



# Informatik II Precomputation



G. Zachmann  
Clausthal University, Germany  
[zach@in.tu-clausthal.de](mailto:zach@in.tu-clausthal.de)

## Definitionen

- **Preconditioning (Vorbehandlung):**  
Aufgabe: löse nacheinander mehrere Probleme, die einen gemeinsamen Anteil haben, d.h., die Eingaben sind von der Form
 
$$\langle x_1; y \rangle \langle x_2; y \rangle \langle x_3; y \rangle \dots$$
 Lösungsansatz: transformiere zunächst den gemeinsamen Anteil, d.h., berechne
 
$$\hat{y} = T(y)$$
 und löse dann die Probleme
 
$$\langle x_1; \hat{y} \rangle \langle x_2; \hat{y} \rangle \langle x_3; \hat{y} \rangle \dots$$
 (hoffentlich effizienter)



G. Zachmann Informatik 2 — SS 10 Preprocessing 2

- **Precomputation (Vorbereitung):**  
 auch wenn die Eingabe nur einmal vorkommt – aber dennoch aus 2 Teilen  $\langle x ; y \rangle$  besteht – so kann es effizienter sein, zunächst einen Teil zu transformieren,  $\hat{y} := T(y)$ , und dann die Lösung aus  $\langle x ; \hat{y} \rangle$  zu berechnen
  - Beispiel: dynamischen Text nach einem Wort durchsuchen
- Bemerkung: im Umgangssprachgebrauch wird "*Precomputation*" oft für beide Arten verwendet
- Wesentliche Frage bei Precomputation / Preconditioning ist immer: Wie hoch ist der Aufwand dafür? macht er evtl. den anschließenden Zeitgewinn kaputt?

G. Zachmann Informatik 2 — SS 10

Preprocessing 3

## Beispiele bisher

- Closest Pairs
  - Vorsortieren der Punkte entlang x-Achse
- Huffman-Kodierung:
  - Bestimmung der (relativen) Häufigkeiten der einzelnen Zeichen

G. Zachmann Informatik 2 — SS 10

Preprocessing 4

## Anwendung dieser Technik im Folgenden

- Auswertung eines Polynoms an vielen Stellen;  
Eingabe ist  $\langle x_i; \underbrace{a_0, \dots, a_n}_{\text{„konstant“}} \rangle$
- Vorgänger-/Nachfolger-Beziehung in einem Baum testen;  
Eingabe ist  $\langle v_1, v_2; \underbrace{\text{Baum}}_{\text{„konstant“}} \rangle$
- Gegebenen Text nach einem Wort durchsuchen

G. Zachmann Informatik 2 — SS 10 Preprocessing 5

## Der Vorgänger im Baum

- Problem:
  - Gegeben ist ein Baum
  - Für viele Paare von Knoten  $(v, w)$  sollen wir entscheiden, ob  $v$  Vorgängerin, d.h. direkte oder indirekte Mutter, von  $w$  ist
- Naive Lösung:
  - Traversiere den Baum von  $w$  aus bis zur Wurzel und vergleiche mit  $v$
  - Aufwand:
    - $O(n)$  im Worst-Case,
    - Untere Schranke für Worst-Case =  $\Omega(\log n)$
- Behauptung: wir können auf dem Baum *Preconditioning* durchführen in Zeit  $\Theta(n)$ , so daß danach das Problem in Zeit  $O(1)$  gelöst werden kann

G. Zachmann Informatik 2 — SS 10 Preprocessing 6

Das Verfahren:

1. Durchlaufe den Baum in *Preorder* und nummeriere dabei die Knoten; bezeichne diese Nummer mit  $N_{\text{pre}}(v)$ 
  - Nummeriere also erst den Knoten selbst mit fortlaufender Nummer, dann nummeriere rekursiv den linken, dann den rechten Teilbaum
2. Nummeriere analog die Knoten in Postorder; Bezeichnung  $N_{\text{post}}(v)$
3. Jetzt gilt für jedes Paar  $v, w$  von Knoten:
  1.  $N_{\text{pre}}(v) \leq N_{\text{pre}}(w) \Leftrightarrow v$  Vorgänger von  $w$  oder  $v$  links von  $w$
  2.  $N_{\text{post}}(v) \geq N_{\text{post}}(w) \Leftrightarrow v$  Vorgänger von  $w$  oder  $v$  rechts von  $w$
4. Zusammen :
 
$$N_{\text{pre}}(v) \leq N_{\text{pre}}(w) \wedge N_{\text{post}}(v) \geq N_{\text{post}}(w) \Leftrightarrow \text{Vorgänger von } w$$

G. Zachmann Informatik 2 — SS 10 Preprocessing 7

Beispiel

```

graph TD
    A((1 A 13)) --> B((2 B 5))
    A --> C((7 C 12))
    B --> D((3 D 1))
    B --> E((4 E 3))
    B --> F((6 F 4))
    E --> I((5 I 2))
    C --> G((8 G 6))
    C --> H((9 H 11))
    H --> B7((10 B 7))
    H --> B8((11 B 8))
    H --> B9((12 B 9))
    H --> B10((13 B 10))
  
```

G. Zachmann Informatik 2 — SS 10 Preprocessing 8

## Wiederholte Auswertung eines Polynoms

- Gegeben: Polynom  $p(x) = a_n x^n + \dots + a_1 x + a_0$
- Gesucht ist:  $p(x_i)$  für viele  $x_1, \dots, x_n$
- **Verboten** ist:
 

```
p = a[0] + a[1]*x + a[2]*x*x + ...
```
- **Noch schlimmer ist!:**

```
p = a[0] + a[1]*x + a[2]*pow(x,2.0)
    + a[3]*pow(x,3.0) ...
```

G. Zachmann Informatik 2 — SS 10 Preprocessing 9

- Einfache Methode: transformiere  $p$  in das sog. **Horner-Schema**

$$p(x) = ((\dots((a_n x + a_{n-1})x + a_{n-2}) \dots)x + a_1) + a_0$$

```
p = a[n]
for j in range( n-1, -1, -1 ):
    p = p*x + a[j]
```
- Reicht für mittellange Polynome völlig aus
- Aufwand:  $n$  Multiplikationen /  $n$  Additionen

G. Zachmann Informatik 2 — SS 10 Preprocessing 10

## Auswertung an äquidistanten Stellen

- Definition: **Vorwärts-Differenzen**

$$\Delta^1 p(x) := p(x+h) - p(x)$$

$$\Delta^i p(x) := \Delta^{i-1} p(x+h) - \Delta^{i-1} p(x)$$

$$\Delta^0 p(x) := p(x)$$
- Behauptung:  $\Delta^i p(x)$  ist ein Polynom vom Grad  $n-i$
- Beweis:
 
$$\begin{aligned} \Delta^1 p(x) &= p(x+h) - p(x) \\ &= a_0 + \dots + a_n(x+h)^n - a_0 - \dots - a_n x^n \\ &= a_0 + \dots + a_n \sum_{i=0}^n \binom{n}{i} x^i h^{n-i} - a_0 - \dots - a_n x^n \\ &= a'_0 + \dots + a'_{n-1} x^{n-1} \end{aligned}$$
- Rest per Induktion

G. Zachmann Informatik 2 — SS 10 Preprocessing 11

- Korollar:
 
$$\Delta^n p(x+jh) \equiv \text{const}$$
- Bemerkung:** auf dieser Beobachtung basiert das ganze Verfahren!
- Bemerkung:  $\Delta^i p(x)$  hängt ab von  $p(x)$ ,  $p(x+h)$ , ...,  $p(x+ih)$
- Beispiel:
 
$$\begin{aligned} \Delta^2 p(x) &= \Delta^1 p(x+h) - \Delta^1 p(x) \\ &= (\Delta^0 p(x+h+h) - \Delta^0 p(x+h)) \\ &\quad - (\Delta^0 p(x+h) - \Delta^0 p(x)) \\ &= p(x+2h) - 2p(x+h) + p(x) \end{aligned}$$

G. Zachmann Informatik 2 — SS 10 Preprocessing 12

## Das Verfahren

- Erstelle eine **Vorwärtsdifferenzen-Pyramide**:
  - Vereinfachende Schreibweise  $\Delta^i p(x + jh) =: \Delta^i p_j$
  - Die Pyramide:

$x+ih$	$p$	$\Delta^1$	$\Delta^2$	...	$\Delta^n$
$x$	$p_0$	$\Delta^1 p_0$	$\Delta^2 p_0$		$\Delta^n p_0$
$x+h$	$p_1$	$\Delta^1 p_1$	$\Delta^2 p_1$		$\Delta^n p_1$
$x+2h$	$p_2$	$\Delta^1 p_2$	$\Delta^2 p_2$		$\Delta^n p_2$
		$\vdots$	$\vdots$		$\vdots$
$x+nh$	$p_n$	$\Delta^1 p_{n-1}$	$\Delta^2 p_{n-1}$		$\Delta^n p_{n-1}$
$x+(n+1)h$	$p_{n+1}$	$\Delta^1 p_n$	$\Delta^2 p_n$		$\Delta^n p_n$

$$\Delta^{i+1} p_j = \Delta^i p_{j+1} - \Delta^i p_j$$

$$\Delta^i p_{j+1} = \Delta^i p_j + \Delta^{i+1} p_j$$

G. Zachmann Informatik 2 — SS 10 Preprocessing 13

- Algo zur Berechnung von vielen äquidistanten Werten von  $p$ :
  1. Initialisieren die Pyramide mit den Punkten  $x, \dots, x + nh$
  2. Berechne den nächsten Punkt an der Stelle  $x + (n+1)h$  durch Fortsetzen der Pyramide um eine Zeile am unteren Ende
- Aufwand: für jeden neuen Punkt braucht man nur  $n$  Additionen (keine Multiplikation!)
- Bemerkungen:
  - Man muß immer nur die unterste Zeile der Pyramide speichern
  - Rundungsfehler akkumulieren sich!  $\rightarrow$  Pyramide ab und zu neu aufsetzen
- Durch Umformung kann man ein Polynom vom Grad  $n$  mit  $\frac{n}{2}$  Multiplikationen an beliebigen Stellen auswerten (s. Knuth)

G. Zachmann Informatik 2 — SS 10 Preprocessing 14

## Suchen in Texten (*String Matching*)

- Aufgabe des String-Matching-Algorithmus
  - kinderleicht
- Naiver Algorithmus
  - Wie würden Sie es tun?
- Knuth-Morris-Pratt-Algorithmus
  - Scharfes Anschauen (= Precomputation) des Musters verbessert die Laufzeit
- Boyer-Moore-Algorithmus
  - "Schlechte" Zeichen erlauben uns, größere Sprünge im Text zu machen
  - Das ist noch besser als "nur" worst-case optimal (in der Praxis)

G. Zachmann Informatik 2 — SS 10 Preprocessing 15

## Die Aufgabe

- Gegeben:
  - Text  $T$  der Länge  $n$  über einem endlichen Alphabet  $\Sigma$ :
 

$T[1]$   $T[n]$   

m	a	n	a	m	a	n	a	p	a	t	i	p	i	t	i	p	i
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---
  - Muster (Pattern)  $P$  der Länge  $m$  über dem selben Alphabet  $\Sigma$ :
 

$P[1]$   $P[m]$   

p	a	t	i
---	---	---	---
- Ausgabe: jedes **Vorkommen** von  $P$  in  $T$ :
 

$T[1+s, \dots, m+s] = P[1..m]$   

m	a	n	a	m	a	n	a	p	a	t	i	p	i	t	i	p	i
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

← Shift  $s$  →

p	a	t	i
---	---	---	---
- Definition: als **(Mis-)Match** wird die (Nicht-)Übereinstimmung von einem Zeichen aus dem Muster mit einem Zeichen im Text bezeichnet

G. Zachmann Informatik 2 — SS 10 Preprocessing 16





## Das naive Verfahren

- Für jede mögliche Verschiebung  $0 \leq i \leq n - m$  prüfe maximal  $m$  Zeichenpaare, bei Mismatch beginne mit nächster Verschiebung

```

# Input: Text T und Muster P
# Output: Liste L mit Verschiebungen i,
#         an denen P in T vorkommt
# Bug: Indizierung von T/P beginnt bei 1
def naive_string_match( T, P ):
    L = []
    for s in range( 0, len(T) - len(P) + 1 ):
        j = 1
        while j <= m and T[s+j] == P[j]:
            j += 1
        if j == m + 1:
            L.append( s )
    return L

```

G. Zachmann Informatik 2 — SS 10 Preprocessing 19

## Beispiel

```

T="A string consisting of 37 characters."
P="sting"

```

A string consisting of 37 characters.

sting

  sting

    sting ...

                  sting

                  sting

                  sting ✓ ...

                  s=14

G. Zachmann Informatik 2 — SS 10 Preprocessing 20

## Aufwand

- Worst-case Beispiel:
 

0	0	...	0	...	0	...	0	0	...
			← s						
			0	...	0	...	0	1	
- Fazit: benötigt im worst-case die Laufzeit
 
$$(n - m + 1) \cdot m \in O(n \cdot m)$$
  - Für  $m \in O(n)$  ist das  $O(n^2)$
  - In der Praxis oft tritt oft ein Mismatch sehr früh auf  $\rightarrow$  Laufzeit  $\sim c \cdot n$
- Untere Schranke für die Textsuche im Worst-Case:**
  - Mind. 1 Vergleich pro  $m$  aufeinander folgende Textzeichen
  - Mind. jedes Zeichen des Patterns muß 1x betrachtet werden
  - $\rightarrow \Omega(m + \frac{n}{m})$

G. Zachmann Informatik 2 — SS 10 Preprocessing 21

## Das Verfahren nach Knuth-Morris-Pratt (KMP)

- Tritt bei einer Verschiebung erstmals bei  $T_i$  und  $P_{j+1}$  ein Mismatch auf, dann gilt:
  - Die zuletzt verglichenen  $j$  Zeichen in  $T$  stimmen mit den ersten  $j$  Zeichen in  $P$  überein
  - $T_i \neq P_{j+1}$

$T_1$	$T_2$	...	...	$T_i$	...	...
						⋈
			$P_1 \quad \dots \quad P_j \quad P_{j+1} \quad \dots \quad P_m$			

G. Zachmann Informatik 2 — SS 10 Preprocessing 22

- Die Idee: bestimme  $j' < j$ , so daß
  - $P[1..j'] = P[j'+1..j] = T[i'..i-1]$

$T_1$	$T_2$	...	...	$T_i$	...	...	
						⋈	
			$P_1 \quad \dots \quad P_j \quad P_{j+1} \quad \dots \quad P_m$				
					⋈		
			$P_1 \dots P_{j'} \quad P_{j'+1} \quad \dots \quad P_m$				

- Vorteil: man kann sofort  $P_{j'+1}$  mit  $T_i$  vergleichen
- M.a.W.: bestimme den **längsten Präfix von  $P$ , der echtes Suffix von  $P[1..j]$  ist**
- Speichere für jedes  $j$  das entsprechende  $j'$  in  $j' = \text{next}[j]$

G. Zachmann Informatik 2 — SS 10
Preprocessing 23

- Beispiel für die Bestimmung von  $\text{next}[j]$ :

$T_1$	$T_2$	...	0 1 0 1 1	0 1 0 1 1	0	...
			0 1 0 1 1	0 1 0 1 1	1	
				0 1 0 1 1	0 1 0 1 1	1

G. Zachmann Informatik 2 — SS 10
Preprocessing 24

- Für  $P = 0101101011$  ist  $next = [0,0,1,2,0,1,2,3,4,5]$

1	2	3	4	5	6	7	8	9	10	
0	1	0	1	1	0	1	0	1	1	← Pattern
		0								}
		0	1							
				0						
				0	1					
				0	1	0				
				0	1	0	1			

next[j] =  
 Länge des längsten Präfixes von  $P$ , das echtes Suffix von  $P[1..j]$  ist

G. Zachmann Informatik 2 — SS 10
Preprocessing 25

### Beispiel zum Vorgehen beim Matching

- Pattern: abrakadabra, next = [0,0,0,1,0,1,0,1,2,3,4]

i	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	
	a	b	r	a	k	a	d	a	b	r	a	b	r	a	b	a	b	r	a	k	...
											≠			≠	≠						
	a	b	r	a	k	a	d	a	b	r	a	k	r	a	k	a	b	k	a	k	
											$jj=11$				$j=0=2$						

next[2] = 0

G. Zachmann Informatik 2 — SS 10
Preprocessing 26

## Implementierung

```

# Input: Text T und Muster P
# Output: Liste L mit Verschiebungen s, an denen P in T vorkommt
def kmp_matcher( T, P ):
    n = len(T)
    m = len(P)
    L = []
    next = compute_next( P )
    j = 0
    for i in range(1,n):
        while j > 0 and T[i] != P[j+1] :
            j = next[j]
        if T[i] == P[j+1]:
            j += 1
        if j == m:
            L.append( i-m )
            j = next[j]
    return L

```

G. Zachmann Informatik 2 — SS 10 Preprocessing 27

## Korrektheit des Algorithmus

- Situation am Beginn der inneren while-Schleife:  
 $P[1..j] = T[i-j .. i-1]$  und  $j \neq m$

$T_1$	$T_2$	...	...	$T_i$	...	...
$P_1$	...	$P_j$	$P_{j+1}$	...	$P_m$	

- Falls  $j = 0 \rightarrow j$  steht vor dem erstem Zeichen von  $P$  (d.h., es gab noch keinen Vergleich zwischen  $P_1$  und  $T_i$ )
- Falls  $j > 0 \rightarrow P$  kann verschoben werden, solange  $j > 0$  und  $T_i \neq P_{j+1}$ 
  - Ist dann  $T[j] = P[j+1]$ , können  $j$  und  $i$  (am Schleifenende) erhöht werden
- Wurde ganz  $P$  verglichen ( $j = m$ ), ist eine Stelle gefunden, und es kann verschoben werden

G. Zachmann Informatik 2 — SS 10 Preprocessing 28

## Laufzeit

- Beobachtungen:
  - Textzeiger  $i$  wird nie zurückgesetzt
  - Textzeiger  $i$  und Musterzeiger  $j$  werden stets gemeinsam inkrementiert
  - Es gilt
 
$$\forall j : \text{next}[j] < j$$
    - $j$  kann, insgesamt über die ganze for-Schleife, nur so oft herabgesetzt werden, wie es heraufgesetzt wurde, also höchstens  $n$  Mal
- Fazit: der KMP-Algorithmus kann in Zeit  $O(n)$  ausgeführt werden, **wenn** das next-Array schon berechnet ist

```

for i in range(1,n):
    while j > 0 and \
          T[i] != P[j+1]:
        j = next[j]
    if T[i] == P[j+1]:
        j += 1
    if j == m:
        L.append( i-m )
        j = next[j]
return L
  
```

G. Zachmann Informatik 2 — SS 10
Preprocessing 29