

# *E2 - Proteine*

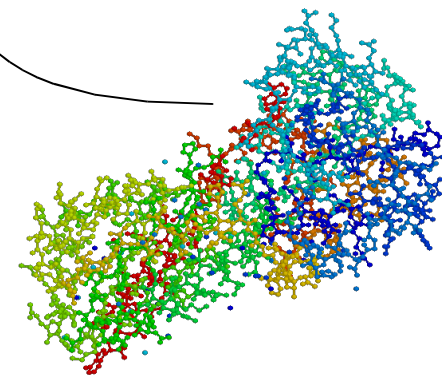
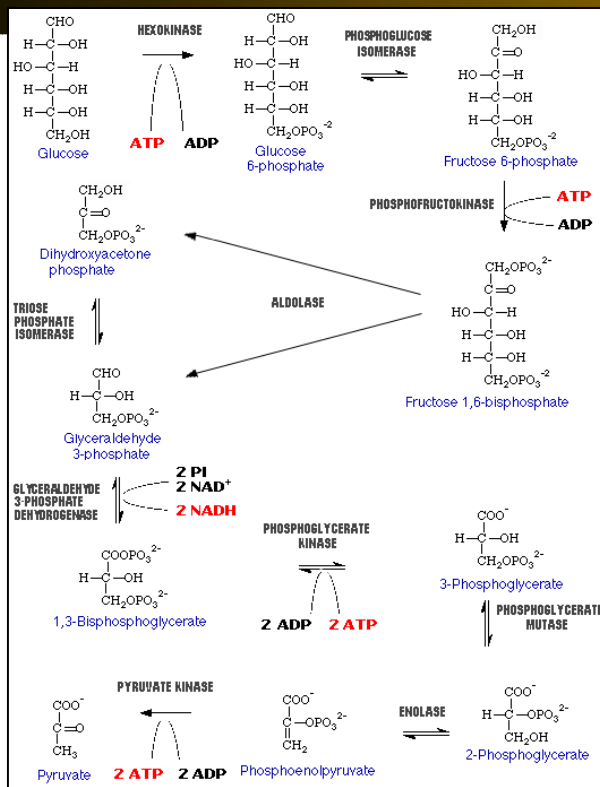


1. Tag: Proteinsequenzen und Sequenzvergleich

Ursula Kummer, Sven Sahle

Femke Mensonides, Irina Surovtsova, Jürgen Zobeley

# Die Zelle - vom Gen zum Enzym



MFKPVDFSETSPVPPDIDLAPTQSPHHVAPSQDSSYDLLS.....

..

.....  
SMLKNKSFLLHGKDYPNNADNNDNEDIRAKTMNRSQSHV

gatccagctg taccattatg taatataata agacacggac gcac.....

# Genetischer Code

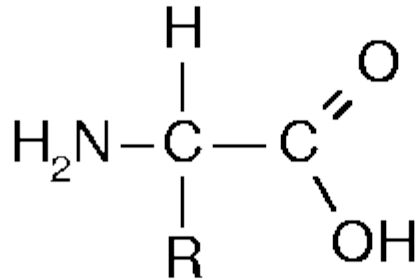
- Code: Abbildung  $A \rightarrow A'$  , wobei  $A$  und  $A'$  jeweils eine Menge von Symbolen sind
- Genetischer Code: Abbildung von Basen-Tripletts (Codone) auf Aminosäuren  
$$\{XYZ \mid X,Y,Z \in \{A,C,G,U\}\}$$
  
 $\rightarrow \{\text{Ala, Cys, Asp, ..., Tyr}\}$
- Eigenschaften des Genetischen Codes:
  - redundant (64 Codone werden auf 20 Aminosäuren und drei STOP-Codone abgebildet)
  - wird bei allen bekannten Lebewesen verwendet ( mit leichten Abweichungen z.B. in Mitochondrien und in Prokaryonten, bisher 15 Code Tabellen bekannt)

# Genetischer Code

| P1 | Position 2 |      |     |     | P3 |
|----|------------|------|-----|-----|----|
|    | G          | A    | C   | U   |    |
| G  | Gly        | Glu  | Ala | Val | G  |
|    | Gly        | Glu  | Ala | Val | A  |
|    | Gly        | Asp  | Ala | Val | C  |
|    | Gly        | Asp  | Ala | Val | U  |
| A  | Arg        | Lys  | Thr | Met | G  |
|    | Arg        | Lys  | Thr | Ile | A  |
|    | Ser        | Asn  | Thr | Ile | C  |
|    | Ser        | Asn  | Thr | Ile | U  |
| C  | Arg        | Gln  | Pro | Leu | G  |
|    | Arg        | Gln  | Pro | Leu | A  |
|    | Arg        | His  | Pro | Leu | C  |
|    | Arg        | His  | Pro | Leu | U  |
| U  | Trp        | STOP | Ser | Leu | G  |
|    | STOP       | STOP | Ser | Leu | A  |
|    | Cys        | Tyr  | Ser | Phe | C  |
|    | Cys        | Tyr  | Ser | Phe | U  |

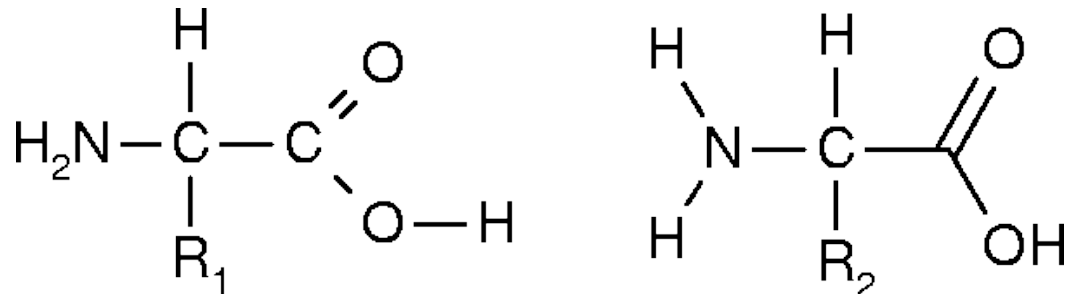
|   |     |              |
|---|-----|--------------|
| A | Ala | Alanin       |
| C | Cys | Cystein      |
| D | Asp | Aspartat     |
| E | Glu | Glutatmat    |
| F | Phe | Phenylalanin |
| G | Gly | Glycin       |
| H | His | Histidin     |
| I | Ile | Isoleucin    |
| K | Lys | Lysin        |
| L | Leu | Leucin       |
| M | Met | Methionin    |
| N | Asn | Asparagin    |
| P | Pro | Prolin       |
| Q | Gln | Glutamin     |
| R | Arg | Arginin      |
| S | Ser | Serin        |
| T | Thr | Threonin     |
| V | Val | Valin        |
| W | Trp | Tryptophan   |
| Y | Tyr | Tyrosin      |

# Aminosäuren: Aufbau

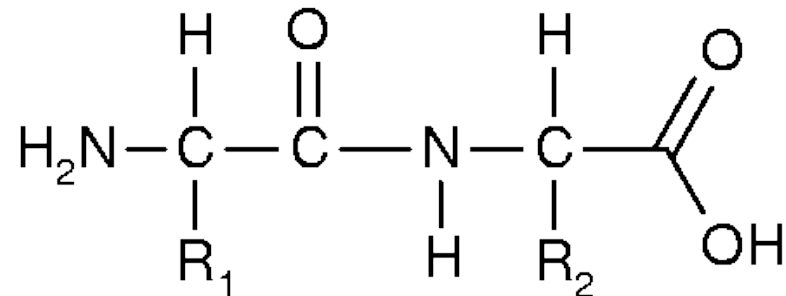


R: Rest

Struktur einer  
Aminosäure

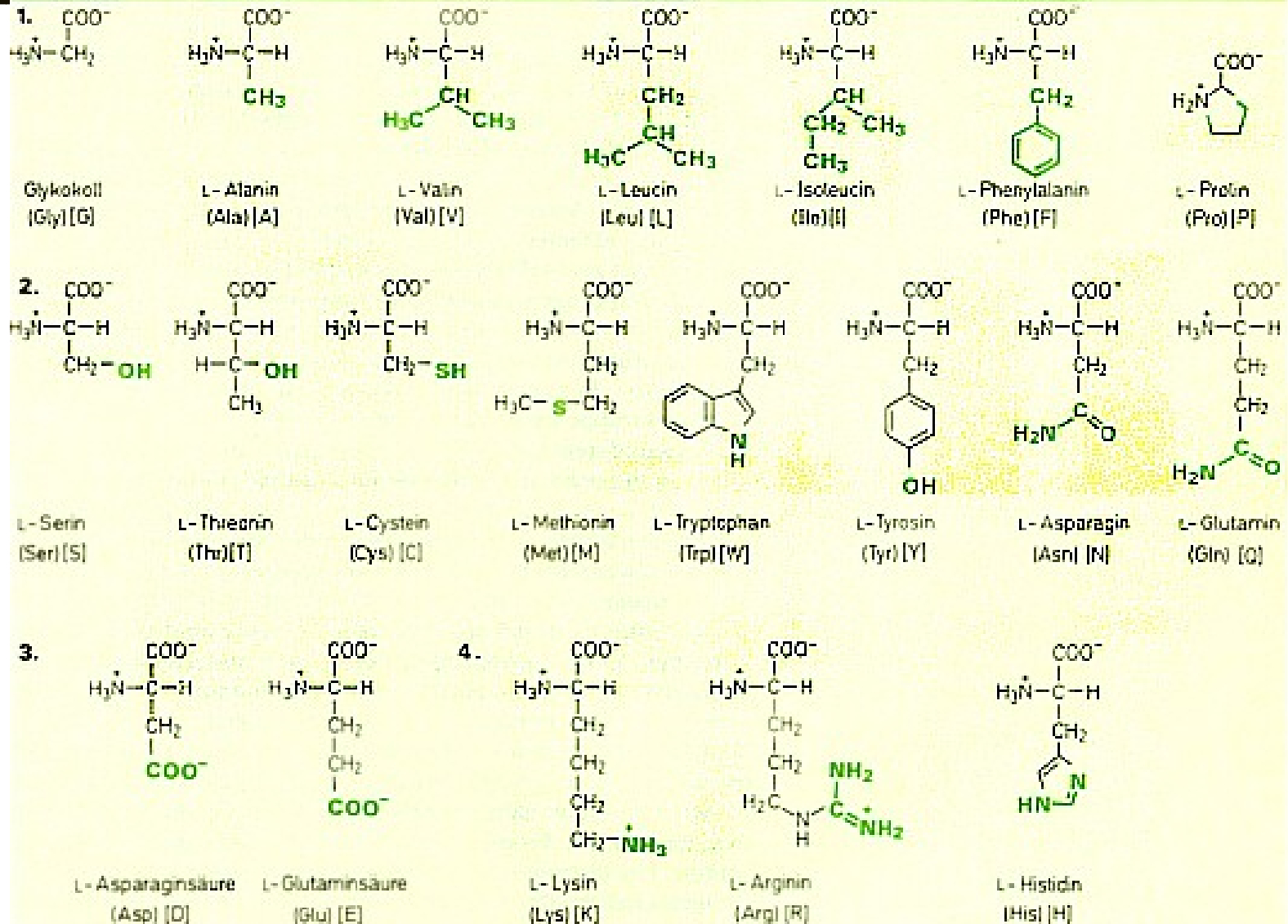


- H<sub>2</sub>O



Entstehung einer  
Peptidbindung zwischen  
zwei Aminosäuren

# Aminosäuren



# Aminosäuren

aliphatisch

klein

winzig

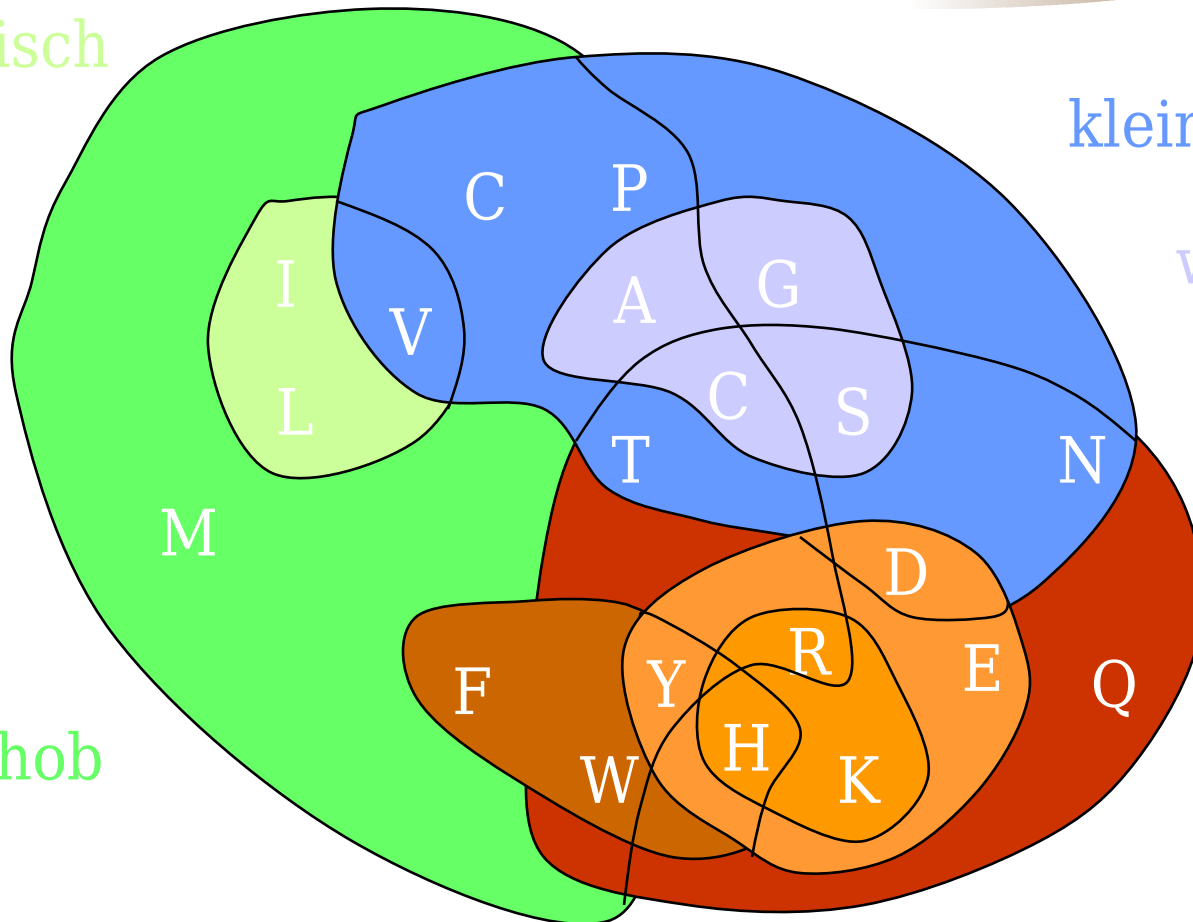
hydrophob

polar

geladen

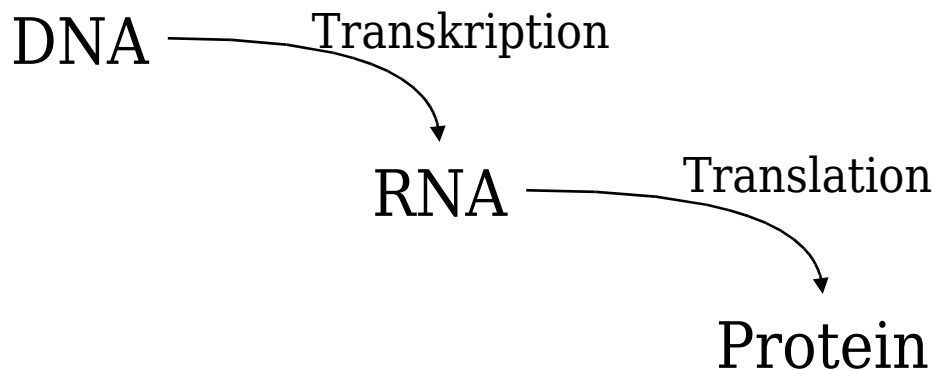
aromatisch

positiv



# *DNA → RNA → Protein*

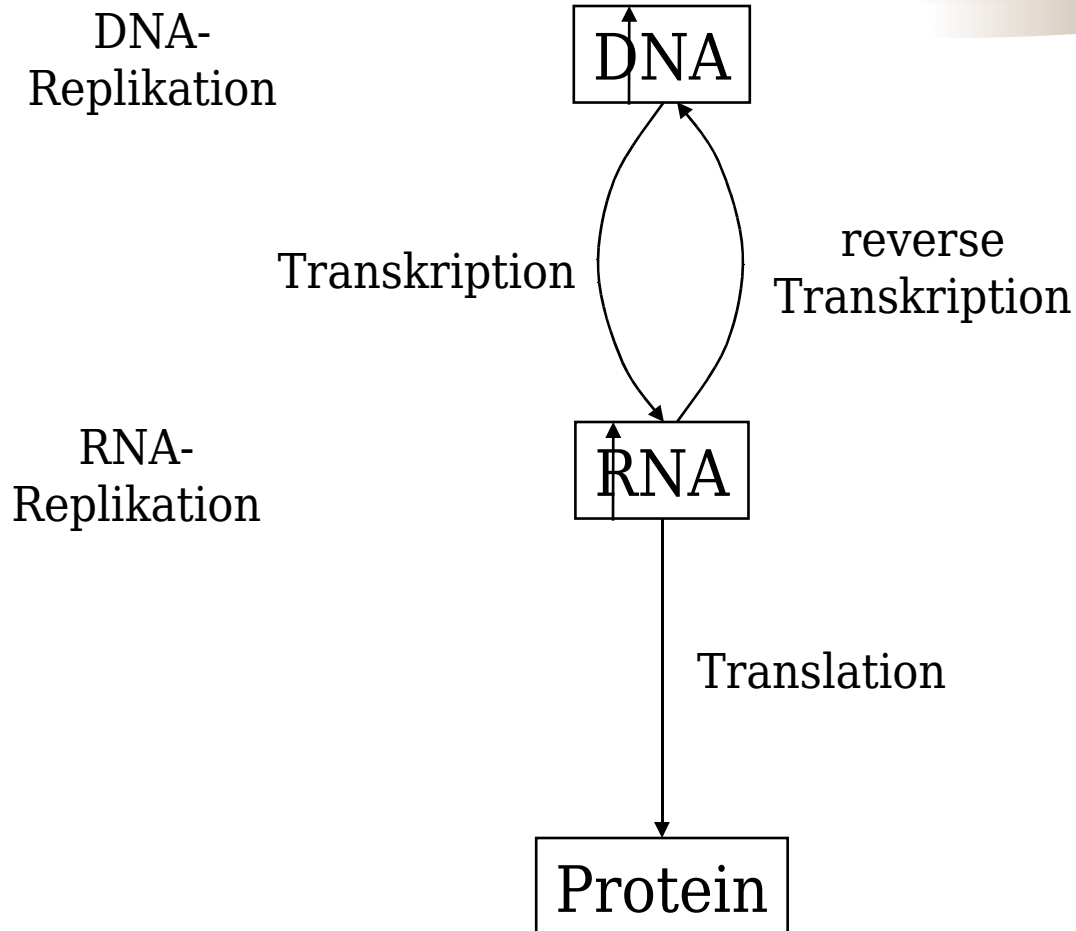
- Zentrales Dogma der Biochemie:
  - Der Fluß der Genetischen Information verläuft von der DNA zur RNA zum Protein



- 1964 kam die Hypothese auf, daß RNA Viren aus ihrer RNA wiederum DNA bilden können.
- 1970 wurde das verantwortliche Enzym gefunden. Folge: Die Lehrbücher mußten umgeschrieben werden

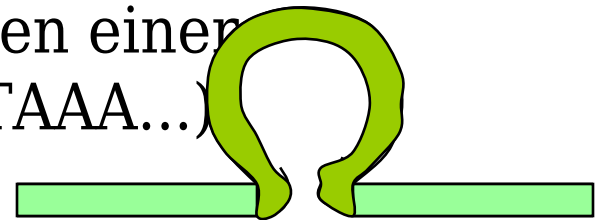


# Erweitertes Dogma

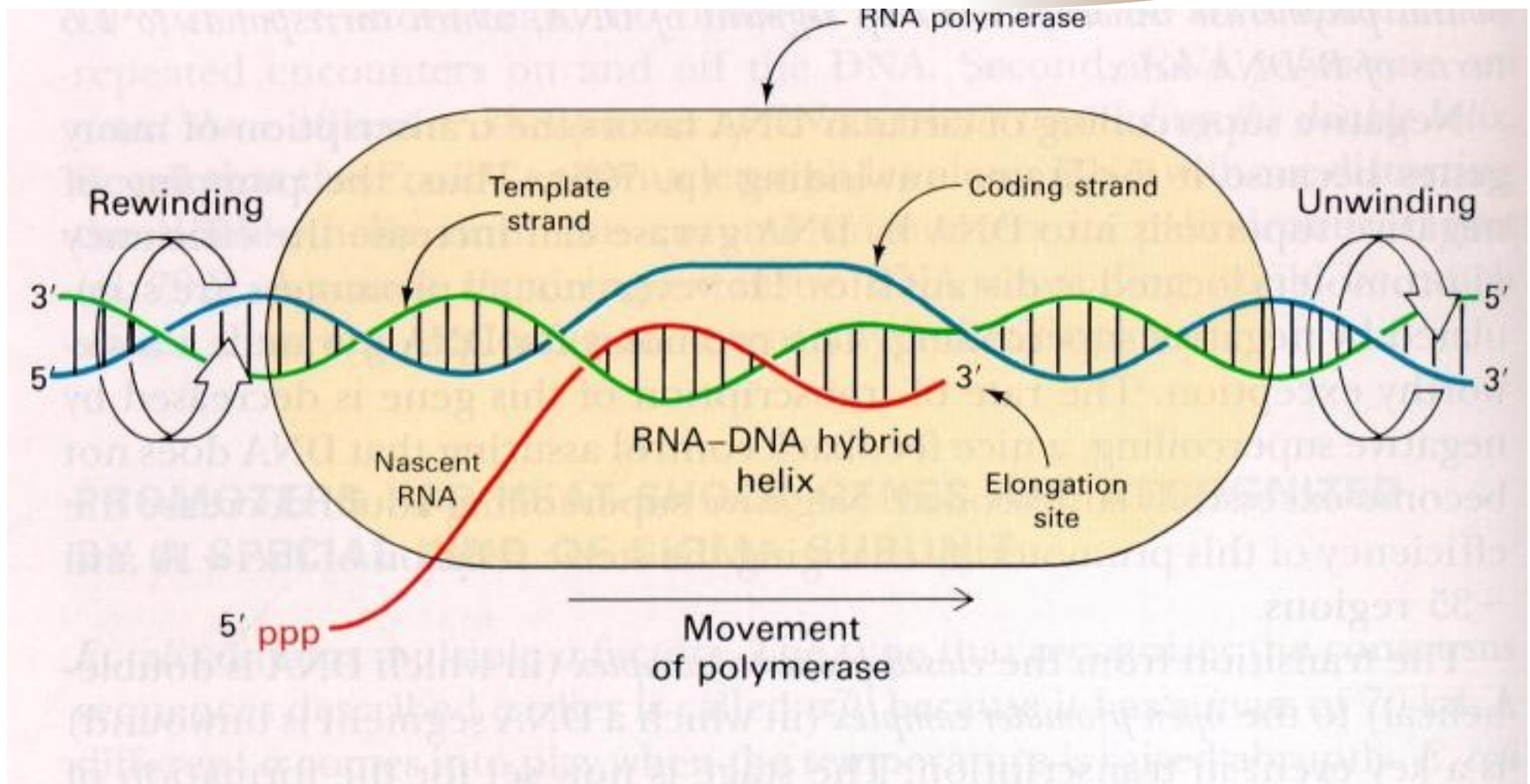


# Transkription

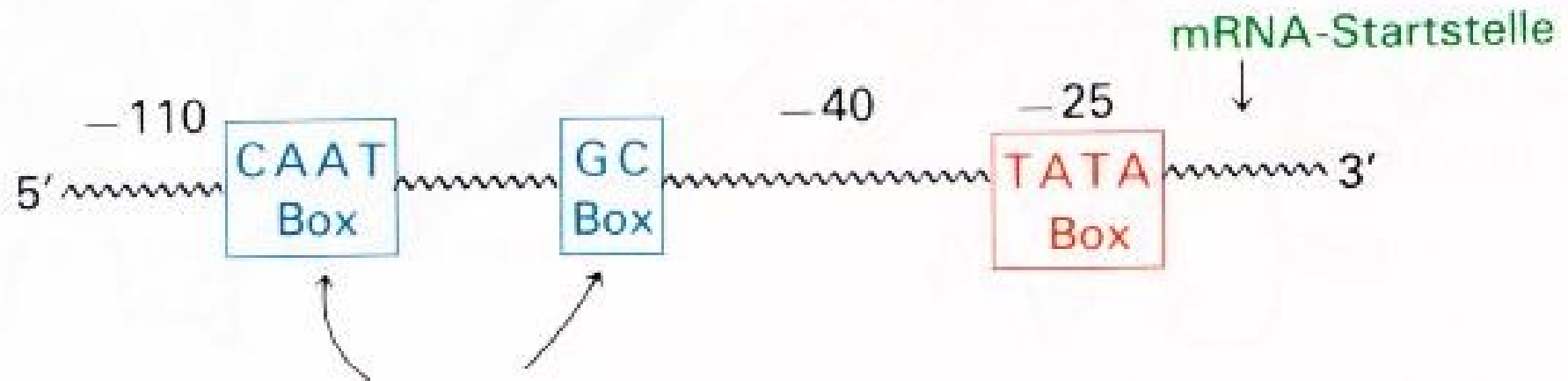
- Umschreibung von DNA in RNA
  - ablesen vom nicht codierenden Strang ( $3' \rightarrow 5'$ ), aber erzeugt Strang ( $5' \rightarrow 3'$ ),
  - ersetzen von Thymin durch Uracil
  - wird durch RNA-Polymerase katalysiert: jeweils auf kurzem Stück wird DNA in Einzelstränge aufgetrennt
  - Start bei Promoter-Sequenz
  - abspalten der RNA nach Erreichen einer bestimmten Sequenz (z.B. ...AATAAA...)
  - danach Splicing der RNA:
    - rausschneiden der Introns
    - alternatives Splicen für unterschiedliche



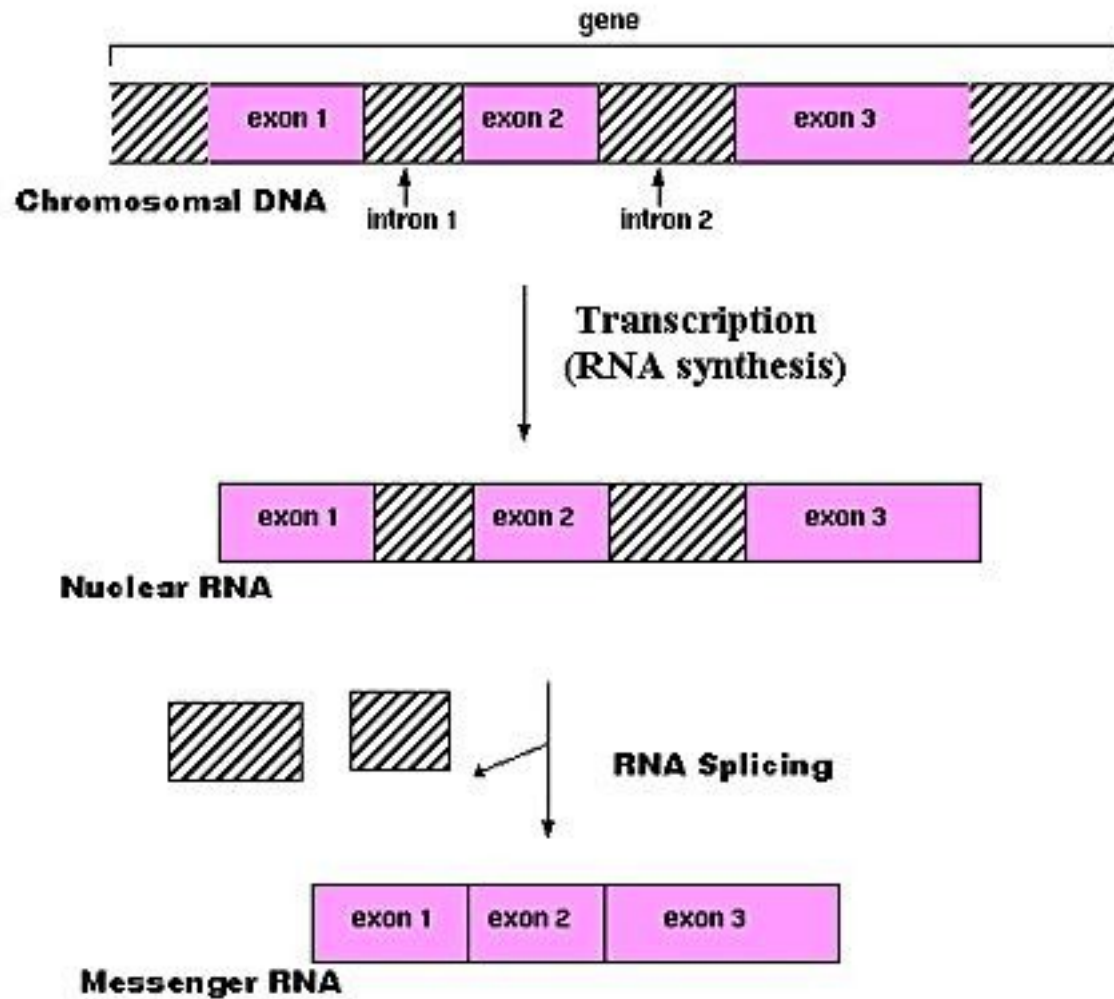
# Transkription



# TATA-Box



# Intron-Extron

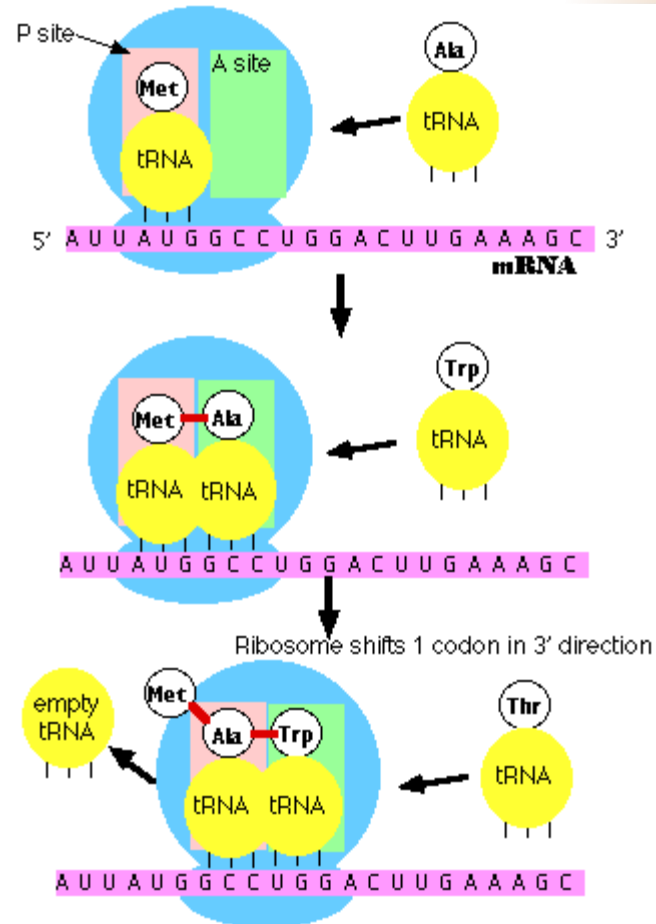


RNA synthesis and processing

# Translation

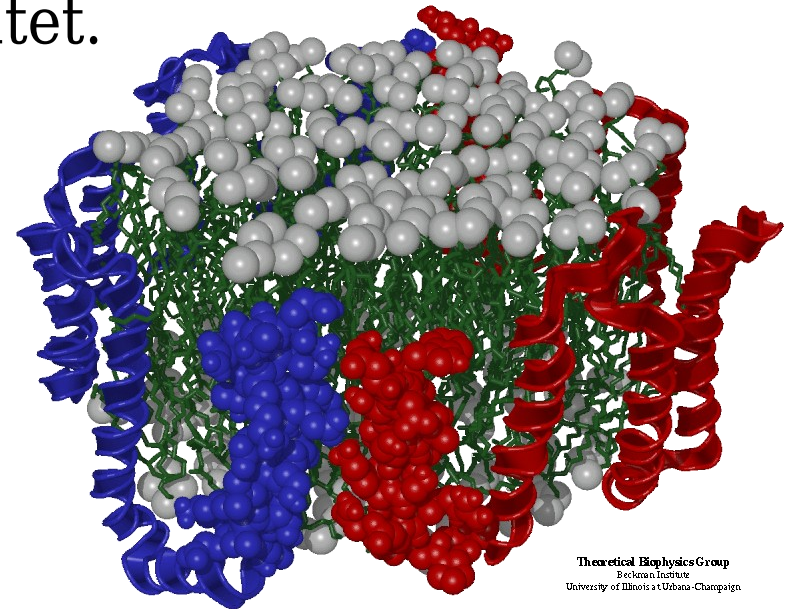
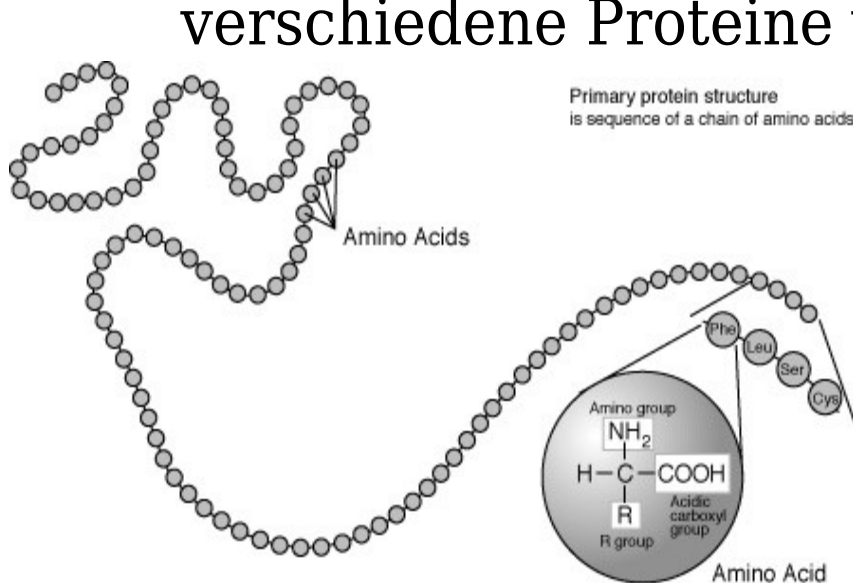
- Proteinsynthese im Ribosom
  - Teilnehmer: mRNA (bestimmen Reihenfolge der Aminosäuren), tRNA (transportieren Aminosäuren), Ribosome (Ort der Translation, Regulation der Bindung von mRNA und tRNA)
  - Initiation : Ribosom bindet an Start-Codon in mRNA 5' AUG 3'
  - Aktivierung : tRNA bindet an Start-Codon in mRNA und ist assoziiert zu A-Untereinheit des Ribosoms
  - Elongation: tRNA wird "weitergereicht" an P-Untereinheit, nächste tRNA mit passendem Anti-Codon bindet an mRNA
  - Aminosäure der ersten tRNA bindet an Aminosäure der zweiten tRNA, Rest verläßt das Ribosom

# Translation



# Proteine

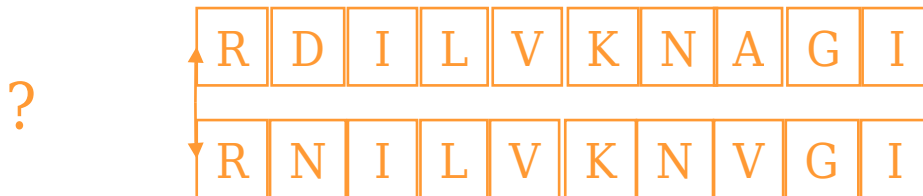
- Proteine sind Makromoleküle, die viele Abläufe in der Zelle bestimmen (Strukturbildung- und Erhaltung, Transport, Schutz und Abwehr, Steuerung und Regelung, Katalyse, Bewegung, Speicherung)
- Im menschlichen Körper werden etwa 100 000 verschiedene Proteine vermutet.





# Paarweiser Sequenzvergleich

- Analyse der evolutionären Beziehung
- Ausgangspunkt für die Vorhersage der Funktion eines Proteines oder seiner Struktur (oder beides)



# Ähnlichkeits-Dogma

- Wenn zwei Sequenzen sehr ähnlich sind, haben sie auch eine
  - eine ähnliche Funktion,
  - eine ähnliche Struktur,und sie haben einen gemeinsamen Vorfahren

**Vorsicht: das stimmt nicht immer !!!**

- Das impliziert, daß
  - die Sequenz eine Syntax bildet, die eine Funktion codiert
  - es gibt auch Redundanz, da einige Elemente ausgetauscht werden können, ohne daß sich die Funktion ändert (robuste Semantik)

# Proteinsequenzen

- Paarweiser Vergleich von je einem Buchstaben aus zwei Sequenzen, d.h. keine Betrachtung statistischer Abhängigkeiten innerhalb einer Sequenz
- Ähnlichkeit von zwei Sequenzen ergibt sich als Summe aus den Einzelähnlichkeiten (Markov-Modell)
- Aufstellung von sogenannten Scoring-Matrizen: Ähnlichkeitswert bezieht sich immer nur auf das dahinterliegende Modell
- Verfahren hauptsächlich für den Vergleich von Aminosäuresequenzen (Proteinen)

# Vergleich von Aminosäuren

- Einfachste Vergleichsmöglichkeit: Identitäts-Matrix
  - gleiche Buchstaben = 1,
  - ungleiche Buchstaben = 0
- Ähnlichkeitsmaße, die über = /  $\neq$  Vergleiche hinausgehen, nutzen
  - chemische oder strukturelle Eigenschaften: polar/unpolar, Form, Größe, Ladung
  - genetische Eigenschaften: minimale Anzahl ausgetauschter Basen in der dazugehörigen DNA
  - evolutionäre Distanz: beobachtete Austausch-Häufigkeiten von Aminosäuren (in bekannten



| A   | B   | C   | D   | E   | F   | G   | H   | I   | K   | L   | M   | N   | P   | Q   | R   | S   | T   | V   | W   | Y   | Z   |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 1.0 | 1.0 | 2.0 | 2.0 | 1.0 | 2.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 2.0 | 1.0 | 1.0 | 2.0 | 2.0 | 2.0 | 1.0 | 1.0 | 2.0 |
| 1.0 | 1.0 | 3.0 | 2.0 | 1.0 | 2.0 | 2.0 | 2.0 | 2.0 | 1.0 | 1.0 | 3.0 | 1.0 | 2.0 | 1.0 | 2.0 | 2.0 | 2.0 | 2.0 | 0.0 | 2.0 | 2.0 |
|     | 3.0 | 1.0 | 0.0 | 2.0 | 2.0 | 1.0 | 1.0 | 0.0 | 1.0 | 0.0 | 1.0 | 1.0 | 0.0 | 2.0 | 2.0 | 1.0 | 1.0 | 2.0 | 2.0 | 2.0 | 0.0 |
|     |     | 3.0 | 2.0 | 1.0 | 2.0 | 2.0 | 1.0 | 1.0 | 1.0 | 0.0 | 2.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 2.0 | 0.0 | 2.0 | 2.0 |
|     |     |     | 3.0 | 0.0 | 2.0 | 1.0 | 1.0 | 2.0 | 1.0 | 1.0 | 1.0 | 1.0 | 2.0 | 1.0 | 1.0 | 1.0 | 1.0 | 2.0 | 1.0 | 1.0 | 3.0 |
|     |     |     |     | 3.0 | 1.0 | 1.0 | 2.0 | 0.0 | 2.0 | 1.0 | 1.0 | 1.0 | 0.0 | 1.0 | 2.0 | 1.0 | 2.0 | 1.0 | 2.0 | 1.0 | 2.0 |
|     |     |     |     |     | 3.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 2.0 | 2.0 | 1.0 | 2.0 | 2.0 | 1.0 | 2.0 | 2.0 |
|     |     |     |     |     |     | 3.0 | 1.0 | 1.0 | 2.0 | 0.0 | 2.0 | 2.0 | 2.0 | 2.0 | 2.0 | 1.0 | 1.0 | 1.0 | 0.0 | 2.0 | 2.0 |
|     |     |     |     |     |     |     | 3.0 | 2.0 | 2.0 | 2.0 | 2.0 | 1.0 | 1.0 | 2.0 | 2.0 | 2.0 | 2.0 | 2.0 | 0.0 | 1.0 | 1.0 |
|     |     |     |     |     |     |     |     | 3.0 | 1.0 | 2.0 | 2.0 | 1.0 | 2.0 | 2.0 | 1.0 | 2.0 | 1.0 | 1.0 | 1.0 | 1.0 | 2.0 |
|     |     |     |     |     |     |     |     |     | 3.0 | 2.0 | 1.0 | 2.0 | 2.0 | 2.0 | 2.0 | 2.0 | 1.0 | 2.0 | 2.0 | 1.0 | 2.0 |
|     |     |     |     |     |     |     |     |     |     | 3.0 | 1.0 | 1.0 | 1.0 | 2.0 | 1.0 | 2.0 | 2.0 | 1.0 | 0.0 | 2.0 | 2.0 |
|     |     |     |     |     |     |     |     |     |     |     | 3.0 | 1.0 | 1.0 | 1.0 | 2.0 | 2.0 | 1.0 | 0.0 | 2.0 | 2.0 | 2.0 |
|     |     |     |     |     |     |     |     |     |     |     |     | 3.0 | 2.0 | 2.0 | 2.0 | 2.0 | 1.0 | 1.0 | 1.0 | 2.0 | 2.0 |
|     |     |     |     |     |     |     |     |     |     |     |     |     | 3.0 | 2.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 3.0 |
|     |     |     |     |     |     |     |     |     |     |     |     |     |     | 3.0 | 2.0 | 2.0 | 1.0 | 2.0 | 1.0 | 2.0 | 2.0 |
|     |     |     |     |     |     |     |     |     |     |     |     |     |     |     | 3.0 | 2.0 | 1.0 | 2.0 | 2.0 | 2.0 | 1.0 |
|     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     | 3.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 |
|     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     | 3.0 | 1.0 | 1.0 | 1.0 | 2.0 |
|     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     | 3.0 | 1.0 | 1.0 | 1.0 |
|     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     | 3.0 | 1.0 | 1.0 |

# Genetische Matrix

$B = D \vee N$

$Z = E \vee Q$

# PAM-Matrix

- basiert auf evolutionärem Modell
  - ähnliche Proteine haben einen gemeinsamen Vorfahren, aus dem beide Sequenzen durch genetische Veränderungen wie z.B. Punktmutationen hervorgegangen sind ( Edit-Distanz)
- empirisch aus Vorkommen von Aminosäuren in ähnlichen (mindestens 85% identischen), homologen Proteinen abgeschätzt
- PAM : Accepted Point Mutation
- PAM 1 - Matrix
  - 1 evolutionärer Schritt
  - 1 Mutation pro 100 Residuen erlaubt (1% Unterschied)  
wie hoch ist Wahrscheinlichkeit, daß sich ein Residuum

# *Berechnung von PAM*

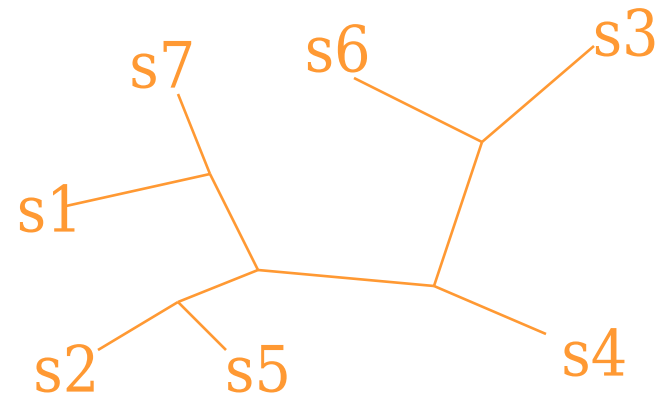


1. Schritt: Evolutionäres Modell aufstellen
2. Schritt: Häufigkeiten der Aminosäuren bestimmen
3. Schritt: Mutationshäufigkeit jeder Aminosäure bestimmen
4. Schritt: Matrix mit Austauschwahrscheinlichkeiten bestimmen
5. Schritt: Evolutionäre Skalierung
6. Schritt: Relative Wahrscheinlichkeit
7. Schritt: ...



# Schritt 1

- Aus nahe verwandten (85 % Identität) Sequenzen wird ein phylogenetischer Baum erzeugt

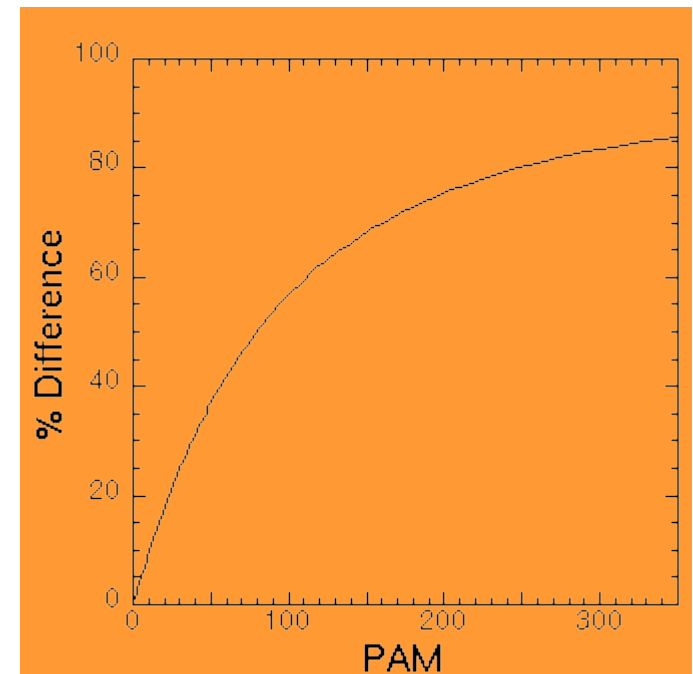


- Accepted Point Mutations:  
1 PAM (Percent Accepted Mutation) ist eine akzeptierte Punkt-Mutation pro 100  $A_{i,j}$  Residuen auf dem Weg zwischen zwei Sequenzen
- $A_{i,j}$  ist die Anzahl der beobachteten

# Evolutionäre Distanz

- Die evolutionäre Distanz entspricht (zumindest bei niedrigen PAM Werten) dem Anteil der unterschiedlichen Aminosäuren

| %Difference | PAM | %Difference | PAM |
|-------------|-----|-------------|-----|
| 1           | 1   | 45          | 67  |
| 5           | 5   | 50          | 80  |
| 10          | 11  | 55          | 94  |
| 15          | 17  | 60          | 112 |
| 20          | 23  | 65          | 133 |
| 25          | 30  | 70          | 159 |
| 30          | 38  | 75          | 195 |
| 35          | 47  | 80          | 246 |
| 40          | 56  | 85          | 328 |



# Schritt 6

M gibt die absolute Wahrscheinlichkeit eines Austauschs an.

Selbst bei zufälligen Sequenzen erhält man aber eine Übereinstimmung von etwa 5%.

Daher müssen die Werte für einen Austausch einer Aminosäure noch in Relation zu dem Wert für einen zufälligen Austausch gesetzt werden.

$p_i^{random} = f_i$  Wahrscheinlichkeit für einen zufälligen Austausch

$R_{i,j} = \frac{M_{i,j}}{p_i^{random}}$  Matrix mit relativen Wahrscheinlichkeiten

# Dayhoff Matrix

|   |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |
|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|---|
| A | 2  |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |
| R | -2 | 6  |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |
| N | 0  | 0  | 2  |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |
| D | 0  | -1 | 2  | 4  |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |
| C | -2 | -4 | -4 | -5 | 4  |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |
| Q | 0  | 1  | 1  | 2  | -5 | 4  |    |    |    |    |    |    |    |    |    |    |    |    |    |   |
| E | 0  | -1 | 1  | 3  | -5 | 2  | 4  |    |    |    |    |    |    |    |    |    |    |    |    |   |
| G | 1  | -3 | 0  | 1  | -3 | -1 | 0  | 5  |    |    |    |    |    |    |    |    |    |    |    |   |
| H | -1 | 2  | 2  | 1  | -3 | 3  | 1  | -2 | 6  |    |    |    |    |    |    |    |    |    |    |   |
| I | -1 | -2 | -2 | -2 | -2 | -2 | -2 | -3 | -2 | 5  |    |    |    |    |    |    |    |    |    |   |
| L | -2 | -3 | -3 | -4 | -6 | -2 | -3 | -4 | -2 | 2  | 6  |    |    |    |    |    |    |    |    |   |
| K | -1 | 3  | 1  | 0  | -5 | 1  | 0  | -2 | 0  | -2 | -3 | 5  |    |    |    |    |    |    |    |   |
| M | -1 | 0  | -2 | -3 | -5 | -1 | -2 | -3 | -2 | 2  | 4  | 0  | 6  |    |    |    |    |    |    |   |
| F | -4 | -4 | -4 | -6 | -4 | -5 | -5 | -5 | -2 | 1  | 2  | -5 | 0  | 9  |    |    |    |    |    |   |
| P | 1  | 0  | -1 | -1 | -3 | 0  | -1 | -1 | 0  | -2 | -3 | -1 | -2 | -5 | 6  |    |    |    |    |   |
| S | 1  | 0  | 1  | 0  | 0  | -1 | 0  | 1  | -1 | -1 | -3 | 0  | -2 | -3 | 1  | 3  |    |    |    |   |
| T | 1  | -1 | 0  | 0  | -2 | -1 | 0  | 0  | -1 | 0  | -2 | 0  | -1 | -2 | 0  | 1  | 3  |    |    |   |
| W | -6 | 2  | -4 | -7 | -8 | -5 | -7 | -7 | -3 | -5 | -2 | -3 | -4 | 0  | -6 | -2 | -5 | 17 |    |   |
| Y | -3 | -4 | -2 | -4 | 0  | -4 | -4 | -5 | 0  | -1 | -1 | -4 | -2 | 7  | -5 | -3 | -3 | 0  | 10 |   |
| V | 0  | -2 | -2 | -2 | -2 | -2 | -2 | -1 | -2 | 4  | 2  | -2 | 2  | -1 | -1 | -1 | 0  | -6 | -2 | 4 |
|   | A  | R  | N  | D  | C  | Q  | E  | G  | H  | I  | L  | K  | M  | F  | P  | S  | T  | W  | Y  | V |

MDM78 PAM250

# Alternative Scoring-Matrizen

- Zahlreiche andere Matrizen
  - anders normalisiert (GCG Version von MDM78 PAM250)
  - neu berechnet aus neueren und umfangreicheren Daten (PET)
  - mit Hilfe von großen Datenmengen (1.7Mio) direkt für größere evolutionäre Distanzen berechnet (GCB, G. Gonnet, M.A. Cohen, & S. Benner (1992))
  - abgeleitet aus Blöcken hoch-konservierter Bereiche in den Proteinen, die mit einem Schwellwert geclustert werden (BLOSUM-t)
  - aus Vergleich der Tertiärstruktur bei Proteinen mit bekannter 3D-Struktur andere Austausch-

| B  | C  | D  | E  | F  | G  | H  | I  | K  | L  | M  | N  | P  | Q  | R  | S  | T  | V  | W  | X  | Y  | Z  |   |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|---|
| -2 | 0  | -2 | -1 | -2 | 0  | -2 | -1 | -1 | -1 | -1 | -2 | -1 | -1 | -1 | 1  | 0  | 0  | -3 | -1 | -2 | -1 | A |
| 6  | -3 | 6  | 2  | -3 | -1 | -1 | -3 | -1 | -4 | -3 | 1  | -1 | 0  | -2 | 0  | -1 | -3 | -4 | -1 | -3 | 2  | B |
|    | 9  | -3 | -4 | -2 | -3 | -3 | -1 | -3 | -1 | -1 | -3 | -3 | -3 | -3 | -1 | -1 | -1 | -2 | -1 | -2 | -4 | C |
|    |    | 6  | 2  | -3 | -1 | -1 | -3 | -1 | -4 | -3 | 1  | -1 | 0  | -2 | 0  | -1 | -3 | -4 | -1 | -3 | 2  | D |
|    |    |    | 5  | -3 | -2 | 0  | -3 | 1  | -3 | -2 | 0  | -1 | 2  | 0  | 0  | -1 | -2 | -3 | -1 | -2 | 5  | E |
|    |    |    |    | 6  | -3 | -1 | 0  | -3 | 0  | 0  | -3 | -4 | -3 | -3 | -2 | -2 | -1 | 1  | -1 | 3  | -3 | F |
|    |    |    |    |    | 6  | -2 | -4 | -2 | -4 | -3 | 0  | -2 | -2 | -2 | 0  | -2 | -3 | -2 | -1 | -3 | -2 | G |
|    |    |    |    |    |    | 8  | -3 | -1 | -3 | -2 | 1  | -2 | 0  | 0  | -1 | -2 | -3 | -2 | -1 | 2  | 0  | H |
|    |    |    |    |    |    |    | 4  | -3 | 2  | 1  | -3 | -3 | -3 | -3 | -2 | -1 | 3  | -3 | -1 | -1 | -3 | I |
|    |    |    |    |    |    |    |    | 5  | -2 | -1 | 0  | -1 | 1  | 2  | 0  | -1 | -2 | -3 | -1 | -2 | 1  | K |
|    |    |    |    |    |    |    |    |    | 4  | 2  | -3 | -3 | -2 | -2 | -2 | -1 | 1  | -2 | -1 | -1 | -3 | L |
|    |    |    |    |    |    |    |    |    |    | 5  | -2 | -2 | 0  | -1 | -1 | -1 | 1  | -1 | -1 | -1 | -2 | M |
|    |    |    |    |    |    |    |    |    |    |    | 6  | -2 | 0  | 0  | 1  | 0  | -3 | -4 | -1 | -2 | 0  | N |
|    |    |    |    |    |    |    |    |    |    |    |    | 7  | -1 | -2 | -1 | -1 | -2 | -4 | -1 | -3 | -1 | P |
|    |    |    |    |    |    |    |    |    |    |    |    |    | 5  | 1  | 0  | -1 | -2 | -2 | -1 | -1 | 2  | Q |
|    |    |    |    |    |    |    |    |    |    |    |    |    |    | 5  | -1 | -1 | -3 | -3 | -1 | -2 | 0  | R |
|    |    |    |    |    |    |    |    |    |    |    |    |    |    |    | 4  | 1  | -2 | -3 | -1 | -2 | 0  | S |
|    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    | 5  | 0  | -2 | -1 | -2 | -1 | T |
|    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    | 4  | -3 | -1 | -1 | -2 | V |
|    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    | 11 | -1 | 2  | -3 | W |
|    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    | -1 | -1 | -1 | X |
|    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    | 7  | -2 | Y |
|    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    | 5  | Z |

# BLOSUM 62

X : any

# Beispiel

|   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|
| R | D | I | L | V | K | N | A | G | I |
| R | N | I | L | V | K | N | V | G | I |

Identitäts-Matrix :  $1 + 0 + 1 + 1 + 1 + 1 + 1 + 0 + 1 + 1 = 8$

Genetische Matrix:  $3 + 2 + 3 + 3 + 3 + 3 + 3 + 2 + 3 + 3 = 28$

AM250:  $6 + 2 + 5 + 6 + 4 + 5 + 2 + 0 + 5 + 5 = 40$

LOSUM62:  $5 + 1 + 4 + 4 + 4 + 5 + 6 + 0 + 6 + 4 = 39$

# Auswahl der Matrix



- Hängt stark von den zu vergleichenden Sequenzen ab
  - je nach (bekanntem) Verwandtschaftsgrad  
120 PAM - 250 PAM
- BLOSUM oft recht gut bei bestimmten Suchverfahren (z.B. BLAST), aber nicht durchgängig für alle Protein-Familien



# Datenbanksuche

- **Ziel:** Aus einer Sequenz-Datenbank alle Einträge bestimmen, die ähnlich zu einer vom Benutzer vorgegebenen Query-Sequenz sind
- **Problem:** Aufwand für ein optimales paarweises Alignment der Query-Sequenz gegen alle Datenbankeinträge viel zu hoch
- **Ansatz:** Heuristische Verfahren
  - BLAST
  - FASTA

- Basic Local Alignment Search Tool
- generiert eine Liste von Segment-Paaren (Teilsequenzen gleicher Länge ohne Gaps) zwischen der Query-Sequenz und Einträgen der Datenbank, die einen Score oberhalb einer vorgegebenen Schwelle besitzen

- Ausgangspunkt: Query-Sequenz, Wortlänge  $w$ , Schwellwerte  $S, T$
- Drei-schrittiger Algorithmus:
  - für eine vorgegebene Wort-Länge  $w$  und eine Score-Matrix werden alle Wörter der Länge  $w$  bestimmt, die bei einem Vergleich mit der Query-Sequenz einen Score  $> T$  ergeben würden
  - Die Datenbank wird auf diese sog.  $w$ -Mere hin durchsucht
  - Jeder Treffer (Sequenz aus DB) wird in beide Richtungen erweitert (ohne Gaps) und es wird geprüft, ob sich ein Score  $> S$  ergibt
  - Ausgabe aller Segmente mit einem Score  $> S$

- Varianten:
  - BLASTn : für Nukleotidsequenzen (DNA)
  - BLASTp: für Aminosäuresequenzen (Proteine)
  - BLASTx: macht eine 6-Frame Translation

# BLAST: Beispiel

http://www.ncbi.nlm.nih.gov/BLAST/

The screenshot shows the NCBI BLAST website in a Mozilla Firefox browser window. The browser's address bar displays the URL <http://www.ncbi.nlm.nih.gov/BLAST/>. The website header includes the BLAST logo and navigation links: Home, Recent Results, Saved Strategies, and Help. A 'My NCBI' section contains links for Sign In and Register. The main content area is titled 'NCBI/BLAST Home' and features a search box with the text 'BLAST finds regions of similarity between biological sequences. [more...](#)'. Below this is a link to 'Learn more about how to use the new BLAST design' and a link to 'Old blast'. The 'BLAST Assembled Genomes' section prompts users to 'Choose a species genome to search, or [list all genomic BLAST databases](#)'. A grid of links lists various species: Human, Mouse, Rat, Arabidopsis thaliana, Oryza sativa, Bos taurus, Danio rerio, Drosophila melanogaster, Gallus gallus, Pan troglodytes, Microbes, and Apis mellifera. The 'Basic BLAST' section asks users to 'Choose a BLAST program to run' and lists options: nucleotide blast, protein blast, blastx, tblastn, and tblastx, each with a brief description and algorithms. On the right side, there are two news sections: 'News' with a link to 'Old BLAST Web Pages to be deleted June 11th 2007' and 'Tip of the Day' with the title 'Using Genomic BLAST' and a link to 'More tips...'. The browser's taskbar at the bottom shows the Start button and several open applications, including Microsoft Outlook Web Access, Microsoft PowerPoint, and the BLAST application itself. The system clock shows the time as 22:46.

# BLAST: Beispiel

UniProtKB/Swiss-Prot entry **Q9GTW9**

[\[Entry info\]](#) [\[Name and origin\]](#) [\[References\]](#) [\[Comments\]](#) [\[Cross-references\]](#) [\[Keywords\]](#) [\[Features\]](#) [\[Sequence\]](#) [\[Tools\]](#)

*Note: most headings are clickable, even if they don't appear as links. They link to the user manual or other documents.*

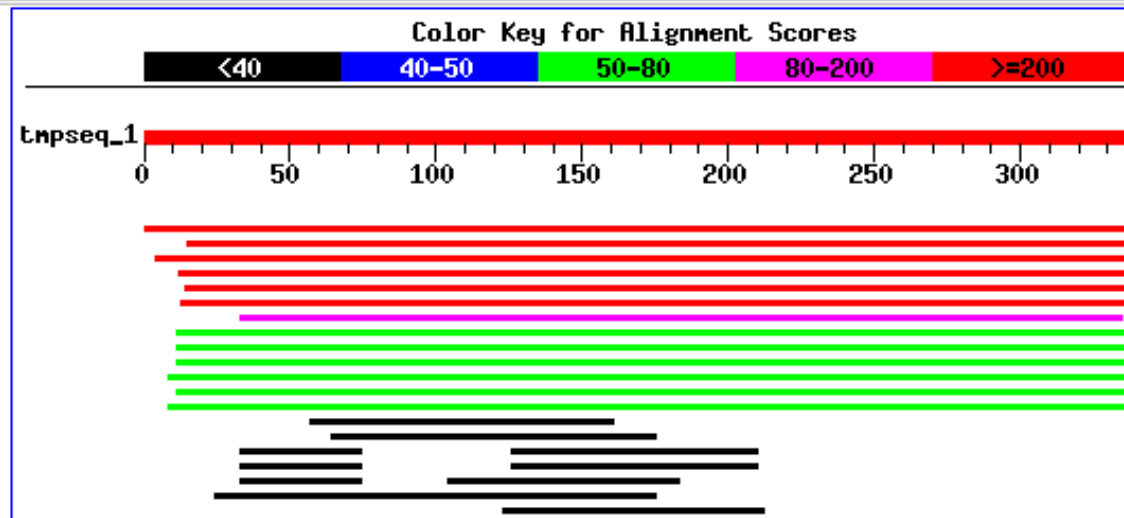
| Entry information                 |                                    |
|-----------------------------------|------------------------------------|
| Entry name                        | GLK1_TRIVA                         |
| Primary accession number          | <b>Q9GTW9</b>                      |
| Secondary accession numbers       | None                               |
| Integrated into Swiss-Prot on     | January 23, 2002                   |
| Sequence was last modified on     | March 1, 2001 (Sequence version 1) |
| Annotations were last modified on | May 1, 2007 (Entry version 24)     |

| Name and origin of the protein |   |
|--------------------------------|---|
| Protein name                   | <b>Glucokinase 1</b>  |
| Synonyms                       | <b>EC 2.7.1.2</b><br><b>Glucose kinase 1</b><br><b>Hexokinase 1</b> |
| Gene name                      | <b>Name: GK1</b>  |

# BLAST: Beispiel

## Distribution of 23 Blast Hits on the Query Sequence

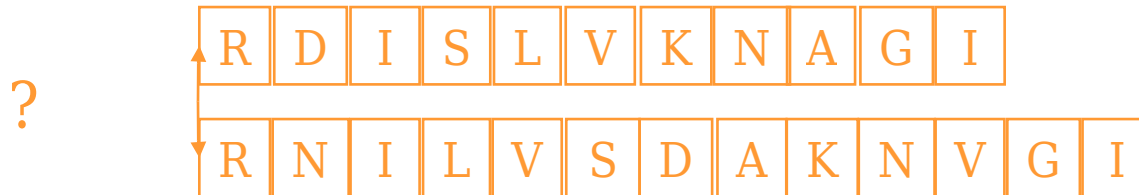
Mouse-over to show defline and scores. Click to show alignments



Sequences producing significant alignments:

|  | Score               | E     |
|--|---------------------|-------|
|  | (bits)              | Value |
| <a href="#">ref NP_011083.1</a> Yer156cp >gi 731528 sp P40093 YEY6_YEAST H...  | <a href="#">660</a> | 0.0   |
| <a href="#">emb CAB71842.1</a> (AL138666) conserved hypothetical protein [...  | <a href="#">309</a> | 2e-83 |
| <a href="#">gb AAF64518.1 AF252871.1</a> (AF252871) GAMM1 protein [Mus musc... | <a href="#">262</a> | 3e-69 |
| <a href="#">pir  T19538</a> hypothetical protein K08H10.8 - Caenorhabditis ... | <a href="#">257</a> | 1e-67 |
| <a href="#">dbj BAB08430.1</a> (AB017067) GAMM1 protein-like [Arabidopsis ...  | <a href="#">255</a> | 6e-67 |

- Neben Substitutionen können auch Einfügungen und Löschungen vorkommen



## Alignment

|   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| R | D | I | S | L | V | - | - | - | K | N | A | G | I |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

|   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| R | N | I | - | L | V | S | D | A | K | N | V | G | I |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|



# Gap Penalty

- Gaps in einem Alignment (Paarung einer Aminosäure mit einer Lücke) werden mit einem schlechten Wert bestraft, z.B. negativer Wert für jedes
  - einzelne Gap, d.h. Funktion linear in der Länge  $k$  der eingefügten bzw. gelöschten Elemente

$$g(k) = q \cdot k$$

- oder zusammengesetzter Wert für Einfügungen/Auslassungen beliebiger Länge

a : Gap Eröffnungsstrafe    b: Gap

Ausweitungsstrafe

$$g(k) = a + b \cdot k$$

affin-lineare Gap Penalty

# Beispiel

R D I S L V - - - K N A G I

R N I - L V S D A K N V G I

titäts-M. :  $1 + 0 + 1 - g(1) + 1 + 1 - g(3) + 1 + 1 + 0 + 1 + 1$

1250 (\*10) :  $6 + 2 + 5 - g(1) + 6 + 4 - g(3) + 5 + 2 + 0 + 5 + 5$

# *Bioinformatik & Internet*



- In den Anfangsjahren: Publikation aller Daten (Sequenzen, Gene, Proteine) in Journalen
- Mittlerweile: Speicherung der Daten in meist öffentlich zugänglichen Datenbanken
- Datenbanken oft sehr einfach strukturiert
- Ziel heute: Integration der Daten durch entsprechende Werkzeuge (Java,

# *Datenbanken im Netz*



- Gensequenzen (DNA + RNA):  
Kollaboration zwischen drei, weltweit verteilten Instituten mit täglichen Austausch der Daten (seit 1980)
  - GenBank at NCBI (National Center for Biotechnology Information)  
[www.ncbi.nlm.nih.gov](http://www.ncbi.nlm.nih.gov)
  - European Molecular Biology Laboratory (EMBL) [www.embl.de](http://www.embl.de)
  - DNA databank of Japan ( DDBJ)  
[www.ddbj.nig.ac.jp](http://www.ddbj.nig.ac.jp)

GenBank Overview - Microsoft Internet Explorer

Adresse <http://www.ncbi.nlm.nih.gov/GenBank/index.html>

NCBI **GenBank Overview**

PubMed Entrez BLAST OMIM Books Taxonomy Structure

Search Entrez for  Go

NCBI  
SITE MAP  
Submit to GenBank  
Updates  
Search GenBank  
Entrez Nucleotide  
BLAST

### International sequence databases exceed 100 gigabases

In August 2005, the INSDC announced the DNA sequence database exceeded 100 gigabases. GenBank is proud of its contributions toward this milestone. We thank all the scientists who have worked through the submission process at GenBank and made their sequence data available to the world. See the related [press release](#).

#### Growth of the International Nucleotide Sequence Database Collaboration

| Date   | GenBank (Billions) | EMBL (Billions) | DDBJ (Billions) | Other (Billions) | Total (Billions) |
|--------|--------------------|-----------------|-----------------|------------------|------------------|
| Aug-00 | ~1                 | ~1              | ~1              | ~1               | ~4               |
| Aug-01 | ~2                 | ~2              | ~2              | ~2               | ~8               |
| Aug-02 | ~5                 | ~5              | ~5              | ~5               | ~20              |
| Aug-03 | ~15                | ~15             | ~15             | ~15              | ~60              |
| Aug-04 | ~40                | ~10             | ~10             | ~10              | ~70              |
| Aug-05 | ~85                | ~10             | ~10             | ~10              | ~115             |

Base Pairs contributed by GenBank® EMBL DDBJ

Internet

# *Datenbanken im Netz*



- Protein Datenbanken
  - SWISSPROT <http://www.expasy.org>
  - PIR <http://pir.georgetown.edu/>
- Suche nach Proteinen
- Muster- und Sequenzvergleich
- Querreferenzen zu anderen Datenbanken
- Literaturreferenzen
- annotierte (und geprüfte Information)

ExPASy Proteomics Server - Microsoft Internet Explorer


Datei Bearbeiten Ansicht Favoriten Extras ?

Zurück Suchen Favoriten Medien

Adresse <http://www.expasy.org/> Wechseln zu Links


[Site Map](#)
[Search ExPASy](#)
[Contact us](#)

Search  for



# ExPASy Proteomics Server

The ExPASy (**Expert Protein Analysis System**) proteomics server of the [Swiss Institute of Bioinformatics](#) (SIB) is dedicated to the analysis of protein sequences and structures as well as 2-D PAGE ([Disclaimer](#) / [References](#)).



In-Silico Analysis of Proteins  
Celebrating the 20th Anniversary of Swiss-Prot

[\[Announcements\]](#)
[\[Job opening\]](#)
[\[Mirror Sites\]](#)

| Databases  | Tools and software packages  |
|--|--|
| <ul style="list-style-type: none"> <li>• <a href="#">UniProt Knowledgebase (Swiss-Prot and TrEMBL)</a> - Protein knowledgebase</li> <li>• <a href="#">PROSITE</a> - Protein families and domains</li> <li>• <a href="#">SWISS-2DPAGE</a> - Two-dimensional polyacrylamide gel electrophoresis</li> <li>• <a href="#">ENZYME</a> - Enzyme nomenclature</li> <li>• <a href="#">SWISS-MODEL Repository</a> - Automatically generated protein models</li> <li>• <a href="#">Links to many other molecular biology databases</a></li> </ul> | <ul style="list-style-type: none"> <li>• <a href="#">Proteomics and sequence analysis tools</a> <ul style="list-style-type: none"> <li>○ <a href="#">Proteomics</a></li> <li>○ <a href="#">DNA -&gt; Protein</a></li> <li>○ <a href="#">Similarity searches (BLAST...)</a></li> <li>○ <a href="#">Pattern and profile searches (ScanProsite...)</a></li> <li>○ <a href="#">Post-translational modification and topology prediction</a></li> <li>○ <a href="#">Primary structure analysis</a></li> <li>○ <a href="#">Secondary and tertiary structure tools (Swiss-PdbViewer...)</a></li> <li>○ <a href="#">Alignment and Phylogenetic analysis</a></li> </ul> </li> <li>• <a href="#">ImageMaster / Melanie</a> - Software for 2-D PAGE analysis</li> <li>• <a href="#">MSight</a> - Mass Spectrometry Imager</li> <li>• <a href="#">Roche Applied Science's Biochemical Pathways</a></li> </ul> |
| Education and services   | Documentation  |
| <ul style="list-style-type: none"> <li>• <a href="#">The ExPASy FTP server</a></li> </ul>  | <ul style="list-style-type: none"> <li>• <a href="#">What's New on ExPASy</a></li> </ul>   |

Internet

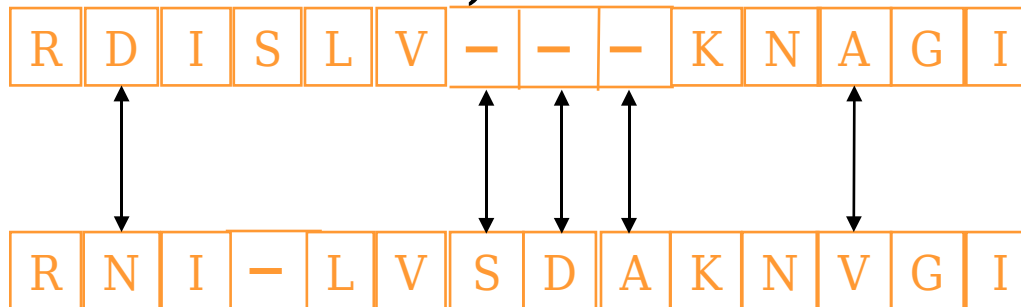
# Datenbankeintrag

```
ID CSAR_HUMAN STANDARD; PRT; 350 AA.
AC P21730;
DT 01-MAY-1991 (Rel. 18, Created)
DT 01-MAY-1991 (Rel. 18, Last sequence update)
DT 15-JUL-1998 (Rel. 36, Last annotation update)
DE C5A ANAPHYLATOXIN CHEMOTACTIC RECEPTOR (C5A-R) (CD88 ANTIGEN).
GN CSR1 OR CSAR.
OS Homo sapiens (Human).
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Mammalia;
OC Eutheria; Primates; Catarrhini; Hominidae; Homo.
RN [1]
RP SEQUENCE FROM N.A.
RX MEDLINE; 91156029.
RA GERARD N.P., GERARD C.;
RT "The chemotactic receptor for human C5a anaphylatoxin.";
RL Nature 349:614-617(1991).
RN [2]
RP SEQUENCE FROM N.A.
RX MEDLINE; 91175748.
RA BOULAY F., MERY L., TARDIF M., BROUCHON L., VIGNAIS P.;
RT "Expression cloning of a receptor for C5a anaphylatoxin on
RT differentiated HL-60 cells.";
RL Biochemistry 30:2993-2999(1991).
CC -!- FUNCTION: RECEPTOR FOR THE CHEMOTACTIC AND INFLAMMATORY
PEPTIDE
CC ANAPHYLATOXIN C5A. THIS RECEPTOR STIMULATES CHEMOTAXIS,
GRANULE
CC ENZYME RELEASE AND SUPEROXIDE ANION PRODUCTION.
CC -!- SUBCELLULAR LOCATION: INTEGRAL MEMBRANE PROTEIN.
CC -!- SIMILARITY: BELONGS TO FAMILY 1 OF G-PROTEIN COUPLED
RECEPTORS.
CC -!- DATABASE: NAME=PROW; NOTE=CD guide CD88 entry;
CC www="http://www.ncbi.nlm.nih.gov/prow/cd/cd88.htm".
CC
```



# Alignment

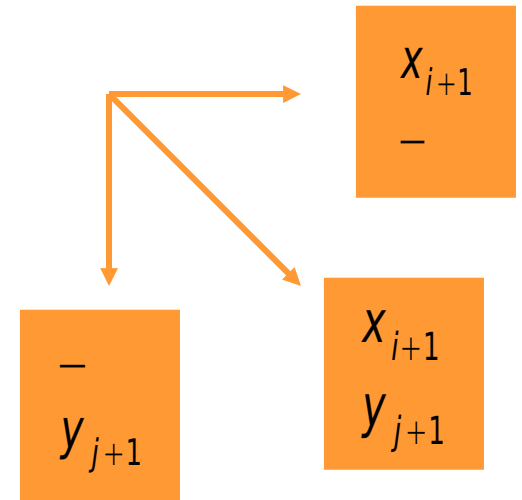
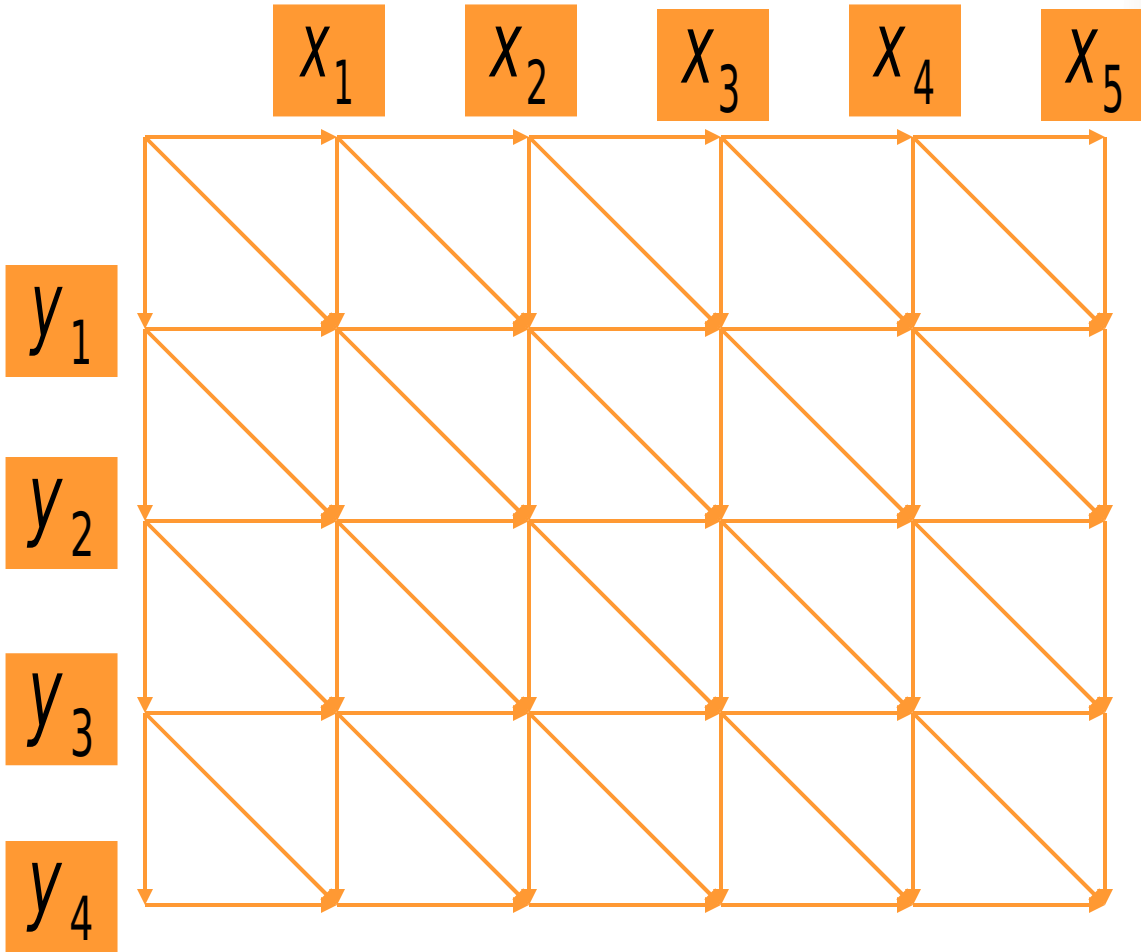
- Suche nach einem Alignment von zwei Sequenzen, so daß der Ähnlichkeitswert maximal ist (oder die Distanz minimal)



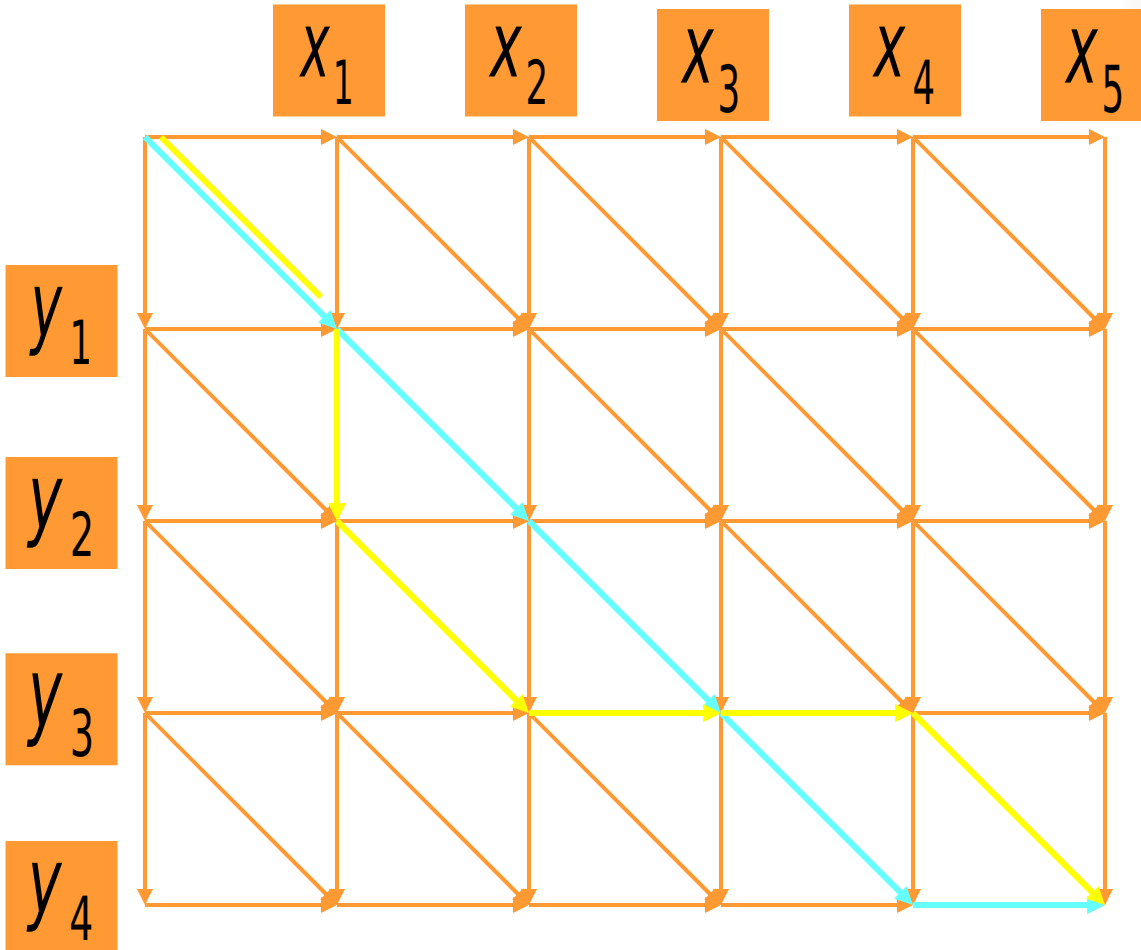
$$\binom{2n}{n} = \frac{(2n)!}{(n!)^2} \approx \frac{2^{2n}}{\sqrt{2\pi n}}$$

Möglichkeiten für zwei Sequenzen der Länge  $n$

# Alignment



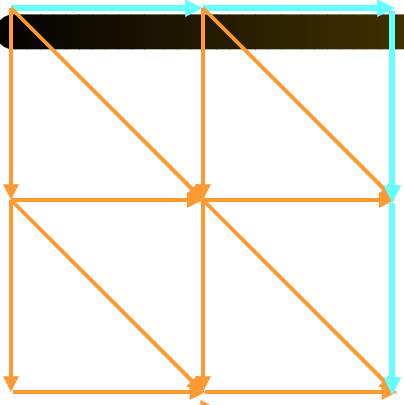
# Alignment



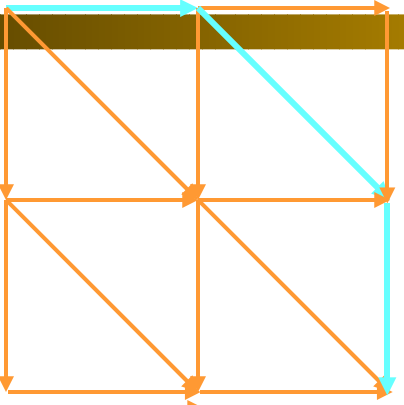
$X_1 X_2 X_3 X_4 X_5$   
 $y_1 y_2 y_3 y_4 -$

$X_1 - X_2 X_3 X_4 X_5$   
 $y_1 y_2 y_3 - - y_4$

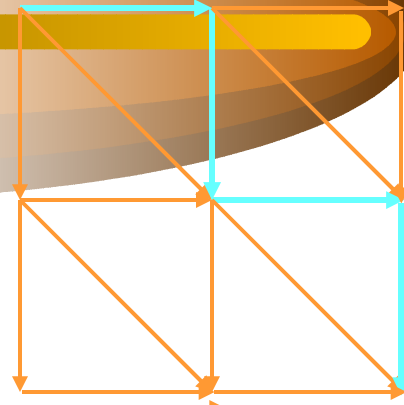
# *Pfade in F*



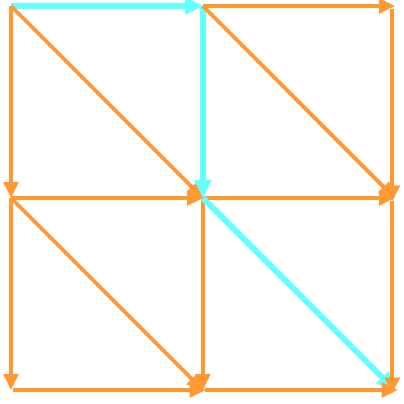
$$\begin{matrix} x_1 x_2 - & - \\ - & y_1 y_2 \end{matrix}$$



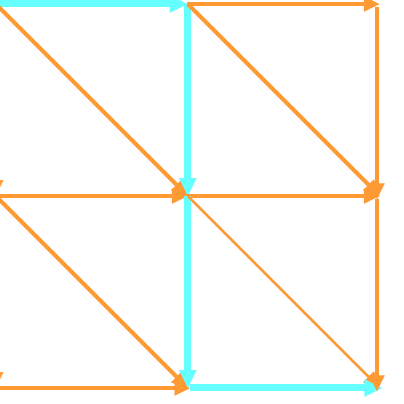
$$\begin{matrix} x_1 x_2 - \\ - & y_1 y_2 \end{matrix}$$



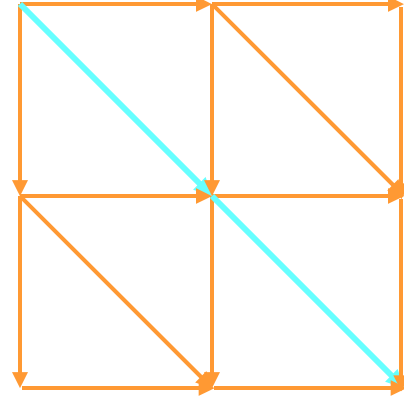
$$\begin{matrix} x_1 - & x_2 - \\ - & y_1 - y_2 \end{matrix}$$



$$\begin{matrix} x_1 - & x_2 \\ - & y_1 y_2 \end{matrix}$$

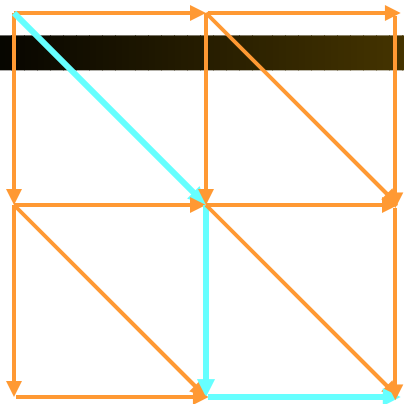


$$\begin{matrix} x_1 - & - x_2 \\ - & y_1 y_2 - \end{matrix}$$



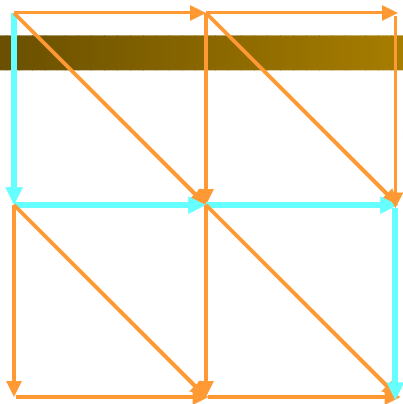
$$\begin{matrix} x_1 x_2 \\ y_1 y_2 \end{matrix}$$

# *Pfade in F*



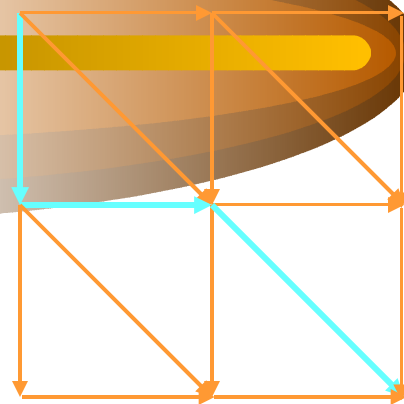
$$x_1 - x_2$$

$$y_1 y_2^-$$



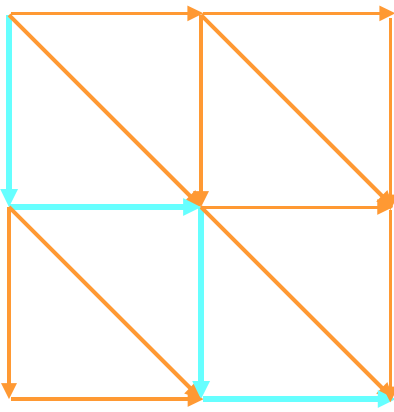
$$- x_1 x_2^-$$

$$y_1^- - y_2$$



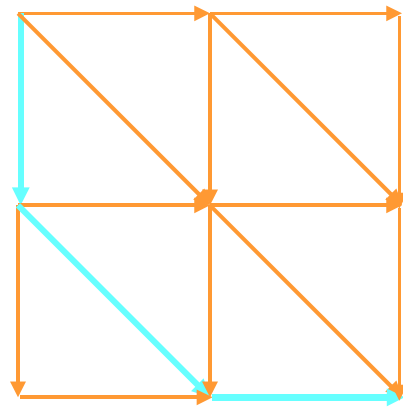
$$- x_1 x_2$$

$$y_1^- y_2$$



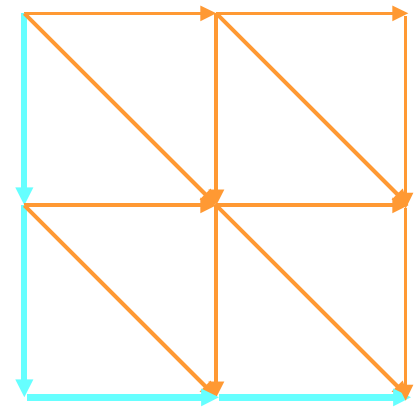
$$- x_1^- x_2$$

$$y_1^- y_2^-$$



$$- x_1 x_2$$

$$y_1 y_2^-$$



$$- - x_1 x_2$$

$$y_1 y_2^- -$$

# *Globales Alignment*

- Needleman-Wunsch Algorithmus
  - benutzt Prinzip der dynamischen Programmierung: optimales Alignment für zwei Sequenzen wird aus optimalen Alignments von Teilsequenzen bestimmt
  - kleinste Einheit: Alignment von zwei Buchstaben (Aminosäuren) bzw. Wert für eine Gap
  - 1. Schritt: Berechnung einer Matrix, die alle möglichen Alignments der Sequenzen repräsentiert. Mit Ausnahme der Initialwerte werden alle Einträge der Matrix mit Hilfe der bereits eingetragenen Werte und einer rekursiven Formel abgeleitet.
  - 2. Schritt: “Ablezen” des besten Alignments aus

# Algorithmus

Sequenz 1:  $x_1 x_2 x_3 \dots x_m$

Sequenz 2:  $y_1 y_2 y_3 \dots y_n$

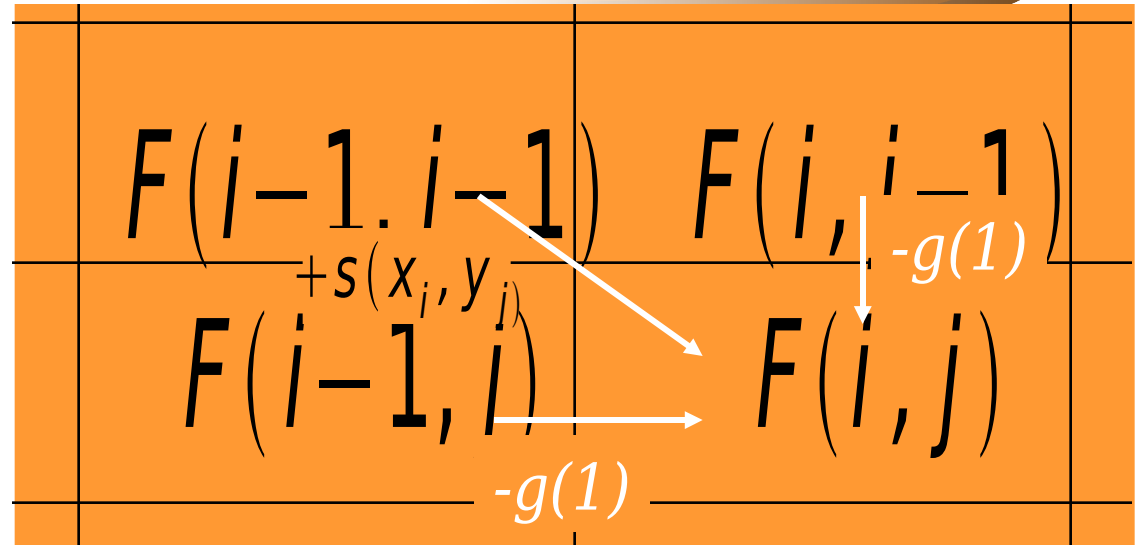
Matrix  $F$  wobei  $F(i, j)$  den Score für das optimale Alignment der Sequenz  $x_1 x_2 \dots x_j$  mit

der Sequenz  $y_1 y_2 \dots y_j$  angibt

| $F$      | - | $x_1$ | $x_2$ | $x_3$ | ... | $x_m$ |
|----------|---|-------|-------|-------|-----|-------|
| -        |   |       |       |       |     |       |
| $y_1$    |   |       |       |       |     |       |
| $y_2$    |   |       |       |       |     |       |
| $y_3$    |   |       |       |       |     |       |
| $\vdots$ |   |       |       |       |     |       |
| $y_n$    |   |       |       |       |     |       |

# Algorithmus

Zerlegungsprinzip:



$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) - g(1) \\ F(i, j-1) - g(1) \end{cases}$$



# *Needleman-Wunsch: Beispiel*



R D I S L V

R N I L V

- PAM 250 (\* 10)
- Lineare Gap Penalty mit  $q = -6$

# *Needleman-Wunsch: Beispiel*

A decorative graphic consisting of a horizontal arrow pointing to the right. The arrow has a black outline and a gradient fill transitioning from dark brown on the left to bright yellow on the right.

*R D I S L V*

0

*R*

*N*

*I*

*L*

*V*

# *Needleman-Wunsch: Beispiel*

|          | <i>R</i> | <i>D</i> | <i>I</i> | <i>S</i> | <i>L</i> | <i>V</i> |     |
|----------|----------|----------|----------|----------|----------|----------|-----|
|          | 0        | -6       | -12      | -18      | -24      | -30      | -36 |
| <i>R</i> | -6       |          |          |          |          |          |     |
| <i>N</i> | -12      |          |          |          |          |          |     |
| <i>I</i> | -18      |          |          |          |          |          |     |
| <i>L</i> | -24      |          |          |          |          |          |     |
| <i>V</i> | -30      |          |          |          |          |          |     |

# Needleman-Wunsch: Beispiel

|   |     | R   | D    | I    | S    | L    | V    |
|---|-----|-----|------|------|------|------|------|
|   | 0   | ←-6 | ←-12 | ←-18 | ←-24 | ←-30 | ←-36 |
| R | -6  | 6   | ←0   | ←-6  | ←-12 | ←-18 | ←-24 |
| N | -12 | 0   | 8    | ←-2  | ←-4  | ←-10 | ←-16 |
| I | -18 | -6  | 2    | 13   | ←-7  | ←-1  | ←-5  |
| L | -24 | -12 | -4   | 7    | 10   | 13   | ←-7  |
| V | -30 | -18 | -10  | 1    | 6    | 12   | 17   |

# *Needleman-Wunsch: Beispiel*



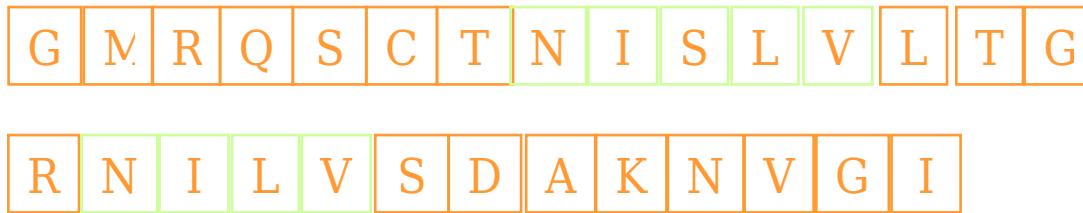
- Ergebnis-Score: 17
- Alignment

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| R | D | I | S | L | V |
|---|---|---|---|---|---|

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| R | N | I | - | L | V |
|---|---|---|---|---|---|

# Optimales Lokales Alignment

- Suche nach lokal bester Übereinstimmung



- Smith-Waterman Algorithmus: wie Needleman-Wunsch, aber
  - Rücksetzen auf den Wert 0, falls  $\text{Score} < 0$
  - Alignment beginnt in Matrix bei einer 0 und endet in Matrix bei maximalem Wert

# *Multiples Alignment*

- Vergleich von mehreren ( $> 2$ ) Sequenzen
- Motivation: Bessere Erkennung von Motiven
  - weniger gut erhaltene Motive gehen im paarweisen Vergleich unter
  - unwichtige Bereiche werden besser ausgeblendet
- Basis für Strukturvorhersage
- Zwei wesentliche Ansätze:
  - simultanes multiples Alignment

# Beispiel

## 8 Fragmente aus Immunoglobulin Sequenzen

VTISCTGSSSNIGAG-NHVKWYQQLPG  
VTISCTGTSSNIGS--ITVNWYQQLPG  
LRLSCSSSGFIFSS--YAMYWVRQAPG  
LSLTCTVSGTSFDD--YYSTWVRQPPG  
PEVTCVVVDVSHEDPQVKFNWYVDG--  
ATLVCLISDFYPGA--VTVAWKADS--  
AALGCLVKDYFPEP--VTVSWNSG---  
VSLTCLVKGFYPSD--IAVEWESNG--

Konservierte Residuen: W, C

Konservierte Regionen: Q.PG

Auffallende Muster: Hydrophobe Residuen: V,L,P,A,I