

# Strings and Languages

## Languages of things

- └ numbers
- └ graphs
- └ polynomials

---

- Each thing should have a finite representation

- Each such representation can be encoded as a string

- Only the two symbols 0 and 1 suffice for all finite encodings as strings

## Alphabet

A finite set of symbols.  
letters

Examples : Binary alphabet  $\{0, 1\}$

Alphabet of decimal digits  $\{0, 1, 2, 3, \dots, 9\}$

Add + and - to this alphabet

Add . to represent floating-point numbers

ASCII alphabet — printable characters

Roman alphabet —  $\{a, b, \dots, z\} \cup \{A, B, \dots, Z\}$

Alphabet for English language

$\{\_ \}$   $\cup$   $\{a, b, \dots, z\}$   $\cup$   $\{A, B, \dots, Z\}$   $\cup$   $\{0, 1, 2, \dots, 9\}$   $\cup$   $\{\text{punctuation symbols}\}$

$\hookrightarrow$  space/blank  $\_ \neq$

Use  $\Sigma$ ,  $\Gamma$ ,  $\Delta$  to stand for alphabets

$\{0, 1\}$  suffices

Symbols:  $a, b, c, 0, 1, 2, \dots$

String: (over an alphabet  $\Sigma$ )

a finite sequence of symbols from  $\Sigma$

$\hookrightarrow$  ordered — first element  
2nd —  
3rd —

Example: 011000110001  $\rightarrow$  binary string

abracadabra  $\rightarrow$  string over the Roman alphabet

Notation:  $u, v, w, x, y, z$   
 $\alpha, \beta, \gamma$  to represent strings

length of a string is the number of symbols in it.

$$|\text{abracadabra}| = 11$$

length of  $w$  is denoted as  $|w|$

String of length 0

"null"

"" ( $\epsilon$ ) — the empty string  $|\epsilon| = 0$ .

Given  $\Sigma$ , we denote by  $\Sigma^*$  the set of all strings over  $\Sigma$ .

$$\Sigma = \{a, b, c\}$$

$$\Sigma^* = \left\{ \epsilon, a, b, c, aa, ab, ac, ba, bb, bc, ca, cb, cc, aaa, aab, aac, \dots \right\}$$

$$a \in \Sigma \quad a^n = \underbrace{aaa \dots a}_{n \text{ times}}$$

$$a^0 = \epsilon$$

$\Sigma^*$  supports an operation called concatenation

$$xy = a_1 a_2 \dots a_m b_1 b_2 \dots b_n$$

$$x = a_1 a_2 \dots a_m$$

$$y = b_1 b_2 \dots b_n$$

# Concatenation on $\Sigma^*$

- closed
- associative
- ~~x~~ - not necessarily commutative
- $\epsilon$  acts as the identity
- ~~x~~ - No inverse

$$(ab)(ba) \neq (ba)(ab)$$

$$w\epsilon = \epsilon w = w$$

monoid

$\Sigma^*$  is a monoid under concatenation

Languages

A language over an alphabet  $\Sigma$  is any subset of  $\Sigma^*$ .

↓  
 $L, L_1, L_2, \dots, A, B \rightarrow$  represent languages  
sets of strings  $\rightarrow$  a lot of things

# Countability vs Uncountability

$\Sigma$  is a finite alphabet ( $\Sigma \neq \emptyset$ )  $|\Sigma| = 1$   
unary alphabet

$\Sigma^*$  is countable. There are  $|\Sigma|^n$  strings of length  $n$   
 $n = 0, 1, 2, 3, \dots$

→ a countable union of countable (finite) sets

$$L \subseteq \Sigma^*$$

The set of all languages over  $\Sigma$

$$= \mathcal{P}(\Sigma^*) = 2^{\Sigma^*}$$

→  
power-set  
theorem

uncountable.

# Representation

Strings  $\rightarrow$  finite sequences (no problem)

Languages  $\rightarrow$  finite  $\rightarrow$  explicit enumeration  
 $\{w_1, w_2, w_3, \dots, w_n\}$

$\searrow$  infinite  $\rightarrow$  explicit enumeration not possible

Require finite description

- English language [the set of all strings of even length]

- Mathematical description [ $\{w \in \{0,1\}^* \mid w \text{ is a multiple of } 3\}$ ]

- Recursive description [palindromes over  $\{a, b\}$ ]

[Grammars] symbols

- [Procedures / Machines]  $\text{symb}_2 \text{symb}_1 \text{core} \text{symb}_1 \text{symb}_2 \dots \text{symb}_n$   $\text{core} \rightarrow \epsilon, a, b$

Languages over  $\Sigma$

Description of a language should be "finite"

uses some alphabet

English/Roman alphabet

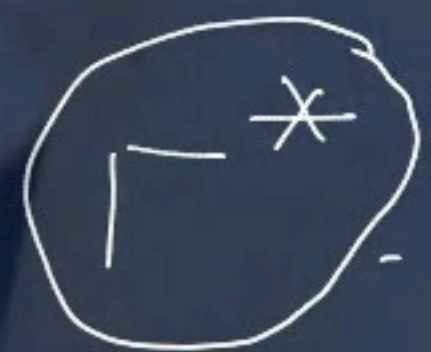
alphabet of math symbols



$|\Gamma|$  is finite (and non-empty)

a string over  $\Gamma$

Descriptions are elements of



countable

cf  $\mathcal{P}(\Sigma^*)$  is uncountable

— Only countably many languages can have finite descriptions

— Most of the languages cannot be described (finitely).

## Computational Problems

Given  $L \subseteq \Sigma^*$  and  $w \in \Sigma^*$ , decide whether  $w \in L$  or not. [Language-membership problem]

Decision problem — This definition is good for all theoretical developments

— Stand the test of time

Problems  $\equiv$  Languages

- There are uncountably many problems.
- Only countably many machines (automata) to solve problems.
- There are unsolvable problems — Difficult to find  
— Proving unsolvability



# Strings and Languages

## • Operations on strings

- concatenation (defined)

- Powers of a string  $w^n, n \geq 0$

$$w^n = \begin{cases} \epsilon & \text{if } n = 0 \\ ww^{n-1} & \text{if } n \geq 1 \end{cases}$$

$$(01)^3 = 010101$$

$$01^3 = 0111$$

$$0^3 1^3 = 000111$$

- Prefix (Suffix):  $w \in \Sigma^*$ .  $u \in \Sigma^*$  is called a prefix (Suffix) of  $w$  if  $w = uv$  for some  $v \in \Sigma^*$ .

All prefixes of

$abbc$  are  $\rightarrow$  proper

$\epsilon, a, ab, abb, abbc,$   
 $abbc$

All suffixes of  $abbc$  are

$\epsilon, c, bc, abc,$   
 $abbc$

$\downarrow$   
proper

## Operations on languages

Languages are subsets of  $\Sigma^*$   $\rightarrow$  universal set

$$A, B \subseteq \Sigma^*$$

$$A \cup B$$

$$A \cap B$$

$$\sim A = \overline{A} = \Sigma^* - A$$

Set-theoretic identities hold

- De Morgan

- Distributivity

$$A = \{a, ab\}$$

$$B = \{\epsilon, b\}$$

$$AB = \{a, b, ab, ab\}$$

## Concatenation

$$A, B \subseteq \Sigma^*$$

$$AB = \{uv \mid u \in A \text{ and } v \in B\}$$

Powers:  $A \subseteq \Sigma^*$

$$\Sigma = \{0, 1\}$$

$$A^n = \begin{cases} \{\epsilon\} & \text{if } n = 0 \\ AA^{n-1} & \text{if } n \geq 1 \end{cases}$$

$$\Sigma^* = \{\epsilon\} \cup \{0, 1\} \cup \{0, 1\}^2 \cup \{0, 1\}^3 \cup \dots$$

$$A = \{a, ab\}$$

$$A^0 = \{\epsilon\}$$

$$A^1 = \{a, ab\}$$

$$A^2 = \{aa, aab, aba, abab\}$$

$A^*$  (asterate of A) (Kleene star)

$$A^* = \bigcup_{n \geq 0} A^n$$
$$A^+ = \bigcup_{n \geq 1} A^n$$

Identities : [easy proofs]

$$A^* A^* = A^*$$

$$(A^*)^* = A^*$$

$$\emptyset^* = \{\epsilon\}$$

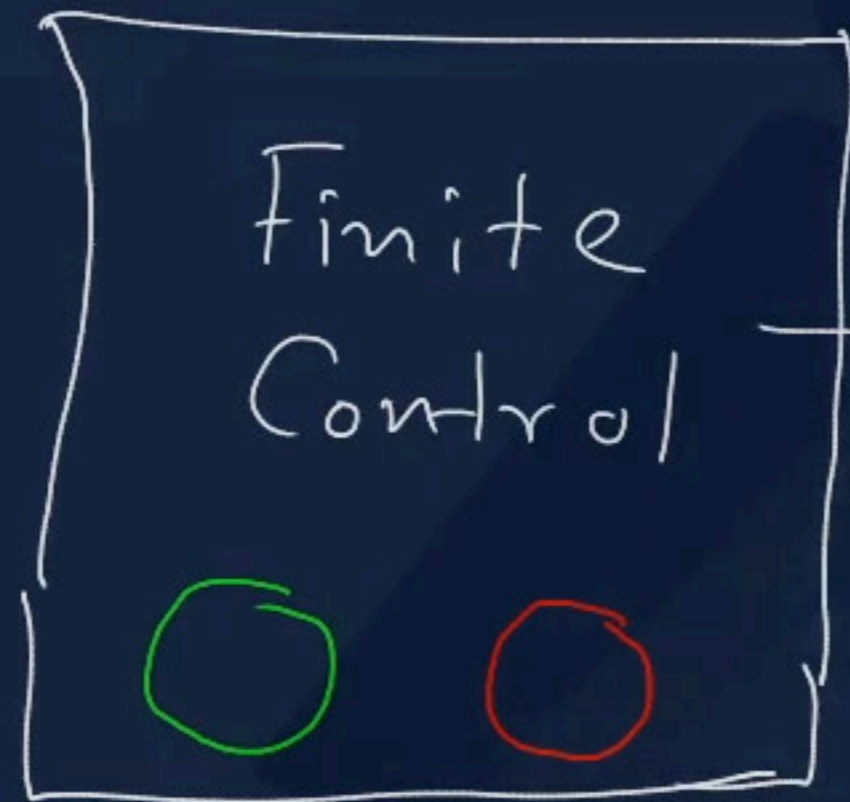
$$[\emptyset^0 = \{\epsilon\}, \emptyset = \emptyset \{\epsilon\} = \emptyset, \emptyset^2, \emptyset^3, \dots = \emptyset]$$

$$\begin{aligned} A^* &= \{\epsilon\} \cup A A^* \\ &= \{\epsilon\} \cup A^* A \end{aligned}$$

# Description of Languages by automata [plural of automaton]

Represents

L



consists of a finite set  $Q$  of states

Input string



(Externally)

$$L \subseteq \Sigma^*$$

Input:  $w \in \Sigma^*$

Output: Accept or Reject

Do computation on w.

If  $w \in L$ , the green lamp glows  
If  $w \notin L$ , the red lamp glows. ] M/c stops.

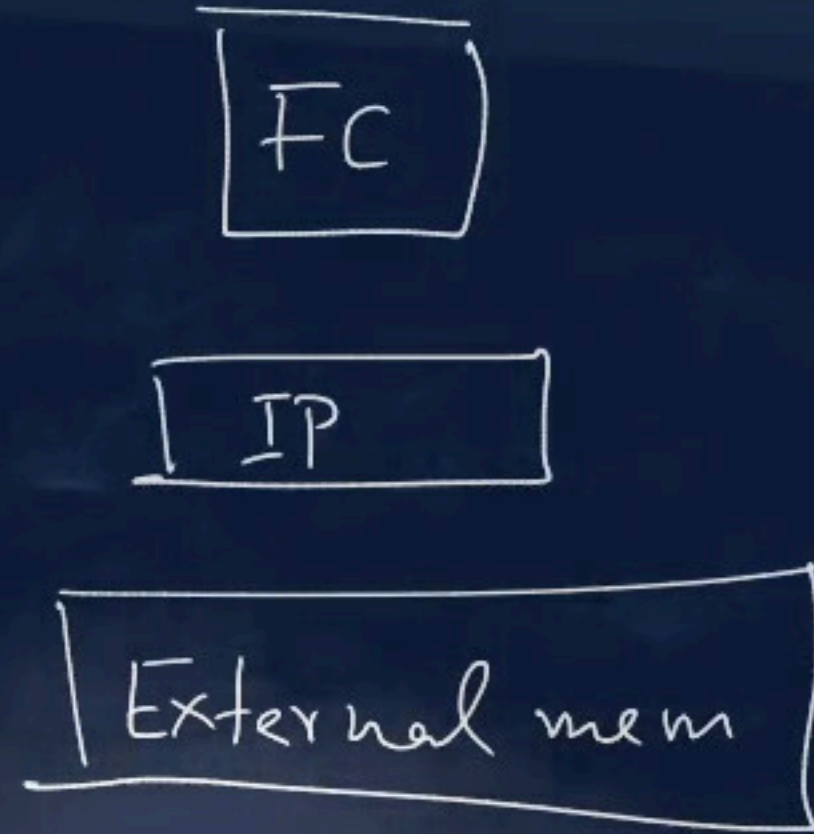
# Different types of automata

External memory

- Finite automata

(no external memory)

Only finite memory (FC)

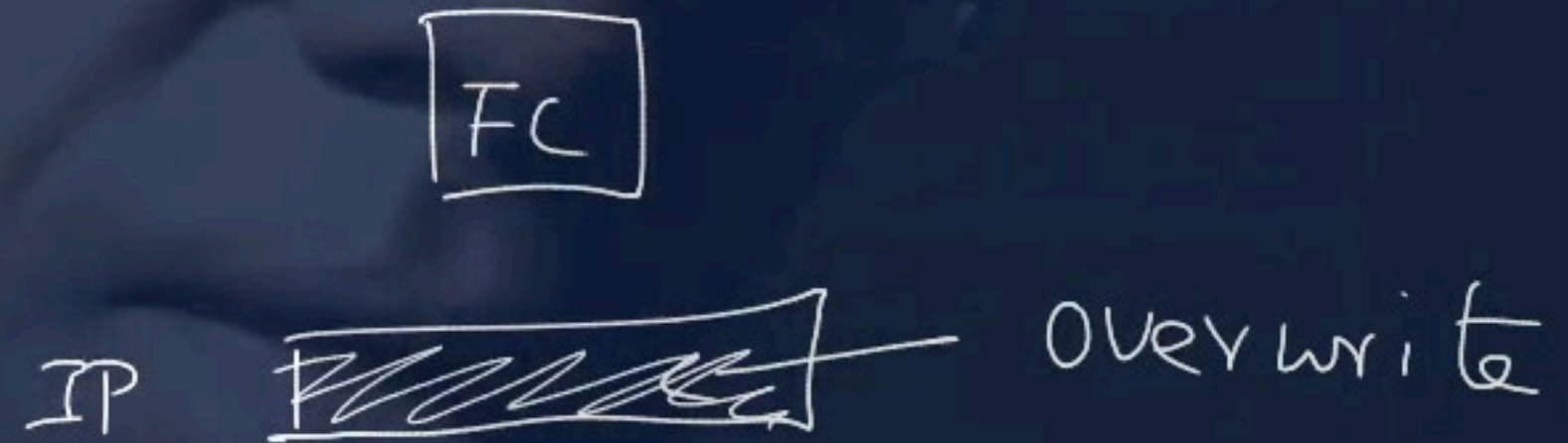


- Pushdown automata

Infinite external memory organized and accessed as a stack

- Linear-bounded automata

IP tape is rewritable.



- Turing machine

Infinite read/write tape



# Chomsky Hierarchy

# Grammars

<u>Regular</u>	—	<u>Finite automata</u>	—	Linear grammars (3)
CFL	—	<u>Pushdown automata</u>	—	<u>Context-free grammars</u> (2)
CSL	—	<u>Linear bounded automata</u>	—	<u>Context-sensitive grammars</u> (1)
	—	<u>Turing machines</u>	—	<u>Unrestricted grammars</u> (Type 0)

