

# String Matching – Applicable Algorithmics

*Esko Ukkonen*  
University of Helsinki  
Esko.Ukkonen@Helsinki.Fi

FDPW 2004 June 8 Petrozavodsk

## Where is string matching?

- 'Every question in theoretical computer science can be formulated as a problem on strings' -*Zvi Galil*
- 'Biology has challenging problems for computer science for one hundred years' -*Donald Knuth*

## Basic question

Text: **KALEVALANKALAVALE**

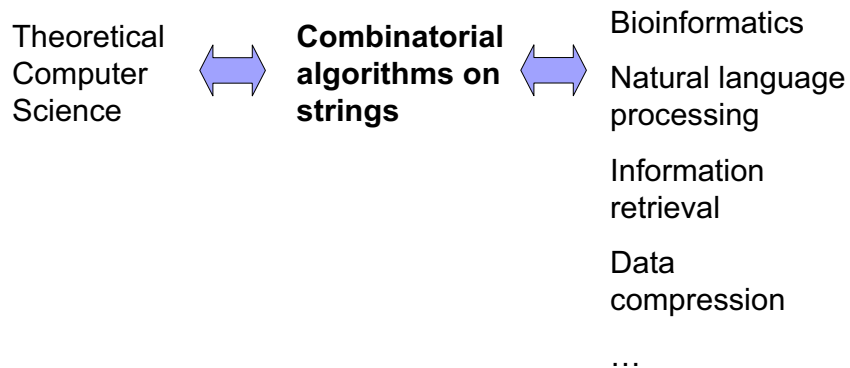
Pattern: **LANKA**

KALEVALANKALAVALE  
LANKA  
LANKA  
LANKA  
LANKA  
LANKA  
LANKA  
LANKA  
LANKA  
LANKA  
LANKA  
LANKA  
LANKA  
LANKA  
LANKA

AAAAAAAAAAAAAAAAAA . .

AAAAAAB

=>  $O(mn)$



## Exact string matching

- pattern  $P = p_1 \dots p_m$  and text  $T = t_1 \dots t_n$  in  $\Sigma^*$
- finite alphabet  $\Sigma$
- $m \ll n$
- find the occurrences of  $P$  in  $T$ :  
Knuth-Morris-Pratt  $O(m+n)$ ,  
multipattern matching Aho-Corasick  $O(m+n)$ ,  
Boyer-Moore  $O(n/m)$ , ...
- more general patterns: gaps, don't-care-symbols, regular expressions, ...

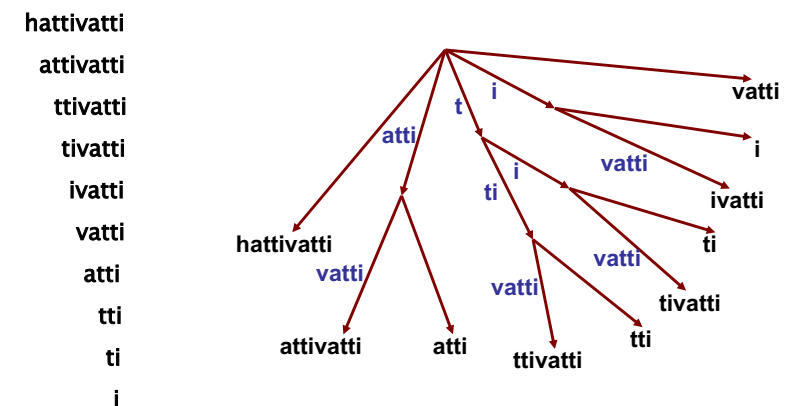
## Text indexing

- assume: text  $T$  stays the same (human genome, ...)
- need to answer queries for several different  $P$
- build an *index* for  $T$  that makes the search for different  $P$  fast: get rid of  $O(|T|)$
- suffix tree, suffix array, suffix automaton

## Suffix tree

- cleverly implemented TRIE that represents all the suffixes of  $T$
- $T = \text{hattivatti}$

hattivatti  
attivatti  
ttivatti  
tivatti  
ivatti  
vatti  
atti  
tti  
ti  
i



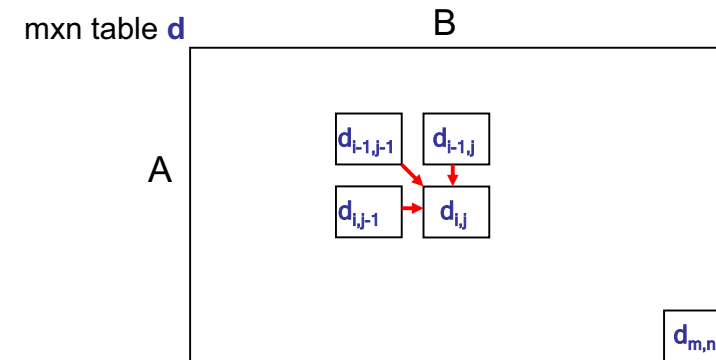


# Distance evaluation problem

- strings A and B, |A| = m, |B| = n, m ≤ n
- evaluate distance d(A,B)
- dynamic programming O(mn); also O(mn/logm) (Crochemore at al SODA2002)
- string alignment

# Dynamic programming

$$d_{i,j} = \min(\text{if } a_i=b_j \text{ then } d_{i-1,j-1} \text{ else } \infty, \\ d_{i-1,j} + 1, \\ d_{i,j-1} + 1) \\ = \text{distance between } i\text{-prefix of } A \text{ and } j\text{-prefix of } B$$



$$d_{i,j} = \min(\text{if } a_i=b_j \text{ then } d_{i-1,j-1} \text{ else } \infty, \\ d_{i-1,j} + 1, \\ d_{i,j-1} + 1)$$

A/B		s	t	o	c	k	h	o	l	m
	0	1	2	3	4	5	6	7	8	9
t	1	2	1	2	3	4	5	6	7	8
u	2	3	2	3	4	5	6	7	8	9
k	3	4	3	4	5	4	5	6	7	8
h	4	5	4	5	6	5	4	5	6	7
o	5	6	5	4	5	6	5	4	5	6
l	6	7	6	5	6	7	6	5	4	5
m	7	8	7	6	7	8	7	6	5	4
a	8	9	8	7	8	9	8	7	6	5

optimal alignment by trace-back

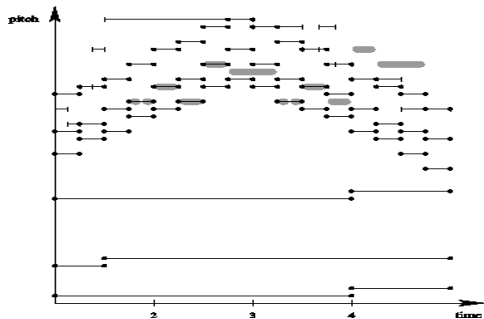
$d_D(A,B)$

# Search problem

- find approximate occurrences of P in T: substrings P' of T such that d(P,P') small
- dyn progr with a small modification: O(mn)
- lots of (practical) improvements:
  - distance bound k => O(kn) search;
  - utilizing regularities of the dp table;
  - filtration approach;
  - bit-parallelism (Gonnet & Baeza-Yates)



Du bist ja so schön, son-der-bar schön;



# Music retrieval and analysis

- symbolically (using notes) encoded music:
  - pitch levels < 300 (typically less)
  - durations < 1000 (typically less)
- sequence of discrete symbols, interpretation in 2D
- not necessarily evenly spaced as traditional strings (different durations)
- many symbols can be 'on' simultaneously

## Music...

- seminal application: **query-by-humming**
  - inaccurate pitches and durations
  - => approximate melody matching (Ghias et al 1995; MELDEX 1997)
- Mongeau & Sankoff (1990)
- **music information retrieval (MIR)**
- monophonic vs polyphonic music
- music analysis: comparative analysis, evolution of musical ideas, characterization of the style of composers, ...

<http://www.cs.helsinki.fi/group/cbrahms/demoengine/>

C-Brahms Retrieval Engine for Melody Searching (Beta)

Help Text version Songlist Mtopia project Main page

LCIS Search 20 Results Show First Matches Sort by Transposition 1 Errors

4C4 4D4 4E4 204 4E4 4D4

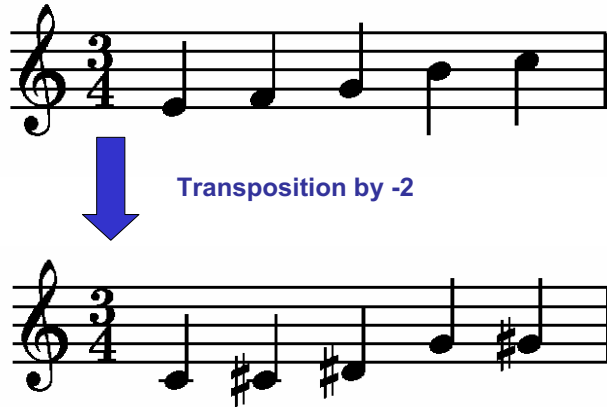
Play Clear Search

Results: 4 matches in 2 songs.

Number	Composer	Title	Opus	Date	Style	Play	Metadata	Score	Appx. Bar #	Matched Notes	Transposition	Errors	Aligns
1	Anonymous	Johany Cock thy Beaver: A Scotch Tune to a Ground	n/a	n/a	Baroque	MIDI	META	PS	22	A#5 C8 D8 F8 C8	22	1	C D E G E D A# C D F - C
2	Anonymous	Johany Cock thy Beaver: A Scotch Tune to a Ground	n/a	n/a	Baroque	MIDI	META	PS	30	A#5 C8 D8 F8 C8	22	1	C D E G E D A# C D F - C
3	Anonymous	Johany Cock thy Beaver: A Scotch Tune to a Ground	n/a	n/a	Baroque	MIDI	META	PS	54	A#5 C8 D8 F8 C8	22	1	C D E G E D A# C D F - C
4	J. Banister (c. 1624-1679)	A Division upon a Ground	n/a	n/a	Baroque	MIDI	META	PS	85	C8 D8 E8 F8 G8 E8 D8	24	1	C D E - G E D C D E F G E D

Pattern: 4C4 4D4 4E4 204 4E4 4D4

## Transposition invariance



## Transposition invariance...

- transposition (systematic shift) of pitch levels keeps the music the same
- the intervals between the notes do not change

## Transposition invariant distances

- alphabet  $\Sigma$ : closed with respect to addition and subtraction
- say,  $\Sigma = \text{integers}$
- $A = a_1 \dots a_m$  translated by  $t$ :  
 $A+t = a_1+t, \dots, a_m+t$
- *transposition invariant distance*  
 $d^T(A,B) = \min_{t \in \Sigma} d(A+t, B)$
- 'longest common hidden melody'

## Exact case is trivial

- interval:  $a_{i+1} - a_i$
- $\text{intervals}(A+t) = \text{intervals}(A)$
- use interval sequences  
 $a_2 - a_1, a_3 - a_2, \dots, a_m - a_{m-1}$   
instead of originals in exact matching
- transposition invariance for free

## Approximate case

- repeat dp for all  $O(mn)$  relevant transpositions:  $O(m^2n^2)$
- sparse dynamic programming (Epstein et al 1992) => transposition invariant distance & search  $O(mn \log \log n)$  (Veli Mäkinen's PhD thesis)

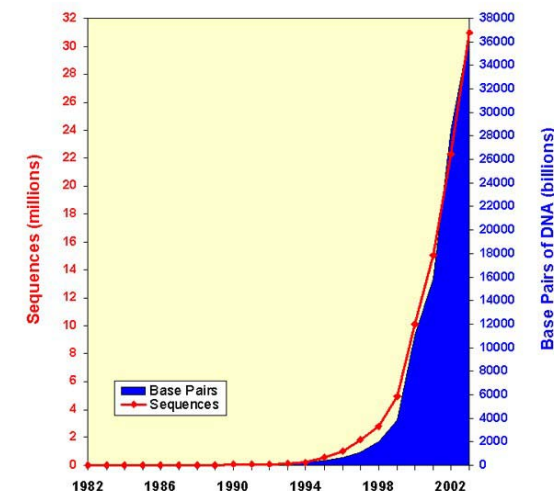
## Applications 2: Bioinformatics/computational biology

- DNA sequencing
- multiple alignment and motif discovery
- lots of more ...

## Genome projects

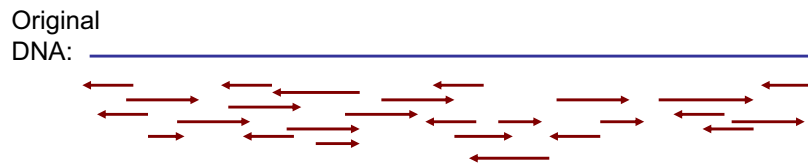
```
cgccgagtgacagagacgctaatacaggc
tgtgttctcaggatgcgtaccgagtggg
agacagcagcacgaccagcgggtggcaga
gacccttgcagacatcaagctctttggg
aacaagtggagcaccgatgatgtacagc
cgatcaatgacatttccctaatacagga
ttacattgcagtgcccaaggagaagtat
gccaagtaatacctccctcacagtg...
```

Growth of GenBank





# Shotgun sequencing



'Gel readings': random location and orientation

# The fragment assembly problem

cctcgagttaagtactgcccgggcttcaacggatctgtcgggagtcg



reconstruction?

cctcgagtttaa  
tacttaactcgag  
cgggcagtacttaa  
aagtactgcccgcg  
gcccgggcttcaacggat  
cccgggcttcaacggatctgtg  
cccgacacagat  
tgtgtcgggagtcg

cctcgagtttaa  
tacttaactcgag  
cgggcagtacttaa  
aagtactgcccgcg  
gcccgggcttcaacggat  
cccgggcttcaacggatctgtg  
cccgacacagat  
tgtgtcgggagtcg



cctcgagtttaa  
ctcgagt-taagta  
t-taagtactgcccg  
aagtactgcccgcg  
gcccgggcttcaacggat  
cccgggcttcaacggatctgtg  
atctgtgtcggg  
tgtgtcgggagtcg



cctcgagt-taagtactgcccgggcttcaacggatctgtgtcgggagtcg

# Fragment assembly (cont.)

- T = DNA sequence (say, the genome of a human), a string over alphabet A, C, G, T
- fragments F<sub>1</sub>, ..., F<sub>m</sub>: pieces of T of length 150 – 800, taken from random locations and in random orientation, some errors can occur
- construct T from its fragments F<sub>1</sub>, ..., F<sub>m</sub>.

## Some numbers (Human Genome)

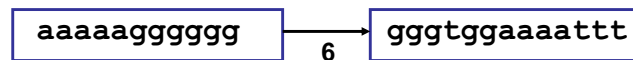
- length about 3 000 000 000 bp
- Celera's fragment data (Feb 2001): 27.27 million fragments, each 150-800 bp
- at least 7-fold coverage by fragments needed => total data length 7 x 3 billion bp

## Overlap graph - Layout - Consensus

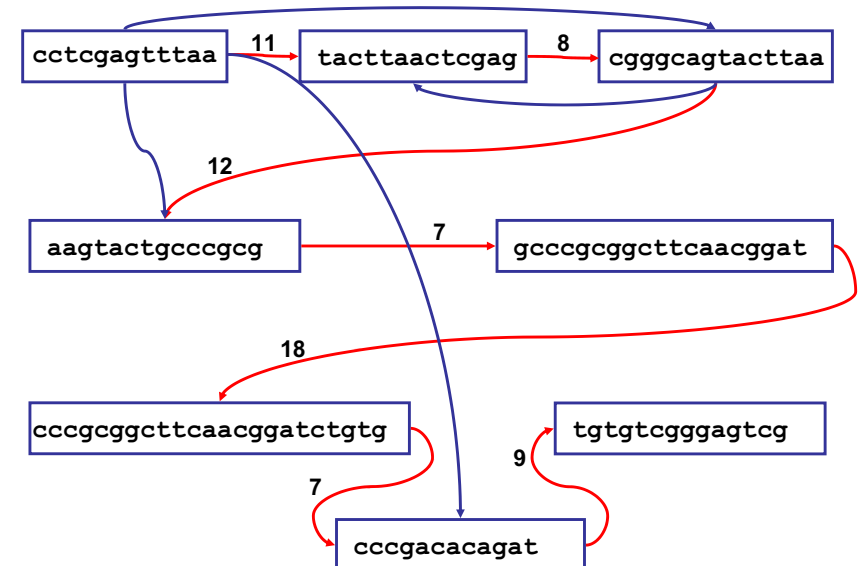
1. For each pair  $(F_i, F_j)$  of fragments, find all acceptable overlaps between  $F_i$  and  $F_j$  and their reverse complements. This gives an overlap graph.
2. Using the overlap graph, build a lay-out that gives global alignment of the fragments.
3. Construct a consensus string from the lay-out. The consensus is the resulting approximation of the original DNA sequence  $T$

## Overlap graph

aaaaagggggg  
gggtggaaaattt



Finding the overlaps: dynamic programming



# Hamiltonian paths $\Leftrightarrow$ Lay-outs

The red path gives our example lay-out:

```

cctcgagtttaa
ctcgagt-taagta
  t-taagtactgcccg
    aagtactgccgcg
      gcccgcggettcaacggat
        cccgcggettcaacggatctgtg
          atctgtgtcggg
            tgtgtcggggagtcg
  
```

# Consensus from lay-out

```

cctcgagtttaa
ctcgagt-taagta
  t-taagtactgcccg
    aagtactgccgcg
      gcccgcggettcaacggat
        cccgcggettcaacggatctgtg
          atctgtgtcggg
            tgtgtcggggagtcg
  
```

cctcgagt-taagtactgcccgcggettcaacggatctgtgtcggggagtcg

## Other problems

- finding homologous sequences: new sequence vs all old ones – the most popular computational task in present-day molecular biology = approximate string matching (filtration, also dynamic programming a la Smith-Waterman). BLAST
- multiple alignment of sequence families to find interesting motifs: NP-hard => heuristics, Hidden Markov models, MCMC

## A multiple alignment

	10	20	30	40	50	60
Hbb_Human.pep	-----VHLTPEEKSAVTALWGKVN--VDEVGGEALGRLLVVYPWTQRFFESFGDLST					
Hbb_Horse.pep	-----VQLSGEKAALVLAALWDKVN---EEVVGGEALGRLLVVYPWTQRFFDSFGDLSN					
Hba_Human.pep	-----VLSPADKTNVKAAGKVGAGHAGEYGAELERMFLSFPTTKTYFPHFDLS--					
Hba_Horse.pep	-----VLSAADKTNVKAAWSKVGGHAGEYGAELERMFLGFPTTKTYFPHFDLS--					
Myg_Phyca.pep	-----VLSEGEWQLVHLVWAKVEADVAGHGQDILIRLFKSHPETLEKFRDRFKHLKT					
Glb5_Petma.pep	PIVDTGSAPLSAAEKTIRSAWAPVYSTYETSGVDILVKFFTSTPAAQEFFPKFKGLTT					
Lgb2_Luplu.pep	-----GALTESQAALVKSSWEEFNANIPKHTHRFFILVLEIAPAAKDLFSFLKGTSE					
	*	*		*	*	
Hbb_Human.pep	PDAVMGNPKVKAHGKKVLAFAFSDGLAHLD----NLKGTFAATLSELHCDKLVDPENFRL					
Hbb_Horse.pep	PGAVMGNPKVKAHGKKVLAHSPGEGVHHLD----NLKGTFAALSELHCDKLVDPENFRL					
Hba_Human.pep	----HGSAQVKGHGKKVADALTNVAHVVD----DMPNALSALSDLHAHKLRLVDPVNFRL					
Hba_Horse.pep	----HGSAQVKAHGKKVGDALTLAVGHLD----DLPGALSNSDLHAHKLRLVDPVNFKL					
Myg_Phyca.pep	EAEMKASEDLKKHGVTVLTALGAILKKKG----HHEAELKPLAQSHATKHKIPIKYLEF					
Glb5_Petma.pep	ADQLKKSADVRWHAERIINAVNDAVASDDT--EKMSMKLRDLGSKHAKSFQVDPQYFKV					
Lgb2_Luplu.pep	VP--QNNPELQAHAGKVFCLVYEAIIQLQVTGVVVVTDATLNKLGSVHVSRG-VADAHFPV					
	..*			*	*	
Hbb_Human.pep	LGNVLVCLVAHHFGKEFTPPVQAAQKVVAGVANALAHKYH-----					
Hbb_Horse.pep	LGNVLVVVLAHHFGKDFTPPELQASQKVVAGVANALAHKYH-----					
Hba_Human.pep	LSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTLSKYR-----					
Hba_Horse.pep	LSHCLLVTLAVHLPNDFTPAVHASLDKFLSSVSTVLTLSKYR-----					
Myg_Phyca.pep	ISEAIIHVLHSRHPGDFGADAQGMNKALELFRKDIAAKYKELGYQG					
Glb5_Petma.pep	LAAVIADTVAAG-----DAGFEKLSMCMICILLRSAY-----					
Lgb2_Luplu.pep	VKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIIVKMEMNDAA---					

```

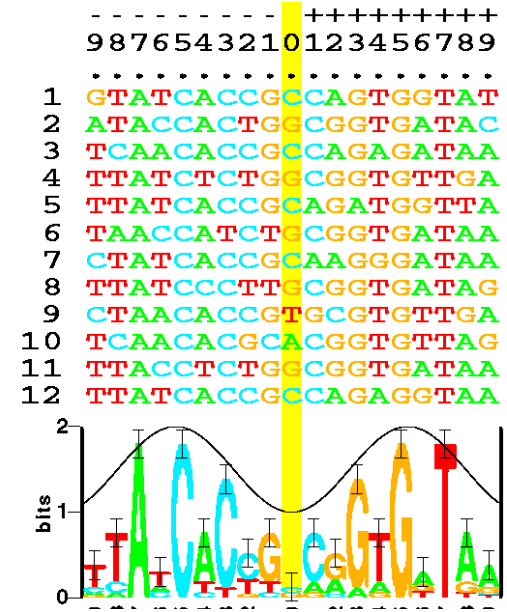
PRETTYBOX of: prettybox.ms(*)  October 26, 1998 11:05:11.41
fa10 .....t.ttttggaa..... 21
fa12 .....t.ttttggaa..... 21
foik .....t.ttttggaa..... 21
e.....t.ttttggaa..... 46
p1m.....t.ttttggaa..... 59
p1s.....t.ttttggaa..... 59
p2s.....t.ttttggaa..... 59
p3s.....t.ttttggaa..... 57
cb3.....t.ttttggaa..... 46
r14.....t.ttttggaa..... 49
r2.....t.ttttggaa..... 49

fa10 .....ffvk..... 60
fa12 .....ffvk..... 60
foik .....ffvk..... 60
e.....ffvk..... 97
p1m.....ffvk..... 110
p1s.....ffvk..... 110
p2s.....ffvk..... 109
p3s.....ffvk..... 92
cb3.....ffvk..... 105
r14.....ffvk..... 100
r2.....ffvk..... 100

fa10 .....t.ttttggaa..... 88
fa12 .....t.ttttggaa..... 88
foik .....t.ttttggaa..... 146
e.....t.ttttggaa..... 159
p1m.....t.ttttggaa..... 159
p1s.....t.ttttggaa..... 159
p2s.....t.ttttggaa..... 157
p3s.....t.ttttggaa..... 143
cb3.....t.ttttggaa..... 152
r14.....t.ttttggaa..... 145
r2.....t.ttttggaa..... 145

fa10 .....t.ttttggaa..... 132
fa12 .....t.ttttggaa..... 132
foik .....t.ttttggaa..... 133
e.....t.ttttggaa..... 204
p1m.....t.ttttggaa..... 212
p1s.....t.ttttggaa..... 212
p2s.....t.ttttggaa..... 210
p3s.....t.ttttggaa..... 196
cb3.....t.ttttggaa..... 204
r14.....t.ttttggaa..... 198
r2.....t.ttttggaa..... 198

```



## Conclusion

- strings are everywhere: DNA is a string of discrete symbols, music ...
- systems biology: cell as a system regulated by the DNA
- pattern discovery rather than pattern matching: statistical techniques, data mining, machine learning

