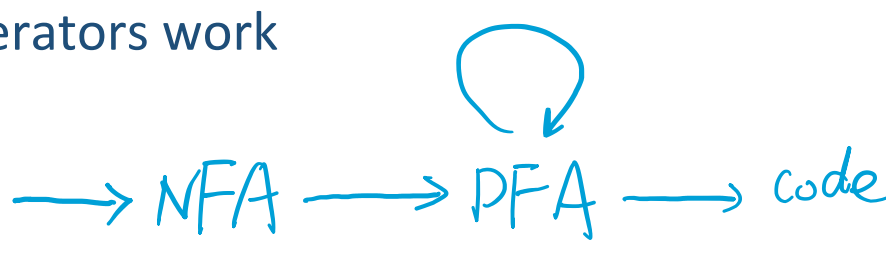


Parsing Overview

Last time: How lexer generators work

Spec = R1 { Action1 }
| R2 { Action2 }
| R3 { Action 3 }
... // in decreasing priority



RE Derivatives

Sometimes **precomputing** NFA/DFA is not worth the trouble or cost => RE libraries

One elegant idea to implement such a RE library: **interpreting** RE directly

Brzozowski

$$RE_0 \xrightarrow{x_1} RE_1 \xrightarrow{x_2} RE_2 \rightarrow \dots \rightarrow \varepsilon$$

$$\partial_x R = R'$$

$$\partial_a a = \varepsilon$$

$$\partial_a \varepsilon = \emptyset$$

$$\partial_a b = \emptyset \quad b \neq a$$

$$\partial_a (R_1 + R_2) = \partial_a R_1 + \partial_a R_2$$

$$\partial_a (R^*) = (\partial_a R) R^*$$

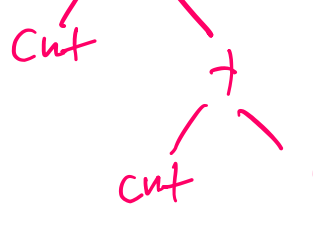
$$\partial_a (R_1 R_2) = (\partial_a R_1) R_2 + v(R_1) (\partial_a R_2)$$
$$v(R_1) = \begin{cases} \varepsilon & \text{if } R_1 \text{ accepts empty str} \\ \emptyset & \text{otherwise} \end{cases}$$

Today: An Overview of Parsing

Where we are:

characters $\xrightarrow{\text{lexing}}$ tokens $\xrightarrow{\text{parsing}}$ AST

$$cnt = cnt + 1$$



Syntax specification: context-free grammar (CFG)

$$S \rightarrow S + S$$

$$S \rightarrow S * S$$

$$S \rightarrow E$$

$$E \rightarrow n$$

$$E \rightarrow (S)$$

- productions
- nonterminal S, E
- terminal symbols $+ * n ()$
- symbols: terminals or nonterminals
- one of the symbols is the start symbol

expressive power comes from recursion (compare with regular expressions)

Some definitions:

- $L(G)$: the language of grammar G
- token stream $\in L(G)$: there is a **derivation** using G

Derivation

Derivation: a sequence of rewrite steps from start symbol to token stream

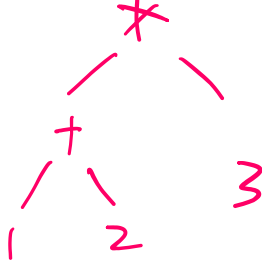
$$\alpha X \beta \rightarrow \alpha \gamma \beta \quad \text{where } X \rightarrow \gamma \text{ is a production}$$

An example derivation:

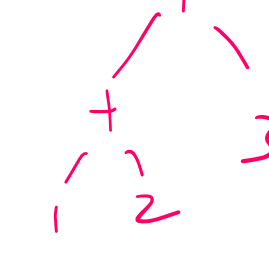
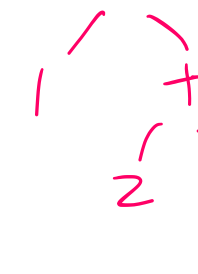
Grammar:

$S \rightarrow S + S$
 $S \rightarrow S * S$
 $S \rightarrow E$
 $E \rightarrow n$
 $E \rightarrow (S)$

$$1 + 2 * 3$$



$$1 + 2 + 3$$



Token stream: $(1 + 4) + 2$

Q: How would we derive the token stream using the grammar?

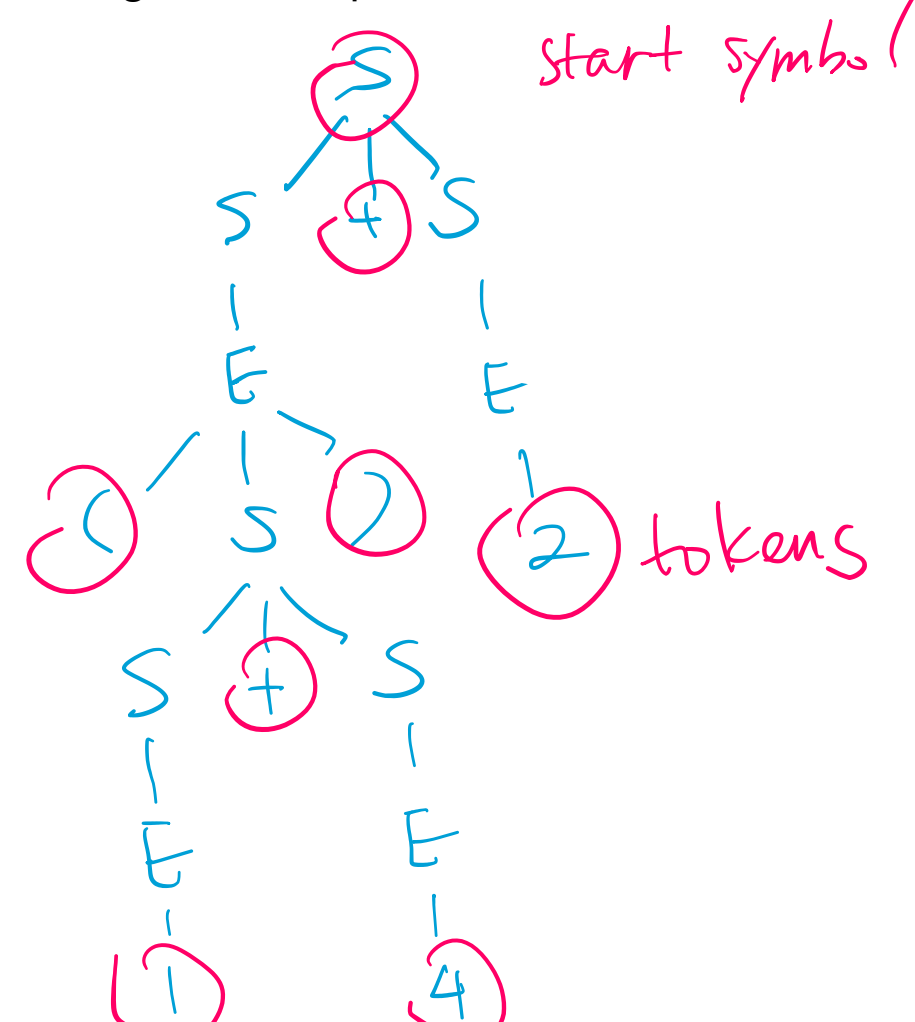
$$\underline{S} \rightarrow \underline{S} + \underline{S} \rightarrow \underline{E} + \underline{S} \rightarrow (\underline{S}) + \underline{S} \rightarrow (\underline{S} + \underline{S}) + \underline{S} \rightarrow \dots \rightarrow (1 + 4) + \underline{S} \rightarrow (1 + 4) + E \rightarrow (1 + 4) + 2$$

Leftmost derivation: expand leftmost nonterminal
turns out to correspond to top-down parsing

Rightmost derivation: expand rightmost nonterminal
turns out to correspond to bottom-up parsing

Parse tree

Derivation generates a parse tree.



Ambiguous grammar

\exists stream of tokens in $L(G)$ s.t. > 1 parse trees for that stream

Fixing ambiguity

Specifying precedence and associativity using additional nonterminals

$$S \rightarrow S + T \mid T$$

$$\text{term } T \rightarrow T * F \mid F$$

$$\text{factor } F \rightarrow n \mid (S)$$

Some parser generators allow precedence/associativity declarations:

CUP precedence Left PLUS;
precedence left TIMES;

Dangling-else problem

$$S \rightarrow \text{if } E \text{ then } S \mid \text{if } E \text{ then } S \text{ else } S \mid \dots$$

$$\text{if } E_1 \text{ then (if } E_2 \text{ then } S_1 \text{ else } S_2)$$

CUP precedence nonassoc IF;
precedence nonassoc ELSE;

$$S \rightarrow U \mid M$$

$$M \rightarrow \text{if } E \text{ then } M \text{ else } M$$

...

$$U \rightarrow \text{if } E \text{ then } S$$

$$\mid \text{if } E \text{ then } M \text{ else } U$$

- statements ending w/ unmatched THEN
- otherwise

Needing context

Example in C++: `int x, k, V, HashTable;`

`HashTable<K,V> x;`

Significant whitespace

Top-Down Parsing.

LL parser

Example S-expressions (McCarthy - Lisp)

$$(\text{let } ((x \ z) \ (+ \ x \ 1))) \quad \text{let } x = z \text{ in } x + 1$$

$$S \rightarrow (L) \mid x$$

$$L \rightarrow \varepsilon \mid S L$$

parse (a b)

Recursive descent

```
parse S () {  
  switch (peek()) {  
    case "(": // S -> (L)  
      consume "(";  
      parse L ();  
      consume ")";  
      return;  
    case ID: // S -> x  
      consume ID;  
      return;  
    default:  
      throw SyntaxError;  
  }  
}
```

```
parse L () {  
  switch (peek()) {  
    case ID: // L -> SL  
      case "(": // L -> SL  
        parse S ();  
        parse L ();  
        return;  
    case ")": // L -> ε  
      return;  
  }  
}
```

default:
throw SyntaxError

