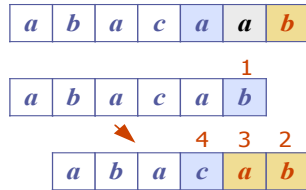


Strengsøgning



Streng



En **streng** er en sekvens af tegn

Lad P være en streng af længde m

Eksempler på streng:

- Java-program
- HTML-dokument
- DNA-sekvens
- Digitaliseret billede

En **delstreng** $P[i..j]$ af P er den delsekvens af P , der består af tegn med indeks imellem i og j

Et **prefix** af P er en delstreng af typen $P[0..i]$

Et **suffix** af P er en delstreng af typen $P[i..m-1]$

Et **alfabet** Σ er mængden af mulige tegn for en mængde af strenge

Lad der være givet en streng T (tekst) og en streng P (mønster).

Eksempler på alfabeter:

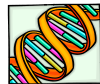
- ASCII
- Unicode
- {A, C, G, T}
- {0, 1}

Strengsøgning (pattern matching) består i at finde en delstreng i T , der er lig med P

Anvendelser:

- Tekstredigeringsprogrammer
- Søgemaskiner
- Biologisk forskning

Rå kraft algoritme



Den rå kraft strengsøgningsalgoritme sammenligner mønsteret P med teksten T for ethvert muligt skift af P i forhold til T , indtil enten

- der er fundet overensstemmelse
- alle mulige placeringer af mønsteret er blevet afprøvet

Algoritmens køretid er $O(nm)$

Eksempel på værste tilfælde:

$T = aaa \dots ah$
 $P = aaah$

forekommer ofte for billeder og DNA-sekvenser, men sjældent for almindelig tekst

Algorithm *BruteForceMatch*(T, P)

Input text T of size n and pattern P of size m

Output starting index of a substring of T equal to P or -1 if no such substring exists

for $i \leftarrow 0$ **to** $n - m$
 { test shift i of the pattern }

$j \leftarrow 0$

while $j < m \wedge T[i+j] = P[j]$

$j \leftarrow j + 1$

if $j = m$

return i {match at i }

return -1 {no match anywhere}

Boyer-Moore's heuristik

Boyer-Moore's strengsøgningsalgoritme er baseret på to heuristikker:

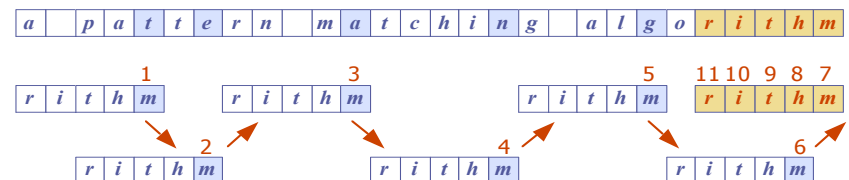
Spejl-heuristikken: Sammenlign P med en delsekvens af T baglæns

Hop-heuristikken: Når en tegnsammenligning mislykkes ved $T[i] = c$:

Hvis P indeholder tegnet c , så skift P således, at den sidste forekomst af c i P kommer ovenpå $T[i]$

Ellers, skift P således, at $P[0]$ kommer ovenpå $T[i+1]$

Eksempel:

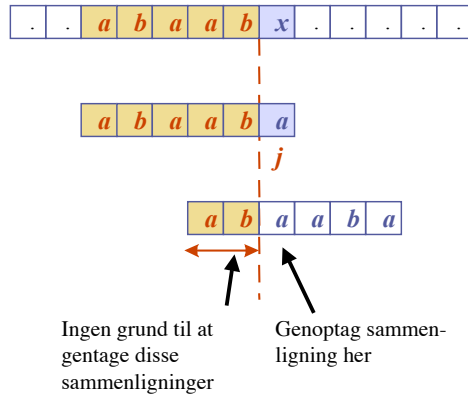


KMP-algoritmen (motivation)

Knuth-Morris-Pratt's algoritme sammenligner mønsteret i teksten fra venstre imod højre, men skifter mønsteret mere intelligent end rå kraft algoritmen (kaster ikke information væk)

Når en tegnsammenligning mislykkes, hvad er da det længste vi kan skifte mønsteret for at undgå overflødige sammenligninger?

Svar: det længste prefix af $P[0..j-1]$, som er suffix af $P[1..j-1]$



9

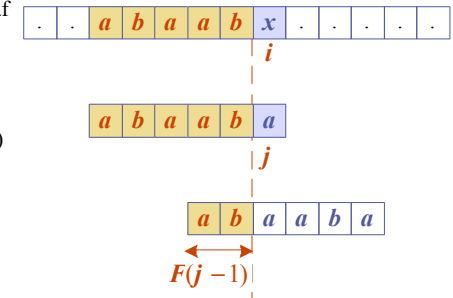
KMP's fejlfunktion

Knuth-Morris-Pratt's algoritme forbehandler mønsteret for at finde prefixer af mønsteret, der findes i mønsteret selv

j	0	1	2	3	4	5
$P[j]$	a	b	a	a	b	a
$F(j)$	0	0	1	1	2	3

Fejlfunktionen $F(j)$ defineres som længden af det længste prefix af $P[0..j]$, som er suffix af $P[1..j]$

Knuth-Morris-Pratt's algoritme ændrer på rå kraft algoritmen, så vi ved en mislykket sammenligning, $P[j] \neq T[i]$, sætter $j \leftarrow F(j-1)$



10

KMP-algoritmen

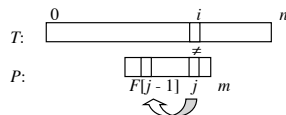
Fejlfunktionen kan repræsenteres ved et array

I hver iteration af while-løkken har vi, at

- i øges med 1, eller
- skiftets størrelse $i-j$ øges med mindst 1 (da $F(j-1) < j$)
- $i-j \leq n$

Der er således højst $2n$ iterationer af while-løkken

KMP-algoritmen kører i optimal tid, $O(m+n)$



Algorithm *KMPMatch*(T, P)

```

 $F \leftarrow failureFunction(P)$ 
 $i \leftarrow 0$ 
 $j \leftarrow 0$ 
while  $i < n$ 
  if  $T[i] = P[j]$ 
    if  $j = m - 1$ 
      return  $i - j$  { match }
    else
       $i \leftarrow i + 1$ 
       $j \leftarrow j + 1$ 
  else
    if  $j > 0$ 
       $j \leftarrow F[j - 1]$ 
    else
       $i \leftarrow i + 1$ 
return -1 { no match }
    
```

11

Beregning af fejlfunktionen



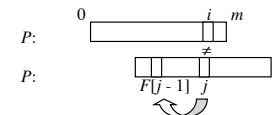
Fejlfunktionen kan repræsenteres ved et array og kan konstrueres i $O(m)$ tid

Konstruktionen svarer til KMP-algoritmen selv

I hver iteration af while-løkken har vi, at

- i øges med 1, eller
- skiftets størrelse $i-j$ øges med mindst 1 (da $F(j-1) < j$)
- $i-j \leq m$

Der er således højst $2m$ iterationer af while-løkken



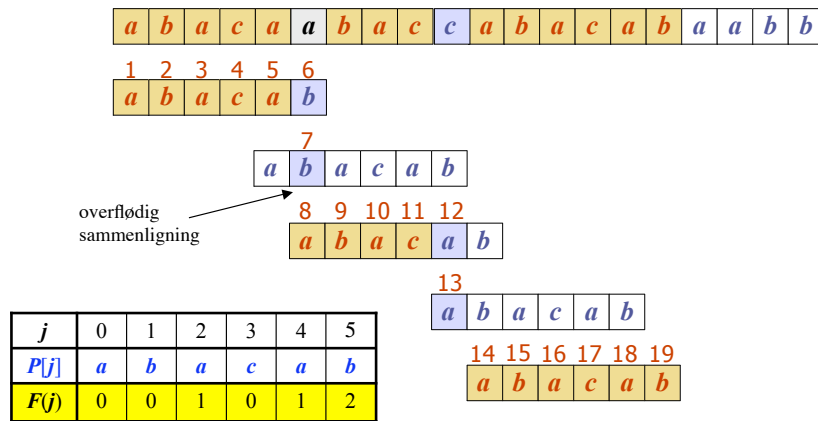
Algorithm *failureFunction*(P)

```

 $F[0] \leftarrow 0$ 
 $i \leftarrow 1$ 
 $j \leftarrow 0$ 
while  $i < m$ 
  if  $P[i] = P[j]$ 
    {we have matched j + 1 characters}
     $F[i] \leftarrow j + 1$ 
     $i \leftarrow i + 1$ 
     $j \leftarrow j + 1$ 
  else if  $j > 0$  then
    {use failure function to shift P}
     $j \leftarrow F[j - 1]$ 
  else
     $F[i] \leftarrow 0$  {no match}
     $i \leftarrow i + 1$ 
    
```

12

Eksempel



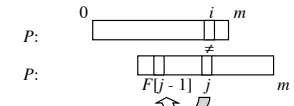
13

Forbedret beregning af fejlfunktionen



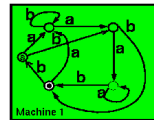
```

Algorithm failureFunction( $P$ )
 $F[0] \leftarrow 0$ 
 $i \leftarrow 1$ 
 $j \leftarrow 0$ 
while  $i < m$ 
  if  $P[i] = P[j]$ 
    {we have matched  $j + 1$  characters}
    if  $P[j] \neq P[j+1]$ 
       $F[i] \leftarrow j + 1$ 
    else
       $F[i] \leftarrow F[j + 1]$ 
     $i \leftarrow i + 1$ 
     $j \leftarrow j + 1$ 
  else if  $j > 0$  then
    {use failure function to shift  $P$ }
     $j \leftarrow F[j - 1]$ 
  else
     $F[i] \leftarrow 0$  {no match}
     $i \leftarrow i + 1$ 
    
```



14

Endelige tilstandsmaskiner



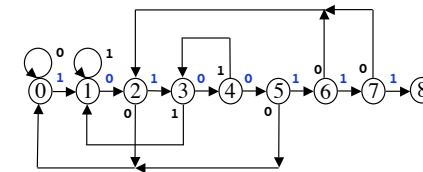
En **endelig tilstandsmaskine** er karakteriseret ved

1. en endelig mængde af tilstande
2. en starttilstand
3. en eller flere sluttillstande
4. en endelig mængde af inputsymboler
5. en funktion, *move*, der afbilder mængden af par bestående af et inputsymbol og en tilstand på mængden af tilstande

$\text{move}(\text{input}, \text{state}) \rightarrow \text{state}$

15

En endelig tilstandsmaskine til genkendelse af mønsteret 10100111



Starttilstand: 0

Sluttilstand: 8

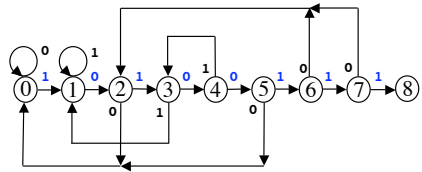
Søgeeksempel:

Tekst: 100111010010100010100111000111

Tilstand: 0120111234562345012345678

16

Repræsentation af en endelig tilstandsmaskine



En endelig tilstandsmaskine kan repræsenteres ved hjælp af en tabel

move-tabel:

input\state	0	1	2	3	4	5	6	7
0	0	2	0	4	5	0	2	2
1	1	1	3	1	3	6	7	8

17

Strengsøgning ved hjælp af en endelig tilstandsmaskine



- (1) Byg en endelig tilstandsmaskine ud fra mønsteret
- (2) Kør maskinen på teksten

```

Algorithm KMP_FSA_Match(T, P)
  move ← moveArray(P)
  i ← 0
  state ← 0
  while i < n and state < m do
    state ← move[T[i], state]
  if state = m then
    return i - m + 1
  return -1 { no match }
    
```

Tidsforbrug: $O(n)$

18

Konstruktion af endelig tilstandsmaskine



Efterfølgertilstande til tilstand i :

- match, så $i + 1$
- mismatch, så den tilstand, der svarer til "længst mulige match" (tag højde for mismatch)

Eksempel: 10100111

Tilstand 6

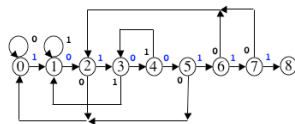
1010011: gå til tilstand 7

1010010: gå til tilstand for 010010 (tilstand 2),
husk tilstand for 010011 ($X = 1$)

Tilstand 7

10100111: gå til tilstand 8

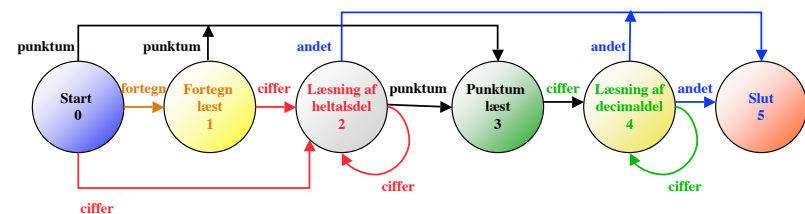
10100110: gå til tilstand for 0100110 (tilstand 2),
husk tilstand 0100111 ($X = 1$)
(disse er efterfølgere af forrige X)



19

Endelig tilstandsmaskine til indlæsning af decimaltal

(f.eks. -3.1415)



20

Rabin-Karp's algoritme



Ide: Brug hashing

- Beregn hashværdien af hver mulig delstreng af længde m i T
- Sammenlign med hashværdien for P
- Hvis to hashværdier er ens, så sammenlign delstrengen i T med P tegn for tegn

Eks: "tabelstørrelse" = 97, $m = 5$.
Find **15926**. Hashværdi = **18** (mod 97)
 $31415926535897932384626433$
 $31415 = 84 \pmod{97}$
 $14159 = 94 \pmod{97}$
 $41592 = 76 \pmod{97}$
 $15926 = 18 \pmod{97}$

Der ikke brug for nogen hashtabel, kun dens størrelse

21

Algoritmedesign



Problem: Hashfunktionen afhænger af m tegn

Løsning: Beregn hashværdi for $i+1$ ud fra hashværdi for i (forbedrer køretiden fra $O(nm)$ til $O(n+m)$)

$$\begin{aligned}31415 &= 84 \pmod{97} \\14159 &= (31415 - 3 \cdot 10000) \cdot 10 + 9 \\&= (84 - 3 \cdot 9) \cdot 10 + 9 \pmod{97} \\&= 579 = 94 \pmod{97} \\41592 &= (94 - 1 \cdot 9) \cdot 10 + 2 = 76 \pmod{97} \\15926 &= (76 - 4 \cdot 9) \cdot 10 + 6 = 18 \pmod{97}\end{aligned}$$

Problem: En fuldstændig sammenligning er nødvendig ved kollision

Afhjælpning: Benyt en meget stor (virtuel) tabelstørrelse

22

Beregning af hashværdi



Antag at en hashværdi er beregnet ud fra

$$h = T[i]d^{m-1} + T[i+1]d^{m-2} + \dots + T[i+m-1],$$

hvor d er antallet af mulige tegn

Den næste værdi af h

$$T[i+1]d^{m-1} + T[i+2]d^{m-2} + \dots + T[i+m-1]d + T[i+m]$$

kan da bestemmes ud fra h :

$$(h - T[i]d^{m-1})d + T[i+m]$$

23

Empirisk undersøgelse af algoritmernes effektivitet



Søgning i "The Oxford English Dictionary" (2nd Edition), cirka 570 millioner tegn

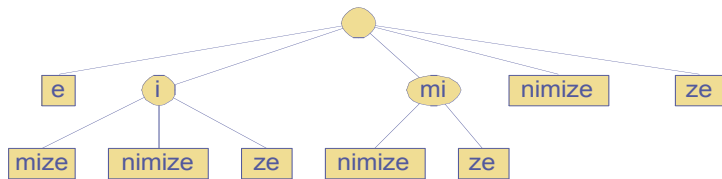
Algoritme	"to be or not to be"	"data"
Rå kraft	1.23	1.74
Knuth-Morris-Pratt	2.16	2.93
Boyer-Moore	1.33	1.16
Rabin-Karp	2.64	3.69
Boyer-Moore-Horspool	1.00	1.00

fra G. H. Gonnet & R. Baeza-Yates:

Handbook of Algorithms and Data Structures
in Pascal and C, Addison-Wesley 1991

24

Tries



25

Forbehandling af strenge

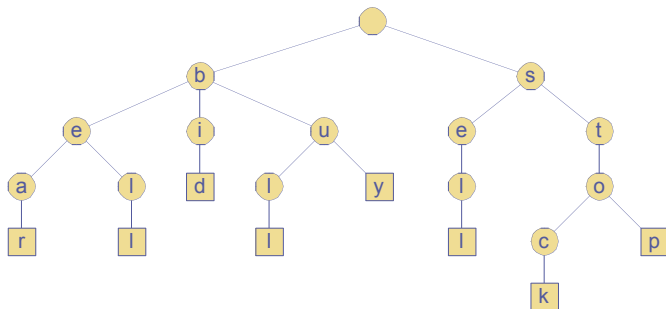
- Forbehandling af mønsteret sætter hastigheden op for strengsøgning
Efter forbehandling af mønsteret udfører KMP-algoritmen strengsøgning i tid, der er proportional med tekstens længde
- Hvis en tekst er lang, uforanderlig og søges i hyppigt (f.eks. Shakespeares værker), kan vi ønske at forbehandle teksten snarere end mønsteret
- Et **trie** (udtales *traj*) er en kompakt datastruktur til repræsentation af strenge, såsom ordene i en tekst
Et trie muliggør strengsøgning i tid, der er proportional med mønsterets længde

26

Standard trie (1)

- Et **standard trie** for en mængde af strenge S er et ordnet træ, hvor
- Hver knude, udtagen roden, indeholder et tegn
 - Børnene af en knude er ordnet alfabetisk
 - Vejene fra roden til de eksterne knuder giver strenge i S

Eksempel: $S = \{ \text{bear, bell, bid, bull, buy, sell, stock, stop} \}$

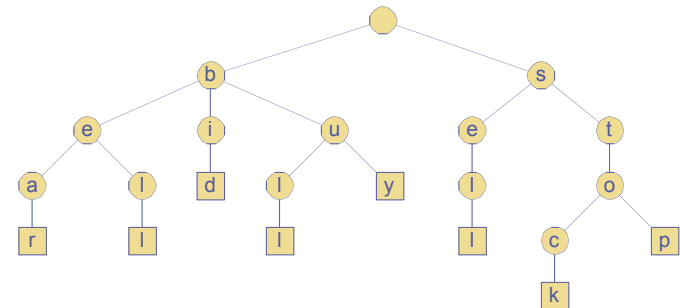


27

Standard trie (2)

Et standard trie bruger $O(n)$ plads og tillader søgning, indsættelse og fjernelse i $O(dm)$ tid, hvor:

- n samlede længde af strengene i S
- m længden af strengparameteren for operationen
- d størrelsen af alfabetet



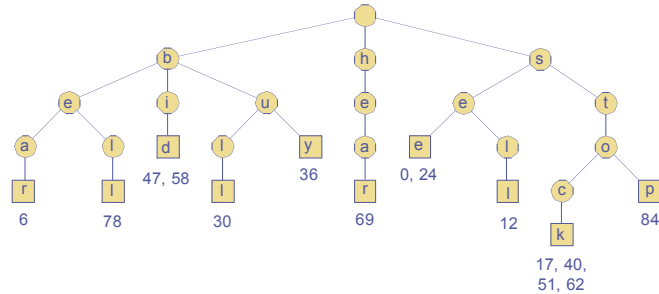
28

Søgning af ord med et trie

Vi indsætter tekstens ord i træet

s	e	e	a	b	e	a	r	?	s	e	l	l	s	t	o	c	k	!						
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	
s	e	e	a	b	u	l	l	?	b	u	y	s	t	o	c	k	!							
24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46		
b	i	d	s	t	o	c	k	!	b	i	d	s	t	o	c	k	!							
47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68			
h	e	a	r	t	h	e	b	e	l	l	?	s	t	o	p	!								
69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88					

Hvert blad indeholder forekomsterne af det tilknyttede ord i teksten

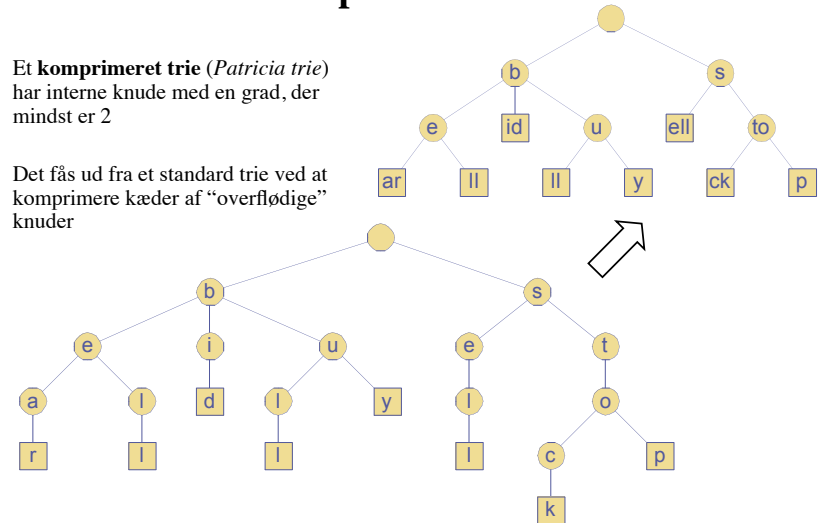


29

Komprimeret trie

Et komprimeret trie (*Patricia trie*) har interne knude med en grad, der mindst er 2

Det fås ud fra et standard trie ved at komprimere kæder af "overflødige" knuder

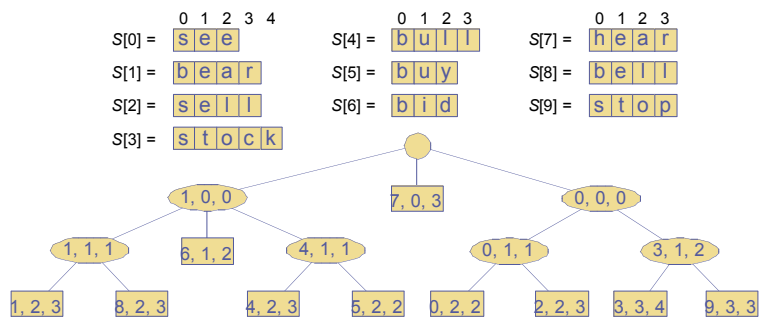


30

Kompakt repræsentation

Kompakt repræsentation af et komprimeret trie for et array af strenge:

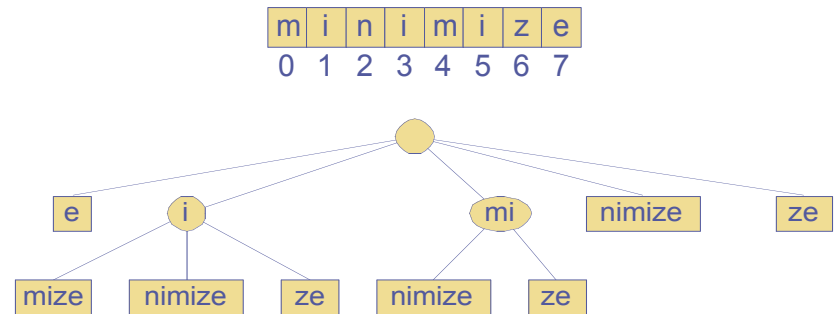
- Knuderne indeholder indeksintervaller i stedet for delstrenge
- Bruger $O(s)$ plads, hvor s er antallet af strenge i arrayet



31

Suffix-trie (1)

Suffix-trie'et for en streng X er det komprimerede trie for alle suffixer i X

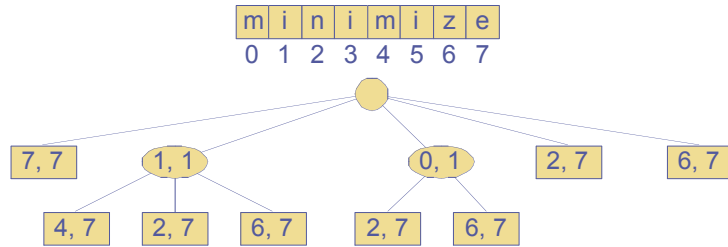


32

Suffix-trie (2)

En kompakt repræsentation af suffix-trie'et for en streng X af længde n fra et alfabet af størrelse d

- bruger $O(n)$ plads
- tillader strengsøgning i X i $O(dm)$ tid, hvor m er mønsterets længde (afhænger ikke af n)



33

Komprimering



34

Filkomprimering



Komprimering reducerer størrelsen af en fil

- for at spare **plads** ved lagring
- for at spare **tid** ved transmission

Komprimering benyttes til

- tekst:** nogle bogstaver er hyppigere end andre
- grafik:** store ensartede områder
- lyd:** gentagne mønstre

35

Run-length encoding



Komprimering ved tælling af gentagelser

Komprimering af **tekst**:

Strengen

AAAABBBAAABBBBBCCCCCDDABCBAABBBBCCDD

omkodes til

4A3BAA5B8CDABC3A4B3CD

Med escape-tegn (her '\'):

\4A\3BAA\5B\8CDABC\3A\4B\3CD

Run-length encoding er normalt ikke særlig effektiv for tekstfiler

36

Trie til kodning (2)

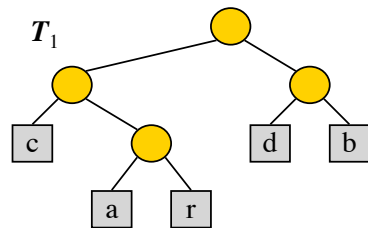
Givet en tekststreng X ønsker vi at finde en prefixkode for tegnene i X , der giver et kort kodeord for X

- Hyppige tegn bør have en korte kodeord
- Sjældne tegn bør have et langt kodeord

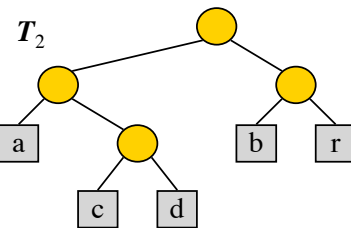
Eksempel:

$X = \text{abracadabra}$

T_1 koder X til 29 bits



T_2 koder X til 24 bits



41

Huffman's algoritme

Givet en streng X konstruerer Huffman's algoritme en prefixkode, der minimerer længden af kodeordet for X

En hob-baseret prioritetskø benyttes som hjælpedatastruktur

Algoritmen er **grådig**, men giver en optimal prefixkode

Køretid:

$O(n + d \log d)$, hvor n er længden af X , og d er antallet af forskellige tegn i X

Algorithm *HuffmanEncoding*(X)

Input string X of size n

Output optimal encoding trie for X

$C \leftarrow \text{distinctCharacters}(X)$

$\text{computeFrequencies}(C, X)$

$Q \leftarrow$ new empty heap

for all $c \in C$

$T \leftarrow$ new single-node tree storing c

$Q.\text{insert}(\text{getFrequency}(c), T)$

while $Q.\text{size}() > 1$

$f_1 \leftarrow Q.\text{minKey}()$

$T_1 \leftarrow Q.\text{removeMin}()$

$f_2 \leftarrow Q.\text{minKey}()$

$T_2 \leftarrow Q.\text{removeMin}()$

$T \leftarrow \text{join}(T_1, T_2)$

$Q.\text{insert}(f_1 + f_2, T)$

return $Q.\text{removeMin}()$

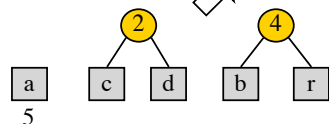
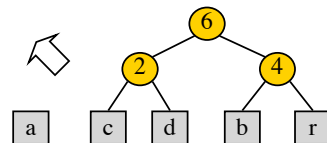
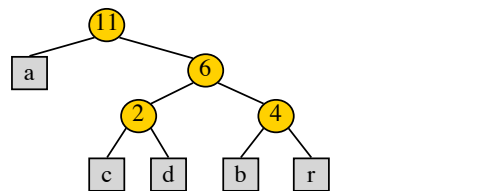
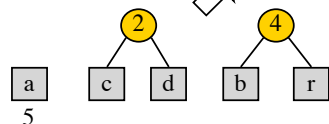
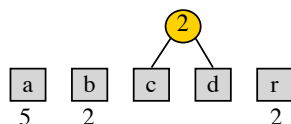
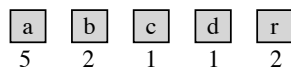
42

Eksempel

$X = \text{abracadabra}$

Frekvenser

a	b	c	d	r
5	2	1	1	2



43

Problemer for Huffmans algoritme



- Kodetræet skal medsendes (typisk 255 bytes)
- To gennemløb af filen (frekvensbestemmelse + kodning)
- Typisk 25% pladsreduktion, men ikke optimal

44

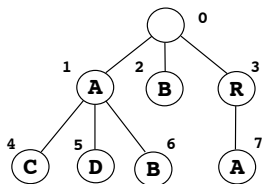
LZW-komprimering

(Lempel, Ziv og Welch, 1977)



Opbyg successivt en ordbog i form af et trie

Eksempel: **ABRACADABRA**



Kodning: **ABR1C1D1B3A**

45

Søgemaskiner



46

Søgemaskiner



En **søgemaskine** er et edb-baseret værktøj til at finde information på nettet

Arkitektur:

Web crawler (spider) indsamler information fra netsiderne ved at forfølge hyperhægter. Stopper aldrig

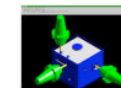
Indexer benytter den indsamlede information til opbygning af datastrukturer, der gør søgning hurtig

Søgemaskine (search engine, retriever) tillader brugere at hente information fra de opbyggede datastrukturer

- Grænseflade til forespørgsler
- Opslag i database med henvisninger til mere end 2 milliarder dokumenter (fylder mere end 4000 GB på disk)
- Rangering af dokumenter

47

Indeksering



Søgningen effektiviseres ved hjælp af en ordbog (et **indeks**)

Ordbogen realiseres som en **inverteret fil**, hvor hver post indeholder et ord og en samling henvisninger til de dokumenter, hvori ordet forekommer

Ordene kaldes **indekstermer**

Samlingen af henvisninger kaldes **forekomstlisten**

bell	10, 22, 23
buy	12
stop	1, 22
bid	24, 27
sell	27
bull	1
stock	33, 102
bear	22, 33

... ..

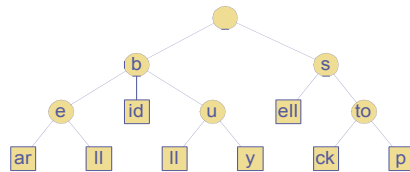
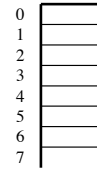
48

Implementering af en inverteret fil



Datastrukturer:

- Et array, der indeholder forekomstlisterne (i vilkårlig orden)
Listerne holdes sorteret (for hurtigt at kunne bestemme fællesmængde)
Gemmes på disk
- Et komprimeret trie for mængden af indekstermer, hvor hver eksterne knude lagrer indekset på forekomstlisten for den tilhørende term
Holdes i arbejdshukommelsen



49

Web crawler



- Henter netsider med henblik på indeksering
- Bredde-først-søgning
 - (1) Begynd med en side, P
 - (2) Find hyperhægtterne (URLs) på P og tilføj dem til en kø
Overgiv P til indexeren
 - (3) Hent P fra køen og gå til (2)

Spørgsmål:

- Hvorledes dybt skal der søges inden for et netsted?
- Hvor hyppigt skal sider besøges?

50

Rangering



Fund skal præsenteres i en rækkefølge

Hvilken?

Relevans, friskhed, popularitet, pålidelighed

Metoder til rangering af nøgleord:

- Tilstedeværelse i dokumenttitel
- Afstand fra start af dokument
- Hyppighed i dokument
- Hægtepopularitet (hvor mange sider peger på siden?)

51