

NOTES ON ON CONTEXT-FREE GRAMMARS
by PRAKASH PANANGADEN

Idea : One can generate sentences by giving rules for rewriting "templates" into actual sentences.

Some Definitions

$\Sigma \rightarrow$ a set of symbols

$\Sigma^* \rightarrow$ all finite sequences of Σ symbols

e.g. $\Sigma = \{a, b\}$ $\Sigma^* = \{\epsilon, a, b, aa, ab, ba, bb, \dots\}$

where " ϵ " means empty sequence

A language over Σ is a subset of Σ^* .

e.g. The language L of all sequences of "a"s and "b"s with an equal number of "a"s and "b"s is a language over $\{a, b\}$.

For example aab is not in L but $bbabbaaa$ is.

To describe languages we use grammars.

Def A grammar has a set Σ of symbols, a set NT of non-terminal symbols, a special symbol (usually S) in NT called the start symbol and a set of rules

or productions of the form

$$A \rightarrow \text{sequence from } (\Sigma \cup N.T.)$$

Discussion : We start with the start symbol and use rules to rewrite the sequences until we get a sequence of symbols from Σ . After this no more rewriting is possible, so we stop. For this reason, symbols in Σ are called terminals.

Ex Here is a grammar for the language L consisting of sequences of a's & b's with an equal number of a's & b's:

$$1. \quad S \rightarrow \epsilon$$

$$3. \quad S \rightarrow bSa$$

$$2. \quad S \rightarrow SS$$

$$4. \quad S \rightarrow aSb$$

An example of generation using this grammar

$$\begin{aligned} S &\xrightarrow{3} \underline{bSa} \xrightarrow{3} b \underline{bSa} a \xrightarrow{4} bb \underline{aSb} aa \xrightarrow{2} \\ &bb a \underline{SS} baa \xrightarrow{4} bb a a \underline{Sb} Sb a a \xrightarrow{3} bb a a Sb \underline{bSa} baa \\ &\xrightarrow{1,1} bb a a, bb a b a a \quad (\text{rule 1 used twice}) \end{aligned}$$

I have written the rule number above each arrow and I have put a dashed underline below the part of the sequence that was just generated.

A grammar must generate all the sequences in the language and only the sequences in the language.

To save space we can put several rules with the same LHS on the same line:

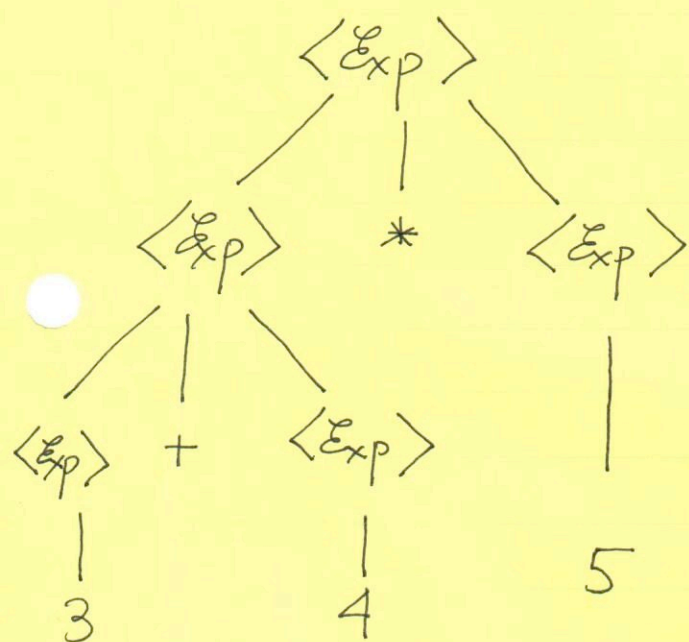
$$S \rightarrow \epsilon / SS / bSa / aSb$$

Grammars for expressions

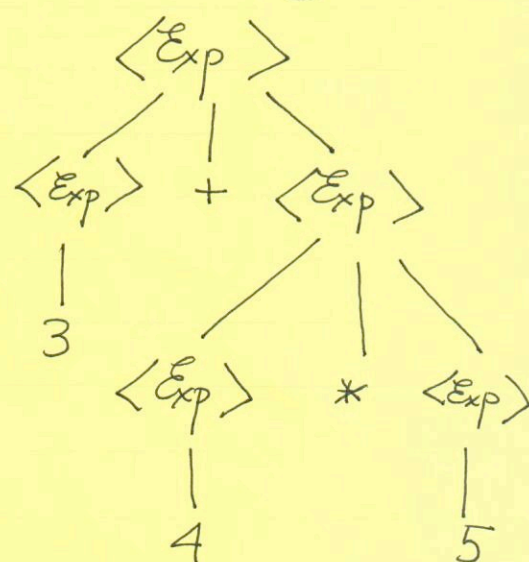
$$NT = \{ \langle \text{Exp} \rangle \} \quad \Sigma = \{ \text{numbers}, +, *, (,) \}$$

$$\langle \text{Exp} \rangle \rightarrow (\langle \text{Exp} \rangle) \mid \langle \text{Exp} \rangle + \langle \text{Exp} \rangle \mid \langle \text{Exp} \rangle * \langle \text{Exp} \rangle \mid \text{number}$$

This correctly generates the language but produces trees with the wrong grouping.

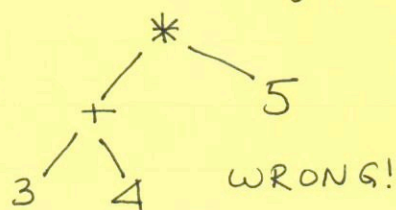


Parse tree for
 $3 + 4 * 5$

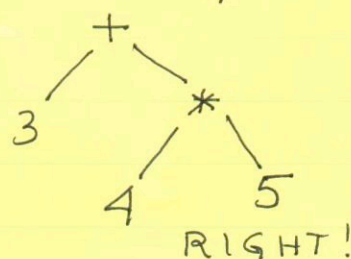


Another parse tree for
 $3 + 4 * 5$

A terrible situation! The same string has 2 parse trees and one of them has the wrong grouping. If we look at the corresponding expression trees we get



and



When a grammar produces 2 or more parse trees for a sequence it is called ambiguous. We need to design unambiguous grammars; but sometimes this is impossible. Here is how we do it for our expressions with generation of numbers thrown in as well.

$NT = \{ \langle N \rangle, \langle D \rangle, \langle E \rangle, \langle F \rangle, \langle T \rangle \}$

Start Symbol is $\langle E \rangle$

$\Sigma = \{ +, *, (,), 0, 1, 2, 3, 4, 5, 6, 7, 8, 9 \}$

$\langle E \rangle \rightarrow \langle E \rangle + \langle T \rangle \mid \langle T \rangle$

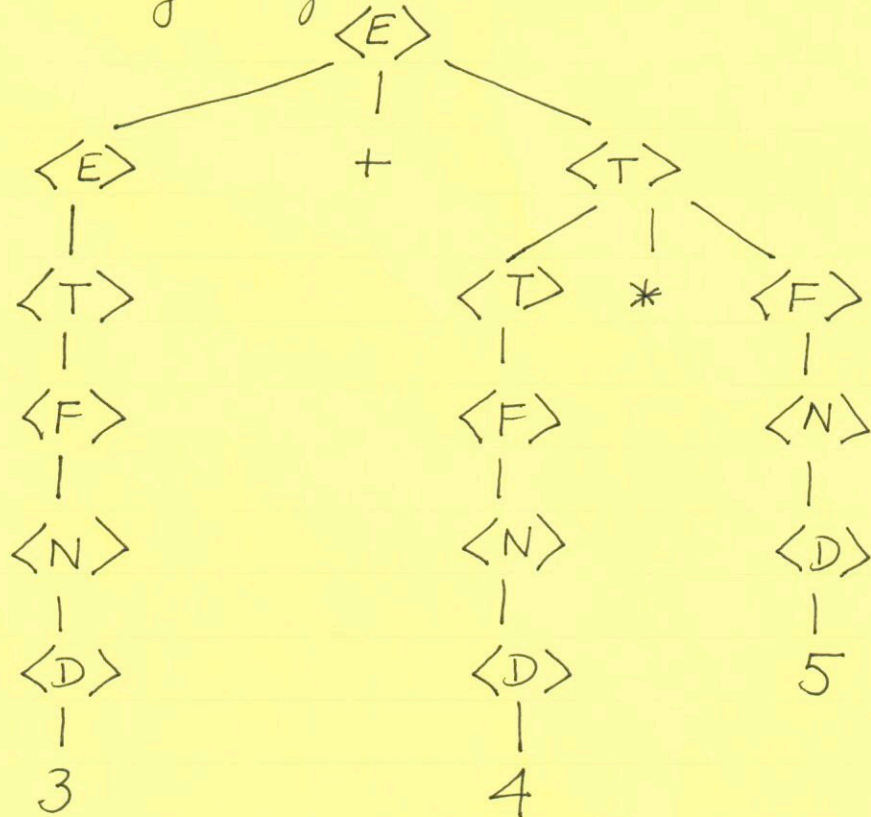
$\langle T \rangle \rightarrow \langle T \rangle * \langle F \rangle \mid \langle F \rangle$

$\langle F \rangle \rightarrow (\langle E \rangle) \mid \langle N \rangle$

$\langle N \rangle \rightarrow \langle N \rangle \langle D \rangle \mid \langle D \rangle$

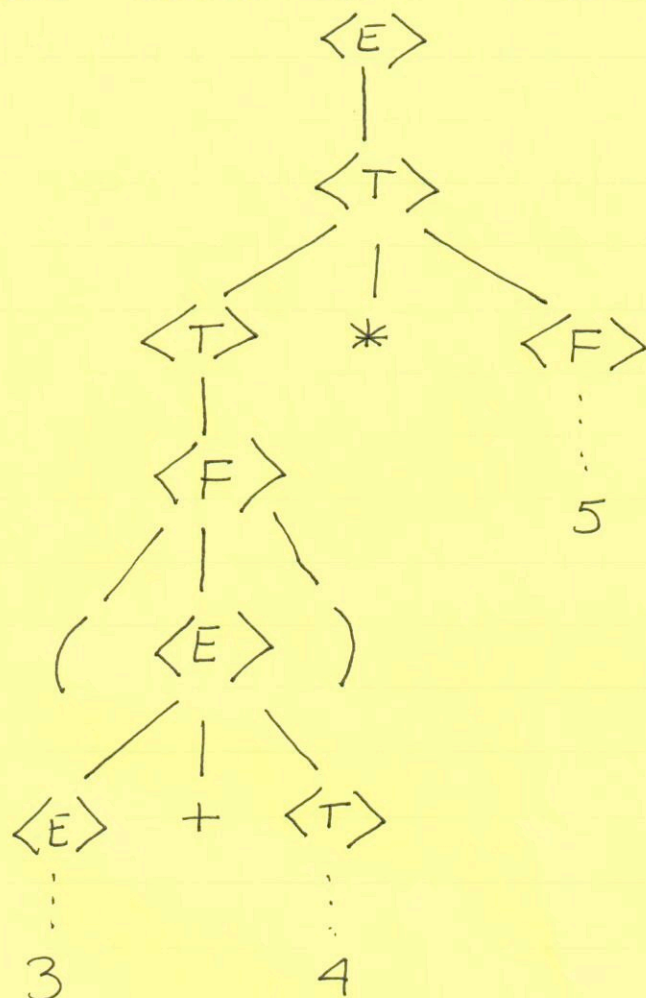
$\langle D \rangle \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

Now we try to generate $3 + 4 * 5$



This is the correct grouping.

Let us try to create the "wrong" grouping.
We want to make the $*$ appear at the top of the tree with the $+$ nested inside it.



The lines indicate the "obvious" steps that I have left out.

Notice how we were forced to put parens around $3 + 4$ when we tried to do this.

Recursive - descent parsing

Here is how you recognize balanced parentheses:

1. Find a left parenthesis
2. Find a balanced string
3. Find a right parenthesis
4. Find a balanced string

This is suggested by the grammar production

$$S \rightarrow (S)S$$

How can this work? It is plainly recursive, just use recursion.

variables first, second : boolean
nextsym : char

function check returns boolean.

Assume a global structure (eg. a file) from which read is getting each character.

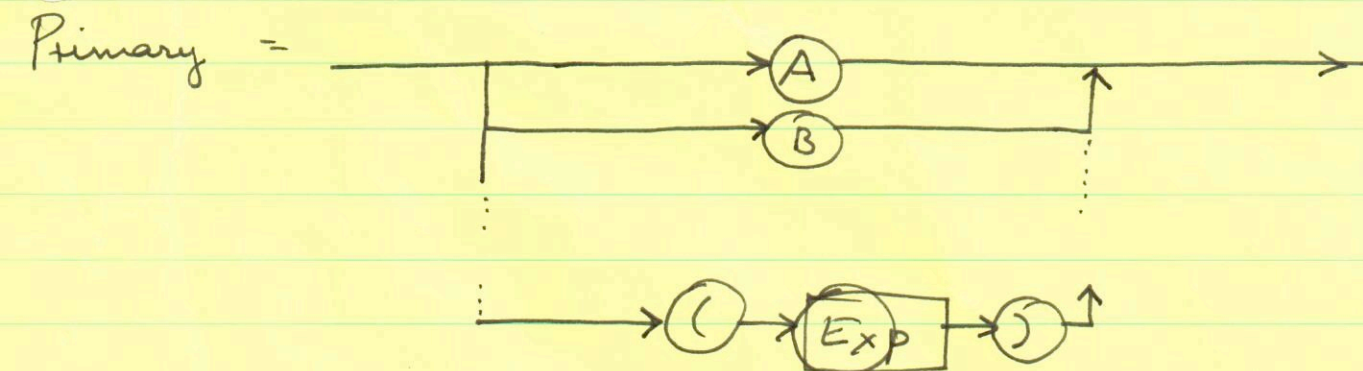
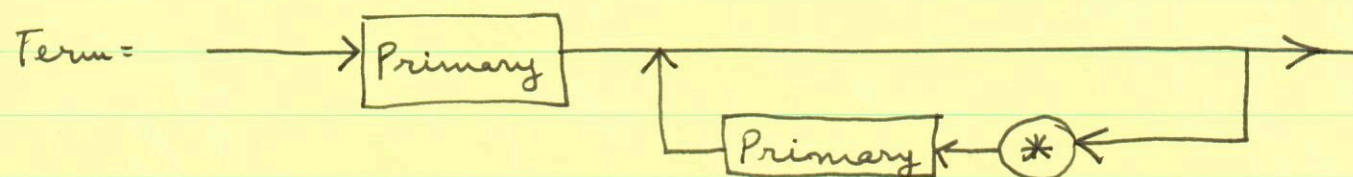
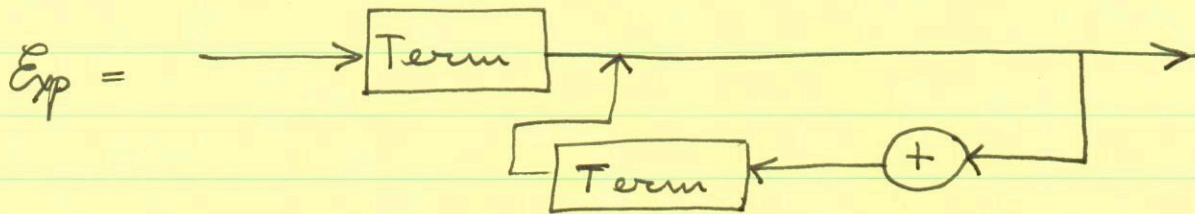
```

if nextsym = '(' then {
    read next char into nextsym;
    first = check(); [Recursion]
    if (first and nextsym = ')') then
        { read into nextsym;
          second = check();
          if second return TRUE
          else return FALSE; }
    else return FALSE
else return FALSE
  
```

Simple Expressions & Syntax Graphs

$+, *, A, B, \dots, Z$

$$\begin{aligned}\langle \text{Exp} \rangle &\rightarrow \langle \text{Term} \rangle \{ \langle \text{Term} \rangle \}^+ \\ \langle \text{Term} \rangle &\rightarrow \langle \text{Primary} \rangle \{ \langle \text{Primary} \rangle \}^* \\ \langle \text{Primary} \rangle &\rightarrow A | B | \dots | Z | (\langle \text{Exp} \rangle)\end{aligned}$$



This is the grammar for your assignment. Assume a procedure called GETSYM which takes the next symbol puts it in a global called "sym" and advances the cursor. The code is best organized using mutual recursion. I will discuss the control flow and not show how the tree is constructed.


```

procedure Exp;
begin

```

```

    Term;

```

```

    while sym = '+' do

```

```

        begin

```

```

            getsym; term

```

```

        end

```

```

    end

```

```

procedure Term;
begin

```

```

    primary;

```

```

    while sym = '*' do

```

```

        begin

```

```

            getsym; primary end

```

```

    end

```

```

procedure primary;
begin

```

```

    if sym in ['A'.. 'Z'] then getsym

```

```

    else if sym = '(' then

```

```

        begin getsym; expression;

```

```

            if sym = ')' then error else getsym

```

```

        end

```

```

    else error

```

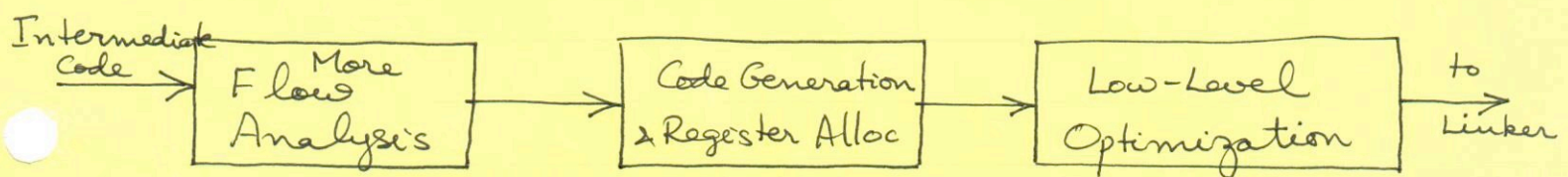
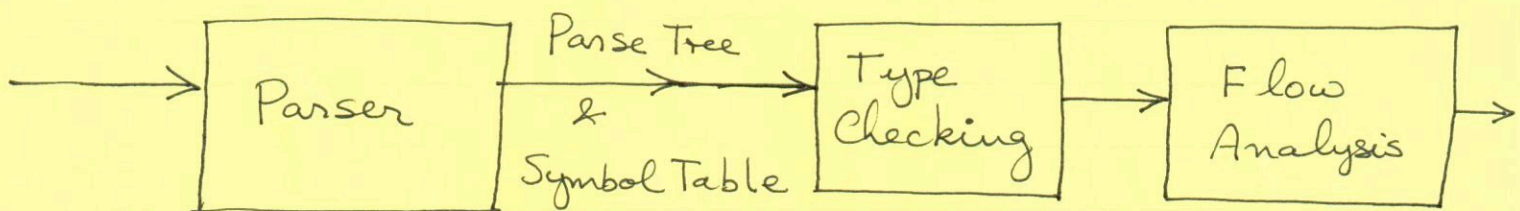
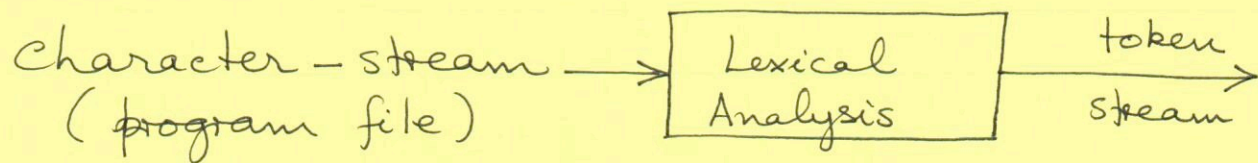
```

end

```

In your code these should be functions that return trees. This code only says whether there is an error or not but it gives the basic control flow. PLEASE produce expression trees not parse trees.

The structure of a compiler



lexical analysis & parsing are so well-understood that we can generate parsers & scanners automatically. The symbol table is the data-structure where variable names and their bindings are stored. The parse-tree is usually simplified to what is called an abstract-syntax tree. For our assignment I want you to generate expression trees directly.