# Tutorial Week 6

**Definition 1.** *The longest common prefix of two sequences $x$ and $y$ is denoted by $\text{lcp}(x, y)$. For some ordered alphabet $A$ such that $a < b$, we say that a string $u$ is smaller lexicographically than a string $v$, if either $u$ is a prefix of $v$, or $\text{lcp}(u, v) = x$ and $u[\|x\|] < v[\|x\|]$.*

**Exercise 1.** *Consider the following list of sequences: anchoritess, ancestral, ancestress, ancestor, ancestorial, ancestorially, anchored, anchoress, anchoritic, and, anchorite, anchor, anchorable, anchorage, ancestries, ancestry, anchorages, anchorate, ancestresses, anchoring, ancestors, anchoretical, anchorites, anchoretish, anchoretism.*

*Consider a random distribution of this list.:*
**a)** *Do a sorting of the list with bucket sort.*
**b)** *Calculate the* lcp *of consecutive elements.*
**c)** *Using the following Lemma compute the rest of the* lcp*'s.*

**Lemma 2.** *Let $L_0 \leq L_1 \leq \ldots \leq L_{n-1}$ and take the integers $d$, $i$, and $f$, with $-1 < d < i < f < n$. Then $\text{lcp}(L_d, L_f) = \min(\text{lcp}(L_d, L_i), \text{lcp}(L_i, L_f))$.*

**d)** *Draw the binary tree of the* lcp*'s for the previous item.*

*Solution:* **a)** The ordered list of sequences is the following: *ancestor, ancestorial, ancestorially, ancestors, ancestral, ancestress, ancestresses, ancestries, ancestry, anchor, anchorable, anchorage, anchorages, anchorate, anchored, anchoress, anchoretical, anchoretish, anchoretism, anchoring, anchorite, anchorites, anchoritess, anchoritic, and.*

To obtain the list by bucket sort, we extend by *nill* positions all short sequences, such that all have the same length. Next, we arrange them by looking at the letters from right to left, keeping them in place (that is, if we compare the $j$th letter of some strings $u$ and $v$, such that $u$ comes before $v$ in the text, then $u$ will occur after $v$ in the new list, only if $u[j] > v[j]$).

**b)** If the sequences are ordered lexicographically, we consider in the following list the element $i$, to be the $i$th element in the list (e.g., *ancestress* will be denoted by 6.) As there are 25 elements in our list, we will denote all these

by the numbers $0, 1, \ldots, 24$. Please recall that a $-1$ has to be added at the beginning and end to ease the computation and follow the algorithm presented in the course.

The list of all the consecutive lcp's is the following:

$\mathrm{lcp}(-1, 0) = \varepsilon$
$\mathrm{lcp}(0, 1) = ancestor$
$\mathrm{lcp}(1, 2) = ancestorial$
$\mathrm{lcp}(2, 3) = ancestor$
$\mathrm{lcp}(3, 4) = ancest$
$\mathrm{lcp}(4, 5) = ancestr$
$\mathrm{lcp}(5, 6) = ancestress$
$\mathrm{lcp}(6, 7) = ancestr$
$\mathrm{lcp}(7, 8) = ancestr$
$\mathrm{lcp}(8, 9) = anc$
$\mathrm{lcp}(9, 10) = anchor$
$\mathrm{lcp}(10, 11) = anchora$
$\mathrm{lcp}(11, 12) = anchorag$
$\mathrm{lcp}(12, 13) = anchora$
$\mathrm{lcp}(13, 14) = anchor$
$\mathrm{lcp}(14, 15) = anchore$
$\mathrm{lcp}(15, 16) = anchore$
$\mathrm{lcp}(16, 17) = anchoreti$
$\mathrm{lcp}(17, 18) = anchoretis$
$\mathrm{lcp}(18, 19) = anchor$
$\mathrm{lcp}(19, 20) = anchori$
$\mathrm{lcp}(20, 21) = anchorite$
$\mathrm{lcp}(21, 22) = anchorites$
$\mathrm{lcp}(22, 23) = anchorit$
$\mathrm{lcp}(23, 24) = an$
$\mathrm{lcp}(24, 25) = \varepsilon$

**c) and d)** We fill directly the balanced binary tree representation of the lcp's with the corresponding values, using Lemma 2 . The presentation of the balanced binary tree is the following:

A tree diagram with the following labeled nodes:

- $\varepsilon$
  - $\varepsilon$
    - $\varepsilon$
      - $\varepsilon$
        - $\varepsilon$
          - $\varepsilon$
          - $ancestor$
        - $ancestor$
          - $ancestorial$
          - $ancestor$
      - $ancest$
        - $ancest$
          - $ancest$
          - $ancestr$
        - $ancestr$
          - $ancestress$
          - $ancestr$
    - $anc$
      - $anc$
        - $anc$
          - $ancestr$
          - $anc$
        - $anchor$
          - $anchor$
          - $anchora$
      - $anchor$
        - $anchora$
          - $anchorag$
          - $anchora$
        - $anchor$
          - $anchor$
          - $anchore$
  - $\varepsilon$
    - $anchor$
      - $anchor$
        - $anchore$
          - $anchore$
          - $anchoreti$
        - $anchor$
          - $anchoretis$
          - $anchor$
      - $anchori$
        - $anchori$
        - $anchorite$
    - $\varepsilon$
      - $anchorit$
        - $anchorites$
        - $anchorit$
      - $\varepsilon$
        - $an$
        - $\varepsilon$