

Computing All Distinct Squares in Linear Time for Integer Alphabets

Hideo Bannai ¹ Shunsuke Inenaga ¹ *Dominik Köppl* ²

¹Department of Informatics, Kyushu University, Japan

²Department of Computer Science, TU Dortmund, Germany

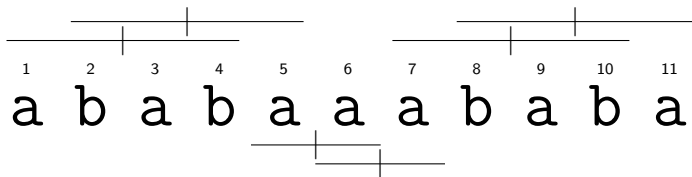
CPM'17

with python a one-liner: `[i**2 for i in range(1,n)]`

squares

1	2	3	4	5	6	7	8	9	10	11
a	b	a	b	a	a	a	b	a	b	a

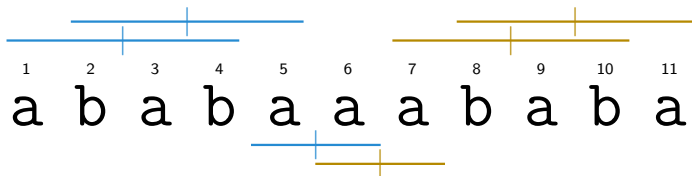
squares



squares

- abab at 1
- baba at 2
- aa at 5
- aa at 6
- abab at 7
- baba at 8

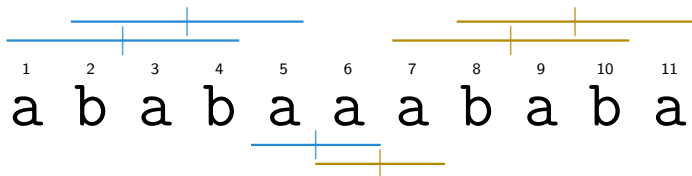
squares



leftmost squares

- abab at 1
- baba at 2
- aa at 5
- ~~aa at 6~~
- ~~abab at 7~~
- ~~baba at 8~~

squares



leftmost squares

- abab at 1
- baba at 2
- aa at 5
- ~~aa at 6~~
- ~~abab at 7~~
- ~~baba at 8~~

fact

leftmost squares \equiv distinct squares

works on distinct squares

#distinct squares

$\leq 2n$	Fraenkel and Simpson'98
$\leq 2n - \Theta(\lg n)$	Ilie'07
$\leq \lfloor 11n/6 \rfloor$	Deza et al.'15
$\leq n$	(yet) unknown

algorithms computing all distinct squares

$\mathcal{O}(\sigma n)$ time	Gusfield and Stoye'04
$\mathcal{O}(n)$ time	Crochemore'14
$\mathcal{O}(n)$ time	this paper

where

- σ : alphabet size
- n : text length

setting

given

- ▀ text T
- ▀ $n := |T|$ text length
- ▀ alphabet of size $\sigma = n^{\mathcal{O}(1)}$

problem

find all distinct squares

goal: $\mathcal{O}(n)$ time

naive solution

- create suffix array and LCP array
- iterate over each text position i
- iterate over all possible periods p
- compare $T[i..]$ with $T[i + p..] \forall 1 \leq i, p \leq n$
- if found a square \Rightarrow check whether already reported

$$\Rightarrow \mathcal{O}\left(\underbrace{n}_{\forall i} \cdot \underbrace{n}_{\forall p} \cdot t_\lambda\right)$$

- t_λ : time for look-up ($t_\lambda = n \lg \sigma$ for a simple trie)

better solutions

idea to get faster

- ▀ check only at certain text position
- ▀ check only periods up to a threshold

sufficient: all borders of Lempel-Ziv factors

idea from

[Gusfield and Stoye'04]

computing all distinct squares in $\mathcal{O}(\sigma n)$ time

Lempel-Ziv parsing

⋮ a b a b a a a b a b a

Lempel-Ziv parsing

⋮ a | b a b a a a b a b a

Lempel-Ziv parsing

a | b | a b a a a b a b a

Lempel-Ziv parsing

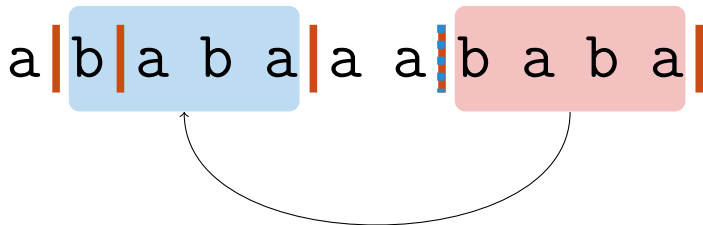


Lempel-Ziv parsing

a | b | a b a a a | b a b a

The diagram shows the string "a b a a a b a b a" with vertical bars separating the characters. A blue box highlights the first 'a' of the fourth segment. A vertical dashed line marks the end of the longest match "a". A purple box highlights the 'a' at the end of the match. A red box highlights the 'a' being added. A curved arrow points from the purple box to the red box.

Lempel-Ziv parsing



Lempel-Ziv parsing

F_1 F_2 F_3 F_4 F_5
a | b | a b a | a a | b a b a |

computing squares

a b a b a a a b a b a

reported squares:

computing squares

F_1 F_2 F_3 F_4 F_5
a | b | a b a | a a | b a b a

reported squares:

computing squares



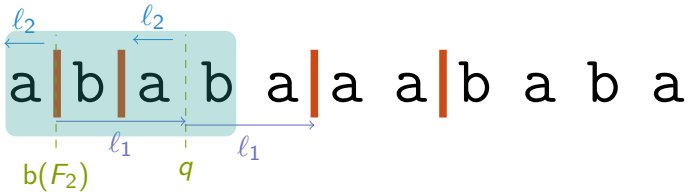
reported squares:

computing squares



reported squares:

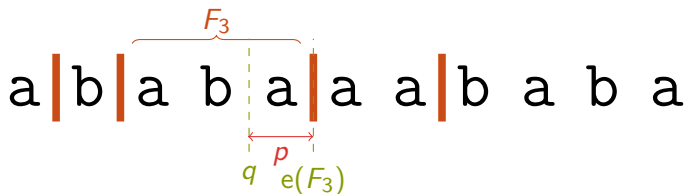
computing squares



reported squares:

▀ abab

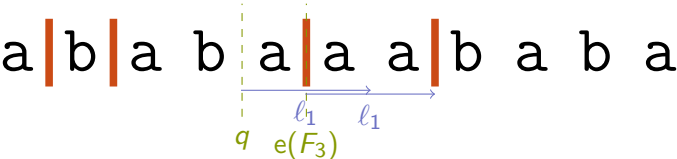
computing squares



reported squares:

▀ abab

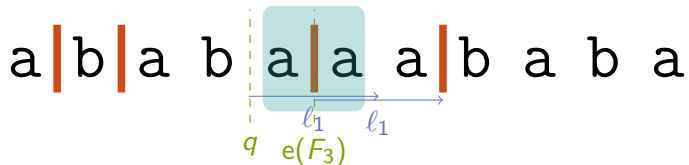
computing squares



reported squares:

- abab

computing squares



reported squares:

- ▀ abab
- ▀ aa

computing squares

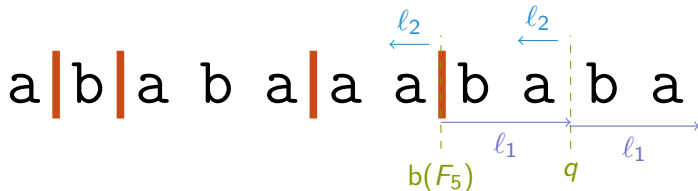


reported squares:

▀ abab

▀ aa

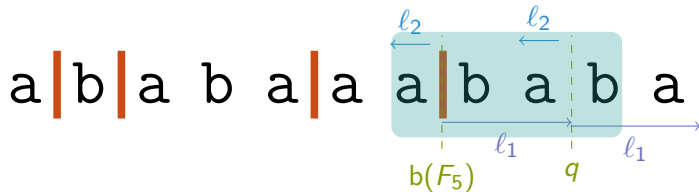
computing squares



reported squares:

- ▀ abab
- ▀ aa

computing squares



reported squares:

- ▀ abab
- ▀ aa
- ▀ abab

computing squares

a | b | a b a | a a | b a b a

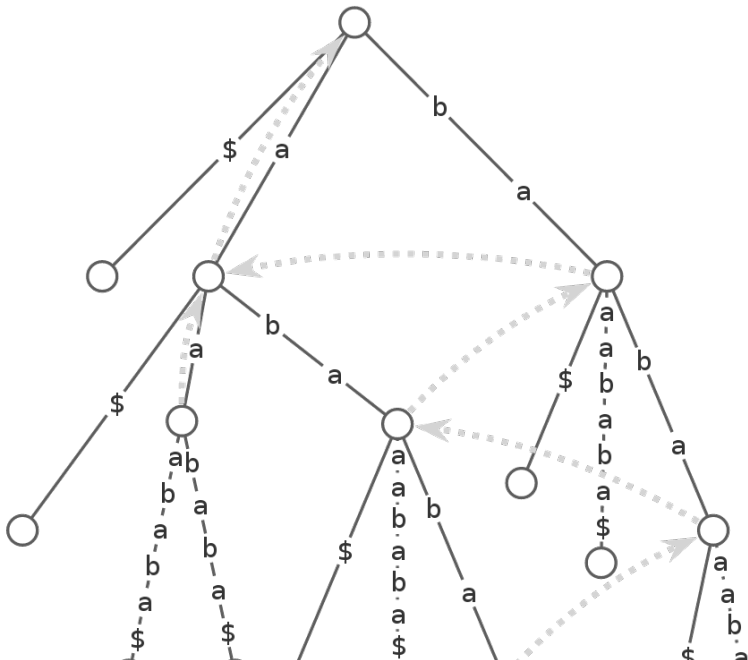
reported squares:

- ▀ abab
- ▀ aa
- ▀ abab

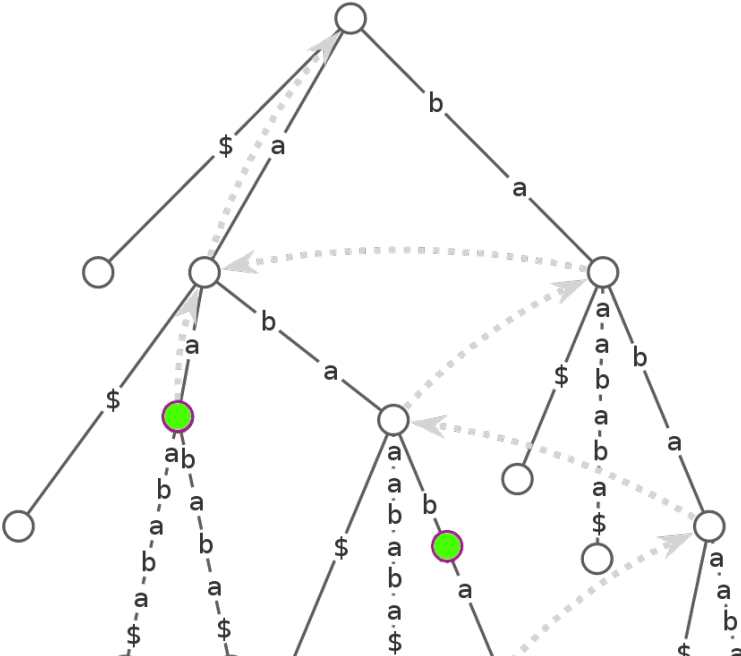
problems:

- ▀ reporting duplicates
- ▀ baba not found

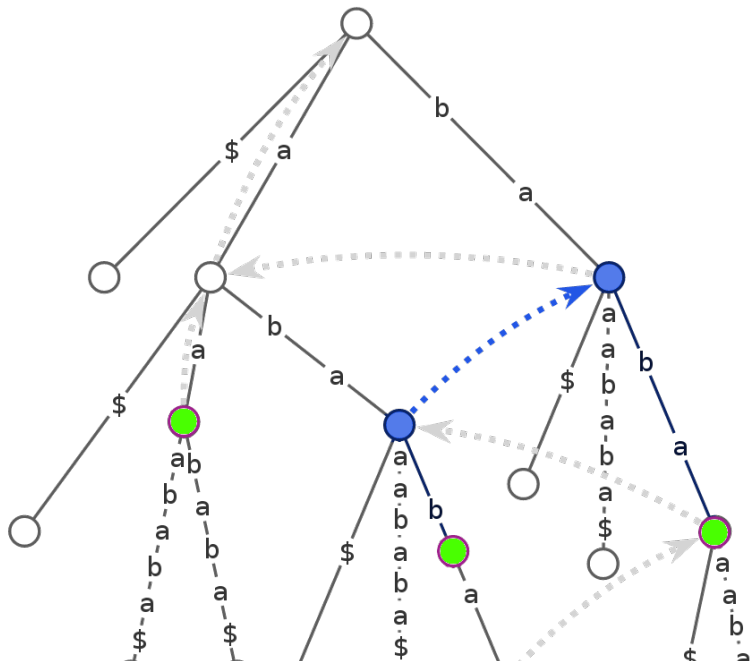
suffix link walk: construct suffix tree



suffix link walk: decorate squares



suffix link walk: rotate



time bounds

- suffix link walk traversals $\mathcal{O}(n)$ nodes
- number of suffix links to a particular node $\leq \sigma$
- \Rightarrow a long traversal can happen $\mathcal{O}(\sigma)$ times

hence $\mathcal{O}(n\sigma)$ time [Gusfield and Stoye'04]

$\mathcal{O}(n)$ time goal

problem

- ▀ \nexists dictionary with $\mathcal{O}(1)$ access/update time
- ▀ store lists, be careful about uniqueness!

solution

use longest previous factor table (LPF)!

longest previous factor table

a b a b a a a b a b a

longest previous factor table

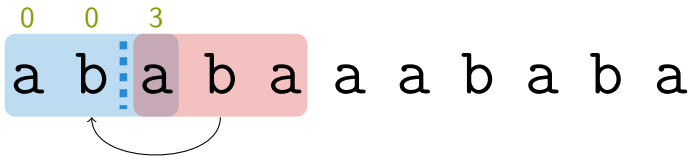
0
a b a b a a a b a b a

The diagram shows the string 'abababababab' with a vertical red bar highlighting the first character 'a'. To the left of this bar are four blue dots, and above it is a yellow '0'. This represents the value of the longest previous factor table at index 0, which is 0 because there are no previous characters to form a factor.

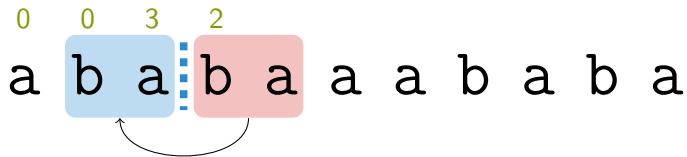
longest previous factor table

0 0
a b a b a a a b a b a

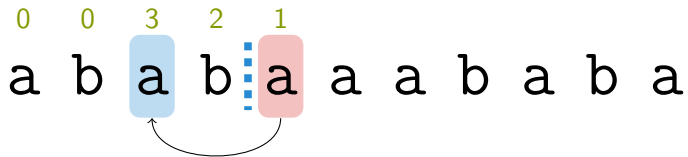
longest previous factor table



longest previous factor table



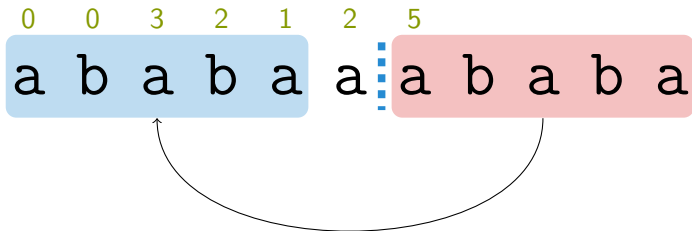
longest previous factor table



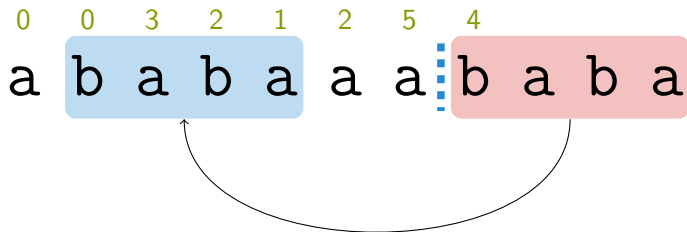
longest previous factor table

0 0 3 2 1 2
a b a b a a a b a b a

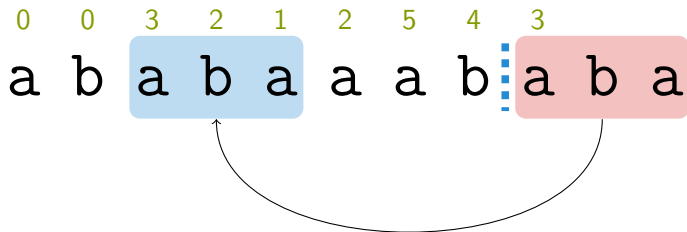
longest previous factor table



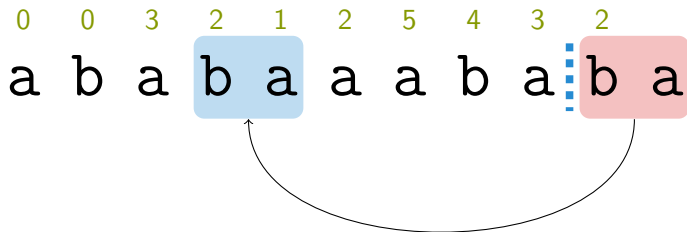
longest previous factor table



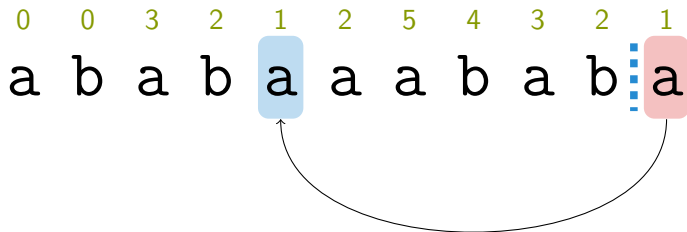
longest previous factor table



longest previous factor table



longest previous factor table

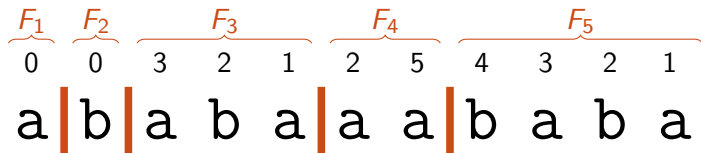


computing squares

0	0	3	2	1	2	5	4	3	2	1
a	b	a	b	a	a	a	b	a	b	a

reported squares:

computing squares



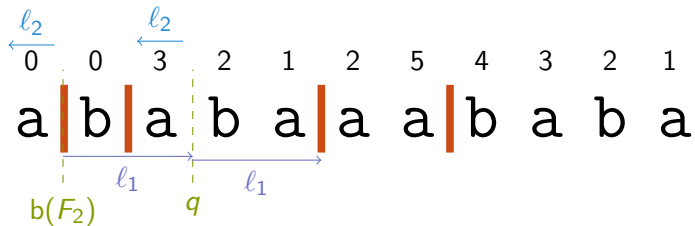
reported squares:

computing squares



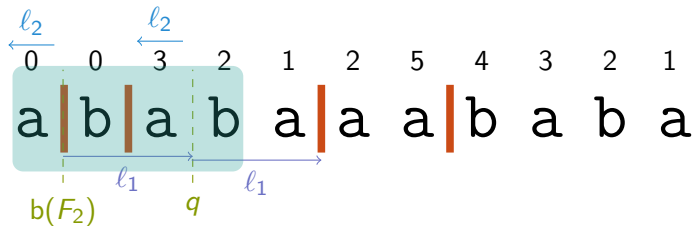
reported squares:

computing squares



reported squares:

computing squares



reported squares:

■ abab

computing squares

0 0 3 2 1 2 5 4 3 2 1
a | b | a b a | a a | b a b a

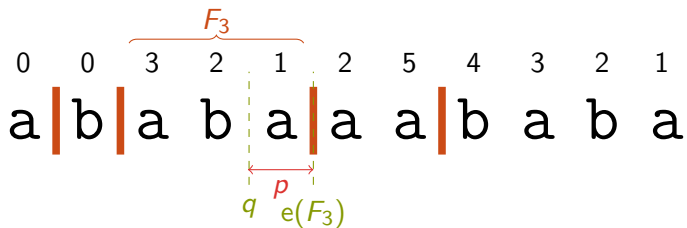
reported squares:

- ▀ abab
- ▀ baba

new techniques:

- ▀ right rotate found squares

computing squares



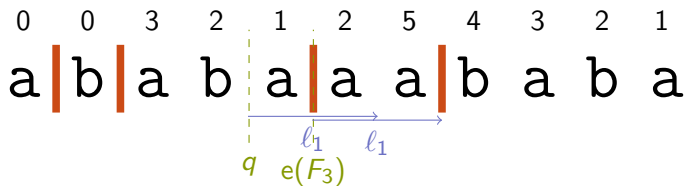
reported squares:

- ▀ abab
- ▀ baba

new techniques:

- ▀ right rotate found squares

computing squares



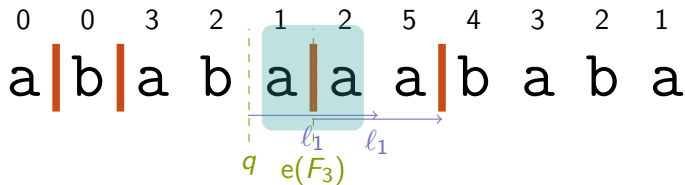
reported squares:

- ▀ abab
- ▀ baba

new techniques:

- ▀ right rotate found squares

computing squares



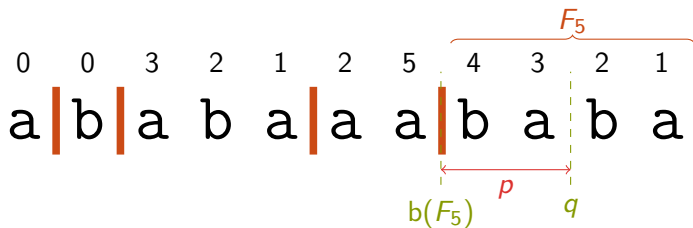
reported squares:

- ▀ abab
- ▀ baba
- ▀ aa

new techniques:

- ▀ right rotate found squares

computing squares



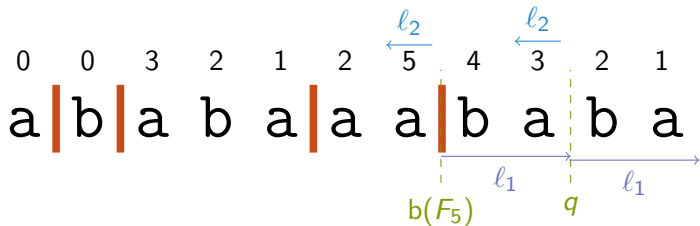
reported squares:

- ▀ abab
- ▀ baba
- ▀ aa

new techniques:

- ▀ right rotate found squares

computing squares



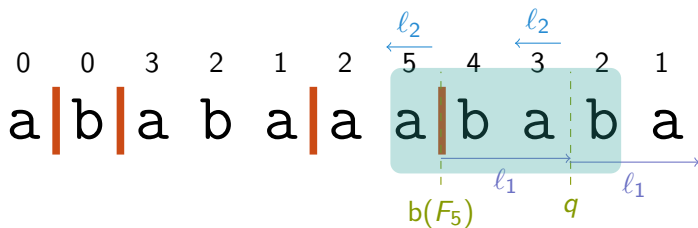
reported squares:

- abab
- baba
- aa

new techniques:

- right rotate found squares

computing squares



reported squares:

- ▀ abab
- ▀ baba
- ▀ aa

new techniques:

- ▀ right rotate found squares
- ▀ skip if $LPF[i] \geq 2p$

RMQs on LPF

naive right rotation of square S takes $\mathcal{O}(|S|)$ time
 $\Rightarrow \mathcal{O}(n^2)$ total time!

our approach

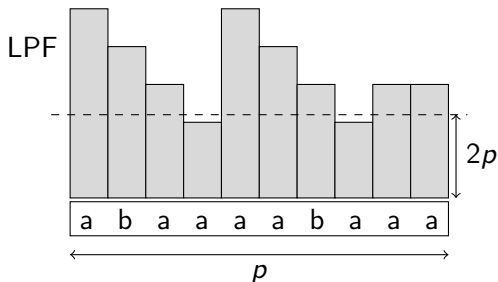
right rotations by RMQ on LPF

RMQs on LPF

naive right rotation of square S takes $\mathcal{O}(|S|)$ time
 $\Rightarrow \mathcal{O}(n^2)$ total time!

our approach

right rotations by RMQ on LPF

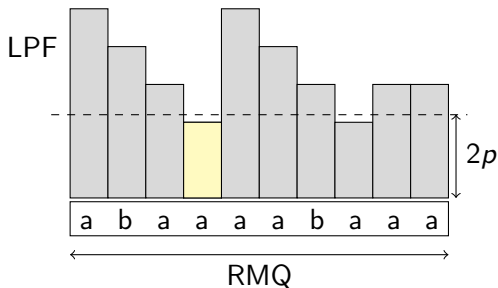


RMQs on LPF

naive right rotation of square S takes $\mathcal{O}(|S|)$ time
 $\Rightarrow \mathcal{O}(n^2)$ total time!

our approach

right rotations by RMQ on LPF

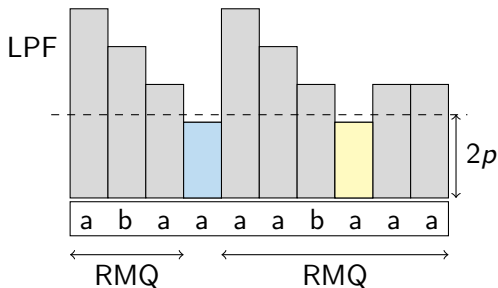


RMQs on LPF

naive right rotation of square S takes $\mathcal{O}(|S|)$ time
 $\Rightarrow \mathcal{O}(n^2)$ total time!

our approach

right rotations by RMQ on LPF

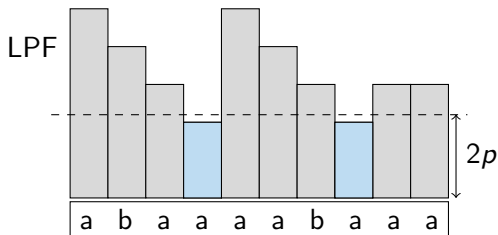


RMQs on LPF

naive right rotation of square S takes $\mathcal{O}(|S|)$ time
 $\Rightarrow \mathcal{O}(n^2)$ total time!

our approach

right rotations by RMQ on LPF



$\#RMQs \leq 3 \times$ number of newly found squares

time analysis

known:

- number of distinct squares $occ = \mathcal{O}(n)$

algorithm:

- build data structures (suffix array, etc.)
- query on $\mathcal{O}(n)$ positions:
 - for a square with LCP
 - find right rotated squares with RMQ

time

$\mathcal{O}(n)$

$\mathcal{O}(1)$

$\mathcal{O}(occ)$ total

$\mathcal{O}(n + occ) = \mathcal{O}(n)$ total time

experiments

collection	σ	z	$\max_x F_x $	$ \text{occ} $	time
dblp.xml	97	7,035,342	1060	7412	70
proteins	26	20,875,097	45,703	3,108,339	245
dna	17	13,970,040	97,966	132,594	310
english	226	13,971,134	987,766	13,408	2639
einstein	125	49,575	906,995	18,192,737	3953

- 200 MiB collections from Pizza&Chili corpus
- σ : alphabet size
- z : # Lempel-Ziv factors
- time in seconds

time bottleneck: RMQs

finding all distinct squares in linear time

techniques

- modification of [Gusfield and Stoye'04]
- using LPF array

open problems

- cope without RMQ
- create MAST in $\mathcal{O}(n)$ time

further results

- computing online in $\mathcal{O}\left(n \lg^2 \lg n / \lg \lg \lg n\right)$ time
- in linear time:
 - decorating suffix tree with information of all squares
 - building topology of the minimal augmented suffix tree (MAST)

finding all distinct squares in linear time

techniques

- modification of [Gusfield and Stoye'04]
- using LPF array

open problems

- cope without RMQ
- create MAST in $\mathcal{O}(n)$ time

further results

- computing online in $\mathcal{O}\left(n \lg^2 \lg n / \lg \lg \lg n\right)$ time
- in linear time:
 - decorating suffix tree with information of all squares
 - building topology of the minimal augmented suffix tree (MAST)

Thank you for listening. Any questions are welcome!