

# ECE 374 B Language Theory: Cheatsheet

## 1 Languages and strings

### Languages

- An *alphabet*  $\Sigma$  is a **finite** set of symbols.

**Definitions** A *string* in  $\Sigma^*$  is a **finite** sequence of symbols in  $\Sigma$ .

- A *language* is  $L$  is a set of strings over some alphabet.

All languages represent mathematical problems.  
 Example: multiplication of two integers:

$$L_{MULT_2} = \left\{ \begin{array}{ccc} 1 \times 1|1, & 1 \times 2|2, & 1 \times 3|3, \dots \\ 2 \times 1|2, & 2 \times 2|4, & 2 \times 3|6, \dots \\ \vdots & \vdots & \vdots \\ n \times 1|n, & n \times 2|2n, & n \times 3|3n, \dots \end{array} \right\} \quad (1)$$

- For languages  $A, B$  the *concatenation* of  $A, B$  is  $AB = \{xy \mid x \in A, y \in B\}$ .
- For languages  $A, B$ , their *union* is  $A \cup B$ , *intersection* is  $A \cap B$ , and *difference* is  $A \setminus B$  (also written as  $A - B$ ).
- For language  $A \subseteq \Sigma^*$  the *complement* of  $A$  is  $\bar{A} = \Sigma^* \setminus A$ .
- $\Sigma^n$  is the set of all strings of length  $n$ .
- $\Sigma^* = \cup_{n \geq 0} \Sigma^n$  is the set of all strings over  $\Sigma$ .
- $\Sigma^+ = \cup_{n \geq 1} \Sigma^n$  is the set of non-empty strings over  $\Sigma$ .

### Language operations

### Strings

- The *length* of a string  $w$  (denoted by  $|w|$ ) is the number of symbols in  $w$ .
- For integer  $n \geq 0$ ,  $\Sigma^n$  is set of all strings over  $\Sigma$  of length  $n$ .  $\Sigma^*$  is the set of all strings over  $\Sigma$ .
- $\Sigma^*$  is the set of all strings of all lengths including empty string.
- $\epsilon$  is a *string* containing no symbols.
- $\emptyset$  is the *empty set*. It contains no strings.

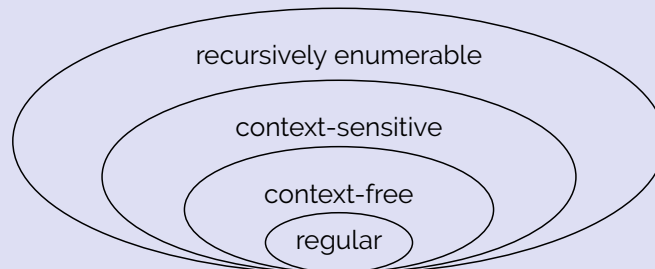
### Definitions

- If  $x$  and  $y$  are strings then  $xy$  denotes their concatenation. Recursively:
  - $xy = y$  if  $x = \epsilon$
  - $xy = a(wy)$  if  $x = aw$
- $v$  is *substring* of  $w \iff$  there exist strings  $x, y$  such that  $w = xvy$ .
  - If  $x = \epsilon$  then  $v$  is a *prefix* of  $w$
  - If  $y = \epsilon$  then  $v$  is a *suffix* of  $w$
- A *subsequence* of a string  $w = w_1w_2 \dots w_n$  is either a subsequence of  $w_2 \dots w_n$  or  $w_1$  followed by a subsequence of  $w_2 \dots w_n$ .
- If  $w$  is a string then  $w^n$  is defined inductively as follows:  $w^n = \epsilon$  if  $n = 0$  or  $w^n = ww^{n-1}$  if  $n > 0$

### String operations

## 2 Overview of language complexity

### Overview



Grammar	Languages	Production Rules	Automaton	Examples
Type-0	recursively enumerable	$\gamma \rightarrow \alpha$ (no constraints)	Turing machine	$L = \{w \mid w \text{ is a TM which halts}\}$
Type-1	context-sensitive	$\alpha A \beta \rightarrow \alpha \gamma \beta$	linear bounded nondeterministic Turing machine	$L = \{a^n b^n c^n \mid n > 0\}$
Type-2	context-free	$A \rightarrow \alpha$	nondeterministic pushdown automata	$L = \{a^n b^n \mid n > 0\}$
Type-3	regular	$A \rightarrow aB$	finite state machine	$L = \{a^n \mid n > 0\}$

Meaning of symbols:

- $a$  - terminal
- $A, B$  - variables
- $\alpha, \beta, \gamma$  - strings in  $\{a \cup A\}^*$  where  $\alpha, \beta$  are maybe empty,  $\gamma$  is never empty

<sup>a</sup>Table borrowed from Wikipedia: [https://en.wikipedia.org/wiki/Chomsky\\_hierarchy](https://en.wikipedia.org/wiki/Chomsky_hierarchy)

### 3 Regular languages

#### Regular language - overview

A language is regular if and only if it can be obtained from finite languages by applying

- union,
- concatenation or
- Kleene star

finitely many times. All regular languages are representable by regular grammars, DFAs, NFAs and regular expressions.

#### Regular expressions

Useful shorthand to denotes a language.

A regular expression  $r$  over an alphabet  $\Sigma$  is one of the following:

**Base cases:**

- $\emptyset$  the language  $\emptyset$
- $\epsilon$  denotes the language  $\{\epsilon\}$
- $a$  denote the language  $\{a\}$

**Inductive cases:** If  $r_1$  and  $r_2$  are regular expressions denoting languages  $L_1$  and  $L_2$  respectively (i.e.,  $L(r_1) = L_1$  and  $L(r_2) = L_2$ ) then,

- $r_1 + r_2$  denotes the language  $L_1 \cup L_2$
- $r_1 \cdot r_2$  denotes the language  $L_1 L_2$
- $r_1^*$  denotes the language  $L_1^*$

**Examples:**

- $0^*$  - the set of all strings of 0s, including the empty string
- $(00000)^*$  - set of all strings of 0s with length a multiple of 5
- $(0 + 1)^*$  - set of all binary strings

#### Nondeterministic finite automata

NFAs are similar to DFAs, but may have more than one transition destination for a given state/character pair.

An NFA  $N$  accepts a string  $w$  iff some accepting state is reached by  $N$  from the start state on input  $w$ .

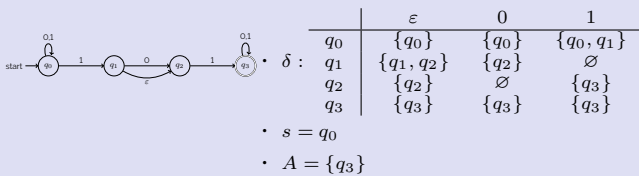
The language accepted (or recognized) by an NFA  $N$  is denoted  $L(N)$  and defined as  $L(N) = \{w \mid N \text{ accepts } w\}$ .

A nondeterministic finite automaton (NFA)  $N = (Q, \Sigma, s, A, \delta)$  is a five tuple where

- $Q$  is a finite set whose elements are called states
- $\Sigma$  is a finite set called the input alphabet
- $\delta : Q \times \Sigma \cup \{\epsilon\} \rightarrow \mathcal{P}(Q)$  is the transition function (here  $\mathcal{P}(Q)$  is the power set of  $Q$ )
- $s$  and  $\Sigma$  are the same as in DFAs

Example:

- $Q = \{q_0, q_1, q_2, q_3\}$
- $\Sigma = \{0, 1\}$



For NFA  $N = (Q, \Sigma, \delta, s, A)$  and  $q \in Q$ , the  $\epsilon$ -reach( $q$ ) is the set of all states that  $q$  can reach using only  $\epsilon$ -transitions.

Inductive definition of  $\delta^* : Q \times \Sigma^* \rightarrow \mathcal{P}(Q)$ :

- if  $w = \epsilon$ ,  $\delta^*(q, w) = \epsilon$ -reach( $q$ )
- if  $w = a$  for  $a \in \Sigma$ ,  $\delta^*(q, a) = \epsilon$ reach( $\bigcup_{p \in \epsilon$ -reach( $q$ )  $\delta(p, a)$ ))
- if  $w = ax$  for  $a \in \Sigma, x \in \Sigma^*$ :  $\delta^*(q, w) = \epsilon$ reach( $\bigcup_{p \in \epsilon$ -reach( $q$ ) ( $\bigcup_{r \in \delta^*(p, a)} \delta^*(r, x)$ ))

#### Regular closure

Regular languages are closed under union, intersection, complement, difference, reversal, Kleene star, concatenation, etc.

#### Deterministic finite automata

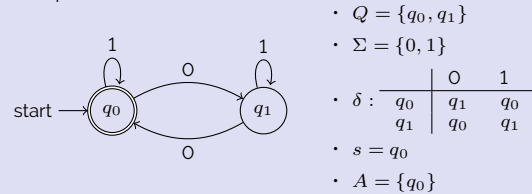
DFAs are finite state machines that can be represented as a directed graph or in terms of a tuple.

The language accepted (or recognized) by a DFA  $M$  is denoted by  $L(M)$  and defined as  $L(M) = \{w \mid M \text{ accepts } w\}$ .

A deterministic finite automaton (DFA)  $M = (Q, \Sigma, s, A, \delta)$  is a five tuple where

- $Q$  is a finite set whose elements are called states
- $\Sigma$  is a finite set called the input alphabet
- $\delta : Q \times \Sigma \rightarrow Q$  is the transition function
- $s \in Q$  is the start state
- $A \subseteq Q$  is the set of accepting/final states

Example:



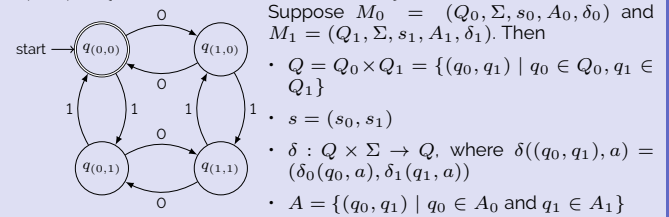
Every string has a unique walk along a DFA. We define the extended transition function as  $\delta^* : Q \times \Sigma^* \rightarrow Q$  defined inductively as follows:

- $\delta^*(q, w) = q$  if  $w = \epsilon$
- $\delta^*(q, w) = \delta^*(\delta(q, a), x)$  if  $w = ax$ .

Can create a larger DFA from multiple smaller DFAs. Suppose

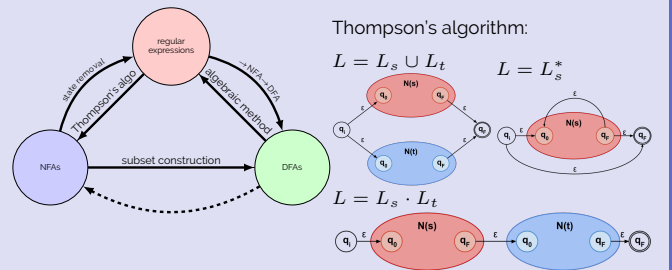
- $L(M_0) = \{w \text{ has an even number of 0s}\}$  (pictured above) and
- $L(M_1) = \{w \text{ has an even number of 1s}\}$ .

$L(M_C) = \{w \text{ has even number of 0s and 1s}\}$



#### Regular language equivalences

A regular language can be represented by a regular expression, regular grammar, DFA and NFA.



**Arden's rule:** If  $R = Q + RP$  then  $R = QP^*$ .

#### Fooling sets

Some languages are not regular (Ex.  $L = \{0^n 1^n \mid n \geq 0\}$ ).

Two states  $p, q \in Q$  are distinguishable if there exists a string  $w \in \Sigma^*$ , such that

$$\delta^*(p, w) \in A \text{ and } \delta^*(q, w) \notin A.$$

or

Two states  $p, q \in Q$  are equivalent if for all strings  $w \in \Sigma^*$ , we have that

$$\delta^*(p, w) \in A \iff \delta^*(q, w) \in A.$$

$$\delta^*(p, w) \notin A \text{ and } \delta^*(q, w) \in A.$$

For a language  $L$  over  $\Sigma$  a set of strings  $F$  (could be infinite) is a fooling set or distinguishing set for  $L$  if every two distinct strings  $x, y \in F$  are distinguishable.

## 4 Context-free languages

### Context-free languages

A language is context-free if it can be generated by a context-free grammar. A context-free grammar is a quadruple  $G = (V, T, P, S)$

- $V$  is a finite set of *nonterminal (variable) symbols*
- $T$  is a finite set of *terminal symbols* (alphabet)
- $P$  is a finite set of *productions*, each of the form  $A \rightarrow \alpha$  where  $A \in V$  and  $\alpha$  is a string in  $(V \cup T)^*$ . Formally,  $P \subseteq V \times (V \cup T)^*$ .
- $S \in V$  is the *start symbol*

Example:  $L = \{ww^R \mid w \in \{0, 1\}^*\}$  is described by  $G = (V, T, P, S)$  where  $V, T, P$  and  $S$  are defined as follows:

- $V = \{S\}$
- $T = \{0, 1\}$
- $P = \{S \rightarrow \varepsilon \mid 0S0 \mid 1S1\}$   
(abbreviation for  $S \rightarrow \varepsilon, S \rightarrow 0S0, S \rightarrow 1S1$ )
- $S = S$

### Pushdown automata

A pushdown automaton is an NFA with a stack.

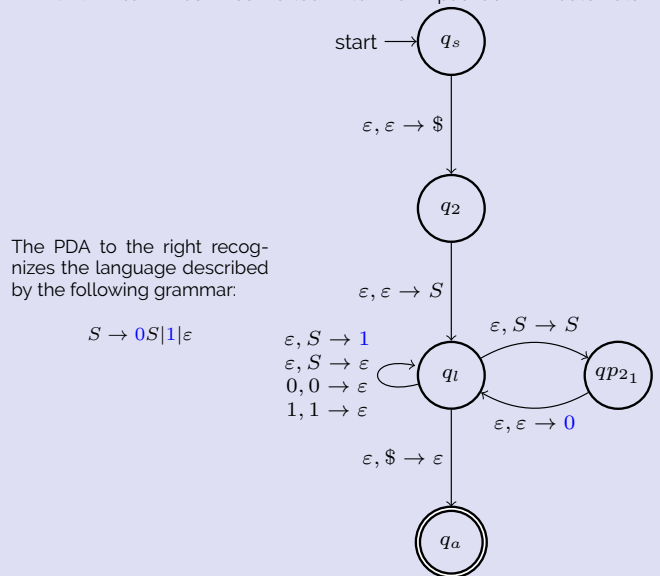
The language  $L = \{0^n 1^n \mid n \geq 0\}$  is recognized by the pushdown automaton:

A *nondeterministic pushdown automaton (PDA)*  $P = (Q, \Sigma, \Gamma, \delta, s, A)$  is a **six** tuple where

- $Q$  is a finite set whose elements are called *states*
- $\Sigma$  is a finite set called the *input alphabet*
- $\Gamma$  is a finite set called the *stack alphabet*
- $\delta : Q \times (\Sigma \cup \{\varepsilon\}) \times (\Gamma \cup \{\varepsilon\}) \rightarrow \mathcal{P}(Q \times (\Gamma \cup \{\varepsilon\}))$  is the *transition function*
- $s$  is the start state
- $A$  is the set of accepting states

In the graphical representation of a PDA, transitions are typically written as  $\langle \text{input read} \rangle, \langle \text{stack pop} \rangle \rightarrow \langle \text{stack push} \rangle$ .

A CFG can be converted to a pushdown automaton.



### Context-free closure

Context-free languages are closed under union, concatenation, and Kleene star.

They are **not** closed under intersection or complement.