

Algorithmen der Sequenzanalyse: verrauschte Motive

AlgSeq – 28/10/2024

Prof. M. Sammeth

Inexakte Vorkommnisse eines Musters

Motive mit einer Abweichung (im *oriC* von *Vibrio cholerae*):

atca**ATGATCAAC**gtaagcttctaagc**ATGATCAAG**gtgctcacacagtttatccacaac
 ctgagtggatgacatcaagataggtcggtgtatctccttctctcgactctcatgacca
 cggaaag**ATGATCAAG**agaggatgatttcttggccatatcgcaatgaatacttgtgactt
 gtgcttccaattgacatcttcagcgccatattgctgctggccaaggtgacggagcgggatt
 acgaaag**CATGATCAT**ggctggttgttctgtttatcttgttttgactgagacttgttagga
 tagacgggtttttcatcactgactagccaaagccttactctgcctgacatcgaccgtaa
 tgataatgaatttacatgcttccgcgcgacgatttacct**CTTGATCAT**cgatccgattgaag
 atcttcaattgttaattctcttgcctcgactcatagccatgatgagct**CTTGATCAT**ggt
 tccttaaccctctattttttacggaaga**ATGATCAAG**ctgctgct**CTTGATCAT**cgtttc

ATGATCAAC und **CATGATCAT**



Acht (statt sechs) Vorkommnisse
eines Nonameres!

$$p_f(N, k, t) \approx \frac{\binom{N-t(k-1)}{t}}{4^{(t-1)k}}$$



$$p_f(540, \mathbf{9}, \mathbf{8}) = 7.242252695405615\mathbf{e-22}$$

k	t	$p_f(540, 9, t)$
7	9	7.451517316341352e-16
8	7	1.6504527748087858e-14
9	6	1.543363450054106e-14
10	4	2.3042467465345107e-09
11	4	3.487107419497277e-11
12	3	7.671120627605887e-08

Mismatches und die Hamming Distanz

Vergleich zweier k -mere p und q

$p:$	p_1	p_2	p_3	\dots	p_k
$q:$	q_1	q_2	q_3	\dots	q_k

$p_i \neq q_i \rightarrow$ ein Mismatch

$\text{HAMMINGDISTANZ}(p, q) :=$ Anzahl aller Mismatches über alle k Positionen

$\text{HAMMINGDISTANZ}(p, q) \leq d \iff$ höchstens d Mismatches im Vergleich p mit q

$\text{Zählen}_d(\text{Text}, \text{Muster}) :=$ Anzahl der Vorkommnisse von Muster in Text
mit höchstens d Mismatches.

Bsp.:

$\text{Zählen}_1(\text{AACAA GCATAACA TTAAGAG, AAAAA}) = 4$, d.h. $\{\text{AACAA}, \text{ATAAA},$

Inexaktes Zählen und Häufige Wörter mit Mismatches

Algorithm: ZÄHLEN(*Text*, *Muster*, *d*)

Zähler $\leftarrow 0$

for *i* $\leftarrow 0$ bis $|Text| - |Muster|$ **do**

Wort $\leftarrow Text(i, |Muster|)$

if HAMMINGDISTANZ(*Wort*, *Muster*) $\leq d$ **then**

Zähler $\leftarrow Zähler + 1$

return *Zähler*

Häufige Wörter mit RC und Mismatches Problem: Finde das/die häufigsten *k*-mer(e) die (zusammen mit dem reversen Komplement und mit höchstens *d* Mismatches) in einer Zeichenkette vorkommen.

Eingabe: Ein String *Text*, ganze Zahlen *k* und *d*.

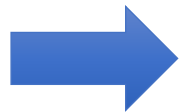
Ausgabe: Alle *k*-mere *Muster*, die $ZÄHLEN_d(Text, Muster) + ZÄHLEN_d(Text, Muster')$ ihres reversen Komplements *Muster'* maximieren (über alle *k*-mere in *Text*).

Unmittelbare Nachbarn ($d \leq 1$)

Unmittelbare Nachbarn := Zwei Zeichenketten p und q , die sich in höchstens einem Mismatch voneinander unterscheiden ($\text{HAMMINGDISTANZ}(p, q) \leq 1$)

Algorithm: UNMITTELBARENACHBARN($Muster$)

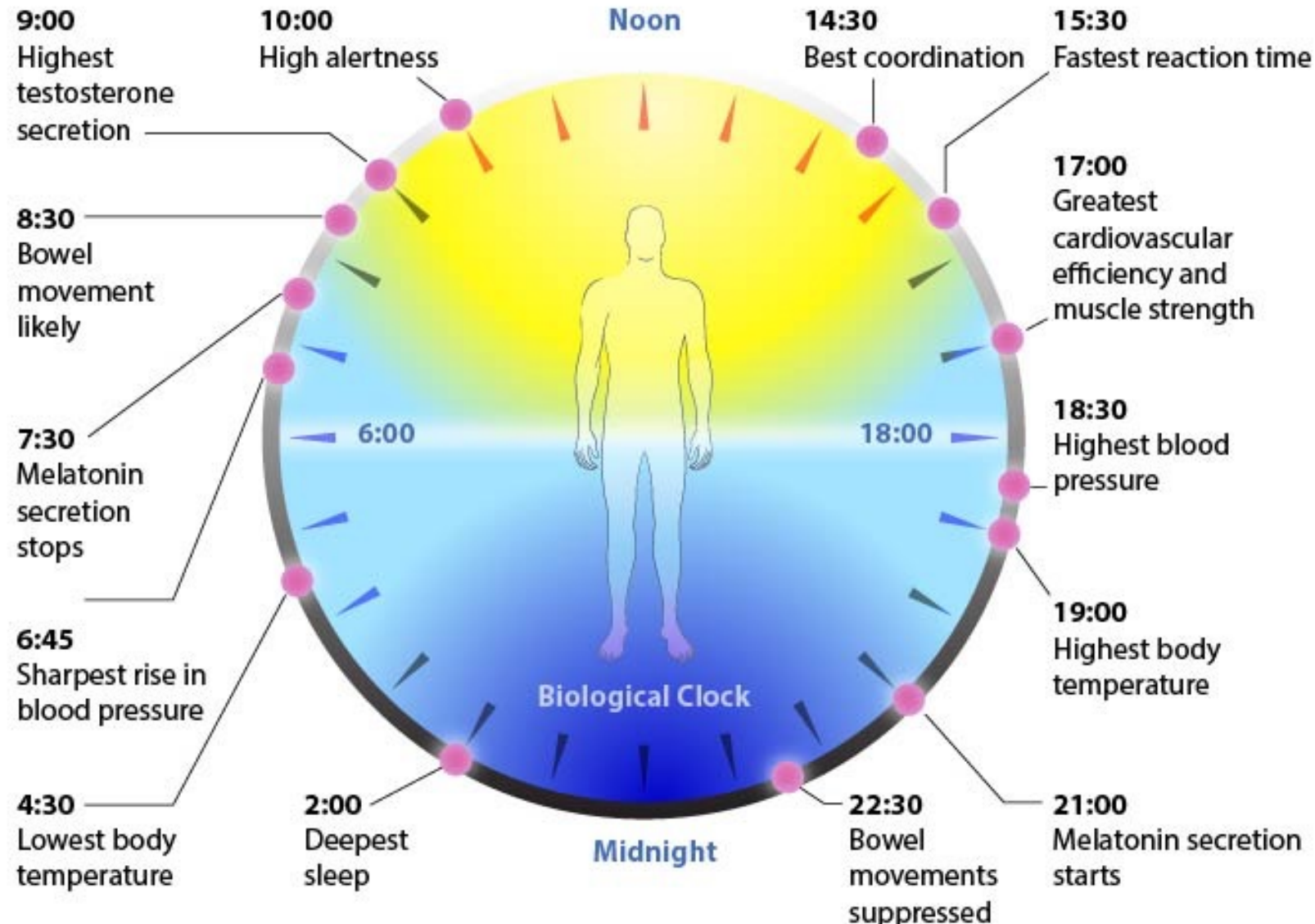
```
Nachbarschaft  $\leftarrow \{Muster\}$ 
for  $i = 1$  to  $|Muster|$  do
     $Symbol \leftarrow i$ -tes Nukleotid von  $Muster$ 
    for jedes Nukleotid  $x$  unterschiedlich von  $Symbol$  do
         $Nachbar \leftarrow Muster$  mit dem  $i$ -ten Nukleotid substituiert durch  $x$ 
         $Nachbarschaft \leftarrow Nachbarschaft \cup \{Nachbar\}$ 
return  $Nachbarschaft$ 
```



Häufige Wörter mit RC und Unmittelbaren Nachbarn Problem:
Finde das/die häufigsten k -mere mit höchstens einem Mismatch

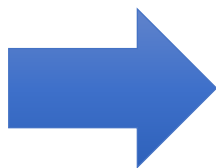
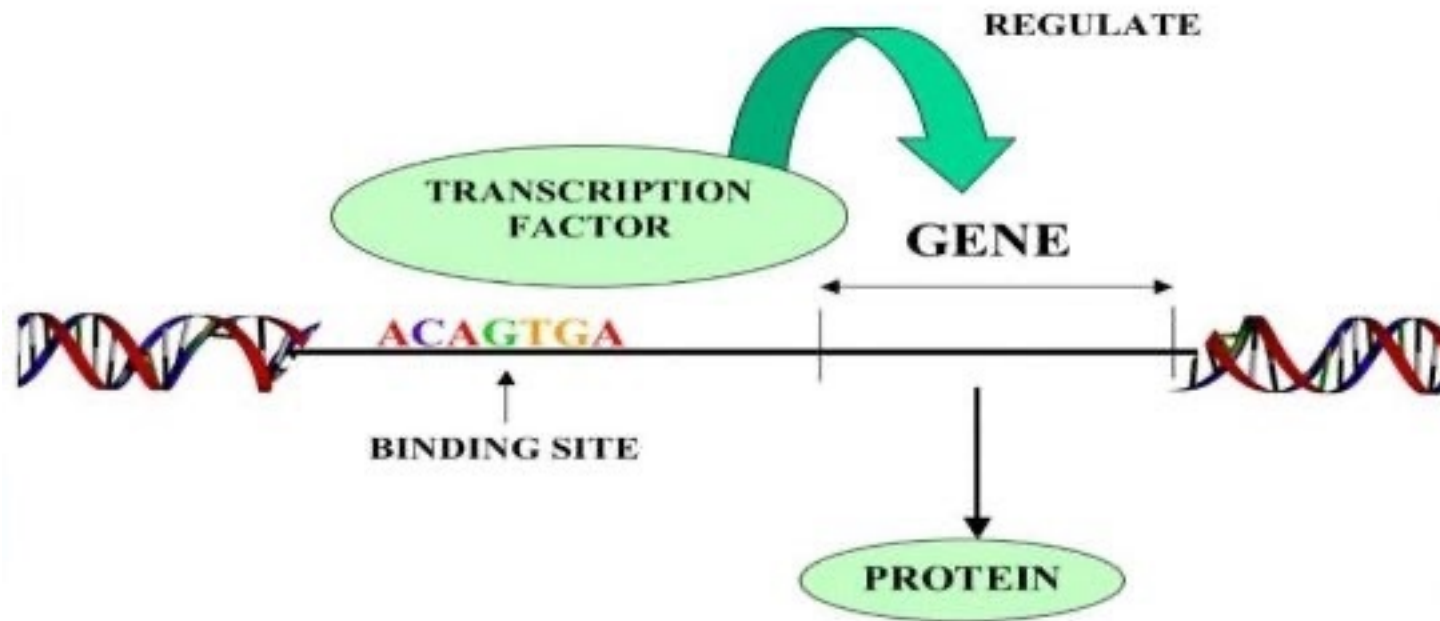
Transkriptionsfaktoren,
z.B. für “Molekulare Uhren”

Gibt es ein “Uhr-Gen”?



Transkriptionsfaktoren und ihre Bindestellen

Transkriptionsfaktoren sind Proteine, die Bindestellen in der DNA erkennen und benachbarte Gene an-/ab-schalten können.



- **mehrere Gene** (mit ähnlichen Funktionen) werden durch den/die gleiche(n) **Transkriptionsfaktor(en)** kontrolliert
- **Motiv** ("regulatorisches Motiv") der Bindestellen des **gleichen** Transkriptionsfaktors ist oft **inexakt**.

z. B. NF- κ B (Transkriptionsfaktor)

NF- κ B kontrolliert Gene der Immunantwort in der Fruchtfliege.

```
1  TCGGGGgTTTtt
2  cCGGtGAcTTaC
3  aCGGGGATTTtC
4  TtGGGGAcTTtt
5  TCGGGGATTTCC
6  TtGGGGAcTTCC
7  TCGGGGATcat
8  TCGGGGATTcCt
9  TaGGGGAacTaC
10 TCGGGtATaaCC
```



Motive von 10 Binstellen für NF- κ B,
inexaktes (nicht konserviertes) Motiv

Verstecken Spielen mit Motiven

(1) 15-mer Motiv **AAAAAAAAAGGGGGG** (exakt) implantiert in 10 Sequenzen:

```
1 atgaccgggatactgatAAAAAAAAAGGGGGGggcgtacacattagataaacgtatgaagtacgttagactcggcgccgccg
2 acccctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTTccgaataAAAAAAAAAGGGGGGa
3 tgagtatccctgggatgacttAAAAAAAAAGGGGGGtgctctcccgattTTTTgaatatgtaggatcattcgccaggggtccga
4 gctgagaattggatgAAAAAAAAAGGGGGGtccacgcaatcgcgaaaccaacgcggacccaaaggcaagaccgataaaggaga
5 tccctTTTgcggtaatgtgccgggaggctggttacgtagggaagccctaacggacttaatAAAAAAAAAGGGGGGcttatag
6 gtcaatcatgttcttgtgaatggatttAAAAAAAAAGGGGGGgaccgcttggcgcacccaaattcagtgtgggcgagcgcaa
7 cggTTTTggcccttgtagaggcccccgAAAAAAAAAGGGGGGcaattatgagagagctaattctatcgcggtgcgtgttcat
8 aacttgagttAAAAAAAAAGGGGGGctggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgta
9 ttggcccatTTggctaaaagcccaacttgacaaatggaagatagaatccttgcatAAAAAAAAAGGGGGGaccgaaaggggaag
10 ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttAAAAAAAAAGGGGGGa
```

(2) 15-mer Motiv **AAAAAAAAAGGGGGG** mit $d=4$ Mismatches implantiert in 10 Sequenzen:

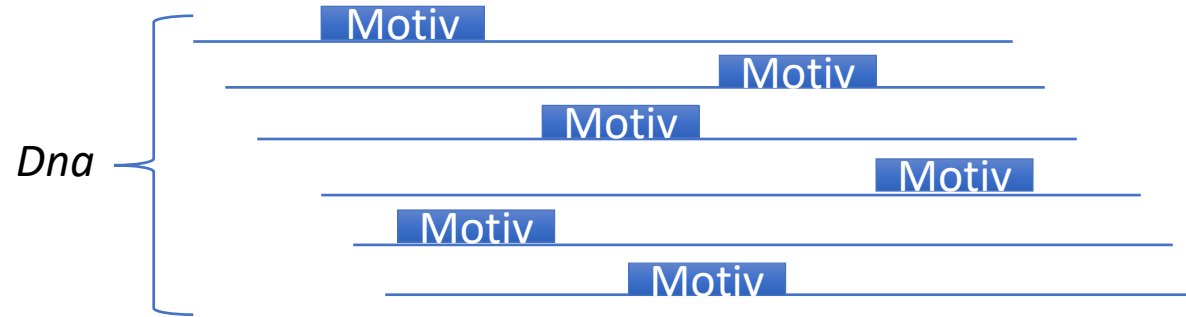
```
1 atgaccgggatactgatAgAAgAAAGGttGGGggcgtacacattagataaacgtatgaagtacgttagactcggcgccgccg
2 acccctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTTccgaatacAAtAAAcGGcGGGa
3 tgagtatccctgggatgacttAAAAtAAtGGAgtGGtgctctcccgattTTTTgaatatgtaggatcattcgccaggggtccga
4 gctgagaattggatgcAAAAAAAGGGattGtccacgcaatcgcgaaaccaacgcggacccaaaggcaagaccgataaaggaga
5 tccctTTTgcggtaatgtgccgggaggctggttacgtagggaagccctaacggacttaatAtAAtAAAGGaaGGGcttatag
6 gtcaatcatgttcttgtgaatggatttAAcAAtAAGGGctGGgaccgcttggcgcacccaaattcagtgtgggcgagcgcaa
7 cggTTTTggcccttgtagaggcccccgAtAAAcAAGGaGGGc caattatgagagagctaattctatcgcggtgcgtgttcat
8 aacttgagttAAAAAAtAGGGaGccctggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgta
9 ttggcccatTTggctaaaagcccaacttgacaaatggaagatagaatccttgcatActAAAAAGGaGcGGaccgaaaggggaag
10 ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttActAAAAAGGaGcGGa
```



(1) lösbar, aber (2) unlösbar mit Algorithmen für HäufigeWörter

Problem Formulierung

(k,d) -Motiv := Sei Dna eine Menge von Strings und d ein Integer, dann ist ein k -mer ein **(k,d) -Motiv**, wenn es in jedem String aus Dna mit höchstens d Mismatches vorkommt.



Implantiertes Motiv Problem:

Finde alle (k,d) -Motive in einer Menge von Zeichenketten.

Input: Eine Menge von Zeichenketten Dna , Integer k und d .

Output: Alle (k,d) -Motive in Dna .

Implantiertes Motiv mit roher Gewalt lösen

hier betrachtete Instanz des Problems (NF-κB Bindemotive): **Subtiles Motif Problem**

*Das 15-mer Motiv **AAAAAAAAAGGGGGG** wurde mit vier zufälligen Mutationen implantiert in zehn 600nt (=typische Länge von upstream regulatorischen Regionen) Sequenzen.*

```
1 atgaccgggatactgatAgAAAgAAGGttGGGggcgtacacattagataaacgtatgaagtacgttagactcggcgccgcccg
2 acccctattttttgagcagatttagtgacctggaaaaaaaatttgagtacaaaactttttccgaatacAAtAAAcGGcGGGa
3 tgagtatccctgggatgacttAAAAtAAtGGaGtGGtgctctcccgatttttgaatatgtaggatcattcgccaggggtccga
4 gctgagaattggatgAAAAAAAAAGGGattGtccacgcaatcgcaaccaacgcggacccaaaggcaagaccgataaaggaga
5 tcccttttgcggtaatgtgccgggaggctggttacgtagggaagccctaacggacttaatAtAAtAAAGGaAGGcttataag
6 gtcaatcatgttcttgtgaatggatttAAcAAtAAGGGctGGgaccgcttggcgcacccaaattcagtgtgggcgagcgcaa
7 cggttttggcccttgtttagaggcccccgtAtAAAcAAGGaGGGccaattatgagagagctaattctatcgcggtgcgtgttcat
8 aacttgagttAAAAAAtAGGGaGccctggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgta
9 ttggcccatgggctaaaagcccaacttgacaaatggaagatagaatccttgcActAAAAAGGaGcGGaccgaaaggggaag
10 ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttActAAAAAGGaGcGGa
```

Brute Force (Rohe Gewalt) := Algorithmische Vorgehensweise, bei der alle
(auch: exhaustive Suche) möglichen Lösungen geprüft werden, und dann
die beste Lösung vorgeschlagen wird.



wie lautet ein “brute force” Ansatz für
das **Implantiertes Motiv** Problem?



Implantiertes Motiv mit roher Gewalt lösen

ein “brute force” Ansatz für das **Implantiertes Motiv** Problem:

Algorithm: MOTIVEAUFZAEHLEN(Dna, k, d)

$Patterns \leftarrow \emptyset$

for jedes k -mer $Muster$ in Dna **do**

for jedes k -mer $Muster'$, das sich von $Muster$ mit höchstens d
 Unterschieden (“Mismatches”) unterscheidet **do**

if $Muster'$ in jeder Zeichenkette von Dna mit höchstens d
 Mismatches vorkommt **then**

$Patterns \leftarrow Patterns \cup \{Muster'\}$

entferne alle Duplikate aus $Patterns$

return $Patterns$



paarweise Vergleiche zwischen Sequenzen helfen kaum (s.u.) \rightarrow
zwei Sprünge im Vergleich: $Muster \rightarrow Muster' \rightarrow ZÄHLEN_2(Muster')$



sehr langsam für große Werte von k und/oder d !

AgAAgAAAGGttGGG
||| | |||
cAAtAAAAcGGGGcG