

## FORMAL LANGUAGES

### Alphabets and Strings

- An **alphabet**  $\Sigma$  is a finite set of **characters** (or **symbols**).
- A **word**, or **sequence**, or **string over**  $\Sigma$  is any group of 0 or more consecutive characters of  $\Sigma$ .
- The **length** of a word is the number of characters in the word.
- The **null string** is the string of length 0. It is denoted  $\epsilon$  or  $\lambda$ .
- A string of length  $n$  is really an ordered  $n$ -tuple of characters written without parentheses or commas.
- Given two strings  $x$  and  $y$  over  $\Sigma$ , the **concatenation** of  $x$  and  $y$  is the string  $xy$  obtained by putting all the characters of  $y$  right after  $x$ .

### Languages over an alphabet

Let  $\Sigma$  be an alphabet. A **formal language over**  $\Sigma$  is a set of strings over  $\Sigma$ .

- $\emptyset$  is the **empty language** (over  $\Sigma$ )
- $\Sigma^n = \{\text{all strings over } \Sigma \text{ that have length } n\}$  where  $n \in \mathbb{N}$
- $\Sigma^+ =$  the **positive closure** of  $\Sigma = \{\text{all strings over } \Sigma \text{ that have length } \geq 1\}$
- $\Sigma^* =$  the **Kleene closure** of  $\Sigma = \{\text{all strings over } \Sigma\}$

### Operations on Languages

Let  $\Sigma$  be an alphabet. Let  $L$  and  $L'$  be two languages defined over  $\Sigma$ .

The following operations define new languages over  $\Sigma$ :

- The **concatenation of  $L$  and  $L'$** , denoted  $LL'$ , is  $LL' = \{xy \mid x \in L \wedge y \in L'\}$
- The **union of  $L$  and  $L'$** , denoted  $L \cup L'$ , is  $L \cup L' = \{x \mid x \in L \vee x \in L'\}$
- The **Kleene closure of  $L$** , denoted  $L^*$ , is  $L^* = \{x \mid x \text{ is a concatenation of any finite number of strings in } L\}$ . Note that  $\epsilon \in L^*$ .

## REGULAR EXPRESSIONS

### Definition

Let  $\Sigma$  be an alphabet. The following are **regular expressions (r.e.)** over  $\Sigma$ :

- I. **BASE**:  $\epsilon$  and each individual symbol of  $\Sigma$  are regular expressions.
- II. **RECURSION**: if  $r$  and  $s$  are regular expressions over  $\Sigma$ , then the following are also regular expressions over  $\Sigma$ :
  - $(rs)$  the concatenation of  $r$  and  $s$
  - $(r | s)$   $r$  or  $s$
  - $(r^*)$  the Kleene closure of  $r$
- III. **RESTRICTION**: The only regular expressions over  $\Sigma$  are the ones defined by I and II above.

### Order of Precedence of Regular Expression Operations

- The order of precedence of r.e. operators are, from highest to lowest:
- Highest:  $()$   $*$  concatenation  $|$  : lowest

### Languages Defined by Regular Expressions

Let  $\Sigma$  be an alphabet. Define a function  $L$  as follows:

$$L: \begin{cases} \{\text{all r.e.'s over } \Sigma\} & \rightarrow \{\text{all languages over } \Sigma\} \\ r & \mapsto L(r) = \text{the language defined by } r \end{cases}$$

- I.  $L(\epsilon) = \{\epsilon\}, \forall a \in \Sigma L(a) = \{a\}$
- II. **RECURSION**: If  $L(r)$  and  $L(s)$  are the languages defined by the regular expressions  $r$  and  $s$  over  $\Sigma$ , then
  - $L(rs) = L(r)L(s)$
  - $L(r|s) = L(r) \cup L(s)$
  - $L(r^*) = (L(r))^*$

### Variations

Some definitions of regular expressions and regular languages define  $\emptyset$  to be a r.e. with  $L(\emptyset) = \emptyset$

## PROPERTIES OF REGULAR EXPRESSIONS

Regular expressions can be simplified by applying the following properties:

For any regular expressions  $r, s, t$ ,

Axiom	Description
$r   s = s   r$	$ $ is commutative
$r   (s   t) = (r   s)   t = r   s   t$	$ $ is associative
$(rs)t = r(st) = rst$	Concatenation is associative
$r(s t) = rs   rt$ and $(s t)r = sr   tr$	Concatenation is distributive over $ $
$r\epsilon = \epsilon r = r$	$\epsilon$ is the identity element for concatenation
$r^{**} = r^*$	$*$ is idempotent
$r^* = (r \epsilon)^*$	

## NOTATIONAL SHORTHANDS

Here are some frequent constructs which have their own notation:

- $(r)^+$  means one or more instances of  $r$ .

$$L((r)^+) = (L(r))^+$$

- $(r)?$  means 0 or 1 instances of  $r$ . i.e.  $(r)? = r|\epsilon$

$$L((r)?) = (L(r|\epsilon)) = L(r) \cup L(\epsilon) = L(r) \cup \{\epsilon\}$$

- Character classes:

$$[abc] = a|b|c$$

$$[a-z] = a|b|\dots|z$$

## REGULAR DEFINITIONS

Regular expressions can be broken down into **regular definitions**: sequences of expressions of the form

$$d_1 \rightarrow r_1$$

...

$$d_n \rightarrow r_n$$

where each  $d_i$  is a distinct name and

$r_i$  is a regular expression over symbols in  $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$