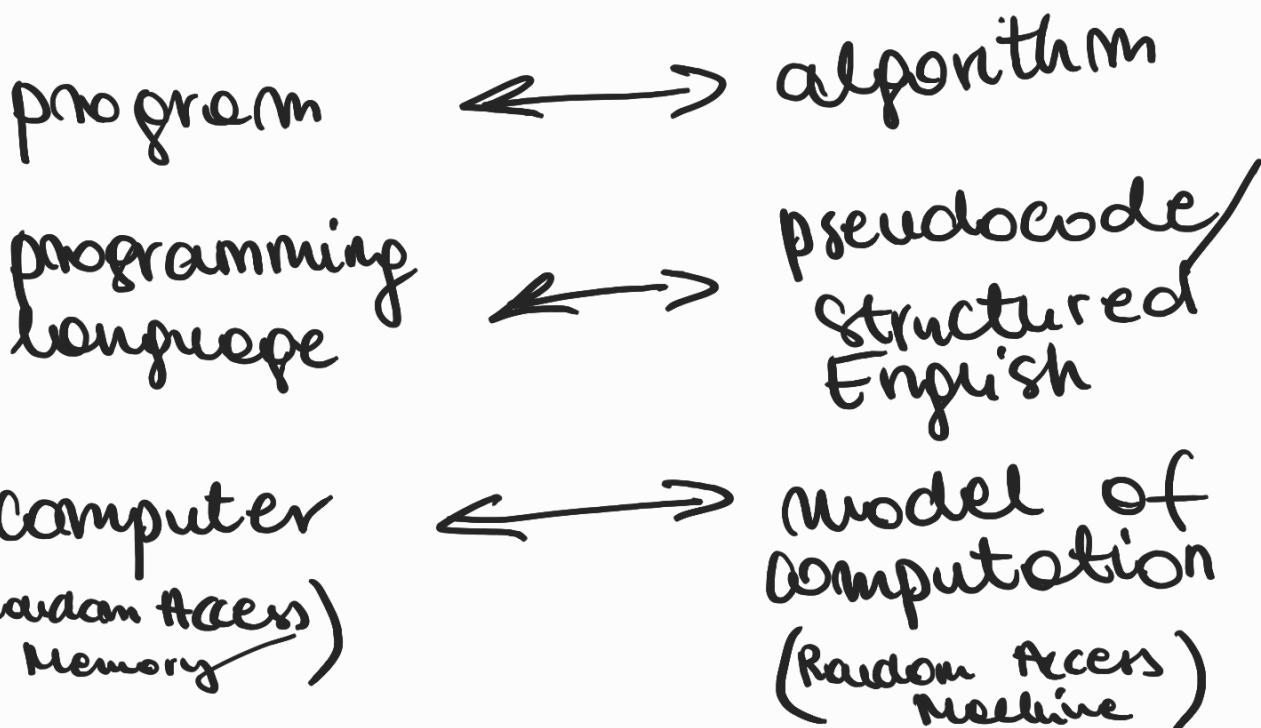


# MODELS OF COMPUTATION, DOCUMENT DISTANCE

What's an algorithm?

- computational procedure for solving a problem
- input  $\rightarrow$  alg  $\rightarrow$  output



## Model of computation

specifies

- what operation an algorithm is allowed
- cost (time) of each operation

① Random Access Machine (RAM)

- Random Access Memory  
modeled by big array
- in  $\Theta(1)$  time an algorithm can
  - load  $\Theta(1)$  words
  - do  $\Theta(1)$  computations
  - store  $\Theta(1)$  words
  - $\Theta(1)$  registers

WORD : w bits

$$\sqrt{w} \lg(\text{size of memory})$$

## ② Pointer Machine

- dynamically allocated objects
- object has  $\Theta(1)$  fields
- field = word (e.g. int)  
or pointer to object

Python model

① "list" = array

$$L[i] = L[j] + 5$$

$\rightarrow \Theta(1)$  time

② object with  $O(1)$  attributes

$x = x.\text{next} \rightarrow O(1)$  time

③  $L.append(x)$

$\hookrightarrow$  table doubling ( $L8$ )

$\hookrightarrow O(1)$  time

④  $L = L1 + L2 \rightarrow$

$\equiv L = []$

for  $x$  in  $L1:$   $\begin{cases} O(1) \\ O(|L1|) \end{cases}$

$L.append(x)$

for  $x$  in  $L2:$   $\begin{cases} O(1) \\ O(|L2|) \end{cases}$

$L.append(x)$

$O(|L1| + |L2|)$

⑤  $\text{len}(L) \rightarrow$  constant

⑥  $x \in L \rightarrow$  constant

⑦  $L.sort() \rightarrow O(|L| \log |L|)$

dictionary

⑧ dict:  $D[\text{key}] = \text{val.} \rightarrow O(1)$

⑨ long:

$x+y \rightarrow O(|x|+|y|)$

log<sub>2</sub>

$$x * y \rightarrow \underbrace{\Theta(|x| + |y|)^2}_{\text{can be } \Theta(|x| * |y|)} \quad \text{better}$$

## Document distance problem:

$$d(D_1, D_2)$$

- document = sequence of words
- word = string of alphanumeric chars
- idea = shared words
- think of document as a vector

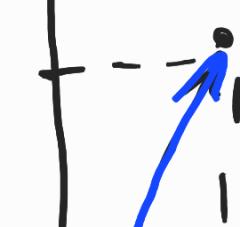
$D[w] = \# \text{occurrences of } w \text{ in } D$

$D_1 =$  "the cat"

$D_2 =$  "the dog" def

$$d'(D_1, D_2) = D_1 \cdot D_2 =$$

$$= \sum_{i=1}^n D_1[i] \cdot D_2[i]$$



$$= \sum_{w_j} D_1(w_j) \cdot D_2(w_j)$$

$$d'(D_1, D_2) = \frac{D_1 \cdot D_2}{|D_1| \cdot |D_2|}$$

$$d'(D_1, D_2) = \arccos$$

Algorithm:

- ① split doc. into words
- ② compute word frequencies
- ③ compute dot product

iterate through the words in doc:

and put in dictionary

20(idea)

found and 1 = M

Count[word] +- 1



$O(1 \text{ word})$

---

OTHER  
OPTION

Sort and  
then count

---

