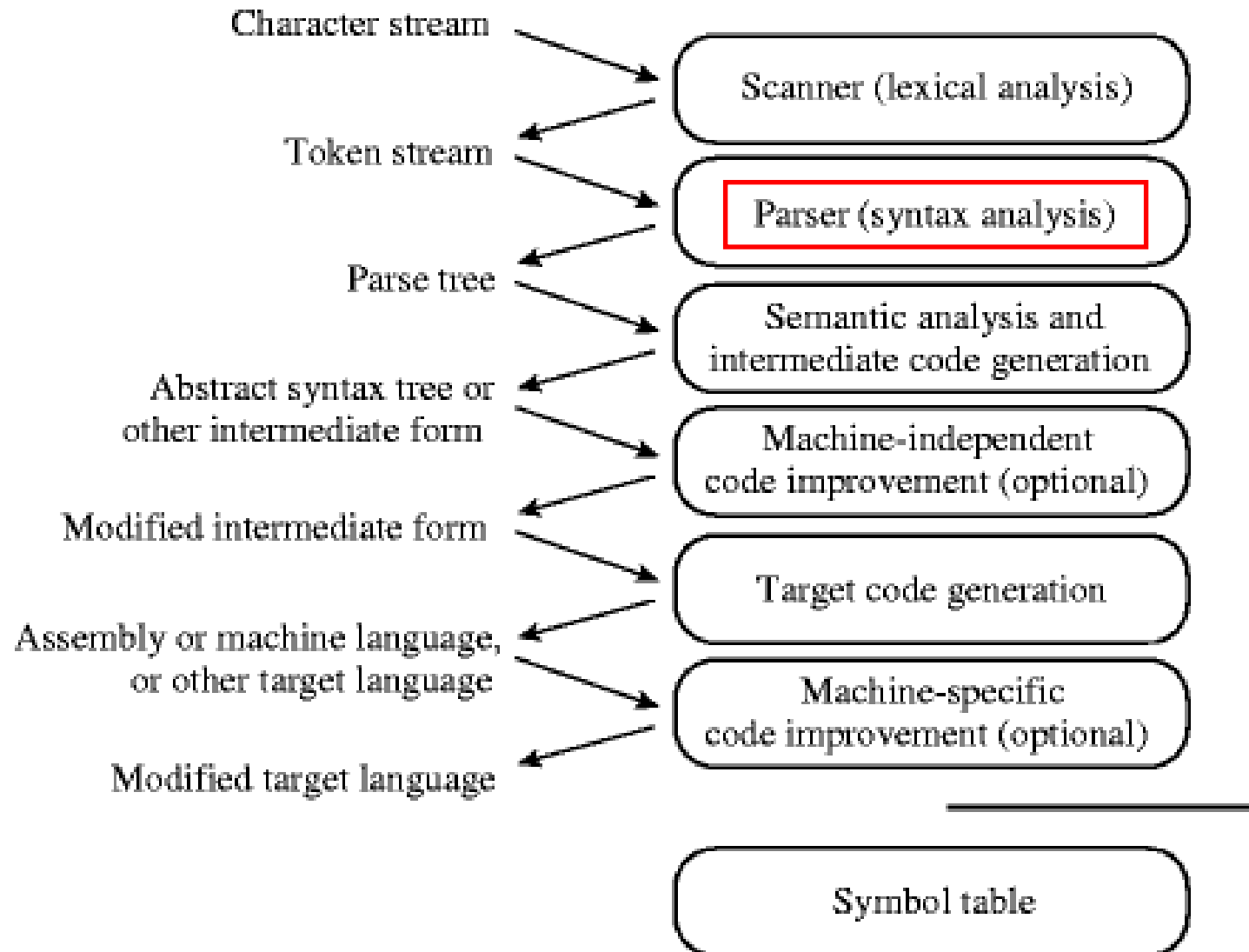

Syntax Specification

Regular Expressions

Phases of Compilation



Syntax Analysis

- Syntax:
 - Webster's definition: *1 a : the way in which linguistic elements (as words) are put together to form constituents (as phrases or clauses)*
- The syntax of a programming language
 - Describes its form
 - » *i.e.* **Organization of tokens (elements)**
 - Formal notation
 - » Context Free Grammars (CFGs)

Review: Formal definition of tokens

- A set of tokens is a set of strings over an alphabet
 - {read, write, +, -, *, /, :=, 1, 2, ..., 10, ..., 3.45e-3, ...}
- A set of tokens is a *regular set* that can be defined by comprehension using a *regular expression*
- For every regular set, there is a *deterministic finite automaton* (DFA) that can recognize it
 - *i.e.* determine whether a string belongs to the set or not
 - Scanners extract tokens from source code in the same way DFAs determine membership

Regular Expressions

- A regular expression (RE) is:
 - A single character
 - The empty string, ϵ
 - The concatenation of two regular expressions
 - » *Notation:* $RE_1 RE_2$ (i.e. RE_1 followed by RE_2)
 - The union of two regular expressions
 - » *Notation:* $RE_1 | RE_2$
 - The closure of a regular expression
 - » *Notation:* RE^*
 - » $*$ is known as the *Kleene star*
 - » $*$ represents the concatenation of 0 or more strings
 - *Non-null enumeration*
 - » *Notation:* RE^+
 - » represents all non-null concatenations of RE (1 or more times)

Regular Expressions Basics

Let alphabet $\Sigma = \{a, b\}$ (means a and b are its only letters)

$$a^* = (\lambda, a, aa, aaa, \dots)$$

$$(ab)^* = (\lambda, ab, abab, ababab, \dots)$$

$$a \cup b = (a, \lambda, b, bb, bb, \dots)$$

$$(a \cup b)^* = \text{all strings containing a's and b's}$$

$$(a^*b^*)^* = (ab^*)^* = \text{all strings containing a's and b's}$$

$$a^*b^* = \{a^i b^j \mid i \geq 0, j \geq 0\}$$

Building Regular Expressions

Regular Expressions as Language

- $*$ while loop
 - iterates 0 or more times
- concatenation uv
 - sequential; first u , then v
- $u \cup v$ OR
 - select from one or the other or both

Description → Regular Expression

Let $\Sigma = \{a, b\}$ – all expressions over this alphabet

Strings with

- exactly one a b^*ab^*
- exactly two a's $b^*ab^*ab^*$
- one or more a's $(b^*ab^*)^*$ or $(a \cup b)^*a(a \cup b)^*$
- even number of a's $(b^*ab^*ab^*)^*$
- even number of a's and exactly one b
 $(aa)^*b(aa)^* \cup (aa)^*ab(aa)^*a$
- odd number of a's $(b^*ab^*ab^*)^*b^*ab^*$
- that don't contain aa $(b \cup ab)^*(\lambda \cup a)$

Regular Expression \rightarrow Description

Same alphabet

- $(aa)^*$ even number of a's
- $(a \cup b)(a \cup b)(a \cup b)(a \cup b)$

all strings of length 4

- $((a \cup b)(a \cup b)(a \cup b)(a \cup b))^*$

strings of length divisible by 4

- $(aa)^* \cap ((a \cup b)(a \cup b)(a \cup b)(a \cup b))^*$

strings of a's of length
divisible by 4

Token Definition Example

- Numeric literals in Pascal
 - Definition of the token *unsigned_number*

$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

$unsigned_integer \rightarrow digit\ digit^* \mid digit^+$

$unsigned_number \rightarrow unsigned_integer\ (\mid \cdot\ unsigned_integer) \mid \epsilon$
 $(\mid e\ (\mid + \mid - \mid \epsilon)\ unsigned_integer) \mid \epsilon$

- **Recursion is not allowed in Regular Expressions!**

Exercise

$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

$unsigned_integer \rightarrow digit\ digit^*$

$unsigned_number \rightarrow unsigned_integer\ ((.\ unsigned_integer) \mid \epsilon)$
 $((e\ (+ \mid - \mid \epsilon)\ unsigned_integer) \mid \epsilon)$

- Regular expression for
 - Decimal numbers

$number \rightarrow \dots$

Exercise

$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

$unsigned_integer \rightarrow digit\ digit^*$

$unsigned_number \rightarrow unsigned_integer\ ((.\ unsigned_integer) \mid \epsilon)$
 $((e\ (+ \mid - \mid \epsilon)\ unsigned_integer) \mid \epsilon)$

- Regular expression for
 - Decimal numbers

$number \rightarrow (+ \mid - \mid \epsilon)\ unsigned_integer\ ((.\ unsigned_integer) \mid \epsilon)$

Exercise

$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

$unsigned_integer \rightarrow digit\ digit^*$

$unsigned_number \rightarrow unsigned_integer\ ((.\ unsigned_integer) \mid \varepsilon)$
 $((e\ (+ \mid - \mid \varepsilon)\ unsigned_integer) \mid \varepsilon)$

- Regular expression for
 - Identifiers

$identifier \rightarrow \dots$

Exercise

$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

$unsigned_integer \rightarrow digit\ digit^*$

$unsigned_number \rightarrow unsigned_integer\ ((.\ unsigned_integer) \mid \epsilon)$
 $((e\ (+ \mid - \mid \epsilon)\ unsigned_integer) \mid \epsilon)$

- Regular expression for
 - Identifiers

$identifier \rightarrow letter\ (letter \mid digit \mid _)^*$

$letter \rightarrow a \mid b \mid c \mid \dots \mid z$