

Source:

1 | The Role of the Lexical Analyzer

1.1 | Lexical vs Syntactic analysis

1. Syntax and raw text are different and can be treated separately
2. it may be more efficient
3. better portability

1.2 | terms: tokens, patterns, lexemes

- #definition token: is a name and a value, where the name like a keyword or an identifier and the value is a section of the source text?
- #definition pattern: basically a regex of what string structures are allowed
- #definition lexeme: part of the source text that is matched by a pattern as an instance of a token

1.3 | common token breakdown

1. keywords (usually one per keyword)
2. operators (sometimes in operator classes)
3. identifiers
4. constants (sometimes one per type)
5. punctuation (usually one per each, including parens, comma, and semicolon)

1.4 | token attributes

- Token name only contains what type of token it is, not the value
 - if the token is "number", then what number actually was it?
- "token name influences parsing decisions, while the attribute value influences translation of tokens after the parse."
- the identifier token **id** needs to associate lots of data, such as it's lexeme, type, and location in memory, etc

1.5 | lexical errors

Sometimes we can modify the source to attempt to fix typos, etc. Such as removing some letters, edit distance, etc.

2 | Input Buffering

#todo-learn

3 | specification of tokens

3.1 | strings and languages (many definitions)

3.1.1 | #definition alphabet

a set of characters. examples include the *binary alphabet* $\{0, 1\}$, ASCII, and Unicode

3.1.2 | #definition string

a string over an alphabet is a "finite sequence of symbols" from that alphabet. It's length $|s|$ is the number of symbols in s . ϵ is the empty string.

3.1.3 | #definition language

countable set of strings over some fixed alphabet. Some languages are abstract, like ϵ or ϵ are boring languages. Also included are the set of C programs and valid english sentences.

3.2 | operations on languages

3.2.1 | union $L \cup M$

standard set union

3.2.2 | concatenation LM

set of pairwise concatenations (anything from the first concat anything from the second)

3.2.3 | Kleene closure L^*

concatenate L zero or more times. $L^0 = \{\epsilon\}$ and $L^n = L^{n-1}L$.

3.2.4 | Positive closure

Kleene closure, but without L^0 .

3.3 | Regular Expressions

this syntax is a little different from "modern" regexes: the vertical bar $|$ represents union instead of "or".

3.3.1 | #definition regular expression (inductive)

1. inductive basis

- (a) ϵ is a regular expression and it's language $L(\epsilon) = \{\epsilon\}$.
- (b) If a is a symbol in the alphabet Σ then a is a regular expression and $L(a) = \{a\}$ (strings of length 1 that are "a").

2. inductive induction (lol)

- (a) union '|'
- (b) concat
- (c) kleene closure
- (d) parens (don't change the value of the internal expression, just used to group things)

3.3.2 | **for ergonomics**

Everything is left associative

- 1. Unary operator * has highest precedence
- 2. concat has second highest precedence
- 3. '|' has lowest precedence

3.3.3 | **#definition regular set**

Any language that can be defined by a regular expression

3.3.4 | **#definition equivalent** $r = s$

If two regular expressions denote the same regular set.

3.4 | **regular definitions**