

# Algorithmen der Sequenzanalyse: Motivanalysen, Sequenzalignment

AlgSeq – 18/11/2024

Prof. M. Sammeth

# Motivanalysen

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**

$BesteMotive \leftarrow Motive$

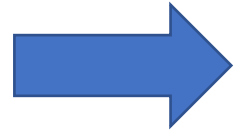
**else**

**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**

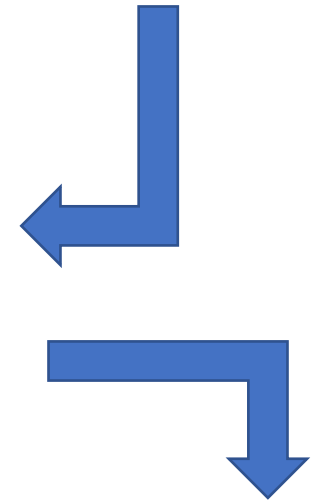


PROFIL(Motive)=

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

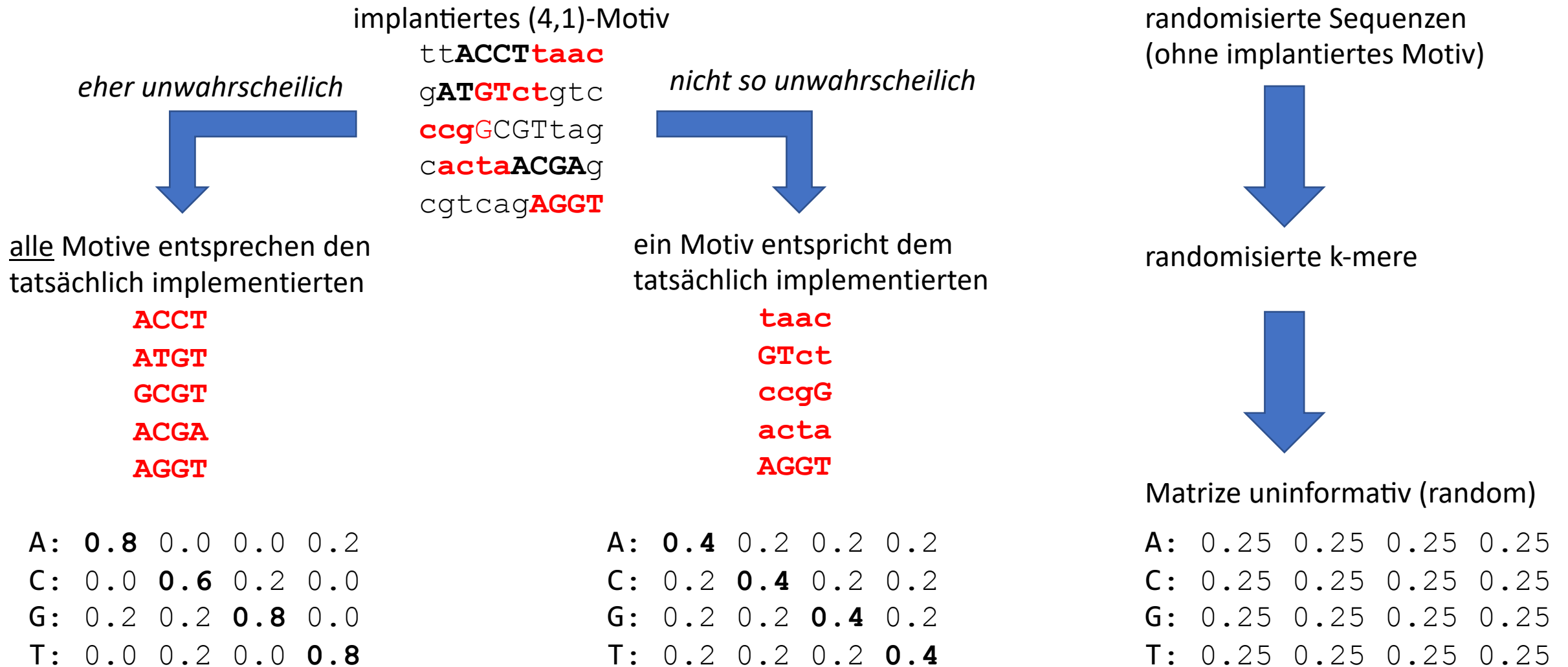
ttAC	tACC	ACCT	CCTt	CTta	Ttaa	taac
.0016	.0016	<b>.0128</b>	.0064	.0016	.0016	.0016
gATG	ATGT	TGTc	GTct	Tctg	ctgt	tgtc
.0016	<b>.0128</b>	.0016	.0032	.0032	.0032	.0016
ccgG	cgGC	gGCG	GCGT	CGTt	GTta	Ttag
.0064	.0032	.0016	<b>.0128</b>	.0032	.0016	.0016
cact	acta	ctaA	taAC	aACG	ACGA	CGAg
.0032	.0064	.0016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

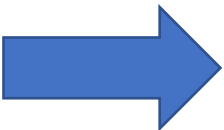
# Evaluierung der Randomisierten Motiv-Suche



# Gibbs Sampling

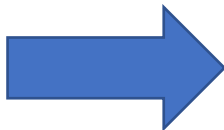
- beginnt auch mit zufällig gewählten  $k$ -meren,  
ist aber vorsichtiger als Randomisierte Motiv-Suche:

→ in jeder “Runde” wird lediglich ein  $k$ -mer betrachtet und entschieden, ob es behalten oder ersetzt werden soll.

ttacctt <b>aac</b>		ttacctt <b>aac</b>
g <b>ata</b> tctgtc		gat <b>atc</b> tgtc
<b>acg</b> gcgttcg		acggcg <b>ttc</b> g
ccct <b>aaa</b> gag		ccctaa <b>aga</b> g
cgtc <b>aga</b> ggt		<b>cgt</b> cagaggt

RANDOMISIERTE MOTIVSUCHE

(kann alle  $k$ -mere in einer Runde austauschen)

ttacctt <b>aac</b>		ttacctt <b>aac</b>
g <b>ata</b> tctgtc		gatatc <b>tgtc</b>
<b>acg</b> gcgttcg		<b>acg</b> gcgttcg
ccct <b>aaa</b> gag		ccct <b>aaa</b> gag
cgtc <b>aga</b> ggt		cgtc <b>aga</b> ggt

GIBBSAMPLER

(tauscht höchstens ein  $k$ -mer pro Runde)

---

**Algorithm:** GIBBSAMPLER( $Dna, k, t, N$ )

---

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

$BesteMotive \leftarrow Motive$

**for**  $j \leftarrow 1$  bis  $N$  **do**

$i \leftarrow \text{RANDOM}(t)$

$Profile \leftarrow$  Profil-Matrix berechnet aus allen  $k$ -meren in  $Motive$  außer  $Motiv_i$

$Motiv_i \leftarrow$  Profil-zufällig gewähltes  $k$ -mer in der  $i$ -ten Sequenz von  $Dna$

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

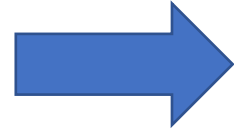
$BesteMotive \leftarrow Motive$

---



implantiertes (4,1)-Motiv

tt**ACCT**taac  
g**ATGTct**gtc  
ccg**GCGT**tag  
cacta**ACGA**g  
cgtcag**AGGT**



zufällig gewählte Motive

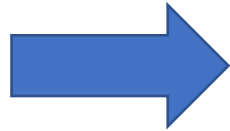
ttACCT**taac**  
gAT**GTct**gtc  
**ccgG**CGTtag  
c**acta**ACGA  
cgtcag**AGGT**



zufällig 3. Sequenz gewählt

tt**ACCTtaac**  
g**ATGTct**gtc  
-----  
c**actaACGA**g  
cgtcag**AGGT**

Motive **taac**  
**GTct**  
**acta**  
**AGGT**



ZAEHLEN(*Motive*)=

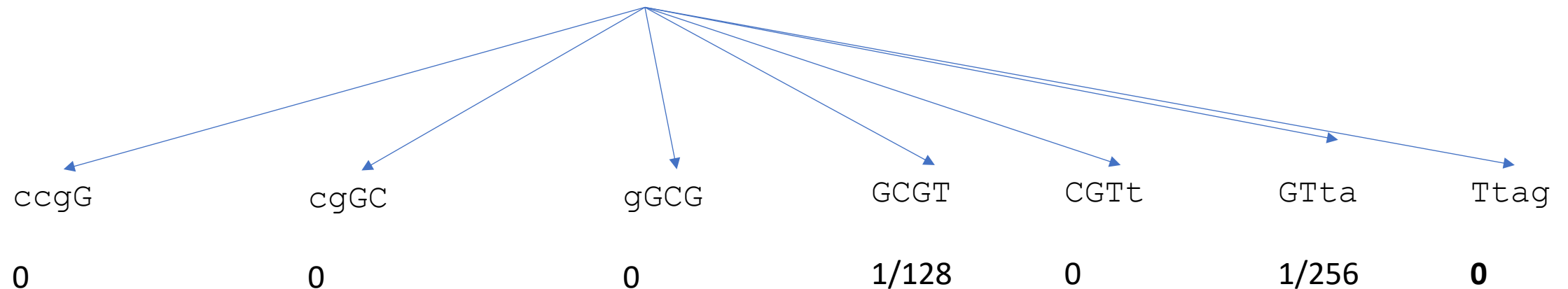
A:	2	1	1	1
C:	0	1	1	1
G:	1	1	1	0
T:	1	1	1	2



PROFIL(*Motive*)=

A:	2/4	1/4	1/4	1/4
C:	0	1/4	1/4	1/4
G:	1/4	1/4	1/4	0
T:	1/4	1/4	1/4	2/4

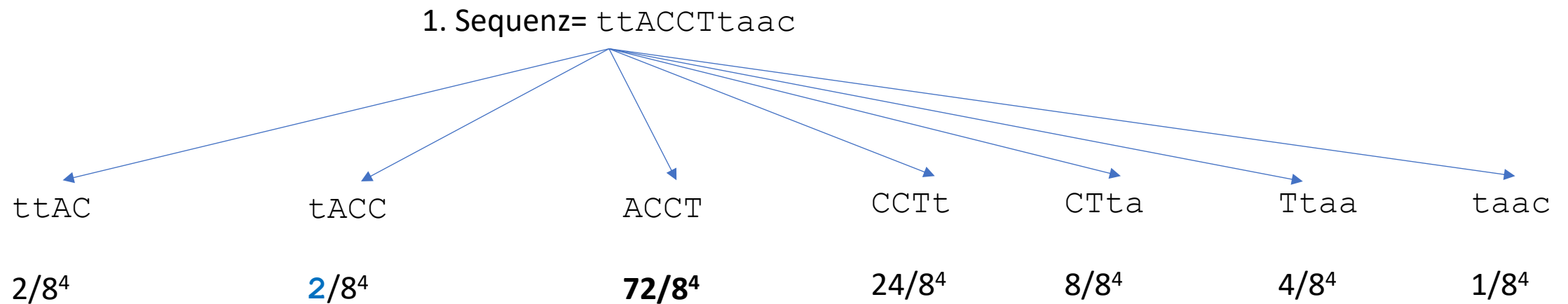
3. Sequenz= ccgGCGTtag



→ kein (fares) gewichtetes Subsampling möglich wo  $Pr() = 0$  !

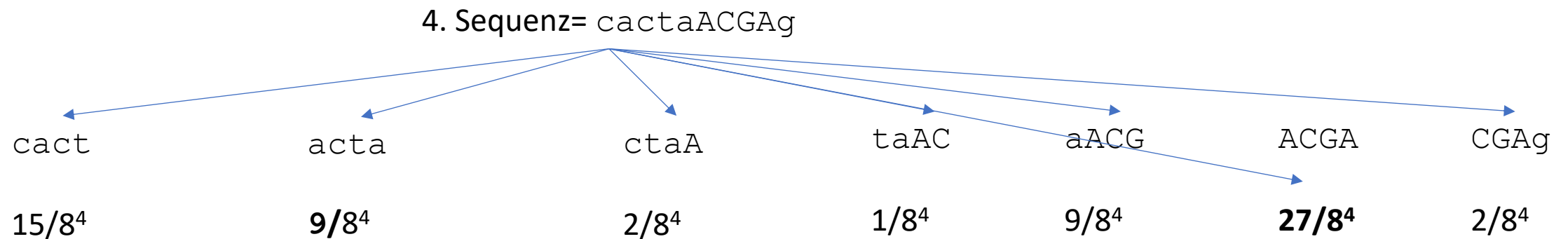


Annahme: **GCGT** wurde gewichtet zufällig in der 3. Sequenz gewählt

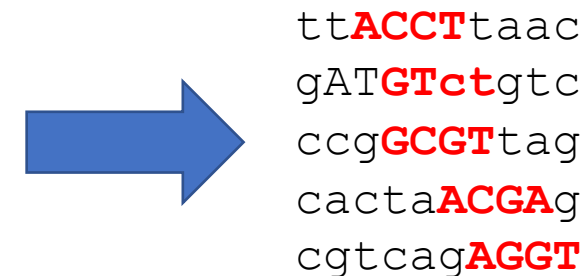


Annahme: **ACCT** wurde gewichtet zufällig in der 1. Sequenz gewählt

Annahme: **ACCT** wurde gewichtet zufällig in der 1. Sequenz gewählt

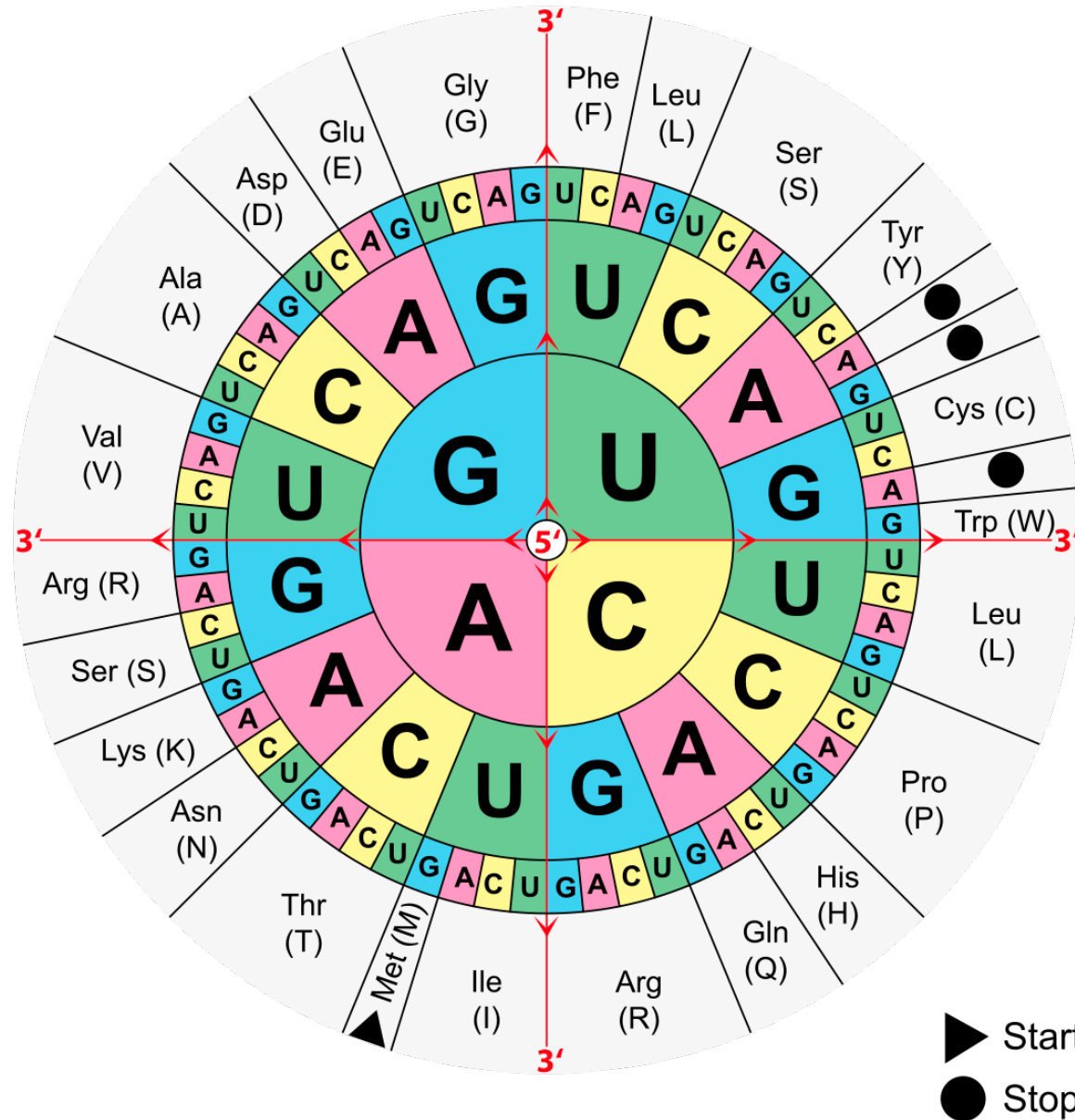


Annahme: **ACGA** wurde gewichtet zufällig in der 4. Sequenz gewählt



# Sequenz- Alignments

# Knacken des Genetischen Codes



# 1961 – Entdeckung der Triplet-Natur von Codons

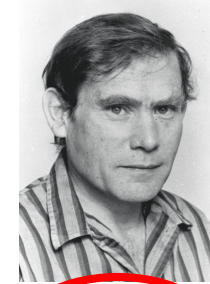
Francis Crick



Alexander Rich

James Watson

Leslie Orgel



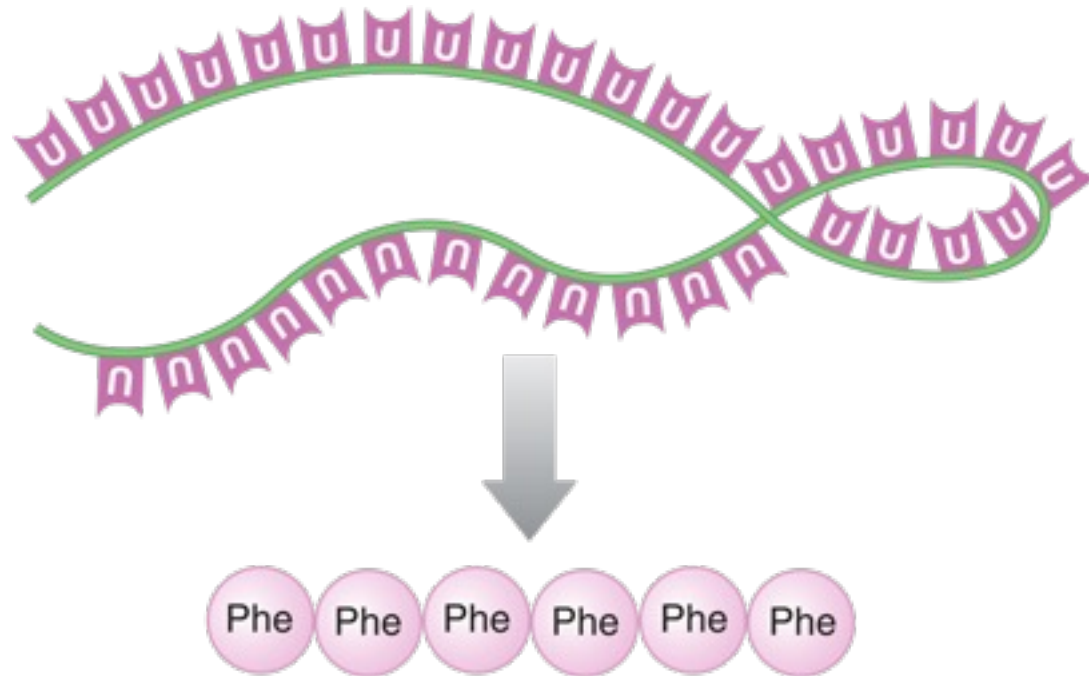
Sidney Brenner

	MUTATION	PHENOTYPE
Wild-type sequence	NONE	rII <sup>+</sup>
FC0 mutant	+	rII <sup>-</sup>
Supression of FC0	+ -	rII <sup>+</sup>
Two base additions	+ +	rII <sup>-</sup>
Three base additions	+ + +	rII <sup>+</sup>

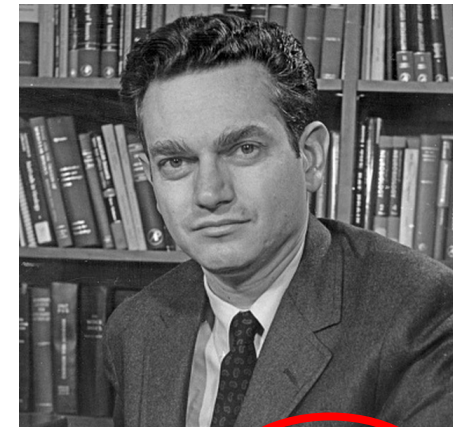
+ Base addition  
- Base deletion

Experimente mit Phagen, Insertion von 3 Mutationen restauriert den Wildtyp!

# 1961 – erstes Codon UUU



*in vitro* Translationssystem für poly-*U* RNAs



Marshall Nirenberg





# 1960s – weitere Experimente mit Homopolymeren



Har Gobind

z.B. UCUCUCUCUCUC...

→ LeuSerLeuSer



Robert W. Holley      Har Gobind Khorana      Marshall W. Nirenberg

- The Nobel Prize in Physiology or Medicine 1968 was awarded jointly to Robert W. Holley, Har Gobind Khorana and Marshall W. Nirenberg "for their **interpretation of the genetic code** and its **function in protein synthesis**" in 1961.

# Der *Nicht*-Ribosomale Code

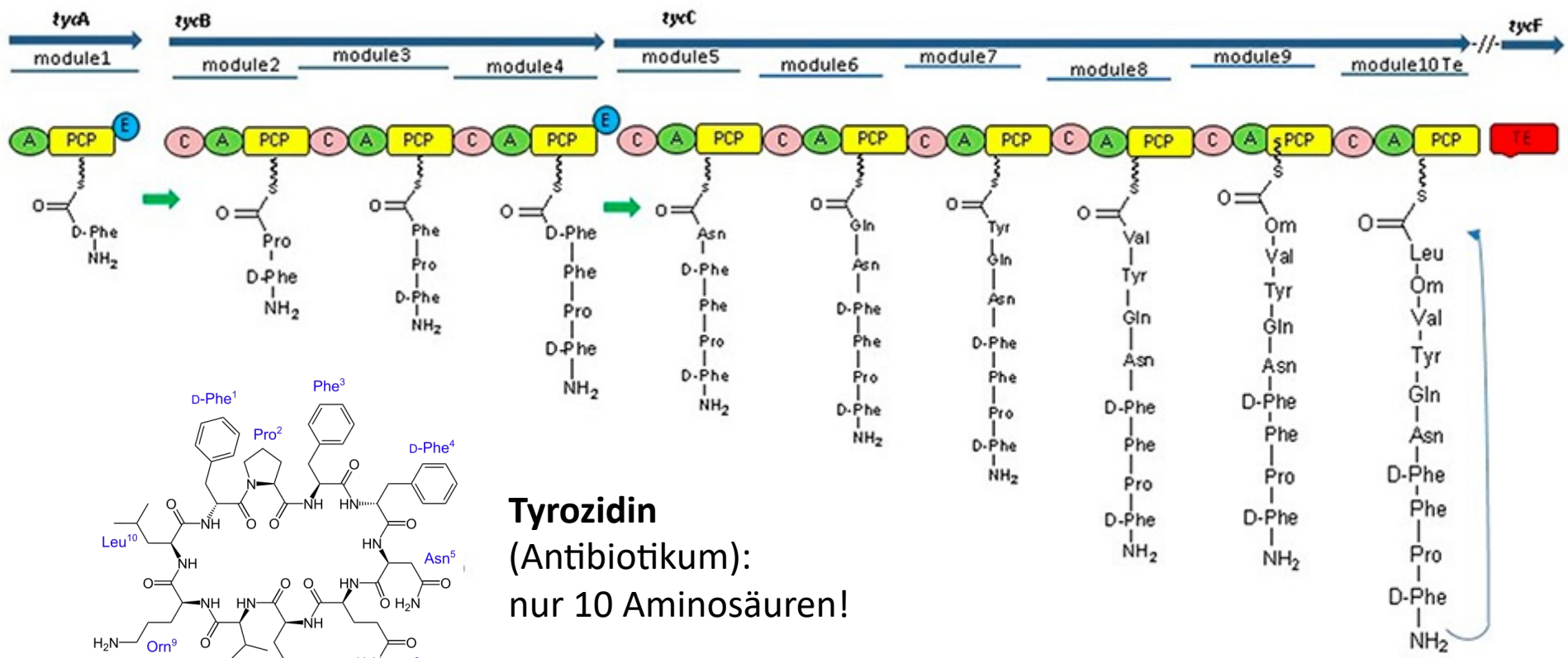
Bakterien und Pilze kodieren Antibiotika und andere *nicht-ribosomale Peptide (NRP)*

Enzym: **NRP Synthetase** (unabhängig von Ribosomen und genet. Code!)

DNA → RNA → **NRP Synthetase** → NRP

NRP Synthetasen bestehen aus Modulen, sog. Adenylierungs-Domänen (A-Domänen),  
wovon jede nur eine Aminosäure synthetisiert.

z.B. **NRP Synthetase für Tyrozin**: 10 A-Domänen mit je 500 ASren → 5.000 AS Protein!



# Knacken des Nicht-Ribosomalen Codes

BIOENGINEERING

## Working Outside the Protein-Synthesis Rules

Proteins built in the ribosome are subject to certain restrictions, so researchers are harnessing a nonribosomal system that might one day make new drugs



Mohamed Marahiel

Sequenzen von drei *A-Domänen* bekannt

**Asp:**

YAFDLGYTCMFPVLLGGGELHIVQKETYTAPDEIAHYIKEHGITYIKLTPSLFHTIVNTASFAFDANFESLRLIVLGGEKIIPIDVIAFRKMYGHTEFINHYGPTEATIGA

**Orn:**

AFDVSAGDFARALLTGGQLIVCPNEVKMDPASLYAIKKYDITIFEATPALVIPLMEYIYEKLDISQLQILIVGSDSCSMEDFKTLVSRFGSTIRIVNSYGVTEACIDS

**Val:**

IAFDASSWEIYAPLLNGGTVCIDYYTTIDIKALEAVFKQHHIRGAMLPPALLKQCLVSAPTMISSLEILFAAGDRLSSQDAILARRAVGSGVYNAYGPTENTVLS

(jeweils nur ca. 110 der ca. 500 Aminosäuren einer A-Domäne gezeigt)



was sind die Gemeinsamkeiten und die (spezifischen) Unterschiede dieser 3 A-Domänen?

# Beobachtungen zum Sequenzvergleich

## 1. Versuch: Direkter Vergleich der Aminosäure → 3 konservierte Spalten

**Asp:**

YAFDLGYTCMFPVLLGGELHIVQKETYTAPDEIAHYIKEHGITYIKLTPSLFHTIVNTASFAFDANFESLRLIVLGGEKIIPIDVIAFRKMYGHTTEFINHYGPTEATIGA

**Orn:**

AFDVSAGDFARALLTGGQLIVCPNEVKMDPASLYAIIKKYDITIFEATPALVIPLMEYIYEQKLDISQLQILIVGSDSCSMEDFKTLVSRFGSTIRIVNSYGVTEACIDS

**Val:**

IAFDASSWEIYAPLLNGGTVVCIDYYTTIDIKALEAVFKQHHIRGAMLPALLKQCLVSAPTMISSLEILFAAGDRLSSQDAILARRAVGSGVYNAYGPTENTVLS

## 2. Versuch: Sequenzähnlichkeiten "mit dem Auge" finden

**Asp:**

YAFDLGYTCMFPVLLGGELHIVQKETYTAPDEIAHYIKEHGITYIKLTPSLFHTIVNTASFAFDANFESLRLIVLGGEKIIPIDVIAFRKMYGHTTEFINHYGPTEATIGA

**Orn:**

AFDVSAGDFARALLTGGQLIVCPNEVKMDPASLYAIIKKYDITIFEATPALVIPLMEYIYEQKLDISQLQILIVGSDSCSMEDFKTLVSRFGSTIRIVNSYGVTEACIDS

**Val:**

IAFDASSWEIYAPLLNGGTVVCIDYYTTIDIKALEAVFKQHHIRGAMLPALLKQCLVSAPTMISSLEILFAAGDRLSSQDAILARRAVGSGVYNAYGPTENTVLS

→ zweite Sequenz "nach rechts" verschieben, durch Einfügen eines Platzhalters "-" am Beginn

**Asp:**

YAFDLGYTCMFPVLLGGELHIVQKETYTAPDEIAHYIKEHGITYIKLTPSLFHTIVNTASFAFDANFESLRLIVLGGEKIIPIDVIAFRKMYGHTTEFINHYGPTEATIGA

**Orn:**

-AFDVSAGDFARALLTGGQLIVCPNEVKMDPASLYAIIKKYDITIFEATPALVIPLMEYIYEQKLDISQLQILIVGSDSCSMEDFKTLVSRFGSTIRIVNSYGVTEACIDS

**Val:**

IAFDASSWEIYAPLLNGGTVVCIDYYTTIDIKALEAVFKQHHIRGAMLPALLKQCLVSAPTMISSLEILFAAGDRLSSQDAILARRAVGSGVYNAYGPTENTVLS

→ 11 konservierte Spalten

# Verschieben durch Einfügen von Platzhaltern (Alignieren)

## 3. weitere fünf (4+1) Platzhalter einfügen

**Asp:**

YAFDLGYTCMFPVLLGGGELHIVQKETYTAPDEIAHYIKEHGIITYIKLTPSLFHTIVNTASFAFDANFESLRLIVLGGEKIIPIDVIAFRKMYGHTEFINHYGPTEATIGA

**Orn:**

-AFDVSAGDFARALLTGGQLIVCPNEVKMDPASLYAIIKKYDITIFEATPALVIPLMEYI-YEQKLDISQLQILIVGSDSCSMEDFKTLVSRFGSTIRIVNSYGVTEACIDS

**Val:**

IAFDASSWEIYAPLLNGGTVCIDYYTTIDIKALEAVFKQHHIRGAMLPALLKQCLVSA-----PTMISSLEILFAAGDRLSSQDAILARRAVGSGVYNAYGPTENTVLS

→ 14 konservierte Spalten

## 4. und noch weitere drei Platzhalter einfügen

**Asp:**

YAFDLGYTCMFPVLLGGGELHIVQKETYTAPDEIAHYIKEHGIITYIKLTPSLFHTIVNTASFAFDANFESLRLIVLGGEKIIPIDVIAFRKMYGHTE-FINHYGPTEATIGA

**Orn:**

-AFDVSAGDFARALLTGGQLIVCPNEVKMDPASLYAIIKKYDITIFEATPALVIPLMEYI-YEQKLDISQLQILIVGSDSCSMEDFKTLVSRFGSTIRIVNSYGVTEACIDS

**Val:**

IAFDASSWEIYAPLLNGGTVCIDYYTTIDIKALEAVFKQHHIRGAMLPALLKQCLVSA-----PTMISSLEILFAAGDRLSSQDAILARRAVGSGV-Y-NAYGPTENTVLS

→ 19 konservierte Spalten (inkl. 2 Duplets)

# Alignments zeigen konservierte/variable Bereiche

Das Alignment zeigt den **konservierten** Kern (**Core**) der A-Domänen:

**Asp:**

YAFDLGYTCMFPVLLGGGELHIVQKETYTAPDEIAHYIKEHGITYIKLTPSLFHTIVNTASFAFDANFESLRLIVLGGEKIIPIDVIAFRKMYGHTE-FINHYGPTEATIGA

**Orn:**

-AFDVSAGDFARALLTGGQLIVCPNEVKMDPASLYAIIKKYDITIFEATPALVIPLMEYI-YEQKLDISQLQILIVGSDSCSMEDFKTLVSRFGSTIRIVNSYGVTEACIDS

**Val:**

IAFDASSWEIYAPLLNGGTVVCIDYYTTIDIKALEAVFKQHHIRGAMLPALLKQCLVSA----PTMISSLEILFAAGDRLSSQDAILARRAVGSGV-Y-NAYGPTENTVLS

→ dadurch konnte Marahiel die mit der kodierten Aminosäure variierende *Signatur* finden:

**Asp:**

YAFDLGYTCMFPVLLGGGELHIVQKETYTAPDEIAHYIKEHGITYIKLTPSLFHTIVNTASFAFDANFESLRLIVLGGEKIIPIDVIAFRKMYGHTE-FINHYGPTEATIGA

**Orn:**

-AFDVSAGDFARALLTGGQLIVCPNEVKMDPASLYAIIKKYDITIFEATPALVIPLMEYI-YEQKLDISQLQILIVGSDSCSMEDFKTLVSRFGSTIRIVNSYGVTEACIDS

**Val:**

IAFDASSWEIYAPLLNGGTVVCIDYYTTIDIKALEAVFKQHHIRGAMLPALLKQCLVSA----PTMISSLEILFAGDRLSSQDAILARRAVGSGV-Y-NAYGPTENTVLS

Signaturen:

LTKVGHIG

→ Asp

VGEIGSID

→ Orn

AWMFAAVL

→ Val

→ das wäre im ursprünglichen Alignment unmöglich gewesen:

**Asp:**

YAFDLGYTCMFPVLLGGGELHIVQKETYTAPDEIAHYIKEHGITYIKLTPSLFHTIVNTASFAFDANFESLRLIVLGGEKIIPIDVIAFRKMYGHTEFINHYGPTEATIGA

**Orn:**

AFDVSAGDFARALLTGGQLIVCPNEVKMDPASLYAIIKKYDITIFEATPALVIPLMEYIYEQKLDISQLQILIVGSDSCSMEDFKTLVSRFGSTIRIVNSYGVTEACIDS

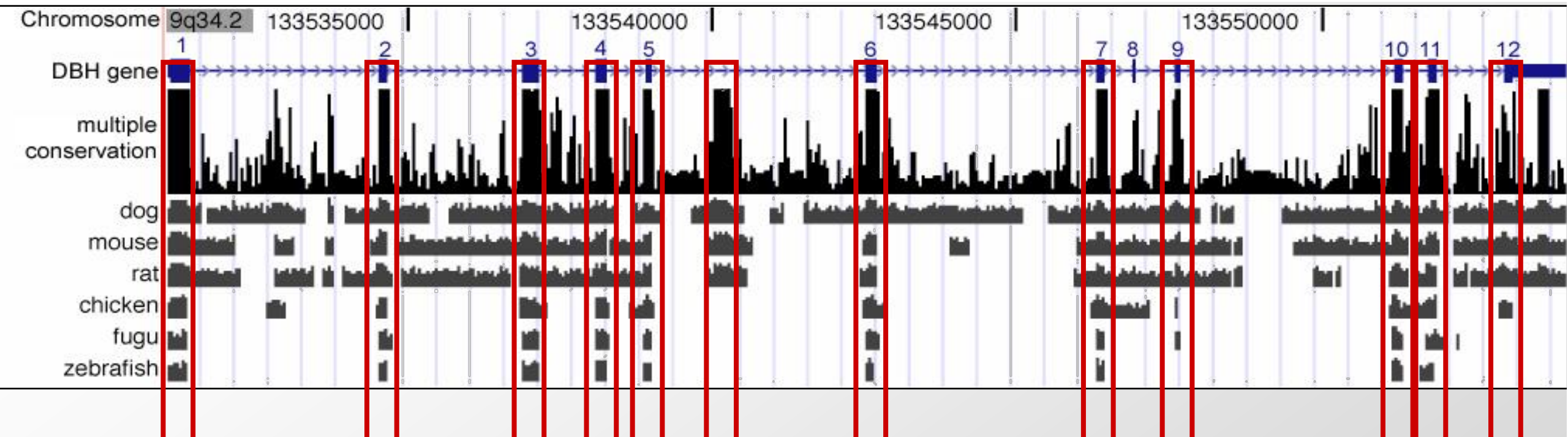
**Val:**

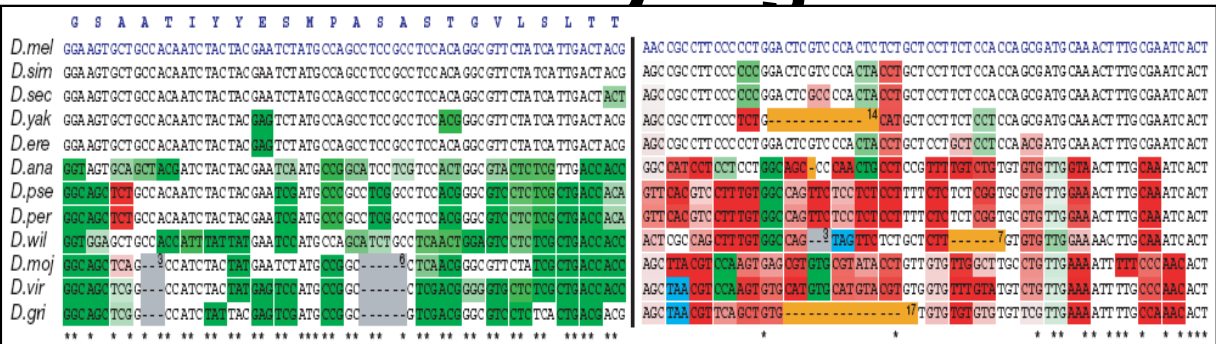
IAFDASSWEIYAPLLNGGTVVCIDYYTTIDIKALEAVFKQHHIRGAMLPALLKQCLVSAPTMISSLEILFAGDRLSSQDAILARRAVGSGVYNAYGPTENTVLS



# Weitere Anwendungen von Alignments:

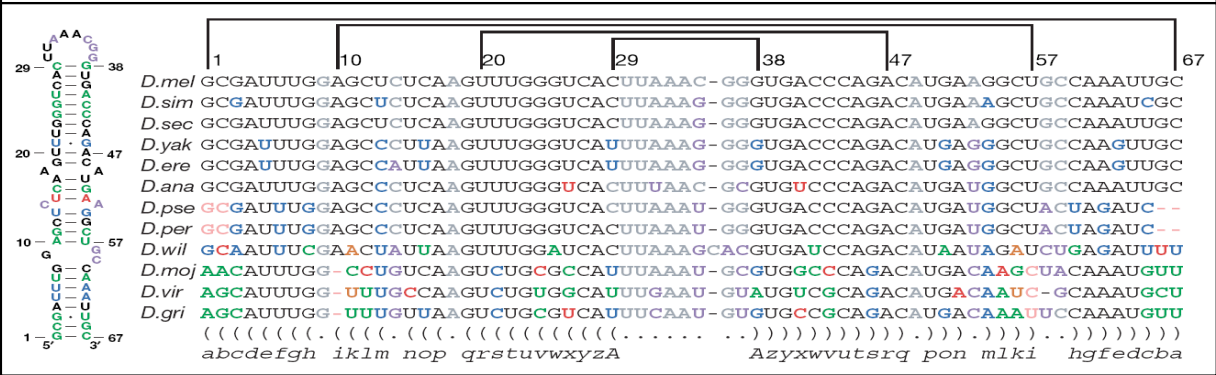
## Konservierung (phylogenet. Kontext)



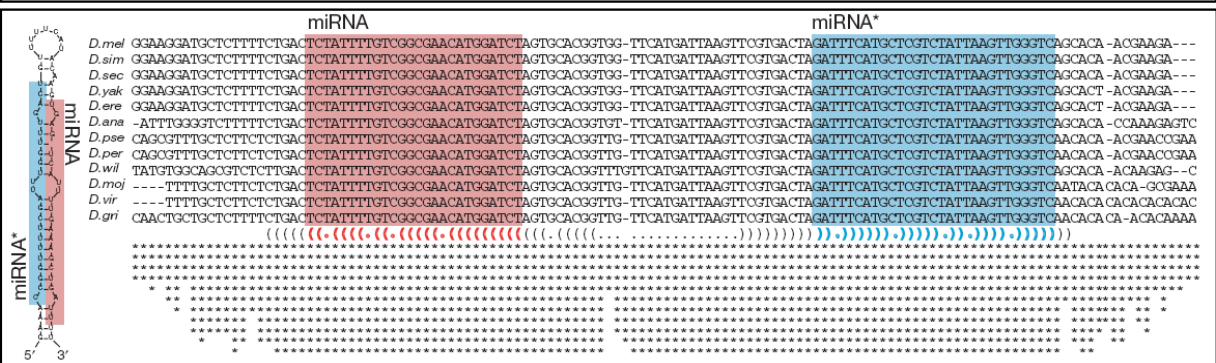


# (Protein-)kodierende Gene

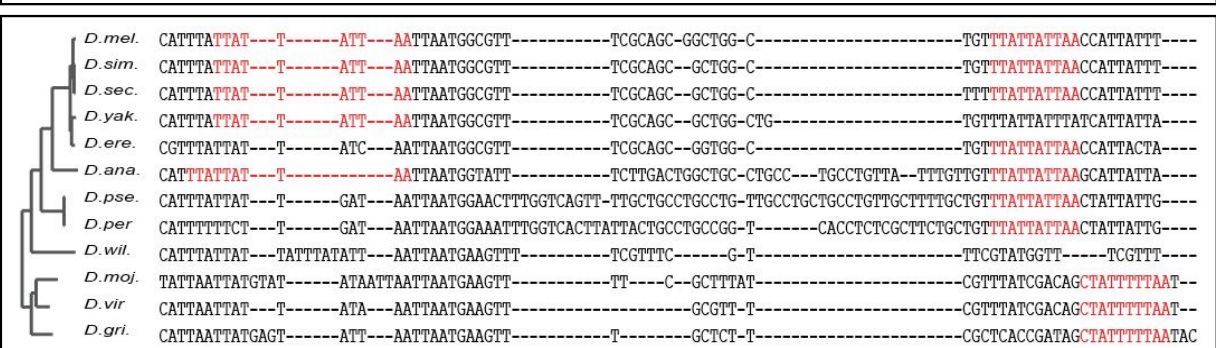
## RNA Strukturen



## miRNAs



## Regulatorische Motive





# Motivation von Alignments

## Vereinfachung: paarweiser Vergleich

$V =$  ATGCATGC

**W=** TGCATGCA

### Beobachtungen:

- nicht alle äquivalenten Sequenzen  $(v,w)$  sind gleich lang
- auch bei gleichlangen Sequenzen  $|v| = |w|$  entspricht das  $i$ -te Symbol des einen String  $v$  manchmal einer (völlig) anderen Position im anderen String  $w$ .

## Alignment:

Einfügen von Platzhalter-Symbolen ("gaps") um eine Matrix zu erhalten,  
in der gleich(-wertig)e Symbole in der gleichen Spalte  $i < |v| + |w|$  stehen.  
→ (Hamming-)Distanz minimiert.

**z.B.**

ATGCATGC → ATGCATGC -  
TGCATGCA - TGCATGCA

ATGCTTA  
TGCATTAA

→

ATGC-TTA-  
-TGCATTAA

HammingDistanz = 8  $\rightarrow$  HammingDistanz = 2

HammingDistanz = 5  $\rightarrow$  HammingDistanz = 3



Zeichen in Strings (verbleibend)	Alignment (wachsend)	P.
A T G T T A T A A T C G T C C		
T G T T A T A T C G T C	A A	+1
G T T A T A C G T C C	A T A T	+1
G T T A T A G T C C	A T - A T C	
T T A T A T C C	A T - G A T C G	+1
T A T A C C	A T - G T A T C G T	+1
A T A C C	A T - G T T A T C G T -	
T A C	A T - G T T A A T C G T - C	
A C	A T - G T T A T A T C G T - C -	
	A T - G T T A T A A T C G T - C - C	

# Alignment als ein Spiel

## Regeln:

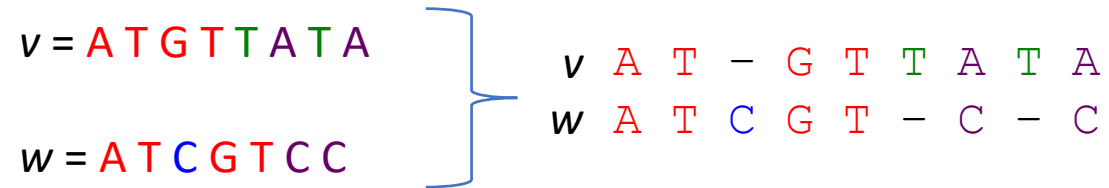
- in jeder Runde nehmen wir das *erste* (linke) Symbol von *einer oder beiden* Strings und fügen es hinten (rechts) ans Alignment
- wenn wir *zwei gleiche Symbole* ans Alignment fügen, bekommen wir +1 Punkt (Greedy!)
- wenn wir nur ein Sybol ans Alignment fügen, wird ein *Platzhalter* "-" am anderen String angefügt

➔ Das Alignieren von ATGTTATA und ATCGTCC gibt +4 Punkte insgesamt!



# Und wie gewinnt man das Spiel?

- Nicht-Ribosomaler Code: "gutes" Alignment zeigt möglichst viele Spalten mit gleichen Symbolen (**Matches**);
- die Matches im Alignment definieren eine **Sequenz von Symbolen**, die so in der gleichen Reihenfolge in beiden Strings  $v$  und  $w$  vorkommen: das ist die "**gemeinsame Subsequenz**";



- ein "gutes" Alignment findet eine lange bzw. die **längste gemeinsame Subsequenz** der alignierten Strings.

---

**Längste Gemeinsame Subsequenz Problem:** *Finde eine längste gemeinsame Subsequenz von zwei Strings.*

**Input:** Zwei Strings,  $v$  und  $w$ .

**Output:** Eine gemeinsame Subsequenz von  $v$  und  $w$  mit maximaler Länge.

---

# Wie viele Alignments gibt es?

$v = \text{ATGTTATA}$   
 $w = \text{ATCGTCC}$

*Abschätzung* nach Alignment-Spiel:

- höchstens 3 Möglichkeiten in jeder Runde
- höchstens  $|v| + |w|$  Runden

$\in \mathcal{O}(3^{|v| + |w|})$  oder  $\mathcal{O}(3^{2L})$  mit  $L = |v| + |w|$

*tatsächlich:*

$$\frac{(|v| + |w|)!}{|v|!|w|!} \quad \text{bzw.} \quad \frac{2L!}{L!^2} \quad \text{mit } L = |v| + |w|$$

$$\sim \frac{2^{2L}}{\sqrt{\pi L}} \quad \text{mit } n! \sim \sqrt{2\pi n} \left(\frac{n}{e}\right)^n \quad \text{Stirlingformel:}$$



exponentiell viele, mögliche Alignments!  
*Brute Force* nur mit (sehr) kurzen Sequenzen möglich



# Übungen

1. Implementieren Sie 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 SCORE(Motiv) nach 20, 200 bzw. nach 2000 Durchläufen?

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

2. Vergleichen Sie die Ergebnisse aus (1) mit entsprechenden Werten nach 20, 200 und 200 Läufen von GIBBSAMPLER, mit je  $N=200$  Iterationen.

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

3. Implementieren Sie einen Greedy-Algorithmus nach dem im Seminar besprochenen Alignment-“Spiel”, mit 1 Punkt für Matches und 0 Punkten für Mismatches bzw. Gaps.

Eingabe: ATGCATGC , TGCATGCA

4. Vergleichen Sie die Ergebnisse aus (3) zu denen einer exhaustiven Suche nach dem besten Alignment zweier kurzer Sequenzen.

Eingabe: ATGCATGC, TGCATGCA

Ausgabe: ATGC-TTA-, -TGCATTAA