

コンパイラ理論 3 字句解析と構文解析

櫻井彰人

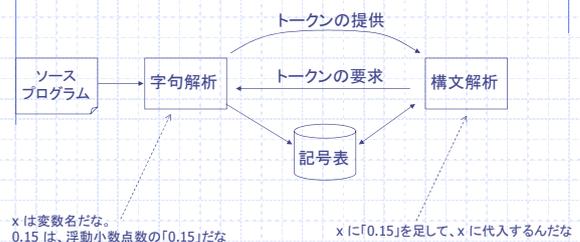
プログラミング言語の定義

- ◆プログラミング言語は次の2階層で定義される。
 - 語彙(単語)とはどういうものか
 - (定義された語彙を用いて)どういう風にプログラムを書くのか?
- 一気に(つまり、一階層で)定義することもできるのだが、分かり難くなる。

コンパイラの構造も

- ◆その2階層に対応するように、コンパイラも作られる
 - 語彙の認識 ⇔ 語彙の定義
 - 字句解析という
 - 文字の並びをみて、単語に区切る
 - 例えば、変数と定数とを区別する
 - プログラムの認識 ⇔ プログラムの定義
 - 構文解析という
 - 例えば、これは「変数xと変数yの値を足してzに代入することだな」と認識する

字句解析と構文解析



なぜ分けるか?

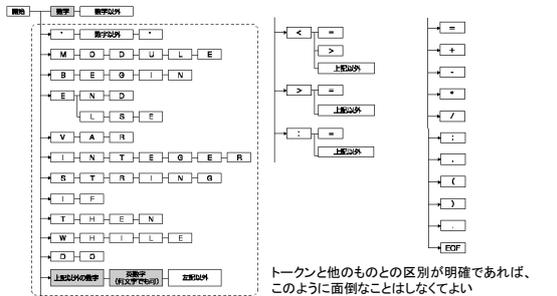
- ◆字句解析を構文解析から分ける理由:
 - 「字句」の定義は、正規文法でできる
 - 簡単な道具で済ませられるところは、簡単にすませよう!
 - 設計が単純になる
 - 効率(速度等)の向上が図れる
 - 可搬性がます
- ◆字句解析・構文解析それぞれによりツールが存在する

トークン・字句・パターン (Tokens, Lexemes, Patterns)

- ◆トークンは、キーワード(if, for, long,...)、演算子(+, *, ...)、識別子、定数、文字列、区切り記号を含む、字句が属すクラスのことをいう
- ◆字句は、文字のある列であって、ソースプログラム内で意味をもつ最小の単位
- ◆パターンは、(Lexで用いるが)あるトークンの生成規則

トークン識別チャート

- 構文解析プログラムが書きやすいようにチャート化



言語の規則の階層1 – 字句

- 字句解析用
- パーサでは使用しない。定義済みとして考える。

<ident> → <letter> (<letter> | <digit>)
 <integer> → <digit>+
 <string> → "<any character except EOF, EOL and " >"

<digit> → 0|1|2|3|4|5|6|7|8|9
 <letter> → a|b|c|d|e|f|g|h|i|j|k|l|m|n|o|p|q|r|s|t|u|v|w|x|y|z|A|B|C|D|E|F|G|H|I|J|K|L|M|N|O|P|Q|R|S|T|U|V|W|X|Y|Z

言語規則の階層2 – 構文

<program> → MODULE IDENT : [<declist>] BEGIN <statlist> END IDENT .
 <declist> → VAR (<identlist> : <type> ;)+
 <statlist> → <statement> (; <statement>)
 <identlist> → IDENT (, IDENT)
 <type> → INTEGER | STRING
 <statement> → IDENT := <expression>
 | IF <relation> THEN <statlist> [ELSE <statlist>] END
 | WHILE <relation> DO <statlist> END
 | IDENT "(" <literal> ")"
 <relation> → <expression> <rel op> <expression>
 <expression> → [<unary op>] <term> (<add op> [<unary op>] <term>)
 <term> → <factor> (<mul op> <factor>)
 <factor> → <literal> | "(" <expression>)
 <literal> → IDENT | NUMBER | STR
 <rel op> → = | < | <= | <> | > | >=
 <unary op> → + | -
 <add op> → + | -
 <mul op> → * | /