# ALGORITHMS OF BIOINFORMATICS

# 7

## Googling Genomes

*18 December 2025*

Prof. Dr. Sebastian Wild

# Outline

# 7 Googling Genomes

# Recall Unit 6

## Application 4: Longest Common Extensions

▶ We implicitly used a special case of a more general, versatile idea:

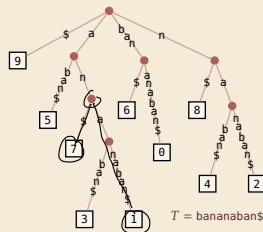The *longest common extension (LCE)* data structure:

   ▶ **Given:** String $T[0..n]$

   ▶ **Goal:** Answer LCE queries, i.e.,
          given positions $i, j$ in $T$,
          how far can we read the same text from there?
          formally: $\text{LCE}(i, j) = \max\{\ell : T[i..i + \ell] = T[j..j + \ell]\}$

⤳ use suffix tree of $T$!

(length of) longest common prefix
of $i$th and $j$th suffix

▶ In $\mathcal{T}$: $\text{LCE}(i, j) = \text{LCP}(T_i, T_j)$ ⤳ same thing, different name!
                   $=$ string depth of
                     *lowest common ancester (LCA)* of
                     leaves $\boxed{i}$ and $\boxed{j}$

▶ in short: $\boxed{\text{LCE}(i, j) = \text{LCP}(T_i, T_j) = \text{stringDepth}(\text{LCA}(\boxed{i}, \boxed{j}))}$



$T = \text{bananaban\$}$

17

# Recall Unit 6

## Efficient LCA

How to find lowest common ancestors?

- ▶ Could walk up the tree to find LCA $\leadsto$ $\Theta(n)$ worst case 👎
- ▶ Could store all LCAs in big table $\leadsto$ $\Theta(n^2)$ space and preprocessing 👎

**Amazing result:** Can compute data structure in $\Theta(n)$ time and space that finds any LCA in **constant(!) time**.

- ▶ a bit tricky to understand
- ▶ but a theoretical breakthrough
- ▶ and useful in practice
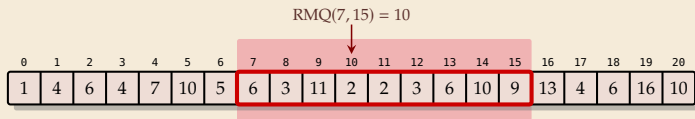
and suffix tree construction inside . . .

$\leadsto$ for now, use $O(1)$ LCA as black box.

$\leadsto$ After linear preprocessing (time & space), we can find LCEs in $O(1)$ time.

18

2

## 7.1 Range-Minimum Queries

# Range-minimum queries (RMQ)

array/numbers don't change

▶ **Given:** Static array $A[0..n)$ of numbers

▶ **Goal:** Find minimum in a range;
    $A$ known in advance and can be preprocessed

RMQ(7, 15) = 10

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| 1 | 4 | 6 | 4 | 7 | 10 | 5 | 6 | 3 | 11 | 2 | 2 | 3 | 6 | 10 | 9 | 13 | 4 | 6 | 16 | 10 |

▶ **Nitpicks:**

▶ Report *index* of minimum, not its value

▶ Report *leftmost* position in case of ties

# Finally: Longest common extensions

- In Unit 6: Left question open how to compute LCA in suffix trees

- But: Enhanced Suffix Array makes life easier!

$$\mathrm{LCE}(i, j) \;=\; \mathrm{LCP}\Big[\mathrm{RMQ}_{\mathrm{LCP}}\big(\min\{R[i], R[j]\} + 1,\; \max\{R[i], R[j]\}\big)\Big]$$



**Inverse suffix array: going left & right**

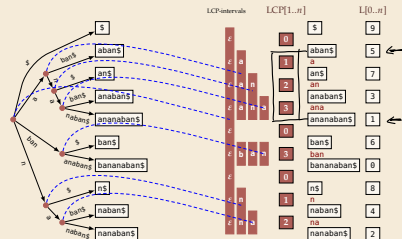- to understand the fastest algorithm, it is helpful to define the *inverse suffix array*:
  - $R[i] = r \iff L[r] = i$   $L = leaf\ array$
    $\iff$ there are $r$ suffixes that come before $T_i$ in sorted order
    $\iff$ $T_i$ has (0-based) *rank r* $\leadsto$ call $R[0..n]$ the *rank array*

| $i$ | $R[i]$ | $T_i$ | | $r$ | $L[r]$ | $T_{L[r]}$ |
|---|---|---|---|---|---|---|
| 0 | $6^{\text{th}}$ | bananaban\$ | | 0 | 9 | \$ |
| 1 | $4^{\text{th}}$ | ananaban\$ | | 1 | 5 | aban\$ |
| 2 | $9^{\text{th}}$ | nanaban\$ | | 2 | 7 | an\$ |
| 3 | $3^{\text{th}}$ | anaban\$ | | 3 | 3 | anaban\$ |
| 4 | $8^{\text{th}}$ | naban\$ | | 4 | 1 | ananaban\$ |
| 5 | $1^{\text{th}}$ | aban\$ | | 5 | 6 | ban\$ |
| 6 | $5^{\text{th}}$ | ban\$ | | 6 | 0 | bananaban\$ |
| 7 | $2^{\text{th}}$ | an\$ | | 7 | 8 | n\$ |
| 8 | $7^{\text{th}}$ | n\$ | | 8 | 4 | naban\$ |
| 9 | $0^{\text{th}}$ | \$ | | 9 | 2 | nanaban\$ |

right $R[0] = 6$

$L[8] = 4$  left

*sort suffixes*

29

**LCP array and internal nodes**

LCP-intervals   LCP$[1..n]$   $L[0..n]$

$\leadsto$ Leaf array $L[0..n]$ plus LCP array LCP$[1..n]$ encode full tree!

39

4

# Rules of the Game

- ▶ For the following, consider RMQ on arbitrary arrays
- ▶ comparison-based    ⇝    values don't matter, only relative order

- ▶ Two main quantities of interest:
    ⇝ space usage $\leq P(n)$
    1. **Preprocessing time**: Running time $P(n)$ of the preprocessing step
    2. **Query time**: Running time $Q(n)$ of one query (using precomputed data)

- ▶ Write $\langle P(n), Q(n) \rangle$ **time solution** for short

# Rules of the Game

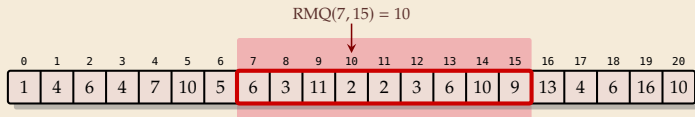- For the following, consider RMQ on arbitrary arrays

- comparison-based $\leadsto$ values don't matter, only relative order

- Two main quantities of interest:

  $\leadsto$ space usage $\leq P(n)$

  1. **Preprocessing time**: Running time $P(n)$ of the preprocessing step
  2. **Query time**: Running time $Q(n)$ of one query (using precomputed data)

- Write $\langle P(n), Q(n) \rangle$ **time solution** for short

### RMQ Implications for LCE
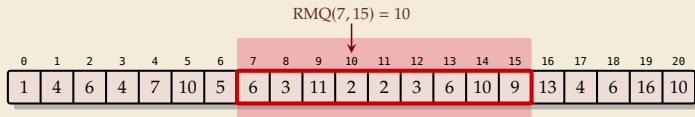
- Recall: Can compute (inverse) suffix array and LCP array in $O(n)$ time

- $\leadsto$ $\langle P(n), Q(n) \rangle$ time RMQ data structure implies
  $\langle P(n) + O(n), Q(n) \rangle$ time LCE data structure

# Trivial Solutions

RMQ(7, 15) = 10

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| 1 | 4 | 6 | 4 | 7 | 10 | 5 | 6 | 3 | 11 | 2 | 2 | 3 | 6 | 10 | 9 | 13 | 4 | 6 | 16 | 10 |

▶ Two easy solutions show extreme ends of scale:

# Trivial Solutions



RMQ(7, 15) = 10

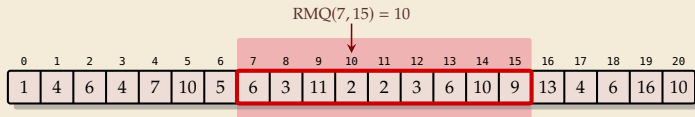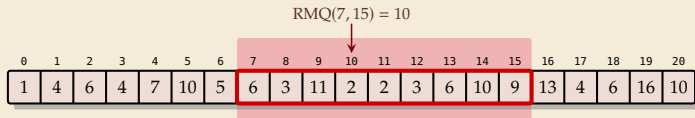| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|---|---|---|---|---|----|---|---|---|----|---|---|---|---|----|---|----|---|---|----|----|
| 1 | 4 | 6 | 4 | 7 | 10 | 5 | 6 | 3 | 11 | 2 | 2 | 3 | 6 | 10 | 9 | 13 | 4 | 6 | 16 | 10 |

▶ Two easy solutions show extreme ends of scale:

1. **Scan on demand**

   ▶ no preprocessing at all

   ▶ answer RMQ($i, j$) by scanning through $A[i..j]$, keeping track of min

   ⤳ $\langle O(1), O(n) \rangle$

6

# Trivial Solutions



RMQ(7, 15) = 10

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| 1 | 4 | 6 | 4 | 7 | 10 | 5 | 6 | 3 | 11 | 2 | 2 | 3 | 6 | 10 | 9 | 13 | 4 | 6 | 16 | 10 |

▶ Two easy solutions show extreme ends of scale:

## 1. Scan on demand

▶ no preprocessing at all

▶ answer $RMQ(i, j)$ by scanning through $A[i..j]$, keeping track of min

⤳ $\langle O(1), O(n) \rangle$

## 2. Precompute all

▶ Precompute all answers in a big 2D array $M[0..n)[0..n)$

▶ queries simple: $RMQ(i, j) = M[i][j]$

⤳ $\langle O(n^3), O(1) \rangle$

6

# Trivial Solutions

RMQ(7, 15) = 10

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| 1 | 4 | 6 | 4 | 7 | 10 | 5 | 6 | 3 | 11 | 2 | 2 | 3 | 6 | 10 | 9 | 13 | 4 | 6 | 16 | 10 |

▶ Two easy solutions show extreme ends of scale:

## 1. Scan on demand

   ▶ no preprocessing at all
   ▶ answer $RMQ(i, j)$ by scanning through $A[i..j]$, keeping track of min
   $\rightsquigarrow \langle O(1), O(n) \rangle$

## 2. Precompute all

   ▶ Precompute all answers in a big 2D array $M[0..n)[0..n)$
   ▶ queries simple: $RMQ(i, j) = M[i][j]$
   $\rightsquigarrow \langle O(n^3), O(1) \rangle$
   ▶ Preprocessing can reuse partial results $\rightsquigarrow \langle O(n^2), O(1) \rangle$

6

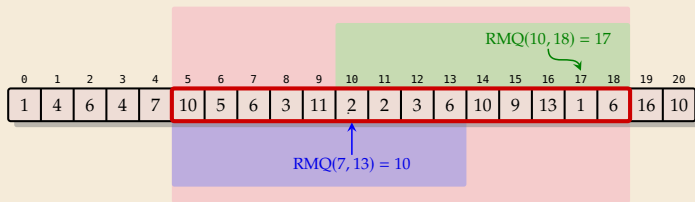# 7.2 RMQ – Sparse Table Solution

# Sparse Table

- **Idea:** Like "precompute-all", but keep only *some* entries

- store $M[i][j]$ iff $\ell = j - i + 1$ is $2^k$.
  - $\rightsquigarrow \leq n \cdot \lg n$ entries
  - $\rightsquigarrow$ Can be stored as $M'[i][k]$

## Sparse Table

- **Idea:** Like "precompute-all", but keep only *some* entries

- store $M[i][j]$    iff    $\ell = j - i + 1$ is $2^k$.
    - $\rightsquigarrow \leq n \cdot \lg n$ entries
    - $\rightsquigarrow$ Can be stored as $M'[i][k]$
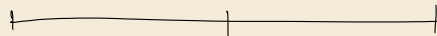
- How to answer queries?

# Sparse Table

▶ **Idea:** Like "precompute-all", but keep only *some* entries

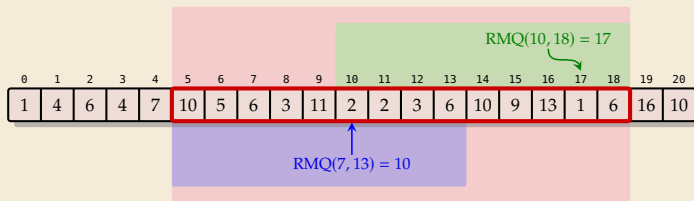▶ store $M[i][j]$ iff $\ell = j - i + 1$ is $2^k$.
  ⤳ $\leq n \cdot \lg n$ entries
  ⤳ Can be stored as $M'[i][k]$

▶ How to answer queries?



RMQ(10, 18) = 17

RMQ(7, 13) = 10

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| 1 | 4 | 6 | 4 | 7 | 10 | 5 | 6 | 3 | 11 | 2 | 2 | 3 | 6 | 10 | 9 | 13 | 1 | 6 | 16 | 10 |

*1.* Find $k$ with $\ell/2 \leq 2^k \leq \ell$

*2.* Cover range $[i..j]$ by
  $2^k$ positions right from $i$ and
  $2^k$ positions left from $j$

*3.* $\mathrm{RMQ}(i, j) =$
  $\arg\min\{A[rmq_1], A[rmq_2]\}$

  with $rmq_1 = \mathrm{RMQ}(i, i + 2^k - 1)$
  $rmq_2 = \mathrm{RMQ}(j - 2^k + 1, j)$

# Sparse Table

▶ **Idea:** Like "precompute-all", but keep only *some* entries

▶ store $M[i][j]$ iff $\ell = j - i + 1$ is $2^k$.
  ↝ $\leq n \cdot \lg n$ entries
  ↝ Can be stored as $M'[i][k]$

▶ How to answer queries?



RMQ(10, 18) = 17

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| 1 | 4 | 6 | 4 | 7 | 10 | 5 | 6 | 3 | 11 | 2 | 2 | 3 | 6 | 10 | 9 | 13 | 1 | 6 | 16 | 10 |

RMQ(7, 13) = 10

▶ Preprocessing can be done in $O(n \log n)$ times

↝ $\langle O(n \log n), O(1) \rangle$ time solution!

1. Find $k$ with $\ell/2 \leq 2^k \leq \ell$

2. Cover range $[i..j]$ by
   $2^k$ positions right from $i$ and
   $2^k$ positions left from $j$

3. $\text{RMQ}(i, j) =$
   $\arg\min\{A[rmq_1], A[rmq_2]\}$

   with $rmq_1 = \text{RMQ}(i, i + 2^k - 1)$
   $rmq_2 = \text{RMQ}(j - 2^k + 1, j)$

## Bootstrapping

▶ We know a $\langle O(n \log n), O(1) \rangle$ time solution

▶ If we use that for $m = \Theta(n/\log n)$ elements, $O(m \log m) = O(n)$!
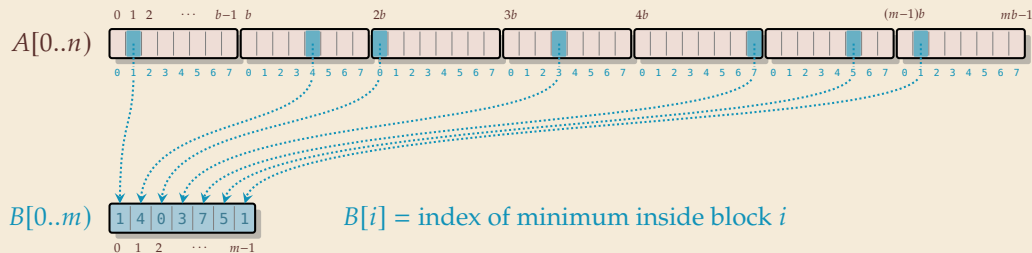
# Bootstrapping

▶ We know a $\langle O(n \log n), O(1) \rangle$ time solution

▶ If we use that for $m = \Theta(n/\log n)$ elements, $O(m \log m) = O(n)$!

▶ Break $A$ into blocks of $b = O(\log n)$ numbers

▶ Create array of block minima $B[0..m)$ for $m = \lceil n/b \rceil = O(n/\log n)$
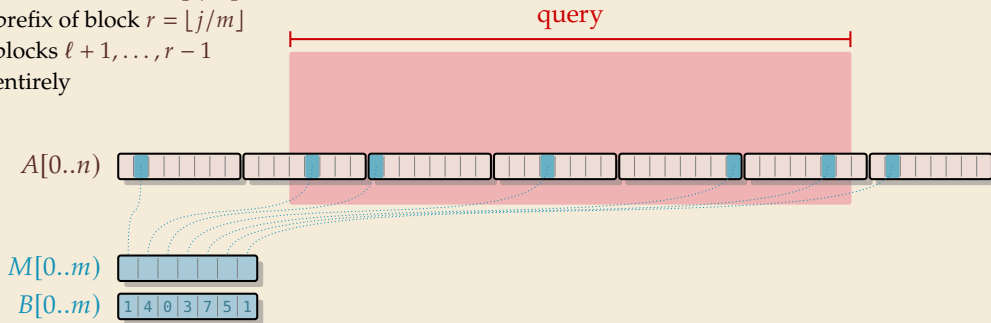


$B[i]$ = index of minimum inside block $i$

# Bootstrapping

- We know a $\langle O(n \log n), O(1) \rangle$ time solution

- If we use that for $m = \Theta(n/\log n)$ elements, $O(m \log m) = O(n)$!

- Break $A$ into blocks of $b = O(\log n)$ numbers

- Create array of block minima $B[0..m)$ for $m = \lceil n/b \rceil = O(n/\log n)$



$B[i]$ = index of minimum inside block $i$

- $\rightsquigarrow$ Use sparse table solution for $B$.

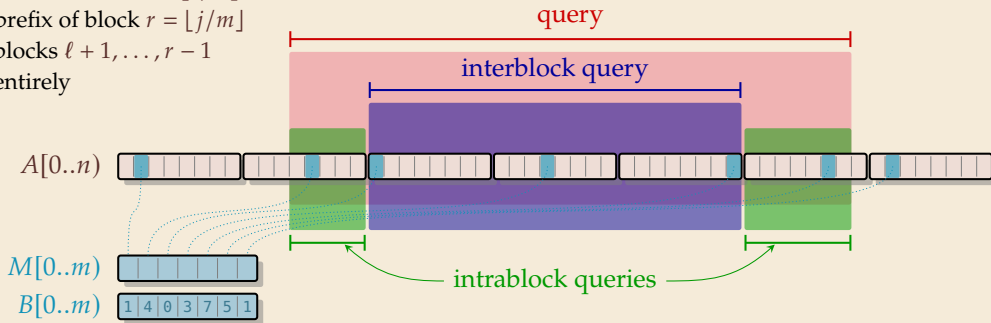$\rightsquigarrow$ Can solve RMQs in $B[0..m)$ in $\langle O(n), O(1) \rangle$ time

# Query decomposition

- Query $\mathrm{RMQ}_A(i, j)$ covers
  - suffix of block $\ell = \lfloor i/m \rfloor$
  - prefix of block $r = \lfloor j/m \rfloor$
  - blocks $\ell + 1, \ldots, r - 1$
    entirely
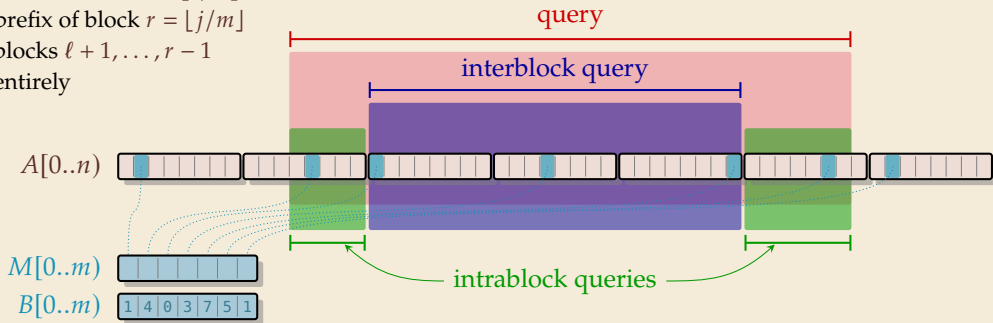
# Query decomposition

- Query $\mathrm{RMQ}_A(i, j)$ covers
  - suffix of block $\ell = \lfloor i/m \rfloor$
  - prefix of block $r = \lfloor j/m \rfloor$
  - blocks $\ell + 1, \ldots, r - 1$ entirely

# Query decomposition
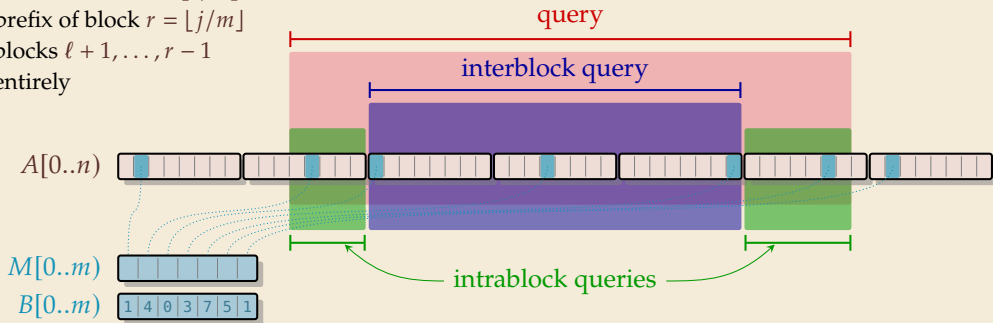
▶ Query $\text{RMQ}_A(i, j)$ covers
  ▶ suffix of block $\ell = \lfloor i/m \rfloor$
  ▶ prefix of block $r = \lfloor j/m \rfloor$
  ▶ blocks $\ell + 1, \ldots, r - 1$
    entirely



▶ $\text{RMQ}_A(i, j) = \underset{k \in K}{\arg\min} \, A[k]$ with $K = \begin{cases} \text{RMQ}_{\text{block } \ell}\big(i - \ell b, \, (\ell + 1)b - 1\big), \\ b \cdot \text{RMQ}_M\big(\ell + 1, r - 1\big) + \\ \qquad B\big[\text{RMQ}_M\big(\ell + 1, r - 1\big)\big], \\ \text{RMQ}_{\text{block } r}\big(rb, \, j - rb\big) \end{cases}$

# Query decomposition

- Query $\mathrm{RMQ}_A(i, j)$ covers
  - suffix of block $\ell = \lfloor i/m \rfloor$
  - prefix of block $r = \lfloor j/m \rfloor$
  - blocks $\ell + 1, \ldots, r - 1$ entirely
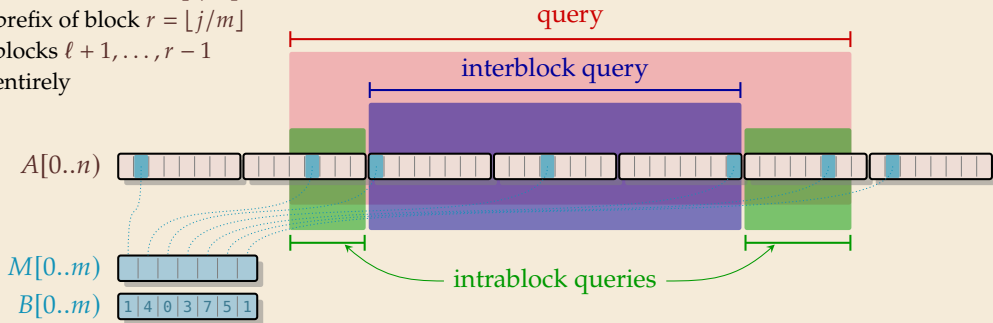


- $\mathrm{RMQ_A}(i, j) \; = \; \underset{k \in K}{\arg\min}\, A[k] \quad \text{with} \; K = \left\{ \begin{array}{l} \mathrm{RMQ}_{\mathrm{block}\,\ell}\big(i - \ell b, (\ell + 1)b - 1\big), \\ b \cdot \mathrm{RMQ}_M\big(\ell + 1, r - 1\big) + \\ \qquad B\big[\mathrm{RMQ}_M\big(\ell + 1, r - 1\big)\big], \\ \mathrm{RMQ}_{\mathrm{block}\,r}\big(rb, j - rb\big) \end{array} \right\}$

- $\rightsquigarrow$ only 3 possible values to check if intrablock and interblock queries known

# Query decomposition

- Query $\text{RMQ}_A(i, j)$ covers
  - suffix of block $\ell = \lfloor i/m \rfloor$
  - prefix of block $r = \lfloor j/m \rfloor$
  - blocks $\ell + 1, \ldots, r - 1$ entirely



query

interblock query

$A[0..n)$

intrablock queries

$M[0..m)$

$B[0..m)$  | 1 | 4 | 0 | 3 | 7 | 5 | 1 |

- $\text{RMQ}_A(i, j) = \underset{k \in K}{\arg\min} A[k]$ with $K = \begin{cases} \text{RMQ}_{\text{block } \ell}(i - \ell b, (\ell + 1)b - 1), \\ b \cdot \text{RMQ}_M(\ell + 1, r - 1) + \\ \quad B[\text{RMQ}_M(\ell + 1, r - 1)], \\ \text{RMQ}_{\text{block } r}(rb, j - rb) \end{cases}$

- ⤳ only 3 possible values to check
  if intrablock and interblock queries known ✓

$O(\log n)$

$\langle O(n), O(\log n) \rangle$

9