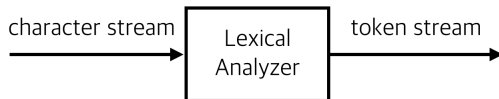


EC3204: Programming Languages and Compilers

Lecture 2 — Lexical Analysis (1) *Specification of Tokens*

Sunbeom So
Fall 2024

Roadmap for Building a Lexical Analyzer



- ① **Specification:** how to express valid lexeme patterns?
 - ▶ Valid C identifiers are strings like `x`, `xy`, `match0`, and `_abc`.⇒ regular expressions
- ② **Recognition:** how to accept (resp., reject) the valid (resp., invalid) lexical patterns?
 - ▶ Accept `match0` as a C identifier, and reject `0match`.⇒ deterministic finite automata (DFA)
- ③ **Implementation:** how to implement such a string recognizer?
 - ⇒ Thompson's construction and subset construction

This Lecture: Specification of Tokens

- Preliminaries for defining a regular expression
 - ▶ alphabets, strings, languages
- Syntax and semantics of regular expressions
- Extensions of regular expressions to simplify the specification.

Preliminary: Alphabet

An alphabet Σ is a finite, non-empty set of symbols.

- $\Sigma = \{0, 1\}$ is a binary alphabet.
- $\Sigma = \{A, B, \dots, Z\}$ is the English alphabet.
- Realistic examples: ASCII, Unicode

ASCII printable characters (character code 32-127)

Codes 32-127 are common for all the different variations of the ASCII table, they are called printable characters, represent letters, digits, punctuation marks, and a few miscellaneous symbols. You will find almost every character on your keyboard. Character 127 represents the command DEL.

DEC	OCT	HEX	BIN	Symbol	HTML Number	HTML Name	Description
32	040	20	00100000	SP	 		Space
33	041	21	00100001	!	!	!	Exclamation mark
34	042	22	00100010	"	"	"	Double quotes (or speech marks)
35	043	23	00100011	#	#	#	Number sign
36	044	24	00100100	\$	$	$	Dollar
37	045	25	00100101	%	%	&percent;	Per cent sign
38	046	26	00100110	&	&	&	Ampersand
39	047	27	00100111	'	'	'	Single quote
40	050	28	00101000	((&lparen;	Open parenthesis (or open bracket)
41	051	29	00101001))	&rparen;	Close parenthesis (or close bracket)
42	052	2A	00101010	*	*	*	Asterisk
43	053	2B	00101011	+	+	+	Plus
44	054	2C	00101100	,	,	,	Comma
45	055	2D	00101101	-	-		Hyphen-minus
46	056	2E	00101110	.	.	.	Period, dot or full stop
47	057	2F	00101111	/	/	/	Slash or divide
48	060	30	00110000	0	0		Zero
49	061	31	00110001	1	1		One
50	062	32	00110010	2	2		Two

(captured from <https://www.ascii-code.com>)

Preliminary: String (Sentence, Word)

A string is a finite sequence of symbols from an alphabet. For example, **1**, **01**, **10110** are strings over $\Sigma = \{0, 1\}$.

The following is the selected list of notations about strings.

- ϵ : the empty string.
- wv : the concatenation of w and v .
 - ▶ Given $w = \text{dog}$ and $v = \text{house}$, $wv = \text{doghouse}$.
 - ▶ The empty string is the identity under the concatenation; for any string s , $\epsilon s = s\epsilon = s$.
- $|w|$: the length of string w .
 - ▶ $|\epsilon| = 0$, $|sa| = |s| + 1$ for some $a \in \Sigma$.
- Regarding concatenation as a product, we can define the “exponentiation” (or “power”) of strings.
 - ▶ Definition: $s^0 = \epsilon$, $s^i = s^{i-1}s$ for all $i > 0$.
 - ▶ $s^1 = s$ since $s^1 = s^0s = \epsilon s = s$.
 - ▶ $s^2 = ss$, $s^3 = sss$, and so on.

Preliminary: Language

A language L is a subset of strings over some alphabet Σ .

The following is the selected list of operations on languages.

- Union: $L_1 \cup L_2 = \{x \mid x \in L_1 \vee x \in L_2\}$
- Concatenation: $L_1 L_2 = \{xy \mid x \in L_1 \wedge y \in L_2\}$
- The power of a language L , denoted L^n :

$$\begin{aligned} L^0 &= \{\epsilon\} \\ L^n &= L^{n-1} L \end{aligned}$$

- The (Kleene) closure of a language L , denoted L^* .

$$L^* = L^0 \cup L^1 \cup L^2 \cup \dots = \bigcup_{i \geq 0} L^i$$

In words: the set of strings obtained by concatenating of L zero or more times.

- The positive closure of a language L , denoted L^+ :

$$L^+ = L^1 \cup L^2 \cup L^3 \cup \dots = \bigcup_{i \geq 1} L^i$$

Regular Expression

A regular expression is a notation to denote a particular type of language.

- Syntax is inductively defined using the production rule.
 - ▶ **Base cases** define the primitive regular expressions.
 - ▶ **Inductive cases** define larger regular expressions from smaller ones.

$$\begin{array}{ll} R \rightarrow \emptyset \mid \epsilon \mid a \in \Sigma & \text{(base cases)} \\ \mid R_1 \mid R_2 \mid R_1 \cdot R_2 \mid R_1^* \mid (R) & \text{(inductive cases)} \end{array}$$

- Semantics (meaning) is also inductively defined.

$$\begin{aligned} L(\emptyset) &= \emptyset \\ L(\epsilon) &= \{\epsilon\} \\ L(a) &= \{a\} \\ L(R_1 \mid R_2) &= L(R_1) \cup L(R_2) \\ L(R_1 \cdot R_2) &= L(R_1)L(R_2) \\ L(R^*) &= (L(R))^* \\ L((R)) &= L(R) \end{aligned}$$

Example

$$\begin{aligned}L(\emptyset) &= \emptyset \\L(\epsilon) &= \{\epsilon\} \\L(a) &= \{a\} \\L(R_1 \mid R_2) &= L(R_1) \cup L(R_2) \\L(R_1 \cdot R_2) &= L(R_1)L(R_2) \\L(R^*) &= (L(R))^* \\L((R)) &= L(R)\end{aligned}$$

Using the semantics above, compute the strings expressed by $a^* \cdot (a \mid b)$.

$$\begin{aligned}L(a^* \cdot (a \mid b)) &= L(a^*)L(a \mid b) \\&= (L(a))^*(L(a) \cup L(b)) \\&= (\{a\})^*(\{a\} \cup \{b\}) \\&= \{\epsilon, a, aa, aaa, \dots\}(\{a, b\}) \\&= \{a, aa, aaa, \dots, b, ab, aab, \dots\}\end{aligned}$$

Exercises

Write regular expressions for the following languages.

- The set of all strings over $\Sigma = \{a, b\}$.
▶
- The set of strings of a 's and b 's, terminated by ab .
▶
- The set of strings with an even number of a 's followed by an odd number of b 's.
▶
- The set of C identifiers.
▶

Exercises

Write regular expressions for the following languages.

- The set of all strings over $\Sigma = \{a, b\}$.
 - ▶ $(a \mid b)^*$
- The set of strings of a 's and b 's, terminated by ab .
 - ▶ $(a \mid b)^*ab$
- The set of strings with an even number of a 's followed by an odd number of b 's.
 - ▶ $(a \cdot a)^*(b \cdot b)^*b$
- The set of C identifiers.
 - ▶ $(A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z \mid _) \cdot (A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z \mid _ \mid 0 \mid 1 \mid \dots \mid 9)^*$

Alternatively, using **regular definition**,

$$\begin{aligned} \textit{letter} &\rightarrow A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z \mid _ \\ \textit{digit} &\rightarrow 0 \mid 1 \mid \dots \mid 9 \\ \textit{id} &\rightarrow \textit{letter}(\textit{letter} \mid \textit{digit})^* \end{aligned}$$

Regular Definitions

For notational convenience, we give names to regular expressions and use the names in subsequent expressions. A **regular definition** is a sequence of definitions of the form:

$$\begin{array}{lcl} d_1 & \rightarrow & r_1 \\ d_2 & \rightarrow & r_2 \\ & \dots & \\ d_n & \rightarrow & r_n \end{array}$$

- 1 Each d_i is a new name such that $d_i \notin \Sigma$.
- 2 Each r_i is a regular expression over $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$.
That is, we restrict r_i to Σ and previously defined regular definitions.

Example

- Unsigned numbers (integers or floating point), e.g., 5280, 0.01234, 6.336E4, or 1.89E-4:

<i>digit</i>	→	0 1 ... 9
<i>digits</i>	→	<i>digit digit*</i>
<i>optionalFraction</i>	→	. <i>digits</i> ϵ
<i>optionalExponent</i>	→	(E (+ - ϵ) <i>digits</i>) ϵ
<i>number</i>	→	<i>digits optionalFraction optionalExponent</i>

Extensions of Regular Expressions

Extensions are used to simplify the specification.

- R^+ : the positive closure of R , i.e., $L(R^+) = L(R)^+$.
- $R?$: zero or one instance of R , i.e., $L(R?) = L(R) \cup \{\epsilon\}$.

To represent character classes, we use the following extensions.

- $[a_1 a_2 \cdots a_n]$: the shorthand for $a_1 \mid a_2 \mid \cdots \mid a_n$.
- $[a_1 - a_n]$: the shorthand for $[a_1 a_2 \cdots a_n]$, where a_1, \dots, a_n are consecutive symbols.
 - ▶ $[abc] = a \mid b \mid c$
 - ▶ $[a-z] = a \mid b \mid \cdots \mid z$.

Examples

- C identifiers:

$letter \rightarrow [A-Za-z_]$
 $digit \rightarrow [0-9]$
 $id \rightarrow letter (letter|digit)^*$

- Unsigned numbers:

$digit \rightarrow [0-9]$
 $digits \rightarrow digit^+$
 $number \rightarrow digits (. digits)? (E [+ -]? digits)?$

Summary

We learned how to express valid lexeme patterns (i.e., define valid inputs of lexical analyzers).

- Preliminaries for defining a regular expression
 - ▶ alphabets, strings, languages
- Syntax and semantics of regular expressions
- Extensions of regular expressions to simplify the specification.

Next class: finite automata for accepting (resp., rejecting) the valid (resp., invalid) lexical patterns.