

# Algorithmen der Sequenzanalyse: Motivanalysen

AlgSeq – 11/11/2024 (Helau)

Prof. M. Sammeth

# Definitionen: Motive bewerten

(i) Distanz eines *Musters* zu einer *Zeichenkette*:

$$d(\text{Muster}, \text{Text}) := \min_{\text{über alle } k\text{-mere } \text{Muster}' \text{ in Text}} \text{HAMMINGDISTANZ}(\text{Muster}, \text{Muster}')$$

z.B.:

$$d(\text{GATTCATCA}, \text{gcaaaGACGCTGAccaa}) = 3$$

(ii) Distanz eines *Musters* zu einer Menge von *Zeichenketten*:

$$d(\text{Muster}, \text{Dna}) := \sum_{i=1}^t d(\text{Muster}, \text{Dna}_i)$$

z.B.:

$$d(\text{AAA}, \text{Dna}) = 1 + 1 + 2 + 0 + 1 = 5$$

$$\text{Dna} = \left\{ \begin{array}{ll} \text{ttaccttAAC} & 1 \\ \text{gATAtctgtc} & 1 \\ \text{ACGgcgttcg} & 2 \\ \text{ccctAAAgag} & 0 \\ \text{cgtcAGAggt} & 1 \end{array} \right.$$

# Definitionen: Motive bewerten

## (iii) Ähnlichstes Muster in einer Zeichenkette

$\text{MOTIV}(\text{Muster}, \text{Text}) := \underset{\text{über alle } k\text{-mere } \text{Muster}' \text{ in Text}}{\text{argmin}} \quad \text{HAMMINGDISTANZ}(\text{Muster}, \text{Muster}')$

z.B.:

$\text{MOTIV}(\text{GATTCTCA}, \text{gacaaaGACGCTGAccaa}) = \text{GACGCTGA}$

(kann mehrdeutig sein!)

## (iv) Ähnlichste Muster in einer Menge von Zeichenketten

$\text{MOTIVE}(\text{Muster}, \text{Dna}) := \bigcup_{i=1}^t \text{MOTIV}(\text{Muster}, \text{Dna}_i)$

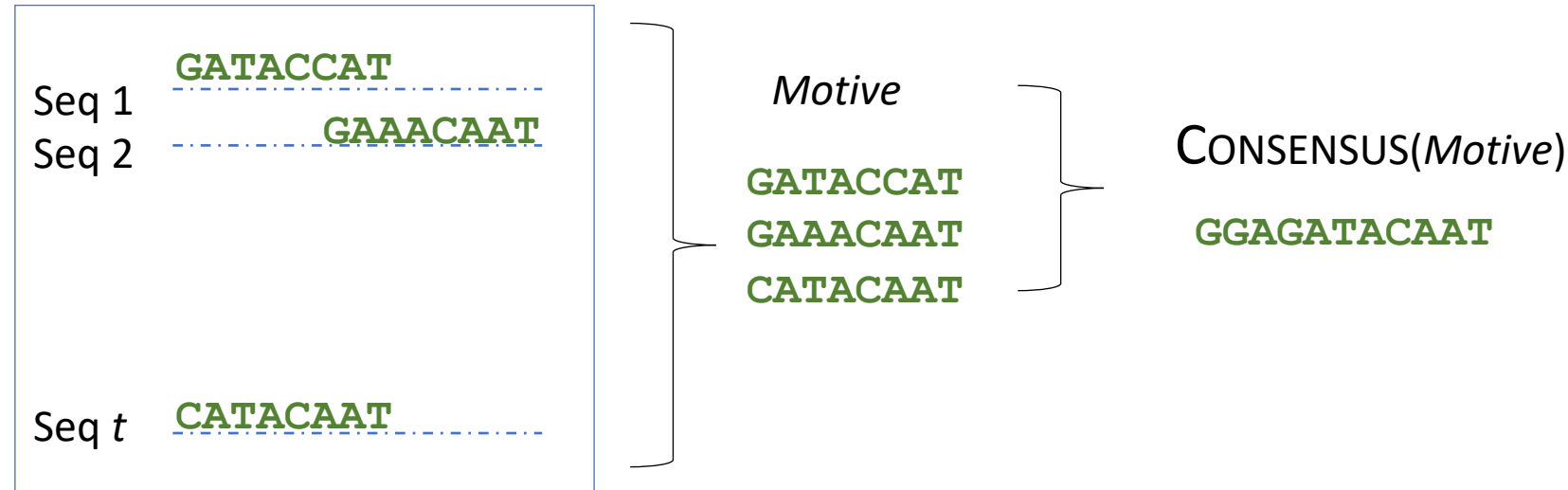
z.B.:  $\text{MOTIVE}(\text{AAA}, \text{Dna}) = \{\text{AAC}, \text{ATA}, \text{ACG}, \text{AAA}, \text{AGA}\}$

$\text{Dna} = \left\{ \begin{array}{l} \text{ttaccttAAC} \\ \text{gATAtctgtc} \\ \text{ACGgcgttcg} \\ \text{ccctAAAgag} \\ \text{cgtcAGAggt} \end{array} \right.$

# Eine neue Perspektive

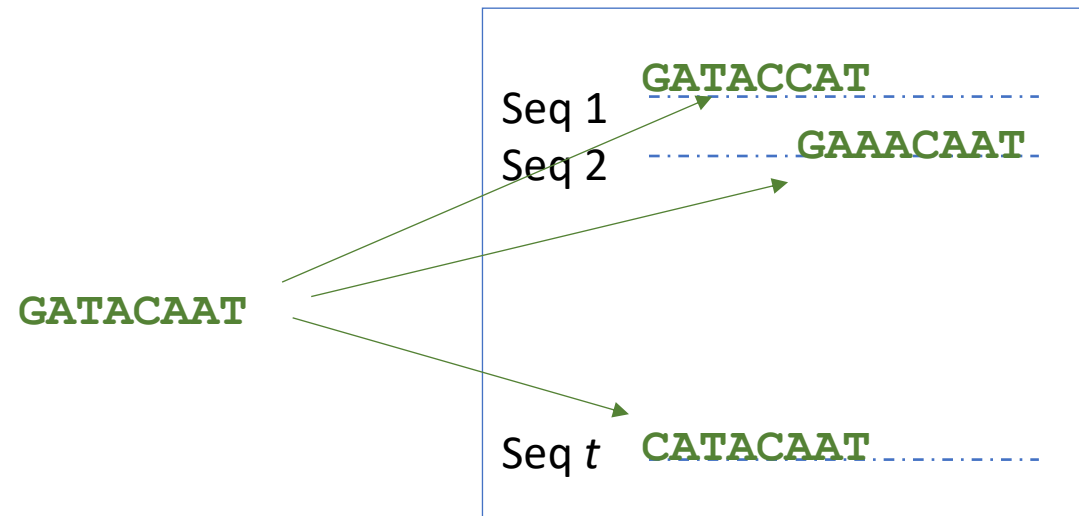
BRUTEFORCEMOTIVSUCHE:

$\forall \text{Motive} \rightarrow \text{suche } \text{CONSENSUS}(\text{Motive})$



Alternative Strategie:

$\forall \text{CONSENSUS}(\text{Motive}) \rightarrow \text{suche } \text{Motive}$



# Re-Definition des Problemes

---

## **Motiv-Findungs Problem:**

*Gegeben eine Menge von Zeichenketten, finde eine Menge von k-meren, mit je einem k-mer von jeder Zeichenkette, welche das Score des intrinsischen Motives minimiert.*

**Eingabe:** Eine Menge von Zeichenketten *Dna* und eine ganze Zahl *k*.

**Ausgabe:** Eine Menge *Motive* mit k-meren, eines von jeder Zeichenkette in *Dna*, welche  $\text{SCORE}(\text{Motive})$  für alle möglichen *k*-mere minimiert.

---

## **Äquivalentes Motiv-Findungs Problem:**

*Gegeben eine Menge von Zeichenketten, finde ein Muster und eine Menge von k-meren, (mit je einem k-mer von jeder Zeichenkette), welches die Distanz zwischen allen möglichen Mustern über alle möglichen Mengen von k-meren minimiert.*

**Eingabe:** Eine Menge von Zeichenketten *Dna* und eine ganze Zahl *k*.

**Ausgabe:** Ein *k*-mer *Muster* und eine Menge *Motive* mit *k*-meren (eines von jeder Zeichenkette in *Dna*), welche  $d(\text{Muster}, \text{Motive})$  für alle möglichen *Muster* und *Motive* minimiert.

---

# Summen von SCORE(Motive) sind kommutativ

$$\text{SCORE}(\text{Motive}) = \sum_{\text{Spalten } j} \sum_{\text{Zeilen } i} I_{i,j} = \sum_{\text{Zeilen } i} \sum_{\text{Spalten } j} I_{i,j} \quad , \text{ mit } I_{i,j} = \begin{cases} 0 \text{ iff } \text{Dna}_i(j, 1) = \text{CONSENSUS}(j, 1) \\ 1 \text{ andersfalls} \end{cases}$$

<i>Motive</i>	T	C	G	G	G	G	g	T	T	T	t	t	3
	c	C	G	G	t	G	A	c	T	T	a	C	+ 4
	a	C	G	G	G	G	A	T	T	T	t	C	+ 2
	T	t	G	G	G	G	A	c	T	T	t	t	+ 4
	a	a	G	G	G	G	A	c	T	T	C	C	+ 3
	T	t	G	G	G	G	A	c	T	T	C	C	+ 2
	T	C	G	G	G	G	A	T	T	c	a	t	+ 3
	T	C	G	G	G	G	A	T	T	c	C	t	+ 2
	T	a	G	G	G	G	A	a	c	T	a	C	+ 4
	T	C	G	G	G	t	A	T	a	a	C	C	+ 3

$$\text{SCORE}(\text{Motive}) \quad 3+ \quad 4+ \quad 0+ \quad 0+ \quad 1+ \quad 1+ \quad 1+ \quad 5+ \quad 2+ \quad 3+ \quad 6+ \quad 4 \quad = \quad 30$$

CONSENSUS(*Motive*)    T   C   G   G   G   G   A   T   T   T   C   C

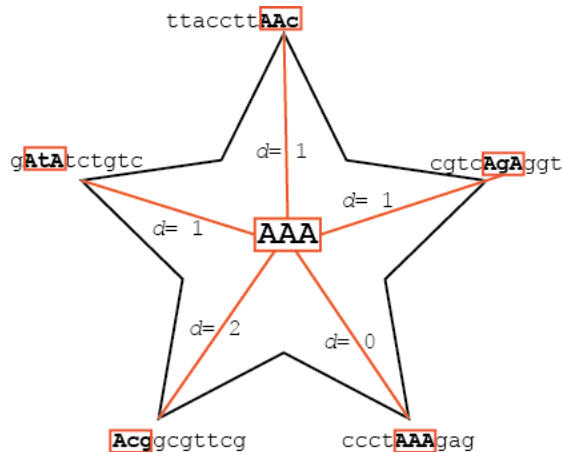
# Die Median-Zeichenkette

## Median-Zeichenketten Problem:

*Finde die Median Zeichenkette.*

**Eingabe:** Eine Menge  $Dna$  von Zeichenketten und eine ganze Zahl  $k$ .

**Ausgabe:** Ein  $k$ -mer *Muster*, das  $d(\text{Muster}, Dna)$  über alle  $k$ -mere Muster minimiert.



## Algorithm: MEDIANZEICHENKETTE( $Dna, k$ )

$Distanz \leftarrow \infty$

**for** jedes  $k$ -mer von AA...AA bis TT...TT **do**

**if**  $Distanz > d(\text{Muster}, Dna)$  **then**

$Distanz \leftarrow d(\text{Muster}, Dna)$

$Median \leftarrow \text{Muster}$

**return**  $Median$



# Laufzeit-Vergleich

$$\mathcal{O}(\text{BRUTEFORCEMOTIVSUCHE}) \\ \in \mathcal{O}(n^t \cdot k \cdot t)$$

vs.

$$\mathcal{O}(\text{MEDIANZEICHENKETTE}) \\ \in \mathcal{O}(4^k \cdot k \cdot n \cdot t)$$

mit  $k \leq 20$  (Motivlänge) und  
 $t \sim 10^3$  (tausende Sequenzen)



2

Mr. BRUTEFORCE



1

Mr. MEDIANZEICHENKETTE



# Praktische Laufzeit

```
1 atgaccgggatactgatAgAAAgAAGGttGGGggcgtacacattagataaacgtatgaagtacgttagactcggcgccgccg
2 acccctatTTTTTgagcagatttagtgacctggaaaaaaaatttgagtacaaaactTTTccgaatacAAtAAAcGGcGGGa
3 tgagtatccctgggatgacttAAAAtAAtGGaGtGGtgctctcccgattTTTtgaatatgtaggatcattcgccaggggtccga
4 gctgagaattggatgAAAAAAAAGGGattGtccacgcaatcgcgaaaccaacgcggacccaaaggcaagaccgataaaggaga
5 tccctTTTgcggtaatgtgccgggaggctggTTacgtagggaagccctaacggacttaatAtAAtAAAGGaGGGcttatag
6 gtcaatcatgttcttgtgaatggatttAAcAAtAAGGGctGGgaccgcttggcgcacccaaattcagtgtgggcgagcgcaa
7 cggTTTTggcccttgTTtagaggcccccgtAtAAAcAAGGaGGGccaattatgagagagctaattctatcgcgTgcgtgttcat
8 aacttgagttAAAAAtAGGGaGccctggggcacatacaagaggagtcttcttatcagttaatgctgtatgacactatgta
9 ttggccccattggctaaaagcccaacttgacaaatggaagatagaatccttgcatActAAAAAGGaGcGGaccgaaaggggaag
10 ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttActAAAAAGGaGcGGa
```

**Subtiles Motif Problem** (upstream regionen):

15-mer Motiv **AAAAAAAGGGGGG** ist implantiert in  
zehn 600nt Sequenzen.

$$\mathcal{O}(\text{MEDIANZEICHENKETTE}) \\ \in \mathcal{O}(4^k \cdot k \cdot n \cdot t)$$



**$4^{15} = 1.073.741.824$**   
**( > 1 Mrd =  $10^9$  !)**

# “Kleines” Subtiles Problem

```
1 atgaccgggatactgatAgAAAgAAGGttGGGggcgtacacattagataaacgtatgaagtacgttagactcggcgccgccg
2 acccctatTTTTTgagcagatttagtgacctggaaaaaaaatttgagtacaaaactTTTccgaatacAAtAAAaGGcGGGa
3 tgagtatccctgggatgacttAAAAtAAtGGaGtGGtgctctcccgattTTTgaatatgtaggatcattcgccaggggtccga
4 gctgagaattggatgAAAAAAAGGGattGtccacgcaatcgcgaaaccaacgcggacccaaaggcaagaccgataaaggaga
5 tccTTTTgCGGtaatgtgccgggaggctggTTacgtagggaagccctaacggacttaatAtAAtAAAGGaGGGcttatag
6 gtcaatcatgttcttTgtgaatggatttAAcAAtAAGGGctGGgaccgcttggcgcacccaaattcagtgtgggCGagcgcaa
7 cggTTTTggcccttgTTtagaggcccccgTAtAAAcAAGGaGGGccaattatgagagagctaattctatcgCGtgCGtgTtcat
8 aacttgagttAAAAAAtAGGGaGccctggggcacatacaagaggagtcttcttatcagTTaatgctgtatgacactatgta
9 ttggcccatTggctaaaagcccaacttgacaaatggaagatagaatccttgcatActAAAAAGGaGcGGaccgaaaggggaag
10 ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttActAAAAAGGaGcGGa
```



MEDIANZEICHENKETTE() für  $k=13$  laufen lassen, und “hoffen”,  
dass es uns einen Teil (d.h. Subsequenz) der Lösung gibt

$4^{13} = 67.108.864$  ( $\sim 10^7$ )



$> \frac{1}{2}$  Tag



gefundenes Motiv:  
**AAAAAtAGaGGG**

# Greedy Algorithmen



## **Brute Force**

**Algorithmus** :=

enumeriert exhaustiv alle  
möglichen Ergebnisse und  
wählt davon dann die beste  
Lösung.

vs.

## **Greedy Algorithmus** :=

iteratives Vorgehen, bei dem  
in jedem Durchlauf die “attraktivste”  
Möglichkeit gewählt wird.



exhaustive Suchraum Analyse:

- findet optimale Lösung(en),  
per Definition
- aber die Berechnung kann  
sehr lange dauern



inexakte Heuristik:

- kann desaströses Ergebnis  
produzieren
- aber für viele (biologische)  
Probleme nützlich

# Mit Profilen Würfeln

*Motive*

```

T C G G G G g T T T t t
c C G G t G A c T T a C
a C G G G G A T T T t C
T t G G G G A c T T t t
a a G G G G A c T T C C
T t G G G G A c T T C C
T C G G G G A T T c a t
T C G G G G A T T c C t
T a G G G G A a c T a C
T C G G G t A T a a C C
    
```

PROFIL(*Motive*)

A:	.2	.2	0	0	0	0	.9	.1	.1	.1	.3	0
C:	.1	.6	0	0	0	0	0	.4	.1	.2	.4	.6
G:	0	0	1	1	.9	.9	.1	0	0	0	0	0
T:	.7	.2	0	0	.1	.1	0	.5	.8	.7	.3	.4



$Pr(\text{Sequenz} \mid \text{PROFIL}(\text{Motive})) :=$

Wahrscheinlichkeit, dass PROFIL das  $k$ -mer *Sequenz* generiert.

$$= \prod_{i=1}^k \text{PROFIL}(\text{Motive})_{\text{Sequenz}(i),i}$$

# Das “Profil-wahrscheinlichste” $k$ -mer

$Pr(\text{Sequenz} \mid \text{PROFIL}(\text{Motive}))$  erlaubt Sequenzen nach PROFIL zu bewerten,  
→ zu vergleichen.

---

## **Profil-wahrscheinlichstes $k$ -mer Problem:**

*Finde das PROFIL-wahrscheinlichste  $k$ -mer einer Zeichenkette.*

**Eingabe:** Eine Zeichenkette *Text*, eine ganze Zahl  $k$ , eine  $4 \times k$  Matrix PROFIL.

**Ausgabe:** Ein  $k$ -mer aus *Text*, das über alle  $k$ -mere in *Text* am Wahrscheinlichsten aus PROFIL generiert wird.

---



wie kann dieses Optimierungskriterium als Zielfunktion einer *greedy* Suche verwendet werden?

# Motiv-Suche “Greedy”

---

**Algorithm:** MOTIVSUCHEGREEDY( $Dna, k, t$ )

---

$BesteMotive \leftarrow$  Matrix von Motiven aus den ersten  $k$ -meren der Zeichenketten in  $Dna$

**for** jedes  $k$ -mer  $Motiv$  in der ersten Zeichenkette von  $Dna$  **do**

- $Motiv_1 \leftarrow Motiv$
- for**  $i \leftarrow 2$  bis  $t$  **do**
  - bilde  $Profile$  aus den Motiven  $Motiv_1, \dots, Motiv_{i-1}$
  - $Motiv_i \leftarrow$   $Profil$ -wahrscheinlichstes  $k$ -mer in der  $i$ -ten Zeichenkette von  $Dna$
- $Motive \leftarrow (Motiv_1, \dots, Motiv_t)$
- if** SCORE( $Motive$ ) < SCORE( $BesteMotive$ ) **then**
  - $BesteMotive \leftarrow Motive$

**return**  $BesteMotive$

---

# Das Problem von MotivSucheGreedy

Minimalbsp: finde das (4,1)-Motiv **ACGT** implantiert in folgenden Strings *Dna*:

*Dna*=

tt	<b>ACCT</b>	taac
g	<b>ATGT</b>	ctgtc
acg	<b>GCGT</b>	tag
cccta	<b>ACGA</b>	g
cgtcag	<b>AGGT</b>	

für  $Motiv_1 = \text{ACCT}$  (das implantierte Motiv) ergibt sich  $\Pr(\text{ATGT}) = 0$

ZAEHLEN(*Motive*)=

A:	<b>1</b>	0	0	0
C:	0	<b>1</b>	<b>1</b>	0
G:	0	0	0	0
T:	0	0	0	<b>1</b>



Matrix ist *sparse*, alle Wahrscheinlichkeiten von **abweichenden Motiven** sind gleich und **gleich 0** !

# Historische Zitate zu Wahrscheinlichkeit 0 und 1

Bsp: Wahrscheinlichkeit, dass die Sonne morgen nicht aufgeht.

5000 Jahre lang wurde überliefert, dass die Sonne aufgeht,  
 $p=0$  dass sie morgen nicht aufgeht

## Cromwell's Regel

*Abgesehen von logischen Zuständen,  
die **wahr** oder falsch sein können,  
sollten bei empirischen Abschätzungen  
**0** und **1** nicht als Wahrscheinlichkeiten  
verwendet werden.*



Oliver Cromwell (vor Angriff auf die Schotten) 1650

## Laplace's Nachfolger Regel (Pseudocounts):

geht davon aus, dass bei der nächsten Beobachtung das bisher uneingetretene Ereignis vorkommt.

( $p=1/1826251$ , dass die Sonne morgen nicht aufgeht)



Pierre-Simon Laplace  
1749-1827  
French Physicist and mathematician

*"I had no need of that hypothesis."* A famous answer to the question from Napoleon about why he didn't mention the name of God in his work



# MotivSucheGreedy mit Laplace's Nachfolger Regel

Dna=

tt**ACCT**taac  
g**ATGT**ctgtc  
acg**GCGT**tag  
cccta**ACGA**g  
cgtcag**AGGT**

Bsp: implantiertes (4,1)-Motiv **ACGT**

erste Sequenz tt**ACCT**taac



Motive= **ACCT**

ZAEHLEN(Motive)=

A:	<b>1</b> +1	0+1	0+1	0+1
C:	0+1	<b>1</b> +1	<b>1</b> +1	0+1
G:	0+1	0+1	0+1	0+1
T:	0+1	0+1	0+1	<b>1</b> +1

Motiv-Matrix mit **ACCT** aus der ersten Sequenz:

jedes Mögliche Ereignis wird einmal mehr als beobachtet gezählt (Pseudocounts)

PROFIL(Motive)=

A:	2/5	1/5	1/5	1/5
C:	1/5	2/5	2/5	1/5
G:	1/5	1/5	1/5	1/5
T:	1/5	1/5	1/5	2/5

Summe der Beobachtungen vergrößert sich entsprechend auch

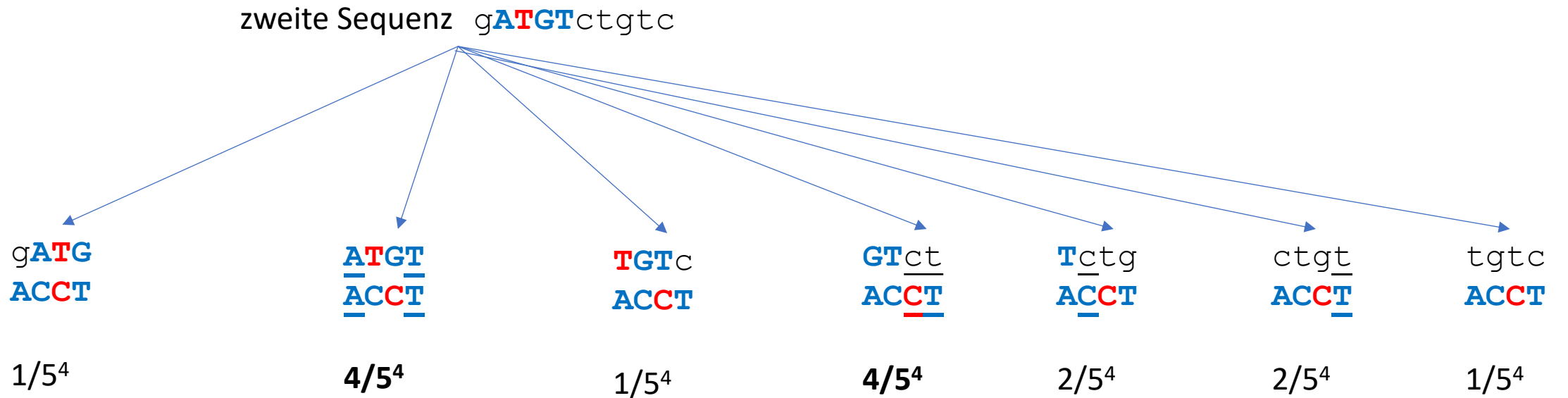
Motive= **ACCT**

ZAEHLEN(Motive)=

A:	<b>1</b> +1	0+1	0+1	0+1
C:	0+1	<b>1</b> +1	<b>1</b> +1	0+1
G:	0+1	0+1	0+1	0+1
T:	0+1	0+1	0+1	<b>1</b> +1

PROFIL(**ACCT**)=

A:	2/5	1/5	1/5	1/5
C:	1/5	2/5	2/5	1/5
G:	1/5	1/5	1/5	1/5
T:	1/5	1/5	1/5	2/5



Zwei profil-  
wahrscheinlichste Motive (!)

Motive= **ACCT**  
**ATGT**

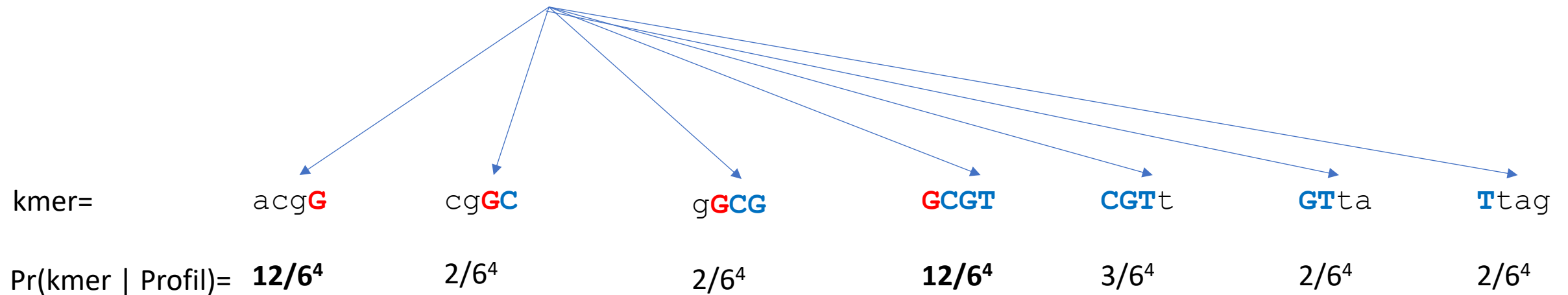
ZAEHLEN(Motive)=

A:	<b>2</b> +1	0+1	0+1	0+1
C:	0+1	<b>1</b> +1	<b>1</b> +1	0+1
G:	0+1	0+1	<b>1</b> +1	0+1
T:	0+1	<b>1</b> +1	0+1	<b>2</b> +1

PROFIL(**ACCT**)=

A:	3/6	1/6	1/6	1/6
C:	1/6	2/6	2/6	1/6
G:	1/6	1/6	2/6	1/6
T:	1/6	2/6	1/6	3/6

dritte Sequenz acg**G****CGT**tag



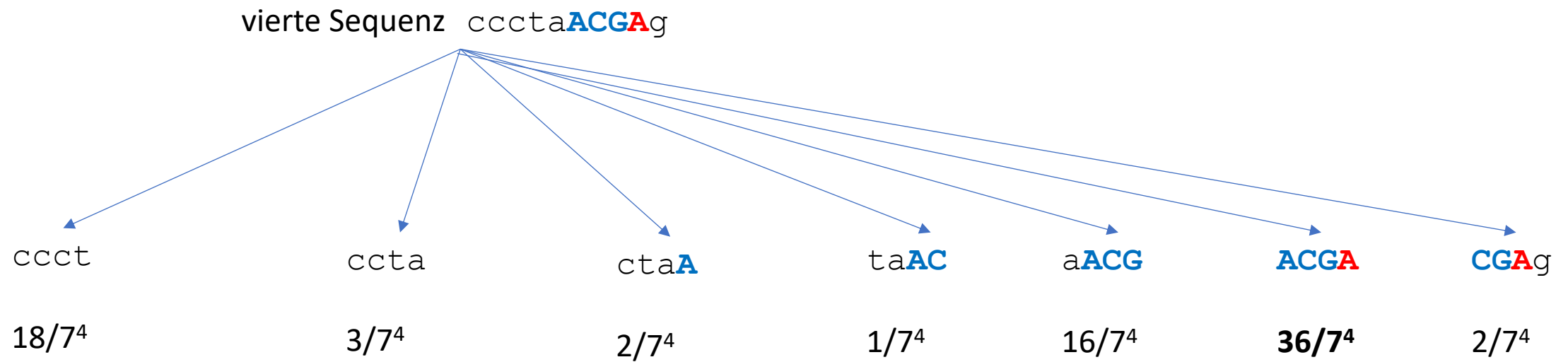
Es kann nur einen geben...



Motive= **ACCT**  
**ATGT**  
 acg**G**

ZAEHLEN(Motive)=  
 A: **3**+1 0+1 0+1 0+1  
 C: 0+1 **2**+1 **1**+1 0+1  
 G: 0+1 0+1 **2**+1 **1**+1  
 T: 0+1 **1**+1 0+1 **2**+1

PROFIL(**ACCT**)=  
 A: 4/7 1/7 1/7 1/7  
 C: 1/7 3/7 2/7 1/7  
 G: 1/7 1/7 3/7 2/7  
 T: 1/7 2/7 1/7 3/7



Korrektes Motiv



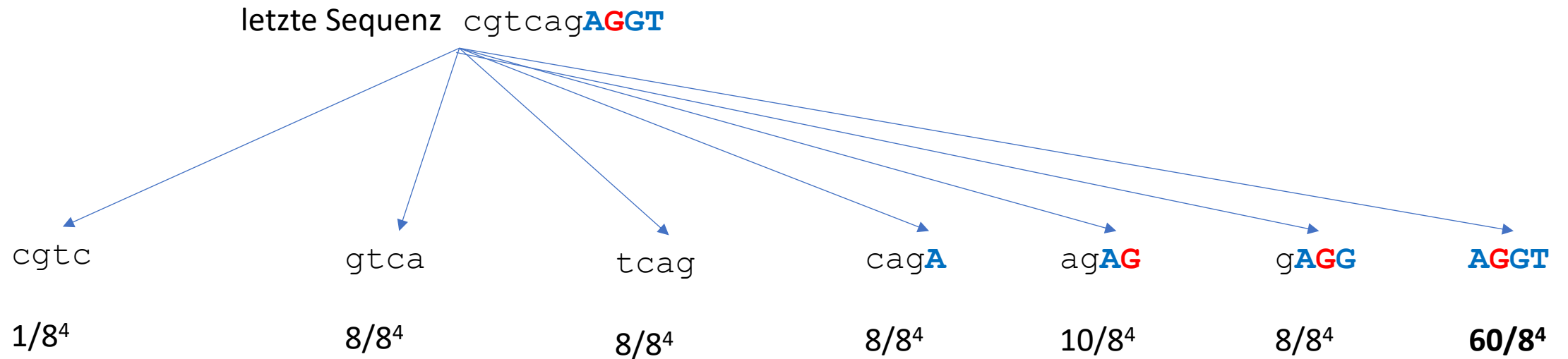
Motive= **ACCT**  
**ATGT**  
acg**G**  
**ACGA**

ZAEHLEN(Motive)=

A:	<b>4</b> +1	0+1	0+1	<b>1</b> +1
C:	0+1	<b>3</b> +1	<b>1</b> +1	0+1
G:	0+1	0+1	<b>3</b> +1	<b>1</b> +1
T:	0+1	<b>1</b> +1	0+1	<b>2</b> +1

PROFIL(**ACCT**)=

A:	5/8	1/8	1/8	2/8
C:	1/8	4/8	2/8	1/8
G:	1/8	1/8	4/8	2/8
T:	1/8	2/8	1/8	3/8



Korrektes Motiv

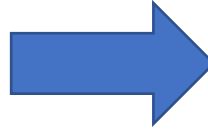
# Greedy Motiv-Suche

implantiertes (4,1)-Motiv **ACGT**



*Dna*=

tt**ACCT**taac  
g**ATGT**ctgtc  
acg**GCGT**tag  
cccta**ACGA**g  
cgtcag**AGGT**



*Motive*= **ACCT**  
**ATGT**  
acg**G**  
**ACGA**  
**AGGT**



*Consensus*= **ACGT**

# Zufalls-basierte Algorithmen

**Randomisierte** Algorithmen  
benutzen zufällige Elemente, um Lösungen zu finden



**Monte Carlo** Algorithmen (hier)

- können sub-optimale Lösungen finden  
(normalerweise mit einer kleinen Fehlerwahrscheinlichkeit)
- sind sehr schnell, können mehrfach laufen



**Las Vegas** Algorithmen:

- Garantie der besten/exakten Lösung
- Laufzeit oft nicht abschätzbar

# Die Profil-wahrscheinlichsten Motive von Sequenzen

$\text{PROFIL}(\text{Motive}) \quad := (4 \times k) \text{ Matrix berechnet über } \text{Motive}.$

$\text{Profil} \quad := \text{eine beliebige } (4 \times k) \text{ Matrix mit Wahrscheinlichkeitsverteilungen.}$

$\text{Pr}(k\text{-mer} \mid \text{Profil}) := \text{Wahrscheinlichkeit, dass ein } k\text{-mer von einem Profil generiert wird}$

- Profil-wahrscheinlichstes Motiv in einem *Text*:

$\text{MOTIV}(\text{Profil}, \text{Text}) := \underset{\text{Sequenz in Text}}{\text{argmin}} \quad \text{Pr}(\text{Sequenz} \mid \text{Profil})$

- Profil-wahrscheinlichste Motive aus jeder Sequenz einer Menge *Dna*:

$\text{MOTIVE}(\text{Profil}, \text{Dna}) \quad := \bigcup_{i=1}^t \text{MOTIV}(\text{Profil}, \text{Dna}_i)$



Beispiel, gegeben:

*Dna*=  
ttaccttaac  
gatgtctgtc  
acggcgttag  
ccctaacgag  
cgtcagaggt

PROFIL=  
A: 4/5 0 0 1/5  
C: 0 3/5 1/5 0  
G: 1/5 1/5 4/5 0  
T: 0 1/5 0 4/5

Profil-wahrscheinliche 4-mere:

MOTIVE(*Profil*, *Dna*) =  
tt**acct**taac  
g**atgt**ctgtc  
acg**gcgt**tag  
cccta**acga**g  
cgtcag**aggt**

PROFIL(*Motive*)=  
A: 4/5 0 0 1/5  
C: 0 3/5 1/5 0  
G: 1/5 1/5 4/5 0  
T: 0 1/5 0 4/5

usw.

PROFIL(MOTIVE(*Profil*, *Dna*)) → MOTIVE(PROFIL(MOTIVE(*Profil*, *Dna*))) →

PROFIL(MOTIVE(PROFIL(MOTIVE(*Profil*, *Dna*)))) → ...

# Randomisierte Motiv-Suche

---

**Algorithm:** RANDOMISIERTEMOTIVSUCHE( $Dna, k, t$ )

---

wähle zufällige  $k$ -mere  $Motive = (Motiv_1, \dots, Motiv_t)$  von je einer Sequenz in  $Dna$

$BesteMotive \leftarrow Motive$

**while** für immer **do**

$Profile \leftarrow \text{PROFILE}(Motive)$

$Motive \leftarrow \text{MOTIVE}(Profile, Dna)$

**if** SCORE( $Motive$ ) < SCORE( $BesteMotive$ ) **then**

$BesteMotive \leftarrow Motive$

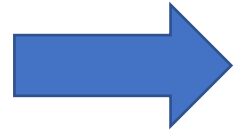
**else**

**return**  $BesteMotive$

---

Beispiel, ein implantiertes  
(4,1)-Motiv **ACGT**:

Dna= ttACCT**taac**  
gAT**GTct**gtc  
**ccgG**CGTtag  
c**acta**ACGAg  
cgtcag**AGGT**



Motive= **taac**  
**GTct**  
**ccgG**  
**acta**  
**AGGT**

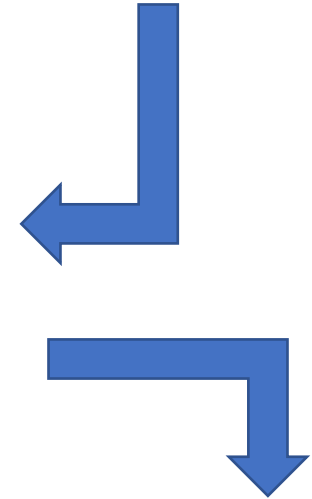


PROFIL(Motive)=

A: 0.4 0.2 0.2 0.2  
C: 0.2 0.4 0.2 0.2  
G: 0.2 0.2 0.4 0.2  
T: 0.2 0.2 0.2 0.4

Wahrscheinlichkeiten  $Pr(kmer, PROFIL(Motive))$ :

ttAC	tACC	ACCT	CCTt	CTta	Ttaa	taac
.0016	.0016	<b>.0128</b>	.0064	.0016	.0016	.0016
gATG	ATGT	TGTc	GTct	Tctg	ctgt	tgtc
.0016	<b>.0128</b>	.0016	.0032	.0032	.0032	.0016
ccgG	cgGC	gGCG	GCGT	CGTt	GTta	Ttag
.0064	.0032	.0016	<b>.0128</b>	.0032	.0016	.0016
cact	acta	ctaA	taAC	aACG	ACGA	CGAg
.0032	.0064	.0016	.0016	.0032	<b>.0128</b>	.0016
cgtc	gtca	tcag	cagA	agAG	gAGG	AGGT
.0016	.0016	.0016	.0032	.0032	.0032	<b>.0128</b>



MOTIVE(PROFIL(Motive))

**ACCT**  
**ATGT**  
**GCGT**  
**ACGA**  
**AGGT**

# Übungen 06

1. Implementieren Sie MEDIANZEICHENKETTE und vergleichen Sie Ihre Lösung mit dem im Seminar besprochenen, “reduzierten” Subtilen Motiv Problem für  $k=13$ .

*Eingabe:* Dna= Algo04\_subtiles\_motiv.txt,  $k=13$

*erwartete Ausgabe:* eine Matrix *Motive* mit  $\text{SCORE}(Motive)=29$  und  $\text{CONSENSUS}(Motive)=$  AAAAAtAGaGGGG.

2. Implementieren Sie GREEDYMOTIVSUCHE ohne sowie mit Pseudocounts, und bewerten Sie die Lösung, die jeder dieser beiden Algorithmen für das “Subtile Motif Problem” findet, indem Sie entsprechend  $\text{SCORE}(Motive)$  berechnen.

*Eingabe:* Dna= Algo04\_subtiles\_motiv.txt,  $k=15$

*erwartete Ausgabe:* das Motive mit dem Consensus des implantierten (15,4)-Motiv AAAAAAAAGGGGGGGG, bzw. Motive deren  $\text{SCORE}(Motive)$  diesem implantierten Muster gleichkommt

3. Implementieren Sie ebenfalls den Algorithmus RANDOMISIERTEMOTIVSUCHE und wenden Sie diesen dann auf das Subtile Motiv Problem an. Wie groß ist das beste, und wie groß der Median von  $\text{SCORE}(Motive)$  nach 20, 200 bzw. nach 2000 Durchläufen?

*Eingabe/erwartete Ausgabe:* s. Aufgabe 2