

Algorithms in Genome Research

4. Dec. 2020

Patterns in non-coding genomic regions

Q.S.

- Standing sites
- repeat regions
- small RNAs (e.g.: rRNAs)
- other motifs that form secondary structures e.g. in tRNA
- horizontally transferred DNA

?

1.

A) (1) Repeat detection on a genomic scale

- tandem repeats



w w w w w w w w

p

approximate X.R.

a) limited edit distance between neighbors

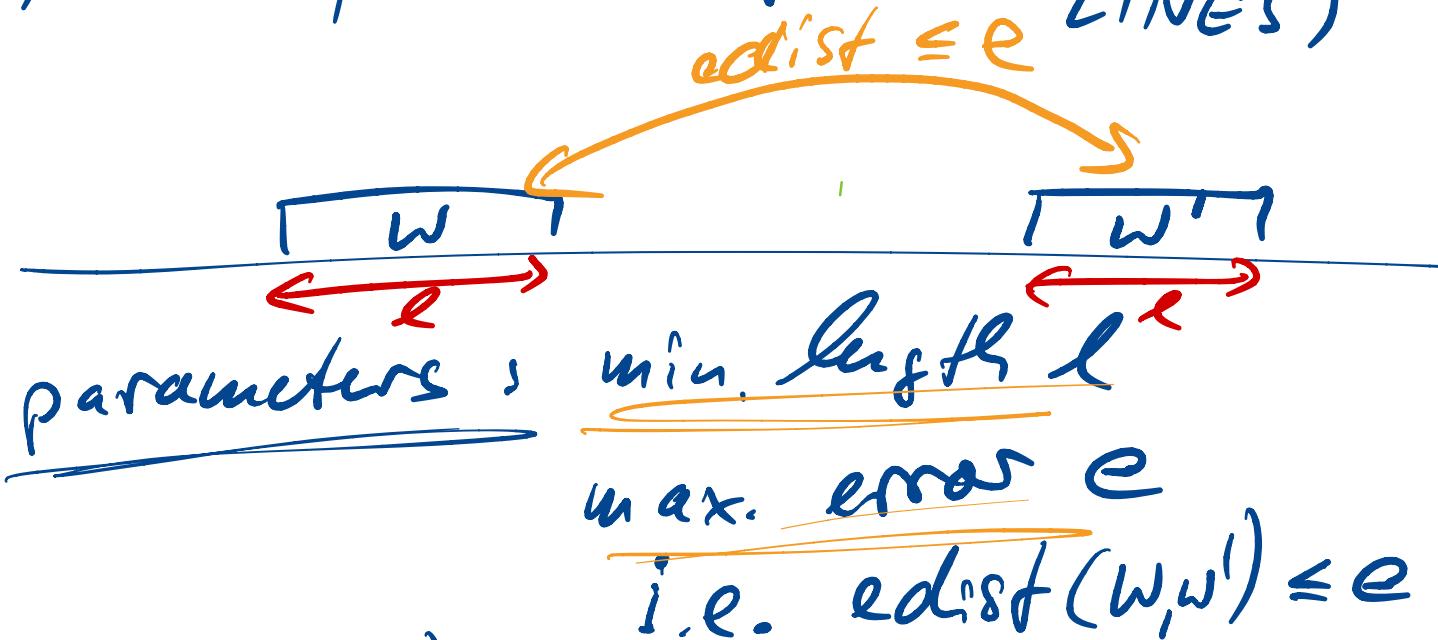
w₁ w₂ w₃ w₄ ...

edit(w_i, w_{i+1}) ≤ t

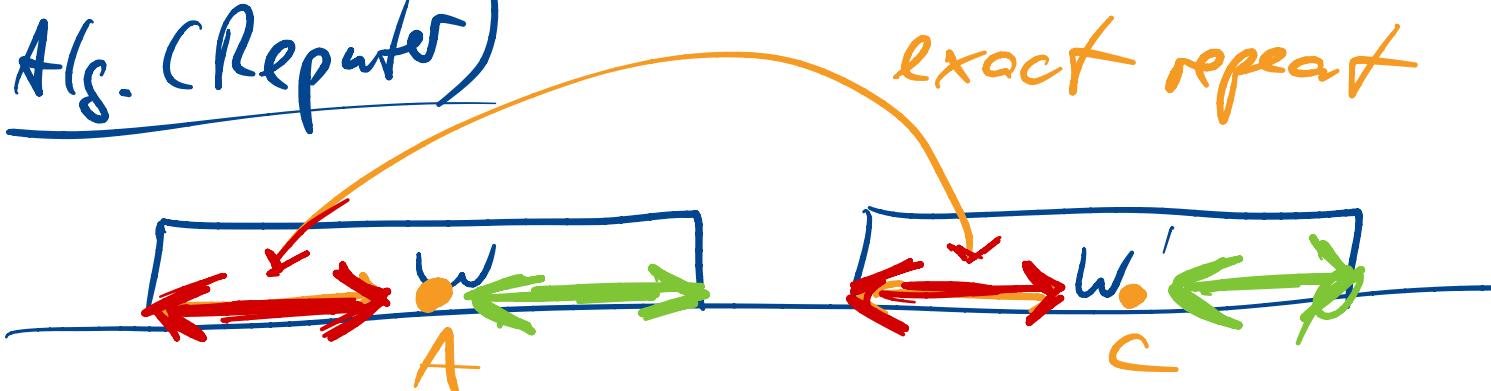
b) global model w

edit(w_i, ĉ) ≤ t

(2) interspersed repeats (SINES, LINES)



Alg. (Repeater)



observe, there is always an exact repeat of seed length

$$s \geq \left\lfloor \frac{l}{e+1} \right\rfloor.$$

Algorithm:

1. find maximal exact repeats

$$\text{of length } s \geq \left\lfloor \frac{l}{e+1} \right\rfloor.$$

\rightarrow using suffix trees / arrays /

2. extend each seed

BWT
 \hookrightarrow up to e errors to the left
and to the right.

$$l=5, e=1$$



3. test if there is an interval of length at least l with e errors.

How many seeds z^2 ?
 $E[z] = O\left(\frac{n^2}{4^s}\right)$

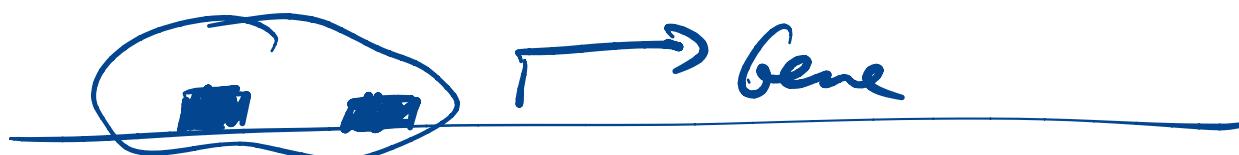
Overall runtime: $O(n + ze)$

in practice: can be applied
 to big genomes
 with reasonable
 parameters
 e.g. $l=100, e=3$

B

Finding Regulatory Regions

Promotor region



discovery

here: ~~detection~~ of unknown binding sites

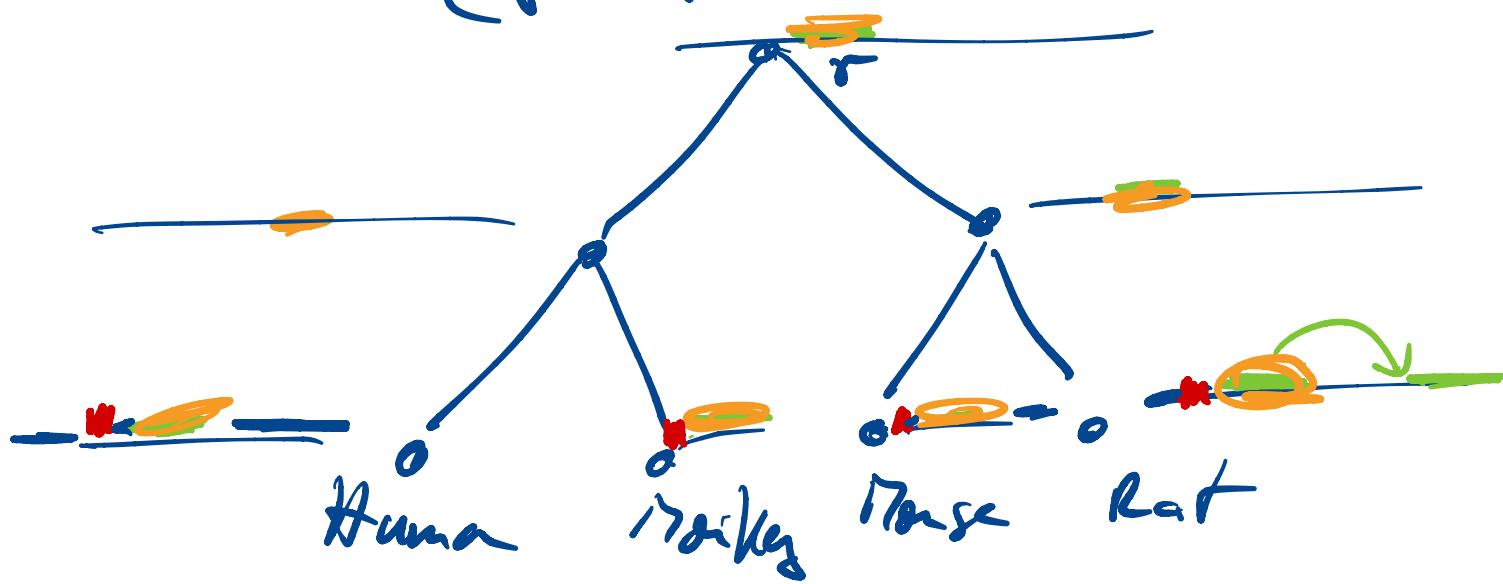
(1) examine co-regulated genes



- find binding site of common regulatory elements

(2) examine orthologous genes
(of different species)

Paralogous



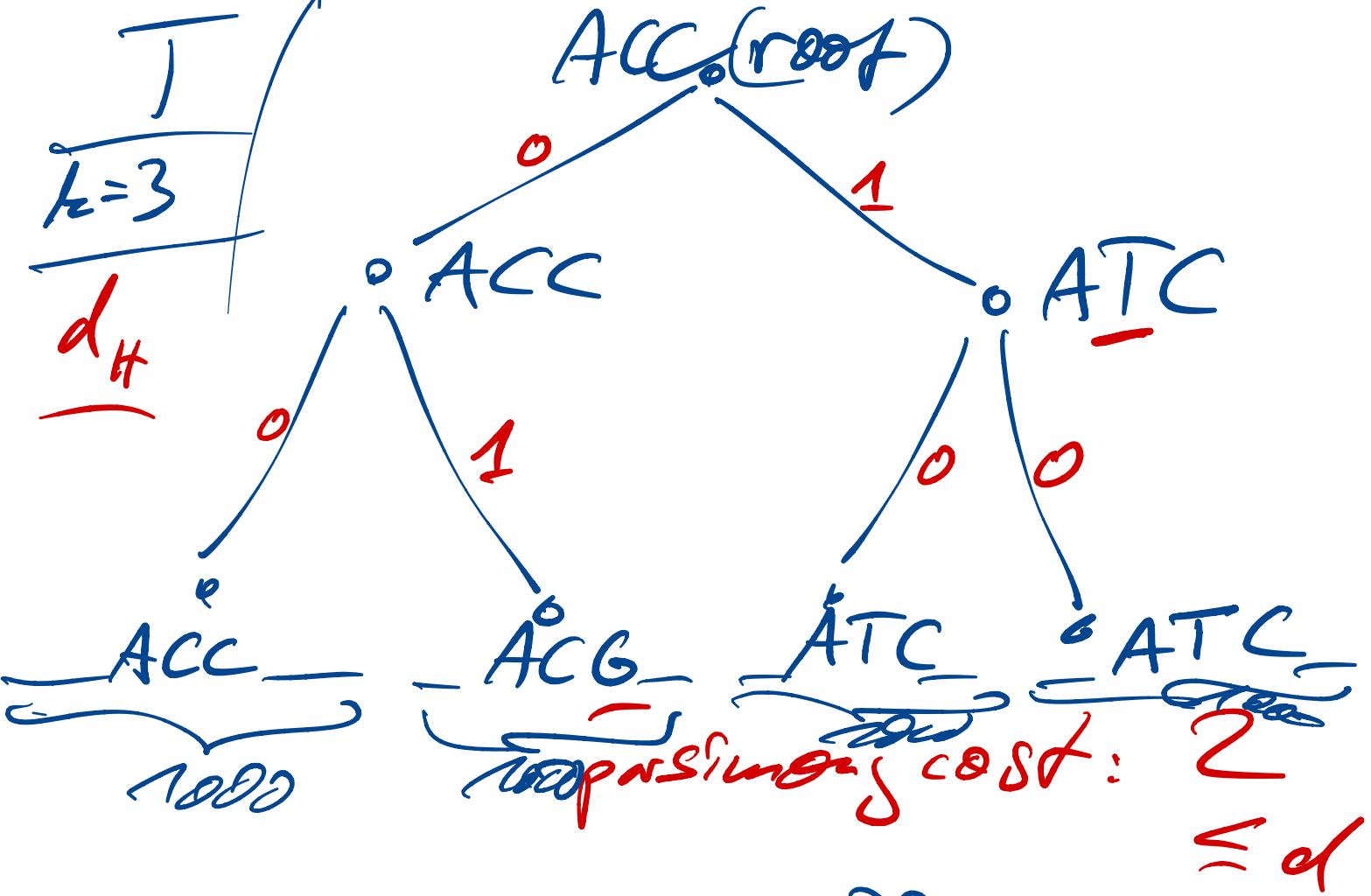
Known: binding sites (8 genes)
 are better conserved
 than the surrounding areas
 binding site, better conserved than
 the rest

Suspecting Parsimony Problem (SPP)

Given: S_1, \dots, S_n set of promoter regions
 of orthologous genes

T phylogenetic tree
 of the n species w/ root
 k length of the motif
 d maximum # of errors allowed.

wanted: All sets of substitutions s_1, \dots, s_n
 of S_1, \dots, S_n , each of length k ,
 such that the parsimony cost
 of s_1, \dots, s_n along T is at most d .



Solution: A simple DP algorithm

bottom up processing of T:

initialisation of leaf u

$$W_u[s] = \begin{cases} 0 & \text{if } s \text{ is a substring of } s_u \\ \infty & \text{otherwise} \end{cases}$$

Recursive Calculation

of internal vertex u
with children $C(u)$

$$w_u[s] = \sum_{v \in C(u)} \min_{t \in \Sigma^k} (w_v[t] + d(s, t))$$

W_u

AA	0
AC	0
AG	0
AT	0
CA	0
:	0
TT	0

$O_{f1} = 1$

$O_{f2} = 1$

