

# Towards a solution to the “runs” conjecture

MAXIME CROCHEMORE<sup>1</sup>, LUCIAN ILIE<sup>2</sup>, LIVIU TINTA<sup>2</sup>

<sup>1</sup>King’s College London, UK

<sup>2</sup>University of Western Ontario, CANADA

# Repetitions

- **repetition** = substring with adjacent occurrences
- **square** – two occurrences
  - ATAGGATAGG = (ATAGG)<sup>2</sup>
  - hotshots = (hots)<sup>2</sup>
- **general (fractional) repetitions**
  - alfalfa = (alf) <sup>$\frac{7}{3}$</sup>
  - otavaltavaltavaltavaltiolta = o(taval) <sup>$\frac{21}{5}$</sup>  iolta
  - CGATATATATATATACGGC = CG(AT) <sup>$\frac{13}{2}$</sup>  CGGC
- for a repetition  $w^\alpha$ 
  - $\alpha$  is the **exponent**
  - the number of letters of  $w$  is the **period**

# Repetitions – Problems

- interesting **combinatorial problems**
  - count the number of various types of repetitions
  - squares already investigated by [Thue, 1906, 1912]
- useful **combinatorial algorithms**
  - find various types of repetitions
  - find **all** repetitions in **linear time**
- many **applications**
  - text algorithms
  - data compression
  - analysis of biological sequences

# Finding all repetitions

- **Goal: find all repetitions in linear time**
  - problem – too many repetitions
  - solution – encode all more compactly
- [Crochemore, 1981, 1983] – linear for one square,  $\mathcal{O}(n \log n)$  for all primitively-rooted maximal integer rep.
- [Apostolico, Preparata, 1983] –  $\mathcal{O}(n \log n)$  for all right-maximal
- [Main, Lorentz, 1985] –  $\mathcal{O}(n \log n)$  for all maximal
- [Main, 1989] – linear for all maximal leftmost
- [Iliopoulos, Moore, Smyth, 1997] – linear for all in Fibonacci

# Runs

- **run** – repetition that is
  - fractional
  - non-extendable (maximal)
  - primitively-rooted
- $o\underline{taval}taval\underline{taval}taval\underline{taval}iota = o(taval)^{\frac{21}{5}}iota$
- $CG\underline{ATATATATATATA}CGGC = CG(AT)^{\frac{13}{2}}CGGC$
- $00\underline{0101}000010 = 00(01)^{\frac{5}{2}}00010$
- $\underline{000}101\underline{0000}10 = 0^31010^410$

# Runs – linear bound

- all repetitions are encoded in runs
- number of runs in a string of length  $n$ 
  - $\mathcal{O}(n)$  – [Kolpakov, Kucherov, 1998]
  - compute all repetitions in linear time
    - modified Main’s algorithm
    - based on Lempel–Ziv factorization – LZ77
    - most important breakthrough – the bound
  - no constant could be derived from the proof
- “runs” conjecture
  - $\text{runs}(n) \leq n$  – supported by computations
- [Rytter, 2006] – first bound:  $5n$

# Runs – simple proof

- Kolpakov and Kucherov's proof – very complicated
  - 8 (large) pages
  - contains case 2.1.2.2
- Rytter's proof – 9 pages
  - highly and weakly periodic runs
- [\[Crochemore, Ilie, 2006\]](#) – simple proof
  - 1.5 pages
  - improved Rytter's idea of “neighbor runs”

# Runs – sum of exponents

- $\mathcal{O}(n)$  – [Kolpakov, Kucherov, 1999]
  - sum of exponents – applications to analysis of algorithms
  - proof very complicated – 9 (large) pages, contains case 2.1.2.3.2
- [Crochemore, Ilie, 2006] – simple proof – 0.5 pages



# Number of runs – upper bounds

- analysis of any algorithm computing all repetitions
- upper bounds
  - [Rytter, 2006] –  $5n$
  - [Puglisi, Simpson, Smyth, 2006] –  $3.48n$
  - [Rytter, 2006] –  $3.44n$
  - [Crochemore, Ilie, 2006] –  $1.6n$

# Sum of exponents – upper bounds

- conjecture
  - [Kolpakov, Kucherov, 1999] –  $\leq 2n$
- [Crochemore, Ilie, 2006] –  $5.6n$ 
  - could be improved by computer verification to  $2.9n$  or lower
  - the first explicit bound for the sum of exponents
  - from Rytter's paper –  $25n$  – “unsatisfactory”

# Proof ideas

- [Rytter, 2006]
  - runs counted at the start – logarithmically many
  - short period runs
  - long period runs
    - weakly periodic
    - highly periodic

abaababaabaababaabaabaabaababaaba

# Proof ideas

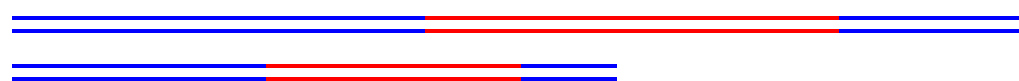
- [Rytter, 2006]
  - runs counted at the start – logarithmically many
  - short period runs
  - long period runs
    - weakly periodic
    - highly periodic

abaababaabaababaabaababaabaa

# Proof ideas

- [Rytter, 2006]
  - runs counted at the start – logarithmically many
  - short period runs
  - long period runs
    - weakly periodic
    - highly periodic

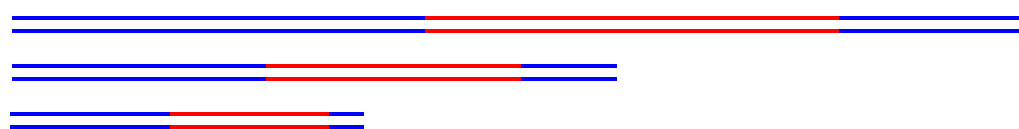
abaababaabaababaabaabaabaabaaba



# Proof ideas

- [Rytter, 2006]
  - runs counted at the start – logarithmically many
  - short period runs
  - long period runs
    - weakly periodic
    - highly periodic

abaababaabaababaabaabaabaababaaba



# Proof ideas

- [Rytter, 2006]
  - runs counted at the start – logarithmically many
  - short period runs
  - long period runs
    - weakly periodic
    - highly periodic

abaababaabaababaabaabaabaaba



# Proof ideas (2)

- [Crochemore, Ilie, 2006]
  - runs counted at the center – linearly many
    - center = beginning of second period
  - short period runs (microruns)
  - long period runs – 10 times better bound
  - for periods  $\geq 87$ 
    - we have  $0.06897n$  compared to  $0.67n$

abaabaabaabaabaabaabaabaabaabaaba



# Proof ideas (2)

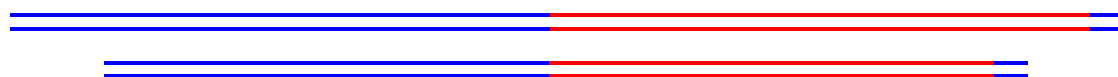
- [Crochemore, Ilie, 2006]
  - runs counted at the center – linearly many
    - center = beginning of second period
  - short period runs (microruns)
  - long period runs – 10 times better bound
  - for periods  $\geq 87$ 
    - we have  $0.06897n$  compared to  $0.67n$

abaabaabaabaabaabaabaabaabaabaaba

# Proof ideas (2)

- [Crochemore, Ilie, 2006]
  - runs counted at the center – **linearly many**
  - center = beginning of second period
  - short period runs (**microruns**)
  - long period runs – 10 times better bound
  - for periods  $\geq 87$
  - we have  $0.06897n$  compared to  $0.67n$

abaabaabaabaabaabaabaabaabaabaaba

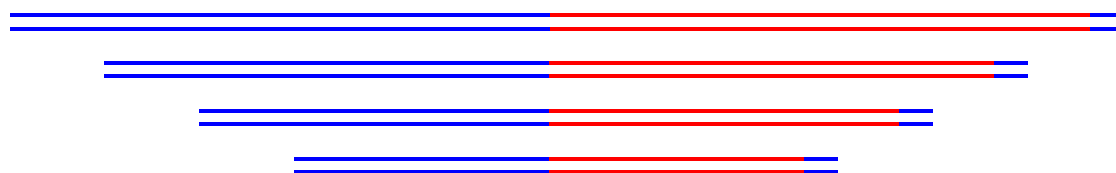




# Proof ideas (2)

- [Crochemore, Ilie, 2006]
  - runs counted at the center – linearly many
    - center = beginning of second period
  - short period runs (microruns)
  - long period runs – 10 times better bound
  - for periods  $\geq 87$ 
    - we have  $0.06897n$  compared to  $0.67n$

abaabaabaabaabaabaabaabaabaabaaba



# Proof ideas (2)

- [Crochemore, Ilie, 2006]
  - runs counted at the center – linearly many
    - center = beginning of second period
  - short period runs (microruns)
  - long period runs – 10 times better bound
  - for periods  $\geq 87$ 
    - we have  $0.06897n$  compared to  $0.67n$

abaabaabaabaabaabaabaabaabaabaaba  
=====

=====

=====

=====

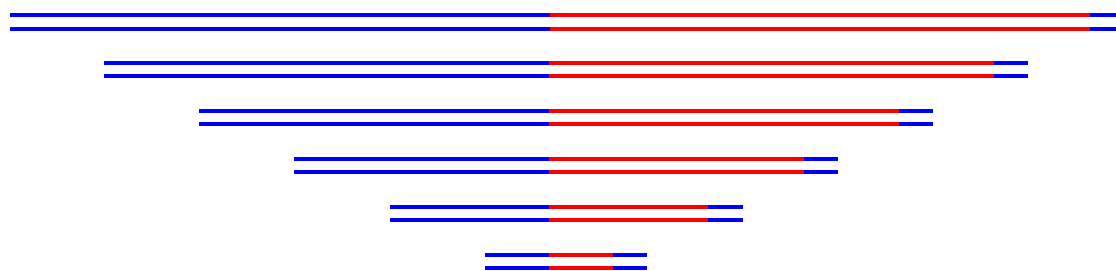
=====

=====

# Proof ideas (2)

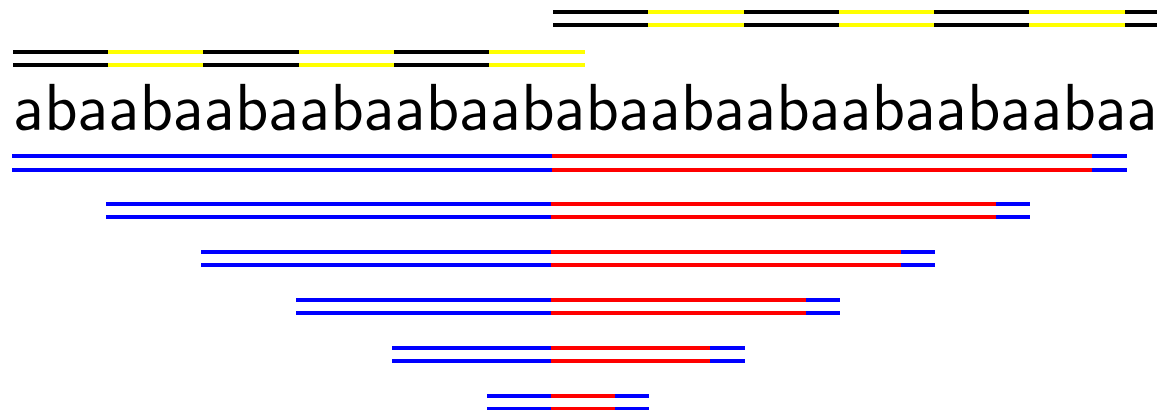
- [Crochemore, Ilie, 2006]
  - runs counted at the center – linearly many
    - center = beginning of second period
  - short period runs (microruns)
  - long period runs – 10 times better bound
  - for periods  $\geq 87$ 
    - we have  $0.06897n$  compared to  $0.67n$

abaabaabaabaabaabaabaabaabaabaaba



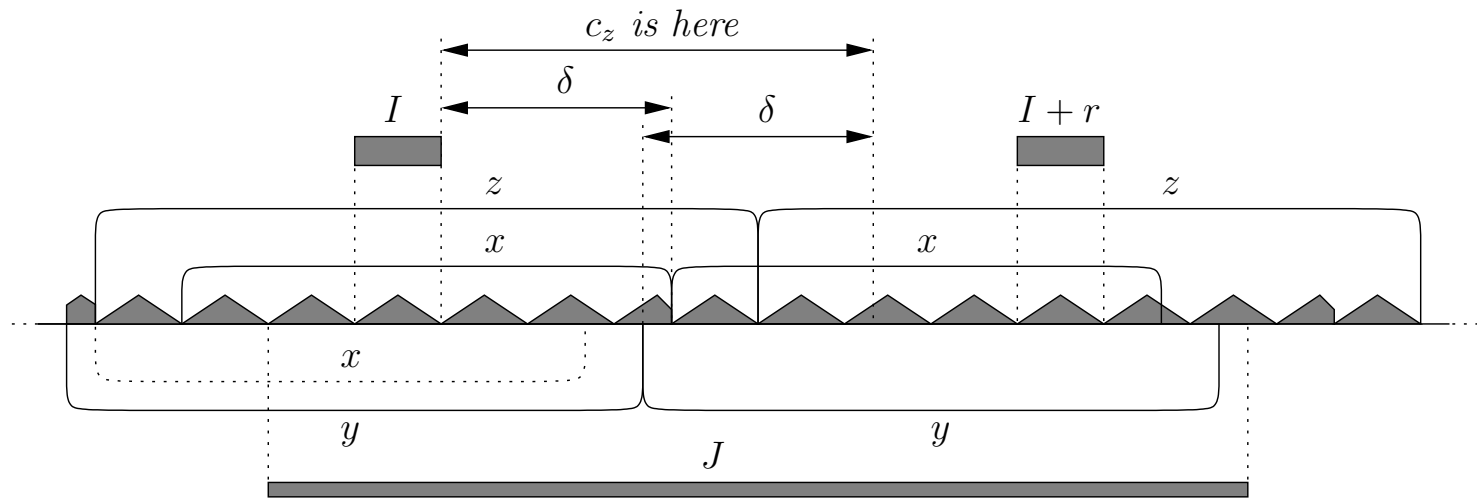
# Proof ideas (2)

- [Crochemore, Ilie, 2006]
  - runs counted at the center – linearly many
    - center = beginning of second period
  - short period runs (microruns)
  - long period runs – 10 times better bound
  - for periods  $\geq 87$ 
    - we have  $0.06897n$  compared to  $0.67n$



# Short periods need separate reasoning

- runs with close centers imply high local periodicity
- for short periods it is not possible





# Long and short periods

- [Crochemore, Ilie, 2006]

- runs with period  $p$  or larger

- **Theorem 1**  $\text{runs}_{\geq p}(n) \leq \frac{6}{p}n$

- runs with period  $p$  or smaller (**microruns**)

- **Theorem 2**  $\text{runs}_{\leq p}(n) \leq bn$

- **Bound:**  $\text{runs}(n) \leq \left(\frac{6}{p+1} + b\right)n$

- $p = 9, b = 1$  gives  $\text{runs}(n) \leq 1.6n$

- higher  $p$  and/or lower  $b$  give better bounds

# Microruns – better upper bounds

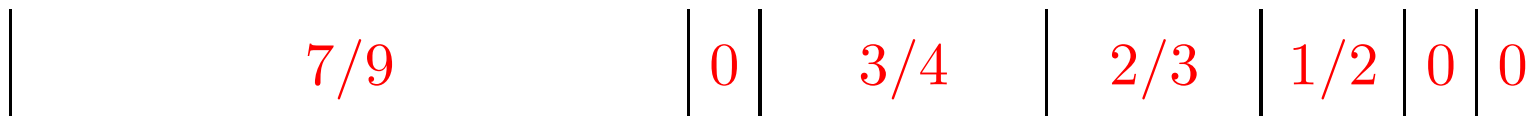
- centers of microruns – not uniformly distributed
- amortized for large enough intervals
- example:  $p = 8, b = 0.8$

0 1 0 0 1 0 1 0 0 1 0 0 1 0 1 0 0 1 0 1 0

1 2 1 1 2 1 2

3 5 3 5 3

8



# Microruns – better upper bounds (2)

- given  $p$  and  $b$ , try all possibilities
  - if all are amortized, then  $\text{runs}_{\leq p}(n) \leq bn$
- trying all strings
  - for arbitrary alphabets – not clear how to do it
  - for fixed alphabets – too much to compute
    - 2 letters, length of amortizing interval 100:  $2^{100}$
- trying all combinations of periods
  - works for arbitrary alphabets
  - too much to compute
    - $p = 100$ , length of amortizing interval 100:  $2^{10000}$
  - **solution: try only the possible combinations**

# Microruns – better upper bounds (3)

- example: if  $\{1, 6\}$  at some position  $i$  in a string  $s$  and  $p = 12$ , then some periods are forbidden

$s_{i-6}$	$s_{i-5}$	$s_{i-4}$	$s_{i-3}$	$s_{i-2}$	$s_{i-1}$	$s_i$
					<del>1</del>	1
					<del>2</del>	<del>2</del>
			<del>3</del>	<del>3</del>	<del>3</del>	<del>3</del>
				<del>4</del>	<del>4</del>	<del>4</del>
					<del>5</del>	<del>5</del>
<del>6</del>	<del>6</del>	<del>6</del>	<del>6</del>	<del>6</del>	<del>6</del>	6
					<del>7</del>	<del>7</del>
				<del>8</del>	<del>8</del>	<del>8</del>
			<del>9</del>	<del>9</del>	<del>9</del>	<del>9</del>
					<del>10</del>	<del>10</del>
					<del>11</del>	<del>11</del>
<del>12</del>	<del>12</del>	<del>12</del>	<del>12</del>	<del>12</del>	<del>12</del>	<del>12</del>

# Microruns – better upper bounds (4)

- all possible sets of periods at the same position for  $p = 9$

- $\emptyset,$   
 $\{1\}, \dots, \{9\},$   
 $\{1, 3\}, \dots, \{1, 9\}, \{2, 5\}, \dots, \{2, 9\}, \{3, 7\}, \dots, \{3, 9\}, \{4, 9\}, \{5, 8\},$   
 $\{1, 3, 5\}, \{1, 3, 7\}, \dots, \{1, 3, 9\}, \{1, 4, 7\}, \{1, 4, 9\}, \{1, 5, 9\}, \{2, 5, 8\},$   
 $\{1, 3, 5, 7\},$   
 $\{1, 3, 5, 7, 9\}$

- there are 37 sets out of potential  $2^9 = 512$

# Current best bound

- obtained using SHARCNET – [www.sharcnet.ca](http://www.sharcnet.ca)

$p$	$b$	solutions	amortize	$\text{runs}(n) \leq$	CPU time (s)
10	0.85	900	100	$1.396n$	40
15	0.89	5275	27	$1.265n$	600
20	0.89	34833	97	$1.176n$	1920
25	0.91	135457	153	$1.141n$	8640
30	0.91	471339	153	$1.104n$	50,400
35	0.93	1455422	82	$1.097n$	230,400
40	0.93	3907110	84	$1.077n$	1,130,400
50	0.93	22635894	139	$1.048n$	27,388,800

# Number of runs – lower bounds

- [Franek, Simpson, Smyth, 2003]
  - $\frac{3}{2\phi}n = 0.927..n$
  - conjectured optimal
- [Kusano, Matsubara, Ishino, Bannai, Shinohara, 2008]
  - $0.944542n$

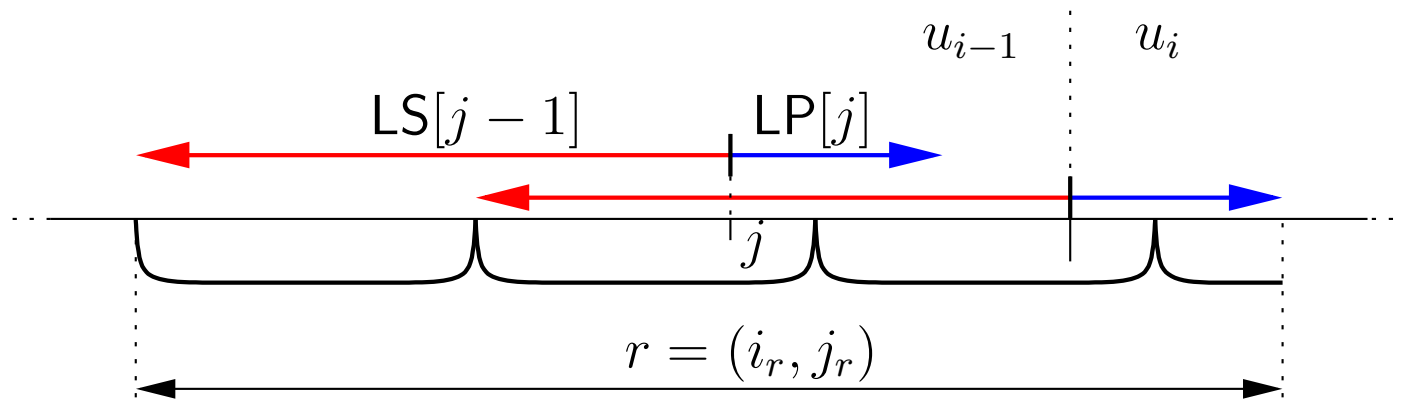
# Further research

- $0.944n \leq \text{runs}(n) \leq 1.048n$
- good enough for practical purposes
- the optimal value seems to be below  $n$ 
  - (the “runs” conjecture is probably true)
- average number of runs
  - [Puglisi, Simpson, 2008] – also linear
- algorithms for computing runs
  - only one idea – use Lempel–Ziv factorization
  - (same for linear on average)
  - different (simpler) algorithms
  - could bring new insight into the structure of strings



# Linear-time algorithm idea

- Lempel–Ziv factorization – `a.b.b.a.abb.baa.ab.ab`
- runs that cross an LZ-factorization point



- those which do not are only translated