# Burrows-Wheeler Transform
## (some of) its properties and applications

Rossano Venturini

Department of Computer Science
University of Pisa

# Basic Concepts in Data Compression

- ***Lossless text data compression***:
  - We would like to design a compressor that, give a text in input, represents is using the smallest possible number of bits. From this representation we must be able to reconstruct the original text without any loss of information.

- **Historical motivations**:
  - Save storage space and/or bandwidth.

S | a | b | r | a | c | a | d | a | b | r | a

C

# 0-th order compressors

S | a | b | r | a | c | a | d | a | b | r | a

C | | | | | | | | | | |

- Build a table: for each symbol stores its frequency

| char | freq |
|------|------|
| a | 5/11 |
| b | 2/11 |
| c | 1/11 |
| d | 1/11 |
| r | 2/11 |

# 0-th order compressors

S | a | b | r | a | c | a | d | a | b | r | a |

C | | | | | | | | | | | |

- Build a table: for each symbol stores its frequency

- Assign a codeword to each symbol. So that,

  - Decompression: Codewords must be uniquely decodable.

| char | freq | code |
|------|------|------|
| a | 5/11 | 0 |
| b | 2/11 | 100 |
| c | 1/11 | 101 |
| d | 1/11 | 110 |
| r | 2/11 | 111 |

# 0-th order compressors

S | a | b | r | a | c | a | d | a | b | r | a |

C | | | | | | | | | | | |

- Build a table: for each symbol stores its frequency

- Assign a codeword to each symbol. So that,

  - Decompression: Codewords must be uniquely decodable.

  - Minimize compress size: Shortest codewords must be assigned to most frequent symbols.

| char | freq | code |
|------|------|------|
| a | 5/11 | 0 |
| b | 2/11 | 100 |
| c | 1/11 | 101 |
| d | 1/11 | 110 |
| r | 2/11 | 111 |

# 0-th order compressors

S | a | b | r | a | c | a | d | a | b | r | a |

C | 0 100 111 0 101 0 110 0 100 111 0 |

- Build a table: for each symbol stores its frequency

- Assign a codeword to each symbol. So that,

  - Decompression: Codewords must be uniquely decodable.

  - Minimize compress size: Shortest codewords must be assigned to most frequent symbols.

- Replace each symbol with its codeword. Compress is C+Table

| char | freq | code |
|------|------|------|
| a | 5/11 | 0 |
| b | 2/11 | 100 |
| c | 1/11 | 101 |
| d | 1/11 | 110 |
| r | 2/11 | 111 |

# 0-th order compressors

S | a | b | r | a | c | a | d | a | b | r | a

C | 0 100 111 0 101 0 110 0 100 111 0

- Decompression is easy:
  - Scan C from left to right

| char | freq | code |
|------|------|------|
| a | 5/11 | 0 |
| b | 2/11 | 100 |
| c | 1/11 | 101 |
| d | 1/11 | 110 |
| r | 2/11 | 111 |

# 0-th order compressors

S | a | b | r | a | c | a | d | a | b | r | a

C | 0 100 111 0 101 0 110 0 100 111 0

- Decompression is easy:
  - Scan C from left to right
  - Every time we identify a codeword, we emit the corresponding symbol.

| char | freq | code |
|------|------|------|
| a | 5/11 | 0 |
| b | 2/11 | 100 |
| c | 1/11 | 101 |
| d | 1/11 | 110 |
| r | 2/11 | 111 |

# 0-th order compressors

S | a | b | r | a | c | a | d | a | b | r | a |

C | 0 100 111 0 101 0 110 0 100 111 0 |

- Decompression is easy:
  - Scan C from left to right
  - Every tim̲e w̲
    codewo̲
    corres̲

**Low compression! we don't exploit regularities in text**

| char | freq | code |
|------|------|------|
| a | 5/11 | 0 |
| b | 2/11 | 100 |
| c | 1/11 | 101 |
| d | 1/11 | 110 |
| r | 2/11 | 111 |

# High order compressors

- We can achieve better compression if the codeword we assign to a symbol also depends on the k symbols preceding it (its context)

# High order compressors

- We can achieve better compression if the codeword we assign to a symbol also depends on the k symbols preceding it (its context)

S | a | b | r | a | c | a | d | a | b | r | a

Build a table for each context
of length k in S

# High order compressors

- We can achieve better compression if the codeword we assign to a symbol also depends on the k symbols preceding it (its context)

S | a | b | r | a | c | a | d | a | b | r | a |

Build a table for each context
of length k in S

k=2

context = ab

# High order compressors

- We can achieve better compression if the codeword we assign to a symbol also depends on the k symbols preceding it (its context)

S | a | b | r | a | c | a | d | a | b | r | a |

Build a table for each context
of length k in S

k=2

context = ab

**char  freq code**

| char | freq | code |
|---|---|---|
| a | 0/2 | - |
| b | 0/2 | - |
| c | 0/2 | - |
| d | 0/2 | - |
| r | 2/2 | 0 |

# High order compressors

- We can achieve better compression if the codeword we assign to a symbol also depends on the k symbols preceding it (its context)

S | a | b | r | a | c | a | d | a | b | r | a

Build a table for each context of length k in S

We need to store a table for each context of size k.

| char | freq | code |
|------|------|------|
| a | 0/2 | - |
| b | 0/2 | - |
| c | 0/2 | - |
| d | 0/2 | - |
| r | 2/2 | 0 |

# The models are the problem!

- The compression improves because we better predict the next symbol.

# The models are the problem!

- The compression improves because we better predict the next symbol.

- Problem:

  - Larger is k, smaller is the compress

# The models are the problem!

- The compression improves because we better predict the next symbol.

- <span style="color:red">Problem:</span>

  - Larger is k, smaller is the compress
  - but we have to store more tables:
    - $O(\sigma^{k+1} \log \sigma)$ bits in the worst case

# The models are the problem!

- The compression improves because we better predict the next symbol.

- <span style="color:red">Problem:</span>

  - Larger is k, smaller is the compress
  - but we have to store more tables:
    - $O(\sigma^{k+1} \log \sigma)$ bits in the worst case
- Since compress size = |C|+ size tables, this approach require a lot of tuning in order to find the best value of k (i.e., the value of k that minimizes compress size)

# The models are the problem!

- The compression improves because we better predict the next symbol.

- <span style="color:red">Problem:</span>

  - Larger is k, smaller is the compress
  - but we have to store more tables:
    - $O(\sigma^{k+1} \log \sigma)$ bits in the worst case

- Since compress size = |C|+ size tables, this approach requires a lot of tuning in order to find the best value of k (i.e., the value of k that minimizes compress size)

- Instead, we would like to have a method that use a 0-th order compressor without care about the length of the contexts

# Rearranging the input

- Idea!
  - Permute the input so that it is more compressible by a 0-th order compressor
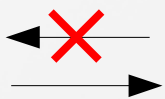
# Rearranging the input

- Idea!
  - Permute the input so that it is more compressible by a 0-th order compressor
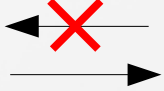- Easiest way: sort the symbol lexicographically

  `abracadabra#` $\longrightarrow$ `#aaaaabbcdrr`

# Rearranging the input

- Idea!
  - Permute the input so that it is more compressible by a 0-th order compressor
- Easiest way: sort the symbol lexicographically

abracadabra#  ⟶  #aaaaabbcdrr

(#,1)(a,5)(b,2)(c,1)(d,1)(r,2)

# Rearranging the input

- Idea!

  – Permute the input ⬚ compressible by a ⬚

- Easiest way: sort the ⬚

Best compression you can achieve!
Decoder must know at least the alphabet distribution.

```
abracadabra#  ⟶  #aa⬚bcdrr
```

```
(#,1)(a,5)(b,2)(c,1)(d,1)(r,2)
```

# Rearranging the input

- Idea!

  - Permute the input so that it is more compressible by a 0-th order compressor

- Easiest way: sort the symbol lexicographically

  abracadabra# $\longrightarrow$ #aaaaabbcdrr

  Which is the problem?

# Rearranging the input

- Idea!

  – Permute the input so that it is more compressible by a 0-th order compressor

- Easiest way: sort the symbol lexicographically

  abracadabra#  ✗→  #aaaaabbcdrr

  **Which is the problem?**

  **The transformation is not reversible!**
  There are 997.920 distinct strings with
  this alphabet distribution!

# Rearranging the input

- Idea!

  – Permute the input so that it is more compressible by a 0-th order compressor

- Easiest way: sort the symbol lexicographically

abracadabra#  ✗ →  #aaaaabbcdrr

- What do you think about this one?

ard#rcaaaabb

# Rearranging the input

- Idea!
  - Permute the input so that it is more compressible by a 0-th order compressor
- Easiest way: sort the symbol lexicographically

abracadabra#  ⟶  #aaaaabbcdrr

- What do you think about this one?

ard#rcaaaabb

Similar, but it is reversible!

Let us given S = abracadabra#

# Burrows-Wheeler Transform

Let us given S = abracadabra#


abracadabra#

# Burrows-Wheeler Transform

Let us given S = abracadabra#

```
abracadabra#
bracadabra#a
```

# Burrows-Wheeler Transform

Let us given S = abracadabra#

```
abracadabra#
bracadabra#a
racadabra#ab
```

Let us given S = abracadabra#

```
abracadabra#
bracadabra#a
racadabra#ab
acadabra#abr
cadabra#abra
adabra#abrac
dabra#abraca
abra#abracad
bra#abracada
ra#abracadab
a#abracadabr
#abracadabra
```
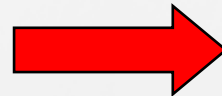
# Burrows-Wheeler Transform

Let us given S = abracadabra#

```
abracadabra#
bracadabra#a
racadabra#ab
acadabra#abr
cadabra#abra
adabra#abrac
dabra#abraca
abra#abracad
bra#abracada
ra#abracadab
a#abracadabr
#abracadabra
```

Let us given S = abracadabra#

```
abracadabra#
bracadabra#a
racadabra#ab
acadabra#abr
cadabra#abra
adabra#abrac
dabra#abraca
abra#abracad
bra#abracada
ra#abracadab
ra#abracadabr
#abracadabra
```

Let us given S = abracadabra#

```
abracadabra#
bracadabra#a
racadabra#ab
acadabra#abr
cadabra#abra
adabra#abrac
dabra#abraca
abra#abracad
bra#abracada
ra#abracadab
a#abracadabr
#abracadabra
```
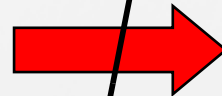
Sort the rows

Let us given S = abracadabra#

```
abracadabra#
bracadabra#a
racadabra#ab
acadabra#abr
cadabra#abra
adabra#abrac
dabra#abraca
abra#abracad
bra#abracada
ra#abracadab
a#abracadabr
#abracadabra
```

**Sort the rows**  →

```
#  abracadabr  a
```

Let us given S = abracadabra#

```
abracadabra#
bracadabra#a
racadabra#ab
acadabra#abr
cadabra#abra
adabra#abrac
dabra#abraca
abra#abracad
bra#abracada
ra#abracadab
a#abracadabr
#abracadabra
```

**Sort the rows**

```
#   abracadabr   a
a   #abracadab   r
```

Let us given S = abracadabra#

```
abracadabra#
bracadabra#a
racadabra#ab
acadabra#abr
cadabra#abra
adabra#abrac
dabra#abraca
abra#abracad
bra#abracada
ra#abracadab
a#abracadabr
#abracadabra
```

Sort the rows

```
#   abracadabr   a
a   #abracadab   r
a   bra#abraca   d
```
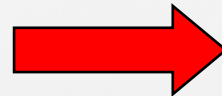
# Burrows-Wheeler Transform

Let us given S = abracadabra#

```
abracadabra#
bracadabra#a
racadabra#ab
acadabra#abr
cadabra#abra
adabra#abrac
dabra#abraca
abra#abracad
bra#abracada
ra#abracadab
a#abracadabr
#abracadabra
```

Sort the rows →

```
#   abracadabr   a
a   #abracadab   r
a   bra#abraca   d
a   bracadabra   #
a   cadabra#ab   r
a   dabra#abra   c
b   ra#abracad   a
b   racadabra#   a
c   adabra#abr   a
d   abra#abrac   a
r   a#abracada   b
r   acadabra#a   b
```
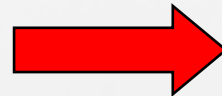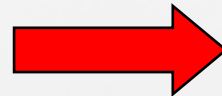
# Burrows-Wheeler Transform

Let us given S = abracadabra#

|  | F |  | L |
|---|---|---|---|
| abracadabra# | # | abracadabr | a |
| bracadabra#a | a | #abracadab | r |
| racadabra#ab | a | bra#abraca | d |
| acadabra#abr | a | bracadabra | # |
| cadabra#abra | a | cadabra#ab | r |
| adabra#abrac | a | dabra#abra | c |
| dabra#abraca | b | ra#abracad | a |
| abra#abracad | b | racadabra# | a |
| bra#abracada | c | adabra#abr | a |
| ra#abracadab | d | abra#abrac | a |
| a#abracadabr | r | a#abracada | b |
| #abracadabra | r | acadabra#a | b |

**Sort the rows** →

# Burrows-Wheeler Transform

Let us given S = abracadabra#

|  | F | discarded | L |
|---|---|---|---|
| abracadabra# | # | abracadabr | a |
| bracadabra#a | a | #abracadab | r |
| racadabra#ab | a | bra#abraca | d |
| acadabra#abr | a | bracadabra | # |
| cadabra#abra | a | cadabra#ab | r |
| adabra#abrac | a | dabra#abra | c |
| dabra#abraca | b | ra#abracad | a |
| abra#abracad | b | racadabra# | a |
| bra#abracada | c | adabra#abr | a |
| ra#abracadab | d | abra#abrac | a |
| a#abracadabr | r | a#abracada | b |
| #abracadabra | r | acadabra#a | b |

Sort the rows ➡

[Burrows-Wheeler, 1994]

Let us given S = abracadabra#

|   | | F | | discarded | | L |
|---|---|---|---|---|---|---|
| abracadabra# | | # | | abracadabr | | a |
| bracadabra#a | | a | | #abracadab | | r |
| racadabra#ab | | a | | bra#abraca | | d |
| acadabra#abr | Sort the rows → | a | | bracadabra | | # |
| cadabra#abra | | a | | cadabra#ab | | r |
| adabra#abrac | | a | | dabra#abra | | c |
| dabra#abraca | | b | | ra#abracad | | a |
| abra#abracad | | b | | racadabra# | | a |
| bra#abracada | | c | | adabra#abr | | a |
| ra#abracadab | | d | | abra#abrac | | a |
| a#abracadabr | | r | | a#abracada | | b |
| #abracadabra | | r | | acadabra#a | | b |

⟷ S

# Burrows-Wheeler Transform:
# why it works

```
F                                    ...   L
ot look upon his like again.    ...   n
ot look upon me; Lest with th...      n
ot love on the wing,-- As I p...      h
ot love your father; But that...      n
ot made them well, they imita...      n
ot madness That I have utter'...      n
ot me? Ham. No, by the rood,    ...   g
ot me; no, nor woman neither,...      n
ot me'? Ros. To think, my lor...      n
ot mend his pace with beating...      n
ot mine. Ham. No, nor mine no...      n
ot mine own. Besides, to be d...      n
ot mock me, fellow-student. I...      n
ot monstrous that this player...      n
ot more like. Ham. But where    ...   n
ot more, my lord. Ham. Is not...      j
ot more native to the heart,    ...   n
ot more ugly to the thing tha...      n
ot move thus. Oph. You must s...      n
ot much approve me.--Well, si...      n
```

Shakespeare's Hamlet

# Burrows-Wheeler Transform: why it works

| $F$ | ... | $L$ |
|---|---|---|
| ot look upon his like again. | ... | n |
| ot look upon me; Lest with th | ... | n |
| ot love on the wing,-- As I p | ... | h |
| ot love your father; But that | ... | n |
| ot made them well, they imita | ... | n |
| ot madness That I have utter' | ... | n |
| ot me? Ham. No, by the rood, | ... | g |
| ot me; no, nor woman neither, | ... | n |
| ot me'? Ros. To think, my lor | ... | n |
| ot mend his pace with beating | ... | n |
| ot mine. Ham. No, nor mine no | ... | n |
| ot mine own. Besides, to be d | ... | n |
| ot mock me, fellow-student. I | ... | n |
| ot monstrous that this player | ... | n |
| ot more like. Ham. But where | ... | n |
| ot more, my lord. Ham. Is not | ... | j |
| ot more native to the heart, | ... | n |
| ot more ugly to the thing tha | ... | n |
| ot move thus. Oph. You must s | ... | n |
| ot much approve me.--Well, si | ... | n |

Shakespeare's Hamlet

**L** is locally homogeneous

Symbols followed by equal k-long contexts are clustered togheter

[Burrows-Wheeler, 1994]



Shakespeare's Hamlet

**L** is locally homogeneous

Symbols followed by equal k-long contexts are clustered togheter

e.g., k=5

For any length k

# Burrows-Wheeler Transform: why it works

Shakespeare's Hamlet

**L** is locally homogeneous

Symbols followed by equal k-long contexts are clustered togheter

thus, L is highly compressible

# Reverse the BWT

# Burrows-Wheeler Transform: reversibility

**F**                                                **L**

<span style="color:red">**unknown**</span>

| F | | L | |
|---|---|---|---|
| # | abracadabr | a | |
| a | #abracadab | r | |
| a | bra#abraca | d | **?** |
| a | bracadabra | # | **?** |
| a | cadabra#ab | r | |
| a | dabra#abra | c | **?** |
| b | ra#abracad | a | |
| b | racadabra# | a | |
| c | adabra#abr | a | |
| d | abra#abrac | a | |
| r | a#abracada | b | |
| r | acadabra#a | b | |

To reobtain S from L we need to map any symbol in **L** to its corresponding occurrence in **F. LF-Mapping**

Simple! For symbols with just one occurrence.

[Burrows-Wheeler, 1994]

F     L

**unknown**

| # | abracadabr | a |
| a | #abracadab | r |
| a | bra#abraca | d | ? |
| a | bracadabra | # |
| a | cadabra#ab | r |
| a | dabra#abra | c | ? |
| b | ra#abracad | a |
| b | racadabra# | a |
| c | adabra#abr | a |
| d | abra#abrac | a |
| r | a#abracada | b |
| r | acadabra#a | b |

To reobtain S from L we need to map any symbol in **L** to its corresponding occurrence in **F. LF-Mapping**

Simple! For symbols with just one occurrence.

[Burrows-Wheeler, 1994]

**F**

**L**

| | unknown | |
|---|---|---|
| # | abracadabr | a ? |
| a | #abracadab | r |
| a | bra#abraca | d |
| a | bracadabra | # |
| a | cadabra#ab | r |
| a | dabra#abra | c |
| b | ra#abracad | a ? |
| b | racadabra# | a ? |
| c | adabra#abr | a ? |
| d | abra#abrac | a ? |
| r | a#abracada | b |
| r | acadabra#a | b |

To reobtain S from L we need to map any symbol in **L** to its corresponding occurrence in **F. LF-Mapping**

Symbols with more than one occurrence?

# Burrows-Wheeler Transform: reversibility

[Burrows-Wheeler, 1994]

F | | L
:-:|:-:|:-:
| **unknown** |
# | abracadabr | **a** ?
a | #abracadab | r
a | bra#abraca | d
a | bracadabra | #
a | cadabra#ab | r
a | dabra#abra | c
b | ra#abracad | **a** ?
b | racadabra# | **a** ?
c | adabra#abr | **a** ?
d | abra#abrac | **a** ?
r | a#abracada | b
r | acadabra#a | b

To reobtain S from L we need to map any symbol in **L** to its corresponding occurrence in **F. LF-Mapping**

**Key Property:**
Equal symbols in L maintain their relative order in F:
First a in L maps to first a in F,...

# Burrows-Wheeler Transform: reversibility

F                           L

**unknown**

| F | | L |
|---|---|---|
| # | abracadabr | a |
| a | #abracadab | r |
| a | bra#abraca | d |
| a | bracadabra | # |
| a | cadabra#ab | r |
| a | dabra#abra | c |
| b | ra#abracad | a  **?** |
| b | racadabra# | a  **?** |
| c | adabra#abr | a  **?** |
| d | abra#abrac | a  **?** |
| r | a#abracada | b |
| r | acadabra#a | b |

To reobtain S from L we need to map any symbol in **L** to its corresponding occurrence in **F. LF-Mapping**

**Key Property:**
Equal symbols in L maintain their relative order in F:
First a in L maps to first a in F,...

# Burrows-Wheeler Transform: reversibility

**F**  **L**

**unknown**

| F | | L |
|---|---|---|
| # | abracadabr | a |
| a | #abracadab | r |
| a | bra#abraca | d |
| a | bracadabra | # |
| a | cadabra#ab | r |
| a | dabra#abra | c |
| b | ra#abracad | a |
| b | racadabra# | a ? |
| c | adabra#abr | a ? |
| d | abra#abrac | a ? |
| r | a#abracada | b |
| r | acadabra#a | b |

To reobtain S from L we need to map any symbol in **L** to its corresponding occurrence in **F. LF-Mapping**

**Key Property:**
Equal symbols in L maintain their relative order in F:
First *a* in L maps to first *a* in F,...

To reobtain S from L we need to map any symbol in **L** to its corresponding occurrence in **F. LF-Mapping**

**Key Property:**
Equal symbols in L maintain their relative order in F:
First a in L maps to first a in F,...

F                        L

**unknown**

| # | abracadabr | a |
| a | #abracadab | r |
| a | bra#abraca | d |
| a | bracadabra | # |
| a | cadabra#ab | r |
| a | dabra#abra | c |
| b | ra#abracad | a |
| b | racadabra# | a |
| c | adabra#abr | a |
| d | abra#abrac | a |
| r | a#abracada | b |
| r | acadabra#a | b |

To reobtain S from L we need to map any symbol in **L** to its corresponding occurrence in **F. LF-Mapping**

**Key Property:**
Equal symbols in L maintain their relative order in F:
First a in L maps to first a in F,...

Proof

[Burrows-Wheeler, 1994]

F          L

**unknown**

| # | abracadabr | **a** |
| a | #abracadab | r |
| a | bra#abraca | d |
| a | bracadabra | # |
| a | cadabra#ab | r |
| a | dabra#abra | c |
| b | ra#abracad | **a** |
| b | racadabra# | a |
| c | adabra#abr | a |
| d | abra#abrac | a |
| r | a#abracada | b |
| r | acadabra#a | b |

To reobtain S from L we need to map any symbol in **L** to its corresponding occurrence in **F. LF-Mapping**

**Key Property:**
Equal symbols in L maintain their relative order in F:
First a in L maps to first a in F,...

Proof

➤ #abracadabra

➤ bra#abracada

[Burrows-Wheeler, 1994]

**F**

**L**

**unknown**

| F | | L |
|---|---|---|
| # | abracadabr | a |
| a | #abracadab | r |
| a | bra#abraca | d |
| a | bracadabra | # |
| a | cadabra#ab | r |
| a | dabra#abra | c |
| b | ra#abracad | a |
| b | racadabra# | a |
| c | adabra#abr | a |
| d | abra#abrac | a |
| r | a#abracada | b |
| r | acadabra#a | b |

To reobtain S from L we need to map any symbol in **L** to its corresponding occurrence in **F. LF-Mapping**

**Key Property:**
Equal symbols in L maintain their relative order in F:
First $a$ in L maps to first $a$ in F,...

Proof

#abracadabra

$<_{lex}$

bra#abracada

[Burrows-Wheeler, 1994]

**F**

**L**

**unknown**

| F | | L |
|---|---|---|
| # | abracadabr | a |
| a | #abracadab | r |
| a | bra#abraca | d |
| a | bracadabra | # |
| a | cadabra#ab | r |
| a | dabra#abra | c |
| b | ra#abracad | a |
| b | racadabra# | a |
| c | adabra#abr | a |
| d | abra#abrac | a |
| r | a#abracada | b |
| r | acadabra#a | b |

To reobtain S from L we need to map any symbol in **L** to its corresponding occurrence in **F. LF-Mapping**

**Key Property:**
Equal symbols in L maintain their relative order in F:
First a in L maps to first a in F,...

Rotate them rightward to obtain their cyclic-rot

Proof

→ #abracadabra

**<lex**

→ bra#abracada

# Burrows-Wheeler Transform: reversibility

[Burrows-Wheeler, 1994]

**F**

**unknown**

**L**

| F | | L |
|---|---|---|
| # | abracadabr | a |
| a | #abracadab | r |
| a | bra#abraca | d |
| a | bracadabra | # |
| a | cadabra#ab | r |
| a | dabra#abra | c |
| b | ra#abracad | a |
| b | racadabra# | a |
| c | adabra#abr | a |
| d | abra#abrac | a |
| r | a#abracada | b |
| r | acadabra#a | b |

To reobtain S from L we need to map any symbol in **L** to its corresponding occurrence in **F. LF-Mapping**

> **Key Property:**
> Equal symbols in L maintain their relative order in F:
> First a in L maps to first a in F,...

> Rotate them rightward i.e., move the a in front

Proof

#abracadabra      a#abracadabr

$<_{lex}$

bra#abracada      abra#abracad

[Burrows-Wheeler, 1994]

F

L

unknown

| # | abracadabr | a |
| a | #abracadab | r |
| a | bra#abraca | d |
| a | bracadabra | # |
| a | cadabra#ab | r |
| a | dabra#abra | c |
| b | ra#abracad | a |
| b | racadabra# | a |
| c | adabra#abr | a |
| d | abra#abrac | a |
| r | a#abracada | b |
| r | acadabra#a | b |

To reobtain S from L we need to map any symbol in **L** to its corresponding occurrence in **F. LF-Mapping**

**Key Property:**
Equal symbols in L maintain their relative order in F:
First a in L maps to first a in F,...

Proof

#abracadabra

$<_{lex}$

bra#abracada

⇒

a#abracadabr

$<_{lex}$

abra#abracad

# Burrows-Wheeler Transform: reversibility

[Burrows-Wheeler, 1994]

F        L

**unknown**

| F | | L |
|---|---|---|
| # | abracadabr | a |
| a | #abracadab | r |
| a | bra#abraca | d |
| a | bracadabra | # |
| a | cadabra#ab | r |
| a | dabra#abra | c |
| b | ra#abracad | a |
| b | racadabra# | a |
| c | adabra#abr | a |
| d | abra#abrac | a |
| r | a#abracada | b |
| r | acadabra#a | b |

To reobtain S from L we need to map any symbol in **L** to its corresponding occurrence in **F. LF-Mapping**

**Key Property:**
Equal symbols in L maintain their relative order in F:
First a in L maps to first a in F,...

Proof

#abracadabra          a#abracadabr

$<_{lex}$                $<_{lex}$

bra#abracada         abra#abracad

[Burrows-Wheeler, 1994]

F

L

**unknown**

| F | L |
|---|---|
| # | a |
| a | r |
| a | d |
| a | # |
| a | r |
| a | c |
| b | a |
| b | a |
| c | a |
| d | b |
| r | b |
| r | |

We reconstruct S backward

**Two Key properties:**

1) LF maps L's to F's symbols

2) L[i] precedes F[i] in S

S=  ------------#

[Burrows-Wheeler, 1994]

F

L

**unknown**

| F | | L |
|---|---|---|
| # | | a |
| a | | r |
| a | | d |
| a | | **#** |
| a | | r |
| a | | c |
| b | | a |
| b | | a |
| c | | a |
| d | | a |
| r | | b |
| r | | b |

We reconstruct S backward

**Two Key properties:**

1) LF maps L's to F's symbols

2) L[i] precedes F[i] in S

$$S = \text{-----------\#}$$

It precedes # in S

F

unknown

a
r
d
#
r
c
a
a
a
a
b
b

#
a
a
a
a
a
b
b
c
d
r
r

We reconstruct S backward

**Two Key properties:**

1) LF maps L's to F's symbols

2) L[i] precedes F[i] in S

S=  - - - - - - - - - - - a#

F

L

**unknown**

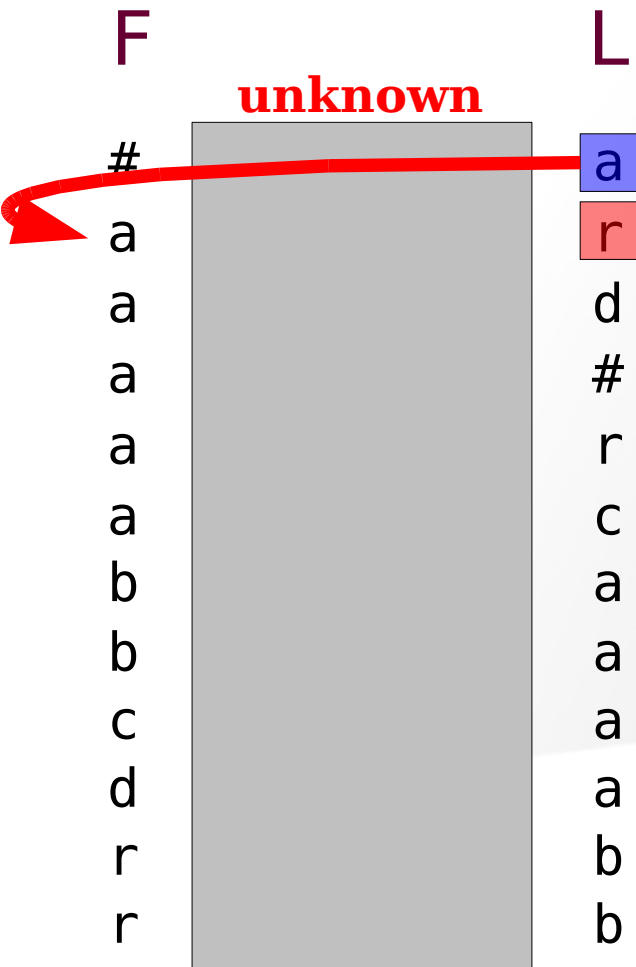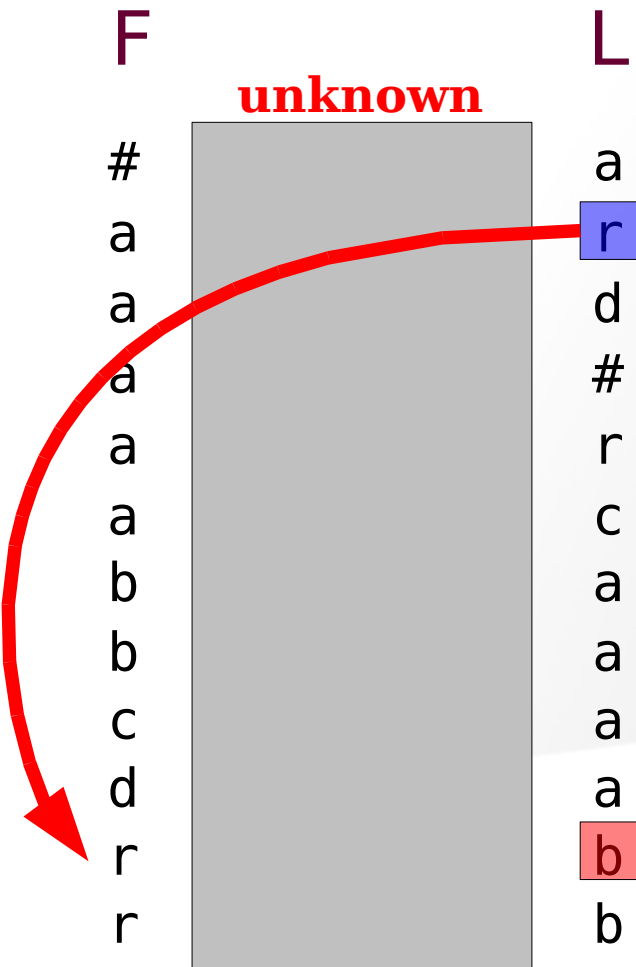| F | | L |
|---|---|---|
| # | | a |
| a | | r |
| a | | d |
| a | | # |
| a | | r |
| a | | c |
| a | | a |
| b | | a |
| b | | a |
| c | | a |
| d | | a |
| r | | b |
| r | | b |

We reconstruct S backward

**Two Key properties:**

1) LF maps L's to F's symbols

2) L[i] precedes F[i] in S

S=  - - - - - - - - - - -**a#**

# Burrows-Wheeler Transform: reversibility

F

L

**unknown**

| F | L |
|---|---|
| # | a |
| a | r |
| a | d |
| a | # |
| a | r |
| a | c |
| b | a |
| b | a |
| c | a |
| d | a |
| r | b |
| r | b |

We reconstruct S backward

**Two Key properties:**

1) LF maps L's to F's symbols

2) L[i] precedes F[i] in S

S= ----------ra#

F          L

**unknown**

| F | L |
|---|---|
| # | a |
| a | r |
| a | d |
| a | # |
| a | r |
| a | c |
| b | a |
| b | a |
| c | a |
| d | a |
| r | b |
| r | b |

We reconstruct S backward

**Two Key properties:**

1) LF maps L's to F's symbols

2) L[i] precedes F[i] in S

S=  --------**bra#**

[Burrows-Wheeler, 1994]

F

unknown

L

# a a a a a a b b c d r r

L a r d # r c a a a b b

We reconstruct S backward

**Two Key properties:**

1) LF maps L's to F's symbols

2) L[i] precedes F[i] in S

S= --------**abra#**

[Burrows-Wheeler, 1994]

F  L

**unknown**

| F | L |
|---|---|
| # | a |
| a | r |
| a | d |
| a | # |
| a | r |
| a | c |
| b | a |
| b | a |
| c | a |
| d | a |
| r | b |
| r | b |

We reconstruct S backward

**Two Key properties:**

1) LF maps L's to F's symbols

2) L[i] precedes F[i] in S

S= **abracadabra#**

F          L

**unknown**

| F | | L |
|---|---|---|
| # | | a |
| a | | r |
| a | | d |
| a | | # |
| a | | r |
| a | | c |
| b | | a |
| b | | a |
| c | | a |
| d | | a |
| r | | b |
| r | | b |

We reconstruct S backward

InvertBWT(L)

Compute LF[0,n-1];
r = 0; i = n;
while (i>0) {
    T[i] = L[r];
    r = LF[r]; i--;
    }

**ies:**

s symbols

in S

S= **abracadabra#**

# Burrows-Wheeler Transform: reversibility

Starting from L, the LF-mapping can be computed in linear time using a simple algorithm

**F**  **unknown**  **L  LF**

| F | unknown | L | LF |
|---|---|---|---|
| # | abracadabr | a | 1 |
| a | #abracadab | r | 10 |
| a | bra#abraca | d | 9 |
| a | bracadabra | # | 0 |
| a | cadabra#ab | r | 11 |
| a | dabra#abra | c | 8 |
| b | ra#abracad | a | 2 |
| b | racadabra# | a | 3 |
| c | adabra#abr | a | 4 |
| d | abra#abrac | a | 5 |
| r | a#abracada | b | 6 |
| r | acadabra#a | b | 7 |

**C**

| # | a | b | c | d | r |
|---|---|---|---|---|---|
| 1 | 6 | 8 | 9 | 10 | 12 |

Auxiliary vector C that stores for each symbol c the number of occs in L of symbols smaller than c

- Scan L
  - Set LF[i]=C[L[i]]
  - Update the counter C[L[i]]++

# Compress the BWT

# Burrows-Wheeler Transform: why it works

Shakespeare's Hamlet

**L** is locally homogeneous

Symbols followed by equal contexts are clustered togheter

thus, L is highly compressible

# Burrows-Wheeler Transform: why it works

Shakespeare's Hamlet

**L** is locally homogeneous

Symbols followed by equal contexts are clustered togheter

thus, L is highly compressible

B&W suggest to compress L in 3 steps:

- Move-To-Front
- Run Length Encoding
- Statistical Encoder (e.g. Huffman)

# Move-to-Front

Transform L which is locally **homogeneous** into a new string that is globally **homogeneous**

L = a r d r c a a a a b b

C =

# Move-to-Front

Transform L which is locally homogeneous into a new string that is globally homogeneous

L = a r d r c a a a a b b

C =

**List**

| |
| --- |
| a |
| b |
| c |
| d |
| r |

Symbols in the alphabet. Initially, they are sorted lexicographically

Transform L which is locally homogeneous into a new string that is globally homogeneous

L = a r d r c a a a a b b

C =

**List**

| |
|:-:|
| a |
| b |
| c |
| d |
| r |

Transform L which is locally homogeneous into a new string that is globally homogeneous

**List**

L = a r d r c a a a a b b

| a |
| b |
| c |
| d |
| r |

Move a in front of List

C = 0
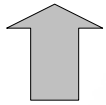
# Move-to-Front

Transform L which is locally homogeneous into
a new string that is globally homogeneous

L = a r d r c a a a a b b

C = 0 4

**List**

| |
|---|
| a |
| b |
| c |
| d |
| r |

Transform L which is locally homogeneous into a new string that is globally homogeneous

L = a r d r c a a a a b b

C = 0 4

**List**

| |
|---|
| a |
| b |
| c |
| d |
| r |

# Move-to-Front

Transform L which is locally homogeneous into a new string that is globally homogeneous

L = a r d r c a a a a b b

C = 0  4

**List**

| |
|:---:|
| r |
| a |
| b |
| c |
| d |

# Move-to-Front

Transform L which is locally homogeneous into a new string that is globally homogeneous

L = a r d r c a a a a b b

C = 0 4 4 1 4 3 0 0 0 4 0
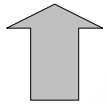
**List**

| |
|---|
| b |
| a |
| c |
| r |
| d |

# Move-to-Front

Transform L which is locally homogeneous into a new string that is globally homogeneous

L = a r d r c a a a a b b

C = 0 4 4 1 4 3 0 0 0 4 0

**List**

| |
|---|
| b |
| a |
| c |
| r |
| d |

Number of distinct symbols since the previous occurrence of a

# Move-to-Front

Transform L which is locally homogeneous into a new string that is globally homogeneous

**List**

L = a r d r c a a a a b b

| b |
|---|
| a |
| c |
| r |
| d |

C = 0 4 4 1 4 3 0 0 0 4 0

If the string is locally homogeneous,
C contain a lot of small numbers.

# Move-to-Front

Transform L which is locally **homogeneous** into a new string that is globally **homogeneous**

**List**

| |
|---|
| b |
| a |
| c |
| r |
| d |

L = a  r  d  r  c  a  a  a  a  b  b

C = 0  4  4  1  4  3  0  0  0  4  0

In particular, runs of equal symbols become runs of 0

# Burrows-Wheeler Transform: why Move-to-Front

```
F                                    ... L
ot look upon his like again.   ... n
ot look upon me; Lest with th ... n
ot love on the wing,-- As I p ... h
ot love your father; But that ... n
ot made them well, they imita ... n
ot madness That I have utter' ... n
ot me? Ham. No, by the rood,   ... g
ot me; no, nor woman neither, ... n
ot me'? Ros. To think, my lor ... n
ot mend his pace with beating ... n
ot mine. Ham. No, nor mine no ... n
ot mine own. Besides, to be d ... n
ot mock me, fellow-student. I ... n
ot monstrous that this player ... n
ot more like. Ham. But where   ... n
ot more, my lord. Ham. Is not ... j
ot more native to the heart,   ... n
ot more ugly to the thing tha ... n
ot move thus. Oph. You must s ... n
ot much approve me.--Well, si ... n
```

Shakespeare's Hamlet
**AGAIN!**

$F$ ... $L$ | **MTF(L)**
---|---

```
ot look upon his like again.   ... n
ot look upon me; Lest with th ... n
ot love on the wing,-- As I p ... h
ot love your father; But that ... n
ot made them well, they imita ... n
ot madness That I have utter' ... n
ot me? Ham. No, by the rood,  ... g
ot me; no, nor woman neither, ... n
ot me'? Ros. To think, my lor ... n
ot mend his pace with beating ... n
ot mine. Ham. No, nor mine no ... n
ot mine own. Besides, to be d ... n
ot mock me, fellow-student. I ... n
ot monstrous that this player ... n
ot more like. Ham. But where  ... n
ot more, my lord. Ham. Is not ... j
ot more native to the heart,  ... n
ot more ugly to the thing tha ... n
ot move thus. Oph. You must s ... n
ot much approve me.--Well, si ... n
```

MTF(L):
+
0
+
1
0
0
+
1
0
0
0
0
0
0
0
+
1
0
0
0

Shakespeare's Hamlet
**AGAIN!**

# After MTF

- Now we have a string with small numbers: lots of 0s, many 1s, ...
- Skewed frequencies: next steps RLE +Huffman!

Symbol frequencies

# Run-Length-Encoding

Encode runs of 0s in C

L = a  r  d  r  c  a  a  a  a  b  b

C = 0  4  4  1  4  3  0  0  0  4  0

C'= 0  5  5  2  5  4  0  0  0  5  0

Replace symbol c>0 with c+1,
i.e. increase the alphabet size by 1

# Run-Length-Encoding

Encode runs of 0s in C

L = a r d r c a a a a b b

C = 0 4 4 1 4 3 0 0 0 4 0

C'= **0** 5 5 2 5 4 **0 0** 5 **0**

Encode runs of 0s with bin(length+1) without the most significant bit. (0/1 are unused)

# Run-Length-Encoding

Encode runs of 0s in C

L = a  r  d  r  c  a  a  a  a  b  b

C = 0  4  4  1  4  3  [0  0  0]  4  0

bin(3+1)= 100

C'= **0**  5  5  2  5  4  [**0  0**]  5  **0**

Encode runs of 0s with bin(length+1)
without the most significant bit. (0/1 are unused)

# Run-Length-Encoding

Encode runs of 0s in C

L = a r d r c a a a a b b

C = 0 4 4 1 4 3 0 0 0 4 0

C' = **0** 5 5 2 5 4 **0 0** 5 **0**

Final step: Huffman to encode C'.
Now it is effective, large number of 0s and 1s

# How to build (efficiently) the BWT?

# How to build the BWT

Let us given S = abracadabra#

|   | F |   | L |
|---|---|---|---|
| abracadabra# | # | abracadabr | a |
| bracadabra#a | a | #abracadab | r |
| racadabra#ab | a | bra#abraca | d |
| acadabra#abr | a | bracadabra | # |
| cadabra#abra | a | cadabra#ab | r |
| adabra#abrac | a | dabra#abra | c |
| dabra#abraca | b | ra#abracad | a |
| abra#abracad | b | racadabra# | a |
| bra#abracada | c | adabra#abr | a |
| ra#abracadab | d | abra#abrac | a |
| a#abracadabr | r | a#abracada | b |
| #abracadabra | r | acadabra#a | b |

**Sort the rows** →

↔ **S**

Let us given S = abracadabra#

F                                L

abracadabra#          #   abracadabr
bracadabra#a          a   #abracadab      r
racadabra#ab          a   bra#abraca      d
acadabra#             b       bra          #
cadabr                        #ab          r
adabra                        bra          c
dabra#                        cad          a
abra#a                        ra#          a
bra#abracada          c    dabra#abr       a
ra#abracadab          d   abra#abrac       a
a#abracadabr          r   a#abracada       b
#abracadabra          r   acadabra#a       b

Obvious inefficiency:
O($n^2$) space

↔ S

# Suffix Array

Let us given S = abracadabra#

**SA**

| 12 | # |
|----|---|
|    |   |
|    |   |
|    |   |
|    |   |
|    |   |
|    |   |
|    |   |
|    |   |
|    |   |
|    |   |
|    |   |

SA stores the starting positions in S
of the suffixes sorted lexicographically

# Suffix Array

Let us given S = abracadabra#

SA

| 12 | # |
|----|---|
|    |   |
|    |   |
|    |   |
|    |   |
|    |   |
|    |   |
|    |   |
|    |   |
|    |   |
|    |   |
|    |   |

SA stores the starting positions in S
of the suffixes sorted lexicographically

# Suffix Array

Let us given S = abracadabra#

## SA

| |
|---|
| 12 | #
| 11 | a#
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |

SA stores the starting positions in S
of the suffixes sorted lexicographically

# Suffix Array

Let us given S = abracadabra#

SA

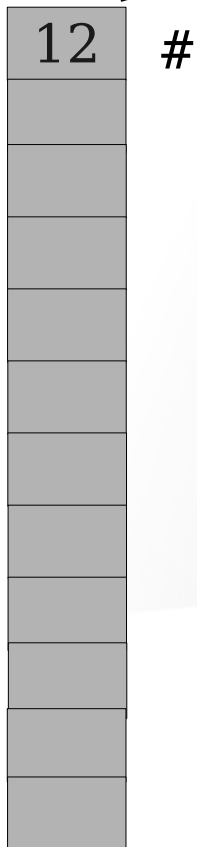| | |
|---|---|
| 12 | # |
| 11 | a# |
| 8 | abra# |
| 1 | abracadabra# |
| 4 | acadabra# |
| 6 | adabra# |
| 9 | bra# |
| 2 | bracadabra# |
| 5 | cadabra#abr |
| 7 | dabra# |
| 10 | ra# |
| 3 | racadabra# |

SA stores the starting positions in S
of the suffixes sorted lexicographically

# Suffix Array

Let us given S = abracadabra#

SA

| | |
|---|---|
| 12 | # |
| 11 | a# |
| 8 | abra# |
| 1 | abracadabra# |
| 4 | acadabra# |
| 6 | adabra# |
| 9 | bra# |
| 2 | bracadabra# |
| 5 | cadabra#abr |
| 7 | dabra# |
| 10 | ra# |
| 3 | racadabra# |

| F | L |
|---|---|
| #abracadabr | a |
| a#abracadab | r |
| abra#abraca | d |
| abracadabra | # |
| acadabra#ab | r |
| adabra#abra | c |
| bra#abracad | a |
| bracadabra# | a |
| cadabra#abr | a |
| dabra#abrac | a |
| ra#abracada | b |
| racadabra#a | b |

# Suffix Array

Let us given S = abracadabra#

SA

| 12 | # |
| 11 | a# |
| 8 | abra# |
| 1 | abracadabra# |
| 4 | acadabra# |
| 6 | adabra# |
| 9 | bra# |
| 2 | bracadabra# |
| 5 | cadabra#abr |
| 7 | dabra# |
| 10 | ra# |
| 3 | racadabra# |

⬅➡

| F | L |
|---|---|
| #abracadabr | a |
| a#abracadab | r |
| abra#abraca | d |
| abracadabra | # |
| acadabra#ab | r |
| adabra#abra | c |
| bra#abracad | a |
| bracadabra# | a |
| cadabra#abr | a |
| dabra#abrac | a |
| ra#abracada | b |
| racadabra#a | b |

# Suffix Array

Let us given S = abracadabra#

SA

| 12 | # | | #abracadabr | a |
| 11 | a# | Given SA, we have L[i] = S[SA[i]-1] racadab | r |
| 8 | abra# | abra#abraca | d |
| 1 | abracadabra# | abracadabra | # |
| 4 | acadabra# | acadabra#ab | r |
| 6 | adabra# | adabra#abra | c |
| 9 | bra# | bra#abracad | a |
| 2 | bracadabra# | bracadabra# | a |
| 5 | cadabra#abr | cadabra#abr | a |
| 7 | dabra# | dabra#abrac | a |
| 10 | ra# | ra#abracada | b |
| 3 | racadabra# | racadabra#a | b |

F          L

# Suffix Array

Let us given S = abracadabra#

SA                     F           L

| SA | | F | L |
|----|-----------|-----------------|---|
| 12 | # | #abracadabr | a |
| 11 | a# | ...racadab | r |
| 8 | abra# | ...#abraca | d |
| 1 | abrac... | ...racadabra | # |
| 4 | acada... | ...adabra#ab | r |
| 6 | adabr... | ...abra#abra | c |
| 9 | bra# | ...a#abracad | a |
| 2 | bracadabra# | bracadabra# | a |
| 5 | cadabra#abr | cadabra#abr | a |
| 7 | dabra# | dabra#abrac | a |
| 10 | ra# | ra#abracada | b |
| 3 | racadabra# | racadabra#a | b |

Given SA, we have L[i] = S[SA[i]-1]

However, the Suffix Array has many other applications in pattern matching!

# How construct SA?

SA

| | |
|---|---|
| 12 | # |
| 11 | a# |
| 8 | abra# |
| 1 | abracadabra# |
| 4 | acadabra# |
| 6 | adabra# |
| 9 | bra# |
| 2 | bracadabra# |
| 5 | cadabra#abr |
| 7 | dabra# |
| 10 | ra# |
| 3 | racadabra# |

Elegant but inefficient

COMPARISON_BASED_CONSTRUCTION(char *$T$, int $n$, char **$SA$)

$\{$ for($i = 0$; $i < n$; $i$++)    $SA[i] = T + i$;

    QSORT($SA$, $n$, sizeof(char *), Suffix_cmp); $\}$

SUFFIX_CMP(char **p, char **q){ return strcmp(*p,*q); }

- $\Theta(n^2 \log n)$ time in the worst-case
- $O(n)$ space

# Suffix Array

Last years, many clever algorithms that compute the Suffix Array in <span style="color:red">linear time:</span>

- Karkkainen-Sanders [J. ACM, 2006]
- Ko-Aluru [J. Discrete Algorithms, 2005

# Suffix Array

Last years, many clever algorithms that compute the Suffix Array in <span style="color:red">linear time:</span>

- Karkkainen-Sanders [J. ACM, 2006]
- Ko-Aluru [J. Discrete Algorithms, 2005

and practical ones (no linear time):

- Manzini-Ferragina [Algorithmica, 2004]
- Maniscalco-Puglisi [ACM JEA, 2006]

# Suffix Array

Last years, many clever algorithms that compute the Suffix Array in <span style="color:red">linear time</span>:

- Karkkai
- Ko-Aluru

and pract

- **Manzini-Ferragina** [Algorithmica, 2004]
- **Maniscalco-Puglisi** [ACM JEA, 2006]

> BWT can be computed
> in linear time!!!
>
> Wheeler discovered
> the BWT in 1983
> but he did not publish it because
> it was considered too slow!

# Suffix Array

Last years, many clever algorithms that compute the Suffix Array in <span style="color:red">linear time:</span>

- Karkkai
- Ko-Aluru

and pract

- **Manzini-Ferragina** [Algorithmica, 2004]
- **Maniscalco-Puglisi** [ACM JEA, 2006]

BWT can be computed
in linear time!!!

Wheeler discovered
the BWT in 1983
but he did not publish it because
it was considered too slow!

In 1994, the
paper was
rejected... and
BWT  is still a
Technical Report.

# A compressor based on BWT: Bzip2

# You find this in your Linux distribution

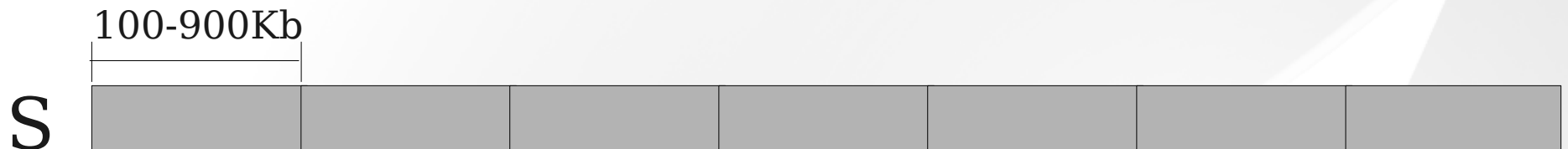It doesn't build the BWT of the whole text but uses blocking approach.

100-900Kb

S

# A real implementation: bzip2

It doesn't build the BWT of the whole text but uses blocking approach.

100-900Kb

S

Just for time performance, the compression it achieves is worse.

# bzip2 vs gzip

| English 5Mb | comp. size | comp. time |
|:---:|:---:|:---:|
| gzip | 2.0 Mb | 1.7 secs |
| bzip2 | 1.5 Mb | 2.2 secs |
| bigbzip | 1.3 Mb | 2.4 secs |

# bzip2 vs gzip

| English 5Mb | comp. size | comp. time |
|---|---|---|
| gzip | 2.0 Mb | 1.7 secs |
| bzip2 | 1.5 Mb | 2.2 secs |
| bigbzip | 1.3 Mb | 2.4 secs |

| English 20Mb | comp. size | comp. time | dec. time |
|---|---|---|---|
| gzip | 7.8 Mb | 7.2 secs | 0.8 secs |
| bzip2 | 5.9 Mb | 11.0 secs | 4.0 secs |
| bigbzip | 4.1 Mb | 20.0 secs | 5.8 secs |

# bzip2 vs gzip

| English 5Mb | comp. size | comp. time |
|---|---|---|
| gzip | 2.0 Mb | 1.7 secs |
| bzip2 | 1.5 Mb | 2.2 secs |
| bigbzip | 1.3 Mb | 2.4 secs |

| English 20Mb | | | |
|---|---|---|---|
| gzip | | | |
| bzip2 | 5.9 Mb | 11.0 secs | 4.0 secs |
| bigbzip | 4.1 Mb | 20.0 secs | 5.8 secs |

... but, now we can increase the block size from 900Kb to 5Mb due to new algorithmic solutions.
Improve compression, same time performace!

# Some other applications

- ***FM-index/CSA:*** *Searching in compressed text.*

  – Given a text S, we can compress in a way that permits us to search any pattern P in S in time proportional to |P|!
  (Notice that |P|<<|S|)

# Some other applications

- **FM-index/CSA***: Searching in compressed text.*

  - Given a text S, we can compress in a way that permits us to search any pattern P in S in time proportional to |P|!
  (Notice that |P|<<|S|)

- ***Booster:*** *Best theoretical class of compressors.*

  - Given an 0-th order compressor C, we can optimally partition BWT(S) in O(|S|) time in order to achieve the best compression with C.

# Some other applications

- **FM-index/CSA:** *Searching in compressed text.*

  - Given a text S, we can compress in a way that permits us to search any pattern P in S in time proportional to |P|!
    (Notice that |P|<<|S|)

- ***Booster:*** *Best theoretical class of compressors.*

  - Given an 0-th order compressor C, we can optimally partition BWT(S) in O(|S|) time in order to achieve the best compression with C.

- Idea used for other data types.

  - e.g. XML. Compression of labeled trees

# Some other applications

- **FM-index/CSA*: *Searching in compressed text.*
  - Given a text S, we can compress in a way that permits us to search any pattern P in S in time proportional to |P|!

  These stuff require another seminar...

  - •

  - Given an 0-th order compressor C, we can optimally partition BWT(S) in O(|S|) time in order to achieve the best compression with C.
- Idea used for other data types.
  - e.g. XML. Compression of labeled trees

# Some other applications

- **FM-index/CSA:** *Searching in compressed text.*
  - Given a text S, we can compress in a way that permits us to search any pattern P in S in time proportional to |P|!

These stuff require another seminar...
This is a threat!

  - Given an 0-th order compressor C, we can optimally partition BWT(S) in O(|S|) time in order to achieve the best compression with C.

- Idea used for other data types.
  - e.g. XML. Compression of labeled trees