# Extracting the Sparse Longest Common Prefix Array from the Suffix Binary Search Tree

Tomohiro I [1]    Robert W. Irving [2]    *Dominik Köppl* [3]    Lorna Love [2]

[1] Department of Artificial Intelligence, Kyushu Institute of Technology, Japan

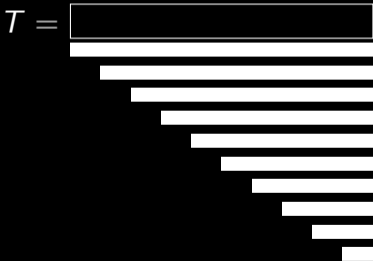[2] School of Computing Science, University of Glasgow, Glasgow, UK

[3] M&D Data Science Center, Tokyo Medical and Dental University, Japan
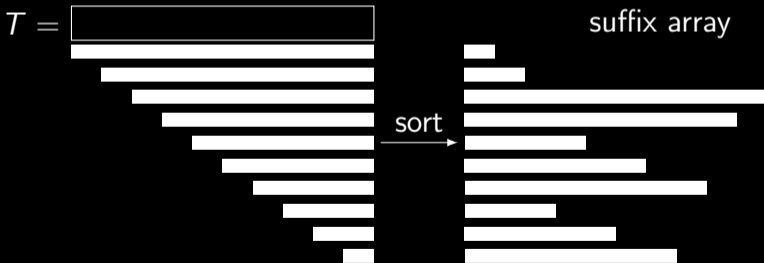
# suffix sorting

$T =$ [                    ]

# suffix sorting

- sort *all* suffixes lexicographically

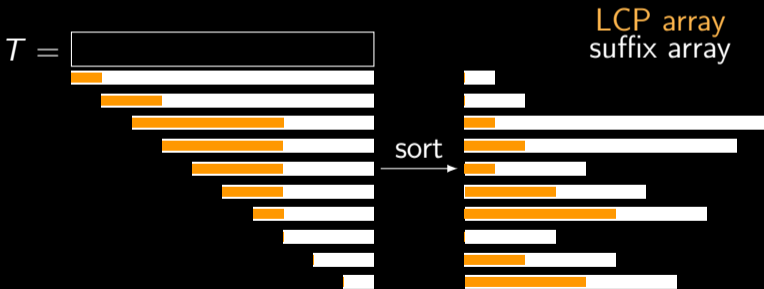$$T = $$

# suffix sorting

- sort *all* suffixes lexicographically



$T =$      suffix array

sort

# suffix sorting

- sort *all* suffixes lexicographically
- lengths of the longest common prefix (LCP) between adjacent suffixes.
- solved in $\mathcal{O}(n)$ time and words of space
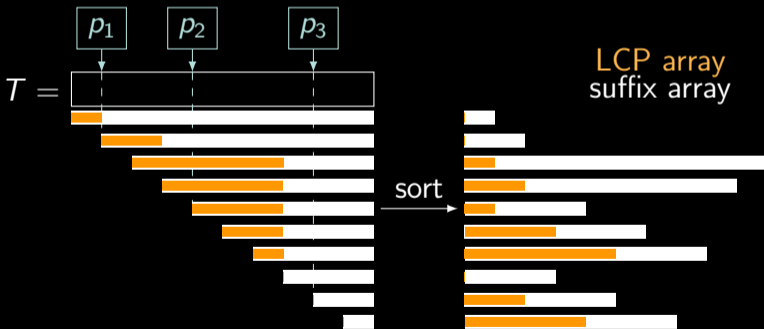
# suffix sorting

- sort *all* suffixes lexicographically
- lengths of the longest common prefix (LCP) between adjacent suffixes.
- solved in $\mathcal{O}(n)$ time and words of space
- sometimes need only suffixes starting at $p_1, \ldots, p_m$

# sparse suffix sorting

- sort *all* suffixes lexicographically
- lengths of the longest common prefix (LCP) between adjacent suffixes.
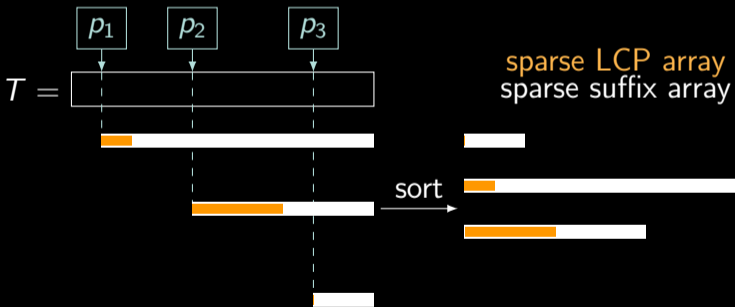- solved in $\mathcal{O}(n)$ time and words of space
- sometimes need only suffixes starting at $p_1, \ldots, p_m$

# dynamic sparse suffix sorting
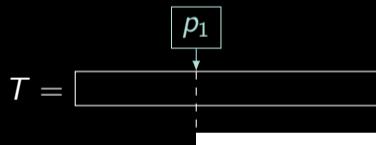
$T = $ ☐

# dynamic sparse suffix sorting

- $p_1, \ldots, p_m$: online, arbitrary order

# dynamic sparse suffix sorting

- $p_1, \ldots, p_m$: online, arbitrary order
- compare two suffixes with LCE query



LCE query $\text{lce}(p_1, p_2)$

# dynamic sparse suffix sorting

- $p_1, \ldots, p_m$: online, arbitrary order
- compare two suffixes with LCE query
- $c := \#$ characters to compare for sorting



compared positions $c$

LCE query $\mathrm{lce}(p_1, p_2)$

# dynamic sparse suffix sorting

- $p_1, \ldots, p_m$: online, arbitrary order
- compare two suffixes with LCE query
- $c := \#$ characters to compare for sorting



compared positions $c$

# dynamic sparse suffix sorting

- $p_1, \ldots, p_m$: online, arbitrary order
- compare two suffixes with LCE query
- $c := \#$ characters to compare for sorting



compared positions $c$

# dynamic sparse suffix sorting

- $p_1, \ldots, p_m$: online, arbitrary order
- compare two suffixes with LCE query
- $c := \#$ characters to compare for sorting
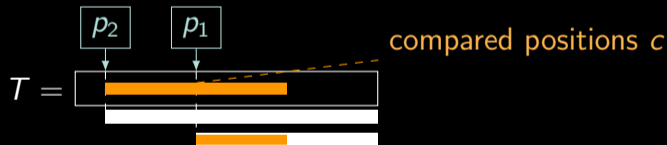


$T =$ compared positions $c$

LCE query

# dynamic sparse suffix sorting

- $p_1, \ldots, p_m$: online, arbitrary order
- compare two suffixes with LCE query
- $c := \#$ characters to compare for sorting
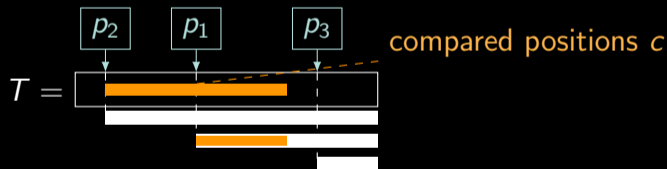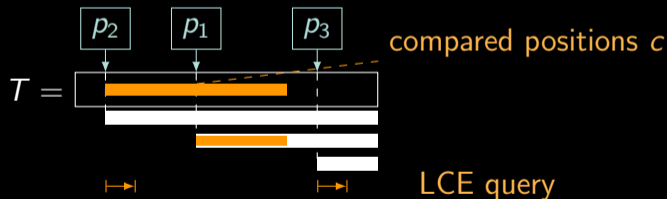


compared positions $c$

$T =$

LCE query

# dynamic sparse suffix sorting

- $p_1, \ldots, p_m$: online, arbitrary order
- compare two suffixes with LCE query
- $c := \#$ characters to compare for sorting
- how to store their order?



compared positions $c$

$T =$

# suffix binary search tree (SBST)



SBST of Irving and Love'03:
binary search tree representation

each node

- ◤ represents a position $p_i$
- ◤ stores a flag $\in \{\mathrm{L}, \mathrm{R}, \perp\}$
- ◤ the LCE with an ancestor

# running example

- ISA : inverse suffix array
- SA : suffix array
- LCP : LCP array

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T[i]$ | c | a | a | t | c | a | c | g | g | t | c | g | g | a | c |
| ISA$[i]$ | 6 | 1 | 4 | 14 | 7 | 3 | 9 | 12 | 13 | 15 | 8 | 11 | 10 | 2 | 5 |

| $r$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SA$[r]$ | 2 | 14 | 6 | 3 | 15 | 1 | 5 | 11 | 7 | 13 | 12 | 8 | 9 | 4 | 10 |
| LCP$[r]$ | 0 | 1 | 2 | 1 | 0 | 1 | 2 | 1 | 3 | 0 | 1 | 2 | 1 | 0 | 2 |

## problem definition

- obtain SA from in-order traversal in $\mathcal{O}(m)$ time.
- how to obtain LCP?

| r | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|
| SA[r] | 2 | 14 | 6 | 3 | 15 | 1 | 5 | 11 | 7 | 13 | 12 | 8 | 9 | 4 | 10 |
| LCP[r] | 0 | 1 | 2 | 1 | 0 | 1 | 2 | 1 | 3 | 0 | 1 | 2 | 1 | 0 | 2 |

# closest left/right ancestors

let $v$ be a node

- $\blacksquare$ $cla_v$ : lowest node having $v$ as a descendant in its left subtree
- $\blacksquare$ $cra_v$ : lowest node having $v$ as a descendant in its right subtree
- $\Rightarrow$ either $cla_v$ or $cra_v$ is $v$'s parent

# LCE value $m_v$ and flag $d_v$

- $ca_v := \text{argmax}_{u \in \{cla_v, cra_v\}} \text{lce}(v, u)$
- if $ca_v = cla_v$, then $m_v = \text{lce}(v, cla_v)$, $d_v = \text{L}$.
- if $ca_v = cra_v$, then $m_v = \text{lce}(v, cra_v)$, $d_v = \text{R}$.
- if $ca_v$ is undefined, then $m_v = 0$, $d_v = \bot$.

| r | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|
| SA[r] | 2 | 14 | 6 | 3 | 15 | 1 | 5 | 11 | 7 | 13 | 12 | 8 | 9 | 4 | 10 |
| LCP[r] | 0 | 1 | 2 | 1 | 0 | 1 | 2 | 1 | 3 | 0 | 1 | 2 | 1 | 0 | 2 |

rules:

e : neither left child nor $\text{cra}_v$ exists
$\Rightarrow \text{LCP}[\text{ISA}[v]] = 0$

| r | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SA[r] | 2 | 14 | 6 | 3 | 15 | 1 | 5 | 11 | 7 | 13 | 12 | 8 | 9 | 4 | 10 |
| rules | e | | | | | | | | | | | | | | |
| | 0 | | | | | | | | | | | | | | |

rules:

e : neither left child nor $\text{cra}_v$ exists
$\Rightarrow \text{LCP}[\text{ISA}[v]] = 0$

r : $d_v = \text{R} \Rightarrow \text{LCP}[\text{ISA}[v]] \geq m_v$

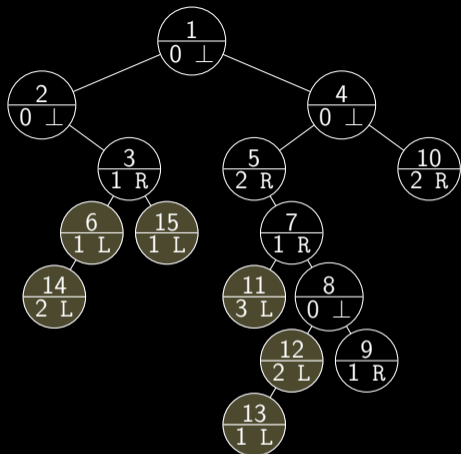| r | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|
| SA[r] | 2 | 14 | 6 | 3 | 15 | 1 | 5 | 11 | 7 | 13 | 12 | 8 | 9 | 4 | 10 |
| rules | e | | | r | | | r | | r | | | | r | | r |
| | 0 | | | 1 | | | 2 | | 1 | | | | 1 | | 2 |

rules:

e : neither left child nor $cra_v$ exists
$\Rightarrow LCP[ISA[v]] = 0$

r : $d_v = R \Rightarrow LCP[ISA[v]] \geq m_v$

l : $d_v = L$ and right subtree of $v$ is
empty $\Rightarrow LCP[ISA[v] + 1] = m_v$

| r | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SA[r] | 2 | 14 | 6 | 3 | 15 | 1 | 5 | 11 | 7 | 13 | 12 | 8 | 9 | 4 | 10 |
| rules | e | l | l | r | l | | r | l | r | l | l | | r | | r |
| | 0 | | 2 | 1 | | 1 | 2 | | 3 | | 1 | 2 | 1 | | 2 |

rules:

e : neither left child nor $\mathrm{cra}_v$ exists
$\Rightarrow \mathrm{LCP}[\mathrm{ISA}[v]] = 0$

r : $d_v = \mathtt{R} \Rightarrow \mathrm{LCP}[\mathrm{ISA}[v]] \geq m_v$

l : $d_v = \mathtt{L}$ and right subtree of $v$ is
empty $\Rightarrow \mathrm{LCP}[\mathrm{ISA}[v]+1] = m_v$

d : $v$ has left child $u \Rightarrow$ rightmost
node in $u$'s subtree determines
$\mathrm{LCP}[\mathrm{ISA}[v]]$

| $r$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SA[$r$] | 2 | 14 | 6 | 3 | 15 | 1 | 5 | 11 | 7 | 13 | 12 | 8 | 9 | 4 | 10 |
| rules | e | l | l | r | l | | r | l | r | l | l | | r | d | r |
| | 0 | | 2 | 1 | | 1 | 2 | | 3 | | 1 | 2 | 1 | 0 | 2 |

9 / 16

rules:

e : neither left child nor $cra_v$ exists $\Rightarrow$ LCP[ISA[$v$]] = 0

r : $d_v = \text{R} \Rightarrow$ LCP[ISA[$v$]] $\geq m_v$

l : $d_v = \text{L}$ and right subtree of $v$ is empty $\Rightarrow$ LCP[ISA[$v$] + 1] = $m_v$

d : $v$ has left child $u \Rightarrow$ rightmost node in $u$'s subtree determines LCP[ISA[$v$]]

a : otherwise: $cra_v$ determines LCP[ISA[$v$]]

| $r$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SA[$r$] | 2 | 14 | 6 | 3 | 15 | 1 | 5 | 11 | 7 | 13 | 12 | 8 | 9 | 4 | 10 |
| rules | e | a | l | r | a | | r | a | r | a | l | | r | d | r |
| LCP[$r$] | 0 | 1 | 2 | 1 | 0 | 1 | 2 | 1 | 3 | 0 | 1 | 2 | 1 | 0 | 2 |

- ◣ rules e, r, l can be computed in constant time per node.
- ◣ how to compute rules d and a?

rules:

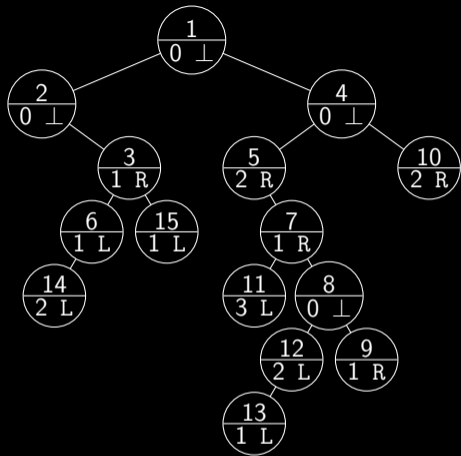  e : neither left child nor $\text{cra}_v$ exists $\Rightarrow \text{LCP}[\text{ISA}[v]] = 0$

  r : $d_v = \text{R} \Rightarrow \text{LCP}[\text{ISA}[v]] \geq m_v$

  l : $d_v = \text{L}$ and right subtree of $v$ is empty $\Rightarrow \text{LCP}[\text{ISA}[v] + 1] = m_v$

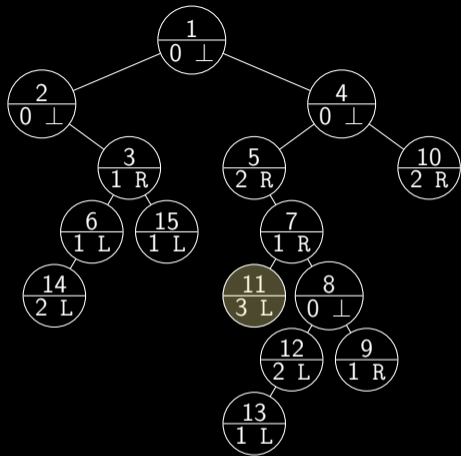  d : $v$ has left child $u \Rightarrow$ rightmost node in $u$'s subtree determines $\text{LCP}[\text{ISA}[v]]$

  a : otherwise: $\text{cra}_v$ determines $\text{LCP}[\text{ISA}[v]]$

| $r$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SA[$r$] | 2 | 14 | 6 | 3 | 15 | 1 | 5 | 11 | 7 | 13 | 12 | 8 | 9 | 4 | 10 |
| rules | e | a | l | r | a | | r | a | r | a | l | | r | d | r |
| LCP[$r$] | 0 | 1 | 2 | 1 | 0 | 1 | 2 | 1 | 3 | 0 | 1 | 2 | 1 | 0 | 2 |

task
compute LCP[11]

## task

compute LCP[11]

$= \text{lce}(\text{cra}_{11}, 11)$ since 11 has no left child

task

compute LCP[11]

- $= \mathrm{lce}(\mathrm{cra}_{11}, 11)$ since 11 has no left child
- $= \mathrm{lce}(\mathrm{cra}_{11}, \mathrm{cla}_{11})$ since $d_{11} = \mathrm{L}$ (proof later)

task

compute LCP[11]

$= \mathrm{lce}(\mathrm{cra}_{11}, 11)$ since 11 has no left child

$= \mathrm{lce}(\mathrm{cra}_{11}, \mathrm{cla}_{11})$ since $d_{11} = \mathrm{L}$ (proof later)

$= \mathrm{lce}(\mathrm{cra}_{11}, 7) = m_7 = 1$ since $\mathrm{cra}_{11} = \mathrm{cra}_7$.

## task

compute LCP[11]

- $= \mathrm{lce}(\mathrm{cra}_{11}, 11)$ since 11 has no left child
- $= \mathrm{lce}(\mathrm{cra}_{11}, \mathrm{cla}_{11})$ since $d_{11} = \mathrm{L}$ (proof later)
- $= \mathrm{lce}(\mathrm{cra}_{11}, 7) = m_7 = 1$ since $\mathrm{cra}_{11} = \mathrm{cra}_7$.
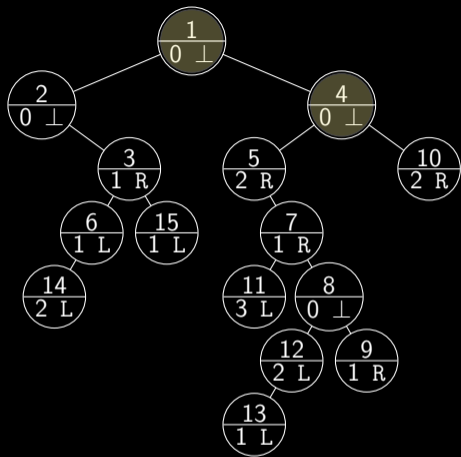- ◣ goal: maintain $\mathrm{lce}(\mathrm{cra}_v, \mathrm{cla}_v)$ for each node $v$ to process

## stack $S$

maintain stack $S$ of LCE values such that, on visiting node $v$, $S$ stores $\text{lce}(\text{cla}_u, \text{cra}_u)$ of all ancestors $u$ of $v$.

$S = \{$

$\quad \text{lce}(\text{cra}_1, \text{cla}_1) = 0,$

$\}$

### stack $S$

maintain stack $S$ of LCE values such that, on visiting node $v$, $S$ stores $\text{lce}(\text{cla}_u, \text{cra}_u)$ of all ancestors $u$ of $v$.

$S = \{$
$\quad \text{lce}(\text{cra}_1, \text{cla}_1) = 0,$
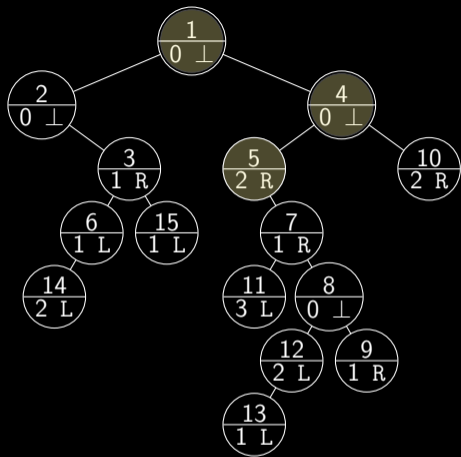$\quad \text{lce}(\text{cra}_4, \text{cla}_4) = 0,$

$\}$

## stack $S$

maintain stack $S$ of LCE values such that, on visiting node $v$, $S$ stores $\text{lce}(\text{cla}_u, \text{cra}_u)$ of all ancestors $u$ of $v$.

$S = \{$
$\quad \text{lce}(\text{cra}_1, \text{cla}_1) = 0,$
$\quad \text{lce}(\text{cra}_4, \text{cla}_4) = 0,$
$\quad \text{lce}(\text{cra}_5, \text{cla}_5) = 0,$

$\}$

## stack $S$

maintain stack $S$ of LCE values such that, on visiting node $v$, $S$ stores $\text{lce}(\text{cla}_u, \text{cra}_u)$ of all ancestors $u$ of $v$.
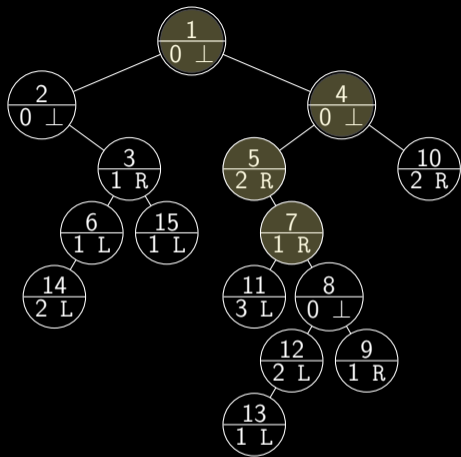
$S = \{$

$\quad \text{lce}(\text{cra}_1, \text{cla}_1) = 0,$

$\quad \text{lce}(\text{cra}_4, \text{cla}_4) = 0,$

$\quad \text{lce}(\text{cra}_5, \text{cla}_5) = 0,$

$\quad \text{lce}(\text{cra}_7, \text{cla}_7) = 0,$

$\}$

## stack $S$

maintain stack $S$ of LCE values such that, on visiting node $v$, $S$ stores $\text{lce}(\text{cla}_u, \text{cra}_u)$ of all ancestors $u$ of $v$.
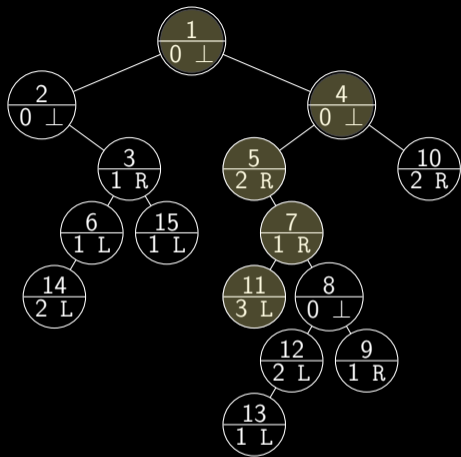
$S = \{$

$\quad \text{lce}(\text{cra}_1, \text{cla}_1) = 0,$

$\quad \text{lce}(\text{cra}_4, \text{cla}_4) = 0,$

$\quad \text{lce}(\text{cra}_5, \text{cla}_5) = 0,$

$\quad \text{lce}(\text{cra}_7, \text{cla}_7) = 0,$

$\quad \text{lce}(\text{cra}_{11}, \text{cla}_{11}) = 1$

$\}$

## stack $S$

maintain stack $S$ of LCE values such that, on visiting node $v$, $S$ stores $\text{lce}(\text{cla}_u, \text{cra}_u)$ of all ancestors $u$ of $v$.

$S = \{$
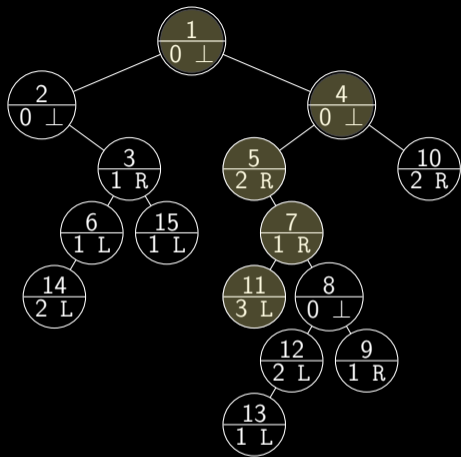
$\quad \text{lce}(\text{cra}_1, \text{cla}_1) = 0,$

$\quad \text{lce}(\text{cra}_4, \text{cla}_4) = 0,$

$\quad \text{lce}(\text{cra}_5, \text{cla}_5) = 0,$

$\quad \text{lce}(\text{cra}_7, \text{cla}_7) = 0,$

$\quad \text{lce}(\text{cra}_{11}, \text{cla}_{11}) = 1$

$\}$

◣ why helpful?

◣ how computable?

### known facts

1. $u, v, w \in [1..n]$ with $T[u..] \prec T[v..] \prec T[w..]$
   $\Rightarrow \mathrm{lce}(u, w) = \min(\mathrm{lce}(u, v), \mathrm{lce}(v, w))$
2. $T[\mathrm{cra}_v..] \prec T[v..] \prec T[\mathrm{cla}_v..]$ (assume $\mathrm{cla}_v$ and $\mathrm{cra}_v$ exist)

### lemma

given

- $\mathrm{lce}(\mathrm{cla}_v, \mathrm{cra}_v)$ and
- $m_v = \mathrm{lce}(v, \mathrm{ca}_v)$,

we can compute

- $\mathrm{lce}(v, \mathrm{cla}_v)$ and
- $\mathrm{lce}(v, \mathrm{cra}_v)$ in constant time.

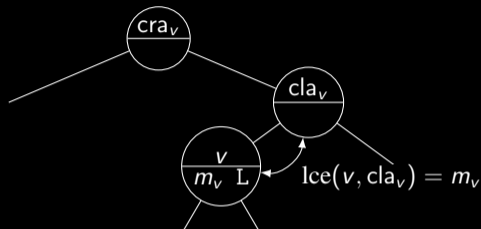## proof of lemma

- wlog., $d_v = L$, and $cla_v$ and $cra_v$ exist

$\Rightarrow ca_v = cla_v$

hence:

- $lce(v, cla_v) = lce(v, ca_v) = m_v$
- $lce(v, cra_v) = lce(cla_v, cra_v)$

the latter is because of Facts 1 and 2:

$$lce(cra_v, cla_v) = \min(lce(v, cra_v), lce(v, cla_v))$$
$$= lce(v, cra_v) \leq lce(v, cla_v)$$



$\square$

# corollary: how to compute stack $S$

given:
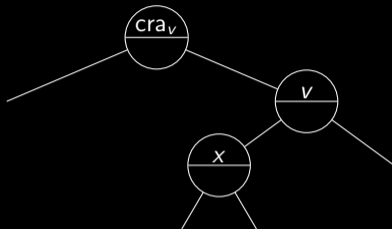
- value $\text{lce}(\text{cla}_v, \text{cra}_v)$
- $x$: $v$'s left child

then:

- $\text{cla}_x = v$ and $\text{cra}_x = \text{cra}_v$
- $\Rightarrow \text{lce}(\text{cla}_x, \text{cra}_x) = \text{lce}(v, \text{cra}_v)$
  computable in constant time by
  lemma

(right child analogously by symmetry)

$\Rightarrow$ can maintain stack $S$ during a top-down traversal in constant time per node.

# subarray extraction

can compute SLCP$[\ell, r]$ in $\mathcal{O}(h + (r - \ell))$ time, where $h$ is the tree's height.

- ◣ augment tree with subtree sizes
- ◣ can find node $\ell$ by top-down traversal (while maintaining $S$)
- ◣ can start in-order traversal at node $\ell$
- ◣ stop traversal when arriving at node $r$
- ◣ number of visited nodes is $\mathcal{O}(h) + r - \ell$, and each node is processed in constant time.

## summary

suffix binary search tree by Irving and Love'03

- maintains ranks of $m$ suffixes
- $\mathcal{O}(m)$ space (each node stores 2 integers + 1 bit)
- construction needs $\mathcal{O}(mh)$ LCE queries ($h$: height)
- can be made balanced ($h = \mathcal{O}(\lg m)$)
- used for sparse suffix sorting by Fischer+'20
  - $\mathcal{O}(c(\sqrt{\lg \sigma} + \lg \lg n) + m \lg m \lg n \lg^* n)$ time
  - $c$: lower bound on number of characters needed to compare
  - $\mathcal{O}(m)$ space                    ($n$ : text length, $\sigma$: alphabet size)

our contribution: can extract

- SSA$[i..i + \ell - 1]$ and
- SLCP$[i..i + \ell - 1]$ in $\mathcal{O}(h + \ell)$ time

any questions are always very welcome!

# open problems

- memory-efficient representations of suffix binary search trees?
- time-efficient implementation via B trees
  - balanced by construction
  - B+ variants have good memory locality
- can we merge two trees efficiently?