

---

# *Permuted Longest-Common-Prefix Array*

**Juha Kärkkäinen**

**University of Helsinki**

**Joint work with Giovanni Manzini and Simon Puglisi**

**CPM 2009**

**Lille, France, June 2009**

# Outline

---

1. **Background**
2. Description
3. Construction
4. Compact Representation

# Suffix Array and LCP Array

---

## Suffix array (SA)

- ▶ Powerful text index
- ▶ Linear time construction

## Longest-common-prefix (LCP) array

- ▶ Often used with suffix array
- ▶ Linear time construction from SA

[Kasai & al., CPM 2001]



# Constructing SA and LCP

---

## SA

- ▶  $5-6n$  bytes, fastest in practice
- ▶  $\sim 2n$  bytes, output to disk sequentially
- ▶ external memory

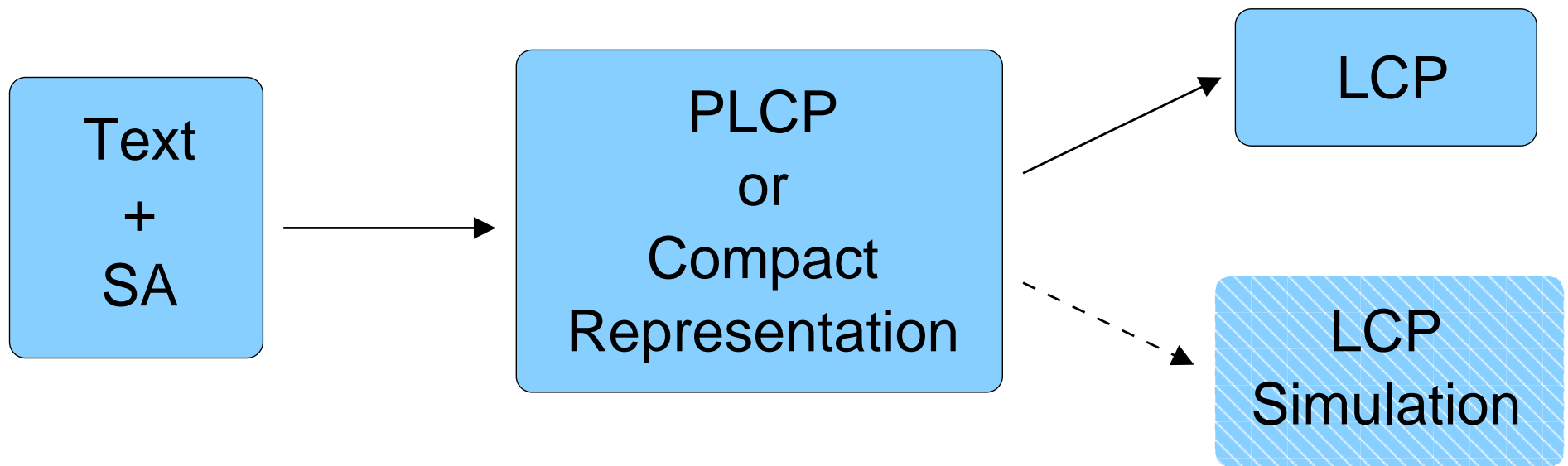
## LCP array from SA

- ▶  $13n$  bytes, **fastest** [Kasai & al., CPM 2001]
- ▶  $9n$  bytes [Mäkinen, F. Inform. 2003]
- ▶  $9n$  or  $6-10n$  bytes [Manzini, SWAT 2004]
- ▶  $\sim 2n$  bytes [Puglisi & Turpin, ISAAC 2008]
  - **semi-external**: sequentially accessed arrays on disk

# Permuted LCP (PLCP) Array

---

- ▶ LCP array in permuted order
- ▶ LCP construction (or simulation) **via PLCP**



- ▶ Faster and more space-efficient

# Constructing LCP from SA

---

## Best previous algorithms

- ▶  $13n$  bytes, fast [Kasai & al., CPM 2001]
- ▶  $\sim 2n$  bytes, semi-external [Puglisi & Turpin, ISAAC 2008]

## New algorithms via PLCP

- ▶ semi-external
- ▶  $5n$  bytes,  $\mathcal{O}(n)$  time, faster
- ▶  $11n$  **bits**,  $\mathcal{O}(n \log n)$  time
- ▶  $(1 + 4/q)n$  bytes,  $\mathcal{O}(nq)$  time

# Outline

---

1. Background
2. **Description**
3. Construction
4. Compact Representation

# Suffix Array (SA)

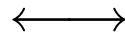
Sorted array of all suffixes

$T = \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 & 6 \\ \text{BANANA} & \# \end{array}$

Pos. order

0	BANANA#
1	ANANA#
2	NANA#
3	ANA#
4	NA#
5	A#
6	#

Position



Lex. order

6	#
5	A#
3	ANA#
1	ANANA#
0	BANANA#
4	NA#
2	NANA#

Suffix array



# LCP Array

Length of **longest common prefix** with preceding suffix

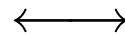
$$T = \begin{array}{cccccc} & 0 & 1 & 2 & 3 & 4 & 5 & 6 \\ & \underline{B} & \underline{A} & \underline{N} & \underline{A} & \underline{N} & \underline{A} & \# \end{array}$$

Pos. order

0	BANANA#
1	ANANA#
2	NANA#
3	ANA#
4	NA#
5	A#
6	#

Lex. order

-	6	#
0	5	A#
1	3	<u>ANA#</u>
<b>3</b>	1	<u>ANANA#</u>
0	0	BANANA#
0	4	NA#
2	2	<u>NANA#</u>



LCP array

SA

# PLCP Array

LCP array **permuted** into the **position order**

$$T = \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 & 6 \\ \text{BANANA}\# \end{array}$$

Pos. order

0	0	BANANA#
3	1	<b>ANANA</b> #
2	2	<b>NANA</b> #
1	3	<b>ANA</b> #
0	4	NA#
0	5	A#
-	6	#

**PLCP**

Lex. order

-	6	#
0	5	A#
1	3	<b>ANA</b> #
3	1	<b>ANANA</b> #
0	0	BANANA#
0	4	NA#
2	2	<b>NANA</b> #

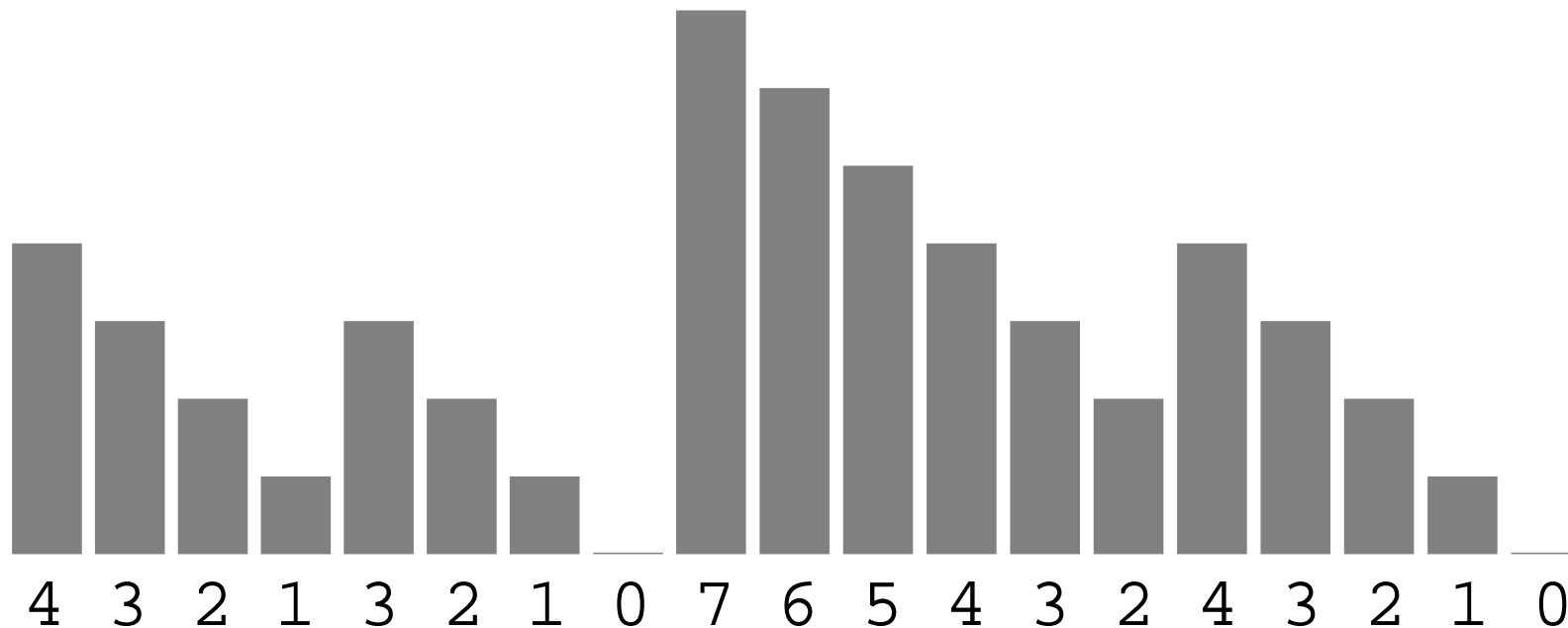
**LCP**



# PLCP Array Structure

---

**Theorem.**  $PLCP[i] \geq PLCP[i - 1] - 1$ .



- ▶ Used in all efficient LCP-from-SA algorithms
- ▶ Enables compact representation

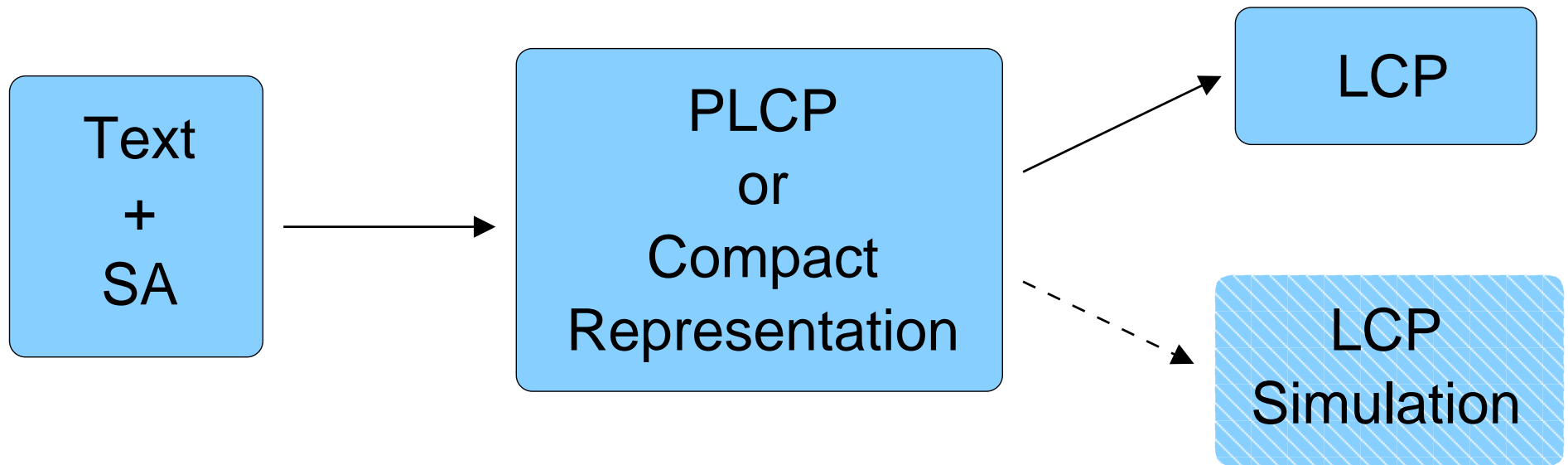
# Outline

---

1. Background
2. Description
3. **Construction**
  - ▶ Irreducible LCP Algorithm
  - ▶  $\Phi$  Algorithm
4. Compact Representation

# PLCP Algorithms

---



- ▶ Irreducible LCP Algorithm
- ▶  $\Phi$  Algorithm

$$\text{▶ } LCP[i] = PLCP[SA[i]]$$

# Irreducible LCP Values

$LCP[i]$  is **irreducible**  $\iff BWT[i] \neq BWT[i - 1]$

$T =$  0 1 2 3 4 5 6  
BANANA#

Pos. order

0	0	BANANA#
3	1	<b>ANANA#</b>
2	2	<b>NANA#</b>
1	3	<b>ANA#</b>
0	4	NA#
0	5	A#
-	6	#

PLCP

$\iff$

Lex. order

-	A	6	#
0	N	5	A#
1	N	3	<b>ANA#</b>
3	B	1	<b>ANANA#</b>
0	#	0	BANANA#
0	A	4	NA#
2	A	2	<b>NANA#</b>

LCP

BWT

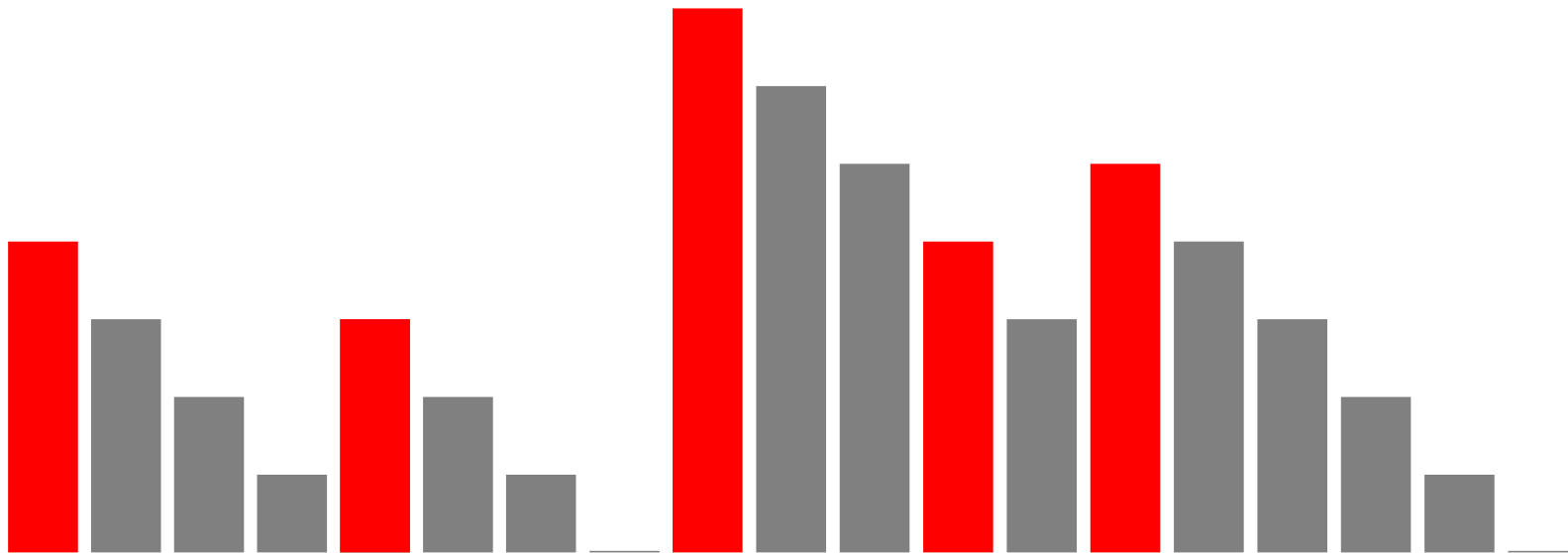
# Irreducible LCP Values

---

**Theorem.** The sum of **irreducible** LCP values is  $\leq 2n \log n$

**Lemma.** For reducible LCP value,  
 $PLCP[i] = PLCP[i - 1] - 1$

[Manzini, SWAT 2004]



# Irreducible LCP Algorithm

---

1. Compute all irreducible PLCP values **naively**
  2. Compute reducible PLCP values:  
 $PLCP[i] := PLCP[i - 1] - 1$
  3. Compute LCP from PLCP
- 
- ▶  $O(n \log n)$  **time**
    - Faster than previous algorithms
  - ▶  $5n$  **bytes**
    - semi-external  $\implies$  10% slower



# $\Phi$ Algorithm

---

1. Compute  $\Phi$ :      **for**  $i := 1$  **to**  $n$  **do**  $\Phi[SA[i]] := SA[i - 1]$
2. Compute PLCP:       $\ell := 0$   
                          **for**  $i := 0$  **to**  $n - 1$  **do**  
                                  **while**  $T[i + \ell] = T[\Phi[i] + \ell]$  **do**  $\ell++$   
   $PLCP[i] := \ell; \ell := \max(\ell - 1, 0)$
3. Compute LCP:      **for**  $i := 1$  **to**  $n$  **do**  $LCP[i] := PLCP[SA[i]]$

Over 30% faster than any previous algorithm

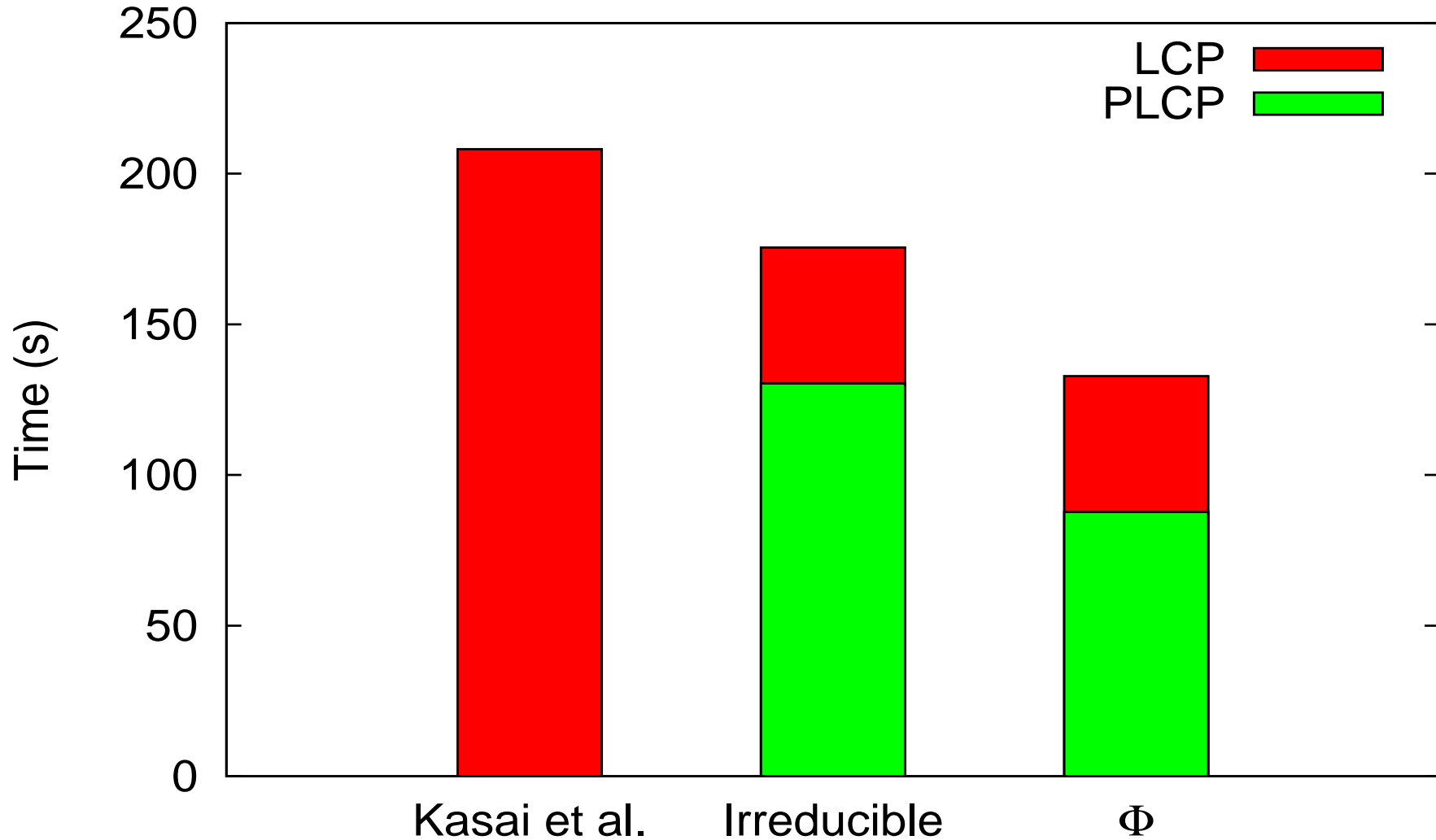
- ▶ Only one random access in each round of main loop

5n bytes

- ▶ semi-external  $\implies$  10% slower

# *Time for Manzini Corpus (708 MiB)*

---



# Outline

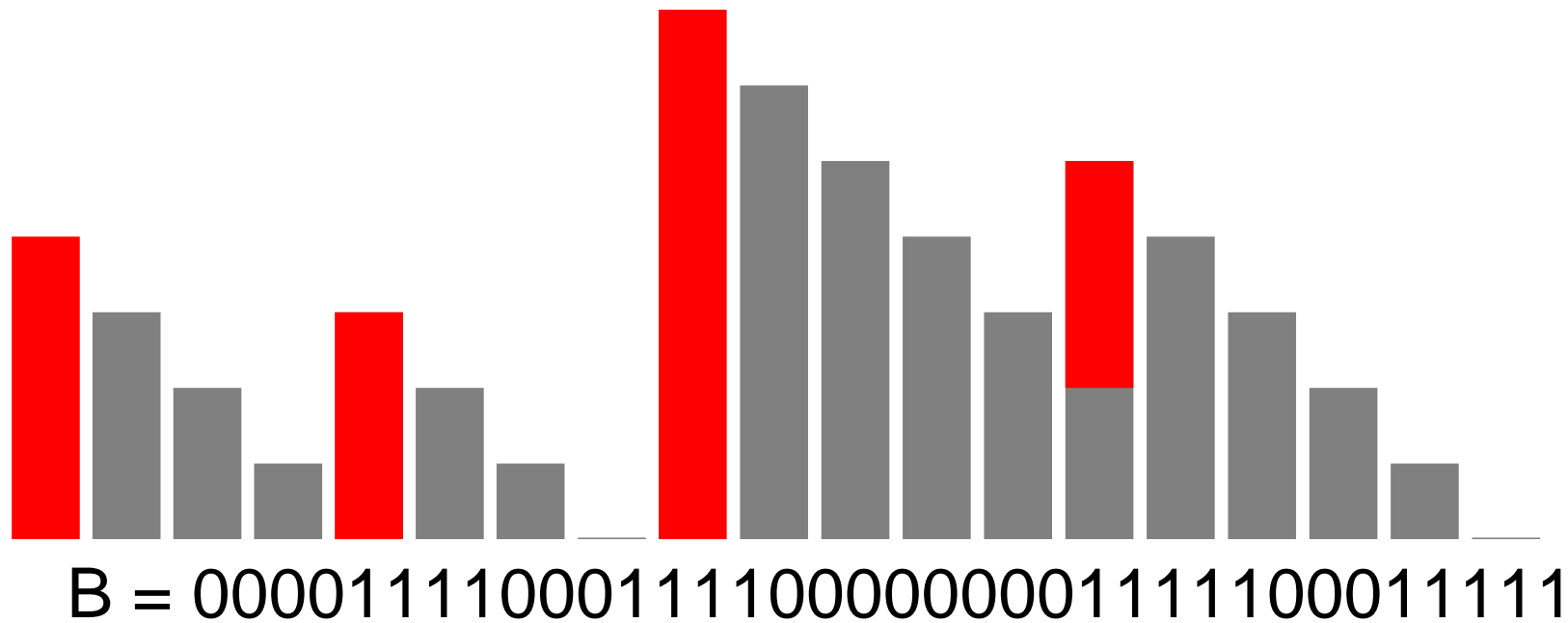
---

1. Background
2. Description
3. Construction
4. **Compact Representation**
  - ▶ Succinct Bitvector
  - ▶ Sparse PLCP Array

# Succinct Bitvector Representation

PLCP stored in bitvector  $B$  of  $2n$  bits

[Sadakane, SODA 2002]

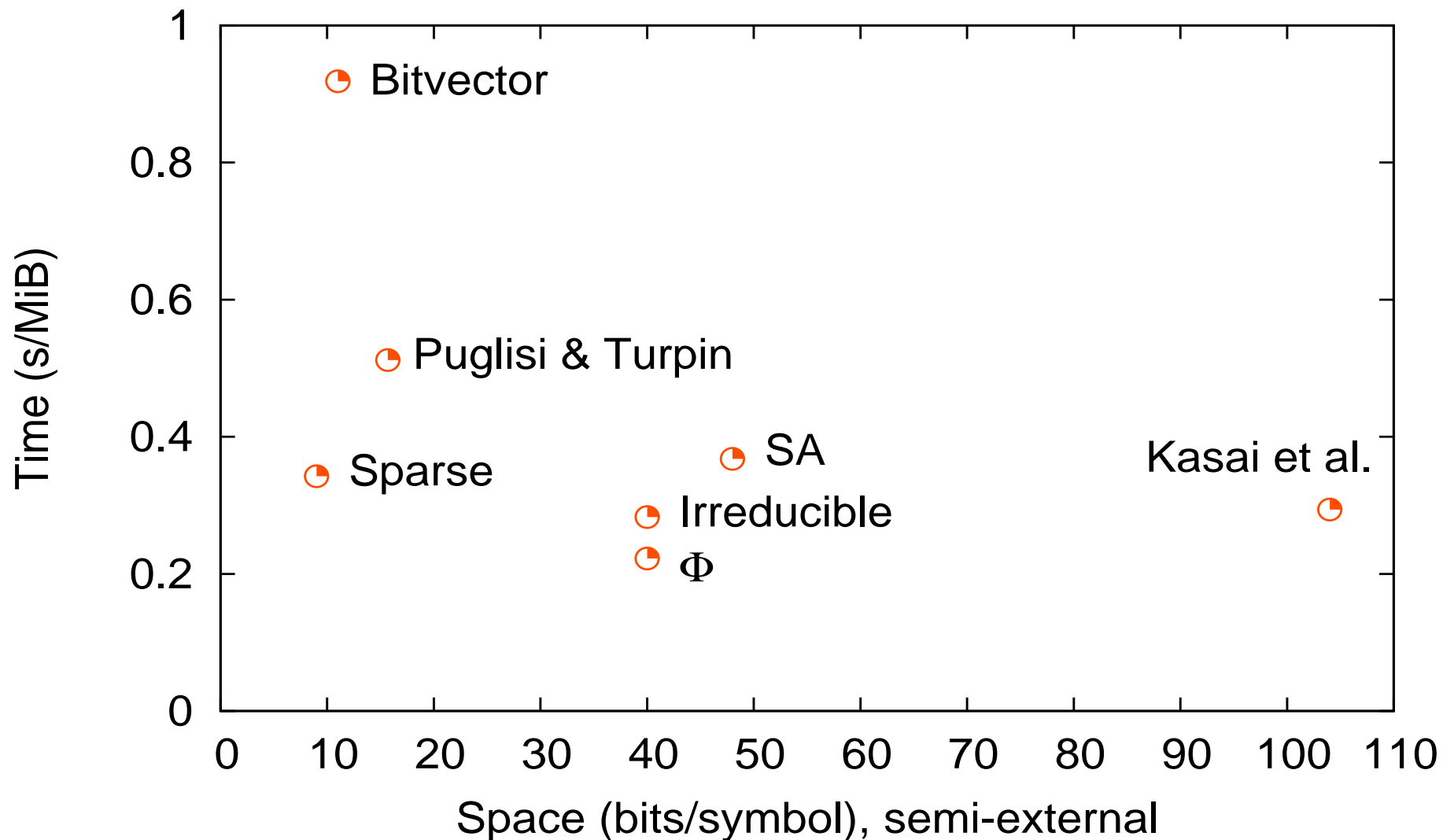


Irreducible LCP Algorithm

- ▶  $\mathcal{O}(n \log n)$  time
- ▶  $11n$  bits (semi-external)



# Space and Time for Manzini Corpus (708 MiB)



# *Conclusion*

---

## Summary

- ▶ Faster and more space-efficient algorithms for (P)LCP array construction
- ▶ Irreducible LCP Theorem

## Open problem

- ▶ External memory algorithm for LCP-from-SA construction

# *Dmitry Khmelev*

---

- ▶ Died at young age in 2004
- ▶ Mathematician doing stringology as hobby
- ▶ Conjectured the Irreducible LCP Theorem
- ▶ LCP construction algorithm
  - Implementation, no publication
  - Uses sparse PLCP array
  - Analysis relies on Irreducible LCP Theorem