

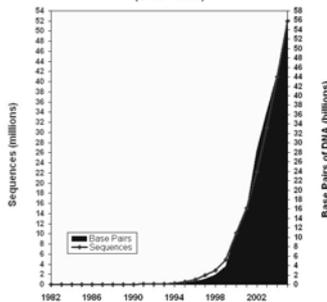
# Datenbanken in der Bioinformatik



## Bioinformatik

Anwendung von Informatikmethoden zur Untersuchung von Problemen der Molekularbiologie, die auf **sehr großen Datenmengen** beruhen

Growth of GenBank (1982 - 2005)



Quelle: <http://www.ncbi.nlm.nih.gov/Genbank/genbankstats.html>, 15.11.07

## Bioinformatikdatenbanken (BIDB)

Es geht um Daten

- eines einzigen (abgeschlossenen) Projektes
- die auf einer weltweiten und andauernden Zusammenarbeit zwischen Forschungsteams beruhen
- eines einzigen Organismus
- über das Vorkommen eines Proteins in allen möglichen Organismen

## Fragestellungen an BIDB

DNS-Analyse und Sequenzierung	Sequenz-DB
Ermittlung phylogenetischer Bäume	Sequenz-DB
Genexpressionsanalyse	→ Spezielle DB
Ermittlung biochemischer Pfade	Sequenz-DB → Spezielle DB
Strukturvorhersage	Proteinsequenz- und Proteinstruktur-DB

## Beispiel-DB\*

Gen bank 	<ul style="list-style-type: none"> <li>■ <a href="http://www.ncbi.nlm.nih.gov">www.ncbi.nlm.nih.gov</a></li> <li>■ <a href="http://www.ebi.ac.uk">www.ebi.ac.uk</a></li> <li>■ <a href="http://www.ddbj.nig.ac.jp">www.ddbj.nig.ac.jp</a></li> </ul>
German Human Genom Project 	<a href="http://www.dhgp.de">www.dhgp.de</a> (2004 eingestellt) <a href="http://www.genome.gov">www.genome.gov</a>

\* von Prof. Cieslik zusammengestellt

## Beispiel-DB

Protein bank	<ul style="list-style-type: none"><li>■ <a href="http://www.expasy.ch">www.expasy.ch</a></li><li>■ <a href="http://www.embl-heidelberg.de">www.embl-heidelberg.de</a></li><li>■ <a href="http://www.rcsb.org/pdb/">www.rcsb.org/pdb/</a></li></ul>
--------------	--



## Beispiel-DB

Phylogeny	<ul style="list-style-type: none"><li>■ <a href="http://www.ucmp.berkeley.edu/exhibit/phylogeny.html">www.ucmp.berkeley.edu/exhibit/phylogeny.html</a></li><li>■ <a href="http://evolution.genetics.washington.edu">http://evolution.genetics.washington.edu</a></li><li>■ <a href="http://tolweb.org/tree">http://tolweb.org/tree</a></li><li>■ <a href="http://awcmee.massey.ac.nz">http://awcmee.massey.ac.nz</a></li></ul>
-----------	--



## Merkmale biologischer Daten

- **Sehr komplex**  
→ mit traditionellen DBMS kann man i. a. nicht alle Datenaspekte erfassen
- **Menge und Variabilität** der Daten sind **sehr hoch**  
→ Handhabung von Datentypen und -werten muss sehr flexibel sein

## Merkmale biologischer Daten

- **Schemata** in biologischen Datenbanken **ändern sich sehr schnell**  
→ heutige Systeme erstellen wenigstens einmal pro Jahr das gesamte Datenbankschema neu
- **Darstellungen** der gleichen Daten durch verschiedene Biologen **unterscheiden sich** meist
- Die meisten Biologen haben **keine Kenntnis über die interne Struktur** der Datenbank

## Merkmale biologischer Daten

- Definition und Darstellung **komplexer Abfragen** sind wichtig  
→ Werkzeuge für die Formulierung von Abfragen müssen bereitgestellt werden
- Benutzer benötigen oft auch den **Zugriff auf „alte“ Datenwerte**  
→ Änderungen der Datenwerte müssen archiviert werden

## Schlussfolgerungen

- herkömmliche DBMS erfüllen nicht alle Anforderungen komplexer biologischer Daten  
→ Weiterentwicklungen der DBMS erforderlich
- **GENOME** (Georgia Tech Emory Network Object Management Environment) ist eine solche Weiterentwicklung (Emory ist eine Universität in Georgia)

## BIDB-Modelle

- weisen Gemeinsamkeiten bezüglich ihrer Datenmodellierung und ihrem Management auf
- BIDB verwenden z. B. folgende Modelle (Prozentangaben beziehen sich auf die von Bry und Kröger untersuchten DB)

## BIDB-Modelle

Modell	Prozent	Bemerkungen
Flat-Files	40 %	ASCII-Texte, unstrukturiert
Relationales	30 %	Herkömmliches Modell, für molekularbiol. Daten wenig geeignet
Objektorientiertes	9 %	Strukturierte Daten, passend für Daten der Molekularbiologie
ACEDB (A C. elegans Database)	4 %	Eigenes, spezielles Modell, Repräsentation von genetischen Daten

## Abfragen

- Die meisten BIDB bieten Webformulare an, mit denen man Abfragen zu den DB erstellen kann. Es sind meist nur begrenzte Abfragetypen.
- Beispiele  
**SWISS-PROT** – derzeit größte Proteindatenb.  
**BLAST** – Basic Local Alignment Search Tool  
**Entrez** – Suchmaschine  
**Tree of Life**

## Literatur

- Bry, F, Kröger, P.: A computational biology database digest: data, data analysis, and data management. Research Report PMS-FB-2002-8
- <http://www.pms.informatik.uni-muenchen.de/publikationen/#PMS-FB-2002-8>
- Bry, F, Kröger, P.: Datenbanken in der Bioinformatik. Informatik Spektrum 17, Okt. 2002
- Elmasri, R., Navathe, S. B.: Grundlagen von Datenbanksystemen. Pearson Studium 2002
- Gibas, C., Jambeck, P.: Praktische Bioinformatik. O'Reilly Verlag 2002

## Literatur

A computational biology database digest: data, data analysis, and data management enthält u. a.:

- 124 Literaturhinweise
- 111 untersuchte BIDB
- Tabelle mit URLs zu Methoden der Datenanalyse (Sequence Alignment, Gene Finding, Gene Expression)