

# Indexing sets of strings: Wheeler languages

## Part I

---

Alberto Policriti & Giovanna D'Agostino



Dpt. of Mathematics, Computer Science, and Physics, University of Udine.

# Why should we care about big data(sets)?

## Earth Biogenome Project



Sequencing 1.5 (!) million species

# Why should we care about big data(sets)?

## Earth Biogenome Project



Sequencing 1.5 (!) million species

### Side (?) observation

coverage increases  $\Rightarrow$  *repetitiveness* increases

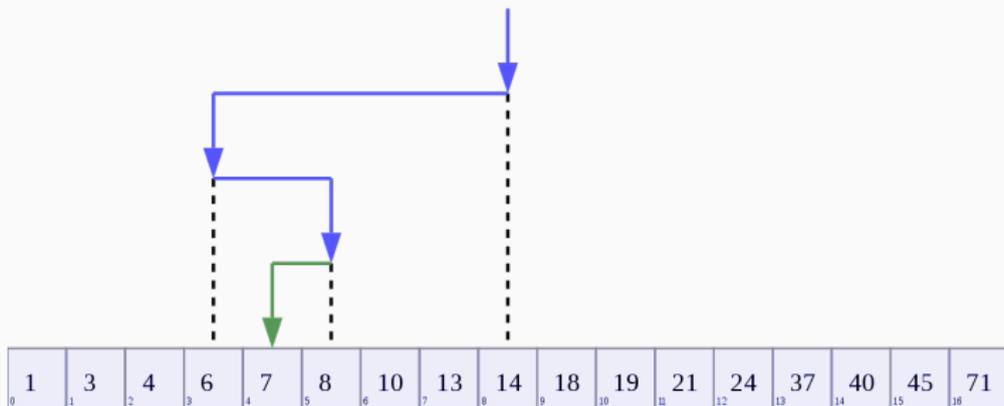
## Example

Given a file  $S$  containing 1T ( $\approx 2^{40}$  numbers) numbers and a number  $n$  answer to the question:

is  $n$  in  $S$ ?

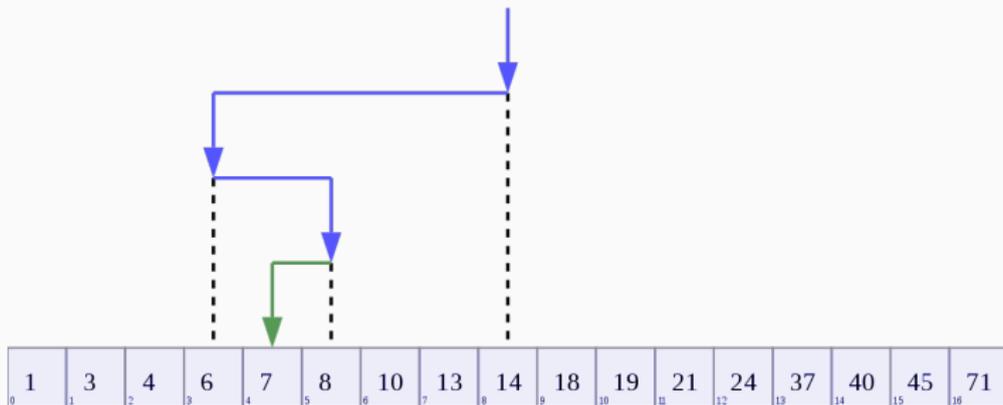
# Binary Search

Search for 7:



# Binary Search

Search for 7:

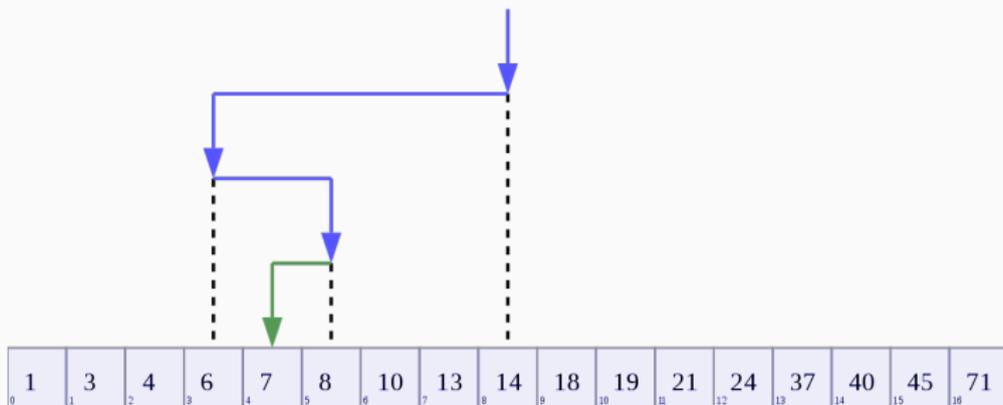


How many questions/steps/instructions?

≈ 40 ... why?

# Binary Search

Search for 7:



$$2^{40} = \underbrace{2 \cdot 2 \cdot 2 \cdot \dots \cdot 2}_{40 \text{ times}}$$

Elements can be *turned into* (encoded by) **binary codes**:  
[indexes](#)

Elements can be *turned into* (encoded by) **binary codes**:

indexes

What about space?

Elements can be *turned into* (encoded by) **binary codes**:

indexes

What about space?

Are there better indexes? (Better codes/encodings?)

## Learned Indexes

### The Case for Learned Index Structures

Tim Kraska\*  
MIT  
Cambridge, MA  
kraska@mit.edu

Alex Beutel  
Google, Inc.  
Mountain View, CA  
alexbeutel@google.com

Ed H. Chi  
Google, Inc.  
Mountain View, CA  
edchi@google.com

Jeffrey Dean  
Google, Inc.  
Mountain View, CA  
jeff@google.com

Neoklis Polyzotis  
Google, Inc.  
Mountain View, CA  
npolyzotis@google.com

30 Apr 2018

Abstract

### The PGM-index: a multicriteria, compressed and learned approach to data indexing

Paolo Ferragina  
University of Pisa, Italy  
paolo.ferragina@unipi.it

Giorgio Vinciguerra  
University of Pisa, Italy  
giorgio.vinciguerra@phd.unipi.it

14 Oct 2019

#### ABSTRACT

The recent introduction of learned indexes has shaken the foundations of the decades-old field of indexing data structures. Combining, or even replacing, classic design elements such as B-tree nodes with machine learning models has proven to give outstanding improvements in the space footprint and time efficiency of data systems. However, these novel approaches are based on heuristics, thus they

formidable results, we still miss proper algorithms and data structures that are flexible enough to work under computational constraints that vary across users, devices and applications, and possibly evolve over time.

In this paper, we restrict our attention to the case of *indexing data structures* for internal or external memory which solve the so-called *fully indexable dictionary* problem. This problem asks to store a multiset  $S$  of real keys in order to

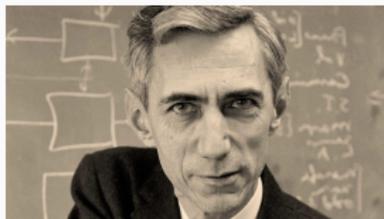
# Entropy of $S$

How difficult is to find an element in  $S$ ?

$$\mathcal{H}(S) = \log |S|$$

How much *information* is carried by each code (of an element of  $S$ )?

BITS and SHANNONS



Probability and entropy of a **random variable distribution**

Probability and entropy of a **random variable distribution**

## **Definition**

For  $X$  r.v. taking values  $x_1, \dots, x_n$

Probability and entropy of a random variable distribution

### Definition

For  $X$  r.v. taking values  $x_1, \dots, x_n$

$$\begin{aligned}\mathcal{H}(X) &= \sum_{i=1}^n p(x_i) \log(1/p(x_i)) \\ &= E[\log(1/p(x_i))] \dots \text{BITS (SHANNONS)}\end{aligned}$$

Probability and entropy of a **random variable distribution**

### **Definition**

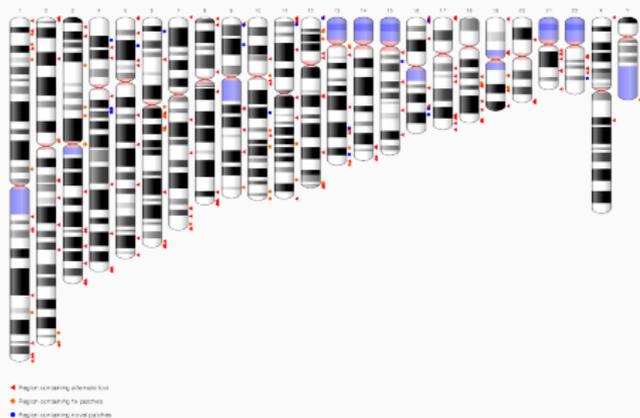
For  $X$  r.v. taking values  $x_1, \dots, x_n$

$$\begin{aligned}\mathcal{H}(X) &= \sum_{i=1}^n p(x_i) \log(1/p(x_i)) \\ &= E[\log(1/p(x_i))] \dots \text{BITS (SHANNONS)}\end{aligned}$$

Then, we cannot encode a sequence  $S$  in less than

$$H(S) = - \sum_{c \in \Sigma} P(c) \cdot \log_2(P(c)) \text{ bits/symbol .}$$

# Entropy of a Genome



```
TTTCCCGTAGGTGAACCTGCGGAAGGATCATTTGTCGAATCCTGCGAAGCAGACGACCTGAGAACATGTTA
AACTCACGGGGCGTTCGGGACGGGGGAGACCTCGGACCGACGCCCCGAGCAGCAGCTCGACCGACGGTTGC
TCGTCGTGCATAAAACTAACCCGGGCGCGGATAAGCGCCAAGGAAAATATGAAAGGATGCCTGCCTC
TCGATGCCCCGTACCGGTTGTGCACGGGAGGATCACGGGCGTCTGGAGAAAACAAAACGACTCTCGGCAA
CGGATATCTCGGCTCTCGCATCGATGAAGAAGTAGCGAACTGCGATACTTGGTGTGAATGCAGAATCC
CGTGAACCATCGAGTCTTTGAACGCAAGTTGCGCCCGAAGCCATTAGGCCGAGGGCACGCTGCCTGGGC
GTCACGCATCGCGTCGCCCCAACACATGTCTCCGTATCGCACGGGTACCCGGGAGGGGGCGGAAACTG
GCTTCCCGTGCCTTTCGTGCGGCTGGCCTAAACACGAGTCCCTAGCGATGGACGTTGTGACAAGTGGTGGT
TGATTTCCCTCAACTAGAGCGCTGTCACGCGATCCTCCATCGATTGAGGAGACTCCCACGACCTGATGCA
```

How do we structure (e.g.) genomes for searching?

- not one string but many strings
- string of characters (not numbers)
- we search for substrings (not single characters)
- ...

How do we structure (e.g.) genomes for searching?

- not one string but many strings
- string of characters (not numbers)
- we search for substrings (not single characters)
- ...

How do we exploit repetitiveness?

We **ordered** data in order to apply binary search.

What about DNA?

## Ordering a string (of DNA)?

### Suffixes

1	AGGTTGCCAGTGT	1	AGGTTGCCAGTGT
2	GGTTGCCAGTGT	9	AGTGT
3	GTTGCCAGTGT	8	CAGTGT
4	TTGCCAGTGT	7	CCAGTGT
5	TGCCAGTGT	6	GCCAGTGT
6	GCCAGTGT	2	GGTTGCCAGTGT
7	CCAGTGT	12	GT
8	CAGTGT	10	GTGT
9	AGTGT	3	GTTGCCAGTGT
10	GTGT	13	T
11	TGT	5	TGCCAGTGT
12	GT	11	TGT
13	T	4	TTGCCAGTGT

## Ordered suffixes (can be compressed): Burrows-Wheeler

AGGTTGCCAGTGT

AGTGT

CAGTGT

CCAGTGT

GCCAGTGT

GGTTGCCAGTGT

GT

GTGT

GTTGCCAGTGT

T

TGCCAGTGT

TGT

TTGCCAGTGT

## Ordered suffixes (can be compressed): Burrows-Wheeler

\$  
AGGTTGCCAGTGT\$  
AGTGT\$  
CAGTGT\$  
CCAGTGT\$  
GCCAGTGT\$  
GGTTGCCAGTGT\$  
GT\$  
GTGT\$  
GTTGCCAGTGT\$  
T\$  
TGCCAGTGT\$  
TGT\$  
TTGCCAGTGT\$

## Ordered suffixes (can be compressed): Burrows-Wheeler

\$AGGTTGCCAGTGT  
AGGTTGCCAGTGT\$  
AGTGT\$AGGTTGCC  
CAGTGT\$AGGTTGC  
CCAGTGT\$AGGTTG  
GCCAGTGT\$AGGTT  
GGTTGCCAGTGT\$A  
GT\$AGGTTGCCAGT  
GTGT\$AGGTTGCCA  
GTTGCCAGTGT\$AG  
T\$AGGTTGCCAGTG  
TGCCAGTGT\$AGGT  
TGT\$AGGTTGCCAG  
TTGCCAGTGT\$AGG

## Ordered suffixes (can be compressed): Burrows-Wheeler

\$AGGTTGCCAGTGT  
AGGTTGCCAGTGT\$  
AGTGT\$AGGTTGCC  
CAGTGT\$AGGTTGC  
CCAGTGT\$AGGTTG  
GCCAGTGT\$AGGTT  
GGTTGCCAGTGT\$A  
GT\$AGGTTGCCAGT  
GTGT\$AGGTTGCCA  
GTTGCCAGTGT\$AG  
T\$AGGTTGCCAGTG  
TGCCAGTGT\$AGGT  
TGT\$AGGTTGCCAG  
TTGCCAGTGT\$AGG

## Definition (Burrows-Wheeler Transform)

$\text{BWT}(\text{AGGTTGCCAGTGT\$}) = \text{T\$CCGTATAGGTGG}$

---

## Proposition

*The BWT of a text can be (easily) **inverted** and **compressed**.*

How much space are we using?

**booster**

The BWT is a *booster* for compression

## An Analysis of the Burrows-Wheeler Transform

Giovanni Manzini

Dipartimento di Informatica, Università del Piemonte Orientale, Italy.

JACM 2001

run-length encoding

How much space are we using?

**booster**

The BWT is a *booster* for compression

## An Analysis of the Burrows-Wheeler Transform

Giovanni Manzini

Dipartimento di Informatica, Università del Piemonte Orientale, Italy.

JACM 2001

run-length encoding

How much space are we using?

**booster**

The BWT is a *booster* for compression

## An Analysis of the Burrows-Wheeler Transform

Giovanni Manzini

Dipartimento di Informatica, Università del Piemonte Orientale, Italy.

JACM 2001

run-length encoding

Long repetitions?

## Searching the $BWT(T)$ ...

... is like searching the suffix array (via the L-F property)

## Searching the $BWT(T)$ ...

... is like searching the suffix array (via the L-F property)

The  $BWT$  orders suffixes.

# What about graphs?

Theoretical Computer Science 698 (2017) 67–78



ELSEVIER

Contents lists available at [ScienceDirect](#)

Theoretical Computer Science

[www.elsevier.com/locate/tcs](http://www.elsevier.com/locate/tcs)



## Wheeler graphs: A framework for BWT-based data structures <sup>☆</sup>



Travis Gagie <sup>a</sup>, Giovanni Manzini <sup>b,c,\*</sup>, Jouni Sirén <sup>d</sup>

<sup>a</sup> *Diego Portales University and CEBIB, Santiago, Chile*

<sup>b</sup> *University of Eastern Piedmont, Alessandria, Italy*

<sup>c</sup> *IIT-CNR, Pisa, Italy*

<sup>d</sup> *Wellcome Trust Sanger Institute, Cambridge, UK*

## A B S T R A C T

---

The famous Burrows–Wheeler Transform (BWT) was originally defined for a single string but variations have been developed for sets of strings, labeled trees, de Bruijn graphs, etc. In this paper we propose a framework that includes many of these variations and that we hope will simplify the search for more.

We first define *Wheeler graphs* and show they have a property we call *path coherence*. We show that if the state diagram of a finite-state automaton is a Wheeler graph then, by its path coherence, we can order the nodes such that, for any string, the nodes reachable from the initial state or states by processing that string are consecutive. This means that even if the automaton is non-deterministic, we can still store it compactly and process strings with it quickly.

We then rederive several variations of the BWT by designing straightforward finite-state automata for the relevant problems and showing that their state diagrams are Wheeler graphs.

© 2017 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

# What about graphs?

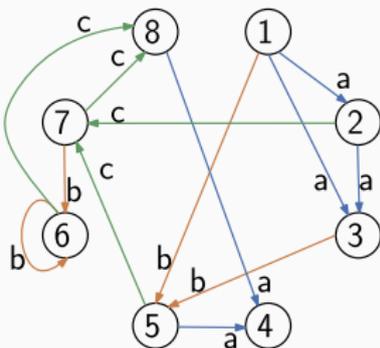
**Definition 1.**  $G$  is a *Wheeler graph* if there is an ordering of the nodes such that nodes with in-degree 0 precede those with positive in-degree and, for any pair of edges  $e = (u, v)$  and  $e' = (u', v')$  labeled  $a$  and  $a'$  respectively, the following monotonicity properties hold:

$$\begin{aligned} a < a' &\implies v < v', \\ (a = a') \wedge (u < u') &\implies v \leq v'. \end{aligned}$$

# What about graphs?

**Definition 1.**  $G$  is a *Wheeler graph* if there is an ordering of the nodes such that nodes with in-degree 0 precede those with positive in-degree and, for any pair of edges  $e = (u, v)$  and  $e' = (u', v')$  labeled  $a$  and  $a'$  respectively, the following monotonicity properties hold:

$$a < a' \implies v < v',$$
$$(a = a') \wedge (u < u') \implies v \leq v'.$$



**Definition 2.**  $G$  is *path coherent* if there is a total order of the nodes such that for any consecutive range  $[i, j]$  of nodes and string  $\alpha$ , the nodes reachable from those in  $[i, j]$  in  $|\alpha|$  steps by following edges whose labels for  $\alpha$  when concatenated, themselves form a consecutive range.

**Lemma 3.** If  $G$  is a *Wheeler graph* by an ordering  $\pi$  on its nodes then it is *path coherent* by  $\pi$ .

- Ordered *sets of strings*

**Definition 2.**  $G$  is *path coherent* if there is a total order of the nodes such that for any consecutive range  $[i, j]$  of nodes and string  $\alpha$ , the nodes reachable from those in  $[i, j]$  in  $|\alpha|$  steps by following edges whose labels for  $\alpha$  when concatenated, themselves form a consecutive range.

**Lemma 3.** If  $G$  is a *Wheeler graph* by an ordering  $\pi$  on its nodes then it is *path coherent* by  $\pi$ .

- Ordered *sets of strings*
- Align to the **right**

**Definition 2.**  $G$  is *path coherent* if there is a total order of the nodes such that for any consecutive range  $[i, j]$  of nodes and string  $\alpha$ , the nodes reachable from those in  $[i, j]$  in  $|\alpha|$  steps by following edges whose labels for  $\alpha$  when concatenated, themselves form a consecutive range.

**Lemma 3.** If  $G$  is a *Wheeler graph* by an ordering  $\pi$  on its nodes then it is *path coherent* by  $\pi$ .

- Ordered *sets of strings*
- Align to the **right**
- *Intervals* in the collection of states
- *Intervals* in the collection of strings

**Definition 2.**  $G$  is *path coherent* if there is a total order of the nodes such that for any consecutive range  $[i, j]$  of nodes and string  $\alpha$ , the nodes reachable from those in  $[i, j]$  in  $|\alpha|$  steps by following edges whose labels for  $\alpha$  when concatenated, themselves form a consecutive range.

**Lemma 3.** If  $G$  is a Wheeler graph by an ordering  $\pi$  on its nodes then it is path coherent by  $\pi$ .

- Ordered sets of strings
- Align to the right
- *Intervals* in the collection of states
- *Intervals* in the collection of strings
- (Infinity)

- minimization  $\Rightarrow$  many function coarsest partition problem (non-deterministic case: **handle with care!**)

- minimization  $\Rightarrow$  many function coarsest partition problem (non-deterministic case: **handle with care!**)
- Myhill-Nerode Theorem  $\Rightarrow$  States of the minimum DFA accepting  $\mathcal{L}$  are sets strings

- minimization  $\Rightarrow$  many function coarsest partition problem (non-deterministic case: **handle with care!**)
- Myhill-Nerode Theorem  $\Rightarrow$  States of the minimum DFA accepting  $\mathcal{L}$  are sets strings
- non-determinism can be eliminated  $\Rightarrow$  power-construction (exponential blow-up)

### Regular Languages meet Prefix Sorting \*

Jarno Alanko<sup>†</sup>

Giovanna D'Agostino<sup>‡</sup>

Alberto Policriti<sup>§</sup>

Nicola Prezza<sup>¶</sup>

1. We show that Wheeler languages are the natural version of regular languages endowed with the co-lexicographic ordering: when sorted, the prefixes of strings belonging to a Wheeler language are partitioned into a *finite* number of *intervals*, each formed by elements from a single Myhill-Nerode equivalence class. In regular languages, those *intervals* are replaced with general *sets*.

2. We show that every Wheeler NFA (WNFA) with  $n$  states admits an equivalent Wheeler DFA (WDFA) with at most  $2n - 1 - |\Sigma|$  states ( $\Sigma$  being the alphabet) that can be computed in  $O(n^3)$  time. This is in sharp contrast with general NFAs (where the blow-up could be exponential).

3. Let  $d$ -NFA denote the class of NFAs with at most  $d$  equally-labeled edges leaving any state. We show that the problem of recognizing and sorting Wheeler  $d$ -NFAs is in P for  $d \leq 2$ . A recent result from Gibney and Thankachan [12] shows that the problem is NP-complete for  $d \geq 5$ . Our result almost completes the picture, the remaining open cases being  $d = 3$  and  $d = 4$ .

4. We provide an online incremental algorithm that, when fed with an acyclic Wheeler DFA's nodes in any topological order, can dynamically compute the co-lexicographic rank of each new incoming node among those already processed with just logarithmic delay.
5. We improve the running time of (4) to linear in the offline setting for arbitrary WDFAs.

6. Given a Wheeler DFA  $\mathcal{A}$  of size  $n$ , we show how to compute, in  $O(n \log n)$  time, the smallest Wheeler DFA recognizing the same language as  $\mathcal{A}$ . If  $\mathcal{A}$  is acyclic, running time drops to  $O(n)$ .
7. Given *any acyclic DFA*  $\mathcal{A}$  of size  $n$ , we show how to compute, in  $O(n + m \log m)$  time, the smallest Wheeler DFA  $\mathcal{A}'$ , of size  $m$ , recognizing the same language as  $\mathcal{A}$ .

### Decidability

Can we decide whether  $\mathcal{L}$  is Wheeler?