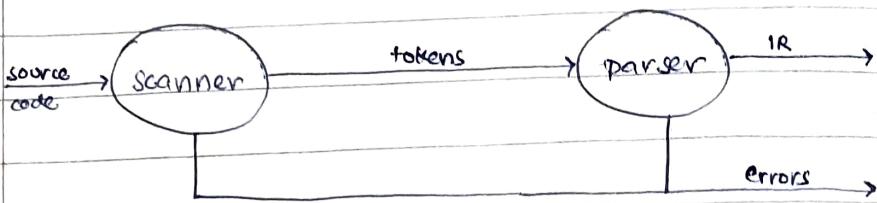


LEXICAL ANALYSISSCANNER

maps character stream to words - the basic unit of syntax
 produces pairs \rightarrow word + part-of-speech.
 (Token) (lexeme) (Token type)

$$x \leftarrow x + y; \quad \langle id, x \rangle = \langle id, x \rangle + \langle id, y \rangle$$

PARSER

checks if the input stream of tokens constitute a syntactically valid program of the language.

if the input program is syntactically correct

\hookrightarrow output a concrete representation of program (AST).

formally grammar $G_1 = (S, N, T, P)$

- S: start symbol
- N: set of non-terminal symbols
- T: set of terminal symbols or words
- P: set of production rules or rewrite rules.

Goal \rightarrow ExprExpr \rightarrow Expr Op Term | TermTerm \rightarrow Number | IdOp \rightarrow + | -

S = Goal (start symbol)

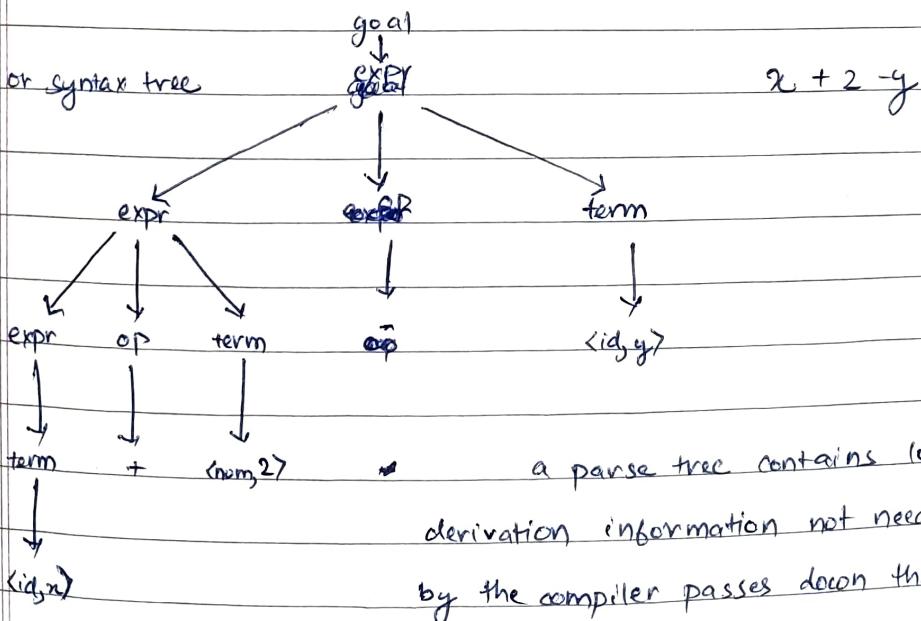
N = {Goal, Expr,

Term, op, }

T = {number, id,
+, -}

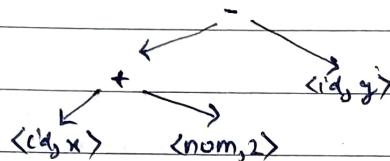
ex given a stream of tokens (read terminals) and the syntax specification in the form of a CFG, how can the parser check the syntactic correctness of the source code?

PARSE TREES

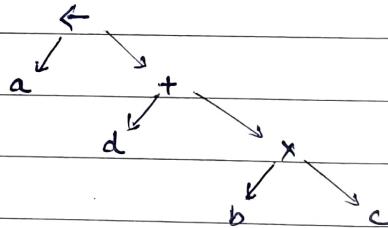


ABSTRACT SYNTAX TREES (AST)

Compilers often use an abstract syntax tree instead of a parse tree. It summarizes grammatical structure, without including detail about the derivation.

 $x + 2 - y$


- much more concise
- one kind of IR

 $a \leftarrow b * c + d$


LEXICAL TOKENS

 $\text{goal} \rightarrow \text{expr}$
 $\text{expr} \rightarrow \text{expr op term} \mid \text{term}$
 $\text{term} \rightarrow \text{number} \mid \text{id}$
 $\text{op} \rightarrow + \mid -$
 $\text{id} \rightarrow \text{alpha id} \mid \text{alpha}$
 $\text{alpha} \rightarrow \text{a} \mid \text{b} \mid \text{c} \mid \dots \mid \text{z}$
 $\text{number} \rightarrow \text{number digit} \mid \text{digit}$
 $\text{digit} \rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9$

goal → expr
 expr → expr op term | term
 term → number | id
 (+, - are possible ops)

goal → expr
 expr → expr + term
 | expr - term | term
 term → number | id

a token is a fundamental (atomic) syntactic unit in a language.
 "it depends upon how the grammar for the language is defined though".

MICRO-SYNTAX AND MACRO SYNTAX

Micro syntax:

the rules that govern the lexical structure of a language.

- specified using regular expressions
- implemented using finite state machines

Macro syntax:

the grammar of the language.

- specified using CFGs
- implemented using pushdown automata

Q: Why can't we encode the micro syntax of the language into the grammar for the macro syntax?

don't lift a needle with a crane?

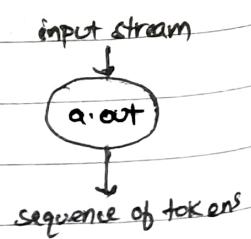
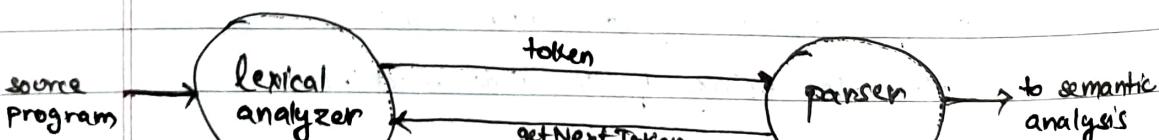
Separation of micro-syntax gives a clean grammar for the core programming language constructs.

grammar need not keep track of nitty gritty details which are taken care of cleanly and efficiently in lexical analyzer.

every grammar rule dropped shrinks the size of parser.

parsing is more harder and slower than lexical analysis.

LEXICAL ANALYSIS AND PARSER



STRUCTURE OF A LEX PROGRAM

declarations

7.7.

translation rules

pattern

{ action }

Y.Y.

auxiliary functions

ex /* regular definitions */

delim [\t\n]

ws { delim } +

letter [A-Za-z]

id letter { letter } { digit } *

digit [0-9]

number { digit } + (. { digit } +) ? (E [+-] ? { digit } +) ?

7.7.

{ ws }

{ * at no action and no return } */

{ if } { return IF; }

{ then } { return THEN; }

{ else } { return ELSE; }

{ id } { yyval = (int) installID(); return ID; }

{ number } { yyval = (int) installNum(); return NUMBER; }

{ < } { yyval = LT; return RELOP; }

{ <= } { yyval = LE; return RELOP; }

{ != } { yyval = NE; return RELOP; }

{ > } { yyval = GT; return RELOP; }

{ >= } { yyval = GE; return RELOP; }

RULES

- ① return longest matching lexeme (token).
- ② token types are prioritized according to their occurrence in the lex file.
- ③ always prefer a longer prefix to a shorter prefix.
- ④ if the longest possible prefix matches two or more patterns, prefer the pattern listed first in the flex program.

LEXICAL ANALYZER DESIGN

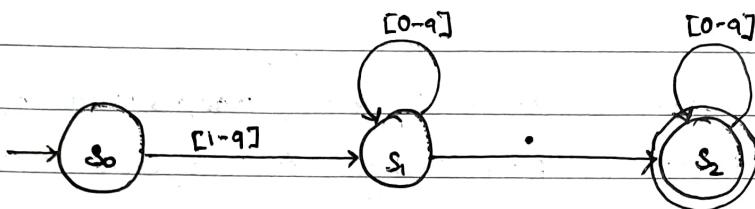
- q given an RE R, build a C function which checks whether an input string is described by the RE R.
1. convert the RE into an NFA.
 2. convert NFA to DFA.
 3. minimize the DFA.
 4. translate the DFA into an equivalent program function.

time complexity: $O(t)$

t = input string length.

ENCODING A DFA INTO A PROGRAM FUNCTION

Ex FLOAT-LITERAL $[1-9][0-9]^*, [0-9]^*$



8/ . 0 1 2 3 4 5 6 7 8 9

s0 s0 s0 s1 s1 s1 s1 s1 s1 s1 s1 s1

s1 s2 s1 s1 s1 s1 s1 s1 s1 s1 s1 s1

s2 s0 s2 s2 s2 s2 s2 s2 s2 s2 s2 s2

s3 s0 s0

note:

we can compress the table by identifying common columns.

TABLE DRIVEN APPROACH FOR ENCODING DFAS

IsFloatLiteral (String str)

Begin

state = s0

for i = 1 to str.length

do

state = nextState (state, str[i])

if (state = s0) then return reject

end for

if (state ∈ F)

then return accept
 else return reject
 end

table

- ① only the transition function, changes used by the next state function changes to recognize another token.
- ② well the set of final states F also changes.

DIRECT CODING APPROACH

isFloatLiteral (String str)

Begin

c = 1

s0: if ($i > \text{str.length}$) then return reject

if ('0' ≤ str[i] ≤ '9') then goto s1

return reject

s1: i = i + 1

if ($i > \text{str.length}$) then return reject

if ('0' ≤ str[i] ≤ '9') then goto s1

if ($\text{str}[i] = \text{!}!$) then goto s2

else return reject

s2: i = i + 1

if ($i > \text{str.length}$) then return accept

if ('0' ≤ str[i] ≤ '9') then goto s2

else return reject

End

TABLE DRIVEN vs DIRECT CODED APPROACH

table driven approach

- one multiplication operation per terminal symbol.
- One branch instruction per terminal. but these branches are highly predictable (loop back).
- good code locality.
- data locality?

direct coded approach

- no multiplication operation.
- one branch instruction per terminal. but predicting the branch targets could be hard.
- code locality?
- data locality?

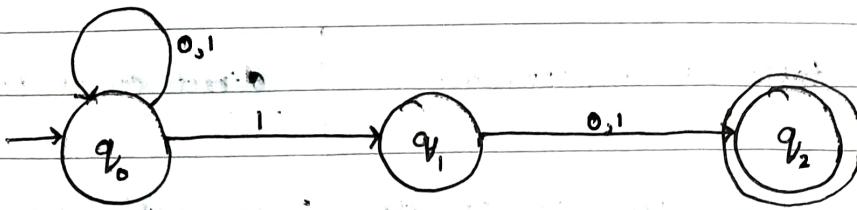
LEXICAL ANALYZER DESIGN

Q given a sequence of REs R_1, R_2, \dots, R_k , construct a function which takes a string as input and outputs the first RE R_i describing the input string.

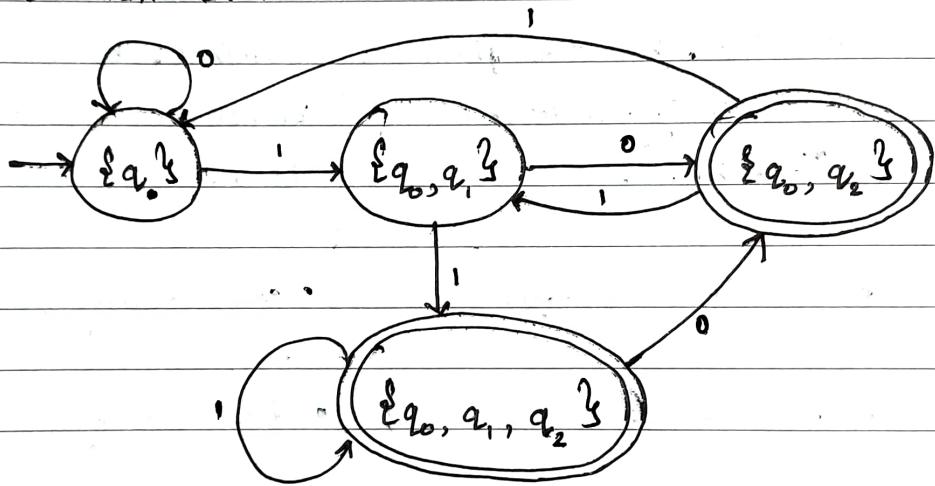
time complexity: $O(kt)$

can we design an $O(t)$ algorithm for the checkString function?

RECALL: NFA TO DFA CONVERSION : ALGORITHM



equivalent DFA



1. convert each of the REs to their respective NFAs.
2. unify the NFAs.
3. convert the unified NFA into a DFA.
4. examine each final state in the DFA (2 possible cases)
 - (a) accepting states are coming from only one NFA.
 - (b) accepting states are coming from different NFAs.
use priority rule.

★ We can use the DFA to do the search over all REs in one shot!

• KEYWORD VS IDENTIFIER

- unify the NFAs for if and Identifier