# Repetition- and linearity-aware rank/select dictionaries

Paolo Ferragina        Giovanni Manzini        Giorgio Vinciguerra

UNIVERSITÀ DI PISA

(ISAAC 2021)

# Compressed rank/select dictionaries

- Given a set $A$ of $n$ elements over an integer universe $0, 1, \dots, u$
    1. Store them in compressed form
    2. Implement $rank(x)$: number of elements in $A$ which are $\leq x$
    3. Implement $select(i)$: return the $i$th smallest element in $A$

- Well-studied building block of succinct data structures

$rank(12) = 3$

| 3 | 6 | 10 | 15 | 18 | 22 | 40 | 43 | 47 | 53 |
|---|---|----|----|----|----|----|----|----|----|
| 1 | 2 | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 |

$\Big\}$ Stored in compressed form

$select(7) = 40$

# Two sources of compressibility

We exploit them both

## Repetitiveness

## Approximate linearity

[Boffa et al., ALENEX '21]



$A$ | 2 | 3 | 5 | 6 | 13 | 14 | 16 | 17 | 20 | 24 | …

Difference between adjacent values

Gap string | 2 | 1 | 2 | 1 | 7 | 1 | 2 | 1 | 3 | 4 | …

Store just a "back reference"

Many nonlinear points
$\Downarrow$
Use piecewise linear $\varepsilon$-approx.

$\mathcal{O}(\log \varepsilon)$ bits

small bit size

Corrections | 3 | 1 | 0 | 0 | -2 | -3
1 2 3 4 5 6

Errors smaller than a given integer $\varepsilon$

$A$ | 3 | 6 | 10 | 15 | 18 | 22
1 2 3 4 5 6

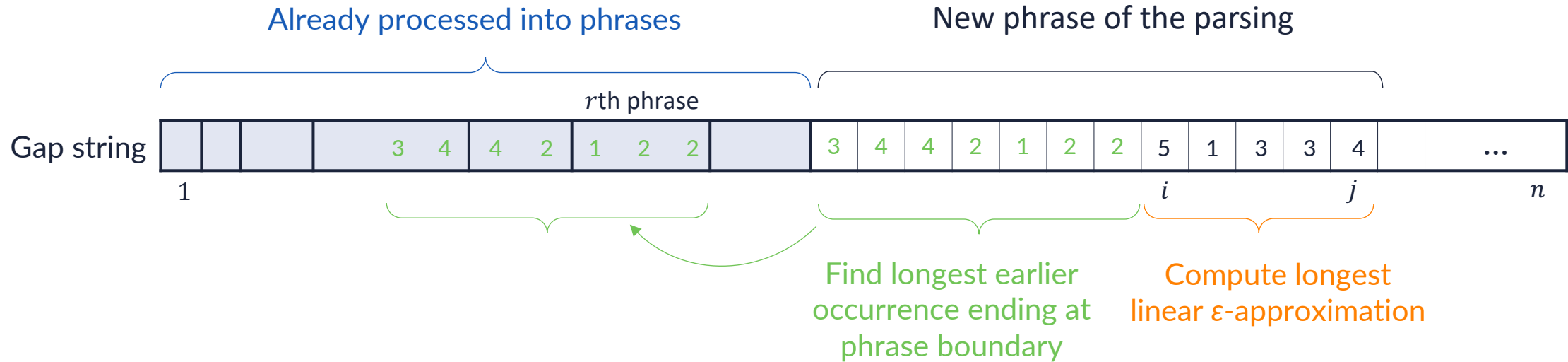# Exploiting repetitiveness and approx. linearity

1. Build on two known repetition-aware methods

   - Lempel-Ziv parsing, LZ-End [Kreft and Navarro, TCS 2013]

   - Block tree [Belazzougui et al., JCSS 2021]

2. Augment them to use linear $\varepsilon$-approximations with corrections

3. Show how to support $rank$ and $select$ in space bounded by the high-order entropy or a repetitiveness measure of the gaps

# The LZ$_\varepsilon$ parsing

Already processed into phrases

New phrase of the parsing

Gap string

| | | | | 3 | 4 | 4 | 2 | 1 | 2 | 2 | | 3 | 4 | 4 | 2 | 1 | 2 | 2 | 5 | 1 | 3 | 3 | 4 | | | ... |

1
$i$
$n$

Find longest earlier occurrence ending at phrase boundary

Compute longest linear $\varepsilon$-approximation



value

$A[i+4]$

$A[i+3]$

$A[i+2]$

$A[i+1]$

$\mathcal{O}(\log \varepsilon)$ bits

Corrections  | 1 | -1 | -1 | -1 | 1 |

$A[i]$

$i$  $i+1$  $i+2$  $i+3$  $i+4$

position

# The LZ$_\varepsilon$ parsing

Already processed into phrases

New phrase of the parsing

$r$th phrase

Gap string

| | | | | 3 | 4 | 4 | 2 | 1 | 2 | 2 | | 3 | 4 | 4 | 2 | 1 | 2 | 2 | 5 | 1 | 3 | 3 | 4 | | | ... |

1            $i$       $j$       $n$

Find longest earlier occurrence ending at phrase boundary
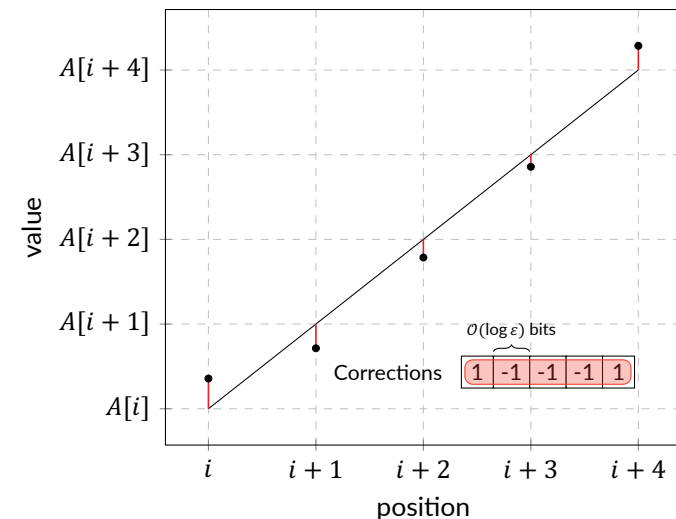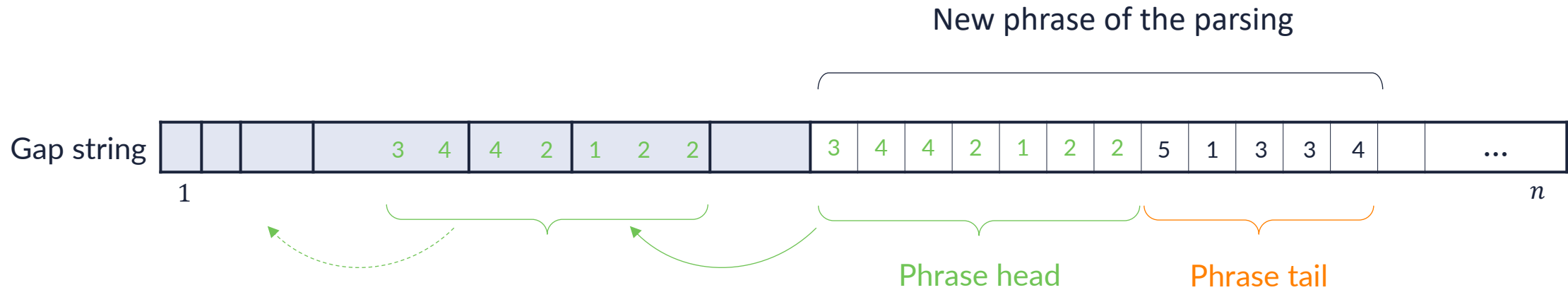
Compute longest linear $\varepsilon$-approximation

For the new phrase we store

- Indexes $i$, $j$, and $r$
- Slope and intercept of the line
- Array of $j - i + 1$ corrections, $\mathcal{O}(\log \varepsilon)$ bits each



value

$A[i+4]$

$A[i+3]$

$A[i+2]$

$A[i+1]$

$A[i]$

$\mathcal{O}(\log \varepsilon)$ bits

Corrections | 1 | -1 | -1 | -1 | 1 |

$i$   $i+1$   $i+2$   $i+3$   $i+4$

position

# Queries in the $\text{LZ}_\varepsilon$ parsing

New phrase of the parsing

Gap string

| | | | | 3 | 4 | 4 | 2 | 1 | 2 | 2 | | 3 | 4 | 4 | 2 | 1 | 2 | 2 | 5 | 1 | 3 | 3 | 4 | | | ... | |

1          $n$

Phrase head          Phrase tail

If $t \in$ Phrase head, then we must recursively unroll the source phrase

If $t \in$ Phrase tail, then we can answer directly with the linear $\varepsilon$-approx

$\text{LZ}_\varepsilon^\rho$: Introduce a trade-off parameter $\rho > 0$ to shorten the phrase head and make queries faster

$select(t)$
$rank(A[t])$

# $\mathbf{LZ}_\varepsilon^\rho$ **bounds**

→ No worse than a traditional LZ-parsing

→ No worse than LA-vector in space

Let $\sigma$ = number of distinct values in the gap string

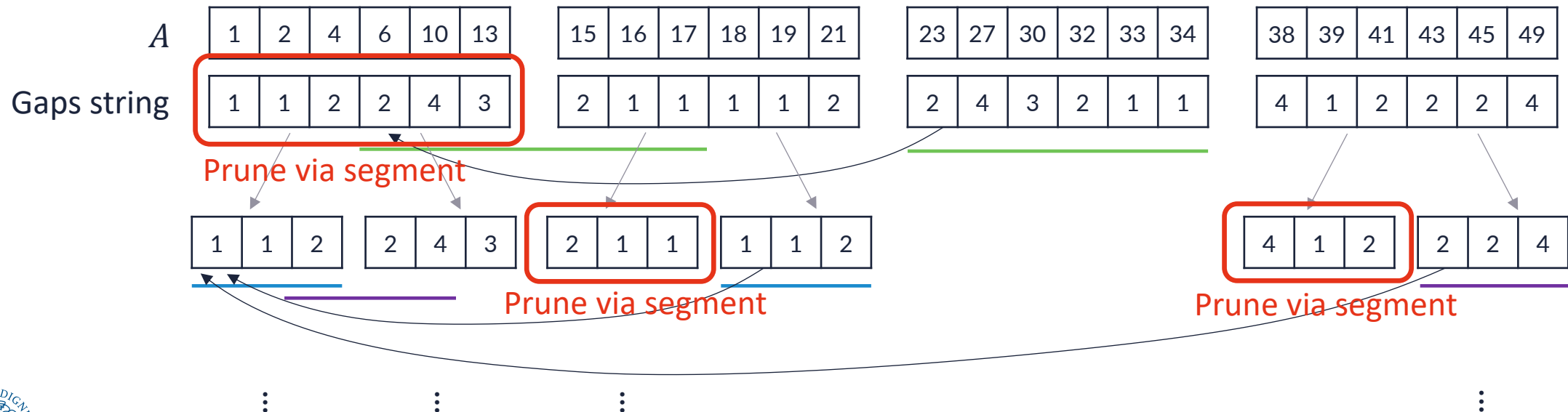| | |
|---|---|
| **Select time** | $\mathcal{O}(\log^{1+\rho} n)$ |
| **Rank time** | $\mathcal{O}(\log^{1+\rho} n + \log \varepsilon)$ |
| **Space in bits** | $nH_k(\text{gap string}) + \mathcal{O}(n/\log^\rho n) + o(n \log \sigma) + \text{space for tails}$ |

Exploit repetitions      Exploit approximate linearity

# The block-$\varepsilon$ tree

- Start with a standard block tree construction on the gap string

$A$

| 1 | 2 | 4 | 6 | 10 | 13 | 15 | 16 | 17 | 18 | 19 | 21 | 23 | 27 | 30 | 32 | 33 | 34 | 38 | 39 | 41 | 43 | 45 | 49 |
|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|

Gap string

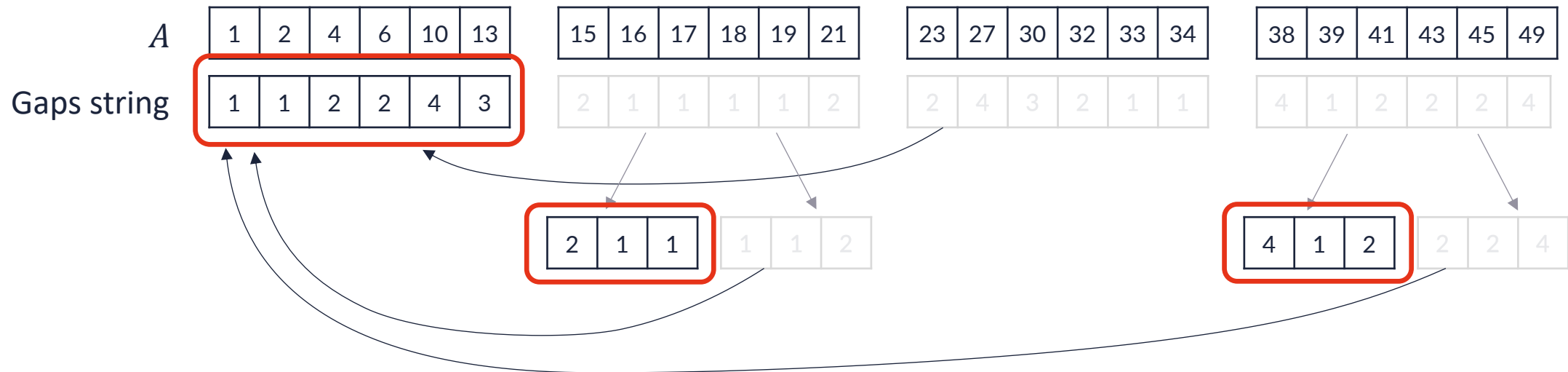| 1 | 1 | 2 | 2 | 4 | 3 | 2 | 1 | 1 | 1 | 1 | 2 | 2 | 4 | 3 | 2 | 1 | 1 | 4 | 1 | 2 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# The block-ε tree

- Start with a standard block tree construction on the gap string
- Assign to each node the bit cost of encoding its subtree
- Prune subtrees that are better compressed by linear ε-approximations



Prune via segment

Prune via segment

Prune via segment

# The block–ε tree

- Start with a standard block tree construction on the gap string
- Assign to each node the bit cost of encoding its subtree
- Prune subtrees that are better compressed by linear ε-approximations
- Store topology, leaf linear ε-approx., and left pointers of copied blocks

# Block-ε tree bounds

- Based on the $\delta$ repetitiveness measure on strings:[1,2,3]

$$\delta = \max\{d_k/k : k = 1, \dots, n\}$$
where $d_k$ = number of distinct substrings of length $k$ in the gap string

- Number of levels is $h = \mathcal{O}\left(\log \frac{n}{\delta}\right)$

| | |
|---|---|
| **Select time** | $\mathcal{O}(h)$ |
| **Rank time** | $\mathcal{O}\left(\log\log \frac{u}{\delta} + h + \log \varepsilon\right)$ |
| **Space in bits** | $\mathcal{O}\left(\delta \log \frac{u}{\delta} \log u\right)$ |

[1] Raskhodnikova et al., Algorithmica (2013)
[2] Christiansen et al., TALG (2020)
[3] Kociumaka et al., LATIN '20

# Experiments with the block-ε tree

- Compared with LA-vector, and a standard block tree

- Datasets: postings lists, positions of symbols in texts (DNA, URLs)

- LA-vector is 10.5× faster in select and 4.7× faster in rank than block tree, but no clear winner in space → Combination of repetitiveness and approximate-linearity makes sense

- Our block-ε tree:
  - wrt LA-vector, it is always slower in select and in rank
  - wrt block tree, it is 2.2× faster in select, either faster (1.3×) or slower (1.3×) in rank
  - has the best space in 2/12 datasets, and the second-best space in 7/12 datasets

  Block-ε tree achieves a good compromise by exploiting both regularities

# Conclusions

- Exploit both repetitiveness and approx. linearity in rank/select dictionaries
- $\text{LZ}_\varepsilon^\rho$ parsing
  - Combine backward copies and linear $\varepsilon$-approximations
  - Space complexity bounded by the $k$th order entropy
- Block-$\varepsilon$ tree
  - Optimise block tree by compressing areas with high approximate linearity
  - Space-time bounds based on the $\delta$ repetitiveness measure
  - Experimentally achieves a good compromise between block trees and LA-vectors

- Future work
  - Implement $\text{LZ}_\varepsilon^\rho$
  - Relation of approximate linearity with other compressibility measures