

dondi.lmu.build/share/pl/syntax-V02.pdf

PROGRAMMING LANGUAGE SYNTAX

- natural language vs programming language.
 - no ambiguity allowed in programming languages - in form (syntax) and meaning (semantics).
 - thus, programming language specifications use formal notations for both syntax and semantics.
- distinction between syntax & semantics:
many programming languages have features that mean the same (shared semantics) but are expressed differently
 - identifying which it's which helps the learning curve.

TWO LEVELS OF SYNTAX

microsyntax

specification

regular expressions

macrosyntax

context-free grammars
expressed in BNF or EBNF

algorithm/
recognition

lexical analysis /
scanning

parsing: LL / top-down
predictive; LR / bottom-up

input

symbol / character stream

token Stream

theoretical
foundation

deterministic finite
automata

deterministic push-down
automata

tools

lex, flex^v

yacc, bison^v

SYNTAX SPECIFICATION

- formalism : set of production rules
- microsyntax rules: concatenation, alternation
(choice among finite alternatives), "kleene closure".
 - the set of strings produced by these 3 rules is a regular set or a regular language.
 - the rules are specified by regular expressions - they "generate" the "regular" language.
 - strings in the regular language are recognized by scanners.
- macrosyntax rules: add recursion
 - the set of strings produced by these rules is a context-free language (CFL).
 - these rules are specified by context-free grammars (CFGs) -

they "generate" the context-free language.

- Strings in the context-free language are recognized by parsers.
- distinction between a full-fledged programming language and a pure formal language:
regular & context-free langs. are purely formal langs - they are just sets of strings, and do not carry meaning.

MICROSYNTAX SPECIFICATION

- a character (in some encoding system:
once ASCII, now unicode).
- the empty string (\emptyset or λ)
- 2 concatenated regular expressions
- 2 regular expressions separated by $|$, denoting a choice between the 2 regexps.
- a regular expression followed by the Kleene star (*), denoting 0 or more instances of that regexp.
- example: numeric literal (unsigned number)

- we use single quotes to identify characters that are expected in the input stream.
- italics indicate regular expressions that are defined elsewhere in the microsyntax

$\text{digit} \rightarrow '0' | '1' | '2' | '3' | '4' | '5' | '6' | '7' | '8' | '9'$

$\text{Unsigned-integer} \rightarrow \text{digit digit}^*$

$\text{Unsigned-number} \rightarrow \text{unsigned-integer}$

$(\text{'.' unsigned-integer}) \mid \epsilon$

$(\text{'e'} (\text{'+'} \mid \text{'-'}) \text{unsigned-integer}) \mid \epsilon$

MICROSYNTAX DESIGN ISSUES

- reserved word or identifier?

reserved words are their own tokens; identifiers get lexed as a single token type with an attached value.

- Case sensitivity in identifiers and keywords?

-

- no: ada, common lisp, fortran 90, pascal

- yes: modula-2/3, C/C++, java, perl, ML, javascript

- other identifier issues!

what characters to accept,

length of identifiers

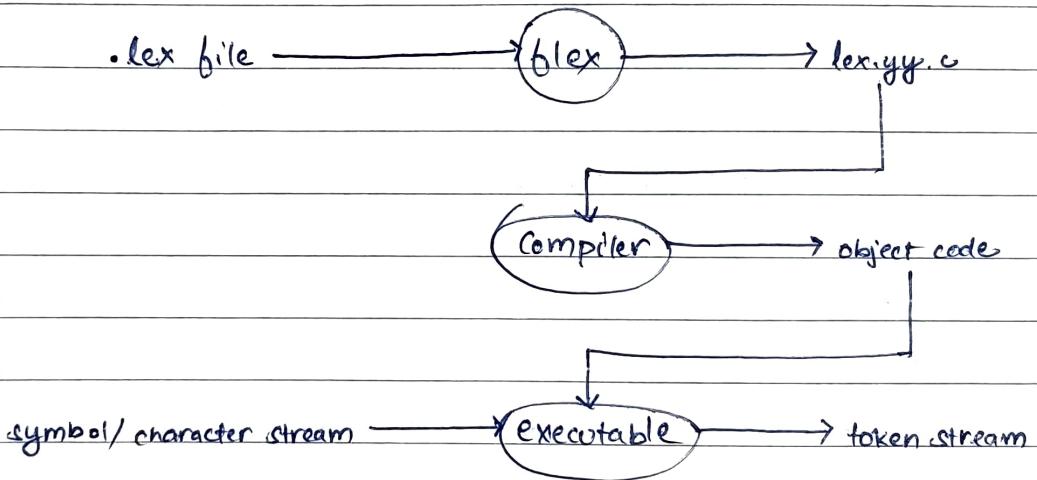
- white space: mostly free format these days — token streams are completely linear, with any amount of white space in between
 - early fortran was not free format: specific character positions had specific functions, with a maximum number of columns (72).
 - some languages allow line breaks to separate statements: haskell, occam, sr
- ditto javascript, but really why? takes semicolons just as well, and makes life simpler.
- comments are processed (essentially, removed) during lexical analysis: can be delimited (nested, or non-nested) or to-end-of-line.

MICROSYNTAX IMPLEMENTATION

- hand coded: e

essentially a "writing out" of a finite state automaton

- details in compiler construction
- really semi-hand coded: lexical analyzers are a well known quantity and follow the same overall pattern regardless of the microsyntax — thus, these days, they are almost always.
- data-driven / table-driven: use a scanner generator; the best known are lex and its newer version, flex.



the token stream provides:

1. what was recognized (the left side of the regular expression) — that is, the token itself.
2. the exact character sequence that was recognized as this token — the token's "value" or "spelling"

RESERVED WORDS VS. IDENTIFIERS

- in Java, `private` is a reserved word, but it is lexically indistinguishable from a variable called, say `large`.
- to handle this, we "cheat" a little bit by maintaining a separate data structure that lists the reserved words in a language; when an "identifier" is found during lexical analysis, it is looked up against the list of known reserved words, and if there is a match, the token for that reserved word is returned instead of the identifier token.
- examples
 - "`500`" is an integer token with value `500`.
 - "`x`" is an identifier token with value "`x`".
 - (in C) "`return`" is a reserved word, so its token is `return` with value as a don't care.
- interesting footnote:
"syntax coloring" features in current programming editors are lexical in nature.

MACROSYNTAX SPECIFICATION

heuristically: "regular expressions with recursions"

- set of productions that define variables or non-terminals in terms of other non terminals and / or terminals (tokens, in the programming languages context).
- the same non-terminal may be defined by more than one production; terminals never appear on the left side of a production.
- a special nonterminal is "blessed" as the start symbol representing the construct defined by the grammar.
- Standard formats:
backus-naur form (BNF) or
extended backus-naur form (EBNF),
named after john backus & peter nau.
- historical tidbit:
BNF first used to specify algol-60
- EBNF is essentially BNF with `l`, `*`, and `()` added.
 - equivalent expressive power with more convenient notation - compare these productions with their strict BNF versions (bold italic == nonterminal; regular italic == token [terminal]):

operator → plus | minus | times | divided-by

identifier-list → identifier (comma identifier)*

MACROSYNTAX DESIGN ISSUES

- we seek a grammar that is unambiguous - given a token stream, we can derive one and only one sequence of production rules that results in that stream from the start symbol.
- general rule:
choose a grammar that reflects the internal structure of the program it produces
- particularly important in arithmetic, with its notions of associativity and precedence of operations.

CASE STUDY: PASCAL'S IF-THEN-ELSE

stmt → if condition then-clause else-clause | other-stmt

then-clause → then stmt

else-clause → else stmt | E

- ambiguous for "if c1 then if c2 then s1 else s2"
- rewrite the grammar (e.g. the productions)?
- implement a disambiguating rule ("the else clause matches the closest unmatched then.")?
- ... or change the syntax? (e.g., explicit end markers such as '3' or addition of an `elsif` keyword)

ANOTHER CASE STUDY: EXPRESSION PARSING

expression → identifier | number | negative expression |

leftparen expression rightparen |

expression operator expression

operator → plus | minus | times | divide

- what derivations in this grammar yield the expression!

Slope * x + intercept

- compare to this grammar, which can parse the same set of token streams:

$\text{expression} \rightarrow \text{term} \mid \text{expression add-op term}$

$\text{term} \rightarrow \text{factor} \mid \text{term mult-op factor}$

$\text{factor} \rightarrow \underline{\text{identifier}} \mid \underline{\text{number}} \mid \underline{\text{negative factor}} \mid$

$\underline{\text{leftparen}} \text{ expression } \underline{\text{rightparen}}$

$\text{add-op} \rightarrow \underline{\text{plus}} \mid \underline{\text{minus}}$

$\text{mult-op} \rightarrow \underline{\text{times}} \mid \underline{\text{divide}}$

- try this for:

$\text{slope} * x + \text{intercept}$

$x - y - z$

MACROSYNTAX IMPLEMENTATION

- when creating a program, we take the start symbol of the language's grammar and progressively replace it with some choice of applicable productions until the result consists entirely of terminals.

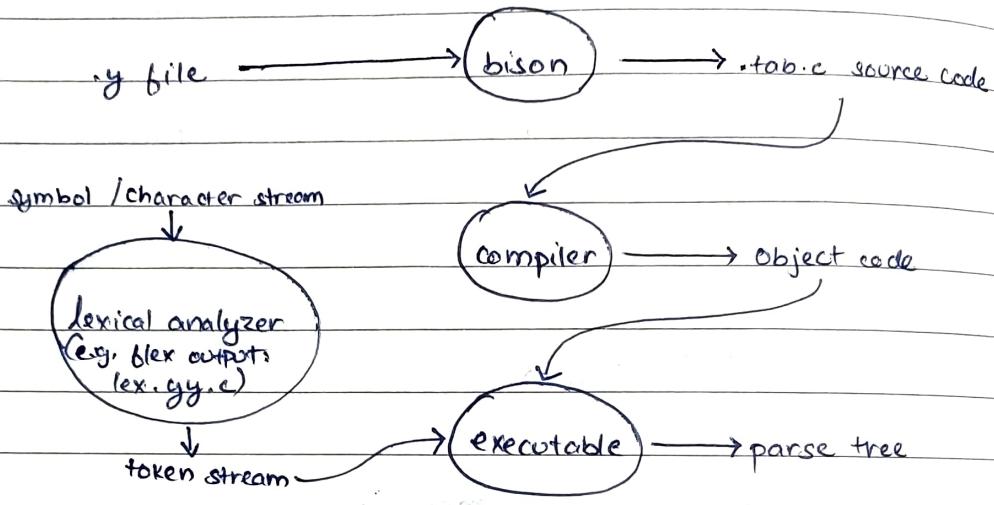
- in reverse (that is, when parsing), we take the stream of terminals (tokens) and determine the sequence of productions that resulted in that stream.
- in general, parsing a context-free grammar is $O(n^3)$.
- 2 CFG categories accommodate $O(n)$ parsing

'L' (L' | R') 'C' n)'

grammar category	meaning	algorithm
LL	left-to-right left-most derivation	top-down, predictive
LR	left-to-right right-most derivation	bottom-up shift-reduce

- we add a number in parentheses to indicate the number of look-ahead tokens required.
- LL grammars allow hand-coded implementation:
top-down / recursive-descent parsers

- LR grammars require a data / table - driven implementation, which use a bottom-up / shift-reduce "parser driver".
 - done by parser generators such as yacc & bison.



- parsing output represents progressively abstract types of data structures, typically best represented a tree (or very similar looking variant).
- in programming languages, the ultimate goal of parser output is an entity that facilitates code-generation & optimization.
- parse trees: a direct mapping from the token stream to the context-free grammar.
- syntax trees: eliminates "helper" tokens and represent the pure syntactic structure of a program.

- abstract syntax trees: static semantics - adds meaning to the symbols of a program, particularly its variables, functions, and other declared entities.

HANDLING SYNTAX ERRORS

- panic mode:

ignore all until we hit a "safe" token.

- phrase-level recovery:

different "safe" sets depending upon the current production.

- implemented with first and follow sets.

- context-sensitive lookahead:

refines follow sets to be context specific.

- historical tidbit:

first attempted by coirth for pascal.

- exception based recovery:

register exception handlers as parser moves through the code; a syntax error "unwinds" to the most recent exception handler.

- error productions:

explicitly define productions in the grammar that represent common or likely syntax errors; when the parser recognizes

that production, a very specific error or warning message can be displayed (and corrections even suggested).