Wavelet Trees for All

Gonzalo Navarro

www.dcc.uchile.cl/gnavarro
gnavarro@dcc.uchile.cl

Department of Computer Science University of Chile



R. Grossi, A. Gupta, J.S. Vitter. High-order entropy-compressed text indexes. *SODA*, Jan 2003.



B. Chazelle. A functional approach to data structures and its use in multidimensional searching. *SIAM J. Comp.* 1988.

イロト イヨト イヨト イヨト

- 32





(日) (종) (종) (종) (종)

- 2



Wavelet + Tree = Wavelet Tree?



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで



◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 臣 のへで



◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 臣 の�?



▲□▶ ▲□▶ ▲三▶ ▲三▶ ▲□▶

Basic Structure

Basic Functionality

Compression



Basic Structure

Basic Functionality

Compression



Basic Structure

Basic Functionality

Compression



Basic Structure

Basic Functionality

Compression



Three Views

Applications as Sequences

Applications as Reorderings

- 2

Three Views

Applications as Sequences

Applications as Reorderings

Three Views

Applications as Sequences

Applications as Reorderings

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Three Views

Applications as Sequences

Applications as Reorderings

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Basic Structure

Basic Functionality

Compression



Basic Wavelet Tree Structure

- Built on a sequence S[1, n] over alphabet $[1, \sigma]$.
- A binary and perfectly balanced tree.
- Each node represents a range [a, b] of $[1, \sigma]$.
- The root represents the whole [1, σ].
- The children of a node v split its alphabet range by half.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

• Each leaf represents some $c \in [1, \sigma]$.

Basic Wavelet Tree Structure

- Associate with each node v a sequence S_v[1, n_v] (not stored).
- Let $[a, b] \subseteq [1, \sigma]$ the alphabet range of v.
- ► Then S_V is the subsequence of S formed by the characters in [a, b].
- At each node v we store a bitmap $B_v[1, n_v]$.
- $B_{\nu}[i] = 0$ iff $S_{\nu}[i]$ belongs to the alphabet of the left child.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

Basic Wavelet Tree Structure



◆□▶ ◆□▶ ▲目▶ ▲目▶ ▲□▶

Basic Space Analysis

- The tree is balanced, so there are $h = \lfloor \lg \sigma \rfloor$ levels.
- We store at most *n* bits per level.
- Thus there are at most $n \lceil \lg \sigma \rceil$ bits.
- Plus $O(\sigma \lg n)$ bits for the tree pointers.
- Total: $n \lg \sigma + O(n) + O(\sigma \lg n)$.
- This is close to the space needed to represent S[1, n] in plain form!

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

Basic Structure

Basic Functionality

Compression



Basic Functionality: Tracking Symbols

- Start at a position *i* at the root bitmap *B_{root}*.
- Where has it gone, left or right?
 - Depends on whether $B_{root}[i] = 0$ or 1.
- And at which position has it been mapped?
 - If it went left, to i_0 = number of 0s up to *i* in B_{root} .
 - If it went right, to $i_1 =$ number of 1s up to *i* in B_{root} .
- Continue recursively downwards.
- When we arrive at the leaf of symbol c, it turns out that S[i] = c.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

Basic Functionality: Tracking Symbols

Values i₀ and i₁ can be defined via operation

 $rank_b(B, i) =$ number of occurrences of bit b in B[1, i]

- Operation rank can be computed in constant time...
- ... by preprocessing *B* and storing o(|B|) further bits.
- Thus we obtain any S[i] in time $O(\lg \sigma)$.
- The space raises to $n \lg \sigma + o(n \lg \sigma) + O(n) + O(\sigma \lg n)$
- Thus the wavelet tree replaces S, more or less within the same space.

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 臣 の�?

Basic Functionality: Tracking Symbols



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

Basic Functionality: Tracking Upwards

- Similarly, we can start at a position *j* of the leaf of a symbol *c* and go upwards.
- Where is the current position *i* in the parent node v?
 - ► If my node is a left child of v, at j_0 = position where the *i*th 0 occurs in B_v .
 - ► If my node is a right child of v, at j_1 = position where the *i*th 1 occurs in B_v .

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

- Continue recursively up to the root.
- At the root, we have found the position in S of the *j*th *c*.

Basic Functionality: Tracking Upwards

Values j₀ and j₁ can be defined via operation

 $select_b(B, j) = position of the jth b in B$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

- Operation select can be also computed in constant time...
- ... by preprocessing **B** and storing o(|B|) further bits.
- Thus we track upwards also in time $O(\lg \sigma)$.
- Indeed, we solve on S the generalization of select_b to sequences.

Basic Functionality: Tracking Upwards



▲ロト ▲母ト ▲ヨト ▲ヨト ヨー のへで

Some Technical Improvements

- The tree pointers can be eliminated by concatenating all the bitmaps of the same level: O(σ lg n) disappears from the space.
- The total space can be reduced to n lg σ + o(n) by using recent low-redudnancy bitmap representations.
- The downward traversal can be sped up to O(lg σ/ lg w), in a machine of w bits, by using multiary trees.
- The upward traversal can be sped up to O(lg^ε σ), at the price of O((1/ε)nlg σ) bits of space, by using long upward pointers.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Basic Structure

Basic Functionality

Compression



- When bitmaps have many more 0s that 1s, or viceversa, they can be compressed.
- If B[1, n] has n_0 0s and n_1 1s, then

$$H_0(B) = \frac{n_0}{n} \log \frac{n}{n_0} + \frac{n_1}{n} \log \frac{n}{n_1}$$

- ► This is lower when $n_0 << n_1$ or viceversa, and maximum $(H_0(B) = 1)$ if $n_0 = n_1$.
- This notion can be extended to sequences S[1, n]:

$$H_0(S) = \sum_{c \in [1,\sigma]} \frac{n_c}{n} \lg \frac{n}{n_c}$$

where n_c is the frequency of c in S.

H₀(S) is lower when some symbols are more frequent than others, and it always holds H₀(S) ≤ lg σ.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

- There exist constant-time rank/select capable bitmap representations that require nH₀(B) + o(n) bits.
- ► What happens if we use them on the bitmaps B_v of a wavelet tree?
- Say $B_{root}[1, n]$ has n_0 0s and n_1 1s, then the first level uses

$$n_0 \lg \frac{n}{n_0} + n_1 \lg \frac{n}{n_1}$$
 bits

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

- Say v_0 is the left child of the root and v_1 its right child.
- Says $B_{v_0}[1, n_0]$ has n_{00} 0s and n_{01} 1s.
- Then the bitmap B_{v_0} uses

$$n_{00} \lg \frac{n_0}{n_{00}} + n_{01} \lg \frac{n_0}{n_{01}}$$
 bits.

Similarly, B_{v1} uses

$$n_{10} \log \frac{n_1}{n_{10}} + n_{11} \log \frac{n_1}{n_{11}}$$
 bits.

 Adding up the three bitmaps, we have that the first two levels use

$$n_{00} \lg \frac{n}{n_{00}} + n_{01} \lg \frac{n}{n_{01}} + n_{10} \lg \frac{n}{n_{10}} + n_{11} \lg \frac{n}{n_{11}}$$
 bits

• This would be $nH_0(S)$ if $\sigma = 4!$

- It is not hard to see that this holds in general: the sum of all n_vH(B_v) is nH₀(S).
- The total space can be as good as $nH_0(S) + o(n)$ bits.
- Previous time complexities are maintained.
- Thus the wavelet tree represents S in compressed form.

Basic Structure

Basic Functionality

Compression


- Zero-order entropy can also be achieved with a radically different technique.
- Instead of a balanced shape, give the wavelet tree a Huffman shape.
- Use the frequencies of the symbols in S.
- One can use plain (and faster) bitmap representations.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

• The total number of bits represented is less than $n(H_0(S) + 1)$.



- Average times are better than O(lg σ): they can be O(1 + H₀(S)) under some conditions.
- Worst case times can worsen, but they can be brought back to O(lg σ) at no asymptotic price, by rebalancing deep subtrees.
- Huffman shape gives a more practical way to reduce the redundancy, from o(n lg σ) to o(n(H₀(S) + 1), even if in theory o(n) can be obtained.

• It requires $O(\sigma \lg n)$ bits of space again.

- Note that this can be used on any variable-length encoding, not only Huffman.
- For example, a sequence of δ -codes.
- The Wavelet Tree would require roughly the same number of bits...
- ... yet it would allow extracting the code of the *i*th symbol in time proportional to its length.

Part II: Applications

Three Views

Applications as Sequences

Applications as Reorderings

Applications as Grids

Three Views

A wavelet tree on a sequence of symbols can be regarded in three ways:

- As representing a sequence.
- As representing a reordering.
- As representing a grid of points.

As a Sequence

- As we have presented it till now.
- It represents S[1, n] within compressed space.
- It supports operations
 - access, downward tracking.
 - select, upward tracking.
 - rank, defined as

 $rank_c(S, i) =$ number of occurrences of symbol c in S[1, i]

As a Sequence: Rank Operation

- Start at a position *i* at the root bitmap *B_{root}*.
- Where do cs go, left or right?
 - If left, continue to left child at position $i \leftarrow i_0$.
 - ▶ If right, continue to right child at position $i \leftarrow i_1$.
- ▶ When we arrive at the leaf of symbol *c*, *i* is the answer.

As a Sequence: Rank Operation



・ロト ・ 日 ・ ・ 田 ・ ・ 日 ・ うへで

As a Reordering

- We have a sequence on the top...
- ... and a stable sorting of its elements at the bottom.
- The wavelet tree describes the sorting process.
- For example, it can represent a permutation π of [1, n] using n lg n + o(n) bits.
- It supports both $\pi(i)$ and $\pi^{-1}(i)$ in time $O(\lg n / \lg \lg n)$.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

As a Reordering: Permutations



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三 つへで

As a Grid of Points

- Simplest case: $n \times n$ grid with *n* points.
- Exactly one point per row and one per column.
- More general cases require (easy) mappings from real coordinates.
- The main operation in these applications: count/report the points in a rectangle.

As a Grid of Points



◆□▶ ◆□▶ ◆□▶ ◆□▶ 三 のへの

As a Grid of Points

Let (x₁, y₁), (x₂, y₂), ..., (x_n, y_n) be the points sorted by x-coordinate.

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○□ のへで

- Consider the wavelet tree of string $S = y_1 y_2 \dots y_n$.
- Then the *i*th point in *x*-order is (i, S[i]).
- And the *i*th point in *y*-order is (*select_i(S*), *i*).

- To count or report the points within $[x_1, x_2] \times [y_1, y_2]$.
- Start at the root with interval $[x_1, x_2]$.
- Project it into the left and into the right children.
- Stop when
 - The current interval $[x_1, x_2]$ is empty.
 - The current "alphabet" interval does not intersect [y1, y2].
 - The current alphabet interval is contained in $[y_1, y_2]$.
- In the latter case, count or report all the elements in $[x_1, x_2]$.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

- Counting time: $O(\lg n / \lg \lg n)$.
- Reporting time $O((k + 1) \lg^{\epsilon} n)$.



= 990



▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のQで



▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のへの

Part II: Applications

Three Views

Applications as Sequences

Applications as Reorderings

(日) (图) (문) (문) (문)

Applications as Grids

- ► A data structure over a text *T*[1, *n*] supporting pattern searches.
- Most classical: suffix trees and suffix arrays.
- The latter are simpler and less space-consuming.
- (Still they do not support all the functionality of suffix trees.)



イロト イヨト イヨト イヨト

- A way to simulate the suffix array A[1, n] in little space uses the Burrows-Wheeler Transform (BWT) of T, T^{bwt}.
- It concatenates the symbols preceding each suffix of A.
- It is a reversible permutation of the string T.
- To reverse it, one needs a limited version of rank on T^{bwt}.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

Ъ

alabar a la alabarda\$ labar a la alabarda\$a abar a la alabarda\$al bar a la alabarda\$ala ar a la alabarda\$alab r a la alabarda\$alaba a la alabarda\$alabar a la alabarda\$alabar la alabarda\$alabar a la alabarda\$alabar a a alabarda\$alabar a l alabarda\$alabar a la alabarda\$alabar a la labarda\$alabar a la a abarda\$alabar a la al barda\$alabar a la ala arda\$alabar a la alab rda\$alabar a la alaba da\$alabar a la alabar a\$alabar a la alabard \$alabar a la alabarda \$alabar a la alabarda a la alabarda\$alaba alabarda\$alabar a la la alabarda\$alabar a\$alabar a la alabard a alabarda\$alabar a a la alabarda\$alaba abar a la alabarda\$al abarda\$alabar a la alabar a la alabarda alabarda\$alabar a la ar a la alabarda\$alab arda\$alabar a la alab bar a la alabarda\$ala barda\$alabar a la ala 2nd "r" da\$alabar a la alabar la alabarda\$alabar a labar a la alabarda\$a labarda\$alabar a la a r a la alabarda\$alaba rda\$alabar a la alaba

- **7** alabar a la alabarda\$
- T^{bwt} araadl ll\$ bbaar aaaa

c = L[i] LF(i) = C[c]+rank_c(L,i)

(日) (四) (三) (三) (三)

크

- The search on the suffix array can be simulated via the so-called backward search.
- It uses (the general version of) rank queries on T^{bwt}.
- The suffix array interval of pattern P is identified in at most 2|P| applications of rank.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ



▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のQで

- This index is called the FM-index.
- In its modern version, it is basically a wavelet tree on T^{bwt}.
- It was proved that a wavelet tree,
 - Built on *T^{bwt}*.
 - Using zero-order compressed bitmaps.
 - So that such compression is "local" (not rare).
- reaches high-order compression of T.
- This is the compression achieved by ppm, bzip2, etc.
- It has to do with the predictability of the next symbol given k previous symbols.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

 This simplifies the implementation of space-efficient FM-indexes.



◆□▶ ◆□▶ ◆目▶ ◆目▶ 目 うへで

Positional Inverted Indexes

- Store the word offsets of each distinct word in a text.
- Used to display snippets and to solve phrase and proximity searches.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

- Can be compressed to $nH_0(T)$ bits (word-based).
- Plus the compressed text, makes $2nH_0(T)$.

Positional Inverted Indexes

- A wavelet tree on the word identifiers uses $\approx nH_0(T)$ bits.
- It represents T: T[i] = access(i).
- It represents the inverted index: the *i*th entry of the list of word w is select_w(i).
- With rank queries, intersections (for phrases) can be done more efficiently.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

Positional Inverted Indexes



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

- ► Directed graph G(V, E) with n = |V| nodes and e = |E| edges.
- An adjacency list takes $n \log e + e \log n$ bits.
- Gives the neighbors of any node $v \in V$ efficiently.
- Reverse neighbors (i.e., who points to v)?
 - Represent the transposed G (and double the space).

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - つへで

- Is v connected to u?
 - Binary search on the list of *v*.

Concatenate all the adjacency lists

$$L = L(v_1) : L(v_2) : \ldots : L(v_n)$$

Add sparse bitmap

$$B = 10^{|L(v_1)|} 10^{|L(v_2)|} \dots 10^{|L(v_n)|}$$

- Space: nlog(e/n) + e log n + O(n) + o(e log n) (without compression).
- Neighbors of v_i, L(v_i):

$$p$$
 = select₁(B, i) - i
 $L(v_i)[j]$ = access(L, p + j)

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

• Reverse neighbors of v_i , $R(v_i)$:

$$p = select_i(L, j)$$

 $R(v_i)[j] = select_0(B, p) - p$

• Is there a link from v_i to v_j ?

$$p_1 = select_1(B, i) - i$$

$$p_2 = select_1(B, i+1) - (i+1)$$

$$v_i \rightarrow v_j \Leftrightarrow rank_j(L, p_2) - rank_j(L, p_1) > 0$$

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 臣 の�?





◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

1: 2,4,5 2: 3,5 3: 2,4 4: 5: 3

Part II: Applications

Three Views

Applications as Sequences

Applications as Reorderings

(日) (图) (문) (문) (문)

Applications as Grids

Permutations

- Consider a permutation with few increasing runs.
- A wavelet-tree-like structure describes a mergesort process.
- This can be arranged in a Hu-Tucker tree shape.
- The total number of bits is less than n(H+2).
- Here $H = \sum (n_i/n) \lg(n/n_i) \le \lg \rho$,
- where n_i is the length of run *i* and ρ the number of runs.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

Permutations



▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のQで
Permutations

- Any $\pi(i)$ can be computed with an upward tracking from *i*.
- Any $\pi^{-1}(i)$ can be computed with a downward tracking from *i*.
- Both can take O(H + 1) on average and O(lg ρ) in the worst case.
- But... do these permutations with few runs arise in practice?
 - E.g. Function Ψ has σ runs of entropy $\leq H_0(T)$.
 - This encoding enables a bidirectional compressed suffix array.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Permutations

- Function Ψ(i) = A⁻¹[A[i] + 1] is used for compressed suffix arrays (CSAs).
- It is the inverse of $LF(i) = A^{-1}[A[i] 1]$.
- A suffix array search is simulated in O(m log n) computations of Ψ.
- It has σ increasing runs, with entropy $H \leq H_0(T)$.
- A representation as a permutation enables, for example, bidirectional CSAs.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• The space is at most $n(H_0(T) + 2)$ bits.

Permutations



▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 ろくぐ

- Basic algorithmic problems on numeric sequences.
- Range quantile query:
 - Preprocess array A[1, n] on [1, U] so that, later,
 - given [l, r] and k,
 - retrieve the kth smallest element in A[I, r]
- Best solution (after much effort): O(n lg n) bits and O(lg n/ lg lg n) time.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

• We can assume $U \leq n$.

Range quantile query on wavelet trees:

- Start at the root v with interval [l, r].
- If $z = rank_0(B, I, r) \ge k$, the answer is on the left child:
- ▶ \rightarrow continue on left child, mapping [*I*, *r*].
- Otherwise, the answer is on the right child:
- ► → continue on the right with k = k z, mapping [l, r].

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

• $n \log U + o(n)$ bits, replacing the array.

O(lg U) time.

- Other problems:
 - Given [l, r] and v, give the smallest number $\geq v$ in [l, r]
 - Given [l, r] and v, give the leftmost number $\geq v$ in [l, r]
- They have clear geometric counterparts (e.g., list range points in order).
- Geometric problems like dominance and visibility reduce to these.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ



◆□▶ ◆□▶ ◆三▶ ◆三▶ ○□ のへの

For each term they store the documents where it appears and its frequency.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

- Useful to have them sorted by document identifier:
 - For conjunctive queries (intersections).
- Useful to have them sorted by decreasing frequency:
 - For ranked queries (bag of words).

- A wavelet tree concatenating the lists in frequency order...
- ... can support both orderings simultaneously:
 - By decreasing frequency: access to the sequence.
 - By increasing doc id: quantile queries.





▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のへで

- It can also simulate inverted lists of ranges of terms.
- Useful for on-the-fly stemming, prefix searches, thesaurus expansions.



◆□▶ ◆□▶ ◆三▶ ◆三▶ ○□ のへで





▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のへで

- It can intersect lists with a fast native algorithm.
- Backtracking as long as no list is empty.
- Easy to generalize to thresholded queries.





◆□▶ ◆□▶ ◆ ≧▶ ◆ ≧ ◆ ● ● ● ● ●



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

- Document retrieval on general string collections.
- Given a pattern *P* and a collection of texts T_1, \ldots, T_D ,
 - In which documents does P appear?
 - With which frequencies?
 - Give me only the k highest-frequency documents.
- One of the main tools: the document array.



(日) (문) (문) (문)

- Document listing (with frequencies) reduces to backtracking on the wavelet tree.
- $O(d \log(D/d))$ time to report the *d* documents.



▲□▶ ▲□▶ ▲注▶ ▲注▶ 注 のへで

- ► Top-*k* document retrieval via range quantile queries.
- Find consecutively more refined quantiles...
- ... until the kth document seen occurs more times than the inter-quantile distance explored up to now.



▲□▶ ▲□▶ ▲□▶ ▲□▶ = のへで

A B C D B D B D A A B B B A C D D D

- Top-k document retrieval via prioritized wavelet tree traversal.
- Set up a priority queue sorted by interval lengths...

D A B C D B D A C A B C D B D B D A C A C B D B D B A C D A C A C Λ 1 0 0 1 0 1 0 0 1 0 0 1 1 0 0 1 0 0 BBAAA CDDCCDDDCCDDCC A B B A A B B B A A B 0 0 1 1 1 0 1 0 1 1 0 0 0 1 0 0

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三回 ● ○○

 Extract each new interval and insert its two children in the queue.



The leaves come out in the proper order, stop after getting k.



◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ● ●

Part II: Applications

Three Views

Applications as Sequences

Applications as Reorderings

Applications as Grids

- Apart from obvious computational geometry applications...
- ... discrete grids model many subproblems in other areas.
- A typical one: pair prefixes and suffixes in grammar-based and Lempel-Ziv-based compressed indexes.

a.l.ab.ar. .a .la. a.lab.ard.a\$



More complex geometric problems: points have values in [1, V].

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

- Find, in rectangular ranges:
 - Sum of values.
 - Average of values.
 - Variance of values.
 - Minima or maxima of values.
 - Quantiles of values.
 - Majorities.

- Some are solved by adding one-dimensional structures to the wavelet tree nodes.
- For example, consider finding the minimum value in a range, 2D-RMQ.
- We can assume *n* points on an $n \times n$ grid, and $V \le n$.
- There exists a (1D) RMQ structure that, on array A[1, n] of numbers,

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

- Requires 2n + o(n) bits of space.
- Answers queries in constant time.
- Does not access A.

- Use a wavelet tree for the points.
- Store the values aligned to the root bitmap, in *V*[1, *n*].
- Subsequences $V_{v}[1, n_{v}]$ of V correspond to $S_{v}[1, n_{v}]$.
- ► We store a one-dimensional RMQ structure per wavelet tree node, for V_v[1, n_v].

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

Overall space is O(n lg n)

- A two-dimensional query range translates into O(lg n) intervals inside wavelet tree nodes.
- The minimum in the two-dimensional range is the minimum of the minima of those intervals of V_{ν} sequences.
- Use the one-dimensional RMQ structure on each involved interval.
- Project each one towards the root to find the values, and choose the minimum.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

• Overall time $O(\lg^{1+\epsilon} n)$.





▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のQで

- A more complicated problem: two-dimensional range quantiles.
- We can build a wavelet tree of grids.
- These have only points, not values.
- Each wavelet tree node includes the points whose values are in a range.
- Two bitmaps tell which x and y coordinates belong to the left/right child.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

Overall space is O(n lg² n)

- The algorithm is similar to that for one dimension.
- We count how many points are in the range mapped to the left child.
- ▶ If there are *k* or more, we go to the left child.
- Else, we go to the right child, subtracting the count from *k*.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• $O(\lg n)$ counting queries: $O(\lg^2 n / \lg \lg n)$ time.



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

- Top-k queries on two dimensions.
- First find the O(lg n) wavelet tree intervals.
- Store them in a priority queue, sorted by minimum value.
- Take the minimum, report it, and insert the two subintervals it splits.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

- Until having extracted the k minima.
- $O((k + \lg n) \lg^{\epsilon} n)$ time.

Binary Relations

- A generalization of graphs.
- We have *t* pairs relating *n* objects with σ labels.
- We can ask who is related with an object, or with a label, or if a pair is related.
- We can also ask, e.g., how many connections are there between a range of objects and a range of labels.
- Or list objects/labels related to a range of labels/objects.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

Binary Relations

1 2 3 4 5 6 7 8 910

- **a** 0000101000
- **b** 1000101000
- **c** 000100001
- **d** 100000000
- **e** 000010001
- **f** 0000110000

(日) (四) (문) (문) (문)

Binary Relations

- On Web graphs, this gives queries over whole domains.
- On inverted indexes (labels are terms and objects are documents):
 - ► Ranges of labels: stemming, prefixes, thesaurus expansion.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

- Ranges of objects: temporal, versioned, hierarchical collections.
- Vertical stripes: vocabulary of (ranges of) documents.
- Style and plagiarism analysis?
- Horizontal stripes: document frequencies.
- Simulate inverted list of a range of words.
- List vocabulary of a range of documents.
Binary Relations



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Documents

Dynamism

- Can support insertions and deletions of symbols.
- All operations take time $O((\lg n / \lg \lg n)(1 + \lg \sigma / \lg \lg n))$.
- Probably near-optimal since dynamic range counting is $\Omega((\lg n / \lg \lg n)^2).$
- However [N. and Nekrich, arXiv 2012]
 - If all you want is access, rank, select, and indels...

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

- You can get the optimal $O(\lg n / \lg \lg n)$ time.
- Worst-case for queries, amortized for indels.
- Zero-order compressed space, as usual.

Conclusions

- Simple and versatile data structure.
- Supports different views of the data within the same space.

▲ロト ▲団ト ▲ヨト ▲ヨト 三目 つんぐ

- Enables succinct and compressed representations.
- Good (near-log) query times for many problems.
- Lots of applications.
- Wavelet trees for all!



| ◆ □ ▶ ◆ ■ ▶ ◆ ■ ◆ ● ● ● ● ● ●



▲□▶ ▲□▶ ▲三▶ ▲三▶ 三回 めんの