

CS 315 – Programming Languages Parsing

CFG (context-free grammar)

- Parser uses CFG rules to perform a *derivation* and create a *parse tree*

For any CFG, we can construct a parser that runs in $O(n^3)$ time, where n is the input size (# of tokens).

We are interested in more scalable parsers, but these parsers cannot work with arbitrary CFGs

- LL parsers work with LL-grammars
 - LL: **L**eft-to-right input order, **L**eftmost derivation
 - These are *top-down* parsers
- LR parsers work with LR-grammars
 - LR: **L**eft-to-right input order, **R**ightmost derivation (in reverse)
 - These are *bottom-up* parsers

Top-down: the parse tree is constructed from top (root first) to bottom (leaves last)

Bottom-up: the parse tree is constructed from bottom (leaves first) to top (root last)

The class of LR grammars is larger than the class of LL grammars, i.e., more CFG grammars are in the former class.

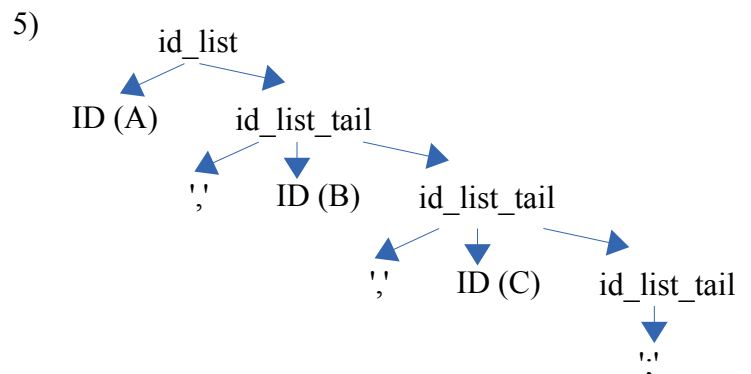
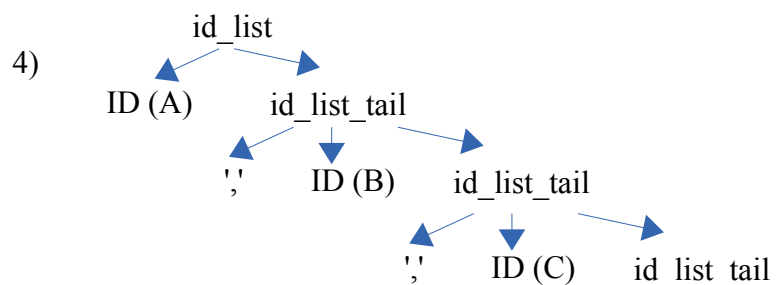
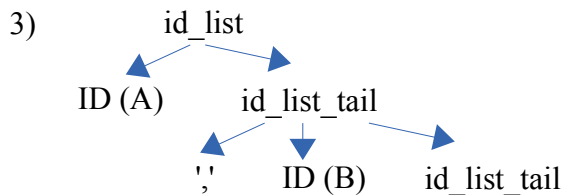
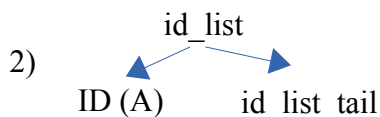
Let us see an example. Consider the following grammar, which is both an LL and an LR grammar:

1. $\text{id_list} \rightarrow \text{ID id_list_tail}$
2. $\text{id_list_tail} \rightarrow ', \text{ID id_list_tail}$
3. $\text{id_list_tail} \rightarrow ',$

Also consider the input: A, B, C;

Top-down parsing

1) id_list



Bottom-up parsing

Either push the next token from the input into the rightmost position in the queue
or

Reduce a set of tokens from the rightmost end of the queue using one of the rules

- | | | |
|----|---|----------|
| 0) | | shift |
| 1) | ID(A) | shift |
| 2) | ID(A) ',' | shift |
| 3) | ID(A) ',' ID(B) | shift |
| 4) | ID(A) ',' ID(B) ',' | shift |
| 5) | ID(A) ',' ID(B) ',' ID(C) | shift |
| 6) | ID(A) ',' ID(B) ',' <u>ID(C) ','</u> | reduce 3 |
| 7) | ID(A) ',' ID(B) <u>',' ID(C) id_list_tail</u> | reduce 2 |



- | | | |
|----|-------------------------------------|----------|
| 8) | ID(A) <u>',' ID(B) id_list_tail</u> | reduce 2 |
|----|-------------------------------------|----------|
-

- | | | |
|----|---------------------------|----------|
| 9) | ID(A) <u>id_list_tail</u> | reduce 1 |
|----|---------------------------|----------|
-

- | | |
|-----|--|
| 10) | |
|-----|--|

A few things to note:

- We have not learned how to programmatically create a parse tree, given an LL or LR grammar.
- For now, we are doing the construction using our intuition. Sometimes there are multiple choices, if you pick the wrong one, you cannot complete the parse tree creation.
- In the upcoming lectures we will learn programmatically creating a parse tree given an LL grammar. The same for LR grammars is left for a compilers course.
- You might have noticed that our bottom-up parse above resulted in pushing the entire input first, and then reducing it. This is undesirable and is a consequence of the right recursion in the grammar. As we will see, **converting it into left recursion will improve the bottom-up parse, but making the grammar a non-LL grammar.** *→ left recursion is good for LR parsing.*

LL-parser: Also called predictive parser

- It predicts what to see next in the input based on the current partial derivation and uses the appropriate rule to expand the leftmost nonterminal

→ look at the next char in the input, and use the CFG accordingly.

LR-parser: Also called the shift-reduce parser

- It pushes tokens into a queue until it sees a complete RHS of a production rule at the end of the queue, which is then reduced to the LHS of the rule. Sometimes there is a choice to be made between pushing and reducing (**pushing now may make it possible to reduce later**) *→ Example of this???*

→ shift until you see a fit that looks like one of the productions in your CFG.

LL-parser: from top to bottom, we get a leftmost derivation

LR-parser: From bottom to top, we get a rightmost derivation

There are different classes of LR parsers, such as (from more powerful to less powerful):

- LR(k)
- LALR(k)
- SLR

} This is how you waste my time

LL parsers can also be generalized as LL(k) parsers. Here k is the number of tokens in the input you can use for looking ahead before you make a decision.

Consider the following grammar for the problem of identifier lists we used before. The grammar is changed to use left recursion, which makes it a non-LL grammar, but it is still an LR grammar. Note that we have still not learned what makes a grammar an LL-grammar.

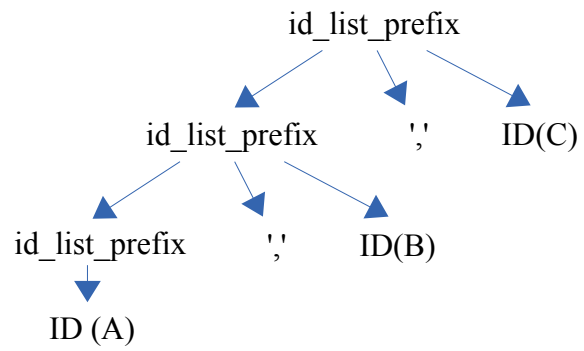
1. $\text{id_list} \rightarrow \text{id_list_prefix};$ *→ recalls the LHS!*
2. $\text{id_list_prefix} \rightarrow \text{id_list_prefix}; \text{ID}$
3. $\text{id_list_prefix} \rightarrow \text{ID}$

Let us do a bottom-up derivation for this grammar, again using the input:

A, B, C ;

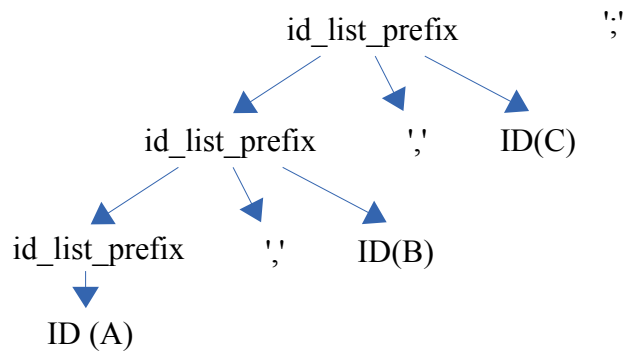
- 0) shift
- 1) ID(A) reduce 3
- 2) id_list_prefix
 ↓
 ID (A) shift
- 3) id_list_prefix ';' shift
 ↓
 ID (A)
- 4) id_list_prefix ';' ID (B) reduce 2
 ↓
 ID (A)
- 5) id_list_prefix shift
 / | \
id_list_prefix ';' ID(B)
 ↓
 ID (A)
- 6) id_list_prefix ';' shift
 / | \
id_list_prefix ';' ID(B)
 ↓
 ID (A)
- 7) id_list_prefix ';' ID(C) reduce 2
 / | \
id_list_prefix ';' ID(B)
 ↓
 ID (A)

8)



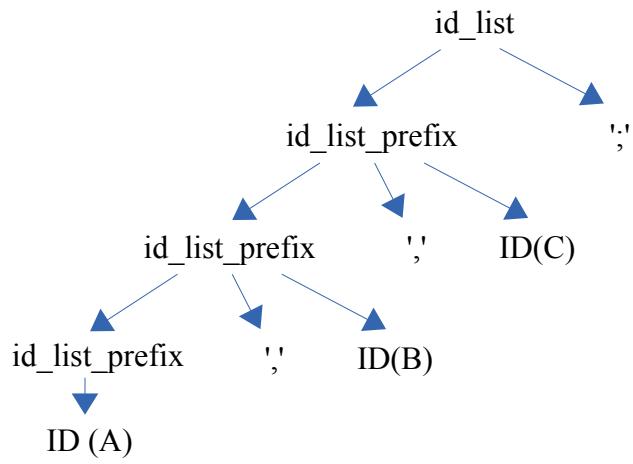
shift

9)



reduce 1

10)



LL Grammars

Consider our famous CFG for arithmetic expressions:

1. $\text{expr} \rightarrow \text{expr add_op term}$
2. $\text{expr} \rightarrow \text{term}$
3. $\text{term} \rightarrow \text{term mul_op factor}$
4. $\text{term} \rightarrow \text{factor}$
5. $\text{factor} \rightarrow \text{ID}$
6. $\text{factor} \rightarrow \text{NUM}$
7. $\text{factor} \rightarrow '(' \text{ expr } ')'$
8. $\text{add_op} \rightarrow '+'$
9. $\text{add_op} \rightarrow '-'$
10. $\text{mul_op} \rightarrow '*'$
11. $\text{mul_op} \rightarrow '/'$

This grammar creates problems for LL(1) parsing:

- Assume we have *expr* as our leftmost nonterminal in our current partial derivation and we have ID as our current input token. Which rule do we apply? Is it rule 1 or rule 2? Note that both rule 1 and rule 2 can expand to eventually produce an ID as the first thing.
- Left recursion: We can expand rule 1 as many times as we want, without impacting the next character from the input that can be consumed.

There are two conditions to ensure that a grammar is LL(1):

1. If two rules have the same LHS, then the set of input characters for which these rules are applicable (called the PREDICT set) must be **non-overlapping**. We will more formally define the PREDICT set a little later.
2. There can be no left recursion.

Now consider the following modified CFG, which is LL(1), for expressions:

1. $\text{stmt} \rightarrow \text{expr}\$$
2. $\text{expr} \rightarrow \text{term term_tail}$
3. $\text{term_tail} \rightarrow \text{add_op term term_tail}$
4. $\text{term_tail} \rightarrow \epsilon$
5. $\text{term} \rightarrow \text{factor factor_tail}$
6. $\text{factor_tail} \rightarrow \text{mul_op factor factor_tail}$
7. $\text{factor_tail} \rightarrow \epsilon$
8. $\text{factor} \rightarrow '(' \text{ expr } ')'$
9. $\text{factor} \rightarrow \text{ID}$
10. $\text{factor} \rightarrow \text{NUM}$
11. $\text{add_op} \rightarrow '+'$
12. $\text{add_op} \rightarrow '-'$
13. $\text{mul_op} \rightarrow '*'$
14. $\text{mul_op} \rightarrow '/'$

Consider *term_tail* with '+' in the input:

- R3 is possible, as *add_op* will eventually expand into '+'
- R4 is not possible, as only input characters it is applicable for are '\$' and ')'

- If we expand R4, we get ϵ , which means empty string
- So what input character can be consumed if we apply R4?
- It is the input character that can follow *term_tail* in the grammar
- We search the RHS of the rules to find *term_tail*
- We find it in two places:
 - R2, where *term_tail* appears as the last thing
 - So we cannot find, right away, what may come after it
 - But we know that whatever can come after *expr* (which is the LHS of R2) can come after *term_tail* as well
 - So we search what can come after *expr*
 - We search the RHS of the rules to find *expr*
 - We find it in two places
 - R1, where only '\$' can come after it
 - R8, where only ')' can come after it
 - R3, where the *term_tail* again appears as the last thing
 - However, LHS is *term_tail* as well, so no new information is learned

There are two ways to construct a top-down parser

- For each non-terminal (LHS), write a procedure

- See the book for complete examples. At this point you should be able to see why left recursion is a problem.

- Create a table, which, for each non-terminal/token pair, specifies the production rule to use.

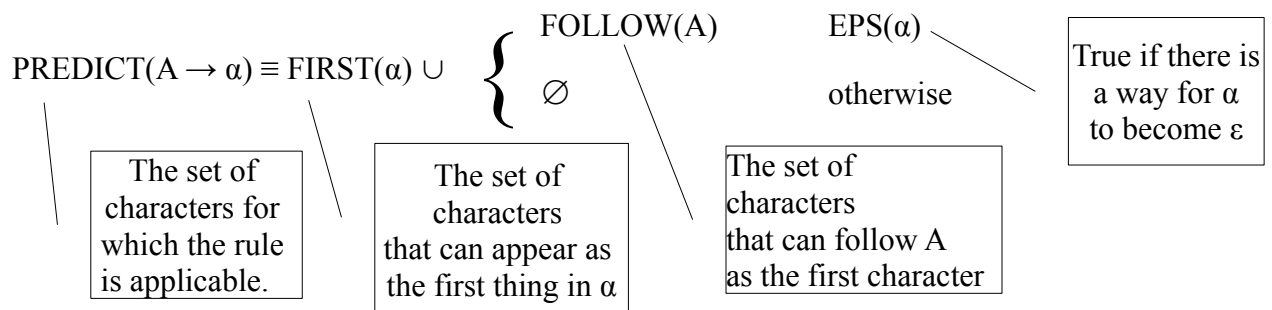
[illegible]

Example input: 3 * (x + y)

input	Production	
NUM(3)	<u>stmt</u>	R1
NUM(3)	<u>expr</u> \$	R2
NUM(3)	<u>term</u> term_tail \$	R5
NUM(3)	<u>factor</u> factor_tail term_tail \$	R10
NUM(3)	<u>NUM</u> factor_tail term_tail \$	Eat NUM
*	<u>factor_tail</u> term_tail \$	R6
*	<u>mul_op</u> factor factor_tail term_tail \$	R13
*	<u>*</u> factor factor_tail term_tail \$	Eat *
(<u>factor</u> factor_tail term_tail \$	R8
(<u>'('</u> expr ')' factor_tail term_tail \$	Eat (
ID(x)	<u>expr</u> ')' factor_tail term_tail \$	R2
ID(x)	<u>term</u> term_tail ')' factor_tail term_tail \$	R5
ID(x)	<u>factor</u> factor_tail term_tail ')' factor_tail term_tail \$	R9
ID(x)	<u>ID</u> factor_tail term_tail ')' factor_tail term_tail \$	Eat ID
+	<u>factor_tail</u> term_tail ')' factor_tail term_tail \$	R7
+	<u>term_tail</u> ')' factor_tail term_tail \$	R3
+	<u>add_op</u> term term_tail ')' factor_tail term_tail \$	R11
+	<u>'+'</u> term term_tail ')' factor_tail term_tail \$	Eat +
ID(y)	<u>term</u> term_tail ')' factor_tail term_tail \$	R5
ID(y)	<u>factor</u> factor_tail term_tail ')' factor_tail term_tail \$	R9
ID(y)	<u>ID</u> factor_tail term_tail ')' factor_tail term_tail \$	Eat ID
)	<u>factor_tail</u> term_tail ')' factor_tail term_tail \$	R7
)	<u>term_tail</u> ')' factor_tail term_tail \$	R4
)	<u> ')'</u> factor_tail term_tail \$	Eat)
\$	<u>factor_tail</u> term_tail \$	R7
\$	<u>term_tail</u> \$	R4
\$	<u>\$</u>	Eat \$

PREDICT Sets

How to find the set of characters for which the rule is applicable.



$$\text{FIRST}(\alpha) \equiv \{c: \alpha \Rightarrow^* c \beta\} \quad \text{EPS}(\alpha) \equiv [\alpha \Rightarrow^* \epsilon] \quad \text{FOLLOW}(A) \equiv \{c: S \Rightarrow^+ \alpha A c \beta\}$$

Here A is a non-terminal, α and β are sentential forms, which contains a mixture of terminals and non-terminals. \Rightarrow^* means 0 or more productions and \Rightarrow^+ means one or more productions. There are simple algorithms to find the PREDICT set programmatically (see the book).

PREDICT(1) $\equiv \{ '(', ID, NUM \}$

PREDICT(2) $\equiv \{ '(', ID, NUM \}$

PREDICT(3) $\equiv \{ '+', '-' \}$

PREDICT(4) $\equiv \{ '$, ')' \}$

FOLLOW(term_tail) \equiv FOLLOW(expr) $\equiv \{ '$, ')' \}$

PREDICT(5) $\equiv \{ '(', ID, NUM \}$

PREDICT(6) $\equiv \{ '*', '/' \}$

PREDICT(7) $\equiv \{ '+', '-', '$, ')' \}$

FOLLOW(factor_tail) \equiv FOLLOW(term)

$$\begin{aligned} & \equiv \text{FIRST}(\text{term_tail}) \cup \begin{cases} \text{FOLLOW}(\text{term_tail}) & \text{EPS}(\text{term_tail}) \\ \emptyset & \text{otherwise} \end{cases} \\ & \equiv \{ '+', '-' \} \cup \text{FOLLOW}(\text{term_tail}) \\ & \equiv \{ '+', '-' \} \cup \{ '$, ')' \} \\ & \equiv \{ '+', '-', '$, ')' \} \end{aligned}$$

Writing an LL(1) Grammar

Problems:

1. Left recursion

$$A \Rightarrow^+ A \alpha$$

E.g.: $\text{id_list} \rightarrow \text{id_list_prefix '}'$
 $\text{id_list_prefix} \rightarrow \text{id_list_prefix '}' \text{ ID}$
 $\text{id_list_prefix} \rightarrow \text{ID}$

2. Common prefixes

$$A \rightarrow \alpha \beta$$

$$A \rightarrow \alpha \gamma$$

E.g.: $\text{stmt} \rightarrow \text{ID} \text{ ":" expr}$
 $\text{stmt} \rightarrow \text{ID} \text{ '(' arguments ')'}$

Techniques:

1. Eliminate left recursion

$$A \rightarrow A \alpha_1 \mid \dots \mid A \alpha_n$$

$$A \rightarrow \beta_1 \mid \dots \mid \beta_m$$

becomes

$$A \rightarrow \beta_1 A' \mid \dots \mid \beta_m A'$$

$$A' \rightarrow \alpha_1 A' \mid \dots \mid \alpha_n A' \mid \epsilon$$

$\text{id_list} \rightarrow \text{id_list_body '}'$

$\text{id_list_body} \rightarrow \text{ID id_list_tail}$

$\text{id_list_tail} \rightarrow \text{' ID id_list_tail} \mid \epsilon$

2. Eliminate common prefixes by left factoring

$$A \rightarrow \alpha \beta$$

$$A \rightarrow \alpha \gamma$$

becomes

$$A \rightarrow \alpha \text{ tail}$$

$$\text{tail} \rightarrow \beta \mid \gamma$$

$$\text{stmt} \rightarrow \text{ID stmt_tail}$$

$$\text{stmt_tail} \rightarrow \text{"::-"} \text{expr}$$

$$\text{stmt_tail} \rightarrow \text{'(' arguments ')}$$