

Leidsche Flesch LunchLezing

KMP

patroonherkennning

Hendrik Jan Hoogeboom

LIACS

26 september 2012

Enterobacteria phage MS2, complete genome (GenBank: EF204940.1)

gggtgggacc cctttcgggg tcctgctcaa cttcctgtcg agctaattgcc atttttaatg tctttagcga gacgctacca tggctatcgc
tgtaggtagc cgaattcca ttcctaggag gtttgacctg tgcgagcttt tagtaccctt gatagggaga acgagacctt cgtcccctcc
gttcgcgttt acgcggacgg tgagactgaa gataactcat tctctttaa atactgcttcg aactggactc ccggtcgttt taactcgact
ggggccaaaa cgaaacagtg gcaactaccc tctccgtatt cacggggggc gtttaagtgtc acatcgatag atcaagggtc ctacaagcga
agtgggtcat cgtggggctg cccgtacgag gagaaagccg gtttcggctt ctccctcgac gcacgctcct gctacagcct cttccctgta
agccagaact tgacttacat cgaagtgccg cagaacgttg cgaaccgggc gtcgaccgaa gtcctgcaaa aggtcaccca gggtaatttt
aaccttggtg ttgctttagc agaggccagg tcgacagcct cacaactcgc gacgcaaacc attgctctcg tgaaggcgta cactgccgct
cgtcgcggta attggcgcca ggcgctccgc taccttgccc taaacgaaga tcgaaagtgtt cgatcaaaac acgtggccgg cagggtggtg
gagttgcagt tcggttggtt accactaatg agtgatatcc aggggtgata tgagatgctt acgaaggttc accttcaaga gtttcttctt
atgagagccg tacgtcaggt cggtactaac atcaagttaa atggccgtct gtcgtatcca gctgcaaac tccagacaac gtgcaacata
tcgcgacgta tcgtgatatg gttttacata aacgatgcac gtttgcatg gttgtcgtct ctaggatctt tgaaccctt aggtatagtg
tgggaaaagg tgcctttctc attcgttgtc gactggctcc tacctgtagg taacatgctc gagggcctta cggccccgt gggatgctcc
tacatgtcag gaacagttac tgacgtaata acgggtgagt ccatcataag cgttgacgct ccctacgggt ggactgtgga gagacagggc
actgctaagg cccaaatctc agccatgcat cgaggggtac aatccgatg gccacaact ggcgctacg taaagtctcc tttctcgatg
gtccatacct tagatgcggt agcattaatc aggcaacggc tctctagata gagccctcaa ccggagtttg aagcatggct tctaacttta
ctcagttcgt tctcgtcgac aatggcggaa ctggcgacgt gactgtcgc ccaagcaact tcgctaacgg ggtcgtgaa tggatcagct
ctaactcgcg ttcacaggct tacaagtaa cctgtagcgt tcgtcagagc tctgcgaga atcgcaaata caccatcaa gtcgagggtc
ctaaagtggc aaccagact gttggtggtg tagagcttc tgtagccga tggcgttcgt acttaaatat ggaactaacc attccaattt
tcgctacgaa ttccgactgc gagcttattg ttaaggcaat gcaaggctct ctaaagatg gaaaccgat tccctcagca atcgcagcaa
actccggcat ctactaatag acgccggcca ttcaaacatg aggattacc atgtcgaaga caacaaagaa gttcaactct ttatgtattg
atcttcctcg cgatctttct ctcgaaattt accaatcaat tgcttctgtc gctactggaa gcggtgatcc gcacagtgc gactttacag
caattgctta ctttaaggac gaattgctca caaagcatcc gacctagggt tctggtaatg acgaggcgac ccgtcgtacc ttagctatcg
ctaagctacg ggaggcgaat gatcggtgcg gtcagataaa tagagaagg tttttacatg acaaatcctt gtcatgggat ccggatgttt
tacaaccag catccgtagc cttattggca acctcctctc tggctaccga tcgtcgttgt ttgggcaatg cacgttctcc aacgggtgct
ctatggggca caagttgcag gatgcagcgc cttacaagaa gttcgtgaa caagcaaccg ttacccccg cgctctgaga gcggctctat
tgggccgaga ccaatgtgcg ccgtggatca gacacgcggt ccgtataac gagtcatatg aatttaggct cgtttaggg aacggagtgt
ttacagttcc gaagaataat aaaatagatc gggctgcctg taaggagcct gatatgaata tgtacctcca gaaaggggtc ggtgcctta
tcagacgccg gctcaaattc gttggtatag atctgaatga tcaatcgatc aaccagcttc tggctcagca gggcagcgt gatggttcgc
ttgcgacgat agacttatcg tctgcatccg attccatctc cgatcgctg gtgtggagtt ttctcccacc tgagctatat tcatatctcg
atcgatccg ctcacactac ggaatcgtag atggcgagac gatacgatgg gaactat tccacaatggg aaatggggtc acgtttgagc
tagagtccat gatattctgg gcaatagtca aagcgacca aatccat tggtaacgccg gaaccatagg catctacggg gacgatatta
tatgtcccag tgagattgca ccccggtgtc tggaggcact tgcctactac ggtttcaaac cgaatctccg taaaacgttc gtgtccgggc
tctttcgcga gagctgcggc gcgcactttt accgtggtgt cgatgtcaa ccgttttaca tcaagaaacc tgttgacaat ctcttcgcc
ttatgctgat attgaaatcg ctacggggtt ggggagttgt cggaggatg tcagatccac gcctttaca ggtgtgggta cgactctct
cccagggtgcc ttcgatgttt ttcgggtggga cggacctcgc tgccgactac tacgtagtca gcccgcccac ggcagtctcg gtatatacca
agactccgta tgggcggcta ctgcggata cccgtacctc gggtttccgt cttgctcgta tcgctcgaga acgcaagttc ttcagcga
agcatgacag tggccgctac atagcgtggt tccatactgg aggtgaagtc accgacagta tgaagtccgc cggcgtgctg attatgcgca
cttcggagtg gctaacgccg gttcccacat tccctcagga gtgtgggcca gcgagctctc ctcggtagct gaccgagga cccccgtaa
cgggggtgggt gtgctcga aa gagcacgggt ccgcgaaagc ggtggctcca ccgaaagggt ggcgggcttc ggcccagga cctcccctg
aagagagggc ccggaattct cccgatttgg taactagctg cttgactagt taccacca

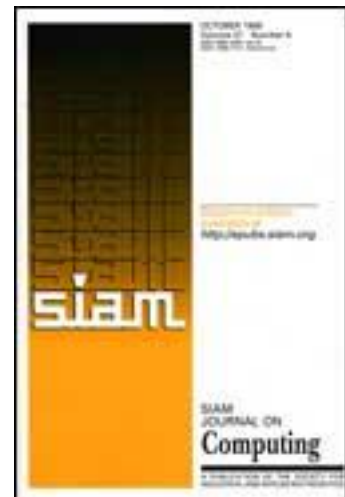
SIAM J. COMPUT.
Vol. 6, No. 2, June 1977

FAST PATTERN MATCHING IN STRINGS*

DONALD E. KNUTH[†], JAMES H. MORRIS, JR.[‡] AND VAUGHAN R. PRATT[¶]

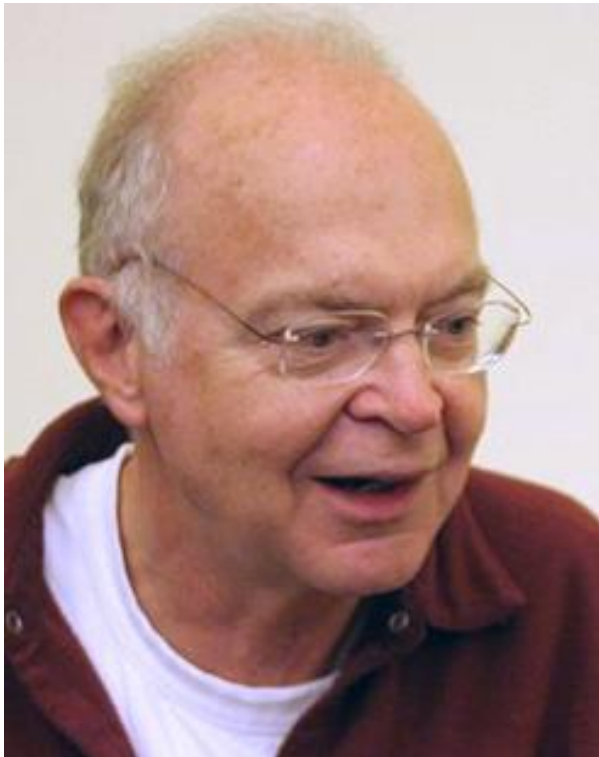
Abstract. An algorithm is presented which finds all occurrences of one given string within another, in running time proportional to the sum of the lengths of the strings. The constant of proportionality is low enough to make this algorithm of practical use, and the procedure can also be extended to deal with some more general pattern-matching problems. A theoretical application of the

Donald Knuth, James H. Morris, Jr & Vaughan Pratt . Fast pattern matching in strings.
SIAM Journal on Computing 6 (1977) 323–350. doi:10.1137/0206024



krasse knarren

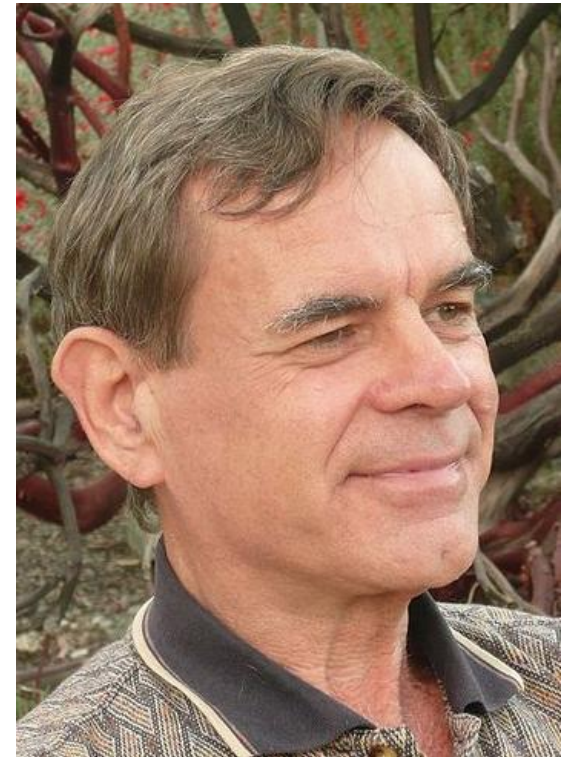
Donald Knuth (1938-)



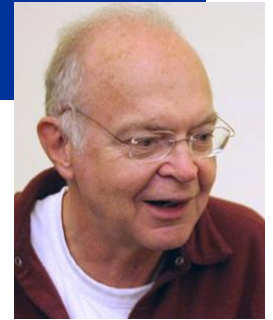
James Morris (1941-)



Vaughan Pratt (1944-)



Donald Knuth

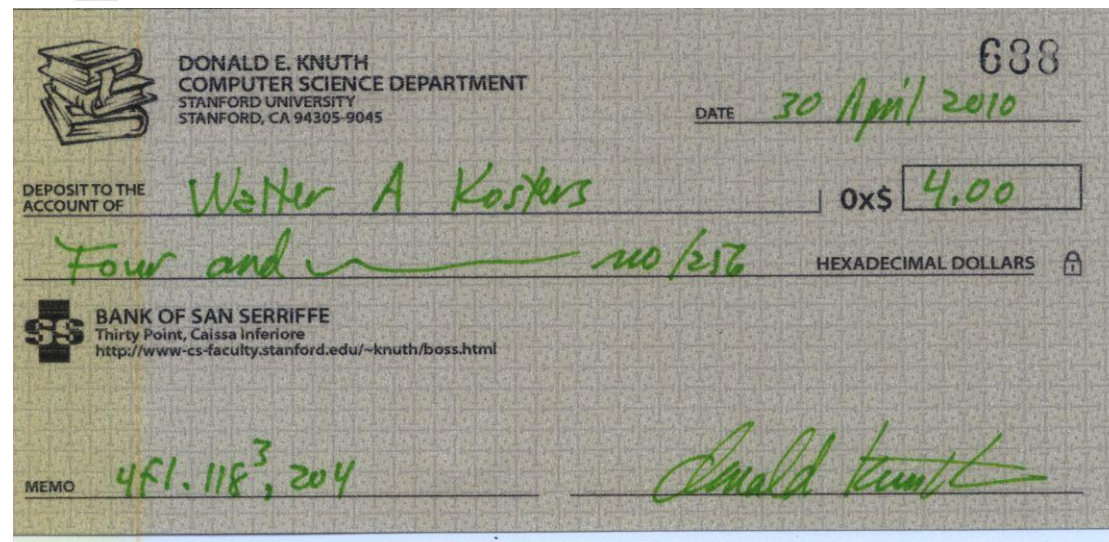
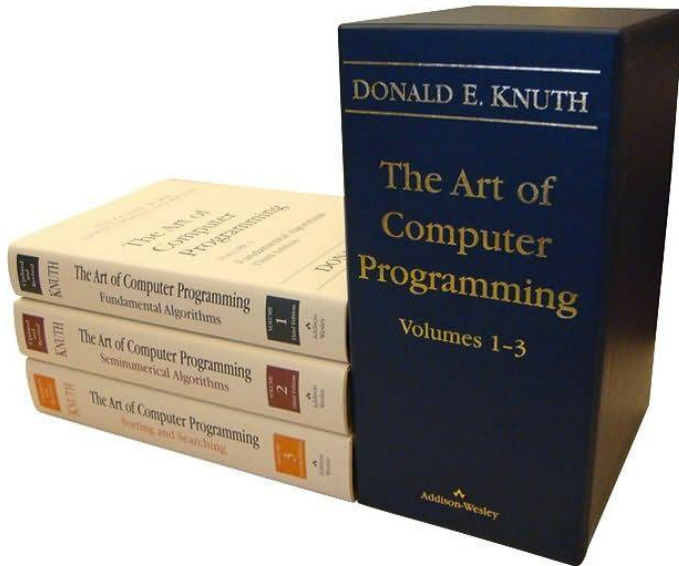


Stanford University

The Art of Computer Programming
géén email

\$2.56

TEX



James Morris



University of California, Berkeley

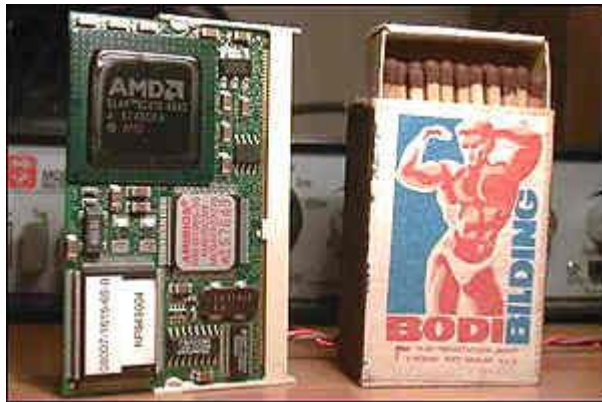
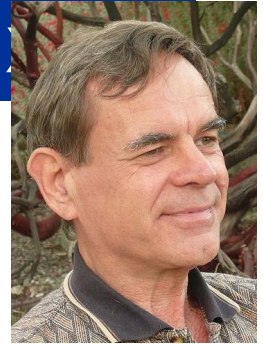
Xerox Alto
ongeveer 40 jaar oud
PC vs. mainframe



muis

“The base machine and one disk were housed in a cabinet about the size of a small refrigerator” (256kB / 2.5MB)

Vaughan Pratt



This computer runs a web site
(1999)

Stanford University

PRIMES \in NP
SUN



pattern matching

Patroon P =
a b c a b c a c a b

Tekst T =
b a b c b a b c a b c a a b c a b c a b c a c a b

naïve algoritme

a₁ b c a b c a c a b
b₁ a b c b a b c a b c a a b c a b c a b c a c a b
↑

naïve algoritme

a₁ b c a b c a c a b
b₁ a b c b a b c a b c a a b c a b c a b c a c a b
↑↑

b a₁ b c a₄ b c a c a b
a₂ b c b₅ a b c a b c a a b c a b c a b c a c a b
↑↑ x

naïve algoritme

a_1 b c a b c a c a b
 b_1 a b c b a b c a b c a a b c a b c a b c a c a b
↑↑

a_1 b c a_4 b c a c a b
b a_2 b c b_5 a b c a b c a a b c a b c a b c a c a b
↑↑ x

a_1 b c a b c a c a b
b a b_3 c b a b c a b c a a b c a b c a b c a c a b
↑↑

naïve algoritme

a_1 b c a b c a c a b
 b_1 a b c b a b c a b c a a b c a b c a b c a c a b
↑↑

a_1 b c a_4 b c a c a b
b a_2 b c b_5 a b c a b c a a b c a b c a b c a c a b
↑↑ x

a_1 b c a b c a c a b
b a b_3 c b a b c a b c a a b c a b c a b c a c a b
↑↑

...

a_1 b c a b c a c_8 a b
b a b c b a_6 b c a b c a a_{13} b c a b c a b c a c a b
↑↑

naïeve algoritme

a a a a a a a a b
a b
↑

kwadratische tijd

gebruik informatie!

$a_1 b c a b c a c_8 a b$
 $b a b c b a_6 b c a b c a a_{13} b c a b c a b c a c a b$
 \uparrow X

hoever moeten/kunnen we P doorschuiven?

T =	...	a	b	c	a	b	c	a	x		
P =	a_1	b	c	a	b	c	a	c_8	a	b			
P' =		a	b	c	a	b	c	a	c	a	b		
			a	b	c	a	b	c	a	c	a	b	
				a_1	b	c	a	b_5	c	a	c	a	b

dat weten we zonder kennis van T !

preprocessing

we gebruiken x niet
(die is 'onbekend')

... a b c a b c a **x** ...
P = a₁ b c a b c a c₈ a b
P' = a₁ b c a b₅ c a c a b

								fout			
								op positie			
j	=	1	2	3	4	5	6	7	8	9	10
P[j]	=	a	b	c	a	b	c	a	c	a	b
next[j]		0	1	1	0	1	1	0	5	0	1
									vervolg		
									op positie		

nul: geen match
schuif helemaal voorbij
huidige positie

positie j: verschuif over $j - \text{next}[j]$

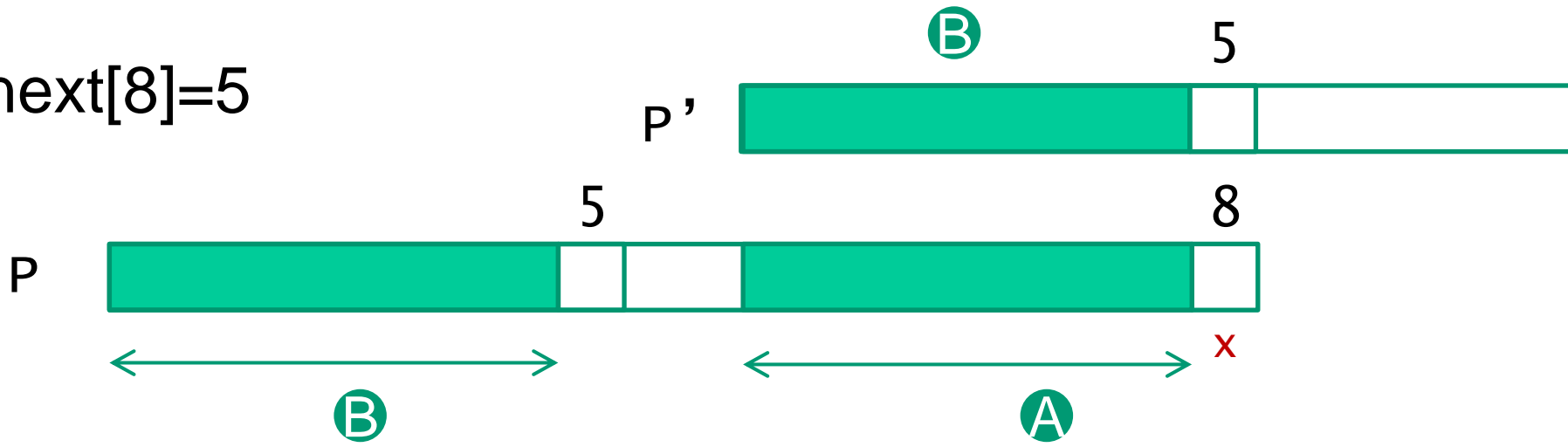
KMP algoritme

```
place pattern at left;  
while pattern not fully matched  
  and text not exhausted do  
  begin  
    while pattern character differs from  
      current text character  
      do shift pattern appropriately;  
        advance to next character of text;  
    end;
```

```
 $j := k := 1;$   
while  $j \leq m$  and  $k \leq n$  do  
  begin  
    while  $j > 0$  and  $text[k] \neq pattern[j]$   
      do  $j := next[j];$   
       $k := k + 1; j := j + 1;$   
    end;
```

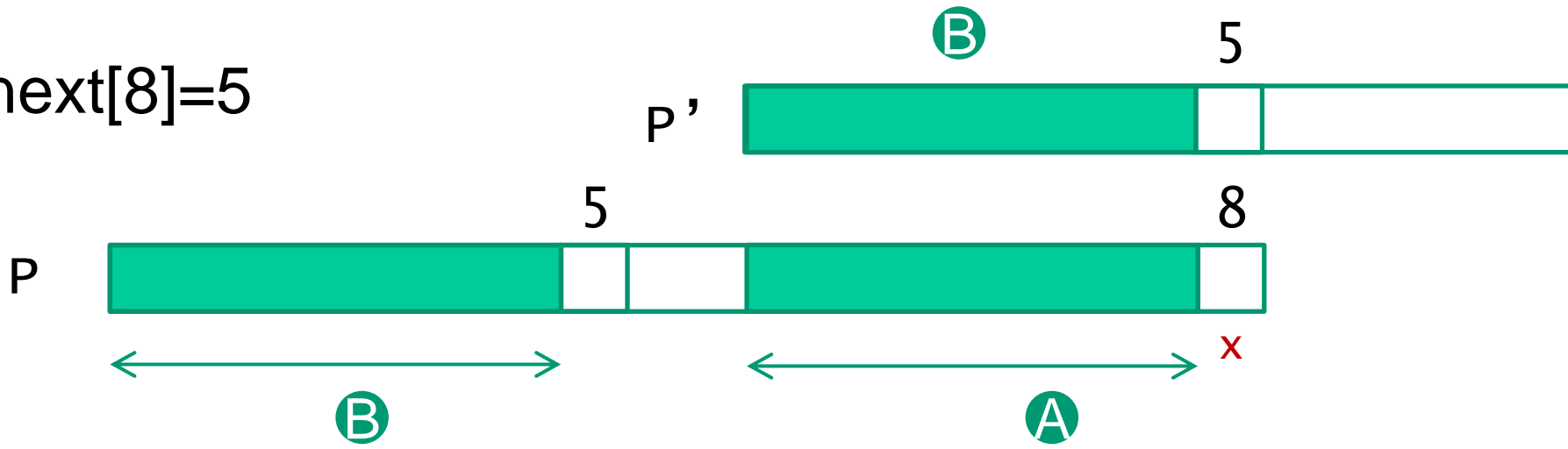
next[k] bepalen? overlap

next[8]=5

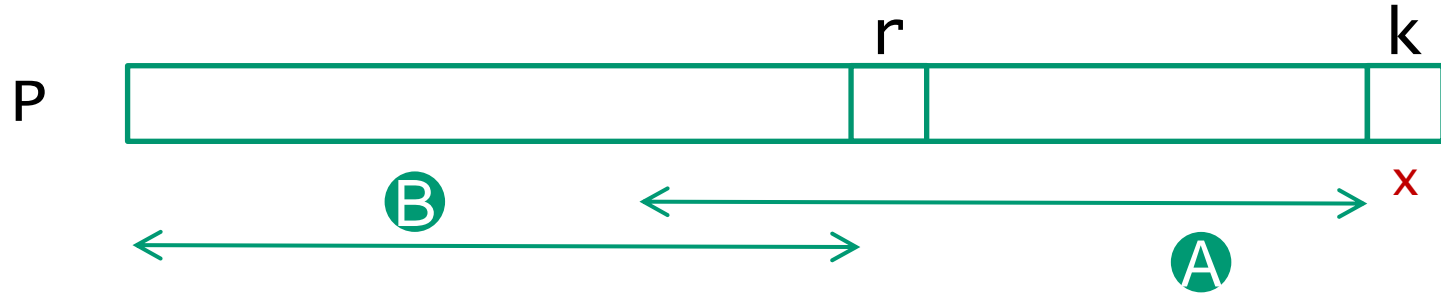


next[k] bepalen? overlap

next[8]=5



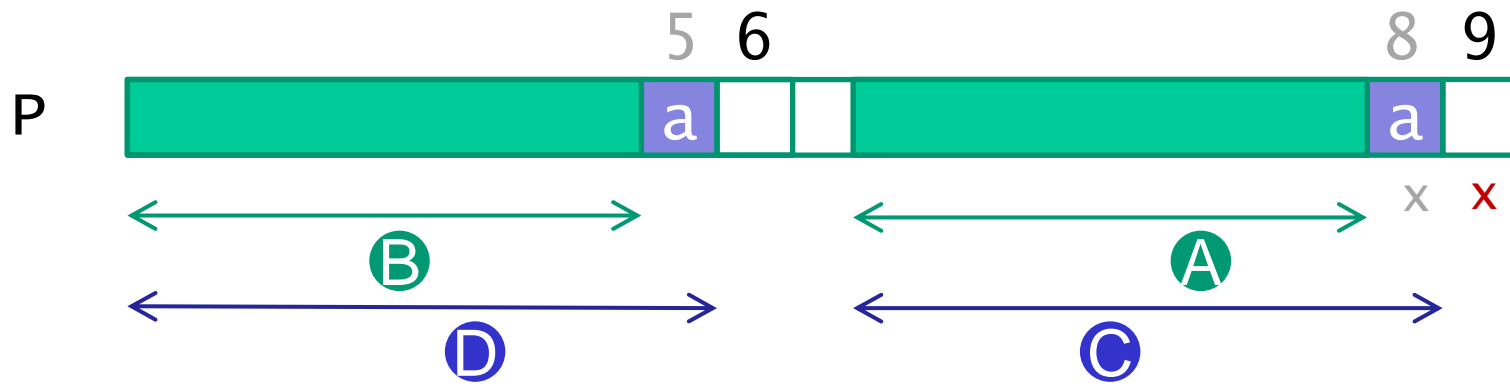
next[k]=r



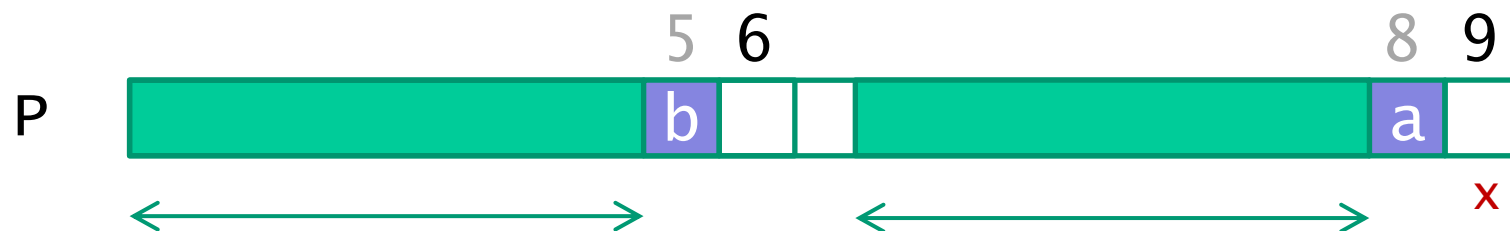
folgende overlap

$\text{next}[8]=5$

$\text{next}[8+1]=5+1$ mits $P[8] = P[5]$

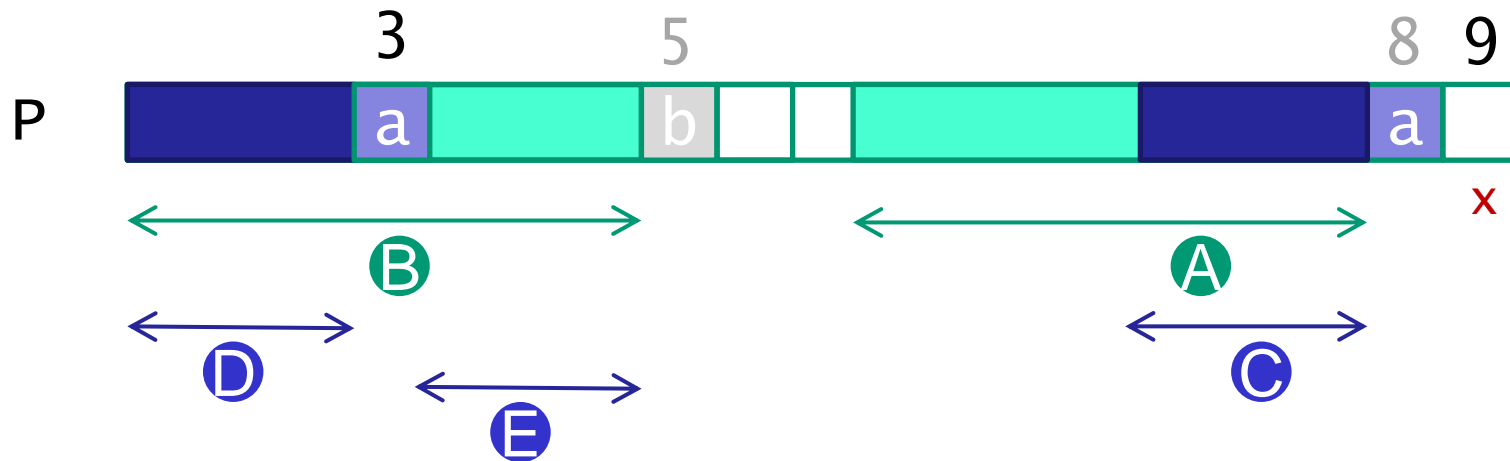


als $P[8] \neq P[5]$



$$\text{next}[8+1]=3+1$$

$$\text{mits } P[8] = P[3]$$



$$3 = \text{next}[5]$$

patroonherkenning op zichzelf toegepast!

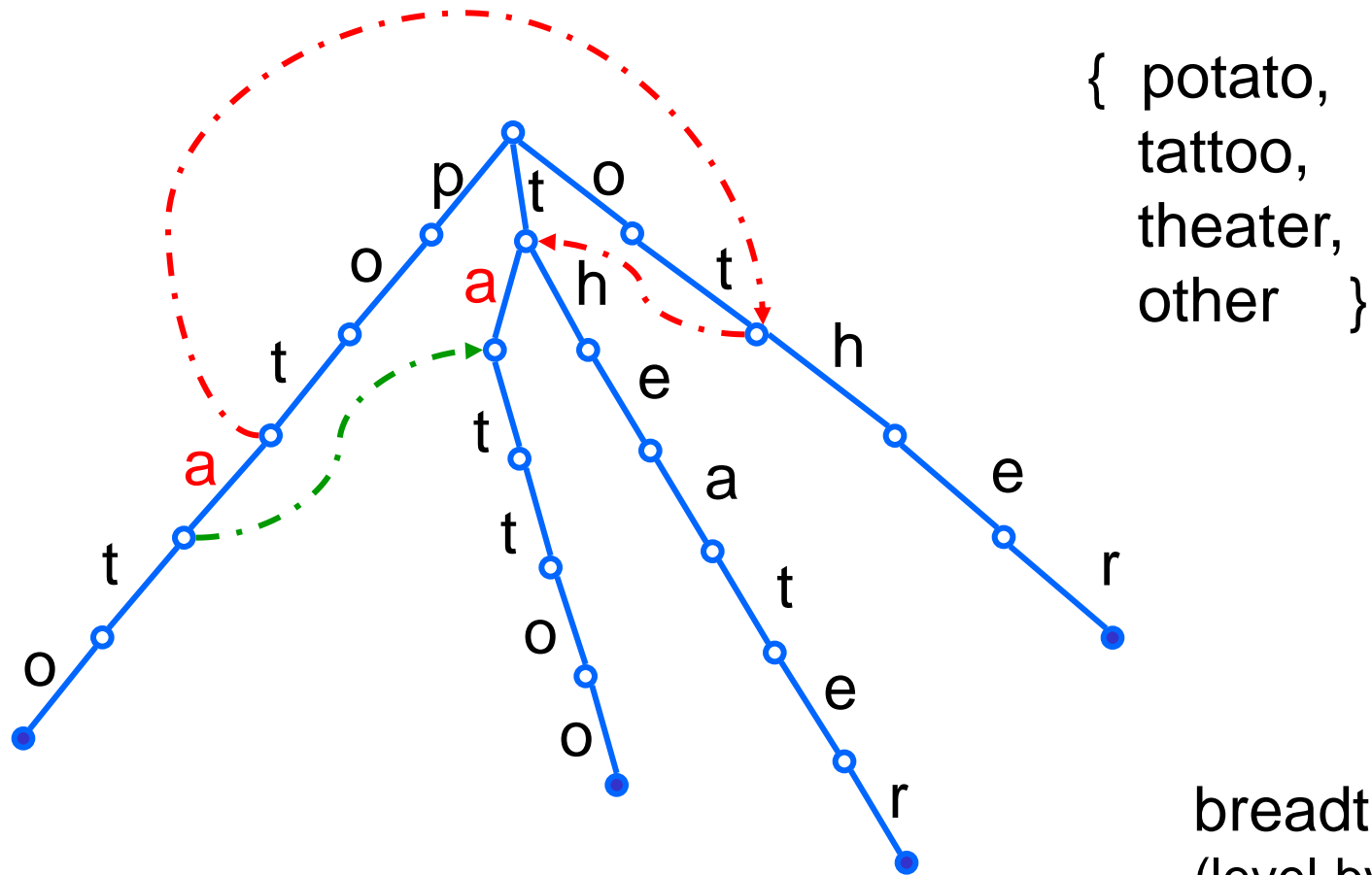
andere methoden

T= marktkoopman

P= schoenveter

schoe...

meerdere woorden: boom



{ potato,
tattoo,
theater,
other }

breadth first
(level-by-level)

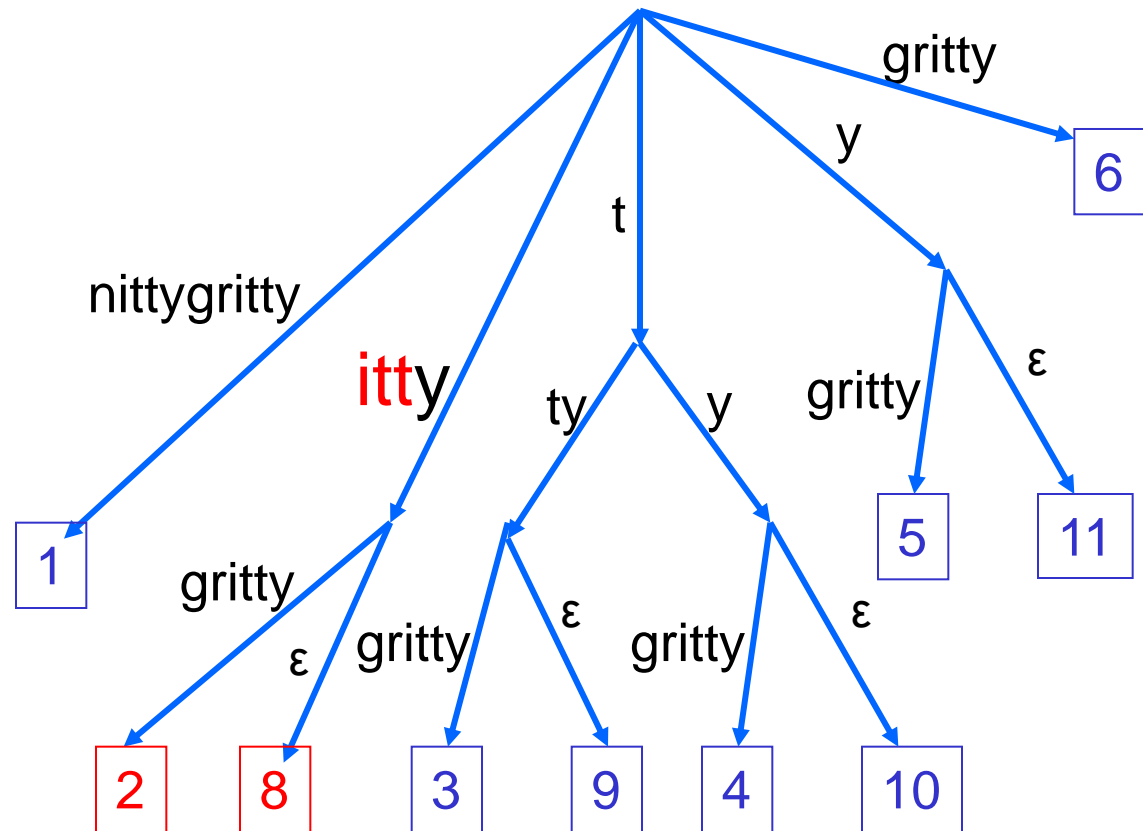
pot	ato
othe	r
theater	tattoo

pot	a	to
ota	☹	
t	a	ttoo

suffix-tree 'nittygritty'

"Algorithm of the Year 1973"

nittygritty	1
itty gritty	2
ttygritty	3
tygritty	4
ygritty	5
gritty	6
ritty	7
itty	8
tty	9
ty	10
y	11

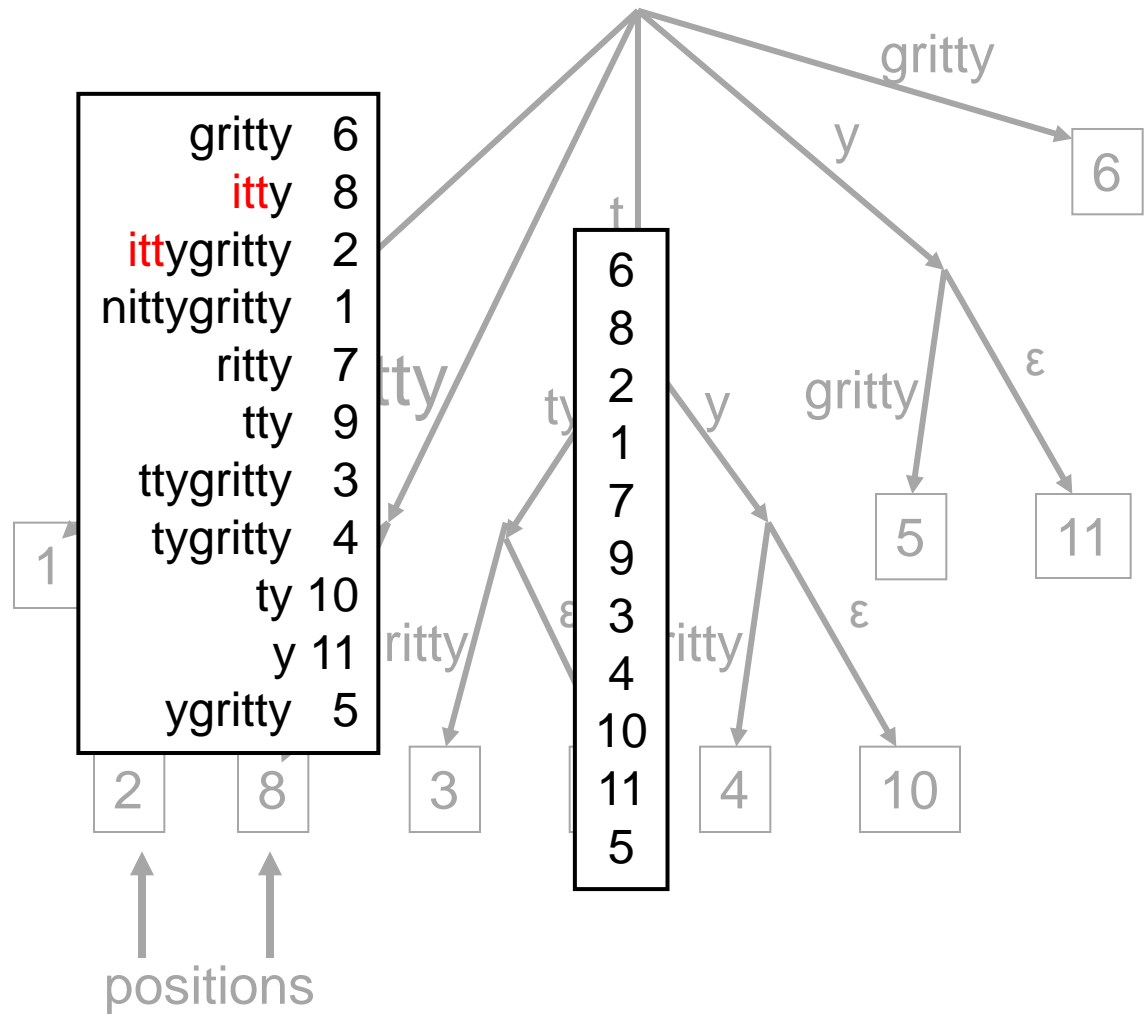


↑ ↑
positions

Weiner (1973),
McCreight (1976), Ukkonen (1995)

suffix-array 'nittygritty'

nittygritty	1
ittygritty	2
ttygritty	3
tygritty	4
ygritty	5
gritty	6
ritty	7
itty	8
tty	9
ty	10
y	11



KMP - Historical remarks

Morris text-editor (1969)

... was too complicated for other implementors of the system to understand, and he discovered several months later that gratuitous "fixes" had turned his routine into a shambles.

Knuth (1970)

- Cook two-way deterministic push-down automata
- Chester strings starting with palindromes

$v v w$

$v \# u v w$

This was the first time in Knuth's experience that automata theory had taught him how to solve a real programming problem better than he could solve it before

dankuwel ...



klaar