#### Compressed String Dictionary Look-up with Edit Distance One

Djamal Belazzougui

LIAFA - Paris 7

Rossano Venturini

University of Pisa



## Problem

- Design a (compressed) index for a set  $D=\{S_1, S_2, ..., S_d\}$ of d strings of total length n drawn from an alphabet of size  $\sigma$
- Given any pattern P, report all the occ strings in D having Edit Distance at most one from P
  - edit(P, S) = minimum # of edit operations (insertion, deletion, substitution) to transform P into S

### Motivations

### Motivations

• Information Retrieval

### Motivations

#### Information Retrieval

- e.g., users may misspell terms in their queries
- Database
  - e.g., you may do not remember the correct spelling of a key

#### • Data Mining

• e.g, clean a noised dataset

- G. Brodal and L. Gasieniec [CPM'96]
  - Space: O(σn log n) bits
  - Time: O(|P| + occ)

- G. Brodal and L. Gasieniec [CPM'96]
  - Space: O(σn log n) bits
  - Time: O(|P| + occ)

Reduced to membership queries in a super set containing all the strings that have edit distance one from any string in D.

For each  $S \in D$ , we add  $O(\sigma |S|)$ 

- G. Brodal and L. Gasieniec [CPM'96]
  - Space: O(σn log n) bits
  - Time: O(|P| + occ)
- G. Brodal and L. Gasieniec [CPM'96]
  - Space: O(n log n) bits
  - Time:  $O(|P| + \log n + occ)$

- G. Brodal and L. Gasieniec [CPM'96]
  - Space: O(σn log n) bits
  - Time: O(|P| + occ)
- G. Brodal and L. Gasieniec [CPM'96]
  - Space: O(n log n) bits
  - Time:  $O(|P| + \log n + occ)$

For each S∈D, consider its factorizations of the form S[1,i-1] S[i] S[i+1,s]

- G. Brodal and L. Gasieniec [CPM'96]
  - Space: O(σn log n) bits
  - Time: O(|P| + occ)
- G. Brodal and L. Gasieniec [CPM'96]
  - Space: O(n log n) bits
  - Time:  $O(|P| + \log n + occ)$

For each S∈D, consider its factorizations of the form S[1,i-1] S[i] S[i+1,s]If P[1,i-1] == S[1,i-1] and P[i,p] == S[i+1,s], and s=p+1 then, S has edit distance one from P

- G. Brodal and L. Gasieniec [CPM'96]
  - Space: O(σn log n) bits
  - Time: O(|P| + occ)
- G. Brodal and L. Gasieniec [CPM'96]
  - Space: O(n log n) bits
  - Time:  $O(|P| + \log n + occ)$

For each S∈D, consider its factorizations of the form S[1,i-1] S[i] S[i+1,s] If P[1,i-1] == S[1,i-1] and P[i,p] == S[i+1,s], and s=p+1 then, S has edit distance one from P Use Balanced Search Trees to index these factorizations

- G. Brodal and L. Gasieniec [CPM'96]
  - Space: O(σn log n) bits
  - Time: O(|P| + occ)
- G. Brodal and L. Gasieniec [CPM'96]
  - Space: O(n log n) bits
  - Time:  $O(|P| + \log n + occ)$
- D. Belazzougui [CPM'09]
  - Space:  $O(n \log \sigma)$
  - Time: O(|P| + occ)

- G. Brodal and L. Gasieniec [CPM'96]
  - Space: O(σn log n) bits
  - Time: O(|P| + occ)
- G. Brodal and L. Gasieniec [CPM'96]
  - Space: O(n log n) bits
  - Time:  $O(|P| + \log n + occ)$
- D. Belazzougui [CPM'09]
  - Space:  $O(n \log \sigma)$
  - Time: O(|P| + occ)



#### Solution based on Patricia tries and Perfect hashing

Time O(|P| + occ)

**Space**  $2nH_k + n \cdot o(\log \sigma) + 2d\log d$ 

#### Solution based on Patricia tries and Perfect hashing

**Time** O(|P| + occ)

**Space**  $2nH_k + n \cdot o(\log \sigma) + 2d\log d$ 

#### Solution based on Compressed Permuterm index

Time  $O(|P| \log \log \sigma + occ)$ Space  $nH_k + n \cdot o(\log \sigma)$  for  $\sigma = \omega(\log^c n)$ 

#### Solution based on Patricia tries and Perfect hashing

Time O(|P| + occ)

**Space**  $2nH_k + n \cdot o(\log \sigma) + 2d\log d$ 

**Top-k strings in**  $O(|P| + k \log k)$ 

Solution based on Compressed Permuterm index

**Time**  $O(|P| \log \log \sigma + occ)$ **Space**  $nH_k + n \cdot o(\log \sigma)$  for  $\sigma = \omega(\log^c n)$ 

#### $D = \{\texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc}\}$

 $D = \{\texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc}\}$ 

 $D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc} \}$ 

 $cmp(D) = {\tt abcc\ accb\ baca\ caac\ cbcc}$  -

access time: O(1)space:  $H_k$  + n o(log  $\sigma$ )

$$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc} \}$$



$$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc} \}$$

5



$$D = \{\texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc}\}$$



$$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc} \}$$

5

С

5



$$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc} \}$$



$$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc} \}$$





#### Factorizations of the form

$$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc} \}$$





-



$\langle 0,$	0,	a,	3,	$5\rangle$			
$\langle 1,$	1,	b,	2,	$3\rangle$			
$\langle 1,$	2,	с,	1,	$3\rangle$			
$\langle 1,$	3,	с,	0,	0 angle			
abcc							



$$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc} \}$$



$$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc} \}$$











$$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc} \}$$







		S[i]					
$\langle 0,$	0,	a,	3,	$5\rangle$			
$\langle 1,$	1,	b,	2,	$3\rangle$			
$\langle 1,$	2,	с,	1,	$3\rangle$			
$\langle 1,$	3,	с,	0,	0 angle			
abcc							

Factorizations of the form

$$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc} \}$$







S[i+1,s] in PT<sup>r</sup>  $\langle 0, 0, a, 3, 5 
angle$  $\langle 1, 1, b, 2, 3 \rangle$  $\langle 1, 2, c, 1, 3 \rangle$  $\langle 1, 3, \mathsf{c}, 0, 0 \rangle$ abcc

Factorizations of the form

$$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc} \}$$







s-i-1  $\langle 0, 0, a, 3, 5 
angle$  $\langle 1, 1, b, 2, 3 \rangle$  $\langle 1, 2, \mathbf{c}, 1, 3 \rangle$  $\langle 1, 3, \mathsf{c}, 0, 0 \rangle$ abcc



$$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc} \}$$



$$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc} \}$$

3







$\langle 0,$	0,	a,	3,	$5\rangle$			
$\langle 1,$	1,	b,	2,	$3\rangle$			
$\langle 1,$	2,	с,	1,	$3\rangle$			
$\langle 1,$	3,	с,	0,	0 angle			
abcc							












$$\langle 0, 0, a, 3, 5 \rangle$$
  
 $\langle 1, 1, b, 2, 3 \rangle$   
 $\langle 1, 2, c, 1, 3 \rangle$   
 $\langle 1, 3, c, 0, 0 \rangle$   
abcc

Index tuples so that, given  $np_i, i, sp_i, p-i+1$ , return the symbols, if any, in tuples of the form  $\langle np_i, i, \star, p-i+1, ns_i \rangle$ .

Define the function F, so that, F(np, i, j, sp) = ciff there exists tuple  $\langle np, i, \mathbf{c}, j, sp \rangle$ 

> $\langle 0, 0, a, 3, 5 \rangle$  $\langle 1, 1, b, 2, 3 \rangle$  $\langle 1, 2, c, 1, 3 \rangle$  $\langle 1, 3, c, 0, 0 \rangle$ abcc

Index tuples so that, given  $np_i, i, sp_i, p-i+1$ , return the symbols, if any, in tuples of the form  $\langle np_i, i, \star, p-i+1, ns_i \rangle$ .

Define the function F, so that,

F(np, i, j, sp) = c

iff there exists tuple

 $\langle np, i, \mathbf{c}, j, sp \rangle$ 

A function F from a subset  $S \subseteq U$  of size n to an alphabet of size  $\sigma$  can be represented within  $nH_0 + n \ o(H_0)$  bits, so that if  $x \in S$ , F(x) evaluated in constant time if  $x \notin S$ , an arbitrary symbol is returned

> Index tuples so that, given  $np_i, i, sp_i, p-i+1$ , return the symbols, if any, in tuples of the form  $\langle np_i, i, \star, p-i+1, ns_i \rangle$ .

 $\langle 0, 0, a, 3, 5 \rangle$  $\langle 1, 1, b, 2, 3 \rangle$  $\langle 1, 2, c, 1, 3 \rangle$  $\langle 1, 3, c, 0, 0 \rangle$ abcc



Define the function F, so that,

F(np, i, j, sp) = c

iff there exists tuple

 $\langle np, i, \mathbf{c}, j, sp \rangle$ 

A function F from a subset  $S \subseteq U$  of size n to an alphabet of size  $\sigma$  can be represented within  $nH_0 + n o(H_0)$  bits, so that if  $x \in S$ , F(x) evaluated in constant time if  $x \notin S$ , an arbitrary symbol is returned

> Index tuples so that, given  $np_i, i, sp_i, p-i+1$ , return the symbols, if any, in tuples of the form $\langle np_i, i, \star, p-i+1, ns_i \rangle$ .

 $\langle 0, 0, a, 3, 5 \rangle$  $\langle 1, 1, b, 2, 3 \rangle$  $\langle 1, 2, c, 1, 3 \rangle$  $\langle 1, 3, c, 0, 0 \rangle$ abcc

Define the function F, so that,

F(np, i, j, sp) = c

iff there exists tuple

 $\langle np, i, \mathbf{c}, j, sp \rangle$ 

A function F from a subset  $S \subseteq U$  of size n to an alphabet of size  $\sigma$  can be represented within  $nH_0 + n \ o(H_0)$  bits, so that if  $x \in S$ , F(x) evaluated in constant time if  $x \notin S$ , an arbitrary symbol is returned

 $\langle 0, 0, a, 3, 5 \rangle$  $\langle 1, 1, b, 2, 3 \rangle$  $\langle 1, 2, c, 1, 3 \rangle$  $\langle 1, 3, c, 0, 0 \rangle$ abcc

We may err!  $P[1,i] \subset P[i+1,p]$ may not be in D. We need to check if it belongs to D!

tuples so that,  $p_i, i, sp_i, p-i+1$ , in the symbols, in tuples of the  $i, \star, p-i+1, ns_i \rangle$ .

Define the function F, so that,

F(np, i, j, sp) = c

iff there exists tuple

 $\langle np, i, \mathbf{c}, j, sp \rangle$ 



Define the function F, so that,

F(np, i, j, sp) = c

iff there exists tuple

 $\langle np, i, \mathbf{c}, j, sp \rangle$ 

Space:  $2nH_k + n \cdot o(\log \sigma) + 2d \log d$ Time: O(|P| + occ)

We may err!  $P[1,i] \ c \ P[i+1,p]$ may not be in D. We need to check if it belongs to D if it belongs to D  $n \ O(1)$ . see the paper  $\langle 1, 3, c, 0, 0 \rangle$ 

abcc



Subset of strings having the same length p+1

a b c c #
b c c # a
c c # a b
c # a b c
# a b c c

a	b	С	С	#
b	С	С	#	a
с	С	#	a	b
С	#	a	b	С
#	a	b	С	с
а	С	С	b	#
С	С	b	#	a
С	b	#	a	С
b	#	a	С	с
#	a	С	С	b
b	a	С	a	#
a	С	a	#	b
с	a	#	b	a
a	#	b	a	С
#	b	a	С	a
С	a	a	С	#
a	a	С	#	С
a	С	#	С	a
С	#	С	a	a
#	С	a	a	С
С	b	С	С	#
b	С	С	#	С
С	С	#	С	b
с	#	С	b	с
#	С	b	С	С

a	b	С	С	#
b	С	С	#	a
с	С	#	a	b
с	#	a	b	с
#	a	b	С	с
a	С	С	b	#
с	С	b	#	a
с	b	#	a	с
b	#	a	С	с
#	a	С	С	b
b	a	С	a	#
a	С	a	#	b
с	a	#	b	a
a	#	b	a	с
#	b	a	С	a
с	a	a	С	#
a	a	С	#	с
a	С	#	С	a
с	#	С	a	a
#	С	a	a	с
с	b	С	С	#
b	С	С	#	с
с	С	#	С	b
с	#	С	b	с
#	С	b	С	с



$D = \{\texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc}\}$	Subset of strings having the same length p+1
	F L
abcc#	<b>#</b> a b c c
bcc#a	#accb
cc#ab	#baca
c # a b c	# c a a <b>c</b>
# a b c c	# с b с <b>с</b>
accb#	a#bac
ccb#a	aac#c
c b # a c	abcc#
b#acc	ac#ca
#accb	aca#b
baca#	accb#
aca#b	b # a c c
ca#ba	baca#
a#bac	bcc#a
#baca	<b>b</b> с с <b># с</b>
caac#	c#abc
aac#c	c#caa
ac#ca	с # с b с

c a # b a

caac#

c b # a c

c b c c #

c c # a b

c c # c b

c c b # a

c b c c # b c c # c c c # c b c # c b c

**#**сbс**с** 

c#caa

#caac

F				
#	а	b	С	С
#	а	С	С	b
#	b	a	С	a
#	С	a	a	с
#	С	b	с	с
a	#	b	a	с
a	a	С	#	с
a	b	С	С	#
a	С	#	С	a
a	С	а	#	b
a	С	С	b	#
b	#	a	С	с
b	a	С	a	#
b	С	С	#	a
b	С	С	#	с
С	#	a	b	с
С	#	С	a	a
С	#	С	b	С
С	a	#	b	a
С	a	a	С	#
С	b	#	a	С
С	b	С	С	#
С	С	#	a	b
С	С	#	С	b
С	С	b	#	a

$D = \{\texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc}\}$	Subset of strings having the same length p+1
	F # a b c c a c c b
Burrows-Wheeler Transform	m baca caac
	# bcc a # bac a a c # c
	a b c c # a c # c a
	aca#b accb# b#acc
	baca# bcc#a
	bcc#c c#abc c#caa
	с # с b с с а # b а
	саас # с b # а с с b с с #
	c c # a b c c # c b
	ccb#a

Given P[1,p], identify the rows prefixed by a cyclic rotation of P on the subset with string of len p+1 i.e., P[i,p]#P[1,i-1]

F				L
#	а	b	С	С
#	a	С	С	b
#	b	a	С	a
#	С	a	a	С
#	С	b	С	С
a	#	b	a	С
a	а	С	#	С
a	b	С	С	#
a	С	#	С	a
a	С	a	#	b
a	С	С	b	#
b	#	а	С	С
b	a	С	a	#
b	С	С	#	a
b	С	С	#	С
С	#	а	b	С
С	#	С	а	a
С	#	С	b	С
С	a	#	b	a
С	a	a	С	#
С	b	#	a	С
С	b	С	С	#
С	С	#	a	b
С	С	#	С	b
с	С	b	#	a

$D = \{\texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc}\}$	Subset of strings having the same length p+1
	F L
Given P[1,p],	# a b c c
identify the rows prefixed by a	# a c c b
cyclic rotation of P on the	#baca
subset with string of len p+l	# c a a <b>c</b>
i.e., P[i,p]#P[1,i-1]	<b>#</b> сbс <b>с</b>
	a # b a c
Strings corresponding to	a a c # c
strings corresponding to	abcc#
these rows are at distance one	ac#ca
from P (an insertion)!	aca#b
	accb#
	b#acc
	baca#
	bcc#a
	b c c # c
	c # a b c
	c#caa
	c # c b c
	c a # b a
	caac#
	c b # a c
	с b с с #
	c c # a b
	c c # c b
	ccb#a

$D = \{\texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc}\}$	Subset of strings having the same length p+1
Given P[1,p], identify the rows prefixed by a cyclic rotation of P on the subset with string of len p+1	FL #abcc #accb #baca #caac
i.e., P[i,p]#P[1,i-1]	# c b c c a # b a c a a c # c
Strings corresponding to these rows are at distance one from P (an insertion)!	a b c c # a c # c a a c a # b
P = acc	accb# b#acc baca#
	bcc#a bcc#c c#abc
	c # c a a c # c b c
	ca#ba caac# cb#ac
	c b c c # c c # a b c c # c b

$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac} \}$	z, cbcc} Subset of strings having the same length p+1
Civen P[1 n]	<u>F L</u>
Given F[1,p],	# a b c c
identify the rows prefixed by a	# a c c b
cyclic rotation of P on the	#baca
subset with string of len p+l	#caac
i.e., P[i,p]#P[1,i-1]	# с b с <b>с</b>
	a#bac
Strings corresponding to	aac#c
these rows are at distance one	abcc#
from D (on incontion)	a c # c a
from F (an insertion):	aca#b
	accb#
P = acc	b#acc
One rotation is cc#a	baca#
	bcc#a
	b c c # c
	c # a b c
	c#caa
	c # c b c
	ca#ba
	caac#
	c b # a c
	<b>c</b> b c c <b>#</b>
	cc#ab
	c c # c b
	ссb#а

$D = \{\texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbcc}\}$	Subset of strings having the same length p+1
Civen D[1 n]	<u>F L</u>
Given P[1,p],	<b>#</b> a b c c
identify the rows prefixed by a	#accb
cyclic rotation of P on the	#baca
subset with string of len p+l	#caac
i.e., P[i,p]#P[1,i-1]	# с b с <b>с</b>
	a#bac
Chuin an annuar an dina ta	aac#c
Strings corresponding to	abcc#
these rows are at distance one	ac#ca
from P (an insertion)!	aca#b
	accb#
<b>P</b> = acc	b#acc
One rotation is cc#a	baca#
abaa is in D	bcc#a
	b с с # с
	c # a b c
	c#caa
	с # с b с
	ca#ba
	caac#
	c b # a c
	с b с с #
	c c # ab
	c c # c b
	ссb#а

$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{c} \}$	caac, cbcc} Subset of strings having the same length p+1
Circus D[1 m]	F L
Given P[1,p],	<b>#</b> a b c c
identify the rows prefixed by a	#accb
cyclic rotation of P on the	#baca
subset with string of len p+l	# c a a <b>c</b>
i.e., P[i,p]#P[1,i-1]	<b>#</b> сbс <b>с</b>
	a#bac
Strings corresponding to	aac#c
strings corresponding to	abcc#
these rows are at distance one	ac#ca
from P (an insertion)!	aca#b
	accb#
<b>P</b> = acc	b#acc
One rotation is cc#a	baca#
abcc <b>is in D</b>	bcc#a
	b c c # c
	c # a b c
How to identify rows	c#caa
prefixed by P[i,p]#P[1,i-1]?	c # c b c
	ca#ba
	caac#
	c b # a c
	<u>cbcc</u> #
	c c # ab
	c c # c b
	c c b # a

$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{bac} \}$	a, caac, cbcc} Subset of strings having the same length p+1
Circan D[1 p]	F L
Given P[1,p],	# a b c c
identify the rows prefixed by a	#accb
cyclic rotation of P on the	#baca
subset with string of len p+l	# c a a <b>c</b>
i.e., P[i.p]#P[1.i-1]	<b>#</b> c b c <b>c</b>
	a # b a c
Casiliana a success and line to	a a c # c
Strings corresponding to	abcc#
these rows are at distance one	a c # c a
from P (an insertion)!	aca#b
	a c c b #
P = acc	b#acc
One rotation is cc#a	baca#
abaa is in D	bcc#a
	b c c # c
	c # a b c
How to identify rows	c#caa
prefixed by P[i,p]#P[1,i-1]?	с # с b с
P. ee. e/ . [.,p] [.,,].	c a # b a
) Peolouroud ecouch feu ecoh of	caac#
-) backward search for each of	c b # a c
them? O(p <sup>2</sup> ) time	с b с с #
	c c # ab
	c c # c b
	c c b # a

$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{o} \}$	caac, cbcc} Subset of strings having the same length p+1
Circus D[1 m]	F L
Given P[1,p],	<b>#</b> a b c c
identify the rows prefixed by a	#accb
cyclic rotation of P on the	#baca
subset with string of len p+l	# c a a <b>c</b>
i.e., P[i,p]#P[1,i-1]	<b>#</b> сbс <b>с</b>
	a#bac
Strings corresponding to	aac#c
Strings corresponding to	abcc#
these rows are at distance one	a c # c a
from P (an insertion)!	aca#b
	accb#
P = acc	b#acc
One rotation is cc#a	baca#
abcc is in D	b c c # a
	b c c # c
	c # a b c
How to identify rows	c#caa
prefixed by P[i,p]#P[1,i-1]?	c # c b c
	c a # b a
-) Backward search for each of	caac#
-) Dackward search for each of $(-2)$ time	c b # a c
them: O(p-) time	<u>cbcc</u> #
	c c # a b
-) Move from P[I,p]#P[1,I-1] to	c c # c b
P[i-1,p]#P[1,i-2] in O(log log σ)	c c b # a

$D = \{ \texttt{abcc}, \texttt{accb}, \texttt{baca}, \texttt{caac}, \texttt{cbc} \}$	CC Subset of strings hav
Given P[1,p],	F L # a b c c
identify the rows prefixed by a	# a c c b
cyclic rotation of P on the	#baca
subset with string of len p+l	# c a a <b>c</b>
i.e., P[i,p]#P[1,i-1]	# c b c <b>c</b>
	a#bac
Strings corresponding to	aac#c
Space: $nH_k + n \cdot$ Time: $O( P  \log 1)$	$o(\log \sigma)$ og $\sigma + occ$
Space: $nH_k + n \cdot$ Time: $O( P  \log l$	$\log \sigma + occ$ ) $\log \sigma + occ$
Space: $nH_k + n \cdot$ Time: $O( P  \log l$	$\frac{O(\log \sigma)}{\log \sigma + occ}$
Space: $nH_k + n \cdot$ Time: $O( P  \log l$ How to identify rows	$O(\log \sigma)$ $O(\sigma + occ)$ $\frac{\sigma + occ}{c \# a b c}$ $c \# a b c$
Space: $nH_k + n \cdot$ Time: $O( P  \log l$ How to identify rows prefixed by P[i,p]#P[1,i-1]?	$O(\log \sigma)$ $O(\sigma + occ)$ $\frac{\sigma + occ}{c \# a b c}$ $c \# a b c$ $c \# c a a$ $c \# c b c$ $c \# c b c$
Space: $nH_k + n \cdot$ Time: $O( P  \log 1)$ How to identify rows prefixed by P[i,p]#P[1,i-1]?	$O(\log \sigma)$ $O(\sigma + occ)$ $D = C = C$ $C = a = b = c$ $C = c = b = c$ $C = c = b = c$ $C = a = c = b = a$ $C = a = c = b = a$ $C = a = c = b = a$
Space: $nH_k + n \cdot$ Time: $O( P  \log 1$ How to identify rows prefixed by P[i,p]#P[1,i-1]? -) Backward search for each of them? $O(n^2)$ time	$O(\log \sigma)$ $O(\log \sigma)$ $O(\sigma + occ)$ $D = C = C$ $C = a b c$ $C = c a a$ $C = c b c$ $C = a b c$ $C = a b c$ $C = a c = b c$ $C = b = a c$
Space: $nH_k + n \cdot$ Time: $O( P  \log 1$ How to identify rows prefixed by P[i,p]#P[1,i-1]? -) Backward search for each of them? O(p <sup>2</sup> ) time	$O(\log \sigma)$ $O(\log \sigma)$ $O(\cos \sigma + occ)$ $D \ c \ c \ w \ c \ b \ c \ c \ a \ b \ c \ c \ a \ c \ b \ c \ c \ a \ c \ b \ c \ c \ a \ c \ b \ c \ c \ a \ c \ b \ c \ c \ a \ c \ b \ c \ c \ c \ c \ c \ c \ c \ c$
Space: $nH_k + n \cdot$ Time: $O( P  \log 1$ How to identify rows prefixed by P[i,p]#P[1,i-1]? -) Backward search for each of them? $O(p^2)$ time -) Move from P[i p]#P[1 i-1] to	$O(\log \sigma)$ $O(\log \sigma)$ $D \subset C = C$ $C = a b C$ $C = c a a$ $C = c b C$ $C = a b a$ $C = a c = b a$ $C = a c = c$ $C = b = c$ $C = $
Space: $nH_k + n \cdot$ Time: $O( P  \log 1)$ How to identify rows prefixed by P[i,p]#P[1,i-1]? -) Backward search for each of them? $O(p^2)$ time -) Move from P[i,p]#P[1,i-1] to P[i, 1, p]#P[1, i, 2] in $O(\log \log \sigma)$	$O(\log \sigma)$ $O(\log \sigma)$ $O(\cos \sigma)$

### Solution based on Patricia tries and Perfect

Time O(|P| + occ)

**Space**  $2nH_k + n \cdot o(\log \sigma) + 2d\log d$ 

**Top-k strings in**  $O(|P| + k \log k)$ 

### Solution based on Patricia tries and Perfect

Time O(|P| + occ)

**Space**  $2nH_k + n \cdot o(\log \sigma) + 2d\log d$ 

**Top-k strings in**  $O(|P| + k \log k)$ 

Solution based on Compressed Permuterm index

**Time**  $O(|P| \log \log \sigma + occ)$ **Space**  $nH_k + n \cdot o(\log \sigma)$  for  $\sigma = \omega(\log^c n)$ 

Solution based on Patricia tries and Perfect



**Space**  $2nH_k + n \cdot o(\log \sigma) + 2d \log d$ 

**Top-k strings in**  $O(|P| + k \log k)$ 

Solution based on Compressed Permuterm index

**Time** 
$$O(|P| \log \log \sigma + occ)$$
  
**Space**  $nH_k + n \cdot o(\log \sigma)$  for  $\sigma = \omega(\log^c n)$ 

Optimal time/space complexities?

Time:  $O(|P| \log \sigma / w + occ)$  where w size of a word Space:  $O(n \log \sigma)$
