epsilon$ $\text{mulop} \rightarrow *$ $\text{factor} \rightarrow (\text{exp}) \mid \text{number}$	$\text{First}(\text{exp}) = \{ (, \text{number} \}$ $\text{First}(\text{exp}') = \{ +, -, \epsilon \}$ $\text{First}(\text{addop}) = \{ +, - \}$ $\text{First}(\text{term}) = \{ (, \text{number} \}$ $\text{First}(\text{term}') = \{ *, \epsilon \}$ $\text{First}(\text{mulop}) = \{ * \}$ $\text{First}(\text{factor}) = \{ (, \text{number} \}$	$\text{Follow}(\text{exp}) = \{ \$, ) \}$ $\text{Follow}(\text{exp}') = \{ \$, ) \}$ $\text{Follow}(\text{addop}) = \{ (, \text{number} \}$ $\text{Follow}(\text{term}) = \{ \$, ), +, - \}$ $\text{Follow}(\text{term}') = \{ \$, ), +, - \}$ $\text{Follow}(\text{mulop}) = \{ (, \text{number} \}$ $\text{Follow}(\text{factor}) = \{ \$, ), +, -, * \}$
--	--	---

$M[N, T]$	(	number	)	+	-	*	\$
exp	$\text{exp} \rightarrow \text{term exp}'$	$\text{exp} \rightarrow \text{term exp}'$	error	error	error	error	error
exp'	error	error	$\text{exp}' \rightarrow \epsilon$	$\text{exp}' \rightarrow \text{addop term exp}'$	$\text{exp}' \rightarrow \text{addop term exp}'$	error	$\text{exp}' \rightarrow \epsilon$
addop	error	error	error	$\text{addop} \rightarrow +$	$\text{addop} \rightarrow -$	error	error
term	$\text{term} \rightarrow \text{factor term}'$	$\text{term} \rightarrow \text{factor term}'$	error	error	error	error	error
term'	error	error	$\text{term}' \rightarrow \epsilon$	$\text{term}' \rightarrow \epsilon$	$\text{term}' \rightarrow \epsilon$	$\text{term}' \rightarrow \text{mulop}$	$\text{term}' \rightarrow \epsilon$

						factor term'	
mulop	error	error	error	error	error	mulop $\rightarrow$ *	error
factor	factor $\rightarrow$ (exp)	factor $\rightarrow$ number	error	error	error	error	error

## LR(1) and Bottom-up Parsers

LL(1) parsing is inherently top-down because the second L indicates a leftmost reduction. LR parsing is a bottom-up parser. We don't necessarily need a number for the lookahead as in some cases we just need the parsing table and can consider the current state of the parsing stack without actually "looking ahead." We will examine LR(1) where the input is considered left to right, we do rightmost reductions, and our lookahead is 1. Later we will look at SLR(1) or simplified LR(1) parsing. In this parsing method, we simplify the parsing table.

Generally, an LR(1) parse is of the following form:

Parse Stack	Input	Action
\$	inputstring\$	
.		
.	Apply sequence of shift or reduce actions	
.		
\$startsymbol	\$	accept

We can perform one of two actions:

1. Shift a token from the inputstring onto the parse stack.
2. Reduce a string (of nonterminal or terminal) to a nonterminal A given a production  $A \rightarrow \alpha$ . In this action, we take  $\alpha$  off the parse stack and put A in its place.

Bottom-up parsers are typically called shift/reduce parsers. Shift/reduce conflicts occur when it is ambiguous which action should be performed at a given point in parsing.

In bottom-up parsers, we must perform "augmentation of the grammar." Consider the below grammar.

$$S \rightarrow (S)S \mid \epsilon$$

We must replace the original start rule. Augmentation is simply creating a new start production with the RHS being the LHS of the original start production.

$$S' \rightarrow S$$

$$S \rightarrow (S)S \mid \epsilon$$

Use LR(1) to determine if the sentential form “()” is a member of the grammar.

Stack	Input	Action
\$	()\$	shift “(“
\$(	)\$	reduce $S \rightarrow \epsilon$
\$(S	)\$	shift “)”
\$(S)	\$	reduce $S \rightarrow \epsilon$
\$(S)S	\$	reduce $S \rightarrow (S)S$
\$S	\$	reduce $S' \rightarrow S$
\$S'	\$	accept

An example with another grammar and sentential form “n + n.” We always use the longest possible production when reducing. We must reduce if the input is “\$” and there is no  $\epsilon$  in the productions.

$$E' \rightarrow E$$

$$E \rightarrow E+n \mid n$$

Stack	Input	Action
\$	n+n\$	shift “n“
\$n	+n\$	reduce $E \rightarrow n$
\$E	+n\$	shift “+”
\$E+	n\$	shift “n”
\$E+n	\$	reduce $E \rightarrow E+n$
\$E	\$	reduce $E' \rightarrow E$
\$E'	\$	accept

An example of creating a language to an NFA for LR(1) parsing:

$$S' \rightarrow S$$

$$S \rightarrow (S)S \mid \epsilon$$

Reduce the above language to “LR items” that we use to determine the states in an NFA for the LR(1) parser. The “.” tracks your progress while reading the input string that is the RHS of a production. We write only the “.” for  $\epsilon$  productions.

$$S' \rightarrow .S$$

$$S' \rightarrow S.$$

$$S \rightarrow .(S)S$$

$$S \rightarrow (.S)S$$

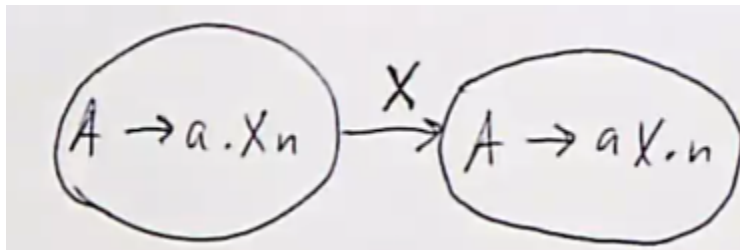
$$S \rightarrow (S.)S$$

$$S \rightarrow (S).S$$

$$S \rightarrow (S)S.$$

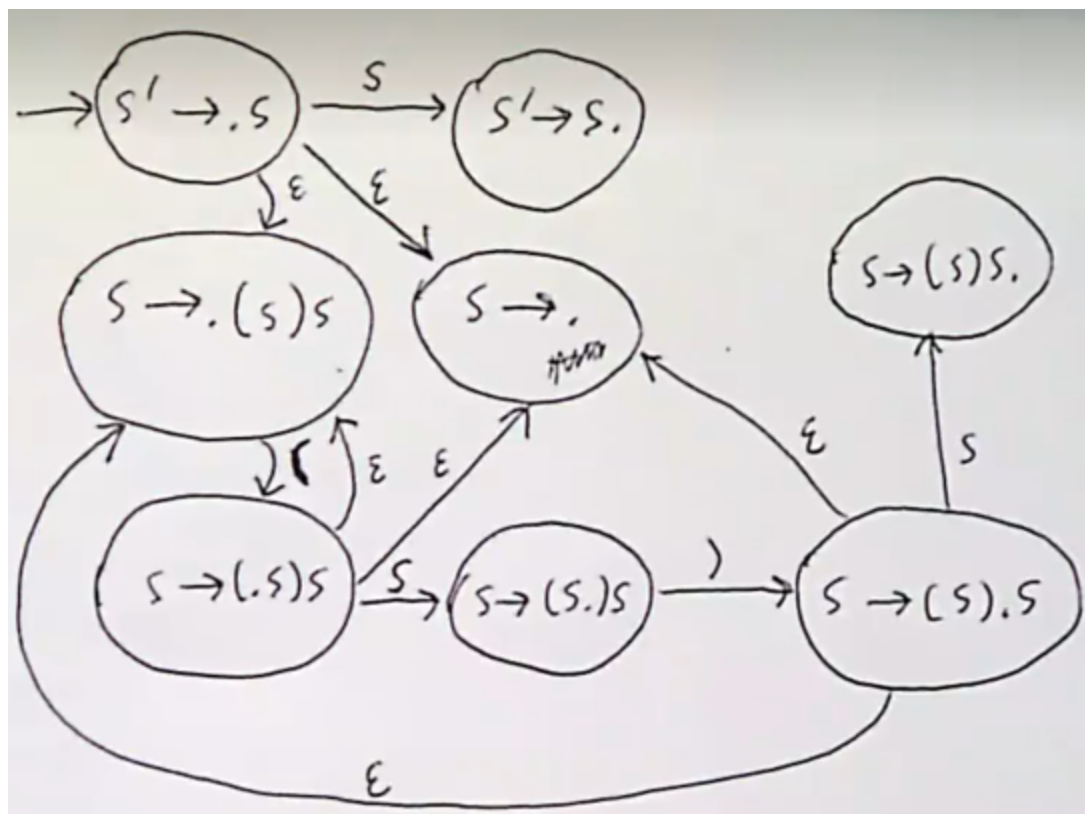
$$S \rightarrow .$$

An NFA for an LR(1) grammar follows the below format. We let the parser decide the accept states and do not draw them as part of the NFA. We use  $\epsilon$  transitions whenever a nonterminal can be reduced by a different production.

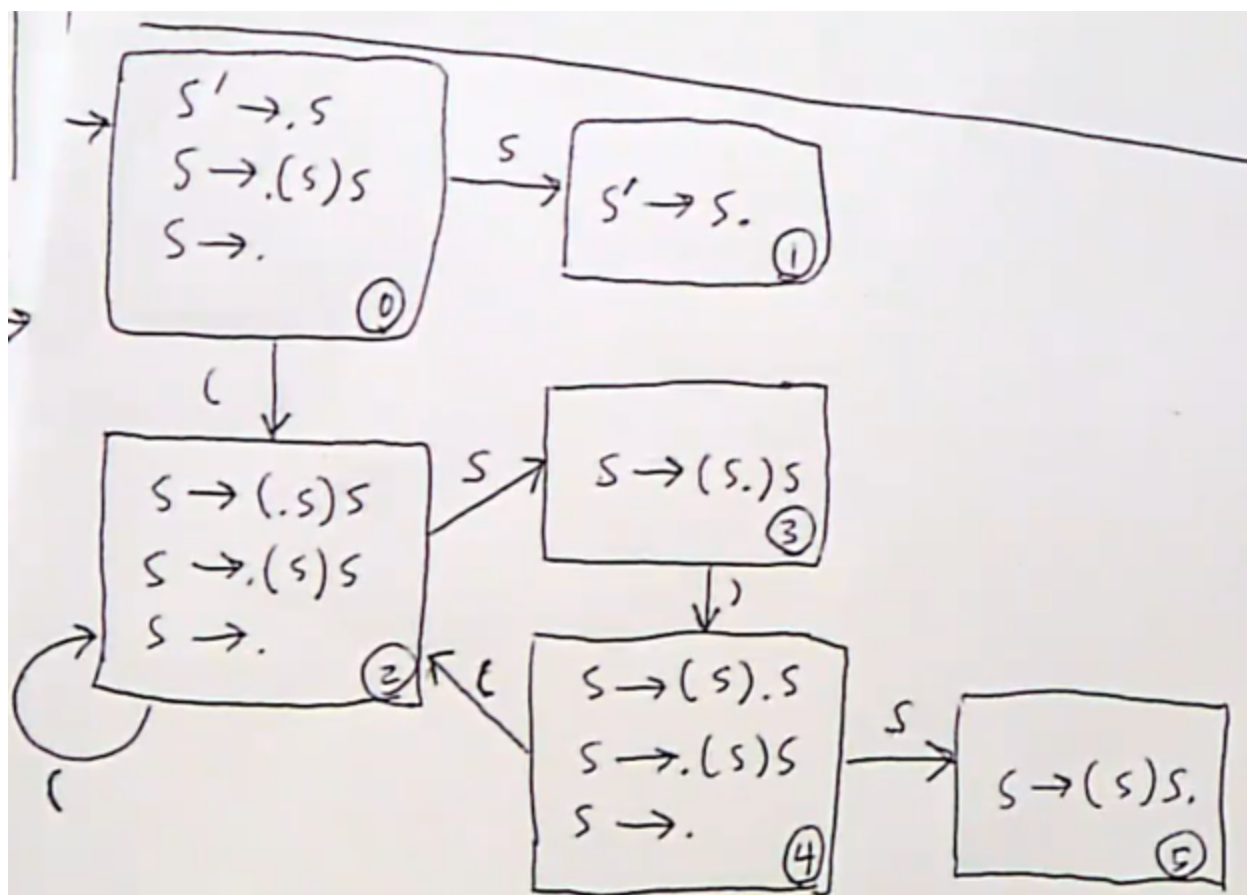


Write the NFA for the language.





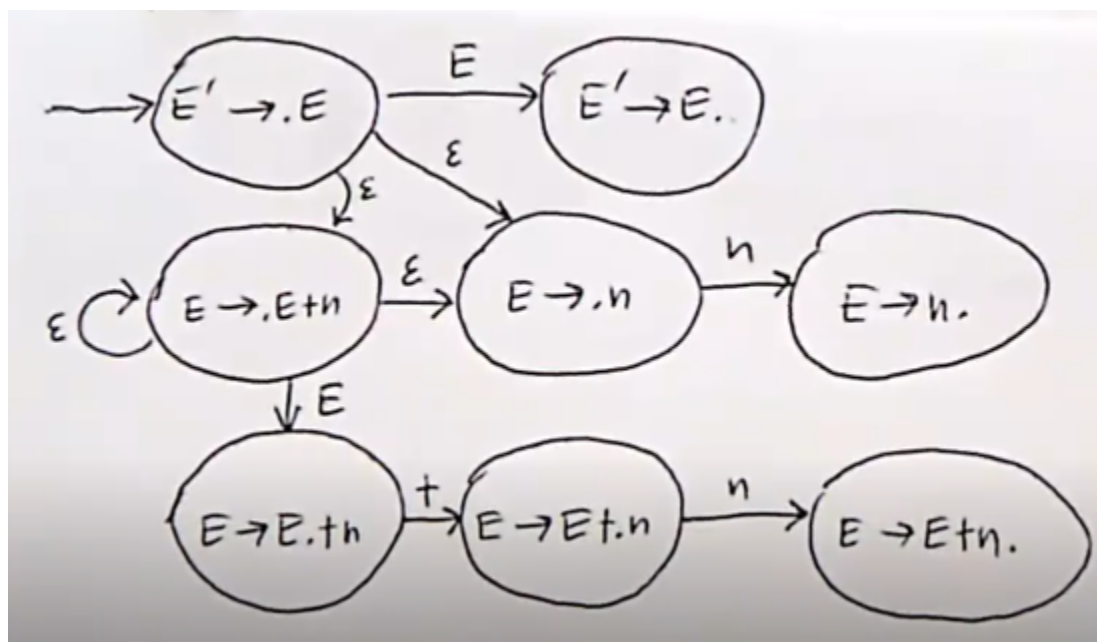
Convert it to a DFA. There are two types of items in each state: kernel items and closure items. Closure items are those productions that are reachable by a  $\epsilon$  transition in the NFA. We often only reference kernel items when discussing the DFA.



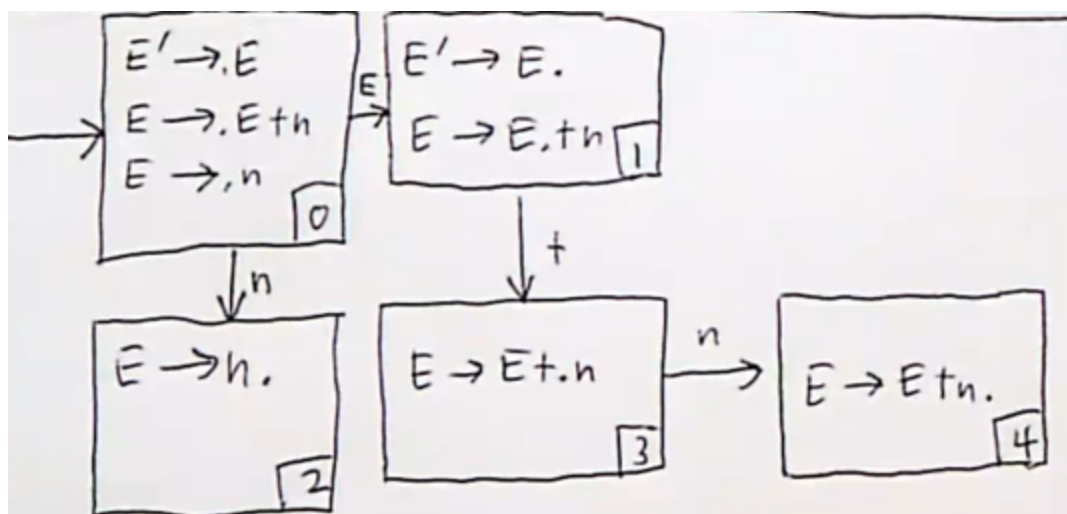
Write the NFA for the below grammar using LR(1) principles.

$E' \rightarrow E$   
 $E \rightarrow E+n \mid n$

$E' \rightarrow .E$   
 $E' \rightarrow E.$   
 $E \rightarrow .E+n$   
 $E \rightarrow E.+n$   
 $E \rightarrow E+.n$   
 $E \rightarrow E+n.$   
 $E \rightarrow .n$   
 $E \rightarrow n.$



Convert to a DFA.



Now, create the LR(1) parsing table. The augmentation is always in state 0. We start the stack with "\$0." For the sentential form "nrestofstring," the parse table looks as follows:

Stack	Input	Action
\$0	\$nrestofstring	shift "n"
\$0n2	\$restofstring	

In LR(0) parsing we do as follows:

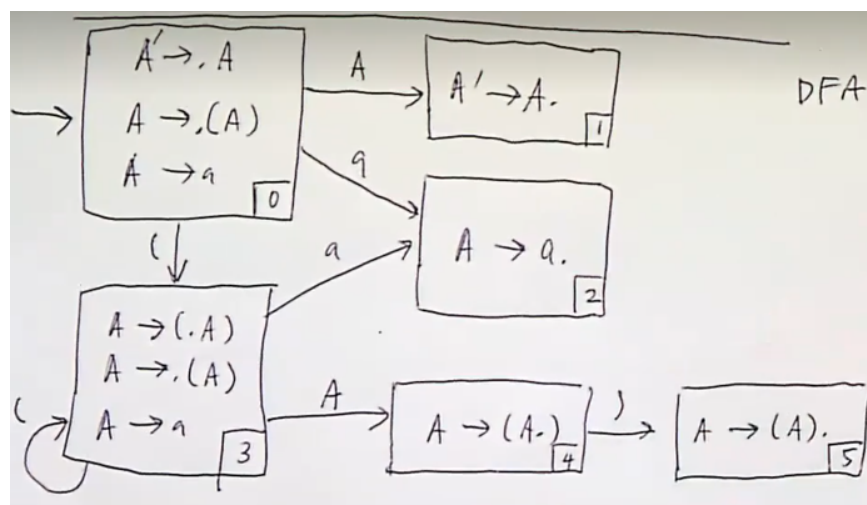
1. If state  $S$  contains an item of the form  $A \rightarrow \alpha.x\beta$  where  $x$  is a terminal, shift  $x$  onto the parse stack. If the terminal is indeed  $x$ , then the new state to be pushed is the state number containing  $A \rightarrow \alpha x.\beta$ . If it is not  $x$ , raise an error.
2. If state  $S$  contains any complete item (e.g.  $A \rightarrow \gamma$ ), then reduce by rule  $A \rightarrow \gamma$ . If we reduce by the rule  $S' \rightarrow S$ , then enter the accept state, provided the input is empty. Error if the input is not empty.

In all other cases, the new state of the parse is computed by:

- a. Remove string  $\gamma$  and all its states from the stack (backing up).
- b. Push  $A$  onto the stack.
- c. Push, as the new state, the number of the state corresponding to a production of the form  $B \rightarrow \alpha A.\beta$ .

Examine the simple grammar and its DFA.

$A \rightarrow (A) \mid a$



Write the parse stack for the sentential form " $((a))$ ." Note: we always want the top of the parse stack to contain the state number we are in.

Stack	Input	Action
\$0	$((a))\$$	shift "("
\$0(3	$(a))\$$	shift "("
\$0(3(3	$a))\$$	shift "a"
\$0(3(3a2	$)$$	reduce $A \rightarrow a$
\$0(3(3A4	$)$$	shift ")"
\$0(3(3A4)5	$)$$	reduce $A \rightarrow (A)$
\$0(3A4	$)$$	shift ")"
\$0(3A4)5	$\$$	reduce $A \rightarrow (A)$
\$0A1	$\$$	reduce $A' \rightarrow A$
\$0A'	$\$$	accept

Write an LR(0) lookup table for the grammar. The "Input" column contains all terminals while the "Goto" column contains all nonterminals.

State	Action	Rule	Input	Goto
			( a )	A

0	shift		3	2	error	1
1	reduce	$A' \rightarrow A$				
2	reduce	$A \rightarrow a$				
3	shift		3	2	error	4
4	shift		error	error	5	error
5	reduce	$A \rightarrow (A)$				

## LR(1) Parsing

“L” indicates we consider the input from left to right, “R” indicates we do a rightmost reduction making it bottom-up, and the “1” indicates we have a lookahead.

In bottom-up parsers we had “LR items,” now we have LR(1) items that are an LR item augmented with a lookahead.

An LR item looks like  $A \rightarrow \alpha.\beta$  where we are parsing through the RHS of the production, we have passed through  $\alpha$ , and we expect a  $\beta$ .

We can augment it to form an LR(1) item by adding a lookahead, in this case “a” (a is the next input in the token):

$[A \rightarrow \alpha.\beta, a]$

We can make DFAs and NFAs composed of LR(1) items like we did for LR items. The lookahead changes the transitions.

LR(1) transitions:

1. Given an LR(1) item  $[A \rightarrow \alpha.X\gamma, a]$  where X is any symbol (terminal or nonterminal), there is a transition on X to the item  $[A \rightarrow \alpha X.\gamma, a]$ . Note: the lookahead remains the same.
2. Given an LR(1) item  $[A \rightarrow \alpha.B\gamma, a]$  where B is a nonterminal, there are  $\epsilon$ -transitions to items  $[B \rightarrow \beta, b]$  for every production  $B \rightarrow \beta$  and for every token b in  $\text{First}(\gamma) \cup \text{First}(a)$  denoted by  $\text{First}(\gamma a)$ . “The transitions keep track of the context in which each structure B needs to appear because we capture the lookahead.”

Consider the grammar:

$A \rightarrow (A) \mid a$

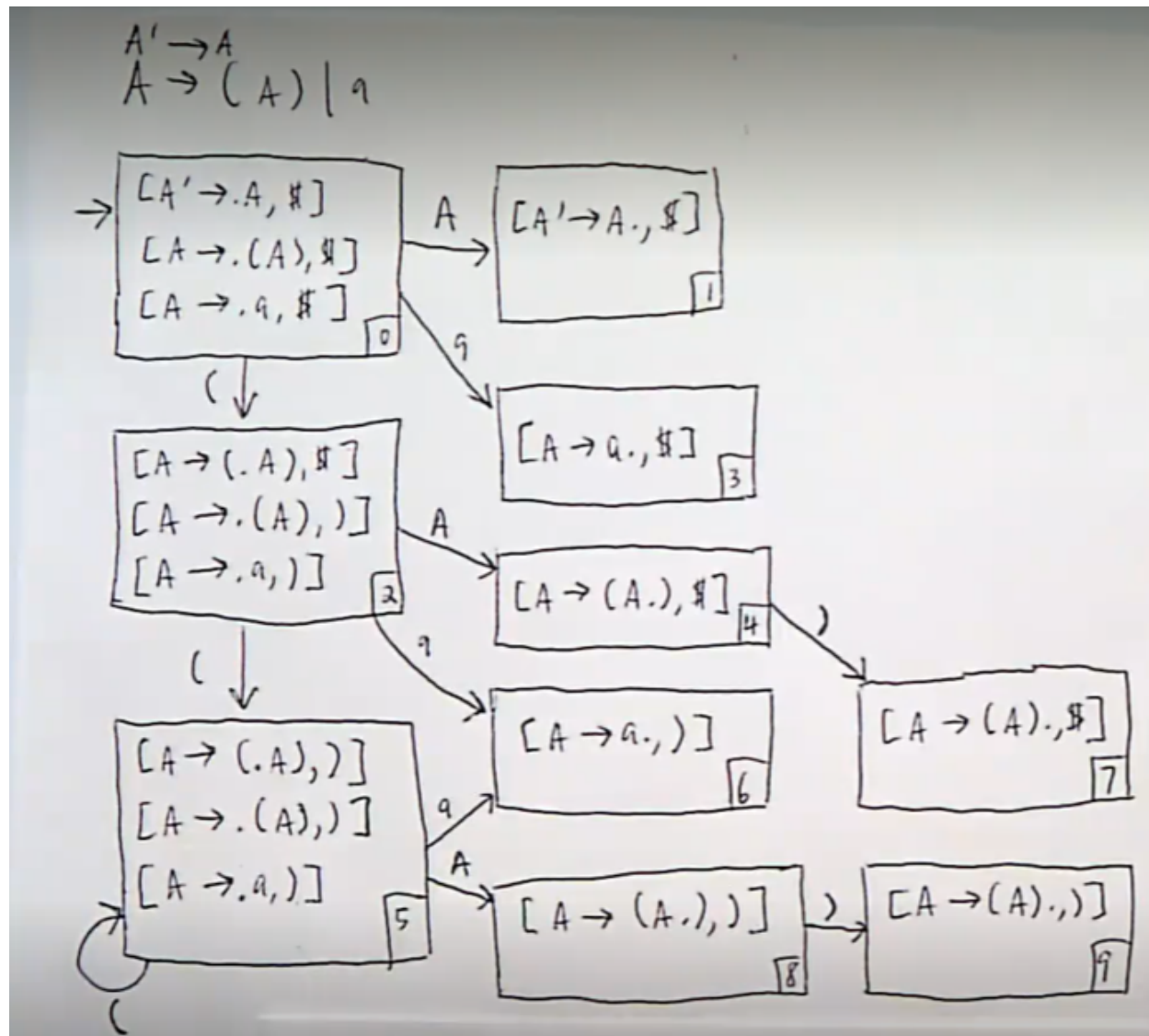
We augment the grammar.

$A' \rightarrow A$

$A \rightarrow (A) \mid a$

We create a DFA by putting the initial LR(1) items in the start state and then subsequent LR(1) items in the rest of the states.

The resulting DFA is below (we went straight to a DFA instead of making the LR items, then an NFA from them, then converting to DFA).



We can use the DFA to build a parse table and attempt to recognize a sentential form to see whether it is a member of the grammar generated by the language.

	Input				Goto
State	(	a	)	\$	A
0	shift; goto state 2	shift; goto state 3			1

1				accept	
2	shift; goto state 5	shift; goto state 6			4
3				reduce; goto state 2	
4			shift; goto state 7		
5	shift; goto state 5	shift; goto state 6			8
6			reduce; goto state 2		
7				reduce; goto state 1	
8			shift; goto state 9		
9			reduce; goto state 1		

Note: reductions are based on what the lookahead is. You should be able to determine whether a sentential form is a member of the language generated by the grammar, given a sentential form and parse table for the grammar.

## LALR(1) Parsing

You achieve LALR(1) in the same way as LR(1): they are LR(0) items with multiple lookaheads.

Given:

$A \rightarrow .B$

We have the LR item with three lookaheads:

$[A \rightarrow .B, \$]$

$[A \rightarrow .B, \#]$

$[A \rightarrow .B, f]$

We can rewrite this as:

$[A \rightarrow .B, \$/\#/f]$