DTU

ALGORITHMIC RESEARCH: COOPERATION AROUND ORESOUND

# Time-Space Trade-Offs for Longest Common Extensions

(To appear at CPM 2012)

Philip Bille[1], Inge Li Gørtz[1], Benjamin Sach[2], and Hjalte Wedel Vildhøj[1]

[1]Technical University of Denmark, DTU Informatics, {phbi,ilg,hwvi}@imm.dtu.dk
[2]University of Warwick, Department of Computer Science, sach@dcs.warwick.ac.uk

University of Copenhagen, April 17, 2012

# The Longest Common Extension Problem
Definition

**Problem:** Preprocess a string $T$ of length $n$ to support LCE queries:

- $\text{LCE}(i,j) =$ The length of the longest common prefix of the suffixes starting at position $i$ and $j$ in $T$.

**Example**

$$T \;=\; \overset{1}{\text{b}}\;\overset{2}{\text{a}}\;\overset{3}{\text{n}}\;\overset{4}{\text{a}}\;\overset{5}{\text{n}}\;\overset{6}{\text{a}}\;\overset{7}{\text{s}} \qquad \text{LCE}(2,4) \;=\; ?$$

# The Longest Common Extension Problem
Definition

**Problem:** Preprocess a string $T$ of length $n$ to support LCE queries:

- $LCE(i, j)$ = The length of the longest common prefix of the suffixes starting at position $i$ and $j$ in $T$.

**Example**

$$T = \overset{\substack{1 \; 2 \; 3 \; 4 \; 5 \; 6 \; 7}}{\text{b a n a n a s}}$$

$$\text{a n a s}$$

$$\text{a n a n a s}$$

$$LCE(2, 4) = ?$$

# The Longest Common Extension Problem
Definition

**Problem:** Preprocess a string $T$ of length $n$ to support LCE queries:

- $\text{LCE}(i,j) = $ The length of the longest common prefix of the suffixes starting at position $i$ and $j$ in $T$.

**Example**

$$T = \overset{\overset{1\ \ 2\ \ 3\ \ 4\ \ 5\ \ 6\ \ 7}{}}{\texttt{b a n a n a s}} \qquad \text{LCE}(2,4) = 3$$

a n a s

a n a n a s

# The Longest Common Extension Problem
Definition

**Problem:** Preprocess a string $T$ of length $n$ to support LCE queries:

- $\text{LCE}(i,j)$ = The length of the longest common prefix of the suffixes starting at position $i$ and $j$ in $T$.

**Example**

$$T \;=\; \overset{1\;\;2\;\;3\;\;4\;\;5\;\;6\;\;7}{\text{b a n a n a s}} \qquad\qquad \text{LCE}(2,5)\;=\;0$$

n a s

a n a n a s

# The Longest Common Extension Problem

Motivation

Longest Common Extensions appear as a subproblem in many string matching problems, including

- Approximate string matching. I.e., find substrings of $T$ such that $T[i \ldots j] \approx P$ (hamming or edit distance).
- Finding palindromes. I.e., find substrings of $T$ such that $T[i \ldots j] = T[i \ldots j]^R$.
- Finding tandem repeats. I.e., find substrings of $T$ such that $T[i \ldots j] = UU$ for some string $U$.

## Example: Palindromes

$$T \;=\; \overset{\text{\scriptsize 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9}}{\texttt{b a b a a b a a c}}$$

# The Longest Common Extension Problem

Motivation

Longest Common Extensions appear as a subproblem in many string matching problems, including

- Approximate string matching. I.e., find substrings of $T$ such that $T[i \ldots j] \approx P$ (hamming or edit distance).
- Finding palindromes. I.e., find substrings of $T$ such that $T[i \ldots j] = T[i \ldots j]^R$.
- Finding tandem repeats. I.e., find substrings of $T$ such that $T[i \ldots j] = UU$ for some string $U$.

## Example: Palindromes

$$T \; = \; \text{b} \; \overset{1}{\text{a}} \; \overset{2}{\text{b}} \; \overset{3}{\text{a}} \; \overset{4}{\text{a}} \; \overset{5}{\text{b}} \; \overset{6}{\text{a}} \; \overset{7}{\text{a}} \; \text{c}$$

center

# The Longest Common Extension Problem

Motivation

Longest Common Extensions appear as a subproblem in many string matching problems, including

- Approximate string matching. I.e., find substrings of $T$ such that $T[i \ldots j] \approx P$ (hamming or edit distance).
- Finding palindromes. I.e., find substrings of $T$ such that $T[i \ldots j] = T[i \ldots j]^R$.
- Finding tandem repeats. I.e., find substrings of $T$ such that $T[i \ldots j] = UU$ for some string $U$.

## Example: Palindromes

$$T = \begin{array}{ccccccccc} \scriptstyle 1 & \scriptstyle 2 & \scriptstyle 3 & \scriptstyle 4 & \scriptstyle 5 & \scriptstyle 6 & \scriptstyle 7 & \scriptstyle 8 & \scriptstyle 9 \\ \texttt{b} & \texttt{a} & \texttt{b} & \texttt{a} & \texttt{a} & \texttt{b} & \texttt{a} & \texttt{a} & \texttt{c} \end{array}$$

center

# The Longest Common Extension Problem

Motivation

Longest Common Extensions appear as a subproblem in many string matching problems, including

- Approximate string matching. I.e., find substrings of $T$ such that $T[i \ldots j] \approx P$ (hamming or edit distance).
- Finding palindromes. I.e., find substrings of $T$ such that $T[i \ldots j] = T[i \ldots j]^R$.
- Finding tandem repeats. I.e., find substrings of $T$ such that $T[i \ldots j] = UU$ for some string $U$.

## Example: Palindromes

$$T = \overset{\substack{1\ \ 2\ \ 3\ \ 4\ \ 5\ \ 6\ \ 7\ \ 8\ \ 9}}{\text{b a b a a b a a c}}$$

center

All maximal palindromes in $P$ can be reported by performing $2n - 1$ LCE queries (one for each possible center).

# Two Simple Solutions

### #1: Store nothing

$$T \; = \; \underset{\underset{i}{\uparrow}}{\overset{1}{\text{b}}} \, \underset{}{\overset{2}{\text{a}}} \, \overset{3}{\text{n}} \, \underset{\underset{j}{\uparrow}}{\overset{4}{\text{a}}} \, \overset{5}{\text{n}} \, \overset{6}{\text{a}} \, \overset{7}{\text{s}}$$

$\text{LCE}(i,j) =$

# Two Simple Solutions

### #1: Store nothing

$$T \;=\; \overset{\scriptsize 1}{\text{b}}\; \overset{\scriptsize 2}{\text{a}}\; \overset{\scriptsize 3}{\text{n}}\; \overset{\scriptsize 4}{\text{a}}\; \overset{\scriptsize 5}{\text{n}}\; \overset{\scriptsize 6}{\text{a}}\; \overset{\scriptsize 7}{\text{s}}$$

$$\uparrow \qquad \uparrow$$
$$i \qquad\quad j$$

$\text{LCE}(i,j) = 1$

# Two Simple Solutions

### #1: Store nothing

$$T \; = \; \overset{1}{\text{b}} \; \overset{2}{\text{a}} \; \overset{3}{\text{n}} \; \overset{4}{\text{a}} \; \overset{5}{\text{n}} \; \overset{6}{\text{a}} \; \overset{7}{\text{s}}$$

$$\underset{i}{\uparrow} \qquad \underset{j}{\uparrow}$$

$\text{LCE}(i,j) = 2$

# Two Simple Solutions

### #1: Store nothing

$$T \;=\; \overset{1}{\text{b}} \; \overset{2}{\text{a}} \; \overset{3}{\text{n}} \; \overset{4}{\text{a}} \; \overset{5}{\text{n}} \; \overset{6}{\text{a}} \; \overset{7}{\text{s}}$$

$i$      $j$

$\text{LCE}(i,j) = 3$

# Two Simple Solutions

### #1: Store nothing

$$T = \begin{matrix} & 1 & 2 & 3 & 4 & 5 & 6 & 7 \\ & b & a & n & a & n & a & s \end{matrix}$$

$i$ (points to position 2), $j$ (points to position 4)

$LCE(i,j) = 3$

Time:   $O(n)$
Space:  $O(1)$

# Two Simple Solutions

### #1: Store nothing

$$T = \underset{i}{\underset{\uparrow}{\overset{1}{b}\ \overset{2}{a}}}\ \overset{3}{n}\ \underset{j}{\underset{\uparrow}{\overset{4}{a}}}\ \overset{5}{n}\ \overset{6}{a}\ \overset{7}{s}$$

$\text{LCE}(i,j) = 3$

| | |
|---|---|
| Time: | $O(n)$ |
| Space: | $O(1)$ |

### #2: Store the suffix tree

# Two Simple Solutions

#1: Store nothing

$$T = \overset{\substack{1\ 2\ 3\ 4\ 5\ 6\ 7}}{\texttt{b a n a n a s}}$$

$i$    $j$

LCE$(i,j) = 3$

| Time: | $O(n)$ |
|-------|--------|
| Space: | $O(1)$ |

#2: Store the suffix tree



LCE$(i,j) = |\text{NCA}(i,j)| = 3$

# Two Simple Solutions

### #1: Store nothing

$$T = \overset{1}{\text{b}} \overset{2}{\text{a}} \overset{3}{\text{n}} \overset{4}{\text{a}} \overset{5}{\text{n}} \overset{6}{\text{a}} \overset{7}{\text{s}}$$

$\uparrow$ $\quad$ $\uparrow$
$i$ $\qquad$ $j$

$\text{LCE}(i,j) = 3$

| Time: | $O(n)$ |
|---|---|
| Space: | $O(1)$ |

### #2: Store the suffix tree



NCA(2, 4)

$\text{LCE}(i,j) = |\text{NCA}(i,j)| = 3$

| Time: | $O(1)$ |
|---|---|
| Space: | $O(n)$ |

# Two Simple Solutions

## #1: Store nothing

$$T = \overset{\overset{1}{b}}{} \overset{2}{a} \overset{3}{n} \overset{4}{a} \overset{5}{n} \overset{6}{a} \overset{7}{s}$$

$i$ points to position 2, $j$ points to position 4

$\text{LCE}(i,j) = 3$

| Time: | $O(n)$ |
|-------|--------|
| Space: | $O(1)$ |

## #2: Store the suffix tree



$\text{NCA}(2,4)$

Trade-off?

| Time: | $O(1)$ |
|-------|--------|
| Space: | $O(n)$ |

$\text{LCE}(i,j) = |\text{NCA}(i,j)| = 3$

# Our Results

Store nothing

Time: $O(n)$
Space: $O(1)$

Trade-off?

Time: $O(1)$
Space: $O(n)$

Store suffix tree

Less space ——→

Faster ——→

# Our Results

Store nothing

| Trade-off parameter $\tau$, $1 \leq \tau \leq n$ |
| --- |

Time:    $O(n)$
Space:   $O(1)$

Randomized

Time:    $O\left(\tau \log\left(\frac{\text{LCE}(i,j)}{\tau}\right)\right)$
Space:   $O\left(\frac{n}{\tau}\right)$

Time:    $O(\tau)$
Space:   $O\left(\frac{n}{\sqrt{\tau}}\right)$

Deterministic

Trade-off?

Time:    $O(1)$
Space:   $O(n)$

Store suffix tree

Less space     Faster

# A Deterministic Solution

**Idea:** Store a subset of the $n$ suffixes in a compacted trie.

$$T = \begin{array}{cccccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \\ d & b & c & a & a & b & c & a & b & c & a & a & b & c & a & c \end{array}$$

# A Deterministic Solution

**Idea:** Store a subset of the *n* suffixes in a compacted trie.

$$T = \quad \begin{array}{cccccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \\ \text{d} & \text{b} & \text{c} & \text{a} & \text{a} & \text{b} & \text{c} & \text{a} & \text{b} & \text{c} & \text{a} & \text{a} & \text{b} & \text{c} & \text{a} & \text{c} \end{array}$$

# A Deterministic Solution

**Idea:** Store a subset of the $n$ suffixes in a compacted trie.

$$T = \begin{array}{cccccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \\ \text{d} & \text{b} & \text{c} & \text{a} & \text{a} & \text{b} & \text{c} & \text{a} & \text{b} & \text{c} & \text{a} & \text{a} & \text{b} & \text{c} & \text{a} & \text{c} \end{array}$$

# A Deterministic Solution

**Idea:** Store a subset of the *n* suffixes in a compacted trie.

$$T = \begin{array}{cccccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \\ d & b & c & a & a & b & c & a & b & c & a & a & b & c & a & c \end{array}$$

# A Deterministic Solution

**Idea:** Store a subset of the *n* suffixes in a compacted trie.

$$T = \quad \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \\ d & b & c & a & a & b & c & a & b & c & a & a & b & c & a & c \end{matrix}$$

# A Deterministic Solution

**Idea:** Store a subset of the $n$ suffixes in a compacted trie.

$$T = \quad \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \\ d & b & c & a & a & b & c & a & b & c & a & a & b & c & a & c \end{matrix}$$

# A Deterministic Solution

**Idea:** Store a subset of the *n* suffixes in a compacted trie.

$$T = \quad \overset{1}{d} \ \overset{2}{b} \ \overset{3}{c} \ \overset{4}{a} \ \overset{5}{a} \ \overset{6}{b} \ \overset{7}{c} \ \overset{8}{a} \ \overset{9}{b} \ \overset{10}{c} \ \overset{11}{a} \ \overset{12}{a} \ \overset{13}{b} \ \overset{14}{c} \ \overset{15}{a} \ \overset{16}{c}$$

# A Deterministic Solution

**Idea:** Store a subset of the $n$ suffixes in a compacted trie.

$$T = \begin{array}{cccccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \\ d & b & c & a & a & b & c & a & b & c & a & a & b & c & a & c \end{array}$$

# A Deterministic Solution

**Idea:** Store a subset of the *n* suffixes in a compacted trie.

$$T = \begin{array}{cccccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \\ \text{d} & \text{b} & \text{c} & \text{a} & \text{a} & \text{b} & \text{c} & \text{a} & \text{b} & \text{c} & \text{a} & \text{a} & \text{b} & \text{c} & \text{a} & \text{c} \end{array}$$

## Difference Covers

A *difference cover modulo* $\tau$ is a set of integers $D \subseteq \{0, 1, \ldots, \tau - 1\}$ such that for any distance $d \in \{0, 1, \ldots, \tau - 1\}$, $D$ contains two elements separated by distance $d$ modulo $\tau$.
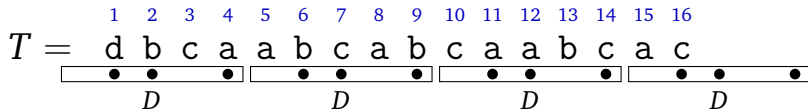
Ex: The set $D = \{1, 2, 4\}$ is a difference cover modulo 5.

| $d$ | 0 | 1 | 2 | 3 | 4 |
|-----|------|------|------|------|------|
| $i,j$ | 1, 1 | 2, 1 | 1, 4 | 4, 1 | 1, 2 |

# A Deterministic Solution

**Idea:** Store a subset of the *n* suffixes in a compacted trie.

$$T = \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \\ \text{d} & \text{b} & \text{c} & \text{a} & \text{a} & \text{b} & \text{c} & \text{a} & \text{b} & \text{c} & \text{a} & \text{a} & \text{b} & \text{c} & \text{a} & \text{c} \end{matrix}$$

$$\underbrace{\phantom{d\ b\ c\ a}}_{D} \quad \underbrace{\phantom{a\ b\ c\ a}}_{D} \quad \underbrace{\phantom{b\ c\ a\ a}}_{D} \quad \underbrace{\phantom{b\ c\ a\ c}}_{D}$$

## Difference Covers

A *difference cover modulo* $\tau$ is a set of integers $D \subseteq \{0, 1, \ldots, \tau - 1\}$ such that for any distance $d \in \{0, 1, \ldots, \tau - 1\}$, $D$ contains two elements separated by distance $d$ modulo $\tau$.

Ex: The set $D = \{1, 2, 4\}$ is a difference cover modulo 5.

| $d$ | 0 | 1 | 2 | 3 | 4 |
|-----|-----|-----|-----|-----|-----|
| $i,j$ | 1, 1 | 2, 1 | 1, 4 | 4, 1 | 1, 2 |

# A Deterministic Solution

**Idea:** Store a subset of the *n* suffixes in a compacted trie.

$$T = \begin{array}{ccccccccccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \\ & d & b & c & a & a & b & c & a & b & c & a & a & b & c & a & c \end{array}$$

## Lemma (Colbourn and Ling[1])

*For any $\tau$, a difference cover modulo $\tau$ of size at most $\sqrt{1.5\tau} + 6$ can be computed in $O(\sqrt{\tau})$ time.*

---

[1] C. J. Colbourn and A. C. Ling. Quorums from difference covers. Inf. Process. Lett. 75(1-2):9–12, 2000

# A Deterministic Solution

**Idea:** Store a subset of the *n* suffixes in a compacted trie.

$$T = \begin{array}{ccccccccccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \\ & d & b & c & a & a & b & c & a & b & c & a & a & b & c & a & c \end{array}$$

## Lemma (Colbourn and Ling[1])
*For any $\tau$, a difference cover modulo $\tau$ of size at most $\sqrt{1.5\tau} + 6$ can be computed in $O(\sqrt{\tau})$ time.*

## Analysis
**Time:** $O(\tau)$

**Space:** $O(\#\text{stored suffixes}) = O\left(\frac{n}{\tau}|D|\right) = O\left(\frac{n}{\sqrt{\tau}}\right)$

---

[1]C. J. Colbourn and A. C. Ling. Quorums from difference covers. Inf. Process. Lett. 75(1-2):9–12, 2000

# A Randomized Solution

## Rabin-Karp Fingerprints

Let $p$ be a sufficiently large prime and choose $b \in \mathbb{Z}_p$ uniformly at random.

$$\phi(S) = \sum_{k=1}^{|S|} S[k] b^k \bmod p \,.$$

$$
\begin{array}{ccccccccccccccccc}
 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \\
T = & \text{d} & \text{b} & \text{c} & \text{a} & \text{a} & \text{b} & \text{c} & \text{a} & \text{b} & \text{c} & \text{a} & \text{a} & \text{b} & \text{c} & \text{a} & \text{c} \\
 = & 3 & 1 & 2 & 0 & 0 & 1 & 2 & 0 & 1 & 2 & 0 & 0 & 1 & 2 & 0 & 2
\end{array}
$$

$$\phi(T[2\ldots 7]) = 120012 \bmod 31 = 11$$

**Crucial property:** With high probability $\phi$ is collision-free on substrings of $T$, i.e., $\phi(S_1) = \phi(S_2)$ iff $S_1 = S_2$.

**Also important:** $\phi(T[i\ldots j+1])$ can be computed from $\phi(T[i\ldots j])$ in $O(1)$ time.

# A Randomized Solution

**Idea:** Store fingerprints of suffixes starting at every $\tau$'th position in $T$.



Blocks of $\tau$ chars

$T =$
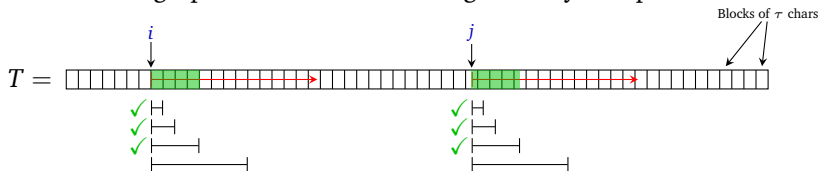
# A Randomized Solution

**Idea:** Store fingerprints of suffixes starting at every $\tau$'th position in $T$.



**Observation:** If $S$ is block aligned we can compute $\phi(S)$ in $O(1)$ time. Otherwise, the time needed is $O(\tau)$.

# A Randomized Solution

**Idea:** Store fingerprints of suffixes starting at every $\tau$'th position in $T$.



**Observation:** If $S$ is block aligned we can compute $\phi(S)$ in $O(1)$ time. Otherwise, the time needed is $O(\tau)$.
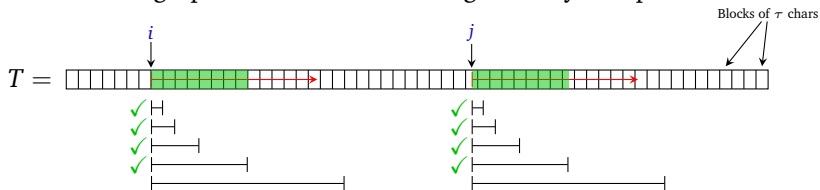
# A Randomized Solution
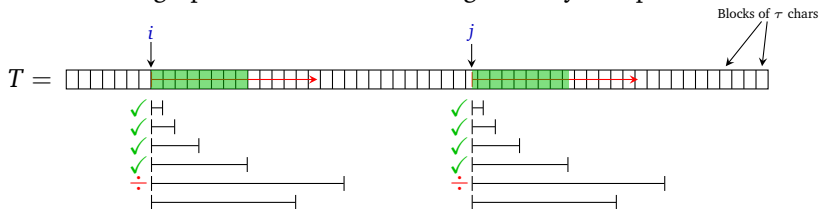
**Idea:** Store fingerprints of suffixes starting at every $\tau$'th position in $T$.

# A Randomized Solution

**Idea:** Store fingerprints of suffixes starting at every $\tau$'th position in $T$.

# A Randomized Solution
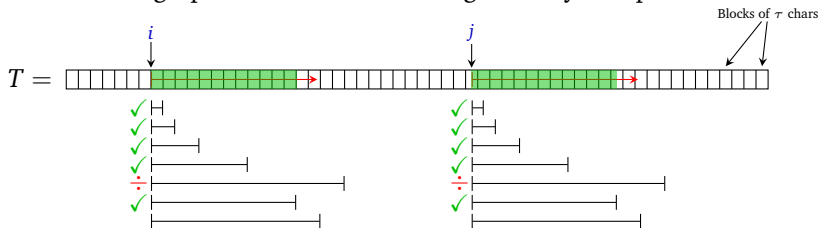
**Idea:** Store fingerprints of suffixes starting at every $\tau$'th position in $T$.

# A Randomized Solution

**Idea:** Store fingerprints of suffixes starting at every $\tau$'th position in $T$.

# A Randomized Solution

**Idea:** Store fingerprints of suffixes starting at every $\tau$'th position in $T$.

# A Randomized Solution

**Idea:** Store fingerprints of suffixes starting at every $\tau$'th position in $T$.

# A Randomized Solution

How to answer a query

**Idea:** Store fingerprints of suffixes starting at every $\tau$'th position in $T$.
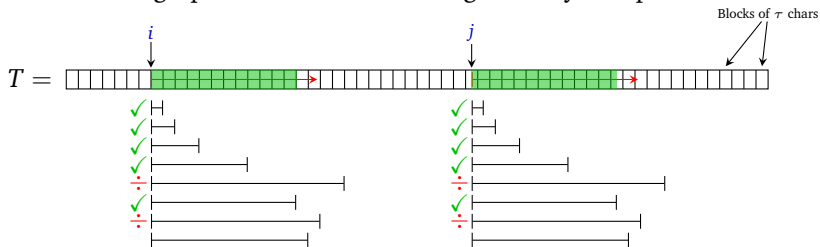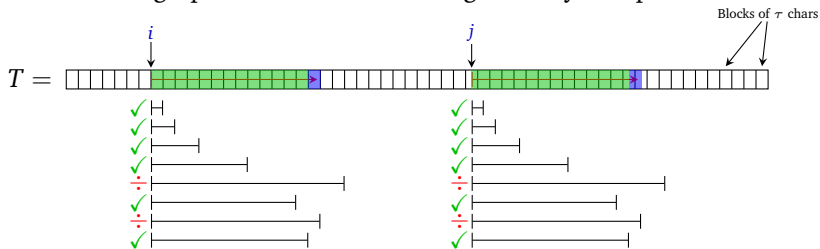
# A Randomized Solution

How to answer a query

**Idea:** Store fingerprints of suffixes starting at every $\tau$'th position in $T$.

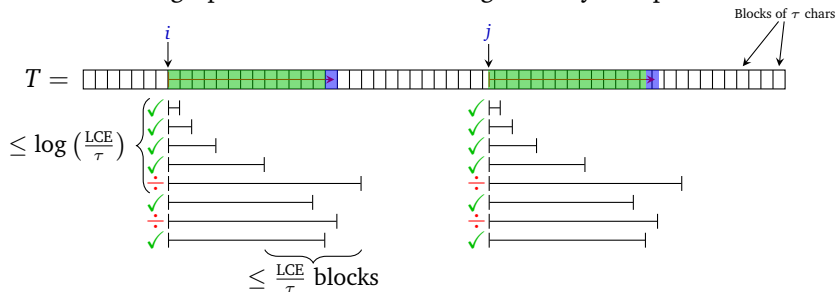# A Randomized Solution

**Idea:** Store fingerprints of suffixes starting at every $\tau$'th position in $T$.

# A Randomized Solution

**Idea:** Store fingerprints of suffixes starting at every $\tau$'th position in $T$.



## Analysis

**Time:** At most $2\log(\frac{\text{LCE}}{\tau})$ fingerprint comparisons each taking time $O(\tau)$. Hence query time $O\left(\tau \log\left(\frac{\text{LCE}}{\tau}\right)\right)$.

**Space:** $O\left(\frac{n}{\tau}\right)$.

# A Randomized Solution

**Idea:** Store fingerprints of suffixes starting at every $\tau$'th position in $T$.



## Analysis

**Time:** At most $2\log(\frac{\text{LCE}}{\tau})$ fingerprint comparisons each taking time $O(\tau)$. Hence query time $O\left(\tau\log\left(\frac{\text{LCE}}{\tau}\right)\right)$.

**Space:** $O\left(\frac{n}{\tau}\right)$.