

Computer Science 420: TruPL - The **T**ru**m**an **P**rogramming Language, Lexical Conventions

This list enumerates the legal lexemes of TruPL. The input alphabet consists of all lowercase alphabetic ASCII characters, the digits '0' through '9', and the following non-alphanumeric characters:

`; : () , = < > + - * / # space tab newline`

All other characters are illegal in TruPL lexemes.

We use the following shorthand notation when describing the lexemes: `[a..z]` is a regular expression that matches any single lowercase ASCII character. The regular expression `[0..9]` matches any digit.

Lexeme	Regular Expression for Lexeme	Attribute
keyword	<i>program</i> <i>procedure</i> <i>int</i> <i>bool</i> <i>begin</i> <i>end</i> <i>if</i> <i>then</i> <i>else</i> <i>while</i> <i>loop</i> <i>print</i> <i>not</i>	<code>kw_program</code> <code>kw_procedure</code> <code>kw_int</code> <code>kw_bool</code> <code>kw_begin</code> <code>kw_end</code> <code>kw_if</code> <code>kw_then</code> <code>kw_else</code> <code>kw_while</code> <code>kw_loop</code> <code>kw_print</code> <code>kw_not</code>
punctuation	<code>;</code> <code>:</code> <code>,</code> <code>:=</code> <code>(</code> <code>)</code>	<code>punc_semi</code> <code>punc_colon</code> <code>punc_comma</code> <code>punc_assign</code> <code>punc_open</code> <code>punc_close</code>
relational operator	<code>=</code> <code><></code> <code>></code>	<code>relop_eq</code> <code>relop_neq</code> <code>relop_gt</code>

	>= < <=	relop_ge relop_lt relop_le
additive operator	+ - or	addop_add addop_sub addop_or
multiplicative operator	* / and	mulop_mul mulop_div mulop_and
identifier	[a..z] ([a..z] [0..9])*	the identifier name
num	[0..9] [0..9]*	the number

Assumptions:

- Comments begin with a '#' and continue until the end of the current line. Any ASCII character may appear in a comment. Your lexical analyzer should skip over all comments.
- Whitespace consists of one or more repetitions of the space character, the tab character ('\t') or the newline character ('\n'). Whitespace may not appear inside any token.
- If a keyword, identifier, **and**, or **or** is followed by another keyword, identifier, num, **and**, or **or**, then the tokens will be separated by whitespace. Thus, "**and**nota23" should be recognized as a single token of type **identifier**, but "**and** nota23" should be recognized as two tokens, and "**1and not a23**" should be recognized as four.
- If two relational operators appear next to each other, the two tokens will be separated by whitespace. Thus, ">=" should be recognized as a single token, but "> =" should be recognized as two tokens.
- If the : token is followed by the = token, the two tokens will be separated by whitespace. ":@" should be recognized as the **punctuation** token.