

# On Compression and Indexing: two sides of the same coin

Paolo Ferragina

Dipartimento di Informatica, Università di Pisa

Paolo Ferragina, Università di Pisa

## What do we mean by "Indexing" ?

### Types of data

Linguistic or tokenizable text  
Raw sequence of characters or bytes

DNA sequences  
Audio-video files  
Executables

### Types of query

Word-based query  
Character-based query

Arbitrary substring  
Complex match

### Two indexing approaches :

- Word-based indexes, here a notion of "word" must be devised !
  - » Inverted files, Signature files, Bitmaps.
- Full-text indexes, no constraint on text and queries !
  - » Suffix Array, Suffix tree, String B-tree [Ferragina-Grossi, JACM 99].

Paolo Ferragina, Università di Pisa

## What do we mean by "Compression" ?

☺ Compression has two positive effects:

- ✓ Space saving (or, double memory at the same cost)
- ✓ Performance improvement
  - Better use of memory levels close to processor
  - Increased disk and memory bandwidth
  - Reduced (mechanical) seek time

» CPU speed nowadays makes (de)compression "costless" !!



### Moral:

More economical to store data in compressed form than uncompressed

- ☑ From March 2001 the Memory eXpansion Technology (MXT) is available on IBM eServers x330MXT
  - ✓ Same performance of a PC with double memory but at half cost

Paolo Ferragina, Università di Pisa

## Compression and Indexing: Two sides of the same coin !

🔍 Do we witness a paradoxical situation ?

- ✓ An index injects redundant data, in order to speed up the pattern searches
- ✓ Compression removes redundancy, in order to squeeze the space occupancy

☺ NO, new results proved a mutual reinforcement behaviour !

- ✓ Better indexes can be designed by exploiting compression techniques
- ✓ Better compressors can be designed by exploiting indexing techniques

In terms of space occupancy

Also in terms of compression ratio

### Moral:

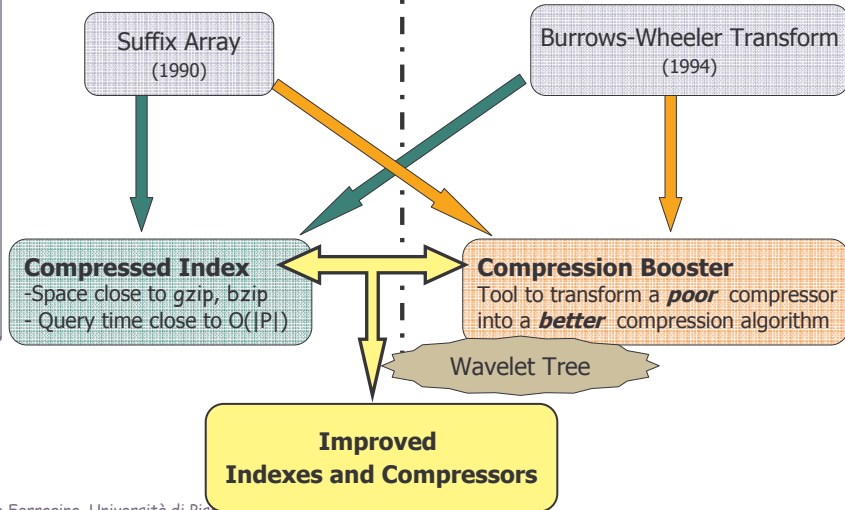
CPM researchers must have a multidisciplinary background, ranging from Data structure design to Data compression, from Architectural Knowledge to Database principles, till Algorithmic Engineering and more...

Paolo Ferragina, Università di Pisa

## Our journey, today...

Index design (Weiner '73)

Compressor design (Shannon '48)

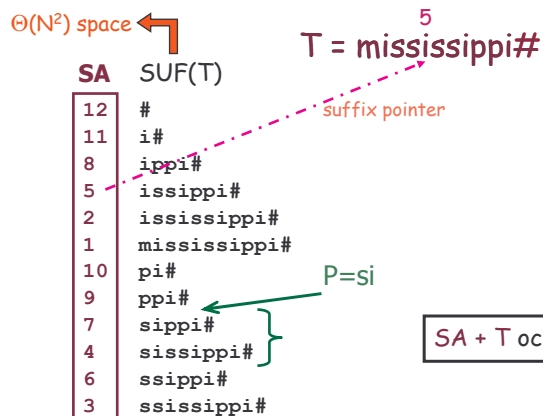


Paolo Ferragina, Università di Pisa

## The Suffix Array [BaezaYates-Gonnet, 87 and Manber-Myers, 90]

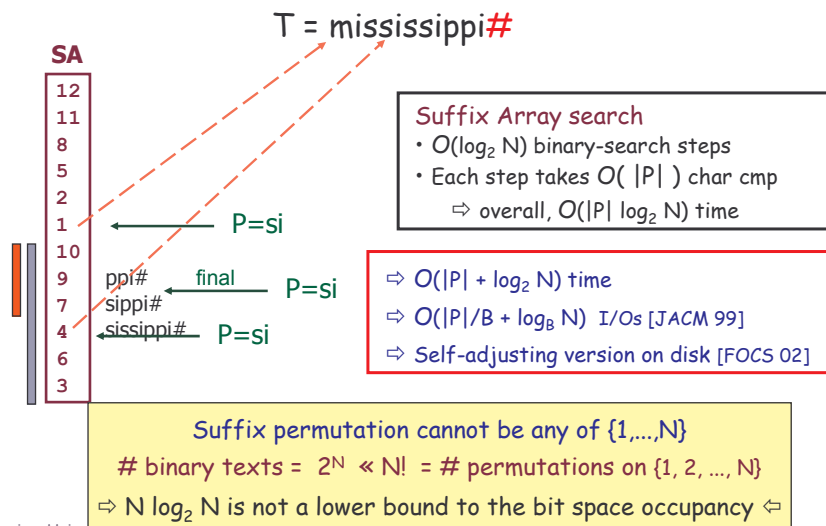
Prop 1. All suffixes in  $SUF(T)$  having prefix  $P$  are contiguous

Prop 2. These suffixes follow  $P$ 's lexicographic position



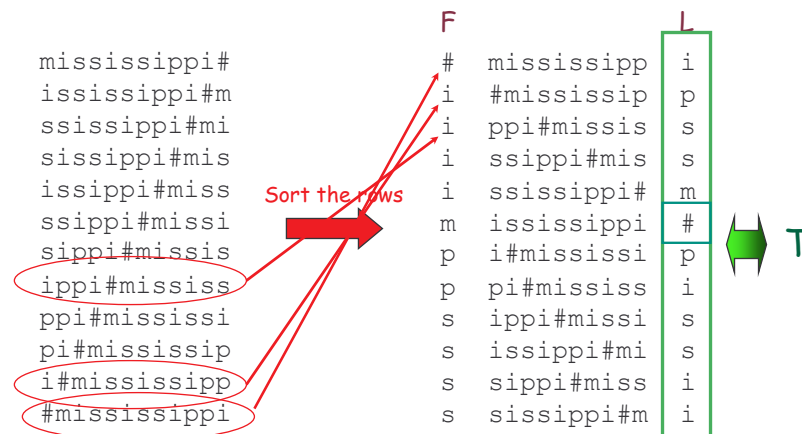
Paolo Ferragina, Università di Pisa

## Searching in Suffix Array [Manber-Myers, 90]



## The Burrows-Wheeler Transform (1994)

Let us given a text  $T = \text{mississippi}\#$



## Why L is so interesting for compression ?

| F | unknown    | L |
|---|------------|---|
| # | mississipp | i |
| i | #mississip | p |
| i | ppi#missis | s |
| i | ssippi#mis | s |
| i | ssissippi# | m |
| m | ississippi | # |
| p | i#mississi | p |
| p | pi#mississ | i |
| s | ippi#missi | s |
| s | issippi#mi | s |
| s | sippi#miss | i |
| s | sissippi#m | i |

A key observation:

- L is locally homogeneous

➔ L is highly compressible

Algorithm Bzip :

- 1 Move-to-Front coding of L
- 2 Run-Length coding
- 3 Statistical coder: Arithmetic, Huffman

T = mississippi#

☑ Bzip vs. Gzip: 20% vs. 33%, but it is s

Paolo Ferragina, Università di Pisa

Building the BWT ⇔ SA construction  
Inverting the BWT ⇔ array visit  
...overall  $O(N)$  time...

## Suffix Array vs. BW-transform

| SA | Rotated text | L |
|----|--------------|---|
| 12 | #mississipp  | i |
| 11 | i#mississip  | p |
| 8  | ippi#missis  | s |
| 5  | issippi#mis  | s |
| 2  | ississippi#  | m |
| 1  | mississippi  | # |
| 10 | pi#mississi  | p |
| 9  | ppi#mississ  | i |
| 7  | sippi#missi  | s |
| 4  | sissippi#mi  | s |
| 6  | ssippi#miss  | i |
| 3  | ssissippi#m  | i |

L includes SA and T. Can we search within L ?



Paolo Ferragina, Università di Pisa

# A compressed index

[Ferragina-Manzini, IEEE Focs 2000]

## Bridging data-structure design and compression techniques:

- ✓ Suffix array data structure
- ✓ Burrows-Wheeler Transform

## The theoretical result:

- ✓ Query complexity:  $O(p + occ \log^\epsilon N)$  time
  - ✓ Space occupancy:  $O(N H_k(T)) + o(N)$  bits  $\rightarrow o(N)$  if  $T$  compressible
- ↪  $k$ -th order empirical entropy

Index does not depend on  $k$   
Bound holds for all  $k$ , simultaneously

## The corollary is that:

- ✓ The Suffix Array is compressible
- ✓ It is a self-index

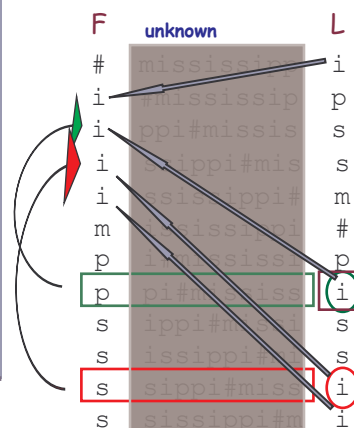
$O(n)$  space: A plethora of papers  
 $H_k$ : Grossi-Gupta-Vitter (03), Sadakane (02),...  
Now, more than 20 papers with more than 20 authors on related subjects

## In practice, the index is much appealing:

- ✓ Space close to the best known compressors, ie. bzip
- ✓ Query time of few millisecs on hundreds of MBs

Paolo Ferragina, Università di Pisa

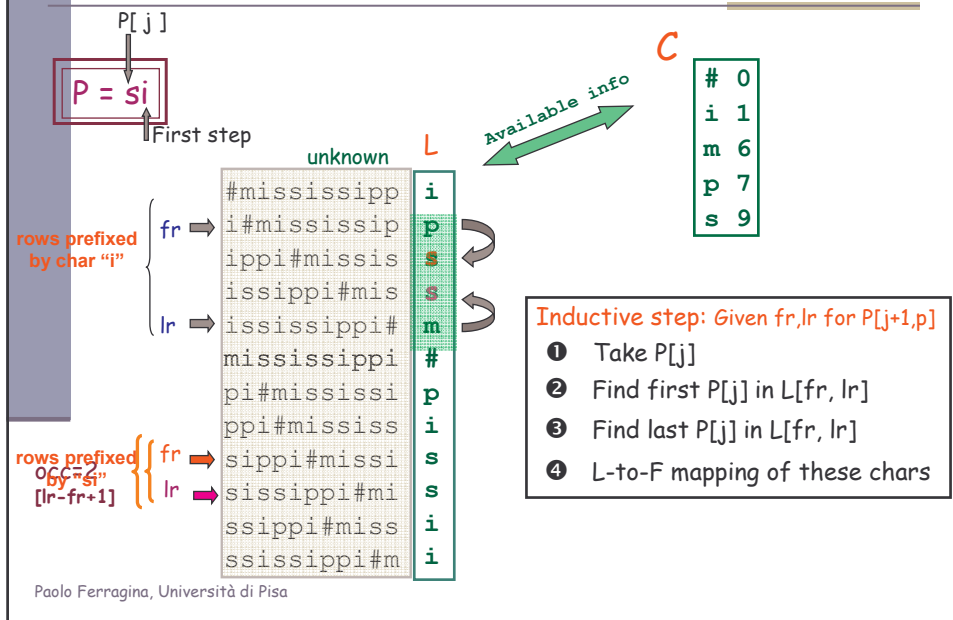
# A useful tool: $L \rightarrow F$ mapping



How do we map L's onto F's chars?  
... Need to distinguish equal chars...

Paolo Ferragina, Università di Pisa

## Substring search in T (Count occurrences)



## Many details are missing...

- ☺ The column L is actually kept compressed:
  - Still guarantee  $O(p)$  time to count the P's occurrences
- ☺ The **Locate** operation takes  $O(\log^E N)$  time
  - Some additional data structure, in  $o(n)$  bits
- ☺ Efficient and succinct index construction [Hon et al., Focs 03]
- ☺ Bio-application: fit Human-genome index in a PC [Sadakane et al., 02]

### Interesting issues:

- ☑ What about arbitrary alphabets? [Grossi et al., 03; Ferragina et al., 04]
- ☑ What about disk-aware, or cache-oblivious, or self-adjusting versions?
- ☑ What about challenging applications: bio, db, data mining, handheld PC, ...

## Where we are ...

We investigated the reinforcement relation:

Compression ideas ↔ Index design

Let's now turn to the other direction

Indexing ideas ↔ Compressor design

Booster

Paolo Ferragina, Università di Pisa

## What do we mean by "boosting" ?

It is a technique that takes a *poor* compressor  $A$  and turns it into a compressor with *better performance guarantee*

A *memoryless* compressor is *poor* in that it assigns codewords to symbols according only to their frequencies (e.g. Huffman)

It incurs in some obvious limitations:

$$T = a^n b^n$$

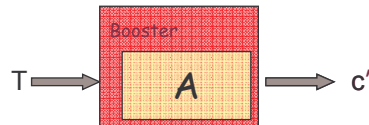
$T =$  random string of length  $2n$  and same number of 'a,b'

Paolo Ferragina, Università di Pisa



## What do we mean by “boosting” ?

It is a technique that takes a *poor* compressor **A** and turns it into a compressor with *better performance guarantee*



Qualitatively, we would like to achieve

- $c'$  is shorter than  $c$ , if  $T$  is compressible
- $\text{Time}(A_{\text{boost}}) = O(\text{Time}(A))$ , i.e. no slowdown
- $A$  is used as a black-box

The more compressible is  $T$ , the shorter is  $c'$

Paolo Ferragina, Univer

Two Key Components: Burrows-Wheeler Transform and Suffix Tree

## The empirical entropy $H_0$

$$H_0(T) = -\sum_i (n_i/n) \log_2 (n_i/n) \quad \text{Frequency in } T \text{ of the } i\text{-th symbol}$$

- ❖  $|T| H_0(T)$  is the best you can hope for a memoryless compressor
- ❖ E.g. **Huffman** or **Arithmetic** coders come close to this bound

We get a better compression using a codeword that depends on the  $k$  symbols preceding the one to be compressed (**context**)

Paolo Ferragina, Università di Pisa

# The empirical entropy $H_k$

For any  $k$ -long context

✓ Compress  $T$  up to  $H_k(T)$   
 $\cong$  compress all  $T[\omega]$  up to their  $H_0$

Use Huffman or Arithmetic

$$H_k(T) = (1/|T|) \sum_{|\omega|=k} |T[\omega]| H_0(T[\omega])$$

❖  $T[\omega]$  = string of symbols that precede  $\omega$  in  $T$

Example: Given  $T = \text{"mississippi"}$ , we have  $T[i] = \text{mssp}$ ,  $T[is] = \text{ms}$

Problems with this approach:

- How to go from all  $T[\omega]$  back to the string  $T$ ?
- How do we choose efficiently the best  $k$ ?

BWT

Suffix Tree

Paolo Ferragina, Università di Pisa

# Use BWT to approximate $H_k$

| unknown     | bwt(T) |
|-------------|--------|
| #mississipp | i      |
| i#mississip | p      |
| ippi#missis | s      |
| issippi#mis | s      |
| issippi#    | m      |
| mississippi | #      |
| pi#mississi | p      |
| ppi#mississ | i      |
| sippi#missi | s      |
| sissippi#mi | s      |
| ssippi#miss | i      |
| ssissippi#m | i      |

Remember that...

$$H_k(T) = (1/|T|) \sum_{|\omega|=k} |T[\omega]| H_0(T[\omega])$$

✓ Compress  $T$  up to  $H_k(T)$   
 $\cong$  compress all  $T[\omega]$  up to their  $H_0$   
 $\cong$  compress pieces of  $bwt(T)$  up to  $H_0$

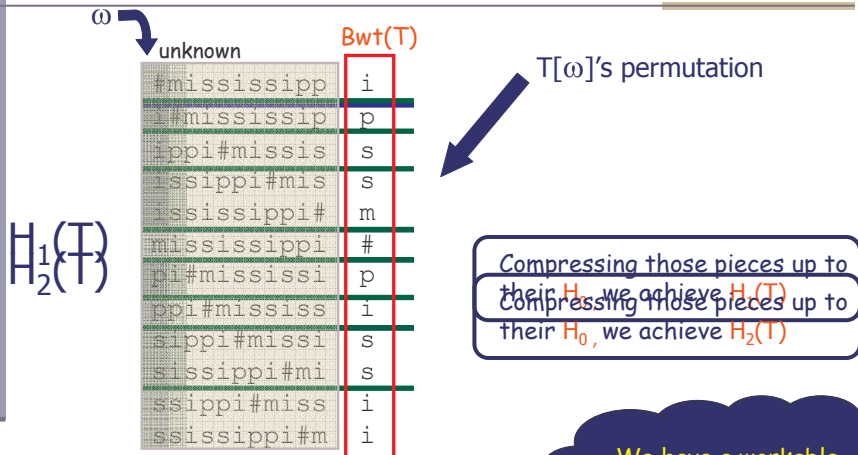
$T[is] = \text{ms}$

$T = \text{mississippi\#}$

$T[\omega]$  is a permutation of a piece of  $bwt(T)$

Paolo Ferragina, Università di Pisa

# What are the pieces of BWT to compress ?

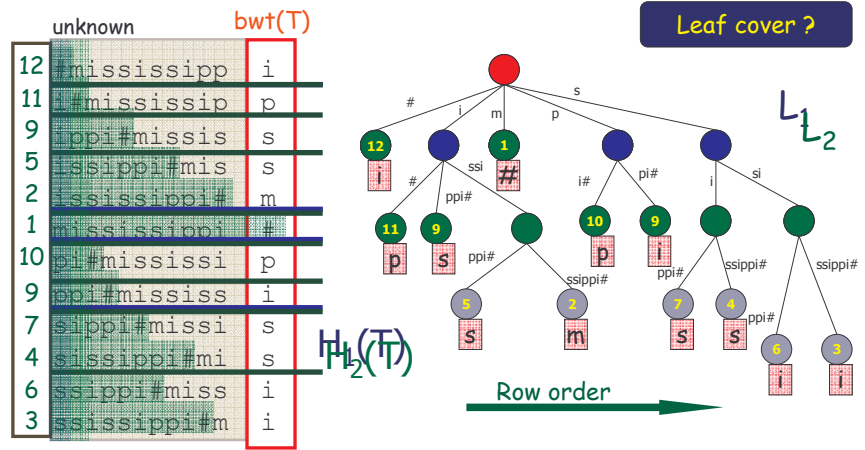


Recall that

- ✓ Compress  $T$  up to  $H_k$
- ≅ compress pieces of  $bwt(T)$  up to  $H_0$

We have a workable way to approximate  $H_k$

# Finding the "best pieces" to compress...



Some leaf covers are "related" to  $H_k$  !!!

## A compression booster [Ferragina-Manzini, SODA 04]

- Let  $A$  be the compressor we wish to boost
- Let  $\text{bwt}(T) = t_1, \dots, t_r$  be the partition induced by the leaf cover  $L$ , and let us define  $\text{cost}(L, A) = \sum_j |A(t_j)|$

**Goal:** Find the leaf cover  $L^*$  of minimum cost

- ✓ It suffices a post-order visit of the suffix tree, hence linear time
- ✓ We have:  $\text{Cost}(L^*, A) \leq \text{Cost}(L_k, A) \approx H_k(T), \forall k$

Technically, we show that

$$|c'| \leq \lambda |s| H_{\frac{1}{\lambda}}(s) + f(|s|) + \log_2 |s| + \gamma_k \quad \forall k$$

Researchers may now concentrate on the "apparently" simpler task of designing 0-th order compressors

Paolo Ferragina, Università di Pisa

[further results joint with Giancarlo-Manzini-Sciortino]

## May we close "the mutual reinforcement cycle" ?

- The Wavelet Tree [Grossi-Gupta-Vitter, Soda 03]
- Using the WT within each piece of the optimal BWT-partition, we get:
  - A compressed index that scales well with the alphabet size  
[joint with Manzini, Makinen, Navarro]
  - Reduce the compression problem to achieve  $H_0$  on binary strings  
[joint with Giancarlo, Manzini, Sciortino]

### Interesting issues:

- ✓ What about space construction of BWT ?
- ✓ What about these tools for "XML" or "Images" ?
- ✓ Other application contexts: bio, db, data mining, network, ...
- ✓ From Theory to Technology ! Libraries, Libraries, .... [e.g. LEDA]

Paolo Ferragina, Università di Pisa

DIMACS Working Group on The Burrows - Wheeler Transform: Ten Years Later - Microsoft Internet Explorer

File Modifica Visualizza Preferiti Strumenti ?

Indietro - - - - - Cerca Preferiti Multimedia - - - - - Vai Collegamenti >>

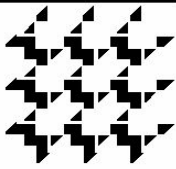
Ingritzzo <http://dimacs.rutgers.edu/Workshops/BWT/>

Google - - - - - Cerca nel Web - - - - - PageRank 10 bloccati - - - - - Vivisimo - - - - -

---

**DIMACS**

*Center for Discrete Mathematics & Theoretical Computer Science  
Founded as a National Science Foundation Science and  
Technology Center*



---

**DIMACS Working Group on The Burrows - Wheeler Transform: Ten Years Later**

August 19 - 20, 2004  
DIMACS Center, CoRE Building, Rutgers University

Organizers:

- Paolo Ferragina, University of Pisa
- Giovanni Manzini, University of Piemonte Orientale
- S. Muthukrishnan, Rutgers University, [muthu@cs.rutgers.edu](mailto:muthu@cs.rutgers.edu)

Presented under the auspices of the Special Focus on [Special Focus on Data Analysis and Mining](#)

---

- [Workshop Announcement](#)
- [Call for Participation](#)

Operazione completata Internet