Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

School of Computer Science and En
University of New South Wales

9. STRING MATCHING ALGORITHMS

# Assignment Project Exam Help

- Assume that you want to find out if a string $B = b_0 b_1 \ldots b_{m-1}$ appears as a (contiguous) substring of a much longer string $A = a \ a_1 \ldots a_{n-1}$.

https://eduassistpro.github.i

Add WeChat edu_assist_pr

# Assignment Project Exam Help

- Assume that you want to find out if a string $B = b_0 b_1 \ldots b_{m-1}$ appears as a (contiguous) substring of a much longer string $A = a\ a_1 \ldots a_{n-1}$.

- The https://eduassistpro.github.i

Add WeChat edu_assist_pr

Assignment Project Exam Help

- Assume that you want to find out if a string $B = b_0 b_1 \ldots b_{m-1}$ appears as a (contiguous) substring of a much longer string $A = a \; a_1 \ldots a_{n-1}$.

- The https://eduassistpro.github.i

- We now show how hashing can be combined with re an efficient string matching algorithm.

Add WeChat edu_assist_pr

- We compute a hash value for the string $B = b_0 b_1 b_2 \ldots b_m$ in the following way.

# Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- We compute a hash value for the string $B = b_0 b_1 b_2 \ldots b_m$ in the following way.

- We will assume that strings $A$ and $B$ are in an alphabet $\mathcal{A}$ with $d$ many symbols in total

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- We compute a hash value for the string $B = b_0 b_1 b_2 \ldots b_m$ in the following way.

- We will assume that strings $A$ and $B$ are in an alphabet $\mathcal{A}$ with $d$ many symbols in total.

- Thus, we can identify each string with a sequence of integers by mapping each sym

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- We compute a hash value for the string $B = b_0 b_1 b_2 \ldots b_m$ in the following way.

- We will assume that strings $A$ and $B$ are in an alphabet $\mathcal{A}$ with $d$ many symbols in total.

- Thus, we can identify each string with a sequence of integers by mapping each sym

- To any string $B = b_0 b_1 \ldots b_{m-1}$ we can now associate an integer whose digits in base $d$ are integers corresponding to each symb

$$ B = h(b_0 b_1 b_2 \ldots b_m) = a $$

- We compute a hash value for the string $B = b_0 b_1 b_2 \ldots b_m$ in the following way.

- We will assume that strings $A$ and $B$ are in an alphabet $\mathcal{A}$ with $d$ many symbols in total.

- Thus, we can identify each string with a sequence of integers by mapping each sym

- To any string $B = b_0 b_1 \ldots b_{m-1}$ we can now associate an integer whose digits in base $d$ are integers corresponding to each symb

$$h(B) = h(b_0 b_1 b_2 \ldots b_m) = d^{m-1} b_0 +$$

- This can be done efficiently using the Horner's rule:

$$h(B) = b_{m-1} + d(b_{m-2} + d(b_{m-3} + d(b_{m-4} + \ldots + d(b_1 + d \cdot b_0))) \ldots)$$

- We compute a hash value for the string $B = b_0 b_1 b_2 \ldots b_m$ in the following way.

- We will assume that strings $A$ and $B$ are in an alphabet $\mathcal{A}$ with $d$ many symbols in total.

- Thus, we can identify each string with a sequence of integers by mapping each sym

- To any string $B = b_0 b_1 \ldots b_{m-1}$ we can now associate an integer whose digits in base $d$ are integers corresponding to each symb

$$h(B) = \ldots b_0 b_1 b_2 \ldots b_m = d^{\ldots} b_0 +$$

- This can be done efficiently using the Horner's rule:

$$h(B) = b_{m-1} + d(b_{m-2} + d(b_{m-3} + d(b_{m-4} + \ldots + d(b_1 + d \cdot b_0)))\ldots)$$

- Next we choose a large prime number $p$ such that $(d+1)\,p$ still fits into a single register and define the hash value of $B$ as $H(B) = h(B) \bmod p$.

- Recall that $A = a_0 a_1 a_2 a_3 \ldots \ldots a_s a_{s+1} \ldots a_{s+m-1} \ldots \ldots a_{N-1}$ where $N >> m$.

# Assignment Project Exam Help

## https://eduassistpro.github.i

## Add WeChat edu_assist_pr

- Recall that $A = a_0 a_1 a_2 a_3 \ldots \ldots a_s a_{s+1} \ldots a_{s+m-1} \ldots \ldots a_{N-1}$ where $N >> m$.
- We want to find efficiently all $s$ such that the string of length $m$ of the form $a_s a_{s+1} \ldots a_{s+m-1}$ and string $b_0 b_1 \ldots b_{m-1}$ are equal.

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- Recall that $A = a_0 a_1 a_2 a_3 \ldots \ldots a_s a_{s+1} \ldots a_{s+m-1} \ldots \ldots a_{N-1}$ where $N >> m$.
- We want to find efficiently all $s$ such that the string of length $m$ of the form $a_s a_{s+1} \ldots a_{s+m-1}$ and string $b_0 b_1 \ldots b_{m-1}$ are equal.
- For each contiguous substring $A_s = a_s a_{s+1} \ldots a_{s+m-1}$ of string $A$ we also compute its hash value as

$$\ldots \ldots \ldots ) \bmod p$$

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- Recall that $A = a_0 a_1 a_2 a_3 \ldots \ldots a_s a_{s+1} \ldots a_{s+m-1} \ldots \ldots a_{N-1}$ where $N >> m$.
- We want to find efficiently all $s$ such that the string of length $m$ of the form $a_s a_{s+1} \ldots a_{s+m-1}$ and string $b_0 b_1 \ldots b_{m-1}$ are equal.
- For each contiguous substring $A_s = a_s a_{s+1} \ldots a_{s+m-1}$ of string $A$ we also compute its hash value as

$$\ldots\ldots ) \bmod p$$

- We c
  symbol-by-symbol matching only if $H(B) = \ldots_s$

- Recall that $A = a_0a_1a_2a_3\ldots\ldots a_sa_{s+1}\ldots a_{s+m-1}\ldots\ldots a_{N-1}$ where $N >> m$.
- We want to find efficiently all $s$ such that the string of length $m$ of the form $a_sa_{s+1}\ldots a_{s+m-1}$ and string $b_0b_1\ldots b_{m-1}$ are equal.
- For each contiguous substring $A_s = a_sa_{s+1}\ldots a_{s+m-1}$ of string $A$ we also compute its hash value as

$$) \bmod p$$

- We c
  symbol-by-symbol matching only if $H(B) = $ $_s$
- Clearly, such an algorithm would be faster than the hai
  comparison only if we can compute the hash values of su ster
  than what it takes to compare strings $B$ and $A_s$ character by character.

- Recall that $A = a_0 a_1 a_2 a_3 \ldots \ldots a_s a_{s+1} \ldots a_{s+m-1} \ldots \ldots a_{N-1}$ where $N >> m$.
- We want to find efficiently all $s$ such that the string of length $m$ of the form $a_s a_{s+1} \ldots a_{s+m-1}$ and string $b_0 b_1 \ldots b_{m-1}$ are equal.
- For each contiguous substring $A_s = a_s a_{s+1} \ldots a_{s+m-1}$ of string $A$ we also compute its hash value as

$$) \bmod p$$

- We c
  symbol-by-symbol matching only if $H(B) =$                                        $_s$

- Clearly, such an algorithm would be faster than the naive
  comparison only if we can compute the hash values of su             ster
  than what it takes to compare strings $B$ and $A_s$ character by character.

- This is where recursion comes into play: we do not have compute the hash
  value $H(A_{s+1})$ of $A_{s+1} = a_{s+1} a_{s+2} \ldots a_{s+m}$ "from scratch", but we can
  compute it efficiently from the hash value $H(A_s)$ of $A_s = a_s a_{s+1} \ldots a_{s+m-1}$ as
  follows.

Since

Assignment Project Exam Help

$$H(A_s) = (d^{m-1}a_s + d^{m-2}a_{s+1} + \ldots d^1 a_{s+m-2} + a_{s+m-1}) \bmod p$$

https://eduassistpro.github.i

Add WeChat edu_assist_pr

Since

$$H(A_s) = (d^{m-1}a_s + d^{m-2}a_{s+1} + \ldots d^1 a_{s+m-2} + a_{s+m-1}) \bmod p$$

by multipl

$(d \cdot H(A_s)) \bmod$

$$= (d^m a_s + d^{m-1}a_{s+1} + \ldots d \cdot a_{s+m-1}) \bmod p$$
$$= (d^m a_s + (d^{m-1}a_{s+1} + \ldots d^2 a_{s+m-2} + d\, a_{s+m}) \bmod p$$
$$= (d^m a_s + H(A_{s+1}) - a_{s+m}) \bmod p$$

- Consequently,

$$H(A_{s+1}) = (d \cdot H(A_s) - d^m a_s + a_{s+m}) \bmod p.$$

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- Consequently,

$$H(A_{s+1}) = (d \cdot H(A_s) - d^m a_s + a_{s+m}) \bmod p.$$

Assignment Project Exam Help

- Note that

$$(d^m a_s) \bmod p = ((d^m \bmod p)a_s) \bmod p$$

and t

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- Consequently,

$$H(A_{s+1}) = (d \cdot H(A_s) - d^m a_s + a_{s+m}) \bmod p.$$

# Assignment Project Exam Help

- Note that

$$(d^m a_s) \bmod p = ((d^m \bmod p) a_s) \bmod p$$

and t

- Als https://eduassistpro.github.i

Add WeChat edu_assist_pr

- Consequently,

$$H(A_{s+1}) = (d \cdot H(A_s) - d^m a_s + a_{s+m}) \bmod p.$$

- Note that

$$(d^m a_s) \bmod p = ((d^m \bmod p) a_s) \bmod p$$

and t

- Als https://eduassistpro.github.i

- Thus, since $H(A_s) < p$ we obtain

$$d \cdot H(A_s) + ((d^m a_s + a_{s+}$$

Add WeChat edu_assist_pr

- Consequently,

$$H(A_{s+1}) = (d \cdot H(A_s) - d^m a_s + a_{s+m}) \bmod p.$$

- Note that

$$(d^m a_s) \bmod p = ((d^m \bmod p) a_s) \bmod p$$

and t

- Als

- Thus, since $H(A_s) < p$   we obtain

- Thus, since we chose $p$ such that $(d+1)\,p$ fits in a register, all the values and the intermediate results for the above expression also fit in a single register.

- Consequently,

$$H(A_{s+1}) = (d \cdot H(A_s) - d^m a_s + a_{s+m}) \bmod p.$$

- Note that

$$(d^m a_s) \bmod p = ((d^m \bmod p) a_s) \bmod p$$

and t

- Als

- Thus, since $H(A_s) < p$ we obtain

- Thus, since we chose $p$ such that $(d+1)p$ fits in a register, all the values and the intermediate results for the above expression also fit in a single register.

- Thus, for every $s$ except $s = 0$ the value of $H(A_s)$ can be computed in constant time independent of the length of the strings $A$ and $B$.

- Thus, we first compute $H(B)$ and $H(A_0)$ using the Horner's rule.

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- Thus, we first compute $H(B)$ and $H(A_0)$ using the Horner's rule.

- Subsequent values of $H(A_s)$ for $s > 0$ are computed in constant time using the above recursion.

- Thus, we first compute $H(B)$ and $H(A_0)$ using the Horner's rule.

- Subsequent values of $H(A_s)$ for $s > 0$ are computed in constant time using the above recursion.

- $H($ ... $A_s$ and $B$ a ... equ ...

- Thus, we first compute $H(B)$ and $H(A_0)$ using the Horner's rule.

- Subsequent values of $H(A_s)$ for $s > 0$ are computed in constant time using the above recursion.

- $H($ ... $A_s$ and $B$ a ... equ ...

- Since $p$ was chosen large, the false positives whe ... but $A_s \neq B$ are very unlikely, which makes the algo ...

- Thus, we first compute $H(B)$ and $H(A_0)$ using the Horner's rule.

- Subsequent values of $H(A_s)$ for $s > 0$ are computed in constant time using the above recursion.

- $H($ ... $A_s$ and $B$ a ... equ ...

- Since $p$ was chosen large, the false positives when ... but $A_s \neq B$ are very unlikely, which makes the algo ...

- However, as always when we use hashing, we cann ... case performance.

- Thus, we first compute $H(B)$ and $H(A_0)$ using the Horner's rule.

- Subsequent values of $H(A_s)$ for $s > 0$ are computed in constant time using the above recursion.

- $H($ ... $A_s$ and $B$ a ... equ ...

- Since $p$ was chosen large, the false positives whe ... but $A_s \neq B$ are very likely, which makes the algo ...

- However, as always when we use hashing, we cann... case performance.

- So we now look for algorithms whose worst case performance can be guaranteed.

- A string matching finite automaton for a string $S$ with $k$ symbols has $k+1$ many states $0, 1, \ldots k$ which correspond to the number of characters matched thus far and a transition function $\delta(s, c)$ where $s$ is a state and $c$ is a character read at the moment.

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- A string matching finite automaton for a string $S$ with $k$ symbols has $k+1$ many states $0, 1, \ldots k$ which correspond to the number of characters matched thus far and a transition function $\delta(s, c)$ where $s$ is a state and $c$ is a character read at the moment.

- We first look at the case when such $\delta(s, c)$ is given by a pre-constructed table.

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- A string matching finite automaton for a string $S$ with $k$ symbols has $k+1$ many states $0, 1, \ldots k$ which correspond to the number of characters matched thus far and a transition function $\delta(s, c)$ where $s$ is a state and $c$ is a character read at the moment.

- We first look at the case when such $\delta(s, c)$ is given by a pre-constructed table.

- To make things easier to describe, we consider the string $S = ababaca$. The table defin

| state | a | b | c |   |
|-------|---|---|---|---|
| 0     | **1** | 0 | 0 | a |
| 1     | 1 | **2** | 0 | b |
| 2     | **3** | 0 | 0 | a |
| 3     | 1 | **4** | 0 | b |
| 4     | **5** | 0 | 0 | a |
| 5     | 1 | 4 | **6** | c |
| 6     | **7** | 0 | 0 | a |
| 7     | 1 | 2 | 0 |   |



state transition diagram for string $ababaca$

- How do we compute the transition function $\delta$, i.e., how do we fill the table?

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- How do we compute the transition function $\delta$, i.e., how do we fill the table?

Assignment Project Exam Help

- Let $B_k$ denote the prefix of the string $B$ consisting of the first $k$ cha

https://eduassistpro.github.i

Add WeChat edu_assist_pr

- How do we compute the transition function $\delta$, i.e., how do we fill the table?

- Let $B_k$ denote the prefix of the string $B$ consisting of the first $k$ cha

- If we
  $B_k$; f
  that

- How do we compute the transition function $\delta$, i.e., how do we fill the table?

- Let $B_k$ denote the prefix of the string $B$ consisting of the first $k$ cha

- If we $B_k$; if that $_m$

- Thus, if $a$ happens to be $B[k+1]$, then $k+1$ and $B_k$ $B_{k+1}$.

- How do we compute the transition function $\delta$, i.e., how do we fill the table?

- Let $B_k$ denote the prefix of the string $B$ consisting of the first $k$ cha

- If we $B_k$; if that $m$ $k$ .

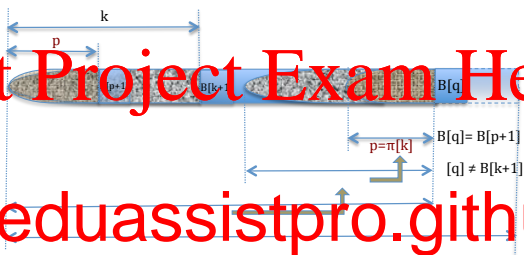- Thus, if $a$ happens to be $B[k+1]$, then $k+1$ and $B_k$ $B_{k+1}$.

- We do that by matching the string against itself: we can recursively compute a function $\pi(k)$ which for each $k$ returns the largest integer $m$ such that the prefix $B_m$ of $B$ is a proper suffix of $B_k$.

```
1:  function
2:      m ← length[B]
3:      let π[1..m] be a new array
4:      π[1] = 0
5:      k = 0
6:      for q
7:          w
               B[k + 1] ≠ B[q]
8:          k = π[k]
9:          if B[k + 1] == B[q]
10:             k = k + 1
11:         π[q] = k
12:     end for
13:     return π
14: end function
```

Assume that length
we have already f

$k$; to compute $\pi[q]$ we check if $B[q] = B[k + 1]$; if true then $\pi[q] = k + 1$; if not true then we find $\pi[k] = p$; if now $B[q] = B[p + 1]$ then $\pi[q] = p + 1$.

- We can now do our search for string $B$ in a longer string $A$:

```
 1: function KMP − Matcher(A, B)
 2:     n    length[A]
 3:     m
 4:     π
 5:     q
 6:     for
 7:         while q > 0 and B[q + 1] ≠ A[i]
 8:         q = π[q]
 9:         if B[q + 1] === A[i]
10:         q           1
11:         if q == m
12:             print pattern occurs with shift i − m
13:             q = π[q]
14:     end for
15: end function
```

- Sometimes we are not interested in finding just the prefect matches, but also in matches that might have a few errors, such as a few insertions, deletions and replacements.

- So ass
  $A =$
  $B =$
  in fin

- Idea: split $B$ into $k + 1$ consecutive subsequen
  length. Then any match in $A$ with at most
  subsequence which is a perfect match for a subsequenc k
  for all perfect matches for all of $k + 1$ subseq
  test by brute force if the remaining parts of $B$ have sufficient number of
  matches in the appropriate parts of $A$.

Assignment Project Exam Help

On a rectangular table there are 25 non-overlapping round coins of equal size

coin with https://eduassistpro.github.i

falling off t

within the table). Show that it is possible to complete

with 100 coins (of course with overlapping of coins).

Add WeChat edu_assist_pr