

Compressed Range-Minimum Queries

Average-Case Analysis of Search Trees

Meets Space-Efficient Data Structures

Sebastian Wild

wild@uwaterloo.ca

based on joint work with
Ian Munro & Patrick Nicholson

UNIVERSITY OF
WATERLOO



My goals (for you!)

1

Showcase new applications of average-case analysis

You know analysis.

→ You will like this.

2

Give flavor of supporting queries on compressed data

- very rough ideas & black boxes
- bit fiddling
- I will treat $n \cdot \frac{\log \log n}{\log n}$ as "much smaller" than n ...

→ tougher ride, but cool stuff

2

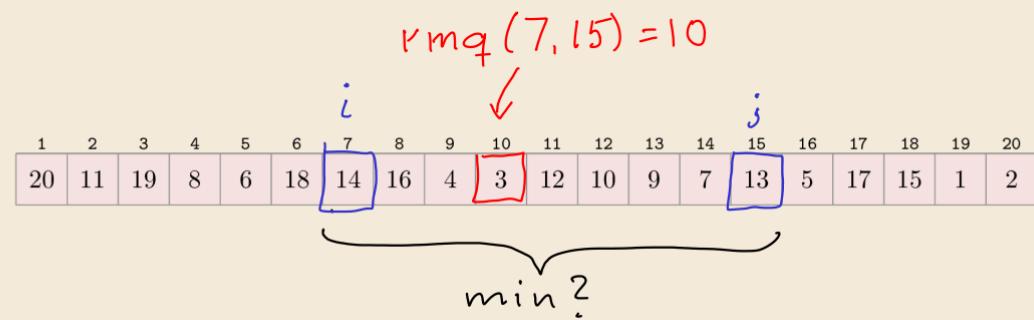
Overview

- ① Range minimum queries & Cartesian trees
- ② Random RMQ & random BSTs
- ③ Tree covering
- ④ Entropy trees

1 Range-Minimum Queries

The RMQ problem:

Given: array $A[1..n]$ of numbers



Task: ① Preprocess A

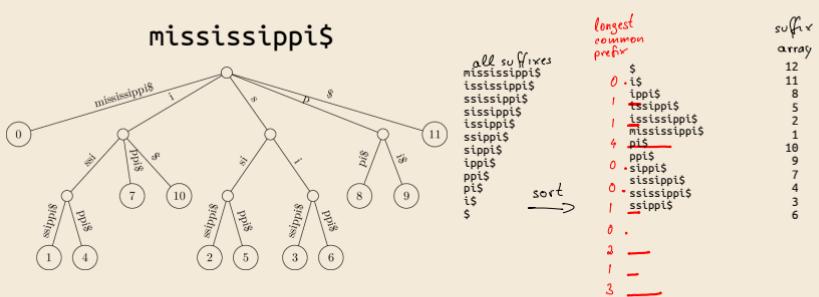
- ② Answer $\text{rmq}(i, j) = \text{position of minimum in } A[i..j]$
 $= \arg \min_{i \leq k \leq j} A[k]$

Applications of RMQ

longest common extensions in strings



mississippi\$

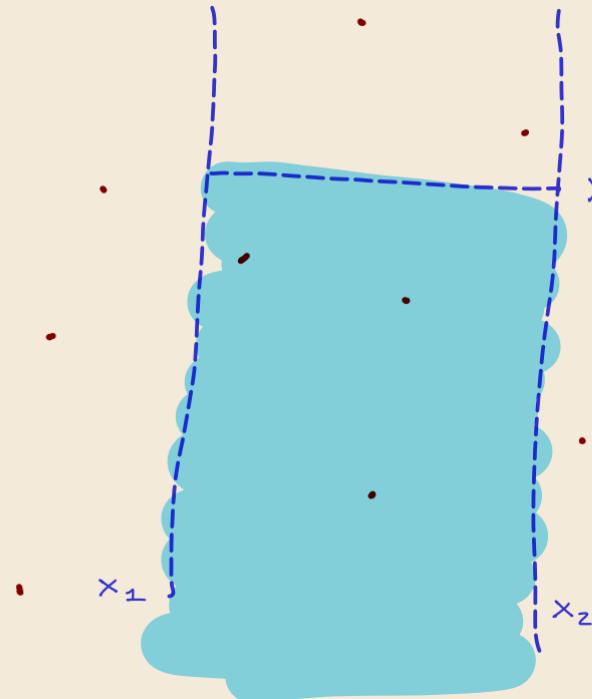


longest common extension of two positions

$$\text{lce}(i, j) = \text{LCP}[\text{rmq}_{\text{LCP}}(i+1, j)]$$

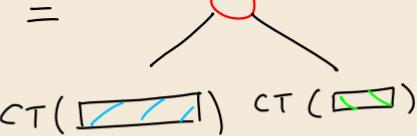
$$\text{LCP}[i] = \text{lcp}(i^{\text{th}} \text{ min } (i-1)^{\text{th}}, \dots, 1^{\text{th}})$$

3-sided orthogonal range queries in 2D



Cartesian Trees

$CT(\boxed{\text{min}})$

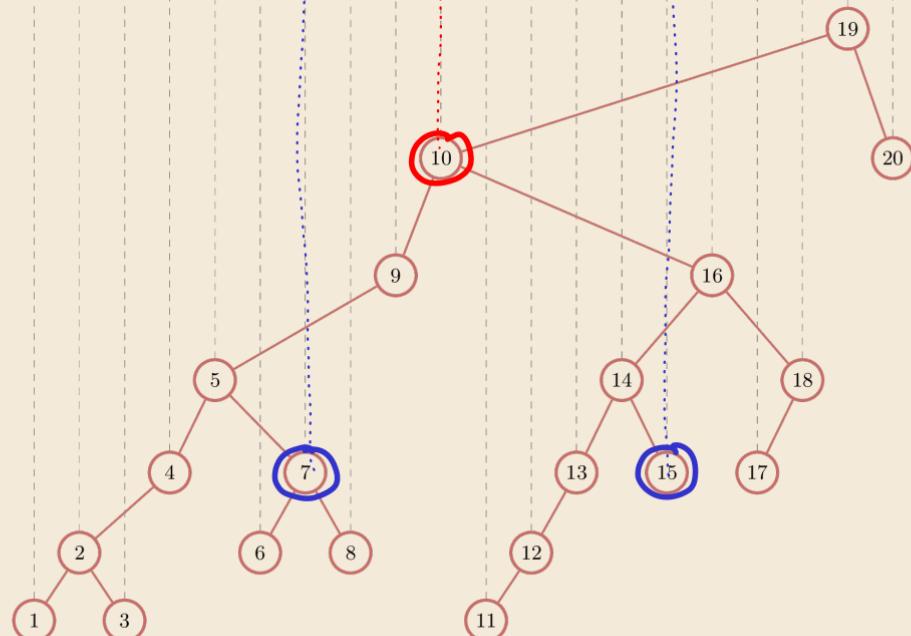


$\text{rmq}(i, j) =$

lowest common ancestor (lca)
of nodes with inorder i and j

$\text{rmq}(7, 15) = 10$

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
20	11	19	8	6	18	14	16	4	3	12	10	9	7	13	5	17	15	1	2



Elementary Solutions

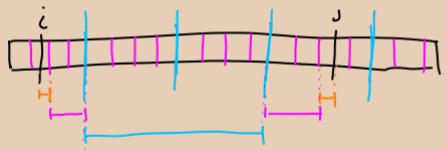
query
time

trivial: store nothing
scan $A[i..j]$ $O(n)$

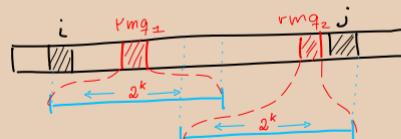


still need original array A

\Rightarrow must store
 n integers



blocking, "4 Russians"
& sparse tables
 $O(1)$ $\Theta(n)$



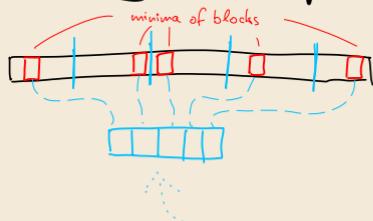
sparse tables
 $O(1)$ $\Theta(n \log n)$

trivial:
precompute all
 $O(1)$ $\Theta(n^2)$

Space usage
(integers)

Details on $O(n)$ Space RMQ

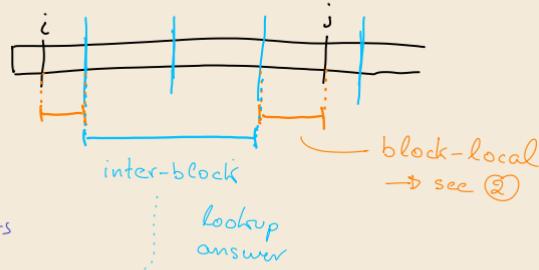
① Blocking / Input Reduction



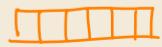
reduced instance $O\left(\frac{n}{\log^3 n}\right)$

↳ use $O(n \log^2 n)$ bit solution:
(sparse tables)

$$O\left(\frac{n}{\log^3 n} \cdot \log^2\left(\frac{n}{\log^3 n}\right)\right) = O\left(\frac{n}{\log n}\right) \text{ bits}$$



② 4-Russians base case



block-local queries: $\leq 4^{\frac{1}{4} \lg n} = \sqrt{n}$ different RMQ answer set

$$\leftarrow \frac{1}{4} \lg n$$

↳ store { pre-computed answers for all possible blocks $O(n)$ }

{ the type of each block (using $2 \cdot \frac{1}{4} \lg n$ bits each) $\frac{n}{\frac{1}{4} \lg n} \cdot 2 \cdot \frac{1}{4} \lg n = 2n$

Cartesian trees

Block type	Range answers
00000000	00000000
00000001	...
00000010	...
00000011	...
:	:
11111110	...
11111111	...

③ Bootstrapping

① needs block size $\omega(\log^2 n)$

② needs block size $< \frac{1}{8} \lg n$

⇒ 2-level blocking



middle level: $\frac{n}{\log^3 n}$ blocks each size $\frac{\log^3 n}{\frac{1}{8} \lg n}$

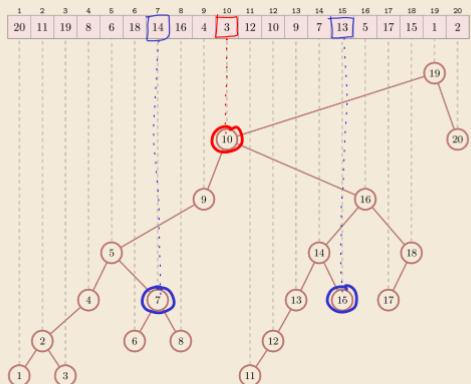
$$O\left(\frac{n}{\log^3 n} \cdot \log^2 n \log\left(\log^2 n\right)\right) = O\left(n \frac{\log \log n}{\log n}\right)$$

Succinct RMQ

goal : minimal-space
asymptotically

data structure with $O(1)$ -time rmq
on word-RAM

no access to A at query time
(encoding model)



- ⇒ can use any data structure for binary trees with
- access to nodes by **inorder index**
 - **lca** of 2 nodes
 - $\lg(C_n) = 2n - \Theta(\log n)$
 - $2n + o(n)$ bits of space optimal
- $\} O(1)$

E level-order unary degree sequence (LOUDS)

E balanced parentheses (BP)

E depth-first unary degree sequence (DFUDS)

E tree covering (TC)

Excursion: Machine Models

logarithmic-cost model

think: Turing machine

numbers $\gg 2^{64}$

machine fixed,
unbounded input

vs

?

vs

vs

unit-cost model

#arithmetic operations

64 bit numbers

fixed-size input

$w = \Theta(1)$

word-RAM model

transdichotomous model

constant-time ops on w bit numbers

w grows with n

$w = \Theta(\log n)$

the machine
grows with
the input!

$w \geq$ size of
numbers
in input

standard model
in algorithm theory

less powerful

more powerful

2 Random RMQ



Can we save space for "typical" RMQ instances?

without compromising query time

on average

natural model:

random-permutations model

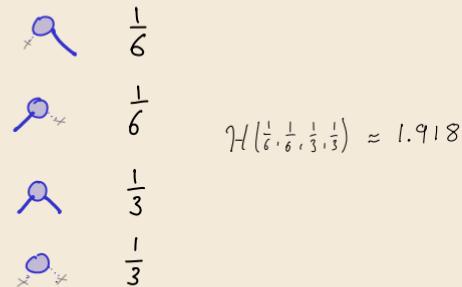
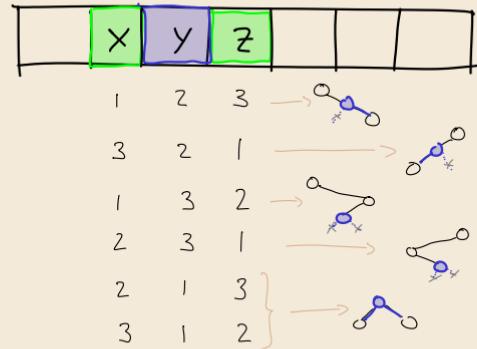
all orderings equally likely

equivalent: $A[i]$ iid real random variables

Previous Work

- local order / node degrees

Davoodi, Navarro, Raman, Rao 2014



→ solves RMQ with average space $1.918n + o(n)$

Can we do better?

- compressible array
empirical entropy

Fischer, Heun 2011

- few runs in array

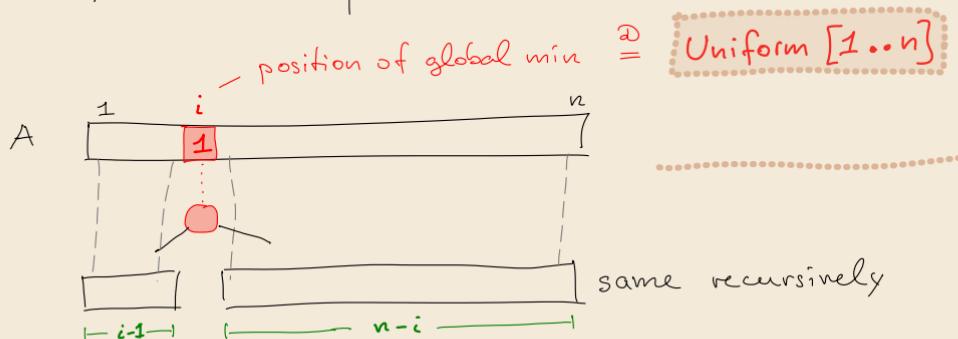
Barbay, Fischer, Navarro 2012

}

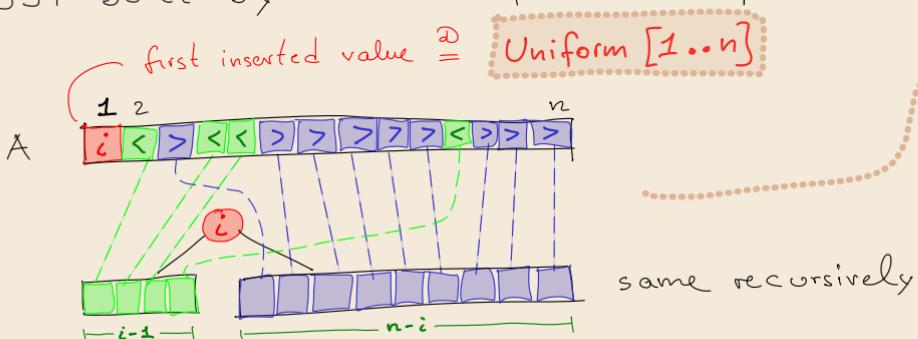
won't help on random arrays

RMQ = BST!

CT for random permutation:



BST built by insertions from random permutation



$$\text{CT}(A) = \text{BST}(A^{-1})$$

inverse function

same shape
distribution
for random perm.

Structural Complexity
of Random Binary Trees

J. C. Kieffer, E-H. Yang, and W. Szpankowski

entropy
 $\sim 1.736n$ bits

Subtree-Size Entropy

$$\Pr_v[T=t] = \prod_{v \in t} \frac{1}{st(v)}$$

prob that among all values in v's subtree
v is inserted first

$$\begin{aligned}
 \text{entropy of } T \text{'s distribution} &= \sum_t \Pr[T=t] \cdot \lg \frac{1}{\Pr[T=t]} \\
 &= \sum_t \lg \left(\prod_{v \in t} st(v) \right) \\
 &= \sum_t \lg \left(\sum_{v \in t} st(v) \right) =: H_{st}(t) \text{ "subtree-size entropy"} \\
 &= \mathbb{E}[H_{st}(t)].
 \end{aligned}$$



If we can encode trees using $H_{st}(t) + l.o.t.$ bits,
we use optimal average space $\stackrel{D}{\sim}$
→ matches entropy

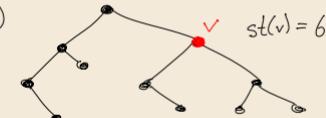
Notation:

t binary tree (shape)

v node in t

$st(v)$ size of v 's subtree

T random BST (shape)



Remark:

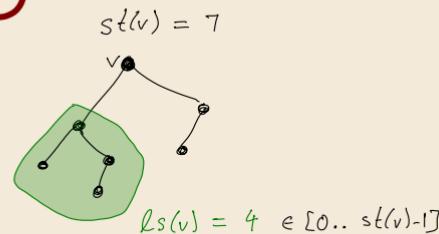
subtree sizes are **global** property
previous methods exploit
local properties

Left-subtree-size encoding

Idea: store the sizes of all left subtrees

for node v : $ls(v) \in [0..st(v)-1]$

uniform dist. in random BSTs?
↳ local entropy $\lg(st(v))$



encode nodes in preorder traversal $\Rightarrow st(v)$ already encoded / known

store $ls(v)$ with arithmetic coding

and compression scheme with $H_{st}(t) + 2$ bits



How to access specific nodes?

How to compute lca?

Bounding the worst case

⚠ H_{st} can be $\sim n \lg n$

use 1 bit to specify:

0 → above encoding

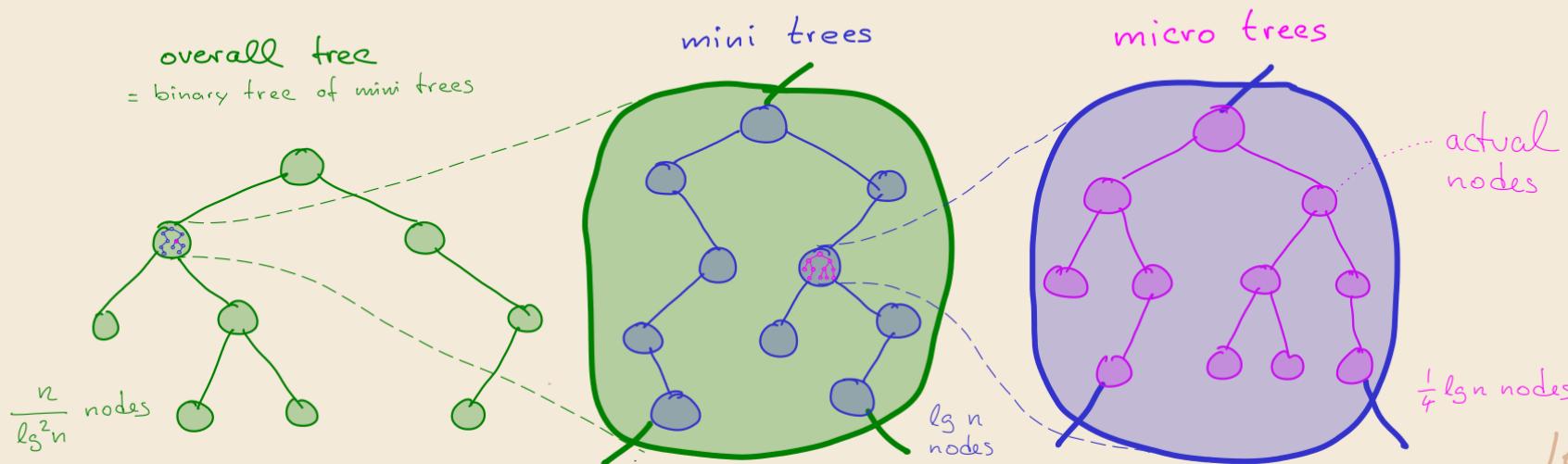
1 → standard $2n$ bit encoding

$\min\{H_{st} + 3, 2n + 1\}$ bit

3 Tree Covering

technique for succinct tree data structure
using 3-level hierarchical decomposition

Geary, Raman, Raman 2006
Farzan, Munro 2012

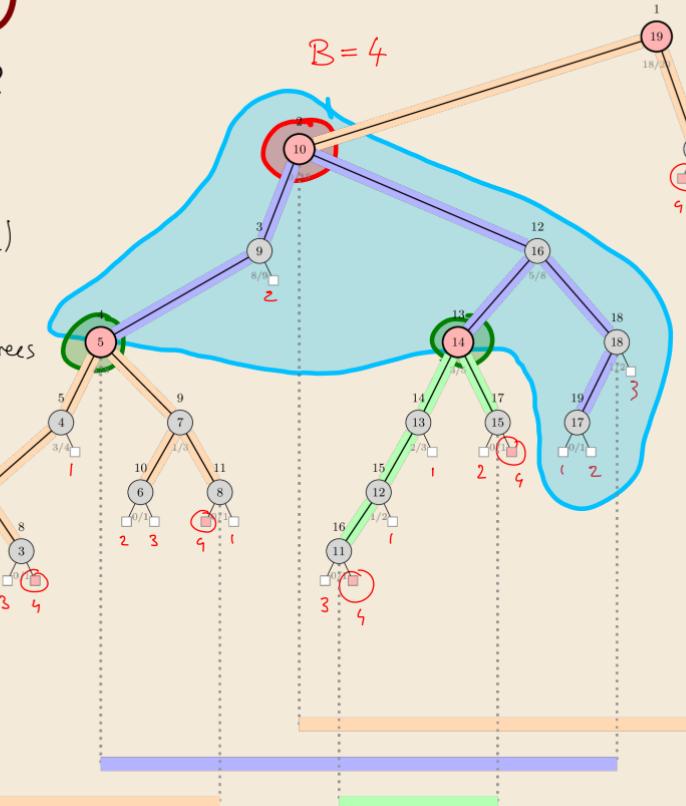
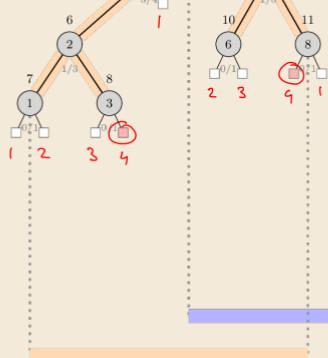


Partitioning Binary Trees

How to define subtrees?

fix parameter B

- ① mark every B th leaf (in inorder)
- ② mark all internal nodes with marked leaf in both subtrees
(+ always mark root)
- ③ marked nodes = roots of blocks

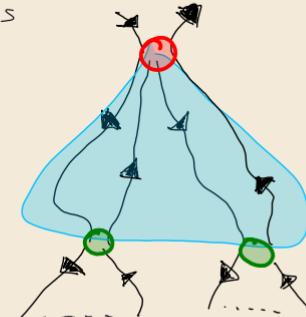


$B=4$

extended binary tree

Properties:

- (a) $\approx \frac{n}{B}$ subtrees
- (b) between 2 and $3B$ nodes
- (c) complete subtree
except for two **portal nodes**
- (d) traversals visit each subtree
 ≤ 3 times



Operations on TC representation



node ids "τ-names"

(τ_1, τ_2, τ_3)

mini tree id

micro tree id

Geary, Raman, Raman 2006

lowest common ancestor

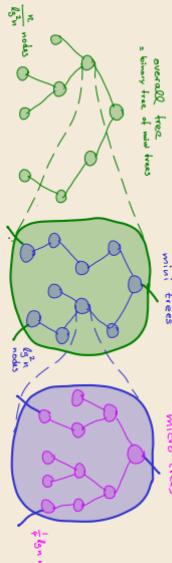


He, Munro, Rao 2012

simplifications

Farzan, Munro 2012

all operations possible in $O(1)$ time w/ $O(n)$ indices



$O(\frac{n}{\log^2 n})$ nodes

↳ can afford pointer/word based representation
 $\rightarrow O(n \log n)$ bit

$O(\frac{n}{\log^2 n})$ trees of $O(\log n)$ nodes each

$O(n)$

↳ as above!

$\leq \frac{1}{4} \log n$ nodes each to encode in $\leq \frac{1}{2} \log n$ bit each

↳ use lookup table by type
 for all local operations

\rightarrow only $2^{\frac{1}{2} \log n} = \sqrt{n}$
 different types

⊕ store type for all microtrees ← dominant space



node id \leftrightarrow (global) preorder

Geary, Raman, Raman 2006



node id \leftrightarrow (global) inorder

Davoodi, Navar, Raman, Rao 2014

4 Entropy Trees

use left-subtree-size encoding for micro trees

+ compressed bitvector to find
starting position of i th entry

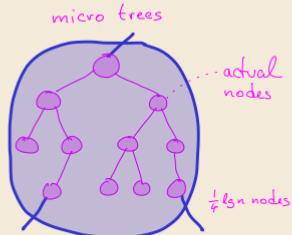


$$\sum_{\text{micro tree } \mu} |\text{type}(\mu)| \leq \sum_{\mu} \left(\sum_{v \in \mu} st_{\mu}(v) + 2 \right) \leq \sum_{\mu, v \in \mu} st(v) + o(n)$$
$$\leq H_{st}(t) + o(n).$$

Theorem: We can support all TC operations on static binary trees in constant time using $\min\{H_{st}, 2n\} + o(n)$ bits.

Corollary: There is a RMQ datastructure with constant query time that occupies the optimal $1.736n + o(n)$ bits on average.

Hypersuccinct Trees



dominant space: type for all microtrees

construction works with any (bounded) code for micro trees

so why not use an instance-optimal one?

Huffman code for micro-tree types



never worse, but unclear how good

→ analysis?

Conclusion

a new AoFA playground?



1 Average-case analysis helps for average-succinct data structures.

Range-min queries is one example. There are more!

optimal avg space ✓

2 Known data-structures techniques work for compressed trees.

How about other combinatorial structures? Other input distributions?

3 Theoretically optimal data structures are often impractical.

Better trade-offs? Simpler solutions?