

## # Cartesian trees. (Vuillemin '80)

Consider a string  $w$  over the Integer alphabet.

A Cartesian tree  $C_T(w)$  of  $w$  is recursively defined as:

> if  $|w| = 0$ ,  $C_T(w)$  is empty.

> otherwise, let  $t$  be the smallest index with the minimum  $w[i]$ .

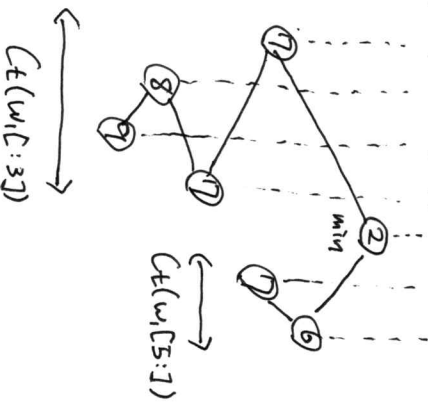
- The root of  $C_T(w)$  is labeled  $w[t]$

- The left subtree of the root is  $C_T(w[1:t-1])$

- The right subtree  $C_T(w[t+1:n])$

Two Cartesian trees are equivalent if they have the same structure.

Example  $w_1 = (7, 8, 9, 7, 2, 7, 6)$



$w_2 = (2, 8, 16, 4, 1, 3, 2) \Rightarrow C_T(w_1) = C_T(w_2)$

## # Known facts.

$C_T(w)$  can be constructed in  $O(|w|)$  time. (Gabow et al. '84)

$C_T(w_1) \stackrel{?}{=} C_T(w_2)$  can be tested in  $O(|w_1| + |w_2|)$  time.

The pattern matching problem w.r.t. Cartesian trees

can be ~~solved~~ solved in  $O(|T_1| + |P|)$  time. (Park et al. '19)

## # Problem. [Approximate Cartesian tree pattern matching]

(Informal) Let the Cartesian edit distance from string  $u$

to string  $v$  be the min. total cost of edits on  $u$  to

make another string  $u'$  allowing  $C_T(u') = C_T(v)$ .

Denote as  $\text{Cdist}(u \rightarrow v)$ .

Given a text  $T$ , ~~and~~ pattern  $P$  and a threshold  $t$ ,

compute all sub strings  $w$  of  $T$  that satisfy

$\text{Cdist}(w \rightarrow P) \leq t$ .

↳ ~~Compute this~~ Compute this

## # The Cartesian edit distance is asymmetric.

Consider  $w = (3, 4, 5, 5, 4, 3)$  and  $u = (4, 5, 5, 5, 4)$ .

$\text{Cdist}(w \rightarrow u) = 1$  but  $\text{Cdist}(u \rightarrow w) = 2$ .

(assume unit cost for insert/delete/substitute.)

# Greedy does not work.

Consider  $w = (5, 4, 4, 3, 2)$  and  $u = (10, 9, 8, 7, 6)$

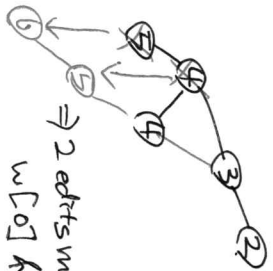
$Ct(w[1:2])$  and  $Ct(u[1:3])$ :



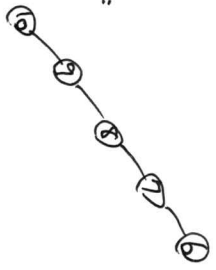
$\Rightarrow$  1 edit is enough

However,

$Ct(w)$ :



and  $Ct(u)$ :



$\Rightarrow$  2 edits made on  $w[0] \& w[1]$

Observe the root value of the "edited" Cartesian tree may differ between edit sequences. This may disallow some future edit sequences if the root value is fixed to the current best ~~version of edits~~ edit sequence. w.r.t.

# Naive upper-bound on  $Cdist(w \rightarrow u)$

Consider deleting all chars in  $w$  except one, and inserting the  $|u|-1$  other chars to make the same Cartesian tree.

$\Rightarrow Cdist(w \rightarrow u)$  is  $O(|w| + |u|)$ .

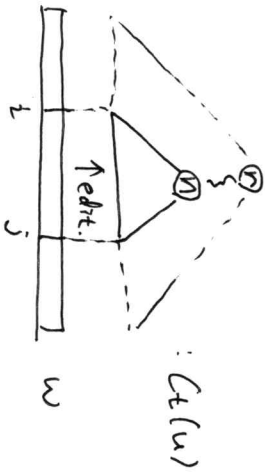
# Idea for computing  $Cdist(w \rightarrow u)$

$\triangleright$  Consider every 3-tuple  $(i, j, n)$ :

-  $0 \leq i \leq j \leq |w|$  and  $n$  is a node in  $Ct(u)$

~~$(i, j, n)$~~

Call each 3-tuple a subproblem; edit  $w[i:j]$  to match the subtree rooted at  $n$  in  $Ct(u)$



Recall that more edits in this subproblem may ~~lead to~~ lead to

higher root labels in the edited  $Ct$ .

$\Rightarrow$  Save (max root label, edit cost) tuples for all

edit costs under the naive upper bound for  $Cdist(w \rightarrow u)$ .

i.e.,  ~~$(i, j, n)$~~  we can define a function that maps

$(i, j, n, \pi)$  to a max root label.

$\hookrightarrow$  target edit cost.

Call that function OPT.

# Base cases. ~~Let~~ Let  $n$  be a leaf node.

1.  $\text{opt}(i, i, n, 1) = \infty$  (1 insert)
2.  $\text{opt}(i, j, n, j-i-1) = \max_{k \in [i, j-1]} w[k]$   $i \leq k \leq j-1$
3.  $\text{opt}(i, j, n, j-i) = \infty$  ( $j-i-1$  deletes, 1 substitution)
4.  $\text{opt}(i, j, n, \infty) = -\infty$  for all undefined  $\infty$ 's.  
↳ "impossible"

# Recurrence. Let  $n$  be an internal node.

W.l.o.g. let  $n$  have two children by placing dummy nodes  $n'$  satisfying  $\text{opt}(i, i, n', 0) = \infty$  and  $\text{opt}(i, j, n', \infty) = -\infty$  if  $i < j$ .

$$1. \text{insOpt}(i, j, n, \alpha) = \max_{y: z: y.z=0} \left( \min(\text{opt}(i, \alpha, y, y)^{-1}, \text{opt}(\alpha, j, r, z)) \right)$$

$$2. \text{subOpt}(i, j, n, \alpha, a) = \max_{y: z: y.z=0} \left( \min(\text{opt}(i, \alpha, y, y)^{-1}, \text{opt}(\alpha, j, r, z)) \right)$$

$$w[a] \text{ if } w[a] \leq \max_{y: z: y.z=0} w[a] \\ \text{if } w[a] > \max_{y: z: y.z=0} w[a] \text{ then } \left( \min(\text{opt}(i, \alpha, y, y)^{-1}, \text{opt}(\alpha, j, r, z)) \right)$$

$$3. \text{opt}(i, j, n, \infty) = \max_{a \in [i, j]} \left( \max(\text{insOpt}(i, j, n, \alpha), \text{subOpt}(i, j, n, \alpha)) \right)$$

⇒ Computing  $\text{opt}(0, \text{root}, \infty)$  for all possible  $\infty$ 's

take  $O(|w|^3 |u| D^2)$  time ( $D$  is the upper bound for  $\text{dist}(u \rightarrow v)$ .)

# Problem solution.

> Note:  $\text{opt}(i, j, \text{root}, \infty)$ 's are all computed through a single run of the recurrence.

⇒ Run the alg that computes  $\text{opt}$ , then return all  $(i, j)$  pairs with  $\infty$  that satisfy  $\text{opt}(i, j, \text{root}, \infty) \leq t$ .

If  $t$  is constant, we have  $D \leq t$ ; therefore the overall RT is  $O(|w|^3 |P|)$ .