# 5 Parallel String Matching

*10 March 2021*

Sebastian Wild

# 5 Parallel String Matching

# Parallelizing string matching

- ► We have seen a plethora of string matching methods

- ► But all efficient methods seem inherently sequential
  *Indeed, they became efficient only after building on knowledge from previous steps!*

  Sounds like the *opposite* of parallel!

⤳ This unit:
  - ► How well can we parallelize string matching?
  - ► What new ideas can help?

  Here: string matching = find *all* occurrences of $P$ in $T$     (more natural problem for parallel)
  always assume $m \leq n$

## 5.1 Elementary Tricks

# Embarrassingly Parallel

- ▶ A problem is called *"embarrassingly parallel"*
  if it can immediately be split into *many, small subtasks*
  that can be solved completely *independently* of each other

- ▶ Typical example: sum of two large matrices     (all entries independent)

- ⇝ best case for parallel computation     (simply assign each processor one subtask)

- ▶ Sorting is not embarrassingly parallel
  - ▶ no obvious way to define many *small* (=efficiently solvable) subproblems
  - ▶ but: some subtasks of our algorithms are, e.g., comparing all elements with pivot

# Clicker Question

Is the string-matching problem "embarrassingly parallel"?

**A** Yes

**B** No

**C** Only for $n \gg m$

**D** Only for $n \approx m$

`sli.do/comp526`

Click on "Polls" tab

# Elementary parallel string matching

**Subproblems in string matching:**

- ▶ string matching = check all guesses $i = 0, \ldots, n - m - 1$
- ▶ checking one guess is a subtask!

# Elementary parallel string matching

**Subproblems in string matching:**

- ▶ string matching = check all guesses $i = 0, \ldots, n - m - 1$
- ▶ checking one guess is a subtask!

**Approach 1:**

- ▶ Check all guesses in parallel _____ all n-m guesses
- ⤳ **Time**: $\Theta(m)$     using sequential checks
         $\Theta(\log m)$ on CREW-PRAM  (⤳ see tutorials)
         $\Theta(1)$     on CRCW-PRAM  (⤳ see tutorials)
- ⤳ **Work**: $\Theta((n - m)m)$  ⤳  not great . . .

# Elementary parallel string matching

**Subproblems in string matching:**

- string matching = check all guesses $i = 0, \ldots, n - m - 1$
- checking one guess is a subtask!

**Approach 1:**

- Check all guesses in parallel
- ⤳ **Time**: $\Theta(m)$     using sequential checks
               $\Theta(\log m)$ on CREW-PRAM ($\leadsto$ see tutorials)
               $\Theta(1)$     on CRCW-PRAM ($\leadsto$ see tutorials)
- ⤳ **Work**: $\Theta((n - m)m)$   ⤳   not great ...



**Approach 2:**

- Divide $T$ into **overlapping** blocks of $\underline{2m}$ characters:
  $T[0..2m), T[m..3m), T[2m..4m), T[3m..5m) \ldots$
- Find matches inside blocks in parallel, using efficient sequential method
      ⤳   $\Theta(2m + m) = \Theta(m)$ each
- ⤳ **Time**: $\underline{\Theta(m)}$      **Work**: $\Theta(\frac{n}{m} \cdot m) = \underline{\Theta(n)}$

# Clicker Question

Is the string-matching problem "embarrassingly parallel"?

**A** Yes

**B** No

**C** Only for $n \gg m$

**D** Only for $n \approx m$

`sli.do/comp526`

Click on "Polls" tab

# Clicker Question

Is the string-matching problem "embarrassingly parallel"?

**A** ~~Yes~~

**B** ~~No~~

**C** Only for $n \gg m$ ✓

**D** ~~Only for $n \approx m$~~

`sli.do/comp526`

Click on "Polls" tab

# Elementary parallel matching – Discussion

👍 very simple methods

👍 could even run distributed with access to part of $T$

👎 parallel speedup only for $m \ll n$

**Goal:**

▶ methods with better parallel time!       ⤳ higher speedup

⤳ must genuinely parallelize the matching process!       (and the preprocessing of the pattern)
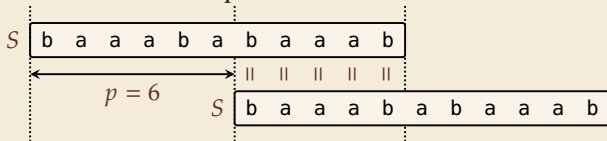
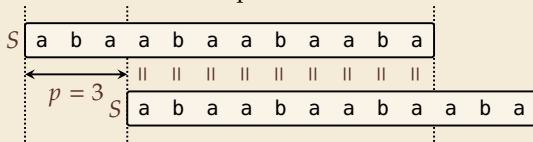⤳ need new ideas

## 5.2 Periodicity

## Periodicity of Strings

▶ $S = S[0..n-1]$ has *period p*  iff  $\forall i \in [0..n-p) : S[i] = S[i+p]$

▶ $p = 0$ and any $p \geq n$ are trivial periods  but these are not very interesting . . .
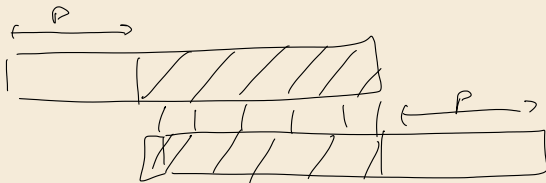
**Examples:**

▶ $S = $ baaababaaab has period 6:

| $S$ | b | a | a | a | b | a | b | a | a | a | b |
|---|---|---|---|---|---|---|---|---|---|---|---|

$\longleftarrow\quad p = 6\quad\longrightarrow$   ‖  ‖  ‖  ‖  ‖

| | | | | | | $S$ | b | a | a | a | b | a | b | a | a | a | b |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

▶ $S = $ abaabaabaaba has period 3:

| $S$ | a | b | a | a | b | a | a | b | a | a | b | a |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

$\longleftarrow p = 3$   ‖  ‖  ‖  ‖  ‖  ‖  ‖  ‖  ‖

| | | | $S$ | a | b | a | a | b | a | a | b | a | a | b | a |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# Periodicity and KMP

### Lemma 5.1 (Periodicity = Longest Overlap)
$p \in [1..n]$ is the *shortest* period in $S = S[0..n-1]$
iff $S[0..n-p)$ is the longest prefix that is also a suffix of $S[p..n)$. ◄

## Periodicity and KMP

**Lemma 5.1 (Periodicity = Longest Overlap)**

$p \in [1..n]$ is the *shortest* period in $S = S[0..n-1]$
iff $S[0..n-p)$ is the longest prefix that is also a suffix of $S[p..n)$. ◄

$S[0..n-1]$ has minimal period $p \iff fail[n] = n - p$