

Space Efficient Data Structures and FM index

Venkatesh Raman

The Institute of Mathematical Sciences, Chennai

NISER Bhubaneswar, February 9, 2019

Overview

Introduction

Data Structures

Goals

Bit Vectors

Strings from a larger alphabet

Sparse Bit Vectors

Trees

Burrows-Wheeler Transform and Indexing

Libraries

Conclusions

Overview

- Plan of the talk

Overview

- Plan of the talk
 - **Why** Space efficient?

Overview

- Plan of the talk
 - **Why** Space efficient?
 - **What** we mean by efficient?

Overview

- Plan of the talk
 - **Why** Space efficient?
 - **What** we mean by efficient? (information theory lower bound)

Overview

- Plan of the talk
 - **Why** Space efficient?
 - **What** we mean by efficient? (information theory lower bound)
 - **How**

Overview

- Plan of the talk
 - **Why** Space efficient?
 - **What** we mean by efficient? (information theory lower bound)
 - **How** Some examples

Overview

- Plan of the talk
 - **Why** Space efficient?
 - **What** we mean by efficient? (information theory lower bound)
 - **How** Some examples
(a binary (or d -ary) vector, subset of a finite universe)

Overview

- Plan of the talk
 - **Why** Space efficient?
 - **What** we mean by efficient? (information theory lower bound)
 - **How** Some examples
(a binary (or d -ary) vector, subset of a finite universe)
 - **Success Story** BWT and FM index

Overview

- Plan of the talk
 - **Why** Space efficient?
 - **What** we mean by efficient? (information theory lower bound)
 - **How** Some examples
(a binary (or d -ary) vector, subset of a finite universe)
 - **Success Story** BWT and FM index
- A recent book
 - *Compact Data Structures: A Practical Approach*, Gonzalo Navarro, Cambridge UP, 2016.

Data Structures

Data Structures

- Pre-process input data so as to answer (long) series of *retrieval* or *update* operations.

Data Structures

- Pre-process input data so as to answer (long) series of *retrieval* or *update* operations.
- Want to minimize:
 1. Query/Update time.
 2. Space usage of data structure.
 3. Time of pre-processing.
 4. Space for pre-processing.

Data Structures

- Pre-process input data so as to answer (long) series of *retrieval* or *update* operations.
- Want to minimize:
 1. Query/Update time.
 2. Space usage of data structure.
 3. Time of pre-processing.
 4. Space for pre-processing.
- In this talk we will worry only about the first two, and

Data Structures

- Pre-process input data so as to answer (long) series of *retrieval* or *update* operations.
- Want to minimize:
 1. Query/Update time.
 2. Space usage of data structure.
 3. Time of pre-processing.
 4. Space for pre-processing.
- In this talk we will worry only about the first two, and our data structures are static.

Space usage of Data Structures

Answering queries on data requires an *index* **in addition to** the data. Index may be much larger than the data. E.g.:

Space usage of Data Structures

Answering queries on data requires an *index* **in addition to** the data. Index may be much larger than the data. E.g.:

- *Range Trees*: data structure for answering 2-D orthogonal range queries on n points.
 - Good worst-case performance but $\Theta(n \log n)$ words of space.

Space usage of Data Structures

Answering queries on data requires an *index* **in addition to** the data. Index may be much larger than the data. E.g.:

- *Range Trees*: data structure for answering 2-D orthogonal range queries on n points.
 - Good worst-case performance but $\Theta(n \log n)$ words of space.
- *Suffix Trees*: data structure for indexing a sequence T of n symbols from an alphabet of size σ .



Space usage of Data Structures

Answering queries on data requires an *index* **in addition to** the data. Index may be much larger than the data. E.g.:

- *Range Trees*: data structure for answering 2-D orthogonal range queries on n points.
 - Good worst-case performance but $\Theta(n \log n)$ words of space.
- *Suffix Trees*: data structure for indexing a sequence T of n symbols from an alphabet of size σ .
 - Supports very complex queries on string patterns quickly but uses $\Theta(n)$ words of space.
 - One word must have at least $\log_2 n$ bits.
 - $\Theta(n)$ words is $\Omega(n \log n)$ bits – raw sequence T is $n \log_2 \sigma$ bits.
 - A good implementation takes 10x to 30x space more than T .

Succinct/Compressed Data Structures

$$\text{Space usage} = \text{“space for data”} + \underbrace{\text{“space for index”}}_{\text{redundancy}}.$$

- Redundancy (working space used by data structure to answer queries) should be small.

Succinct/Compressed Data Structures

$$\text{Space usage} = \text{“space for data”} + \underbrace{\text{“space for index”}}_{\text{redundancy}}.$$

- Redundancy (working space used by data structure to answer queries) should be small. Ideally $o(\text{inputsize})$.
- What should be the space for the data?

Why care about space?

Why care about space?

- While the cost of memory continues to go down, the growth of data is increasing at a much higher rate. (E.g. Search Engines, Genome data)

Why care about space?

- While the cost of memory continues to go down, the growth of data is increasing at a much higher rate. (E.g. Search Engines, Genome data)
- Space is important if we want to pack a lot of data into handheld devices.

Why care about space?

- While the cost of memory continues to go down, the growth of data is increasing at a much higher rate. (E.g. Search Engines, Genome data)
- Space is important if we want to pack a lot of data into handheld devices.
- Sometimes, better space usage increases the amount of data that can be stored in main memory, thereby increasing time efficiency too.

Models of Computation

- Computational model:
 - Unit-cost RAM with word size $\Theta(\log n)$ bits.
 - Operations on $O(\log n)$ bit operands (addition, subtraction, OR, multiplication, ..) in $O(1)$ time.
 - Space counted in terms of bits.

Models of Computation

- Computational model:
 - Unit-cost RAM with word size $\Theta(\log n)$ bits.
 - Operations on $O(\log n)$ bit operands (addition, subtraction, OR, multiplication, ..) in $O(1)$ time.
 - Space counted in terms of bits.
 - There are also other models like Cell-probe model with word size $\Theta(\log n)$ bits (normally used for lower bounds).

“Space for Data”

Definition (Information-theoretic Lower Bound)

If an object x is chosen from a set S then in the worst case we need $\log_2 |S|$ bits to represent x .

“Space for Data”

Definition (Information-theoretic Lower Bound)

If an object x is chosen from a set S then in the worst case we need $\log_2 |S|$ bits to represent x .

- x is a binary string of length n .
- S is the set of all binary strings of length n .
- $\log_2 |S| = \log_2 2^n = n$ bits.

“Space for Data”

Definition (Information-theoretic Lower Bound)

If an object x is chosen from a set S then in the worst case we need $\log_2 |S|$ bits to represent x .

- x is a permutation over $\{1, \dots, n\}$.
- S is the set of all permutations over $\{1, \dots, n\}$.
- $\log_2 |S| = \log_2 n! = n \log_2 n - n \log_2 e + o(n)$ bits.

Note that the standard way to represent a permutation takes $n \lceil \lg n \rceil$ bits.

“Space for Data”

Definition (Information-theoretic Lower Bound)

If an object x is chosen from a set S then in the worst case we need $\log_2 |S|$ bits to represent x .

- x is a binary string of length n with m 1s.
- S is the set of all binary strings of length n with m 1s.
- $\log_2 |S| = \log_2 \binom{n}{m} = m \log_2(n/m) + O(m)$ bits.
 - E.g. if $m = O(n/\log n)$ then the lower bound is $O(m \log \log n) = o(n)$ bits.
 - if we just write down the positions of the 1's, that is $m \lceil \log_2 n \rceil$ bits

“Space for Data”

Definition (Information-theoretic Lower Bound)

If an object x is chosen from a set S then in the worst case we need $\log_2 |S|$ bits to represent x .

- x is a triangulated planar graph of n nodes.
- S is the set of all triangulated planar graphs with n nodes.
- $\log_2 |S| \sim 3.24n$ bits.

There are also bounds for general graphs, chordal graphs, bounded treewidth graphs.

Overview

Introduction

Data Structures

Goals

Bit Vectors

Strings from a larger alphabet

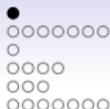
Sparse Bit Vectors

Trees

Burrows-Wheeler Transform and Indexing

Libraries

Conclusions



Succinct Data Structures

Aim is to store using space:

$$\text{Space usage} = \text{“space for data”} + \underbrace{\text{“space for index”}}_{\text{lower-order term}}.$$

and perform operations *directly* on it.

- For static DS, often get $O(1)$ time operations.
- Representation often tightly tied to set of operations.
- They work in practice!



Bit Vectors

Data: Sequence X of n bits, x_1, \dots, x_n .

ITLB: n bits; total space $n + o(n)$ bits.



Bit Vectors

Data: Sequence X of n bits, x_1, \dots, x_n .

ITLB: n bits; total space $n + o(n)$ bits.

Operations:

- $rank_1(i)$: number of 1s in x_1, \dots, x_i .
- $select_1(i)$: position of i th 1.

Also $rank_0$, $select_0$. Ideally all in $O(1)$ time.

Example: $X = 01101001$, $rank_1(4) = 2$, $select_0(4) = 7$.



Bit Vectors

Data: Sequence X of n bits, x_1, \dots, x_n .

ITLB: n bits; total space $n + o(n)$ bits.

Operations:

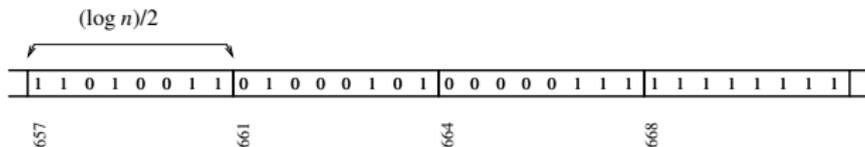
- $rank_1(i)$: number of 1s in x_1, \dots, x_i .
- $select_1(i)$: position of i th 1.

Also $rank_0$, $select_0$. Ideally all in $O(1)$ time.

Example: $X = 01101001$, $rank_1(4) = 2$, $select_0(4) = 7$.

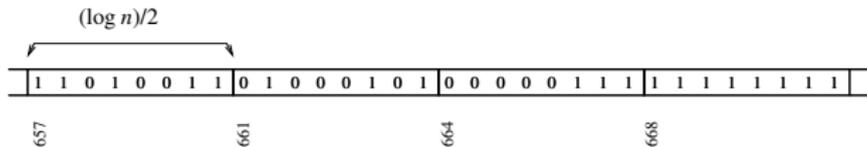
Operations introduced in [Elias, *J. ACM '75*], [Tarjan and Yao, *C. ACM '78*], [Chazelle, *SIAM J. Comput '85*], [Jacobson, *FOCS '89*].

Bit Vectors: Implementing $rank_1$



- Naive solution: store answer to all $rank_1$ queries. Space: $O(n \log n)$ bits.
- Sample: store answer only to every $(\log n)/2$ -th $rank_1$ queries. Space: $O(n)$ bits.

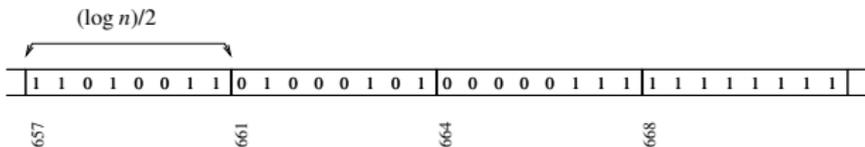
Bit Vectors: Implementing $rank_1$



- Naive solution: store answer to all $rank_1$ queries. Space: $O(n \log n)$ bits.
- Sample: store answer only to every $(\log n)/2$ -th $rank_1$ queries. Space: $O(n)$ bits.
- How to support $rank_1$ in $O(1)$ time?

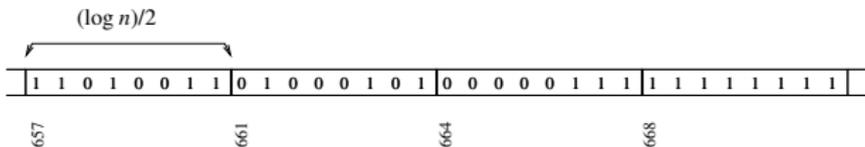
○
○○●○○○
○
○○○○
○○○
○○○○○○○

Bit-Vectors: Implementing $rank_1$



- Scanning the $(\log n)/2$ block takes $O(\log n)$ time.

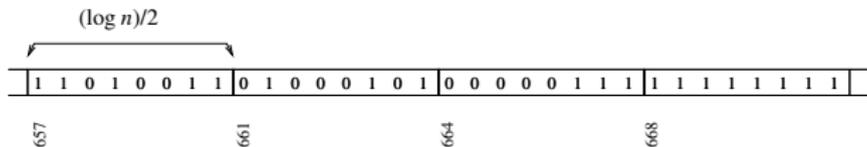
Bit-Vectors: Implementing $rank_1$



- Scanning the $(\log n)/2$ block takes $O(\log n)$ time.
- We will use what is called the “Four Russians trick”.



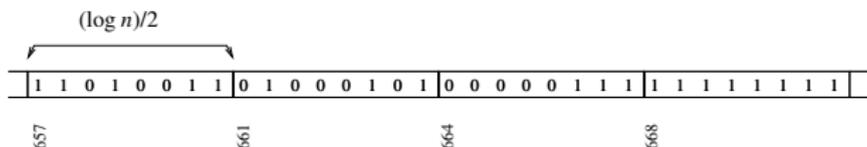
Bit-Vectors: Implementing $rank_1$



- Scanning the $(\log n)/2$ block takes $O(\log n)$ time.
- We will use what is called the “Four Russians trick”.
- Let $k = (\log n)/2$. Create a table A with $2^{k+\log_2 k} = O(\sqrt{n} \log n) = o(n)$ entries.

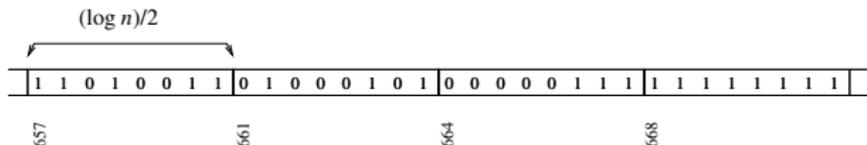


Bit-Vectors: Implementing $rank_1$



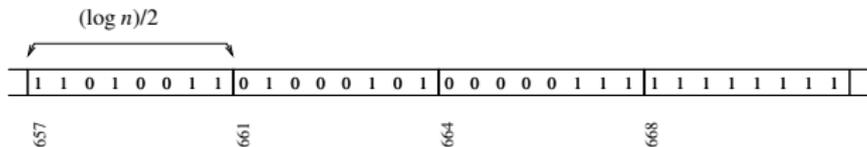
- Scanning the $(\log n)/2$ block takes $O(\log n)$ time.
- We will use what is called the “Four Russians trick”.
- Let $k = (\log n)/2$. Create a table A with $2^{k+\log_2 k} = O(\sqrt{n} \log n) = o(n)$ entries.
- $A[y_1 \dots y_{\log_2 k} x_1 \dots x_k] = \text{number of 1s in } x_1 \dots x_{y+1}$ where $y = y_1 \dots y_{\log_2 k}$. (The “four Russians” trick.)

Bit-Vectors: Implementing $rank_1$

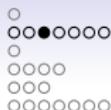


- Scanning the $(\log n)/2$ block takes $O(\log n)$ time.
- We will use what is called the “Four Russians trick”.
- Let $k = (\log n)/2$. Create a table A with $2^{k+\log_2 k} = O(\sqrt{n} \log n) = o(n)$ entries.
- $A[y_1 \dots y_{\log_2 k} x_1 \dots x_k] =$ number of 1s in $x_1 \dots x_{y+1}$ where $y = y_1 \dots y_{\log_2 k}$. (The “four Russians” trick.)
- $rank_1(x) = 657 + \underbrace{A[10111010011]}_3$.

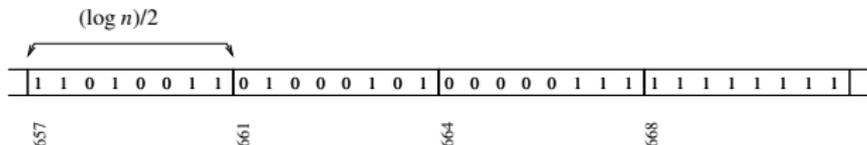
Bit-Vectors: Implementing $rank_1$



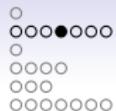
- Scanning the $(\log n)/2$ block takes $O(\log n)$ time.
- We will use what is called the “Four Russians trick”.
- Let $k = (\log n)/2$. Create a table A with $2^{k+\log_2 k} = O(\sqrt{n} \log n) = o(n)$ entries.
- $A[y_1 \dots y_{\log_2 k} x_1 \dots x_k] =$ number of 1s in $x_1 \dots x_{y+1}$ where $y = y_1 \dots y_{\log_2 k}$. (The “four Russians” trick.)
- $rank_1(x) = 657 + A[\underbrace{10111010011}_3]$.
- $O(n)$ bits, $O(1)$ time.



Bit-Vectors: Implementing $rank_1$



- Scanning the $(\log n)/2$ block takes $O(\log n)$ time.
- We will use what is called the “Four Russians trick”.
- Let $k = (\log n)/2$. Create a table A with $2^{k+\log_2 k} = O(\sqrt{n} \log n) = o(n)$ entries.
- $A[y_1 \dots y_{\log_2 k} x_1 \dots x_k] =$ number of 1s in $x_1 \dots x_{y+1}$ where $y = y_1 \dots y_{\log_2 k}$. (The “four Russians” trick.)
- $rank_1(x) = 657 + \underbrace{A[10111010011]}_3$.
- $O(n)$ bits, $O(1)$ time.
- Many theoretical SDS: decompose + sample + table lookup.



Bit-Vectors: Implementing $rank_1$

Improve redundancy by two-level approach.

Bit-Vectors: Implementing $rank_1$

Improve redundancy by two-level approach.

- Store answer for every $\log^2 n$ positions. This takes only $O(n \log n / \log^2 n = n / \log n) = o(n)$ bits.

```
○  
○○○●○○○  
○  
○○○○  
○○○  
○○○○○○○
```

Bit-Vectors: Implementing $rank_1$

Improve redundancy by two-level approach.

- Store answer for every $\log^2 n$ positions. This takes only $O(n \log n / \log^2 n = n / \log n) = o(n)$ bits.
- Then for every $(\log n) / 2$ positions, store answer **within** the block. This takes $O(n(\log \log n) / \log n) = o(n)$ bits.



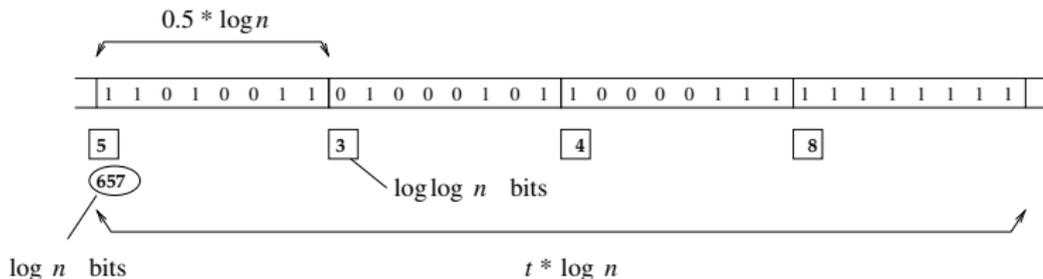
Bit-Vectors: Implementing $rank_1$

Improve redundancy by two-level approach.

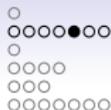
- Store answer for every $\log^2 n$ positions. This takes only $O(n \log n / \log^2 n = n / \log n) = o(n)$ bits.
- Then for every $(\log n)/2$ positions, store answer **within** the block. This takes $O(n(\log \log n) / \log n) = o(n)$ bits.
- Then store, as before, a table to find answers within $(\log n)/2$ positions.



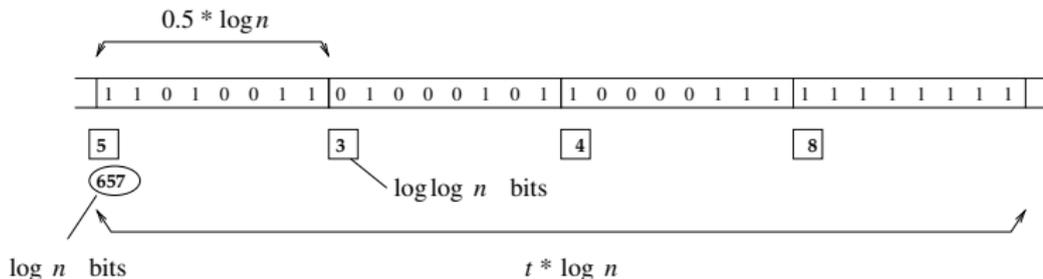
Bit-Vectors: Implementing $rank_1$ Two-level approach



$$\begin{aligned}
 \text{Space} &= n + O\left(\frac{n}{t \lg n} \cdot \lg n + \frac{n}{\lg n} \cdot \lg \lg n\right) + O(\sqrt{n} \cdot \lg n) \\
 &= n + O(n \log \log n / \log n) \text{ bits: choose } t = \Theta(\log n / \log \log n).
 \end{aligned}$$

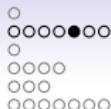


Bit-Vectors: Implementing $rank_1$ Two-level approach

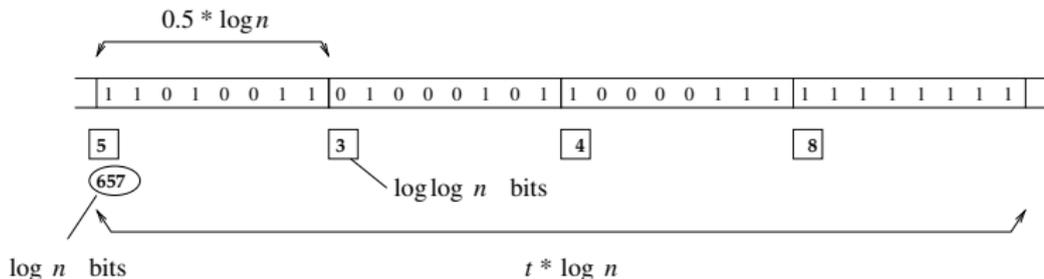


$$\begin{aligned} \text{Space} &= n + O\left(\frac{n}{t \lg n} \cdot \lg n + \frac{n}{\lg n} \cdot \lg \lg n\right) + O(\sqrt{n} \cdot \lg n) \\ &= n + O(n \log \log n / \log n) \text{ bits: choose } t = \Theta(\log n / \log \log n). \end{aligned}$$

- Redundancy $O(n \lg \lg n / \lg n)$ bits, *optimal* for $O(1)$ time operations [Golynski, *TCS'07*].



Bit-Vectors: Implementing $rank_1$ Two-level approach



$$\begin{aligned}
 \text{Space} &= n + O\left(\frac{n}{t \lg n} \cdot \lg n + \frac{n}{\lg n} \cdot \lg \lg n\right) + O(\sqrt{n} \cdot \lg n) \\
 &= n + O(n \log \log n / \log n) \text{ bits: choose } t = \Theta(\log n / \log \log n).
 \end{aligned}$$

- Redundancy $O(n \lg \lg n / \lg n)$ bits, *optimal* for $O(1)$ time operations [Golynski, *TCS'07*].
- Supporting $select_1$ is similar, though a bit complicated.

Bit-Vectors: Implementing $select_1$; the idea



Bit-Vectors: Implementing $select_1$; the idea

- We will try to manage by using extra $O(n/\log \log n)$ bits.

Bit-Vectors: Implementing $select_1$; the idea

- We will try to manage by using extra $O(n/\log \log n)$ bits.
- Store answer for every $\lg n(\lg \lg n)^{th}$ 1, takes space $n/\lg \lg n$ bits.

```
○  
○○○○○●○  
○  
○○○○  
○○○  
○○○○○○○
```

Bit-Vectors: Implementing $select_1$; the idea

- We will try to manage by using extra $O(n/\log \log n)$ bits.
- Store answer for every $\lg n(\lg \lg n)^{th}$ 1, takes space $n/\lg \lg n$ bits.
- If the range r between two consecutive answers stored is of size more than $(\lg n \lg \lg n)^2$, store the positions of all the $\lg n(\lg \lg n)$ 1 in the range;

```
○
○○○○○●○
○
○○○○
○○○
○○○○○○○
```

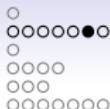
Bit-Vectors: Implementing $select_1$; the idea

- We will try to manage by using extra $O(n/\log \log n)$ bits.
- Store answer for every $\lg n(\lg \lg n)^{th}$ 1, takes space $n/\lg \lg n$ bits.
- If the range r between two consecutive answers stored is of size more than $(\lg n \lg \lg n)^2$, store the positions of all the $\lg n(\lg \lg n)$ 1 in the range; takes $(\lg n)^2(\lg \lg n)$ bits, which is at most $r/\lg \lg n$.

```
○  
○○○○○●○  
○  
○○○○  
○○○  
○○○○○○○
```

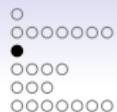
Bit-Vectors: Implementing $select_1$; the idea

- We will try to manage by using extra $O(n/\log \log n)$ bits.
- Store answer for every $\lg n(\lg \lg n)^{th}$ 1, takes space $n/\lg \lg n$ bits.
- If the range r between two consecutive answers stored is of size more than $(\lg n \lg \lg n)^2$, store the positions of all the $\lg n(\lg \lg n)$ 1 in the range; takes $(\lg n)^2(\lg \lg n)$ bits, which is at most $r/\lg \lg n$.
- Otherwise recurse.

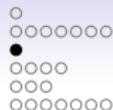


Bit-Vectors: Implementing $select_1$; the idea

- We will try to manage by using extra $O(n/\log \log n)$ bits.
- Store answer for every $\lg n(\lg \lg n)^{th}$ 1, takes space $n/\lg \lg n$ bits.
- If the range r between two consecutive answers stored is of size more than $(\lg n \lg \lg n)^2$, store the positions of all the $\lg n(\lg \lg n)$ 1 in the range; takes $(\lg n)^2(\lg \lg n)$ bits, which is at most $r/\lg \lg n$.
- Otherwise recurse. After a couple of levels, the range will be small enough ($O((\lg \lg n)^4)$) that a table look up can complete the job.



Wavelet Tree – Representing strings from a larger alphabet

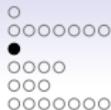


Wavelet Tree – Representing strings from a larger alphabet

Data: Sequence $S[1..n]$ of symbols from an alphabet of size σ .

Operations:

$\left. \begin{array}{l} \mathit{rank}(c, i): \text{ number of } c\text{'s in } S[1..i]. \\ \mathit{select}(c, i): \text{ position of } i\text{-th } c. \\ \mathit{access}(i): \text{ return } S[i]. \end{array} \right\} \text{ in } O(\log \sigma) \text{ time.}$



Wavelet Tree – Representing strings from a larger alphabet

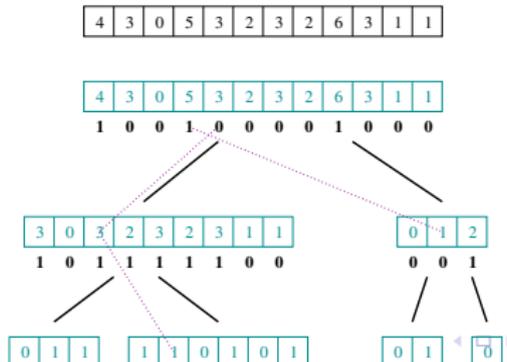
Data: Sequence $S[1..n]$ of symbols from an alphabet of size σ .

Operations:

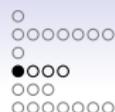
$rank(c, i)$: number of c 's in $S[1..i]$.
 $select(c, i)$: position of i -th c .
 $access(i)$: return $S[i]$.

} in $O(\log \sigma)$ time.

Store $\log_2 \sigma$ BVs: $\underbrace{n \log \sigma}_{\text{raw size}} + o(n \log \sigma)$ bits [Grossi, Vitter, *SJC '05*].



A Bit vector with only m 1s



A Bit vector with only m 1s

Data: Sequence X of n bits, x_1, \dots, x_n with m 1s.

Operations:

- $select_1(i)$.

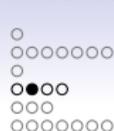
Data: Set $X = \{x_1, \dots, x_m\} \subseteq \{1, \dots, n\}$, $x_1 < x_2 < \dots < x_m$.

Operations:

- $access(i)$: return x_i .

ITLB: $\log_2 \binom{n}{m} = m \log_2(n/m) + O(m)$ bits.

[Elias, *J. ACM*'75], [Grossi/Vitter, *SICOMP*'06], [Raman et al., *TALG*'07].



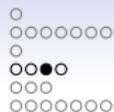
Elias-Fano Representation

Bucket according to most significant b bits.

Example. $b = 3$, $\lceil \log_2 n \rceil = 5$, $m = 7$.

x_1	0	1	0	0	0
x_2	0	1	0	0	1
x_3	0	1	0	1	1
x_4	0	1	1	0	1
x_5	1	0	0	0	0
x_6	1	0	0	1	0
x_7	1	0	1	1	1

Bucket	Keys
000	—
001	—
010	x_1, x_2, x_3
011	x_4
100	x_5, x_6
101	x_7
110	—
111	—



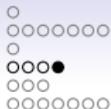
Elias-Fano

- ▷ Store only low-order bits.
- ▷ Keep sizes of all buckets.

Example
select(6)

bkt	sz	data
000	0	—
001	0	—
010	3	$\underbrace{00}_{x_1}, \underbrace{01}_{x_2}, \underbrace{11}_{x_3}$
011	1	$\underbrace{01}_{x_4}$
100	2	$\underbrace{00}_{x_5}, \underbrace{10}_{x_6}$
101	1	$\underbrace{11}_{x_7}$
110	0	—
111	0	—

Elias-Fano



Elias-Fano

- Choose $b = \lfloor \log_2 m \rfloor$ bits. In bucket: $\lceil \log_2 n \rceil - \lfloor \log_2 m \rfloor$ -bit keys.

Elias-Fano

- Choose $b = \lfloor \log_2 m \rfloor$ bits. In bucket: $\lceil \log_2 n \rceil - \lfloor \log_2 m \rfloor$ -bit keys.
- $m \log_2 n - m \log_2 m + O(m) = m \log_2(n/m) + O(m)$ bits for lower part.

```

○
○○○○○○○
○
○○○●
○○○
○○○○○○○

```

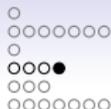
Elias-Fano

- Choose $b = \lfloor \log_2 m \rfloor$ bits. In bucket: $\lceil \log_2 n \rceil - \lfloor \log_2 m \rfloor$ -bit keys.
- $m \log_2 n - m \log_2 m + O(m) = m \log_2(n/m) + O(m)$ bits for lower part.

Encoding Bucket Sizes

Bucket no:	000	001	010	011	100	101	110	111
Bucket size:	0	0	3	1	2	1	0	0

- Use a *unary* encoding: 0, 0, 3, 1, 2, 1, 0, 0 → **110001010010111**.



Elias-Fano

- Choose $b = \lfloor \log_2 m \rfloor$ bits. In bucket: $\lceil \log_2 n \rceil - \lfloor \log_2 m \rfloor$ -bit keys.
- $m \log_2 n - m \log_2 m + O(m) = m \log_2(n/m) + O(m)$ bits for lower part.

Encoding Bucket Sizes

Bucket no:	000	001	010	011	100	101	110	111
Bucket size:	0	0	3	1	2	1	0	0

- Use a *unary* encoding: 0, 0, 3, 1, 2, 1, 0, 0 \rightarrow **110001010010111**.
- z buckets, total size $m \Rightarrow m + z = O(m)$ bits ($z = 2^{\lfloor \log_2 m \rfloor}$).
 - Overall space of E-F bit-vector is $m \log_2(n/m) + O(m)$ bits.
- In which bucket is the 6th key? \triangleright "rank₁ of 6th 0".
 - $select_1$ in $O(1)$ time.

Elias-Fano

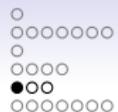
- Choose $b = \lfloor \log_2 m \rfloor$ bits. In bucket: $\lceil \log_2 n \rceil - \lfloor \log_2 m \rfloor$ -bit keys.
- $m \log_2 n - m \log_2 m + O(m) = m \log_2(n/m) + O(m)$ bits for lower part.

Encoding Bucket Sizes

Bucket no:	000	001	010	011	100	101	110	111
Bucket size:	0	0	3	1	2	1	0	0

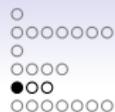
- Use a *unary* encoding: 0, 0, 3, 1, 2, 1, 0, 0 \rightarrow **110001010010111**.
- z buckets, total size $m \Rightarrow m + z = O(m)$ bits ($z = 2^{\lfloor \log_2 m \rfloor}$).
 - Overall space of E-F bit-vector is $m \log_2(n/m) + O(m)$ bits.
- In which bucket is the 6th key? \triangleright "rank₁ of 6th 0".
 - $select_1$ in $O(1)$ time.
 - Redundancy can be made $o(m)$ and membership and Rankone can also be supported (RRR01)

Tree Representations



Tree Representations

Data: n -node binary tree.

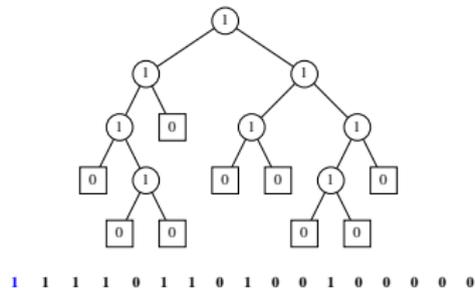


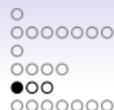
Tree Representations

Data: n -node binary tree.

Operations: Navigation (left child, right child, parent).

- Visit nodes in level-order and output 1 if internal node and 0 if external ($2n + 1$ bits) [Jacobson, *FOCS '89*]. Store sequence of bits as bit vector.



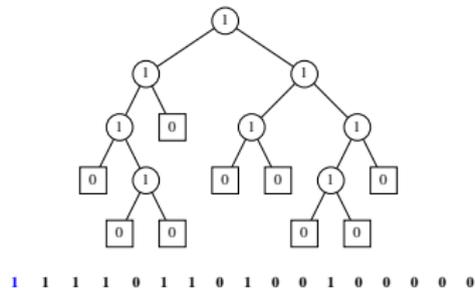


Tree Representations

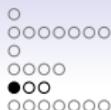
Data: n -node binary tree.

Operations: Navigation (left child, right child, parent).

- Visit nodes in level-order and output 1 if internal node and 0 if external ($2n + 1$ bits) [Jacobson, *FOCS '89*]. Store sequence of bits as bit vector.



- Number internal nodes by position of 1 in bit-string

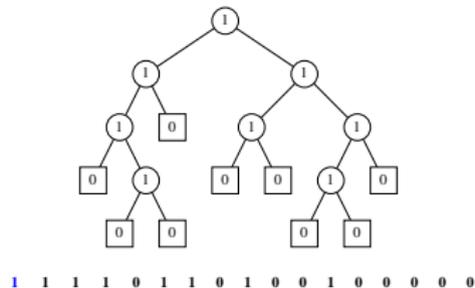


Tree Representations

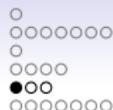
Data: n -node binary tree.

Operations: Navigation (left child, right child, parent).

- Visit nodes in level-order and output 1 if internal node and 0 if external ($2n + 1$ bits) [Jacobson, *FOCS '89*]. Store sequence of bits as bit vector.



- Number internal nodes by position of 1 in bit-string
- Left child = $2 * rank_1(i)$.

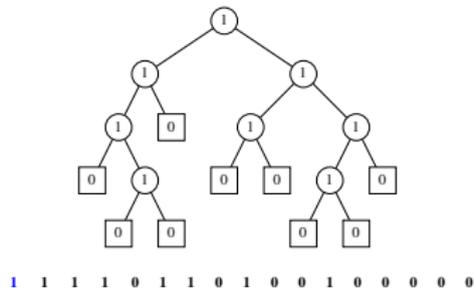


Tree Representations

Data: n -node binary tree.

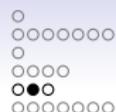
Operations: Navigation (left child, right child, parent).

- Visit nodes in level-order and output 1 if internal node and 0 if external ($2n + 1$ bits) [Jacobson, *FOCS '89*]. Store sequence of bits as bit vector.



- Number internal nodes by position of 1 in bit-string
- Left child = $2 * rank_1(i)$. E.g. Left child of node 7 = $7 * 2 = 14$. Right child = $2 * rank_1(i) + 1$. parent = $\text{select}_1(\lfloor i/2 \rfloor)$.

Tree Representations



Tree Representations

- "Optimal" representations of many kinds of trees e.g. ordinal trees (rooted arbitrary degree (un-)labelled trees, e.g. XML documents), tries.
- Wide range of $O(1)$ -time operations, e.g.:
 - ordinal trees in $2n + o(n)$ bits [Navarro, Sadakane, *TALG'12*].

Tree Representations

Pattern Matching – Compressed Text Indexing

Pattern Matching – Compressed Text Indexing

Data: Sequence T ("text") of m symbols from alphabet of size σ .

ITLB: $n \log_2 \sigma$ bits.

Operation: Given pattern P , determine if P occurs (exactly) in T (and report the number of occurrences, starting positions etc).



Pattern Matching – Compressed Text Indexing

Data: Sequence T ("text") of m symbols from alphabet of size σ .

ITLB: $n \log_2 \sigma$ bits.

Operation: Given pattern P , determine if P occurs (exactly) in T (and report the number of occurrences, starting positions etc).

- For a human genome sequence, m is about 3 billion (3×10^9) characters, and $\sigma = 4$.



Pattern Matching – Compressed Text Indexing

Data: Sequence T ("text") of m symbols from alphabet of size σ .

ITLB: $n \log_2 \sigma$ bits.

Operation: Given pattern P , determine if P occurs (exactly) in T (and report the number of occurrences, starting positions etc).

- For a human genome sequence, m is about 3 billion (3×10^9) characters, and $\sigma = 4$.
- Standard data structure is *suffix tree*, which answers this query in $O(|P|)$ time but takes $O(n \log n)$ bits of space.
- In practice, a ST is about 10-30 times larger than the text.



Pattern Matching – Compressed Text Indexing

Data: Sequence T ("text") of m symbols from alphabet of size σ .

ITLB: $n \log_2 \sigma$ bits.

Operation: Given pattern P , determine if P occurs (exactly) in T (and report the number of occurrences, starting positions etc).

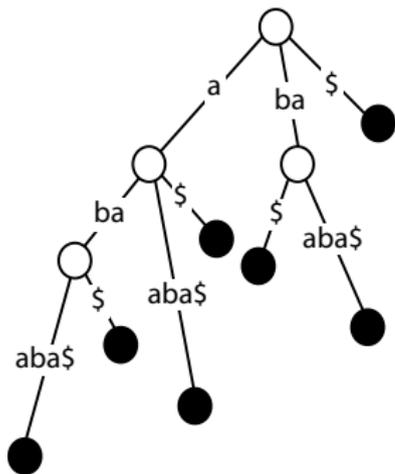
- For a human genome sequence, m is about 3 billion (3×10^9) characters, and $\sigma = 4$.
- Standard data structure is *suffix tree*, which answers this query in $O(|P|)$ time but takes $O(n \log n)$ bits of space.
- In practice, a ST is about 10-30 times larger than the text.
- A number of SDS have been developed: we'll focus on the FM-Index [Ferragina, Manzini, *JACM '05*].

Previous Popular Solution – Suffix Trees

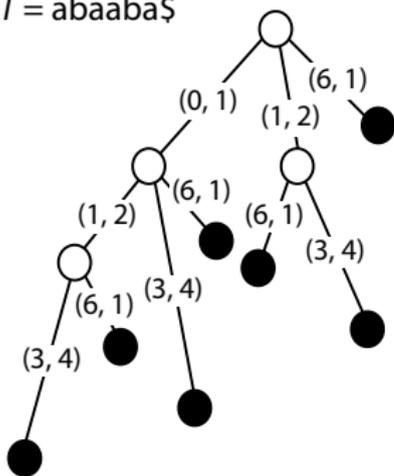
Suffix tree

$T = \text{abaaba}\$$

Idea 2: Store T itself in addition to the tree. Convert tree's edge labels to (offset, length) pairs with respect to T .



$T = \text{abaaba}\$$



Space required for suffix tree is now $O(m)$

Previous Popular Solution – Suffix Trees



Previous Popular Solution – Suffix Trees

- A (compressed) **trie** containing all the suffixes of T . The tree contains $m + 1$ leaves and at most m other nodes.

```
○  
○○○○○○○  
○  
○○○○  
○○○  
○●○○○○○
```

Previous Popular Solution – Suffix Trees

- A (compressed) **trie** containing all the suffixes of T . The tree contains $m + 1$ leaves and at most m other nodes.
- Each leaf is labelled with the starting position of the suffix ending at that leaf.



Previous Popular Solution – Suffix Trees

- A (compressed) **trie** containing all the suffixes of T . The tree contains $m + 1$ leaves and at most m other nodes.
- Each leaf is labelled with the starting position of the suffix ending at that leaf.
- Each edge has a string, that can be represented by the starting and ending position of the substring in the text.
- Overall, naive implementation takes about $4n$ words or $4n \lg n$ bits.
- Progress in succinct data structures has brought the space down to $m \lg m + O(m)$ bits (in addition to the text).
- P exists in T if and only if P is a prefix of a suffix of T . So, follow from the root matching P . If success, the leaves in the entire subtree gives the list of occurrences.



Previous Popular Solution – Suffix Trees

- A (compressed) **trie** containing all the suffixes of T . The tree contains $m + 1$ leaves and at most m other nodes.
- Each leaf is labelled with the starting position of the suffix ending at that leaf.
- Each edge has a string, that can be represented by the starting and ending position of the substring in the text.
- Overall, naive implementation takes about $4n$ words or $4n \lg n$ bits.
- Progress in succinct data structures has brought the space down to $m \lg m + O(m)$ bits (in addition to the text).
- P exists in T if and only if P is a prefix of a suffix of T . So, follow from the root matching P . If success, the leaves in the entire subtree gives the list of occurrences.
- $O(n + occ)$ to find all occurrences

Previous popular solution - Suffix Arrays

Suffix array

$T\$ = \text{abaaba}\$$ ← As with suffix tree,
 T is part of index

SA(T) =
(SA = "Suffix Array")

6	\$
5	a \$
2	a a b a \$
3	a b a \$
0	a b a a b a \$
4	b a \$
1	b a a b a \$

$m + 1$
integers

Suffix array of T is an array of integers in $[0, m]$ specifying the lexicographic order of $T\$$'s suffixes

Suffix array: querying

Is P a substring of T ?

1. For P to be a substring, it must be a prefix of ≥ 1 of T 's suffixes
2. Suffixes sharing a prefix are consecutive in the suffix array

Use binary search

6	\$
5	a \$
2	a a b a \$
3	a b a \$
0	a b a a b a \$
4	b a \$
1	b a a b a \$

Suffix array: querying

Is P a substring of T ?

Do binary search, check whether P is a prefix of the suffix there

How many times does P occur in T ?

Two binary searches yield the range of suffixes with P as prefix; size of range equals # times P occurs in T

Worst-case time bound?

$O(\log_2 m)$ bisections, $O(n)$ comparisons per bisection, so $O(n \log m)$

6	\$
5	a \$
2	a a b a \$
3	a b a \$
0	a b a a b a \$
4	b a \$
1	b a a b a \$

Previous popular solution - Suffix Arrays

Previous popular solution - Suffix Arrays

- A permutation of $\{1, 2, \dots, m\}$. $S[i]$ is the starting position of the i -th suffix in the lexicographic order.

Previous popular solution - Suffix Arrays

- A permutation of $\{1, 2, \dots, m\}$. $S[i]$ is the starting position of the i -th suffix in the lexicographic order.
- Takes $m \lg m$ bits.

Previous popular solution - Suffix Arrays

- A permutation of $\{1, 2, \dots, m\}$. $S[i]$ is the starting position of the i -th suffix in the lexicographic order.
- Takes $m \lg m$ bits. Naive binary search takes $O(n \lg m)$ time.

```
○  
○○○○○○○  
○  
○○○○  
○○○  
○○●○○○
```

Previous popular solution - Suffix Arrays

- A permutation of $\{1, 2, \dots, m\}$. $S[i]$ is the starting position of the i -th suffix in the lexicographic order.
- Takes $m \lg m$ bits. Naive binary search takes $O(n \lg m)$ time.
- With what is called an LCP array taking another $m \lg m$ bits, the search time can be brought down to $O(n + \lg m)$ bits.

The FM-Index

```

○
○○○○○○○
○
○○○○
○○○
○○○●○○○

```

The FM-Index

Based on the Burrows-Wheeler transform of the text T .

Example: $T = \text{mississippi}$

<i>F</i>											<i>L</i>
i	m	i	s	s	i	s	s	i	p		p
i	p	p	i	m	i	s	s	i	s		s
i	s	s	i	p	p	i	m	i	s		s
i	s	s	i	s	s	i	p	p	i		m
m	i	s	s	i	s	s	i	p	p		i
p	i	m	i	s	s	i	s	s	i		p
p	p	i	m	i	s	s	i	s	s		i
s	i	p	p	i	m	i	s	s	i		s
s	i	s	s	i	p	p	i	m	i		s
s	s	i	p	p	i	m	i	s	s		i
s	s	i	s	s	i	p	p	i	m		i

$\text{BWT}(T) = \text{pssmi} \text{pissii}$

Burrows-Wheeler Transform

Text transform that is useful for compression & search.

banana

banana\$

anana\$b

nana\$ba

ana\$ban

na\$bana

a\$banan

\$banana

sort



\$banana

a\$banan

ana\$ban

anana\$b

banana\$

nana\$ba

na\$bana

BWVT(banana) =

annb\$aa

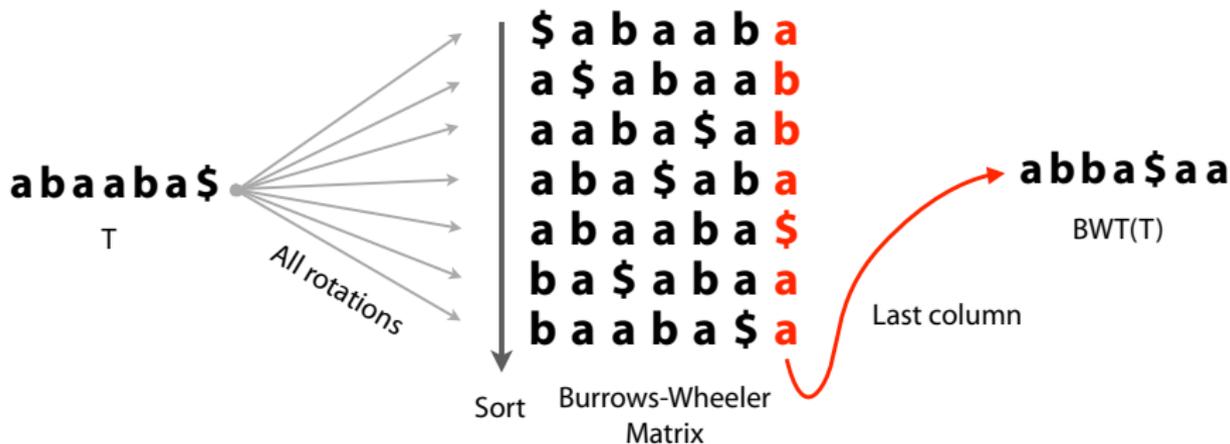
Tends to put runs of the same character together.

Makes compression work well.

“bzip” is based on this.

Burrows-Wheeler Transform

Reversible permutation of the characters of a string, used originally for compression



How is it useful for compression?

How is it reversible?

How is it an index?

Burrows-Wheeler Transform

BWM bears a resemblance to the suffix array

\$ a b a a b a
a \$ a b a a b
a a b a \$ a b
a b a \$ a b a
a b a a b a \$
b a \$ a b a a
b a a b a \$ a

BWM(T)

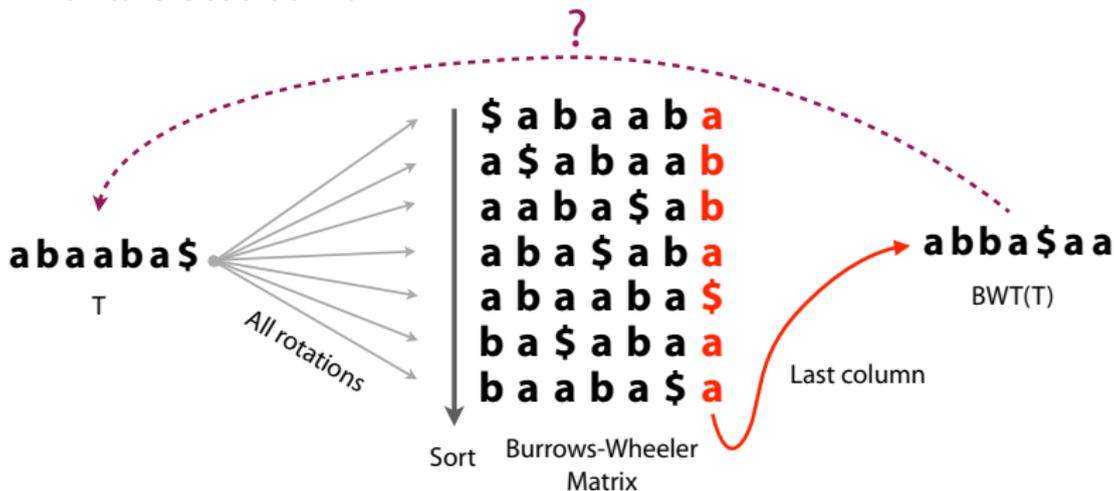
6	\$
5	a \$
2	a a b a \$
3	a b a \$
0	a b a a b a \$
4	b a \$
1	b a a b a \$

SA(T)

Sort order is the same whether rows are rotations or suffixes

Burrows-Wheeler Transform

How to reverse the BWT?



BWM has a key property called the *LF Mapping*...

Burrows-Wheeler Transform: T-ranking

Give each character in T a rank, equal to # times the character occurred previously in T . Call this the T -ranking.

a₀ b₀ a₁ a₂ b₁ a₃ \$

Now let's re-write the BWM including ranks...

Burrows-Wheeler Transform

BWM with T-ranking:

<i>F</i>							<i>L</i>
\$	a ₀	b ₀	a ₁	a ₂	b ₁	a ₃	a ₃
a ₃	\$	a ₀	b ₀	a ₁	a ₂	b ₁	b ₁
a ₁	a ₂	b ₁	a ₃	\$	a ₀	b ₀	b ₀
a ₂	b ₁	a ₃	\$	a ₀	b ₀	a ₁	a ₁
a ₀	b ₀	a ₁	a ₂	b ₁	a ₃	\$	\$
b ₁	a ₃	\$	a ₀	b ₀	a ₁	a ₂	a ₂
b ₀	a ₁	a ₂	b ₁	a ₃	\$	a ₀	a ₀

Look at first and last columns, called *F* and *L*

And look at just the **a**s

as occur in the same order in *F* and *L*. As we look down columns, in both cases we see: **a₃, a₁, a₂, a₀**

Burrows-Wheeler Transform

BWM with T-ranking:

<i>F</i>							<i>L</i>
\$	a ₀	b ₀	a ₁	a ₂	b ₁	a ₃	
a ₃	\$	a ₀	b ₀	a ₁	a ₂	b₁	
a ₁	a ₂	b₁	a ₃	\$	a ₀	b₀	
a ₂	b₁	a ₃	\$	a ₀	b ₀	a ₁	
a ₀	b ₀	a ₁	a ₂	b₁	a ₃	\$	
b₁	a ₃	\$	a ₀	b ₀	a ₁	a ₂	
b₀	a ₁	a ₂	b₁	a ₃	\$	a ₀	

Same with **b**s: **b₁**, **b₀**

Burrows-Wheeler Transform: LF Mapping

BWM with T-ranking:

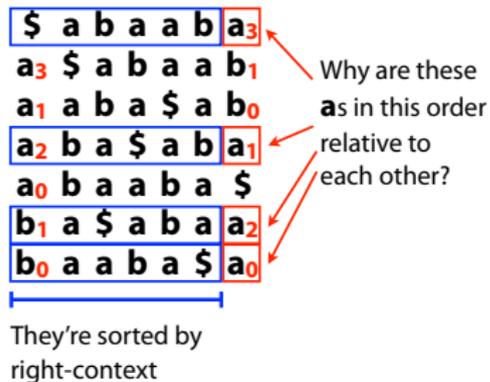
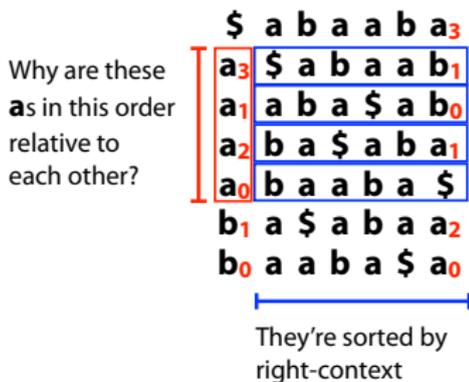
	<i>F</i>							<i>L</i>
	\$	a ₀	b ₀	a ₁	a ₂	b ₁	a ₃	
	a ₃	\$	a ₀	b ₀	a ₁	a ₂	b ₁	
	a ₁	a ₂	b ₁	a ₃	\$	a ₀	b ₀	
	a ₂	b ₁	a ₃	\$	a ₀	b ₀	a ₁	
	a ₀	b ₀	a ₁	a ₂	b ₁	a ₃	\$	
	b ₁	a ₃	\$	a ₀	b ₀	a ₁	a ₂	
	b ₀	a ₁	a ₂	b ₁	a ₃	\$	a ₀	

LF Mapping: The i^{th} occurrence of a character c in L and the i^{th} occurrence of c in F correspond to the *same* occurrence in T

However we rank occurrences of c , ranks appear in the same order in F and L

Burrows-Wheeler Transform: LF Mapping

Why does the LF Mapping hold?



Occurrences of c in F are sorted by right-context. Same for L !

Whatever ranking we give to characters in T , rank orders in F and L will match

Burrows-Wheeler Transform: LF Mapping

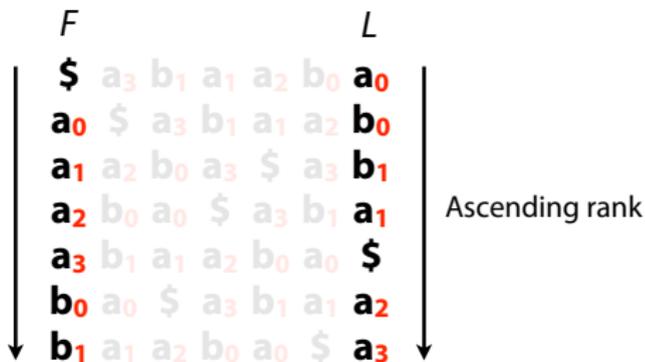
BWM with T-ranking:

<i>F</i>								<i>L</i>
\$	a ₀	b ₀	a ₁	a ₂	b ₁	a ₃		
a ₃	\$	a ₀	b ₀	a ₁	a ₂	b ₁		
a ₁	a ₂	b ₁	a ₃	\$	a ₀	b ₀		
a ₂	b ₁	a ₃	\$	a ₀	b ₀	a ₁		
a ₀	b ₀	a ₁	a ₂	b ₁	a ₃	\$		
b ₁	a ₃	\$	a ₀	b ₀	a ₁	a ₂		
b ₀	a ₁	a ₂	b ₁	a ₃	\$	a ₀		

We'd like a different ranking so that for a given character, ranks are in ascending order as we look down the F / L columns...

Burrows-Wheeler Transform: LF Mapping

BWM with B-ranking:



F now has very simple structure: a \$, a block of **a**s with ascending ranks, a block of **b**s with ascending ranks

Burrows-Wheeler Transform

<i>F</i>	<i>L</i>	
\$	a₀	
a₀	b₀	
a₁	b₁	← Which BWM row <i>begins</i> with b₁ ?
a₂	a₁	Skip row starting with \$ (1 row)
a₃	\$	Skip rows starting with a (4 rows)
b₀	a₂	Skip row starting with b₀ (1 row)
row 6 → b₁	a₃	Answer: row 6

Burrows-Wheeler Transform

Say T has 300 **A**s, 400 **C**s, 250 **G**s and 700 **T**s and $\$ < \mathbf{A} < \mathbf{C} < \mathbf{G} < \mathbf{T}$

Which BWM row (0-based) begins with **G**₁₀₀? (Ranks are B-ranks.)

Skip row starting with **\$** (1 row)

Skip rows starting with **A** (300 rows)

Skip rows starting with **C** (400 rows)

Skip first 100 rows starting with **G** (100 rows)

Answer: row $1 + 300 + 400 + 100 =$ **row 801**

Burrows-Wheeler Transform: reversing

Reverse BWT(T) starting at right-hand-side of T and moving left

Start in first row. F must have $\$$. L contains character just prior to $\$$: a_0

a_0 : LF Mapping says this is same occurrence of a as first a in F . Jump to row beginning with a_0 . L contains character just prior to a_0 : b_0 .

Repeat for b_0 , get a_2

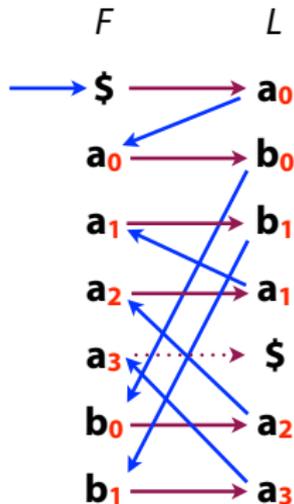
Repeat for a_2 , get a_1

Repeat for a_1 , get b_1

Repeat for b_1 , get a_3

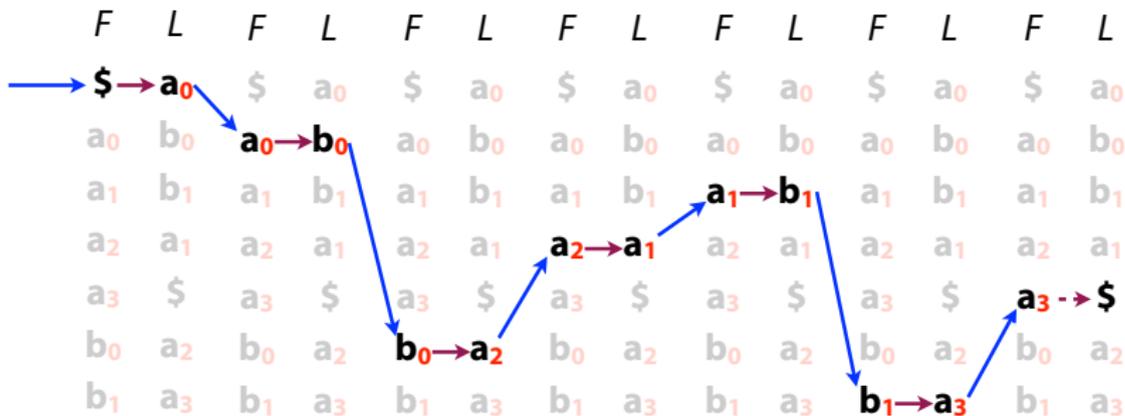
Repeat for a_3 , get $\$$, done

Reverse of chars we visited = $a_3 b_1 a_1 a_2 b_0 a_0 \$ = T$



Burrows-Wheeler Transform: reversing

Another way to visualize reversing BWT(T):



T: a₃ b₁ a₁ a₂ b₀ a₀ \$

Burrows-Wheeler Transform

We've seen how BWT is useful for compression:

Sorts characters by right-context, making a more compressible string

And how it's reversible:

Repeated applications of LF Mapping, recreating T from right to left

How is it used as an index?

FM Index

FM Index: an index combining the BWT *with a few small auxiliary data structures*

“FM” supposedly stands for “Full-text Minute-space.”
(But inventors are named Ferragina and Manzini)

Core of index consists of F and L from BWM:

F can be represented very simply
(1 integer per alphabet character)

And L is compressible

Potentially very space-economical!

F								L
\$	a	b	a	a	b	a		
a	\$	a	b	a	a	b		
a	a	b	a	\$	a	b		
a	b	a	\$	a	b	a		
a	b	a	a	b	a	\$		
b	a	\$	a	b	a	a		
b	a	a	b	a	\$	a		

└──────────┘
Not stored in index

Paolo Ferragina, and Giovanni Manzini. "Opportunistic data structures with applications." *Foundations of Computer Science, 2000. Proceedings. 41st Annual Symposium on*. IEEE, 2000.

FM Index: querying

Though BWM is related to suffix array, we can't query it the same way

\$	a	b	a	a	b	a
a	\$	a	b	a	a	b
a	a	b	a	\$	a	b
a	b	a	\$	a	b	a
a	b	a	a	b	a	\$
b	a	\$	a	b	a	a
b	a	a	b	a	\$	a



6	\$
5	a \$
2	a a b a \$
3	a b a \$
0	a b a a b a \$
4	b a \$
1	b a a b a \$

We don't have these columns; binary search isn't possible

FM Index: querying

Look for range of rows of $BWM(T)$ with P as prefix

Do this for P 's shortest suffix, then extend to successively longer suffixes until range becomes empty or we've exhausted P

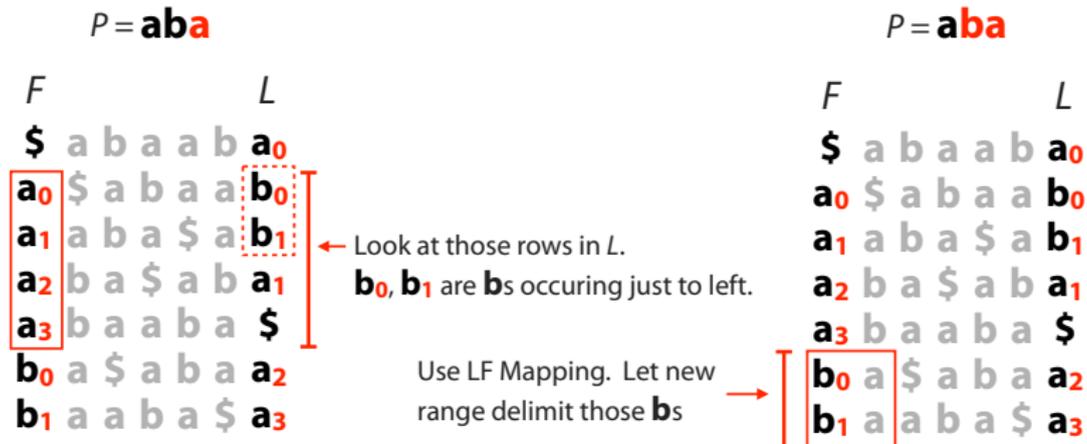
$P = \mathbf{aba}$

Easy to find all the rows beginning with \mathbf{a} , thanks to F 's simple structure

F		L				
$\$$	a	b	a	a	b	a_3
a_0	$\$$	a	b	a	a	b_1
a_1	a	b	a	$\$$	a	b_0
a_2	b	a	$\$$	a	b	a_1
a_3	b	a	a	b	a	$\$$
b_0	a	$\$$	a	b	a	a_2
b_1	a	a	b	a	$\$$	a_0

FM Index: querying

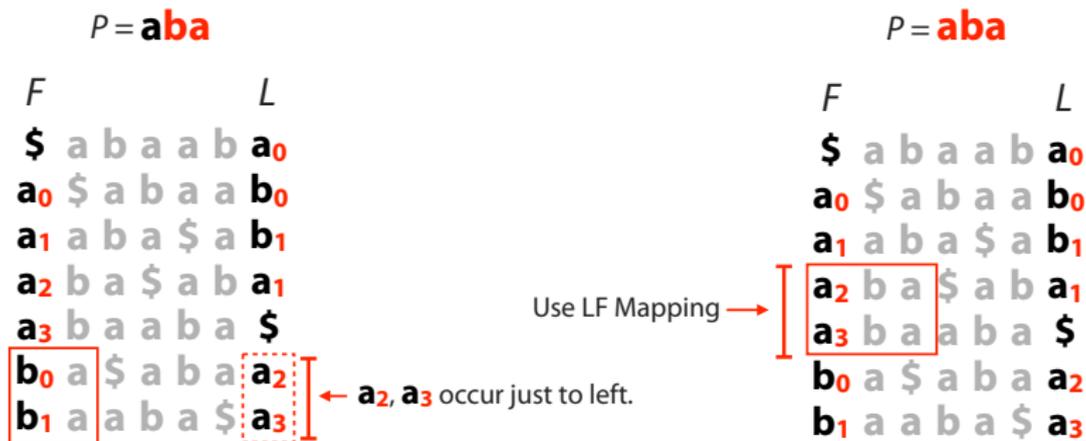
We have rows beginning with **a**, now we seek rows beginning with **ba**



Now we have the rows with prefix **ba**

FM Index: querying

We have rows beginning with **ba**, now we seek rows beginning with **aba**



Now we have the rows with prefix **aba**

FM Index: querying

$P = \mathbf{aba}$

Now we have the same range, $[3, 5)$, we would have got from querying suffix array

	<i>F</i>		<i>L</i>					
	\$	a	b	a	a	b	a	a_0
	a_0	\$	a	b	a	a	b	b_0
	a_1	a	b	a	\$	a	b	b_1
$[3, 5)$	a_2	b	a	\$	a	b	a_1	
	a_3	b	a	a	b	a	\$	
	b_0	a	\$	a	b	a	a_2	
	b_1	a	a	b	a	\$	a_3	

Where are these?

6	\$
5	a \$
2	a a b a \$
3	a b a \$
0	a b a a b a \$
4	b a \$
1	b a a b a \$

Unlike suffix array, we don't immediately know *where* the matches are in T...

FM Index: querying

When P does not occur in T , we will eventually fail to find the next character in L :

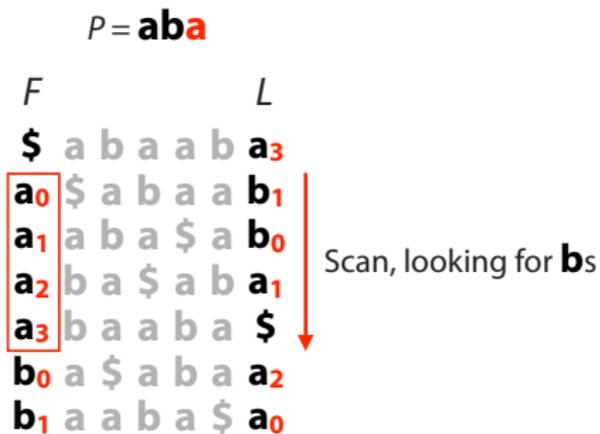
$P = \mathbf{bba}$

F						L	
$\$$	a	b	a	a	b	$\mathbf{a_0}$	
$\mathbf{a_0}$	$\$$	a	b	a	a	$\mathbf{b_0}$	
$\mathbf{a_1}$	a	b	a	$\$$	a	$\mathbf{b_1}$	
$\mathbf{a_2}$	b	a	$\$$	a	b	$\mathbf{a_1}$	
$\mathbf{a_3}$	b	a	a	b	a	$\$$	
Rows with ba prefix	$\mathbf{b_0}$	a	$\$$	a	b	a	$\mathbf{a_2}$
	$\mathbf{b_1}$	a	a	b	a	$\$$	$\mathbf{a_3}$

← No **bs!**

FM Index: querying

If we *scan* characters in the last column, that can be very slow, $O(m)$



FM Index: lingering issues

- (1) Scanning for preceding character is slow

	\$	a	b	a	a	b	a	a₀
a₀	\$	a	b	a	a	b	a	b₀
a₁	a	b	a	\$	a	b	a	b₁
a₂	b	a	\$	a	b	a	a	a₁
a₃	b	a	a	b	a	a	\$	
b₀	a	\$	a	b	a	a	a	a₂
b₁	a	a	b	a	\$	a	a	a₃

$O(m)$
 scan

- (2) Storing ranks takes too much space

```

def reverseBwt(bw):
    """ Make T from BWT(T) """
    ranks, tots = rankBwt(bw)
    first = firstCol(tots)
    rowi = 0
    t = "$"
    while bw[rowi] != '$':
        c = bw[rowi]
        t = c + t
        rowi = first[c][0] + ranks[rowi]
    return t
  
```

m integers

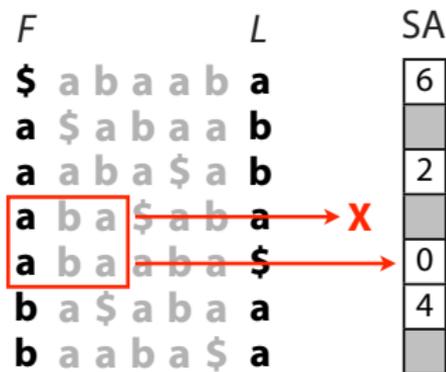
- (3) Need way to find where matches occur in T :

	\$	a	b	a	a	b	a	a₀
a₀	\$	a	b	a	a	b	a	b₀
a₁	a	b	a	\$	a	b	a	b₁
a₂	b	a	\$	a	b	a	a	a₁
a₃	b	a	a	b	a	a	\$	
b₀	a	\$	a	b	a	a	a	a₂
b₁	a	a	b	a	\$	a	a	a₃

Where?

FM Index: resolving offsets

Idea: store some, but not all, entries of the suffix array

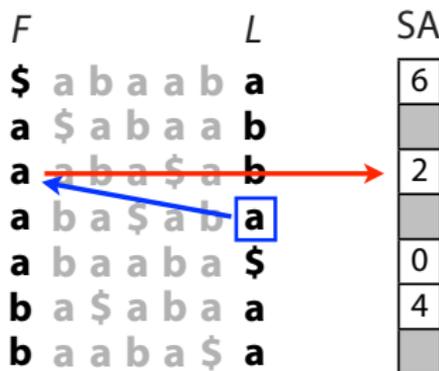


Lookup for row 4 succeeds - we kept that entry of SA

Lookup for row 3 fails - we discarded that entry of SA

FM Index: resolving offsets

But LF Mapping tells us that the **a** at the end of row 3 corresponds to...
...the **a** at the beginning of row 2



And row 2 has a suffix array value = 2

So row 3 has suffix array value = 3 = 2 (row 2's SA val) + 1 (# steps to row 2)

If saved SA values are $O(1)$ positions apart in T , resolving offset is $O(1)$ time

FM Index: problems solved

Solved! At the expense of adding some SA values ($O(m)$ integers) to index
Call this the "SA sample"

- (3) Need a way to find where these occurrences are in T :

\$	a	b	a	a	b	a ₀
a ₀	\$	a	b	a	a	b ₀
a ₁	a	b	a	\$	a	b ₁
a ₂	b	a	\$	a	b	a ₁
a ₃	b	a	a	b	a	\$
b ₀	a	\$	a	b	a	a ₂
b ₁	a	a	b	a	\$	a ₃

With SA sample we can do this in
 $O(1)$ time per occurrence

To Summarize (FM index)

○
○○○○○○○
○
○○○○
○○○
○○○○●○○

To Summarize (FM index)

- Existence of P in T , and
- the number of occurrences (occ) of P in T

To Summarize (FM index)

- Existence of P in T , and
- the number of occurrences (occ) of P in T

can be determined in $O(n)$ time using

```
○  
○○○○○○○  
○  
○○○○  
○○○  
○○○○●○○
```

To Summarize (FM index)

- Existence of P in T , and
- the number of occurrences (occ) of P in T

can be determined in $O(n)$ time using

- $m \lg \sigma$ bits, for BWT (last column)
- $o(m \lg \sigma)$ bits for rank
- $\sigma \lg m$ bits for count of each character (first column)

○
○○○○○○○
○
○○○○
○○○
○○○○●○○

To Summarize (FM index)

- Existence of P in T , and
- the number of occurrences (occ) of P in T

can be determined in $O(n)$ time using

- $m \lg \sigma$ bits, for BWT (last column)
- $o(m \lg \sigma)$ bits for rank
- $\sigma \lg m$ bits for count of each character (first column)

and the position of all occurrences of P in T can be determined in

```
○  
○○○○○○○  
○  
○○○○  
○○○  
○○○○●○○
```

To Summarize (FM index)

- Existence of P in T , and
- the number of occurrences (occ) of P in T

can be determined in $O(n)$ time using

- $m \lg \sigma$ bits, for BWT (last column)
- $o(m \lg \sigma)$ bits for rank
- $\sigma \lg m$ bits for count of each character (first column)

and the position of all occurrences of P in T can be determined in

- additional $O(k \text{ occ})$ time, using
- an additional $(m \lg m)/k$ bits of space (using a sampled suffix array)

```
○  
○○○○○○○  
○  
○○○○  
○○○  
○○○○●○○
```

To Summarize (FM index)

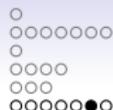
- Existence of P in T , and
- the number of occurrences (occ) of P in T

can be determined in $O(n)$ time using

- $m \lg \sigma$ bits, for BWT (last column)
- $o(m \lg \sigma)$ bits for rank
- $\sigma \lg m$ bits for count of each character (first column)

and the position of all occurrences of P in T can be determined in

- additional $O(k \text{ occ})$ time, using
- an additional $(m \lg m)/k$ bits of space (using a sampled suffix array)
- For example, $O(\text{occ} \lg m)$ time using additional $O(m)$ bits of space.



Contrasting with Suffix Arrays and Suffix Trees

FM Index	$O(m \lg \sigma)$ bits 1.5GB for human genome	$O(n)$ time for finding existence and occ $O(n + occ \lg m)$ for finding all occurrences
Suffix Array	$2m \lg m$ bits + text about 12GB for human genome	$O(n + \lg m)$ time for all operations
Suffix Tree	$3m \lg m$ bits + text about 47GB in MUMmer for human genome; with optimization ($m \lg m + O(m)$ bits)	$O(n)$ time for boolean query $O(n + occ)$ for finding all occurrences useful for many other operations

Introduction

Data Structures

Goals

Bit Vectors

Strings from a larger alphabet

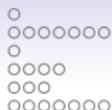
Sparse Bit Vectors

Trees

Burrows-Wheeler Transform and Indexing

Libraries

Conclusions



Libraries

- A number of good implementations of succinct data structures in C++ are available.
- Different platforms, coding styles:
 - `sds1-lite` (Gog, Petri et al. U. Melbourne).
 - `succinct` (Grossi and Ottaviano, U. Pisa).
 - `Sux4J` (Vigna, U. Milan, Java).
 - `LIBCDS` (Claude and Navarro, Akori and U. Chile).
- All open-source and available as Git repositories.

Conclusions

- SDS are a relatively mature field in terms of breadth of problems considered.

Conclusions

- SDS are a relatively mature field in terms of breadth of problems considered.
- Quite practical; FM index has been implemented in BIO software (Bowtie).

Conclusions

- SDS are a relatively mature field in terms of breadth of problems considered.
- Quite practical; FM index has been implemented in BIO software (Bowtie).
- Some foundational questions still not addressed (e.g. lower bounds). at least in dynamic SDS.

Thank You

Thank You

Special thanks to Rajeev Raman (Leicester University) and Ben Langmead (Johns Hopkins) for some of the slides