

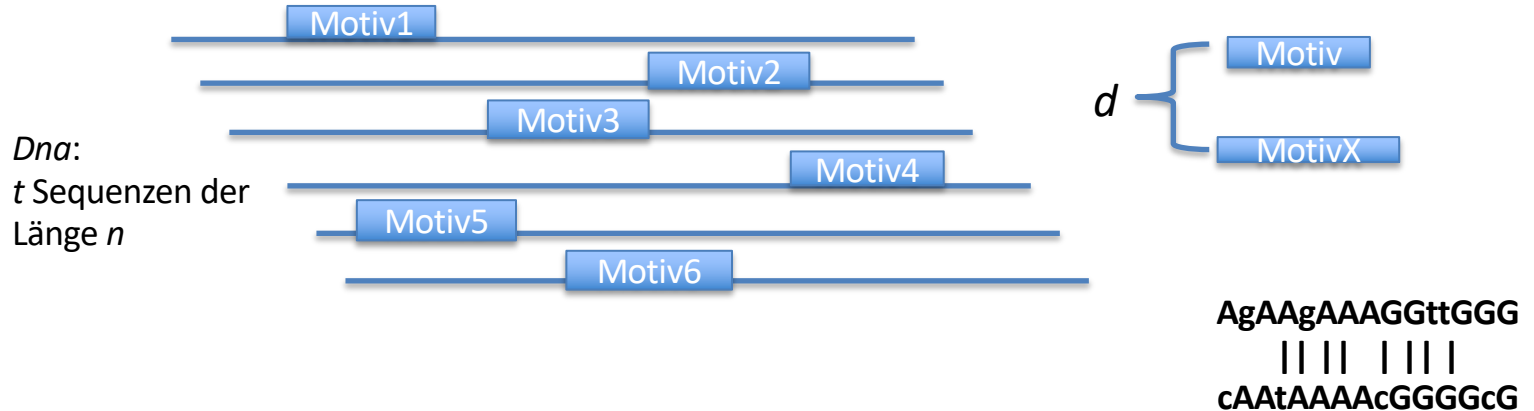
# Algorithmen der Sequenzanalyse: Motivanalysen

AlgSeq – 04/11/2024

Prof. M. Sammeth

# Implantierte Motive

**$(k,d)$ -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

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.) →  
zwei Sprünge im Vergleich:  $Muster \rightarrow Muster' \rightarrow ZÄHLEN_2(Muster')$

# $d$ -Nachbarschaften

Für  $d=1$  kann die Nachbarschaft (1-Nachbarschaft) mit einer Doppelschleife erzeugt werden:

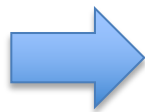
---

**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
```

---



Wie kann uns dieser Algorithmus helfen,  
um  $d$ -Nachbarschaften für  $d>1$  zu generieren?

# Iterative Lösung

---

**Algorithm:** NACHBARNITERATIV( $Muster, Muster, d$ )

---

$Nachbarschaft \leftarrow \{Muster\}$   
**for**  $j \leftarrow 1$  bis  $d$  **do**  
    **for** jede Zeichenkette  $Muster'$  in  $Nachbarschaft$  **do**  
         $Nachbarschaft \leftarrow$   
             $Nachbarschaft \cup \text{UNMITTELBARENACHBARN}(Muster')$   
**return**  $Nachbarschaft$

---

---

**Algorithm:** UNMITTELBARENACHBARN( $Muster$ )

---

$Nachbarschaft \leftarrow \{Muster\}$   
**for**  $i = 1$  to  $|Muster'|$  **do**  
     $Symbol \leftarrow i\text{-tes Nukleotid von } Muster$   
    **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**  $Nachbarschaft$

---

# Rekursiver Ansatz

---

**Algorithm:** NACHBARN( $Muster, Muster, d$ )

---

```
if  $d = 0$  then
  └ return  $Muster$ 
if  $|Muster| = 1$  then
  └ return  $\{A, C, G, T\}$ 
Nachbarschaft  $\leftarrow \emptyset$ 
SuffixNachbarschaft  $\leftarrow$  NACHBARN( $Suffix(Muster), d$ )
for jede Zeichenkette  $Text$  aus  $SuffixNachbarschaft$  do
  if HAMMINGDISTANZ( $SUFFIX(Muster, d)$ )  $< d$  then
    for jedes Nukleotid  $x$  do
      └  $Nachbarschaft \leftarrow Nachbarschaft \cup x \cdot Text$ 
    else
      └  $Nachbarschaft \leftarrow$ 
        └  $Nachbarschaft \cup ERSTESYMBOL(Muster) \cdot Text$ 
   $Nachbarschaft \leftarrow$ 
     $Nachbarschaft \cup UNMITTELBARENACHBARN(Muster')$ 
return  $Nachbarschaft$ 
```

---

# Subtiles Motiv Problem

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

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

```
1 atgaccgggataactgatAgAAgAAAGGttGGGggcgtacacattagataaacgtatgaagtacgttagactcggcgccgcgcg
2 acccctatTTTTTtgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTTccgaataccAAtAAAAcGGcGGGa
3 tgagtatccctgggatgacttAAAAtAAtGGAgtGGtgctctcccgatTTTTgaatatgtaggatcattcgccaggggtccga
4 gctgagaattggatgcAAAAAAGGGattGtccacgcaatcgcgaaaccaacgcggacccaaaggcaagaccgataaaggaga
5 tccctTTTtgcggtaatgtgccgggaggctggttacgtagggaagccctaacggacttaatAtAAtAAAGGaaGGGcttatag
6 gtcaatcatgttcttgtgaatggatttAAcAAtAAGGGctGGgaccgcttggcgcacccaaattcagtgtgggcgagcgcaa
7 cggTTTTggcccttgttagaggcccccgtAtAAAcAAGGaGGGccaattatgagagagctaatactatcgcggtgcgtgttcat
8 aacttgagttAAAAAAtAGGGaGccctggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgta
9 ttggcccatTggctaaaagcccaacttgacaaatggaagatagaatccttgcatActAAAAAGGaGcGGaccgaaaggggaag
10 ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttActAAAAAGGaGcGGa
```

# 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.) →  
zwei Sprünge im Vergleich:  $Muster \rightarrow Muster' \rightarrow ZÄHLEN_2(Muster')$

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



# Motive bewerten

*Motive* := eine  $(t \times k)$  Matrix von Symbolen der  $t$  als Motive betrachteten  $k$ -mere.

*Motive*

(hier Bindestellen NF-κB)

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

$\text{SCORE}(\text{Motive}) :=$  Anzahl der nicht-konservierten Symbole einer Spalte, d.h. Anzahl der Symbole, die nicht dem häufigsten Symbol der Spalte entsprechen.

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

# Motive bewerten (2)

$\text{ZÄHLEN}(\text{Motive}) := (4 \times k)$  Matrix der absoluten Häufigkeiten jedes DNA-Symboles.

$\text{ZÄHLEN}(\text{Motive})$	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	10	10	9	9	1	0	0	0	0	0
	T:	7	2	0	0	1	1	0	5	8	7	3	4

$\text{PROFIL}(\text{Motive}): (4 \times k)$  Matrix der relativen Häufigkeiten jedes DNA-Symboles.

$\text{PROFIL}(\text{Motive})$	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

$\text{CONSENSUSSTRING}(\text{Motive}) :=$  Konkatenation des häufigsten Symboles jeder Spalte  
(zufällige Entscheidung falls mehrere zusammen am häufigsten)

$\text{CONSENSUSSTRING}(\text{Motive})$	<b>T</b>	<b>C</b>	<b>G</b>	<b>G</b>	<b>G</b>	<b>G</b>	<b>A</b>	<b>T</b>	<b>T</b>	<b>T</b>	<b>C</b>	<b>C</b>
-----------------------------------------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

# Das Problem mit dem Konsensus

NF- $\kappa$ B (Transkriptionsfaktor Bindestelle in *D.melanogaster*):

1	2	3	4	5	6	7	8	9	10	11	12
T	C	G	G	G	G	A	T/C	T	T	C	C/T

CSRE (Transkriptionsfaktor Bindestelle in *S.cerevisiae*):

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
C	G/C	G/T	T/A	C/T	G/C	C/G	A	T	G/T	C/G	A	T	C/T	C/T	G/T



Konsensus (Zeichenkette) kann Mehrdeutigkeiten nicht korrekt abbilden.

# Entropie

Die Entropie  $H$ :

ist ein kondensiertes Maß der Unsicherheit über eine Wahrscheinlichkeitsverteilung

$$H(p_1, \dots, p_k) = -\sum_{i=1}^k p_i \log_2(p_i), \text{ mit } \log_2(0) := 0$$

z.B. in PROFIL(Motive) von NF-κB:

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

zweite Spalte:  $H(0.2, 0.6, 0.0, 0.2) \sim 1.371$

konserviertere letzte Spalte:  $H(0.0, 0.6, 0.0, 0.4) \sim 0.971$

5. Spalte – stark konserviert:  $H(0.0, 0.0, 0.9, 0.1) \sim 0.467$



je *größer* die Konservierung, desto *kleiner* die Entropie

# Motiv Logos

Informationsgehalt (IC information content)  $:= 2 - H(p_1, \dots, p_N)$



je *größer* die Konservierung, desto *größer* der Informationsgehalt!

# Das Motiv Findungs Problem

---

## Motiv-Findungs Problem:

Gegeben eine Menge von Zeichenketten, finde *Motive*, eine Menge von k-meren mit je einem k-mer von jeder Zeichenkette, welche das  $\text{SCORE}(\text{Motive})$  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.

---

```
1 atgaccgggatactgatAgAAAgAAGGttGGGggcgtacacattagataaacgtatgaagtacgttagactcggcgccgcgcg
2 acccctatTTTTTgagcagatttagtgacctggaaaaaaaaatttgagtacaaaactTTTTccgaatacAAtAAAAcGGcGGGa
3 tgagtatccctgggatgacttAAAAtAAtGGaGtGGtgctctcccgattTTTTgaatatgtaggatcattcgccaggggtccga
4 gctgagaattggatgAAAAAAAAGGGattGtccacgcaatcgcgaaaccaacgcggacccaaaggcaagaccgataaaggaga
5 tccctTTTTgcggtaatgtgccgggaggctggttacgtagggaagccctaacggacttaatAtAAtAAAGGaaGGGcttatag
6 gtcaatcatgttcttgtgaatggatttAAcAAtAAGGGctGGgaccgcttggcgcacccaaattcagtggtgggcgagcgcaa
7 cggTTTTggcccttgtagaggcccccgtAtAAAcAAGGaGGGccaattatgagagagctaattctatcgcggtgcgtgttcat
8 aacttgagttAAAAAAAtAGGGaGccctggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgta
9 ttggcccatgggctaaaagcccaacttgacaaatggaagatagaatccttgcatActAAAAAGGaGcGGaccgaaaggggaag
10 ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttActAAAAAGGaGcGGa
```

**$(k,d)$ -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  $Dna$ .

---

**Motiv-Findungs Problem:**

Gegeben eine Menge von Zeichenketten, finde *Motive*, eine Menge von  $k$ -meren mit je einem  $k$ -mer von jeder Zeichenkette, welche das  $\text{SCORE}(\text{Motive})$  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.

---

# Motiv-Suche “Brute Force”



“Brute Force” Methode:

BRUTEFORCEMOTIVSUCHE:

- für alle möglichen  $k$ -mer Kompositionen *Motive* in *Dna*
- berechne  $\text{SCORE}(\text{Motive})$
- bestimme das Minimum aller  $\text{SCORE}(\text{Motive})$

$t$  Sequenzen in Dna

jede Sequenz mit  $n$  Nukleotiden

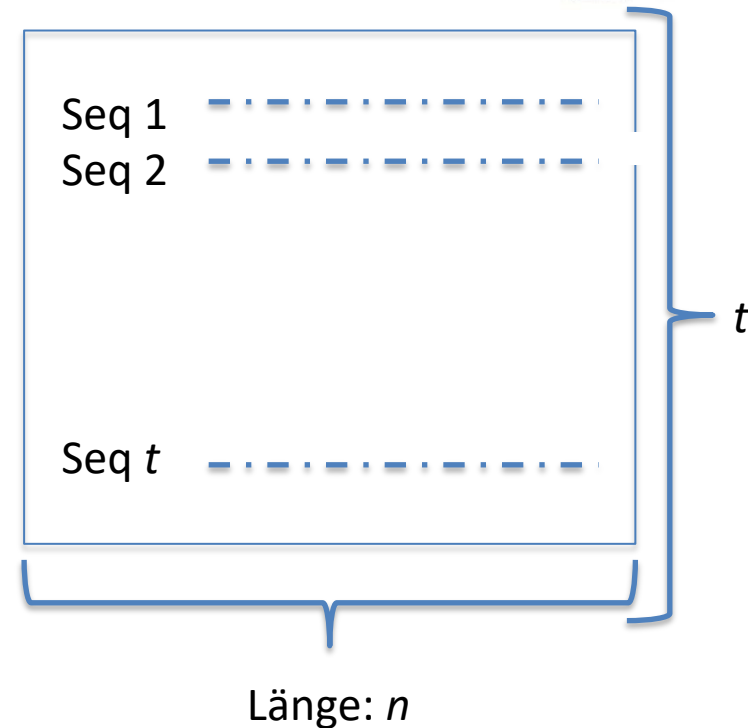
→  $(n - k + 1)$   $k$ -mere in einer Sequenz

→  $(n - k + 1)^t$  Motiv-Matrizen

→  $k \cdot t$  Schritte für  $\text{SCORE}(\text{Motive})$

→ insgesamt:

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





# Übungen

- (1) Beschreiben Sie für eine  $t \times k$  Matrix *Motive* den minimalen und den maximalen Wert von  $\text{SCORE}(\text{Motive})$ .
- (2) Alternativ bzw. analog zur Berechnung von  $\text{SCORE}(\text{Motive})$  können auch die Entropie-Werte der einzelnen Spalten von *Motive* zu einer Gesamt-Entropie der Matrix aufaddiert werden. Bestimmen Sie nach dieser Regel die Entropie der Matrix des NF-κB Motives aus dem Seminar.
- (3) Was ist der *maximale* und was der *minimale* Wert für den *Informationsgehalt*, und *bei welchen Eingaben* tritt jeder dieser Extremwerte auf?
- (4) Implementieren Sie  $\text{BRUTEFORCEMOTIVSUCHE}(\text{Dna}, k)$  und versuchen Sie damit, die in den Folien verwendete “Miniatur” des subtilen Motiv Problemes mit  $n=82$  zu lösen.

Eingabe: Datei `Algo05_subtiles_motiv_mini.txt`

erwartete Ausgabe: das 15-mer `AAAAAAAAAGGGGGGGG`