

Informatik II - SS 2014

(Algorithmen & Datenstrukturen)

Vorlesung 21 (29.7.2014)

String Matching (Textsuche) II



**UNI
FREIBURG**

Fabian Kuhn

Algorithmen und Komplexität

Textsuche / String Matching

Gegeben:

- Zwei Zeichenketten (Strings)
- Text T (typischerweise lang)
- Muster P (engl. pattern, typischerweise kurz)

Ziel:

- Finde alle Vorkommen von P in T

Annahmen:

- Länge Text T : n , Länge Muster P : m

Beispiel:

- Text: dubadubadudadubidubadubidubiduda
- Pattern: dubidu

Naiver Algorithmus

```
TestPosition(s):           // tests if  $T[s, \dots, s + m - 1] == P$   
     $t := 0$   
    while  $t < m$  and  $T[s + t] = P[t]$  do  
         $t := t + 1$   
    return  $(t = m)$             $O(n \cdot m)$ 
```

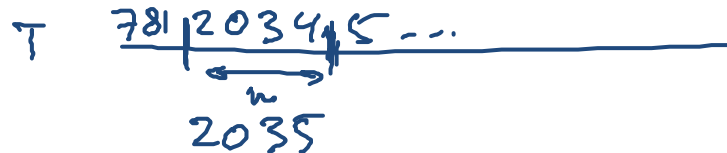
String-Matching:

```
for  $s := \underline{0}$  to  $\underline{n - m}$  do  
    if TestPosition(s) then  
        report found match at position  $s$ 
```

Rabin-Karp Algorithmus

Grundidee

- Wir schieben wieder ein Fenster der Grösse m über den Text und schauen an jeder Stelle, ob das Muster passt



- Zur Einfachheit nehmen wir an, dass der Text nur aus den Ziffern 0, ..., 9 besteht
 - dann können wir das Muster und das Fenster als Zahl verstehen
- Wenn wir das Fenster eins nach rechts schieben, kann die neue Zahl einfach aus der alten berechnet werden

$\overline{2034}5$
m
 $2034 - 2 \cdot 10^3 + 5$
 $\rightarrow 0345$
Zahlen vergleichen kostet $O(m)$ Zeit

Lösung von Rabin und Karp:

- Wir rechnen alles mit den Zahlen modulo M
 - M sollte möglichst gross sein, allerdings klein genug, damit die Zahlen $0, \dots, M - 1$ in einer Speicherzelle (z.B. 32 Bit) Platz haben
- Muster und Textfenster sind dann beides Zahlen aus dem Bereich $(M=100)$ $\{0, \dots, M - 1\}$
 $2034 \bmod 100 = 34$
- Beim Schieben des Fensters um eine Stelle, lässt sich die neue Zahl wieder in $O(1)$ Zeit berechnen
 - Falls das nicht klar ist, siehe spätere Folie...
- Falls das Muster gefunden wird, sind die zwei Zahlen gleich, falls nicht, können sie trotzdem gleich sein
 - Falls die Zahlen gleich sind, dann überprüfen wir nochmals wie beim naiven Algorithmus Buchstabe für Buchstabe

Rabin-Karp Algorithmus: Pseudo-Code

Text $T[0 \dots n - 1]$, Muster $P[0 \dots m - 1]$, Basis b , Modulus M

$h := \underline{b^{m-1}} \underline{\text{mod } M}$

$p := 0; t := 0;$

for $i := 0$ **to** $m - 1$ **do**

$p := (p \cdot b + P[i]) \text{ mod } M$

$t := (t \cdot b + T[i]) \text{ mod } M$

Zahl, welche das Muster repräsentiert
1. Textfenster

$s := 0;$

while $s \leq n - m$ **do**

if $p = t$ **then**

 TestPosition(s)

$t := \underline{((\underline{t} - \underline{T[s]} \cdot \underline{h}) \cdot \underline{b} + \underline{T[s + m]}) \text{ mod } M}$

Fenster um 1 nach rechts schieben

Rabin-Karp Algorithmus: Laufzeit

Vorberechnung:

Muster a^m $\underbrace{a \dots a}_m$
Text a^n $\underbrace{a \dots a}_n$ $\Theta(n \cdot m)$

Im schlechtesten Fall:

- Der schlechteste Fall tritt ein, falls die Zahlen in jedem Schritt übereinstimmen. Dann muss man in jedem Schritt Buchstabe für Buchstabe überprüfen, ob man das Muster wirklich gefunden hat.
 - Sollte bei guter Wahl von M nicht allzu oft geschehen...
 - ausser, wenn das Muster tatsächlich sehr oft ($\Theta(n)$ mal) vorkommt...

Im besten Fall:

- Im besten Fall sind die Zahlen nur gleich, falls das Muster auch wirklich gefunden wird. Die Kosten sind dann $O(\underline{n} + \underline{k} \cdot \underline{m})$, falls das Muster im Text k Mal vorkommt.

Wahl der Parameter...

Zahlendarstellung und Wahl von M

- Wir hätten gerne, dass wenn $x \neq y$, dann ist $h(x) = h(y)$ “unwahrscheinlich” (für $h(x) := x \bmod M$)
- Nehmen wir an, dass die Buchstaben in Muster und Text als Ziffern zur Basis b dargestellt werden
 - in unserem Fall, haben wir $b = 10$
- Falls b und M einen gemeinsamen Teiler haben, ist $h(x) = h(y)$ trotz $x \neq y$ nicht so unwahrscheinlich

Extremfall: M ist ein Teiler von b

$$3 \cdot 10 \bmod 5 = 0$$

Bsp: $b = 10, M = 5$

Pattern $\alpha_{n-1}, \alpha_{n-2}, \dots, \alpha_0 = \sum \alpha_i \cdot 10^i$

$$\left(\sum_{i=0}^{n-1} \alpha_i \cdot 10^i \right) \bmod 5 = \sum_{i=0}^{n-1} \underbrace{(\alpha_i \cdot 10^i \bmod 5)}_{0 \text{ für alle } i \neq 0} = \alpha_0 \bmod 5$$

Zahendarstellung und Wahl von M

- Wir hätten gerne, dass wenn $x \neq y$, dann ist $h(x) = h(y)$ “unwahrscheinlich” (für $h(x) := x \bmod M$)
- Nehmen wir an, dass die Buchstaben in Muster und Text als Ziffern zur Basis b dargestellt werden
 - in unserem Fall, haben wir $b = 10$
- Falls b und M einen gemeinsamen Teiler haben, ist $h(x) = h(y)$ trotz $x \neq y$ nicht so unwahrscheinlich

Wir wählen deshalb

- Die Basis b als genug grosse Primzahl
 - bei ASCII-Zeichen muss $b > 256$ sein
- M kann dann beliebig gewählt werden, am besten als Zweierpotenz
 - Zwischenresultate sind $< M \cdot b$, das sollte also in 32 (64) Bit Platz haben

Rechnen Modulo m

$$x \bmod M = y \Leftrightarrow \exists q \in \mathbb{Z}: y = x + q \cdot m \wedge y \in \{0, \dots, M - 1\}$$

- $x \bmod M$: addiere/subtrahiere M von x bis die Zahl im Bereich $\{0, \dots, M - 1\}$ ist

Rechenregeln:

$$(a \cdot b) \bmod m = ((a \bmod m) \cdot (b \bmod m)) \bmod m$$

$$(a + b) \bmod m = ((a \bmod m) + (b \bmod m)) \bmod m$$

Rechnen Modulo m

$$x \bmod M = y \Leftrightarrow \exists q \in \mathbb{Z}: y = x + q \cdot m \wedge y \in \{0, \dots, M - 1\}$$

- $x \bmod M$: addiere/subtrahiere M von x bis die Zahl im Bereich $\{0, \dots, M - 1\}$ ist

Rechenregeln:

$$\begin{aligned} & (a \cdot b) \bmod m = ((a \bmod m) \cdot (b \bmod m)) \bmod m \\ & (a + b) \bmod m = ((a \bmod m) + (b \bmod m)) \bmod m \end{aligned}$$

Schieben des Fensters:

- Fenster von Stelle s nach Stelle $s + 1$ schieben

$$t := ((\underset{\uparrow}{t} - T[s] \cdot \underset{\uparrow}{h}) \cdot b + T[s + m]) \bmod M$$

$h = b^{m-1} \bmod M$

Rechnen Modulo m

$$17 \bmod 5 = \\ 17 = 2 + 3 \cdot 5$$

$$x \bmod M = y \Leftrightarrow \exists q \in \mathbb{Z}: y = \underline{x} + \underline{q \cdot M} \wedge y \in \underline{\{0, \dots, M - 1\}}$$

Negative Zahlen

- Damit ist $x \bmod M$ immer im Bereich $\{0, \dots, M - 1\}$

Beispiele:

$$24 \bmod 10 = 4, \quad 4 \bmod 10 = 4, \quad -4 \bmod 10 = \underline{6}$$

$$-4 = 6 + (-1) \cdot 10$$

- **Aber:** In Java / C++ / Python ist $(-x) \% m = -(x \% m)$

Beispiele:

$$24 \% 10 = 4, \quad 4 \% 10 = 4, \quad -4 \% 10 = \underline{\underline{-4}}$$

- **Workaround:** Falls das Resultat von $x \% m$ negativ ist, einfach m dazuaddieren, dann kommt man in den richtigen Bereich

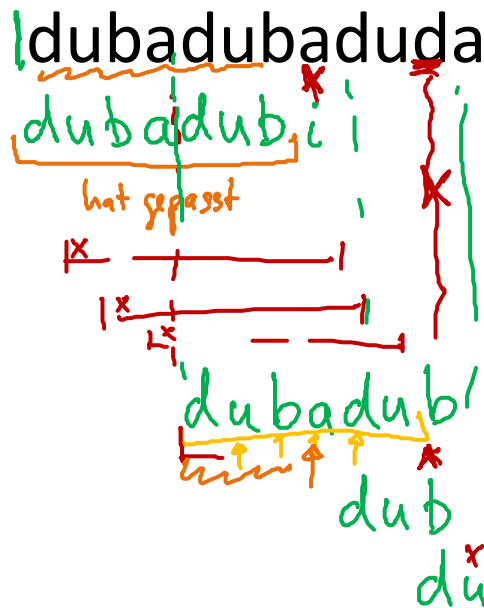
Algorithmus von Knuth, Morris, Pratt

- Kann wir das Problem immer in Zeit $O(n)$ lösen?
 - im schlechtesten Fall...

Schauen wir uns nochmals ein Beispiel an:

- Pattern: dubadubi

Text: |dubadubadudadubidubadubidubiduda



Knuth-Morris-Pratt Algorithmus

Idee:

- Falls wir beim Testen des Musters P an Stelle t feststellen, dass $P[t]$ nicht mit dem Text an der entsprechenden Stelle übereinstimmt, dann wissen wir, dass die Stellen $P[0 \dots t - 1]$ übereingestimmt haben.
- Das können wir bei der weiteren Suche ausnutzen

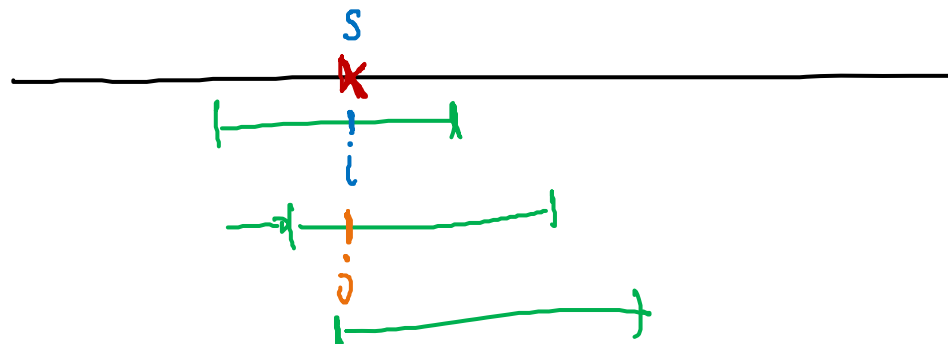
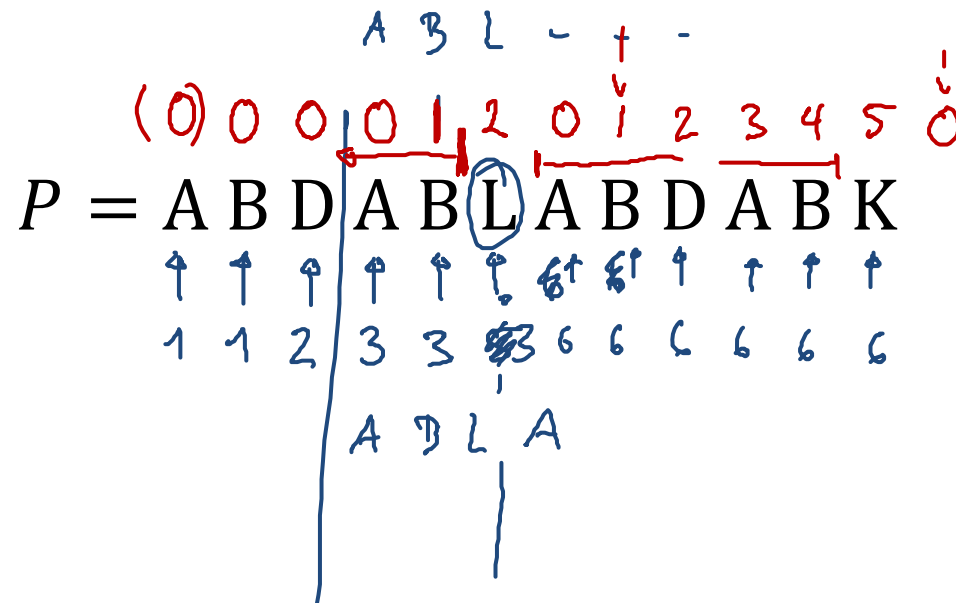
Beispiel: $P = \text{ABDABL} \boxed{\text{ABDABK}}$

ABDABL ABDABK
 Text... ABDABX
 * ABDABL ...
 ↑
 ... ABDABL ...
~~ABDAB~~ ABDA...

A B C D E F
 . . . A B C
 ABDABL ABDABK

Knuth-Morris-Pratt Alg.: Initialisierung

- Wir merken uns an jeder Stelle des Musters, wie weit wir das Suchfenster bei einem "Mismatch" weiterschieben können.



Knuth-Morris-Pratt Algorithmus

Vorbereitung: Array S der Länge $m + 1$

- $S[i]$: Stelle in P , an welcher man die neue Suche beginnt, falls beim Testen der Stelle i im Pattern ein Mismatch auftritt
- $S[0] = -1$, $S[1] = 0$
- $S[m]$: Stelle in P , an welcher man weitersucht, nachdem P erfolgreich gefunden wurde

Beispiel:

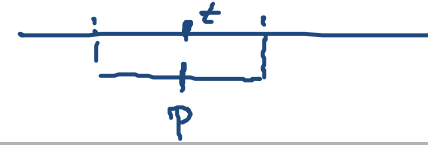
$P = [\underline{A, B, D, A, B, L, A, B, D, A, B}, D]$

$S = [-1, 0, 0, 0, 1, 2, 0, 1, 2, 3, 4, 5, 3]$

S hängt nur von P_{ab}

A B D A B ...
↑ A B D ...

Knuth-Morris-Pratt Algorithmus



t := 0; p := 0 // t: Position in Text, p: Position im Pattern

while $t < n$ do

if $T[t] = P[p]$ then // characters match

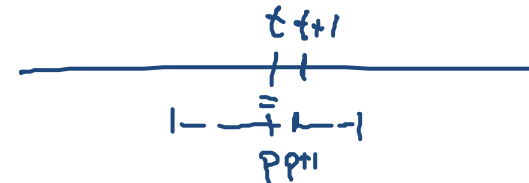
if $p = m - 1$ then // pattern found

pattern found at position $t - m + 1$

$p := S[m]$; $t := t + 1$

else

$p := p + 1$; $t := t + 1$



else // characters don't match

if $p = 0$ then // mismatch at first character

$t := t + 1$

else

$p := S[p]$

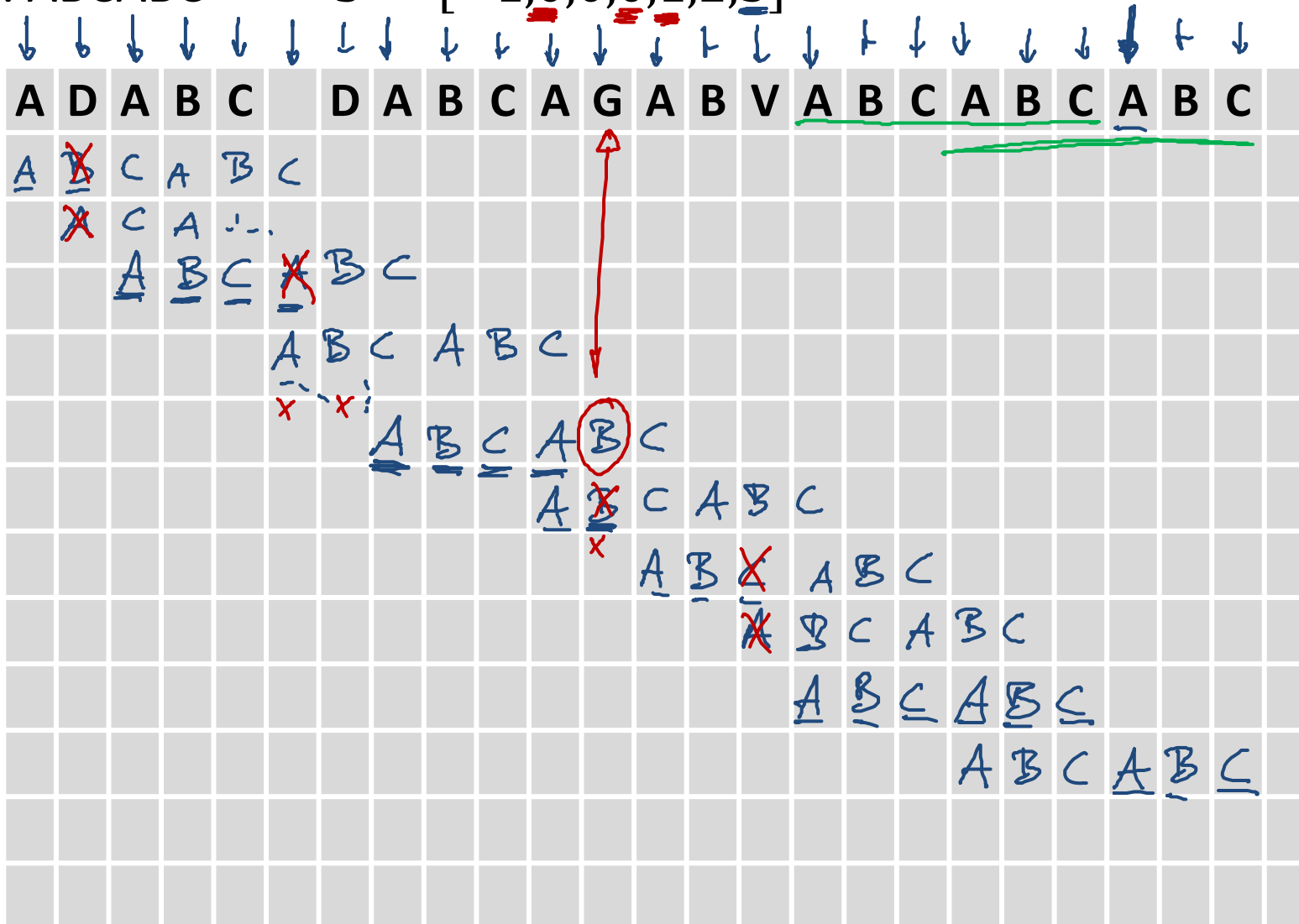


Knuth-Morris-Pratt Alg.: Beispiel

Pattern: ABCABC

$S = [-1, 0, 0, 0, 1, 2, 3]$

Text:



Knuth-Morris-Pratt Alg.: Laufzeit

Laufzeit ohne Initialisierung des Arrays S:

$t := 0$; $p := 0$

while $t < n$ **do**

if $T[t] = P[p]$ **then**

if $p = m - 1$ **then**

pattern found

$p := S[m]$; $t := t + 1$

else

$p := p + 1$; $t := t + 1$

else

if $p = 0$ **then**

$t := t + 1$

else

$p := S[p]$

$S[p] < P$

$\Rightarrow \leq 2n$
Schleifendurchläufe
 \sum
Laufzeit: $O(n)$

$\leq n$ mal



kann auch $\leq n$ mal vorkommen

Fenster wird geschoben (nach rechts)

Initialisierung

Vorbereitung von Array S :

- $P = [A, B, D, A, B, L, A, B, D, A, B, D]$
 $S = [-1, 0, 0, 0, 1, 2, 0, 1, 2, 3, 4, 5, 3]$
- An Position in $S[i]$ (für $i \in \{2, \dots, m\}$) steht

$$S[i] := \max_{k < i} \{P[i - k \dots i - 1] = P[0 \dots k - 1]\}$$

- $S[i]$: Länge des längsten echten Teilstückes von $P[0 \dots i - 1]$, welches an Stelle $i - 1$ endet, und welches auch Anfangsstück von P ist

Berechnung von $S[i]$:

- Falls $P[S[i - 1]] = P[i - 1]$, dann ist $S[i] = S[i - 1] + 1$
- Sonst testen, ob es einen kürzeres, passendes Anfangsstück gibt
 - Wir werden gleich anschauen, wie man das macht...

Berechnung von $S[i]$: Beispiel

$h := S[i - 1]$ $i \geq 2$ $S[0] = -1, S[1] = 0$
 while $h \geq 0$ do

 if $P[i - 1] = P[h]$ then
 $S[i] := h + 1; h := -2$
 else
 $h := S[h]$

if $h = -1$ then $S[i] = 0$ $i=4$ $i=5$ $i=12$

Beispiel: $P = [A, B, D, A, B, L, A, B, D, A, B, D, X]$

$S = [-1, 0, 0, 0, 1, 2, 0, 1, 2, 3, 4, 5, \dots]$

$i=2$ $P[1] \stackrel{?}{=} P[0]$
 (= "B") (= "A")
 $h = S[0] = -1$

$i=4$ $h = S[3] = 0$
 if $P[3] = P[0]$
 "A" "A"
 $\hookrightarrow S[4] = 0 + 1 = 1$

$i=5$ $h = S[4] = 1$
 if $P[4] = P[1]$
 "B" "B"
 $S[5] = 2$

$i=12$
 $h = S[11] = 5$
 $P[11] = P[5]$
 "D" \neq "L"
 $h = S[5] = 2$
 $P[11] = P[2]$
 "D" = "D" $\Rightarrow S[12] = 2 + 1 = 3$

Berechnung von $S[i]$: Laufzeit

```

h := S[i - 1]
while h ≥ 0 do
  if P[i - 1] = P[h] then
    S[i] := h + 1; h := -1
  else
    h := S[h]
if h = -1 then S[i] = 0

```

Beobachtung:

$$\underline{S[i]} \leq \underline{S[i - 1]} + 1$$

Falls $S[i] = S[i - 1] + 1$: 1 Schleifendurchlauf

Falls $S[i] < S[i - 1]$:

- Wert von h nimmt in jedem Schleifendurchlauf ab
- Am Schluss ist $S[i] = h + 1$
- Anzahl Schleifendurchläufe $\leq \Delta h + 1 = S[i - 1] - S[i] + 2$

$$\underbrace{\quad}_{S[i-1] - (S[i] - 1)}$$

Berechnung von $S[i]$: Laufzeit

Falls $S[i] = S[i - 1] + 1$:

- Anzahl Schleifendurchläufe = $1 = S[i - 1] - S[i] + 2$

Falls $S[i] < S[i - 1]$:

- Anzahl Schleifendurchläufe $\leq \Delta h + 1 = S[i - 1] - S[i] + 2$

Gesamtlaufzeit:

$$\sum_{i=2}^m (S[i-1] - S[i] + 2) = (m-1) \cdot 2 + (S[1] - S[2] + S[2] - S[3] + \dots)$$

$$= (m-1) \cdot 2 + \underbrace{S[1]}_{\geq 0} - \underbrace{S[m]}_{\geq 0}$$

$$\leq \underline{2(m-1)}$$

Knuth-Morris-Pratt Algorithmus:

- Berechnet zuerst in Zeit $O(m)$ das Array S der Länge m
 - hängt nur vom Pattern P ab
 - beschreibt an jeder Position im Pattern, wo (im Pattern) man bei einem Mismatch weitersuchen muss
- Mit Hilfe von S werden dann alle Vorkommen von P in T in Zeit $O(n)$ gefunden
 - In jedem Schritt kann man entweder die aktuelle Suchposition in T oder die Position des Suchfensters in T um mindestens 1 nach rechts verschieben

Gesamtlaufzeit: $O(m + n) = O(n)$