

## Teilwortsuche mit dem Algorithmus von Rabin und Karp, Fingerabdrücke

$a = 55318654691459543965671439459743987463261911246129429129987539873497175\dots$

$b = 1415926$

$$a = a_1 a_2 \dots a_n \quad 0 \leq a_i, b_i \leq B-1 \quad \text{Bytes}$$

$$b = b_1 b_2 \dots b_m \quad B = 256$$

In unserem Beispiel  $B=10$  (der Einfachheit halber)

- GRUNDIDEE: Interpretiere die Zeichen  $a_i, b_i$  als Ziffern einer Zahl zur Basis  $B$ .

$$b \rightarrow \hat{b} = b_1 B^{m-1} + b_2 B^{m-2} + \dots + b_{m-1} B + b_m$$

$$a_i a_{i+1} \dots a_{i+m-1} \rightarrow \hat{a}_i = a_i B^{m-1} + a_{i+1} B^{m-2} + \dots + a_{i+m-1} \cdot 1$$

$$a_{i+1} \dots a_{i+m} \rightarrow \hat{a}_{i+1} = a_{i+1} B^{m-1} + \dots + a_{i+m-1} \cdot B + a_{i+m}$$

$$\hat{a}_{i+1} = (\hat{a}_i - a_i B^{m-1}) \times B + a_{i+m}$$

$$a_i a_{i+1} \dots a_{i+m-1} = b \iff \hat{a}_i = \hat{b} \iff \tilde{a}_i = \tilde{b}$$

als Zeichenketten  $\uparrow$   $\uparrow$  als Zahlen

Das Teilwortproblem kann in  $O(m+n)$  Rechenschritten gelöst werden.

Berechne  $\hat{b}, \hat{a}_1 \dots O(m)$  (Horner Schema!)  $\hat{a}_i \rightarrow \hat{a}_{i+1} O(1)$   
 $n \times \text{mal}$

• IDEE: Wir rechnen modulo  $Q$ .

$Q$  wird so gewählt, dass die Rechnungen mit Maschinenarithmetik durchgeführt werden können.

Bsp. long 64 Bit.  $B=256$   $Q \leq 2^{56}$

$$\tilde{b} := \hat{b} \bmod Q \quad \tilde{a}_i := \hat{a}_i \bmod Q$$

$$\tilde{a}_{i+1} = \left( \tilde{a}_i - a_i B^{m-1} \right) \times B + a_{i+m} \pmod{Q}$$

$$\hat{c} := B^{m-1} \bmod Q$$

$\underbrace{\quad\quad\quad}_{8} \quad \underbrace{\quad\quad\quad}_{56 \text{ bits}}$   
 $\quad\quad\quad \approx 64$   
 $\underbrace{\quad\quad\quad}_{56} \quad \quad \quad 8$   
 $\quad\quad\quad \underbrace{\quad\quad\quad}_{64}$

• Wahl von  $Q$ :  $Q = 2^{56}$  SCHLECHT!

$Q$  ... eine zufällige Primzahl mit  $2^{55} \leq Q \leq 2^{56}$

SATZ: Die Anzahl der Primzahlen  $\leq N$  ist  $\sim \frac{N}{\ln N}$ .

$\ln 2^{56} \approx 40$ . Es kommen ca.  $\frac{2^{55}}{40}$  Primzahlen  $Q$  in Frage.

$$\Pr[\text{Fehlalarm an Stelle } i] = \Pr[\hat{a}_i \neq \hat{b}, \text{ aber } \tilde{a}_i = \tilde{b}]$$

$$\hat{a}_i \equiv \hat{b} \pmod{Q} \Leftrightarrow Q \mid \hat{a}_i - \hat{b}$$

$$= \Pr[Q \mid \underbrace{\hat{a}_i - \hat{b}}_{8m \text{ bits}} \neq 0] \leq \frac{m/6}{2^{55}/40} \leq \frac{7m}{2^{55}} \leq \frac{m}{10^{15}}$$

$$\leq 2^{8m}$$

$$\hat{a}_i - \hat{b} \text{ hat } \leq \frac{\log 2^{8m}}{\log 2^{55}} \text{ verschiedene Primfaktoren } \geq 2^{55}$$

$$= \frac{8m}{55} \leq \frac{m}{6}$$

• Es gibt effiziente Primzahl tests. (randomisiert.)

Probiere wiederholt eine Zufallszahl im Bereich  $2^{55} < Q < 2^{56}$ , bis  $Q$  eine Primzahl ist.

Nach 40 Versuchen (im Mittel) hat man Erfolg.

abcacbacbacbacbac  
 bcacbacbacbacbac  
 accnacancnacdna  
 cadbadcbacdabda  
 acandadadcabcdabc  
 adbcdbcbdbdacdadna  
 aaandnaadcancdann

bac  
 caa  
 ncn

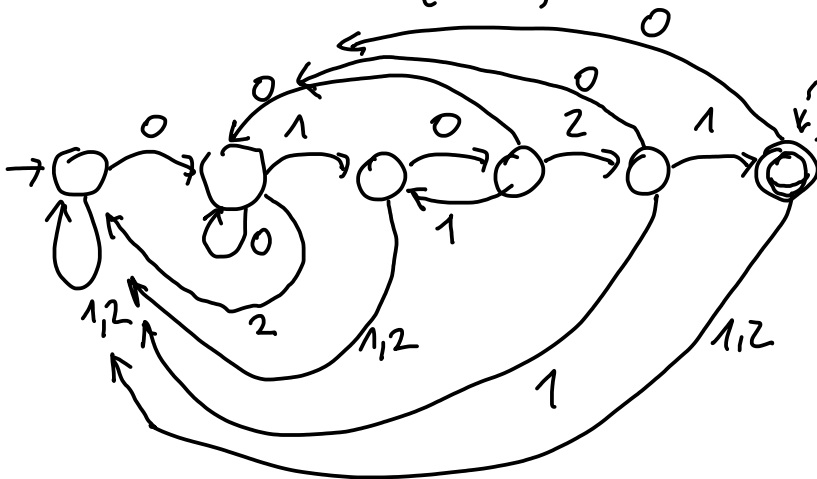
zweidimensionale  
 Musteruche

### Teilwortsuche mit endlichen Automaten

$$b \text{ kommt in } a \text{ vor} \iff a \in \underbrace{\Sigma^* b \Sigma^*}_{\text{regulär}}$$

$\Rightarrow$  Es gibt einen deterministischen endlichen Automaten, der diese Sprache akzeptiert.

$b = 01021$      $\Sigma = \{0,1,2\}$



$\Sigma^* b = \text{Sprache}$   
 Speicher  $|b| \cdot |\Sigma|$

für beliebige reguläre  
 Ausdrücke.