

Was ist Bioinformatik?

Alexander Vilbig

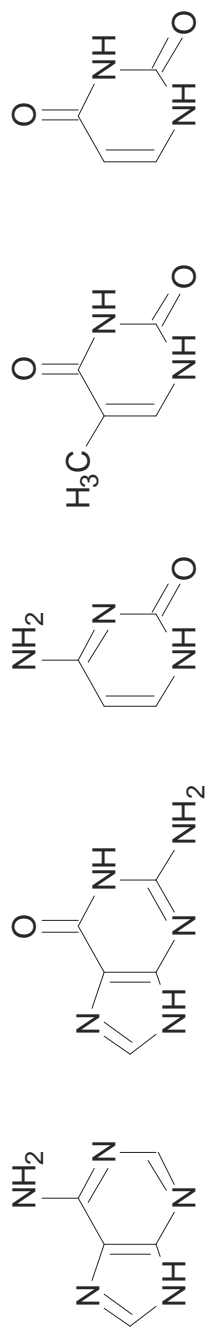
23. Oktober, 2001



Übersicht

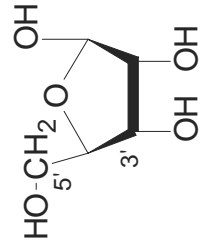
- Biologische Grundlagen
 - DNA, RNA, Gen, Protein, ...
- Informatik in der Biologie
 - Sequenz-Alignment
 - Molekulardynamik-Simulation
- Biologie in der Informatik
 - Evolutionäre Heuristiken
- Zusammenfassung & Ausblick
- Literatur und Ressourcen

Nukleinsäuren I

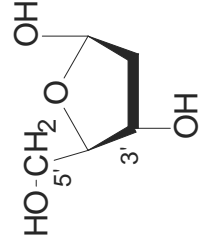


Purin-Basen

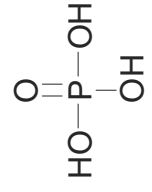
Pyrimidin-Basen



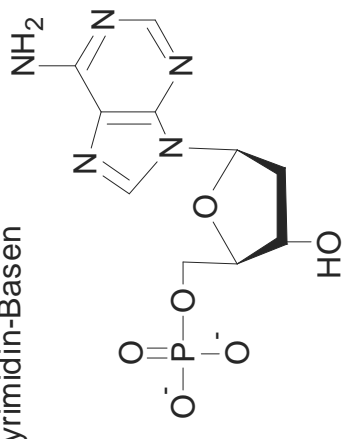
Ribose



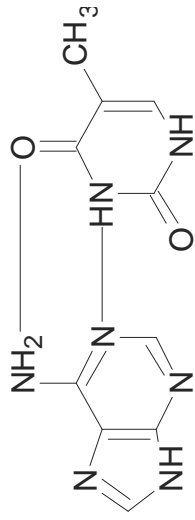
Desoxyribose



Phosphat

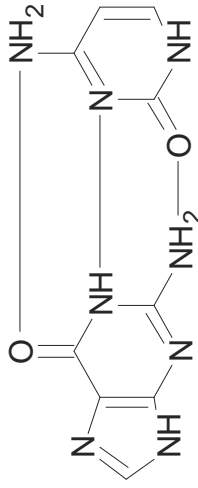


Nukleinsäuren II



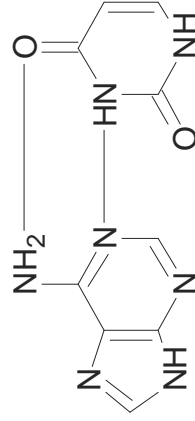
Adenin

Thymin



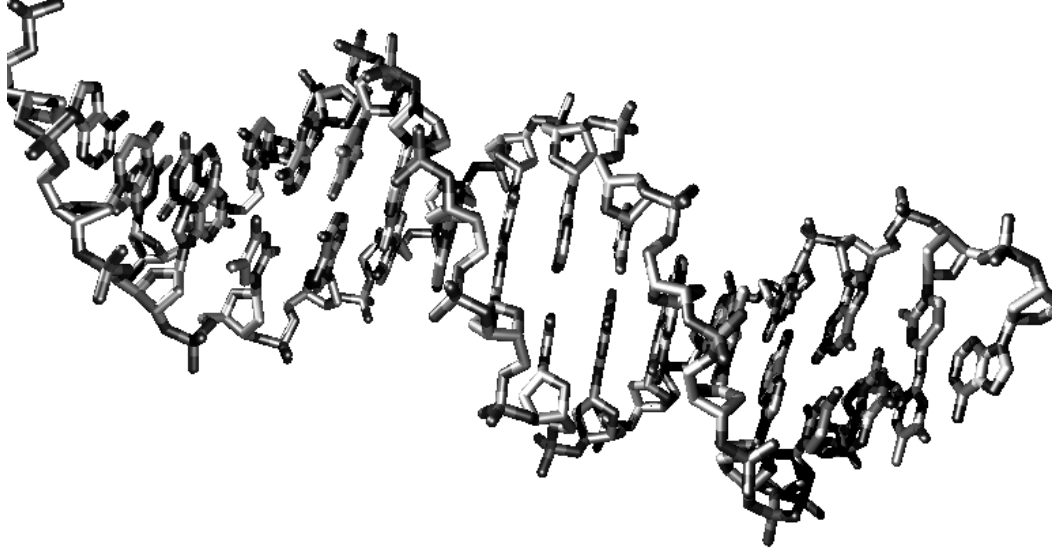
Guanin

Cytosin



Adenin

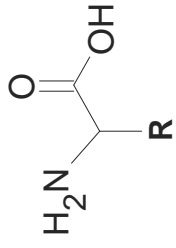
Uracil



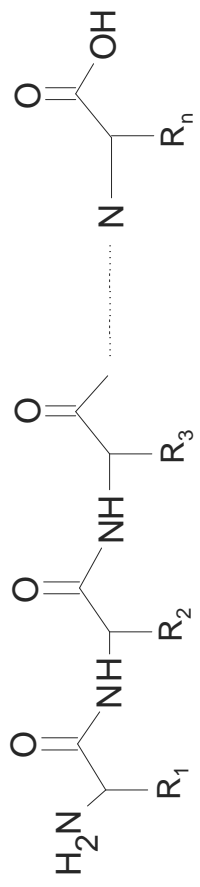
roteine



□ roteine

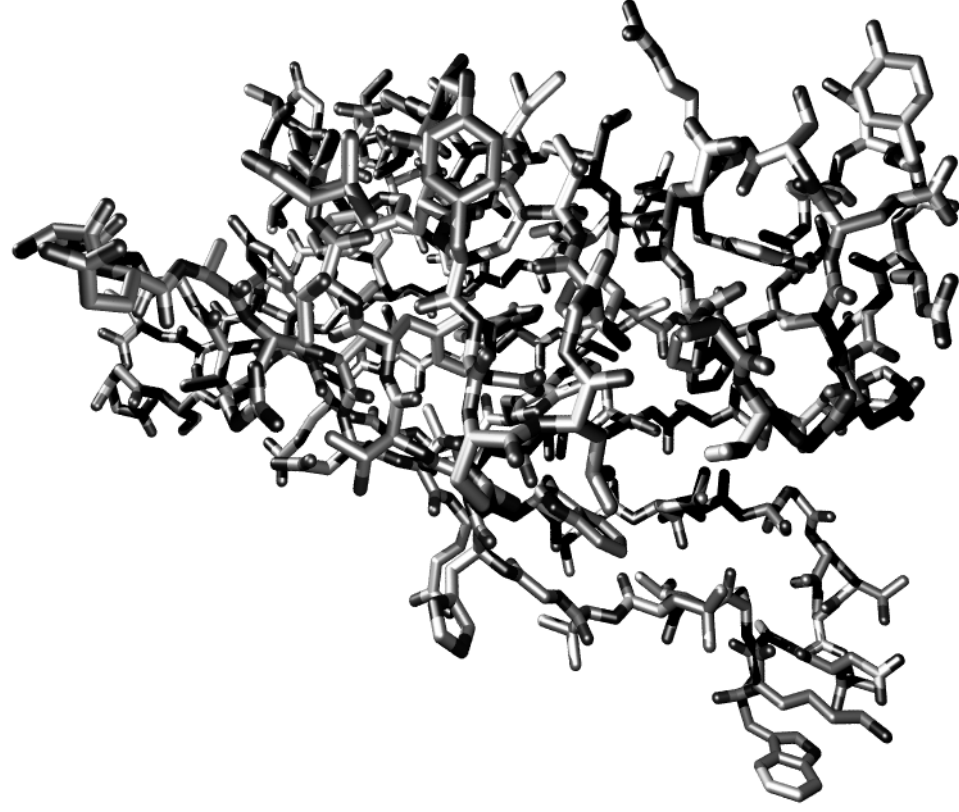


Aminosäure

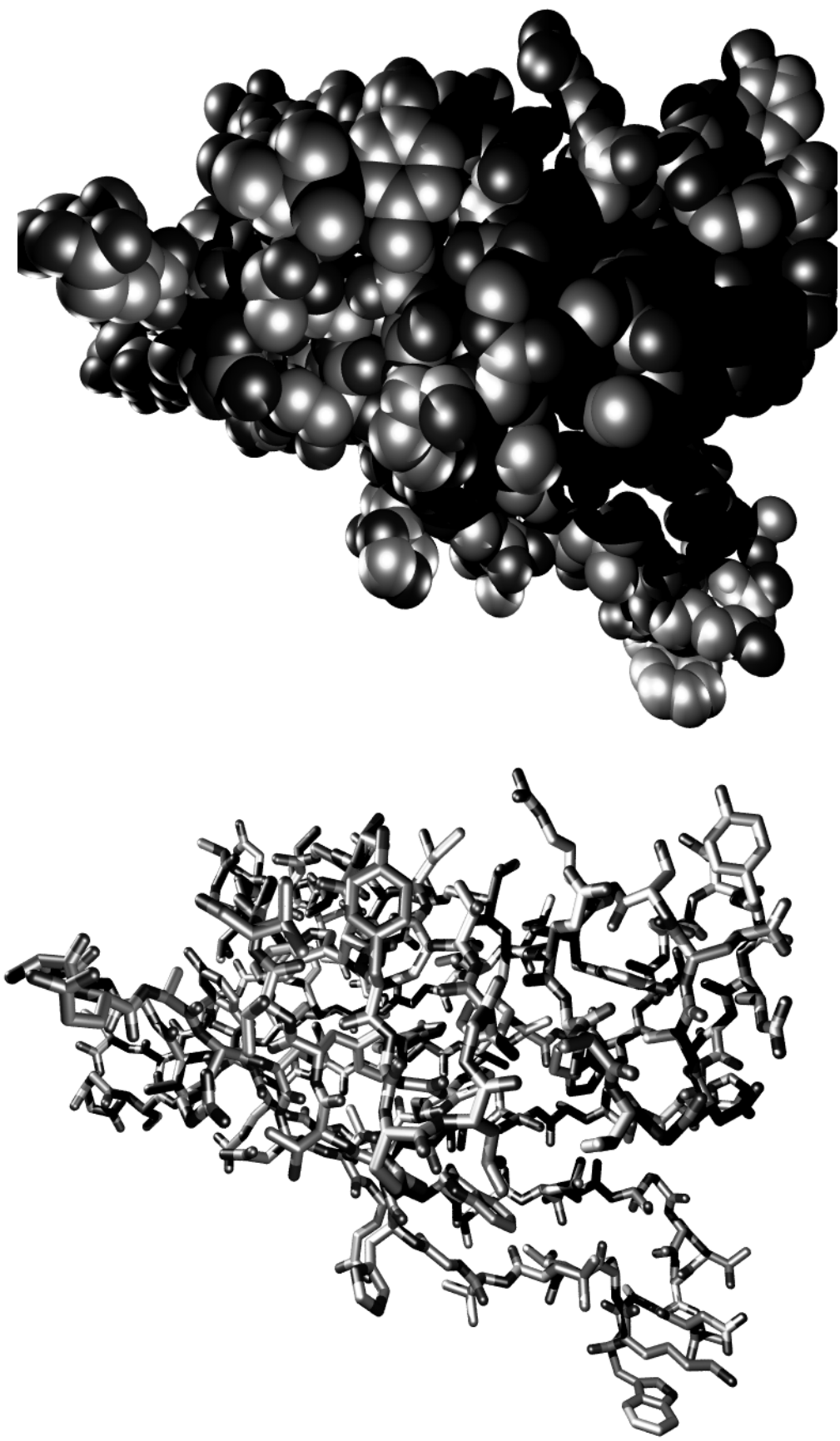


Polypeptid

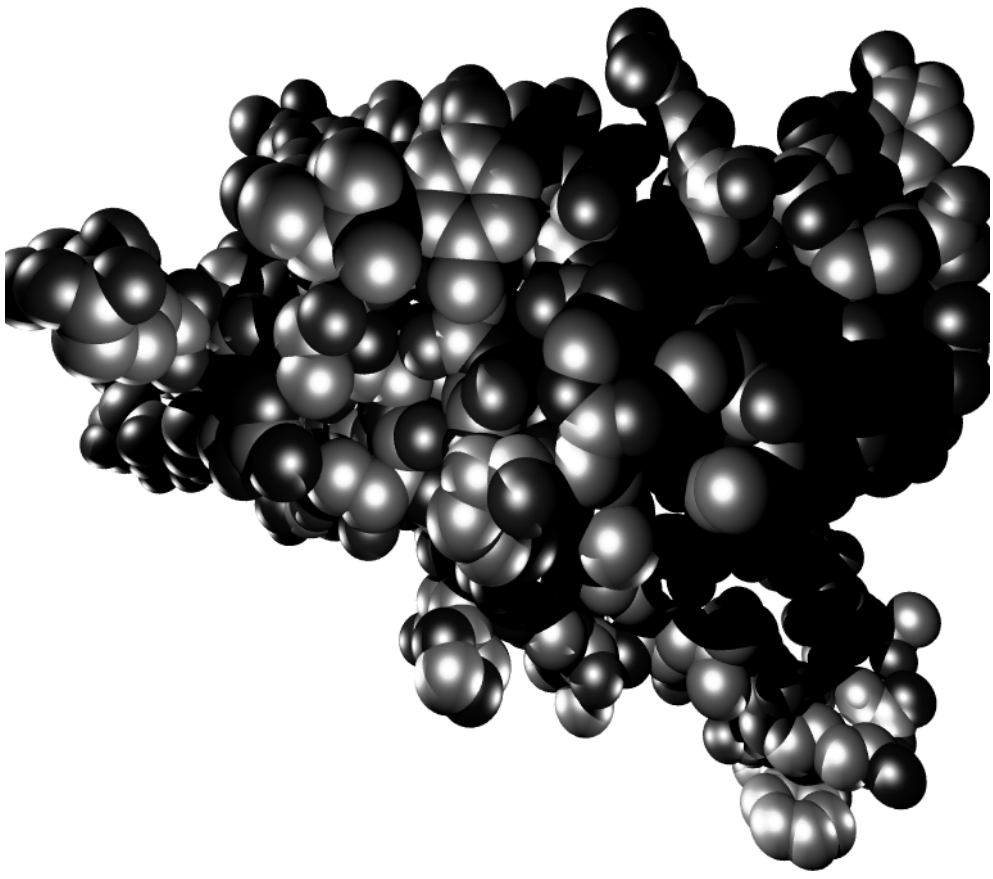
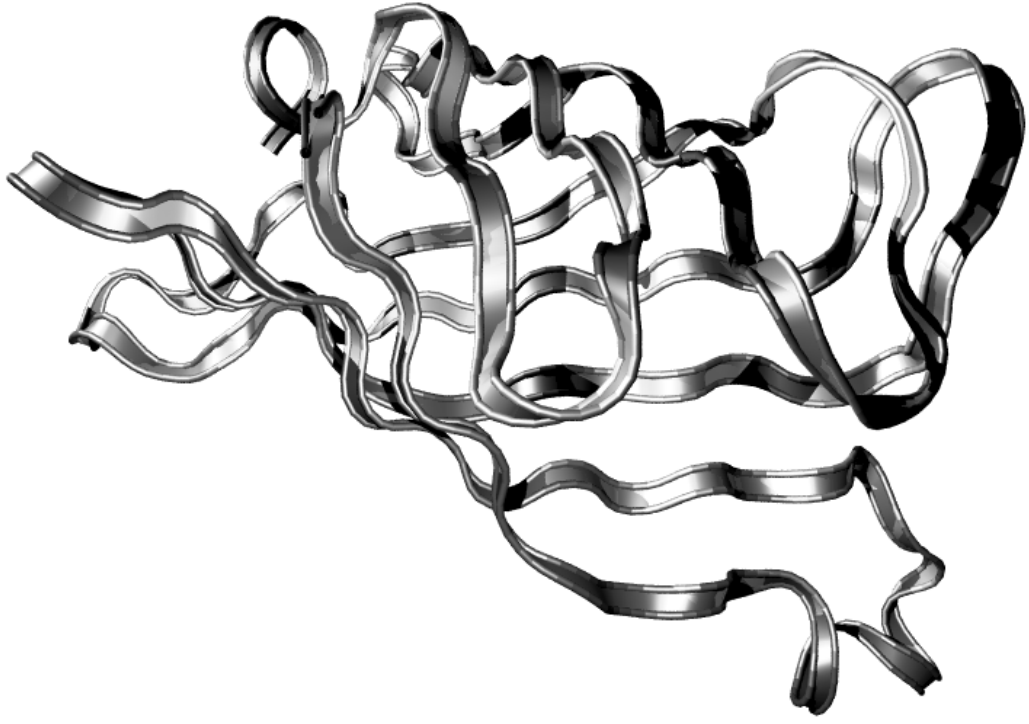
□ roteine



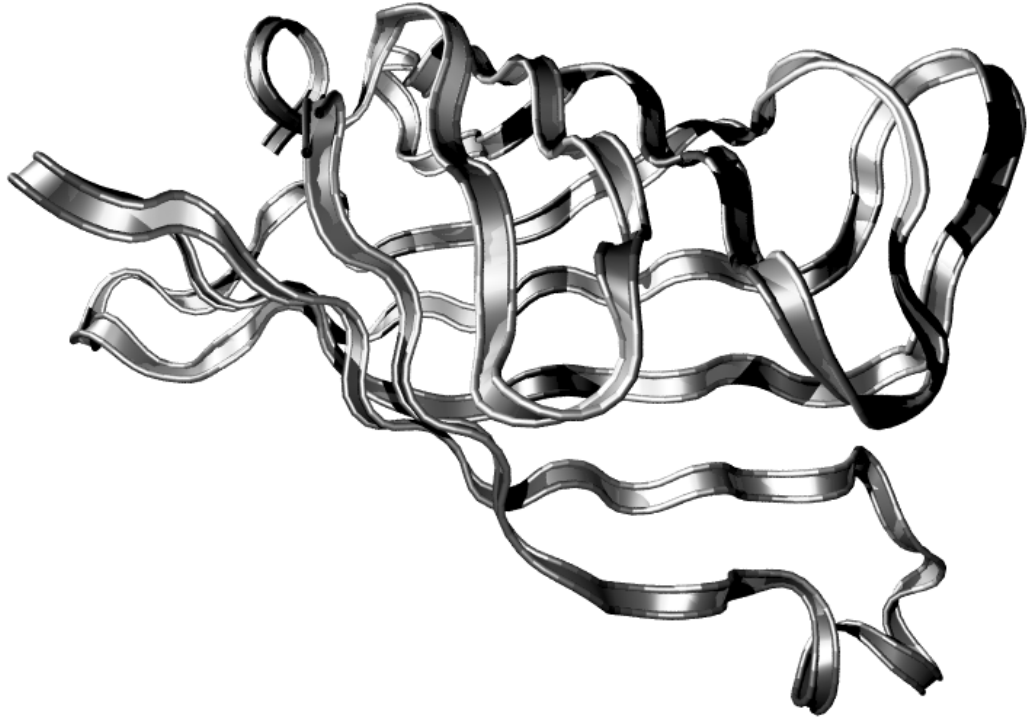
□ roteine



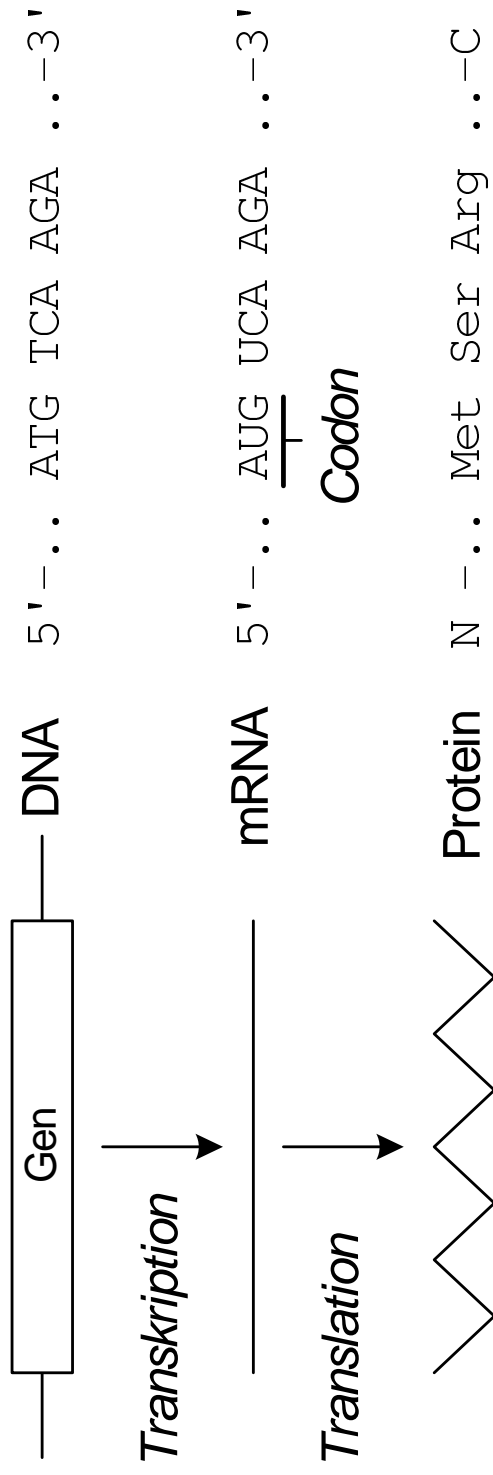
□ roteine



□ roteine



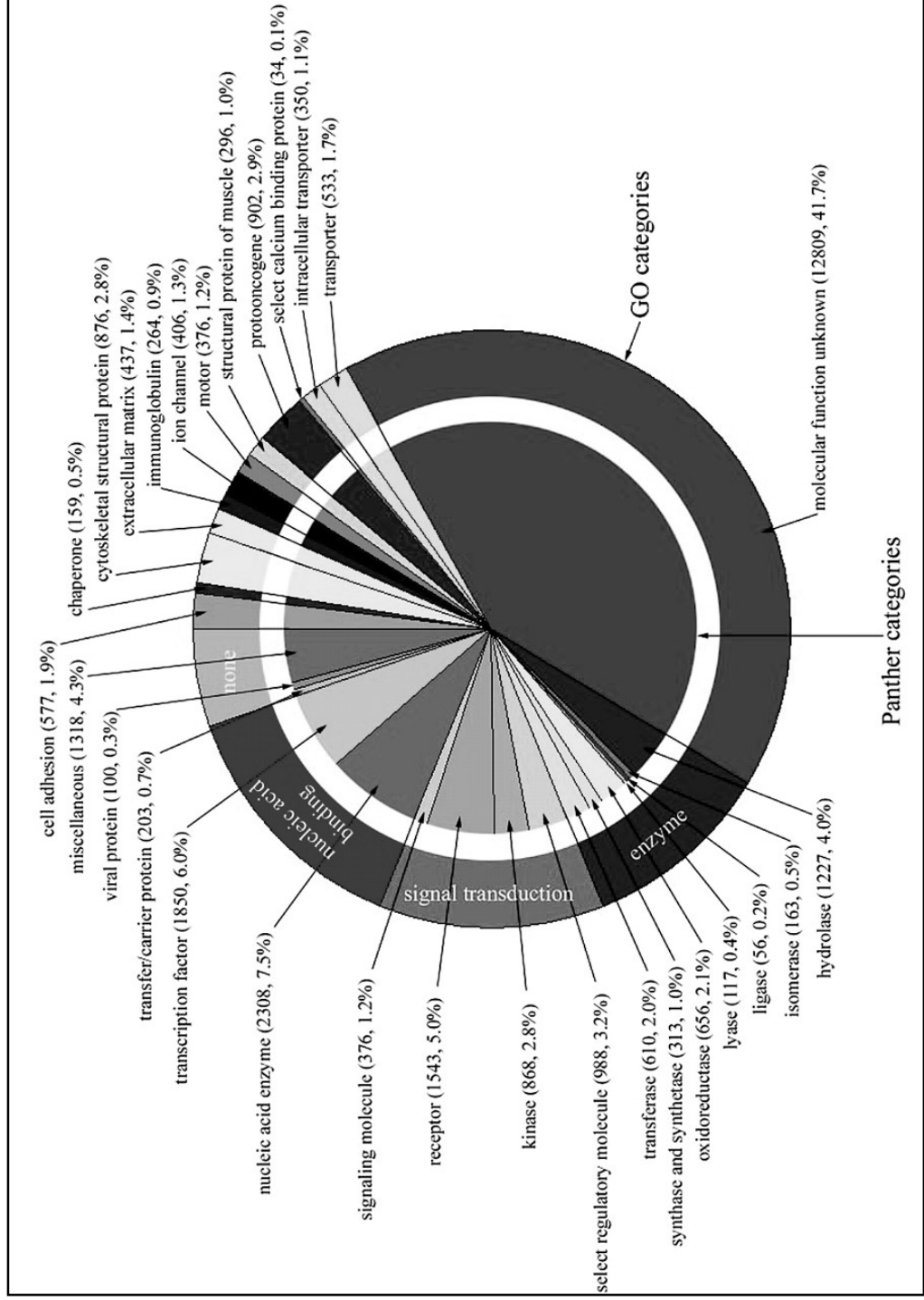
□ □ □ **ression**



■ Was menschliche Genom

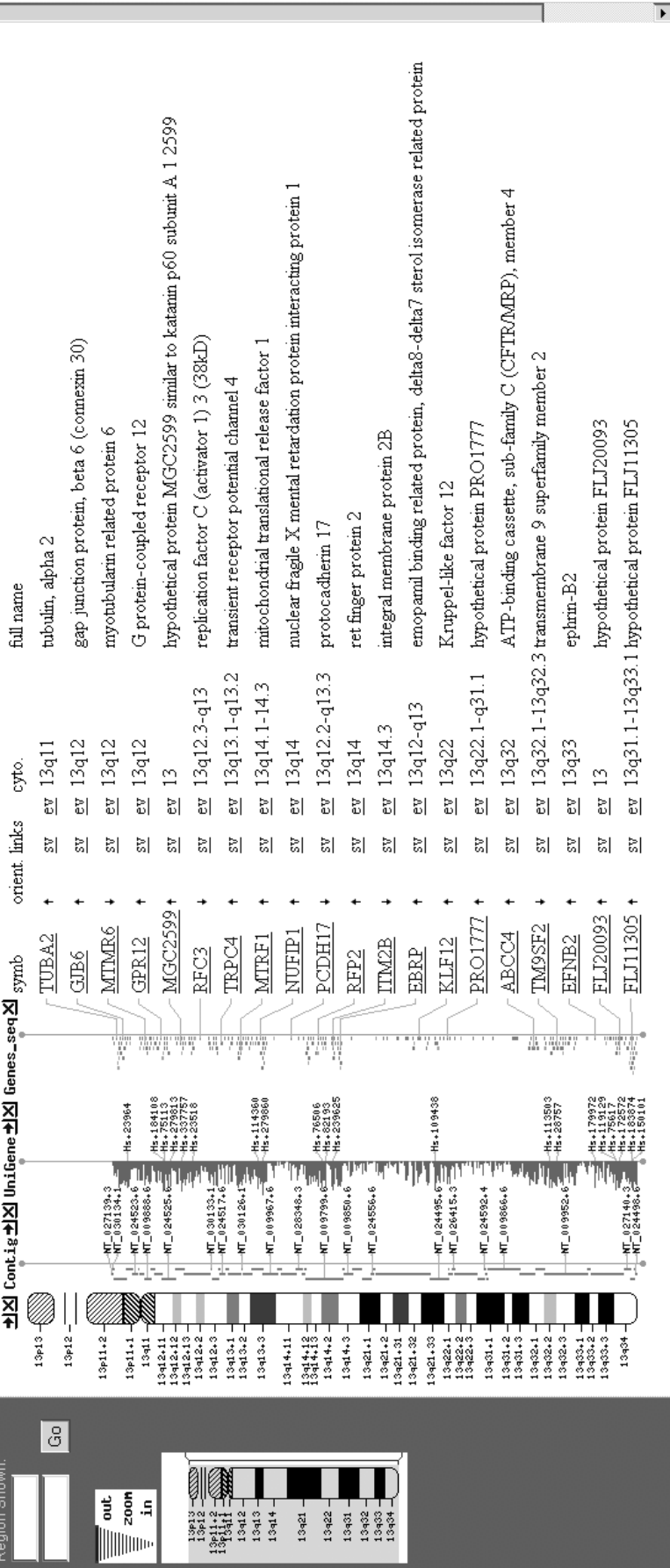
- Umfasst ca. 3,2 Milliarden Basenpaare
- Geschätzte 25.000 Gene
- Physisch verteilt auf 23 Chromosomen
- Sequenzierung in 2003 praktisch abgeschlossen
 - Human Genome Project öffentlich
 - Celera Genomics Inc. kommerziell

□ as menschliche □ enom



Homo sapiens Map View build 26
 Chromosome: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 X Y
 Master: Genes On Sequence Map [Display settings](#)

Total Genes On Chromosome: 326 [13 not localized]
 Region Displayed: 0-1166M bp [Download/View Sequence/Evidence](#)
 Genes Labeled: 20 Total Genes in Region: 313



Informatik in der Biologie

- Datenerfassung und -verwaltung
 - Steuerung von Laborgeräten
 - Labormanagement, Sequence Assembly, ...
 - Datenbanken, Queries, Präsentation, API-Anbindung, ...
- Datenanalyse
 - Sequenzbasiert Alignment, Phylogenie, ...
 - Strukturbasiert Molecular Docking, Structure Prediction, ...
 - Visualisierung
- Simulation theoretischer Modelle
 - Biochemische Reaktionen
 - Metabolismus
 - Ontogenese, ...

□ **n** □ **en** □ **un** □ **s** **beis** □ **iel** □ □ □ **e** □ **uen** □ □ □ **li** □ **nment**

- □ **Bestm** □ **gliche** □ **Ausrichtung biomolekularer Sequenzen**
 - **I** **Insertion, Deletion, Substitution, Match**

A = TTAAGCACATTGGATTCTCGGCAGCGTGG

B = TTAAGGATTCTCCGTAG

Genomsequenzierung

- Bestmögliche Ausrichtung biomolekularer Sequenzen
 - Insertion, Deletion, Substitution, Match

A = TTAAGCAC TTGGATTCTCGGCAGCGTGG
B = TTAAGGATTCTCCGTAG

Global

A = TTAA-GCACTTGGATTCTCGGCAGCGTGG
B = TAAA-----GGATTCT---C--CGTAG

Alignment

Bestmögliche Ausrichtung biomolekularer Sequenzen

- Insertion, Deletion, Substitution, Match

A = TTAAGCACTTGGATTCTCGGCAGCGTGG
B = TTAAGGATTCTCCGTAG

Global

A TTAAGCACTTGGATTCTCGGCAGCGTGG
B TAAA-----GGATTCT--C--CGTAG

Semi-global

A TTAAGCACTTTT---GGATTCTCGGCAG
B -----TTAAAGGATTCTCCGTAG

Alignment

- Bestmögliche Ausrichtung biomolekularer Sequenzen
 - ▮ Insertion, Deletion, Substitution, Match

A = TTAAGCACTTGGATTCTCGGCAGCGTGG
B = TTAAGGATTCTCCGTAG

Global

A TTAAGCACTTGGATTCTCGGCAGCGTGG
B TAAA-----GGATTCT--C--CGTAG

Semi-global

A TTAAGCACTTT---GGATTCTCGGCAG
B -----TTAAAGGATTCTCCGTAG

mit Gaps

A TTAAGCACTTGGATTCTCGGCAGCGTGG
B TTAAG-----GATTCTC-----CGTAG

Berechnung eines globalen Alignments I

$$\text{Distanz } d(A, B) = \min_{i=1}^k \sum_{i=1}^k d(a_i, b_i) \quad d(a_i, b_j) = \begin{cases} 0, & a_i = b_j \\ > 0, & \text{sonst} \end{cases}$$

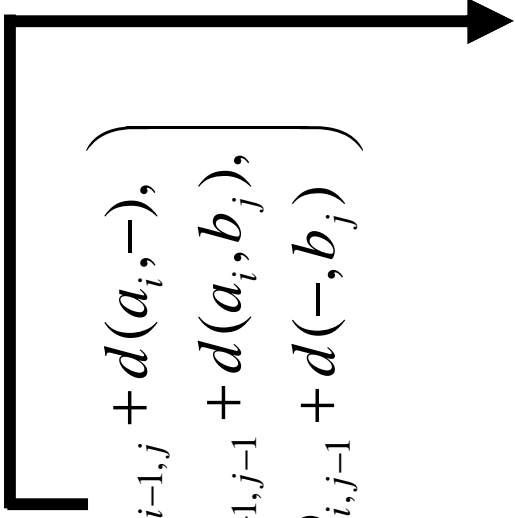
$$D_{i,j} = \min \begin{pmatrix} D_{i-1,j} + d(a_i, -), \\ D_{i-1,j-1} + d(a_i, b_j), \\ D_{i,j-1} + d(-, b_j) \end{pmatrix}$$

$$D_{i-1,j} = \min \begin{pmatrix} D_{i-2,j} + d(a_{i-1}, -), \\ D_{i-2,j-1} + d(a_{i-1}, b_j), \\ D_{i-1,j-1} + d(-, b_j) \end{pmatrix}$$

Berechnung eines globalen Alignments I

$$\text{Distanz } d(A, B) = \min_{i=1}^k \sum_{i=1}^k d(a_i, b_i) \quad d(a_i, b_j) = \begin{cases} 0, & a_i = b_j \\ > 0, & \text{sonst} \end{cases}$$

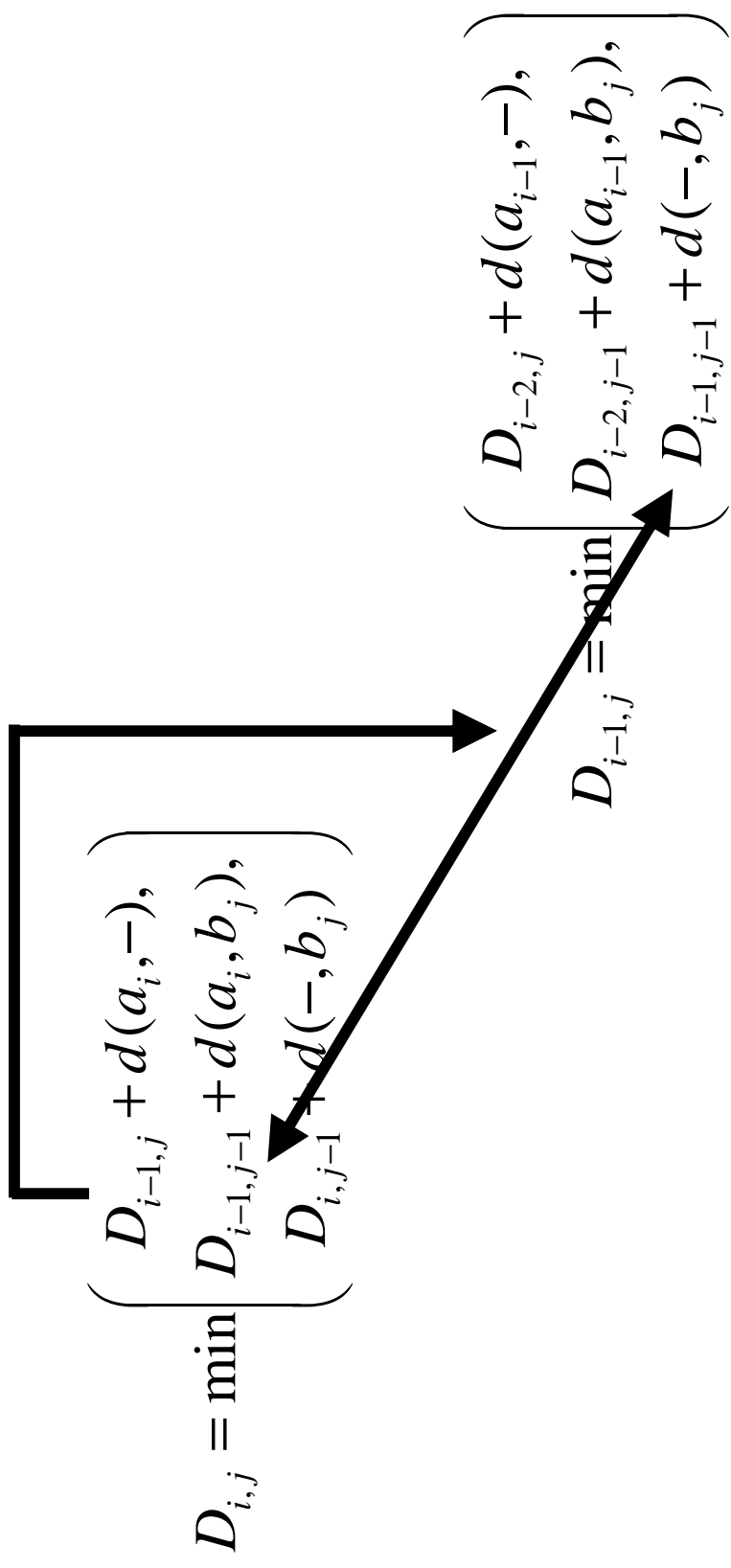
$$D_{i,j} = \min \begin{pmatrix} D_{i-1,j} + d(a_i, -), \\ D_{i-1,j-1} + d(a_i, b_j), \\ D_{i,j-1} + d(-, b_j) \end{pmatrix}$$



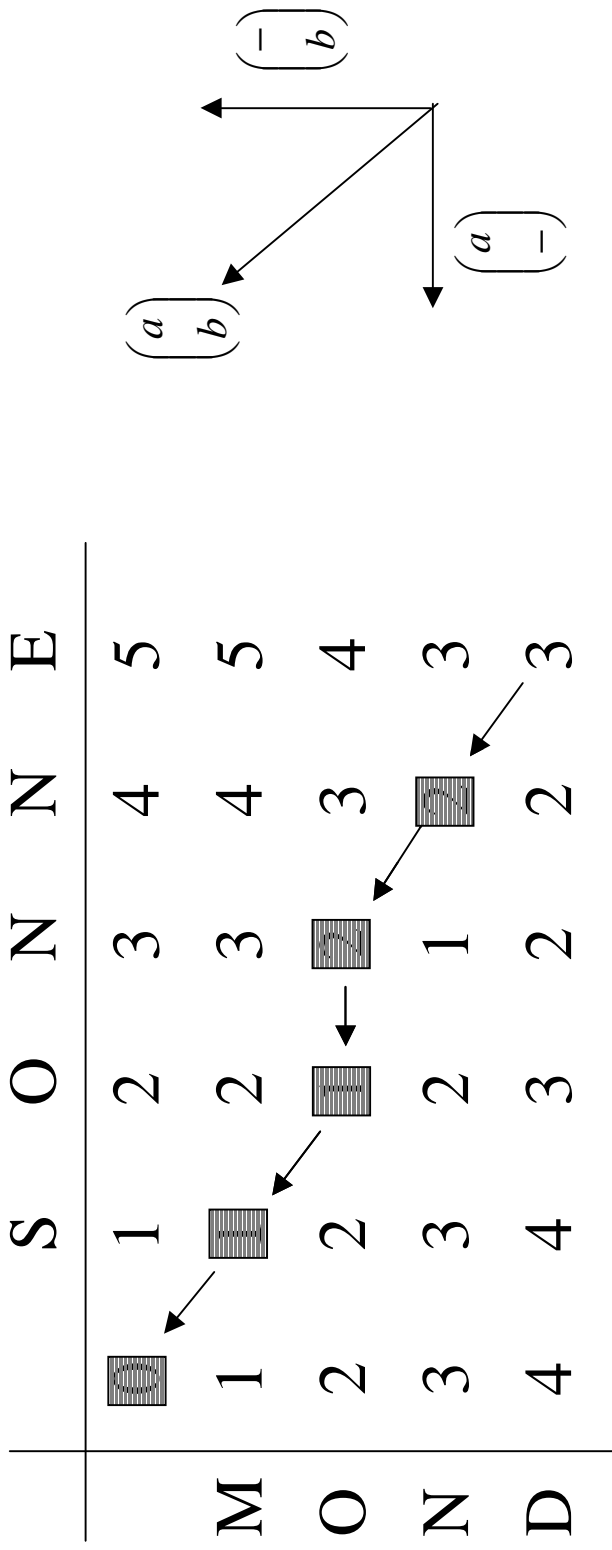
$$D_{i-1,j} = \min \begin{pmatrix} D_{i-2,j} + d(a_{i-1}, -), \\ D_{i-2,j-1} + d(a_{i-1}, b_j), \\ D_{i-1,j-1} + d(-, b_j) \end{pmatrix}$$

Berechnung eines globalen Alignments I

$$\text{Distanz } d(A, B) = \min_{i=1}^k \sum_{i=1}^k d(a_i, b_i) \quad d(a_i, b_j) = \begin{cases} 0, & a_i = b_j \\ > 0, & \text{sonst} \end{cases}$$



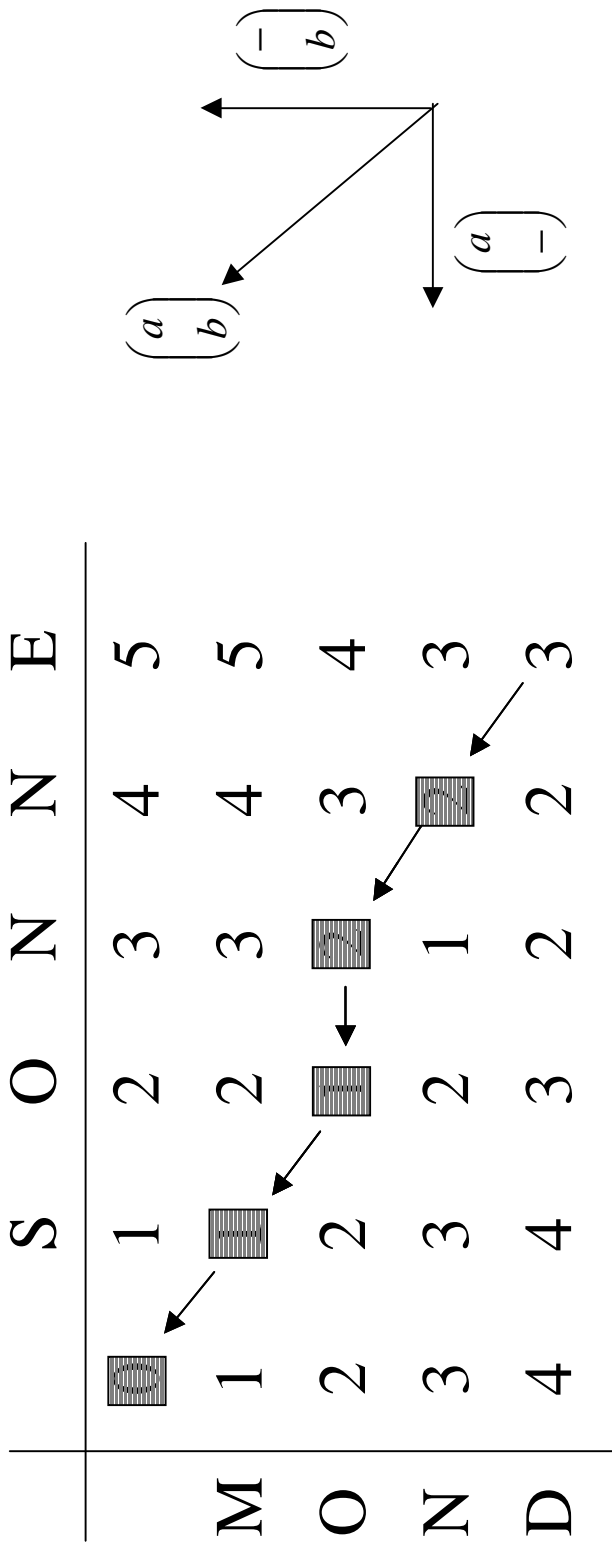
Berechnung eines globalen Alignments II



s o n n e

m o - n d

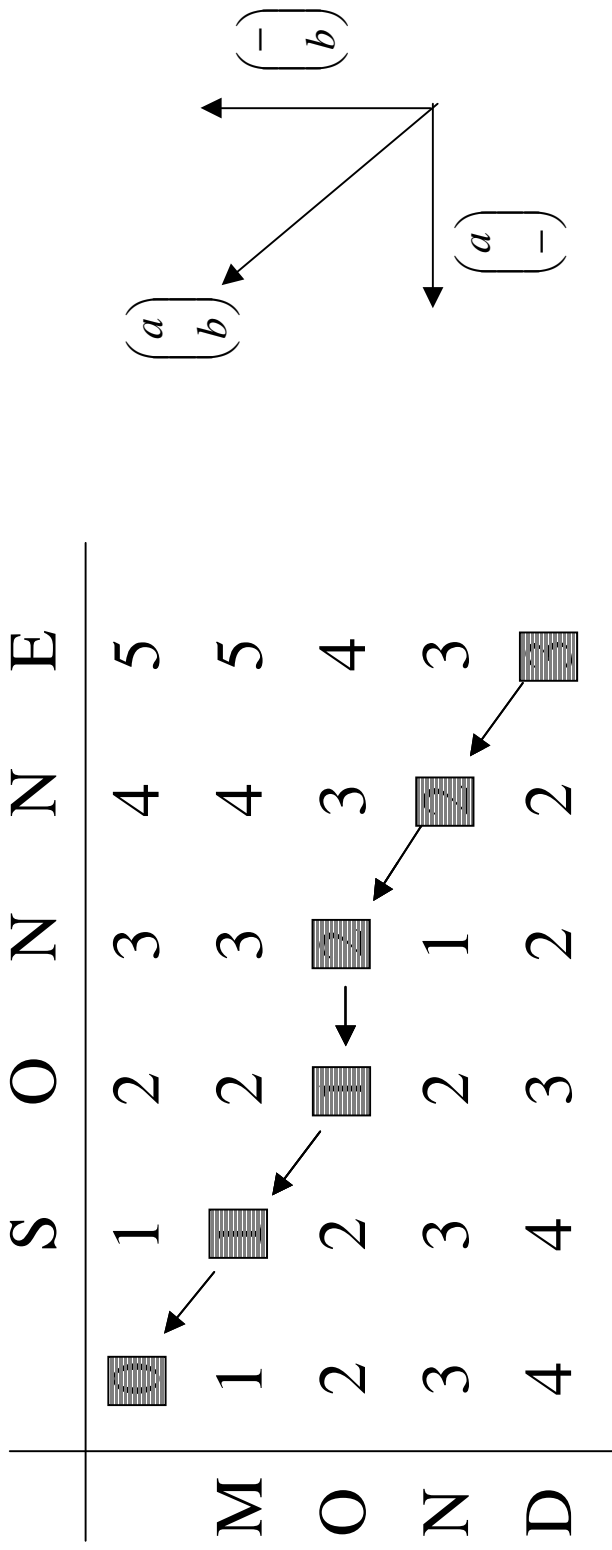
Berechnung eines globalen Alignments II



s o n n e

m o - n d

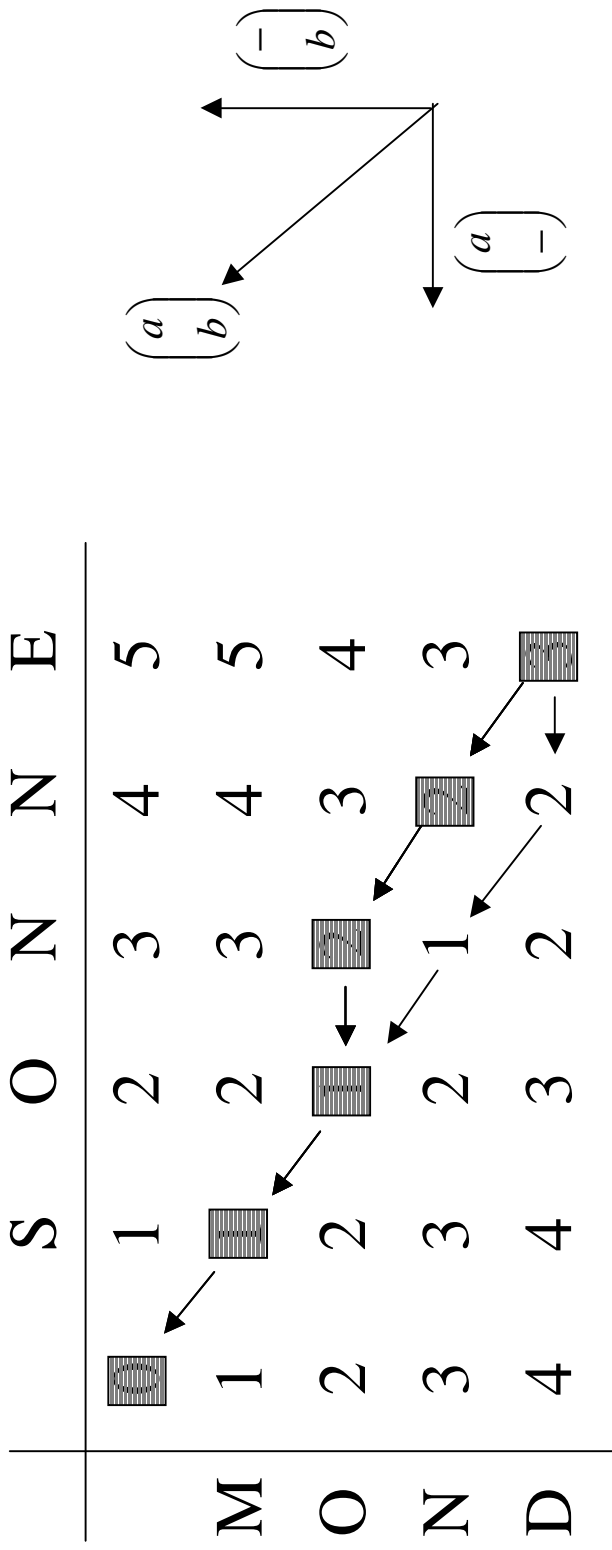
Berechnung eines globalen Alignments II



s o n n e

m o - n d

Berechnung eines globalen Alignments II



s o n n e sonne

m o - n d mond-

□ **n** □ **en** □ **un** □ **s** □ **beis** □ **iel** □ □

□ **olekular** □ □ **namik** □ □ **imulation**

- Simulation des Erhaltens molekularer Systeme
 - Detaillierungsgrad □ Individuelle Atome oder funktionale Gruppen
 - Systemgröße □ □ □ □ - □ □ □ □ Atome □ Lösungsmittelatome
 - Zeiträume □ □ □ □ □ - □ □ □ □ s □ Pico- bzw. Nanosekunden □
- Zielsetzung
 - Untersuchung von Bindungsvorgängen
 - Berechnung thermodynamischer Größen
 - Drug Design, ...
- Theoretisches Modell
 - Klassische Mechanik vs. Quantenmechanik
 - Newton'sche Bewegungsgleichungen für große Systeme

■ molekularmechanische Simulation I

- Algorithmus von Verlet
 - Integration der Bewegungsgleichungen
 - Kombination der Taylor-Entwicklung an den Stellen $t+\Delta t$ und $t-\Delta t$

$$\frac{\partial \mathbf{r}_i(t)}{\partial t} = \mathbf{v}_i(t)$$

$$\frac{\partial \mathbf{v}_i(t)}{\partial t} = \frac{\mathbf{F}_i(t)}{m_i}$$

Molekulardynamik-Simulation I

- Algorithmus von Verlet
- Integration der Bewegungsgleichungen
- Kombination der Taylor-Entwicklung an den Stellen $t+\Delta t$ und $t-\Delta t$

$$\frac{\partial \mathbf{r}_i(t)}{\partial t} = \mathbf{v}_i(t) \quad \frac{\partial \mathbf{v}_i(t)}{\partial t} = \frac{\mathbf{F}_i(t)}{m_i}$$

$$\mathbf{r}_i(t + \Delta t) = 2\mathbf{r}_i(t) - \mathbf{r}_i(t - \Delta t) + \frac{\mathbf{F}_i(t)}{m_i} \Delta t^2 + O(\Delta t^4)$$

Molekulardynamik-Simulation I

- Algorithmus von Verlet
- Integration der Bewegungsgleichungen
- Kombination der Taylor-Entwicklung an den Stellen $t+\Delta t$ und $t-\Delta t$

$$\frac{\partial \mathbf{r}_i(t)}{\partial t} = \mathbf{v}_i(t) \quad \frac{\partial \mathbf{v}_i(t)}{\partial t} = \frac{\mathbf{F}_i(t)}{m_i}$$

$$\mathbf{r}_i(t + \Delta t) = 2\mathbf{r}_i(t) - \mathbf{r}_i(t - \Delta t) + \frac{\mathbf{F}_i(t)}{m_i} \Delta t^2 + O(\Delta t^4)$$

$$\mathbf{F}_i(t) = - \frac{\partial U[\mathbf{r}_1(t), \mathbf{r}_2(t), \dots, \mathbf{r}_N(t)]}{\partial \mathbf{r}_i(t)}$$

Molekulardynamik-Simulation I

- Algorithmus von Verlet
- Integration der Bewegungsgleichungen
- Kombination der Taylor-Entwicklung an den Stellen $t+\Delta t$ und $t-\Delta t$

$$\frac{\partial \mathbf{r}_i(t)}{\partial t} = \mathbf{v}_i(t) \quad \frac{\partial \mathbf{v}_i(t)}{\partial t} = \frac{\mathbf{F}_i(t)}{m_i}$$

$$\mathbf{r}_i(t + \Delta t) = 2\mathbf{r}_i(t) - \mathbf{r}_i(t - \Delta t) + \frac{\mathbf{F}_i(t)}{m_i} \Delta t^2 + O(\Delta t^4)$$

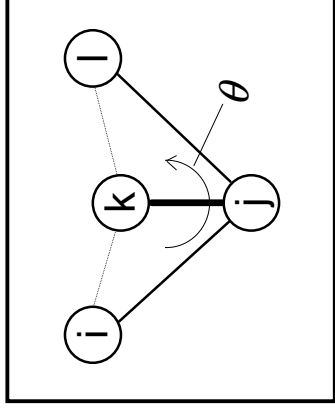
$$\mathbf{F}_i(t) = - \frac{\partial U[\mathbf{r}_1(t), \mathbf{r}_2(t), \dots, \mathbf{r}_N(t)]}{\partial \mathbf{r}_i(t)}$$

$$U(\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_N) = E_B + E_\theta + E_\varphi + E_\psi + E_{VDW} + E_C$$

Molekulardynamik-Simulation I

- Algorithmus von Verlet
- Integration der Bewegungsgleichungen
- Kombination der Taylor-Entwicklung an den Stellen $t+\Delta t$ und $t-\Delta t$

$$\frac{\partial \mathbf{r}_i(t)}{\partial t} = \mathbf{v}_i(t) \quad \frac{\partial \mathbf{v}_i(t)}{\partial t} = \frac{\mathbf{F}_i(t)}{m_i}$$



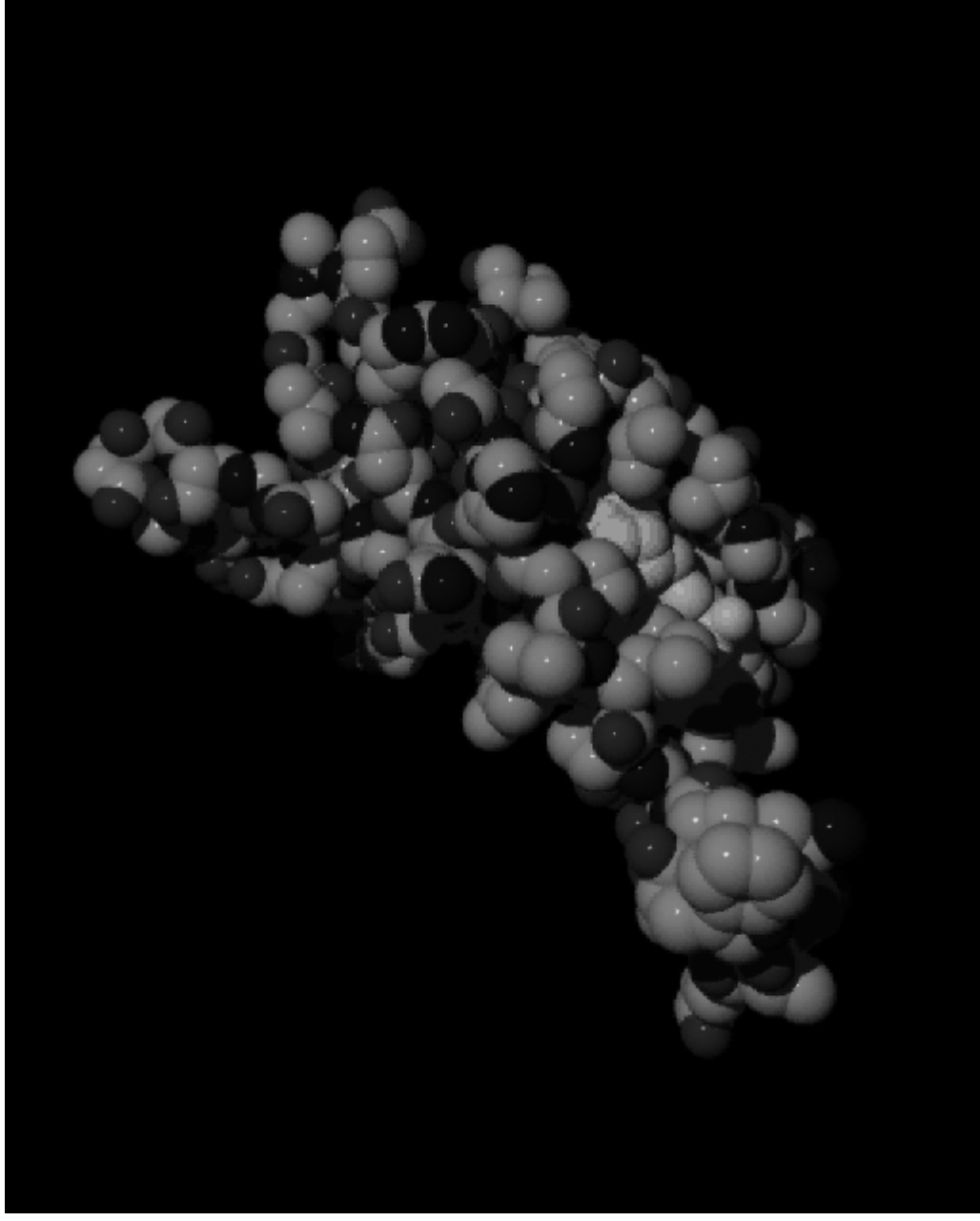
$$\mathbf{r}_i(t + \Delta t) = 2\mathbf{r}_i(t) - \mathbf{r}_i(t - \Delta t) + \frac{\mathbf{F}_i(t)}{m_i} \Delta t^2 + O(\Delta t^4)$$

$$\mathbf{F}_i(t) = - \frac{\partial U[\mathbf{r}_1(t), \mathbf{r}_2(t), \dots, \mathbf{r}_N(t)]}{\partial \mathbf{r}_i(t)}$$

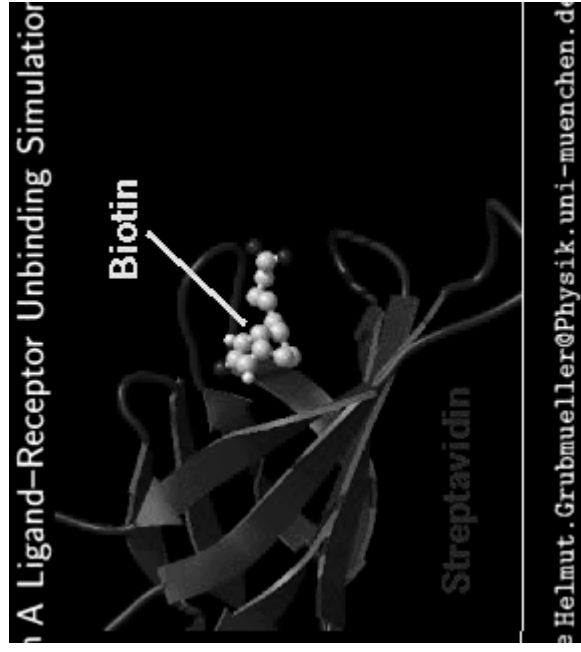
$$U(\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_N) = E_B + E_\theta + E_\varphi + E_\psi + E_{VDW} + E_C$$

Molekulardynamik-Simulation II

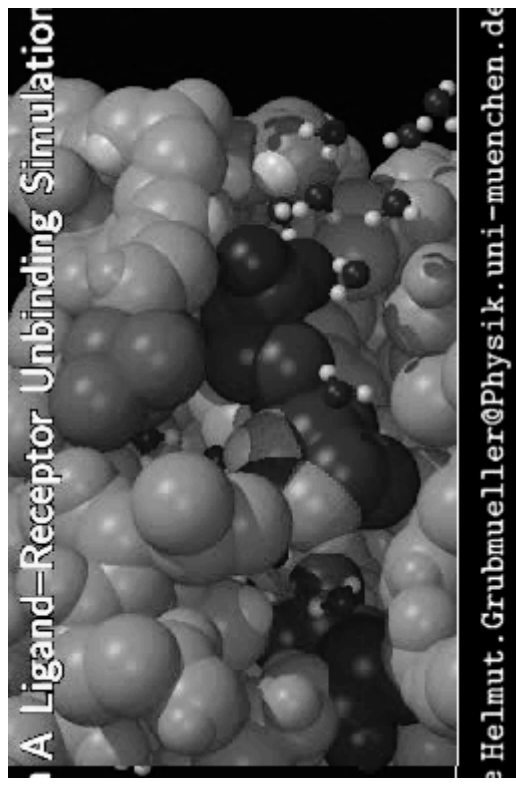
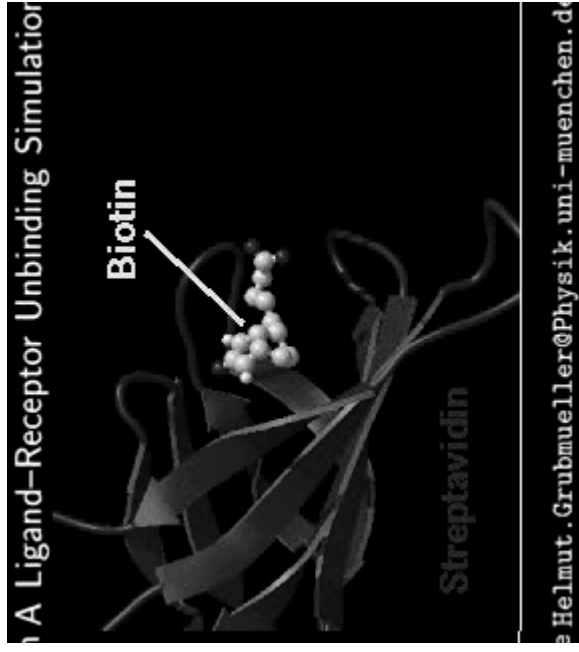
Molekulardynamik-Simulation II



Molekulardynamik-Simulation II

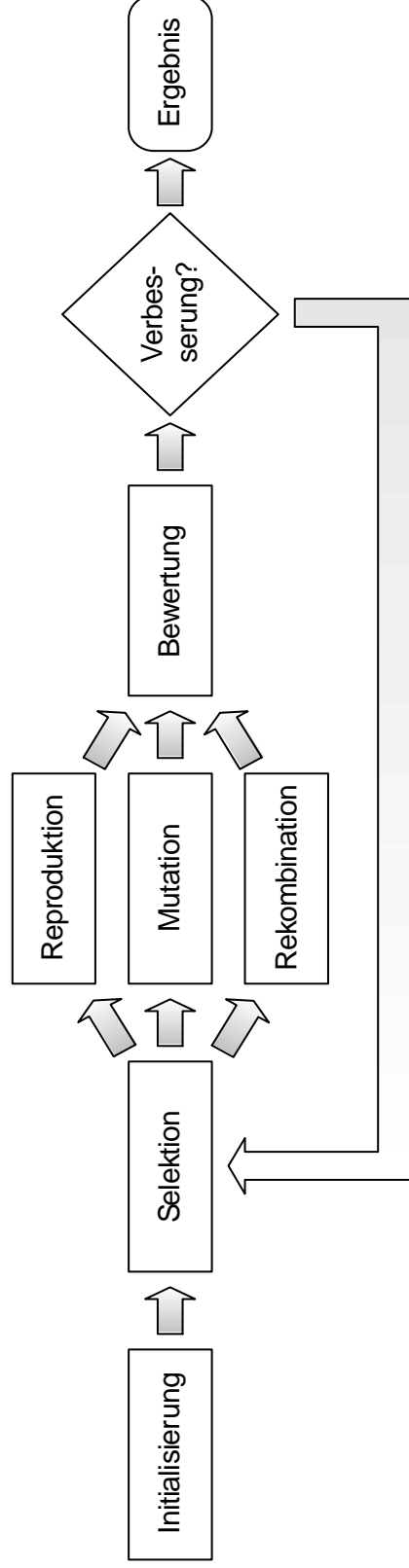


Molekulardynamik-Simulation II

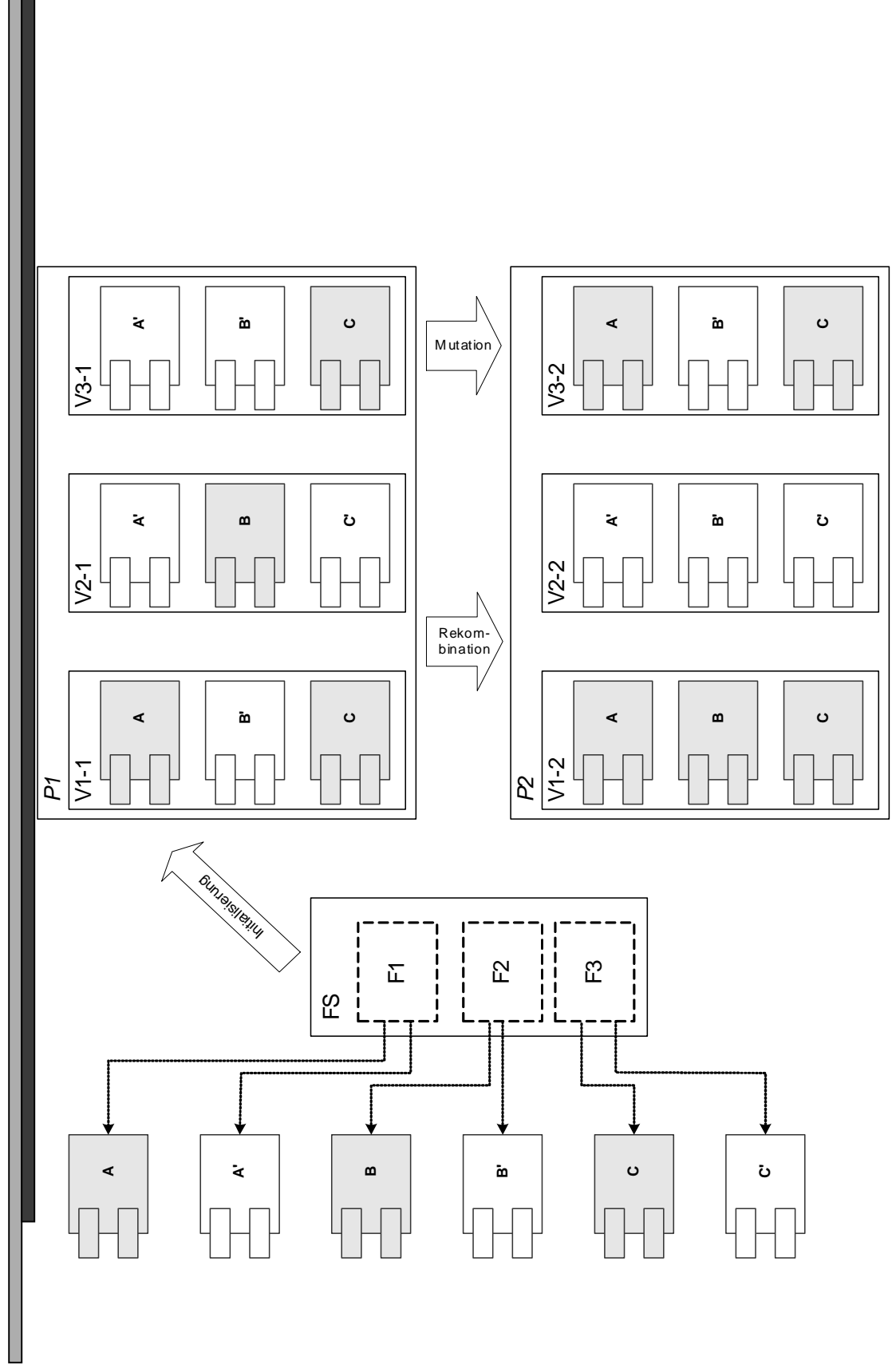


■ Evolutionäre Heuristik zur kombinatorischen Optimierung

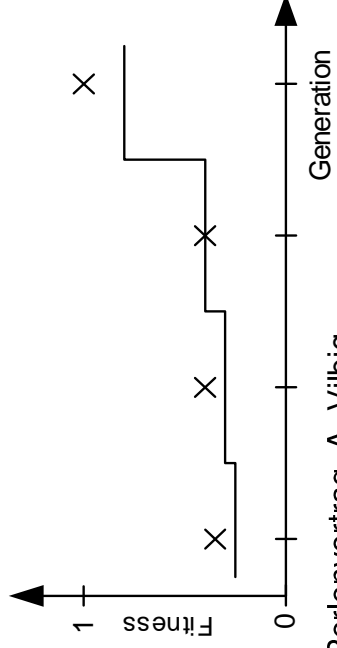
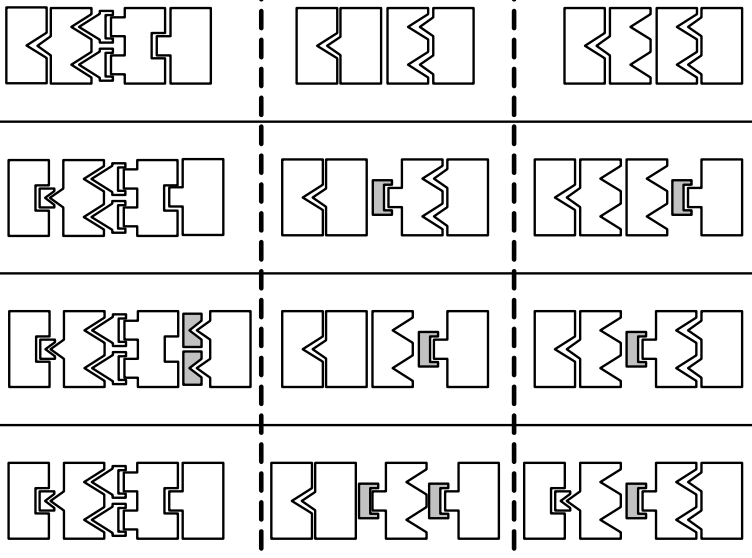
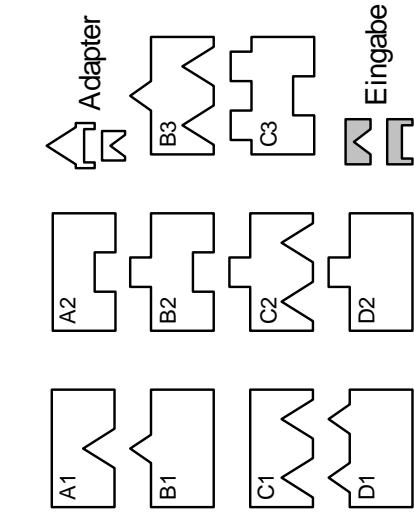
- Evolutionäre Heuristik zur kombinatorischen Optimierung
 - Anwendung bei komplexer, nicht-stetiger Zielfunktion
 - Betrachtung einer Menge von Lösungen als *Individuen* einer *Population*, welche sich durch *Selektion* und fortgesetzte Anwendung der *Genetischen Operatoren* weiterentwickelt
 - Lineare Kodierung der Lösung erforderlich \leftrightarrow *Chromosom*
 - Theorie nur für einfache Kodierung ausgearbeitet („Bitstrings“)
 - Anpassung der *Fitness-Funktion* und Gen. Operatoren erforderlich



Genetische Algorithmen I



Genetische Algorithmen II



Zusammenfassung & Ausblick

- Informatik in der Biologie
 - ┆ Essentiell zur Bewältigung großer Datenmengen
 - ┆ Anwendung „traditioneller“ Algorithmen zur Klärung wichtiger biologischer Fragestellungen
 - ┆ Vielfältige Lösungsansätze
- Biologie in der Informatik
 - ┆ Evolutionäre Heuristiken als interessante Alternative zu bekannten Verfahren
 - ┆ Erfassung und Verarbeitung unscharfer Informationen
 - ┆ Potential bei weitem nicht ausgeschöpft
- Zukünftige Fragestellungen & Ansätze
 - ┆ Genomics, Proteomics: Holistische Betrachtung komplexer biomolekularer Systeme
 - ┆ Protein Folding: Zusammenhang zwischen Sequenz und Struktur
 - ┆ Vision: „Organic Software Engineering“ - Analogie zwischen Software-System und Organismus?

Literatur

- Biochemie
 - L. Stryer, *Biochemistry*. Freeman & Co, New York 1988.
- Genetik
 - B. Lewin, *Genes IV*. Oxford University Press, New York 1990.
- Sequenzanalyse
 - M. Waterman, *Introduction to Computational Biology: Maps, Sequences and Genomes*. Chapman & Hall, London 1995.
- Molekulardynamik-Simulation
 - J. M. Haile, *Molecular Dynamics Simulation: Elementary Methods*. Wiley-Interscience, New York, 1992.
- Genetische Algorithmen
 - D. Goldberg, *Genetic Algorithms in Search, Optimization, and Machine Learning*. Addison-Wesley 1989.

WWW-Ressourcen

- Human Genome Project
 - www.ornl.gov/hgmis/, www.nhgri.nih.gov/HGP/
- Forschungseinrichtungen
 - National Center for Biotechnology Information: www.ncbi.nlm.nih.gov/
 - European Bioinformatics Institute: www.ebi.ac.uk/
 - European Molecular Biology Laboratory: www.embl-heidelberg.de/
- Online-Dokumente und URLs
 - Folding@Home: www.stanford.edu/group/pandegroup/cosm/
 - MD Simulation des Streptavidin-Biotin Komplexes (MPI Göttingen)
www.mpibpc.gwdg.de/abteilungen/071/streptavidin.html
 - BioComputing Hypertext Coursebook der Universität Bielefeld
www.techfak.uni-bielefeld.de/bcd/Curric/welcome.html
 - Vorlesung Scientific Computation an der ETH Zürich
linneus20.ethz.ch:8080/wsrscript.html
 - Lehrstuhl für effiziente Algorithmen der TU München
wwwmayr.informatik.tu-muenchen.de/