

# 50006 - Compilers - (Dr Dulay) Lecture 3

Oliver Killane

09/04/22

## LL Parsing

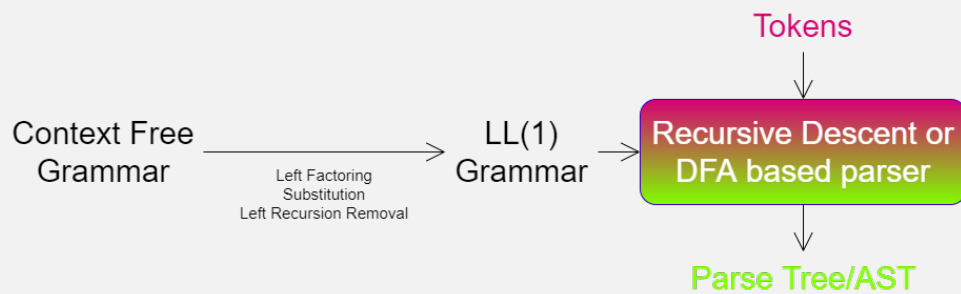
Top-down parsing using either recursive descent (can be hand-coded) or an **LL(1)** pushdown automaton (generated by an **LL(1)** parser-generator (e.g ANTLR)).

### Definition: LL(k) Grammar

A grammar is **LL(k)** if a  $k$ -token lookahead is sufficient to determine which alternative of a rule to use when parsing.

- **LL(1)** uses the current token only.
- **LL(0)** does not exist (would be deciding based on 0 tokens).
- When using **LL/Top Down Parsing** leaf nodes are constructed from the root.

The parser created for the grammar can be implemented as either a **DFA** or to parse by **recursive descent**.



### Definition: LL(1) Grammar

A grammar is **LL(1)** if for a rule  $A \rightarrow \alpha \mid \beta$  ( $A$  is non-terminal).

$$first(\alpha) \cap first(\beta) = \emptyset$$

$$\wedge \epsilon \in first(\alpha) \Rightarrow (first(\beta) \cap follow(A) = \emptyset)$$

$$\wedge \epsilon \in first(\beta) \Rightarrow (first(\alpha) \cap follow(A) = \emptyset)$$

$$A \rightarrow \underbrace{\alpha}_{t \in first(\alpha)} \mid \underbrace{\beta}_{t \in first(\beta)}$$

If both could start with t, if we get t as the current token, we cannot determine if we are parsing a or b.

$$A \rightarrow \underbrace{\alpha}_{\epsilon \in first(\alpha)} \mid \underbrace{\beta}_{t \in first(\beta)}$$

$follow(A) = \{\dots, t, \dots\}$

If a contains  $\epsilon$ , then we could potentially skip over this rule to whatever follows A.

Here with token t we cannot determine if we are skipping past this rule, or parsing b.

And the other order, and with many alternatives

This extends to more than two alternatives.

### Definition: Extended Backus-Naur Form (BNF)

Used for writing context free grammars, and includes several useful features:

$$\begin{aligned} \{\alpha\} & \quad 0 \text{ or more occurrences of } \alpha \\ [\alpha] & \quad 0 \text{ or 1 occurrences of } \alpha \\ (\dots) & \quad \text{A way to group elements together} \end{aligned}$$

For example we can left factor rules with alternatives that have intersecting first sets.

$$Expr \rightarrow Term \text{ '}' Expr \mid Term$$

then becomes

$$Expr \rightarrow Term \text{ ['}' Expr]}$$

Or we can remove left recursion:

$$Sequence \rightarrow Sequence \text{ ';' } Statement \mid Statement$$

then becomes

$$Sequence \rightarrow Statement \{ \text{' ;' } Statement \}$$

### Definition: Recursive Descent Parser

Consists of a set of parse functions for each rule which take in some tokens, and return remaining tokens and a generated **AST**.

We can use a basic function for matching **terminal** tokens:

```
1 # Get the next token. If it matches the expected, pop it from
2 # the list of tokens, else throw an error.
3 def match(expected: token):
4     if lexical_analyser.next_token() == expected:
5         lexical_analyser.pop_token()
6     else:
7         raise error("Unexpected token")
```

We can apply some basic rules for the main patterns:

```
1 # A B
2 A()
3 B()
4
5 # A | B (to be an LL(1) grammar, first sets must be disjoint)
6 if next_token() in first(A):
7     A()
8 elif next_token() in first(B):
9     B()
10
11 # {A}
12 while next_token() in first(A):
13     A()
14
15 # [A]
16 if next_token() in first(A):
17     A()
```

### Example: Statements

$$\begin{aligned} Stat &\rightarrow IfStat \mid BeginStat \mid PrintStat \\ IfStat &\rightarrow \text{'if' } Expr \text{'then' } Stat [\text{'else' } Stat] \text{'fi'} \\ BeginStat &\rightarrow \text{'begin' } Stat \{ \text{';' } Stat \} \text{'end'} \\ PrintStat &\rightarrow \text{'print' } Expr \end{aligned}$$

```
1 def Stat() -> Statement:
2     if next_token == IF:
3         return IfStat()
4     elif next_token == BEGIN:
5         return BeginStat()
6     elif next_token == PRINT:
7         return PrintStat()
8     else:
9         raise error("Expected Statement Starting Token")
10
11 # Parse an if statement, returning the statement if parse succeeds.
12 def IfStat() -> Statement:
13     match(IF)
14     cond = Expr()
15     match(THEN)
16     if_branch = Stat()
17     else_branch = None
18     if next_token() == ELSE:
19         match(ELSE)
20         else_branch = Stat()
21     match(FI)
22     return IfStatement(cond, if_branch, else_branch)
23
24 # Parse a block of statements, returning the block statement if
25   ↳ successful
26 def BeginStat() -> Statement:
27     stats = []
28     match(BEGIN)
29     stats.append(Stat())
30     while next_token() == SEMICOLON:
31         match(SEMICOLON)
32         stats.append(Stat())
33     match(END)
34     return Block(stats)
35
36 # Parse a print statement, returning the statement if successful.
37 def PrintStat() -> Statement:
38     match(PRINT)
39     return PrintStatement(Expr())
```

## AST Construction

Class hierarchies can be used to create/organise nodes into variants of a given type (e.g if statements, print statements and assignments are all statements, so implement/inherit from some Statement class/interface).

### Other languages

While in languages such as Java or Python, class hierarchies are used to implement ASTs other languages have cleaner representations.

- In **Haskell** the `data` keyword can be used to build an **ADT** for an **AST**
- Rust supports enums similar to the haskell `data` keyword.

## CFG $\rightarrow$ LL(1) Conversion

Transformations must be applied to a non-LL(1) context free grammar, which cannot always be automated.

- Left Recursion Removal and Substitution are usually applied first.
- The semantics of the grammar must be maintained.
- Readability is a key consideration, especially if the language is to be extended in future.

### Definition: Left Factorisation

Where two or more alternatives of a rule have a common prefix (first element to be parsed), factor this to be parsed before determining which alternative to parse.

non-LL(1)	EBNF LL(1)	BNF LL(1)
$A \rightarrow B C \mid B D$	$A \rightarrow B (C \mid D)$	$A \rightarrow B X$ $X \rightarrow C \mid D$
$A \rightarrow B C \mid B$	$A \rightarrow B [C]$	$A \rightarrow B X$ $X \rightarrow C \mid \epsilon$

### Definition: Substitution

Substituting a rule/non-terminal with its alternatives.

- Can be used to find indirect conflicts that may require left-factoring.
- Can produce grammars encoding more information in a rule (makes it easier to produce the required AST)

$$\begin{aligned} A &\rightarrow B \mid C \\ B &\rightarrow \text{'hello'} \\ C &\rightarrow \text{'hello'} \text{'there'} \end{aligned}$$

Here we have an indirect conflict as both alternatives for  $C$  start with 'hello'.

$$A \rightarrow \text{'hello'} \mid \text{'hello'} \text{'there'}$$

Now we can left-factor.

$$A \rightarrow \text{'hello'} [\text{'there'}]$$

#### Definition: Left Recursion Removal

Grammars cannot be **LL(1)** with left-recursion. We can use **direct left recursion removal** to eliminate left recursion from rules while leaving the grammar mostly intact.

$$A \rightarrow X \mid A Y \Rightarrow A \rightarrow X \{Y\}$$

For example with left-associative arithmetic expressions:

$$\begin{aligned} Expr &\rightarrow Expr \text{ ('+' | '-' ) } Term \mid Term \\ Term &\rightarrow Term \text{ ('*' | '/' ) } Factor \mid Factor \\ Factor &\rightarrow \text{'(' Expr ')'} \mid \underline{int} \end{aligned}$$

The **parse tree** may no longer represent the associativity, hence we will need to ensure the arithmetic is still associative when we construct the **AST**:

$$\begin{aligned} Expr &\rightarrow Term \{ \text{'+' | '-'} \} Term \\ Term &\rightarrow Factor \{ \text{'*' | '/'} \} Factor \\ Factor &\rightarrow \text{'(' Expr ')'} \mid \underline{int} \end{aligned}$$

## Error Recovery

When detecting an error in parsing...

- Useful error messages need to contain information collected during parsing.
- Error recovery can be used to allow further parse errors to be detected, with as little code as possible being skipped but avoiding nonsense/spurious error messages being created as a result.
- Error correction can be attempted to allow for semantic checks to be performed. The corrections must attempt to emulate what the semantics of the erroneous code likely was.

### Definition: Panic Mode Recovery

Each parse function has a **syncset** of tokens. When an error occur, the parser skips forwards until it encounters one of these tokens.

- Additional tokens can be provided as arguments to the parsing function (e.g add **follow set** of outer non-terminal to inner parse)
- The **follow set** of a rule is often used as the **syncset**
- 

```
1  # match a token, if necessary report and error, return boolean if the
   ↪ token matched expected.
2  def match(expected) -> bool:
3      if next_token() == expected:
4          pop_token()
5          return True
6      else:
7          add_error(next_token(), parser_pos(), [expected],
   ↪ INCORRECT.TOKEN)
8          return False
9
10 # Skip until at the end of the file, or token matches the sync set
11 def skipto(syncset):
12     while next_token() not in syncset and next_token() != EOF:
13         pop_token()
14
15
16 def check(expectset, syncset, error):
17     if next_token() not in expectset:
18         add_error(next_token(), parser_pos(), expectset, error)
19         skipto(expectset + syncset)
```