

Parsing — Part II (Top-down parsing, left-recursion removal)

Parsing Techniques



Top-down parsers (LL(1), recursive descent)

- Start at the root of the parse tree and grow toward leaves
- Pick a production & try to match the input
- Bad "pick" \Rightarrow may need to backtrack
- Some grammars are backtrack-free (predictive parsing)

Bottom-up parsers (LR(1), operator precedence)

- Start at the leaves and grow toward root
- As input is consumed, encode possibilities in an internal state
- Start in a state valid for legal first tokens
- Bottom-up parsers handle a large class of grammars



A top-down parser starts with the root of the parse tree The root node is labeled with the goal symbol of the grammar

Top-down parsing algorithm:

Construct the root node of the parse tree Repeat until the fringe of the parse tree matches the input string

- 1 At a node labeled A, select a production with A on its lhs and, for each symbol on its rhs, construct the appropriate child
- 2 When a terminal symbol is added to the fringe and it doesn't match the fringe, backtrack
- 3 Find the next node to be expanded

(label \in NT)

- The key is picking the right production in step 1
 - → That choice should be guided by the input string

Remember the expression grammar?

ELAWARE 1743 s

Version with precedence derived last lecture

And the input $\underline{x} - \underline{2} \star \underline{y}$

A possible parse



Consider the following parse of x - 2 * y

1				
	Rule	Sentential Form	Input	
	Ι	Goal	↑ <u>×</u> - <u>2</u> * γ	
	1	Expr	↑ <u>×</u> - <u>2</u> *⊻	
	2	Expr + Term	↑ <u>× - 2</u> * ⊻	
	2	Expr + Term + Term	<u>↑× - 2 * y</u>	
	2	Expr+ Term + Term + Term	↑ <u>× - 2 * y</u>	
	2	Expr+Term + Term ++Term	1 × − 2 * y	

consuming no input !

This doesn't terminate

(obviously)

- Wrong choice of expansion leads to non-termination
- Non-termination is a bad property for a parser to have
- Parser must make the right choice



Top-down parsers cannot handle left-recursive grammars

Formally,

A grammar is left recursive if $\exists A \in NT$ such that $\exists a \text{ derivation } A \Rightarrow^{+} A\alpha$, for some string $\alpha \in (NT \cup T)^{+}$

Our expression grammar is left recursive

- This can lead to non-termination in a top-down parser
- For a top-down parser, any recursion must be right recursion
- We would like to convert the left recursion to right recursion

Non-termination is a bad property in any part of a compiler



To remove left recursion, we can transform the grammar

```
Consider a grammar fragment of the form
```

```
Fee \rightarrow Fee \alpha
```

```
β
```

where neither α nor β start with Fee

```
We can rewrite this as
```

```
Fee → β Fie
Fie → α Fie
```

3 |

where Fie is a new non-terminal

This accepts the same language, but uses only right recursion



If it picks the wrong production, a top-down parser may backtrack Alternative is to look ahead in input & use context to pick correctly

How much lookahead is needed?

• In general, an arbitrarily large amount

Fortunately,

- Large subclasses of CFGs can be parsed with limited lookahead
- Most programming language constructs fall in those subclasses

Among the interesting subclasses are LL(1) and LR(1) grammars

Basic idea

Given $A \rightarrow \alpha \mid \beta$, the parser should be able to choose between $\alpha \& \beta$

FIRST sets

For some *rhs* $\alpha \in G$, define FIRST(α) as the set of tokens that appear as the first symbol in some string that derives from α That is, $\underline{x} \in FIRST(\alpha)$ *iff* $\alpha \Rightarrow^* \underline{x} \gamma$, for some γ

We will defer the problem of how to compute FIRST sets until we look at the *LR(1)* table construction algorithm



Basic idea

Given $A \rightarrow \alpha \mid \beta$, the parser should be able to choose between $\alpha \& \beta$

FIRST sets

For some $rhs \alpha \in G$, define FIRST(α) as the set of tokens that appear as the first symbol in some string that derives from α That is, $\underline{\mathbf{x}} \in \text{FIRST}(\alpha)$ iff $\alpha \Rightarrow^* \underline{\mathbf{x}} \gamma$, for some γ

The LL(1) Property

If $A \rightarrow \alpha$ and $A \rightarrow \beta$ both appear in the grammar, we would like

FIRST(α) \cap FIRST(β) = \emptyset

This would allow the parser to make a correct choice with a lookahead of exactly one symbol! This is almost correct



See the next slide





What about $\epsilon\text{-productions?}$

 \Rightarrow They complicate the definition of LL(1)

If $A \rightarrow \alpha$ and $A \rightarrow \beta$ and $\varepsilon \in \text{FIRST}(\alpha)$, then we need to ensure that $\text{FIRST}(\beta)$ is disjoint from $\text{FOLLOW}(\alpha)$, too

Define FIRST⁺(α) as

- FIRST(α) \cup FOLLOW(α), if $\varepsilon \in$ FIRST(α)
- FIRST(α), otherwise

Then, a grammar is *LL(1)* iff $A \rightarrow \alpha$ and $A \rightarrow \beta$ implies

 $\mathsf{FIRST}^{+}(\alpha) \cap \mathsf{FIRST}^{+}(\beta) = \emptyset$

FOLLOW(α) is the set of all words in the grammar that can legally appear immediately after an α

VIVERSITY OF ELAWARE

Given a grammar that has the LL(1) property

- Can write a simple routine to recognize each *lhs*
- Code is both simple & fast

Consider $A \rightarrow \beta_1 \mid \beta_2 \mid \beta_3$, with FIRST⁺(β_1) \cap FIRST⁺ (β_2) \cap FIRST⁺ (β_3) = \emptyset

Grammars with the *LL(1)* property are called predictive grammars because the parser can "predict" the correct expansion at each point in the parse.

Parsers that capitalize on the *LL(1)* property are called predictive parsers.



We build a recursive descent parser for the following grammar:

$$A \rightarrow B \mid CA \mid a$$
$$B \rightarrow bB \mid x$$
$$C \rightarrow c$$

The term descent refers to the direction in which the parse tree is built.

Recursive Descent Parsing

To actually build a parse tree:

- Augment parsing routines to build B() nodes
- Node for each symbol on *rhs*
- Action is to receive all *rhs* nodes, make them children of *lhs* node, and return this new node

To build an abstract syntax tree

- Build fewer nodes
- Put them together in a different order

```
B()
```

if (lookahead() = b)
then return B().addOne();
if (lookahead() = x)
then return new BNode(0);
throw Exception;

This is a preview of Chapter 4



if (lookahead() = b)
then return new BNode(read(),B());
if (lookahead() = x)
then return new BNode(read());
throw Exception;

Left Factoring



What if my grammar does not have the LL(1) property?

 \Rightarrow Sometimes, we can transform the grammar

How would you rewrite the grammar

 $A \rightarrow aab \mid aac \mid aad$



Left Factoring



What if my grammar does not have the LL(1) property?

 \Rightarrow Sometimes, we can transform the grammar

How would you rewrite the grammar

 $A \rightarrow aab \mid aac \mid aad$

Rewrite to

$$A \rightarrow aa A'$$
$$A' \rightarrow b \mid c \mid d$$

(Generality)



<u>Question</u>

By *eliminating left recursion* and *left factoring*, can we transform an arbitrary CFG to a form where it meets the *LL(1)* condition? (and can be parsed predictively with a single token lookahead?)

<u>Answer</u>

Given a CFG that doesn't meet the *LL(1)* condition, it is undecidable whether or not an equivalent *LL(1)* grammar exists.

<u>Example</u>

 $\{a^n 0 b^n \mid n \ge 1\} \cup \{a^n 1 b^{2n} \mid n \ge 1\}$ has no LL(1) grammar

Language that Cannot Be LL(1)



<u>Example</u>

 $\{a^n 0 b^n | n \ge 1\} \cup \{a^n 1 b^{2n} | n \ge 1\}$ has no *LL(1)* grammar

 $G \rightarrow \underline{aAb}$ $\mid \underline{aBbb}$ $A \rightarrow \underline{aAb}$ $\mid \underline{0}$ $B \rightarrow \underline{aBbb}$ $\mid \underline{1}$

Problem: need an unbounded number of <u>a</u> characters before you can determine whether you are in the A group or the B group.

Recursive Descent (Summary)

- 1. Build FIRST (and FOLLOW) sets
- 2. Massage grammar to have *LL(1)* condition
 - a. Remove left recursion
 - b. Left factor it
- 3. Define a procedure for each non-terminal
 - a. Implement a case for each right-hand side
 - b. Call procedures as needed for non-terminals
- 4. Add extra code, as needed
 - a. Perform context-sensitive checking
 - b. Build an IR to record the code

Can we automate this process?





Given an LL(1) grammar, and its FIRST & FOLLOW sets ...

- Emit a routine for each non-terminal
 - \rightarrow Multiple if-then statements to check alternate rhs's
 - \rightarrow Each returns a node on success and throws an error else
 - \rightarrow Simple, working (, *perhaps ugly*,) code
- This automatically constructs a recursive-descent parser

I don't know of a system that does this ...

Improving matters

- Bunch of if-then statements may be slow
 - \rightarrow Good case statement implementation would be better
- What about a table to encode the options?
 - \rightarrow Interpret the table with a skeleton, as we did in scanning

Building Top-down Parsers

Strategy

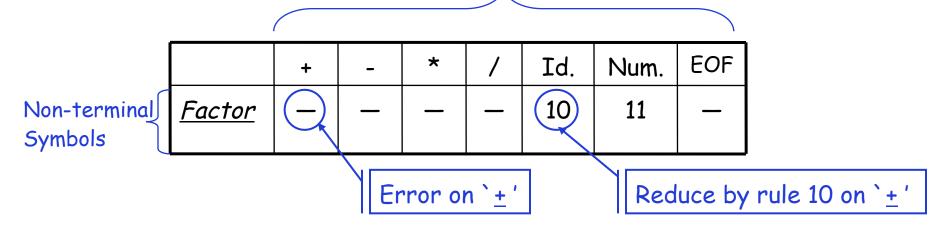
- Encode knowledge in a table
- Use a standard "skeleton" parser to interpret the table

Example

• The non-terminal *Factor* has three expansions

 \rightarrow (*Expr*) or <u>Identifier</u> or <u>Number</u>

• Table might look like: Terminal Symbols





Building Top Down Parsers

Building the complete table

- Need a row for every NT & a column for every T
- Need an algorithm to build the table

Filling in TABLE[X,y], $X \in NT$, $y \in T$

- 1. entry is the rule $X \rightarrow \beta$, if $y \in FIRST(\beta)$
- 2. entry is the rule $X \rightarrow \varepsilon$ if $y \in FOLLOW(X)$ and $X \rightarrow \varepsilon \in G$
- 3. entry is error if neither 1 nor 2 define it

If any entry is defined multiple times, G is not LL(1)

This is the *LL(1)* table construction algorithm





```
word \leftarrow nextWord()
push EOF onto Stack
push the start symbol onto Stack
TOS \leftarrow top of Stack
loop forever
  if TOS = EOF and word = EOF then
                                                             exit on success
      report success and exit
  else if TOS is a terminal or eof then
    if TOS matches word then
       pop Stack
                                         // recognized TOS
       word \leftarrow nextWord()
    else
        report error looking for TOS
  else
                                          // TOS is a non-terminal
    if TABLE[TOS,word] is A \rightarrow B_1 B_2 \dots B_k then
       pop Stack
                                         // get rid of A
       push B_k, B_{k-1}, ..., B_1 on stack
                                        // in that order
    else report error expanding TOS
  TOS \leftarrow top of Stack
```