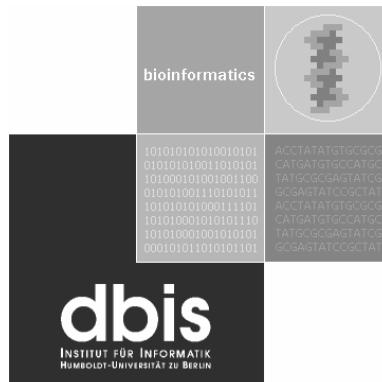
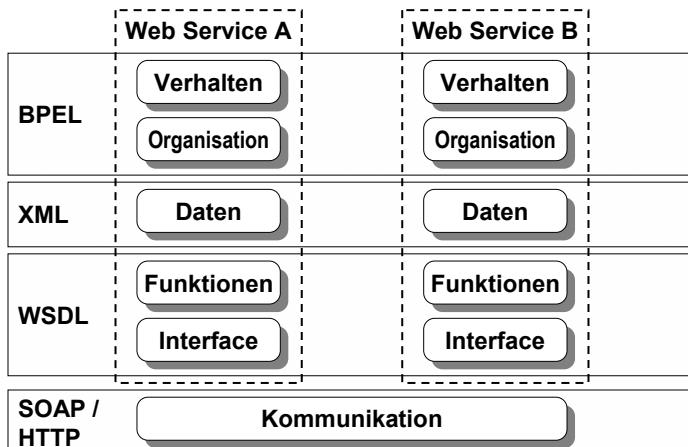


Datenintegration in der Bioinformatik

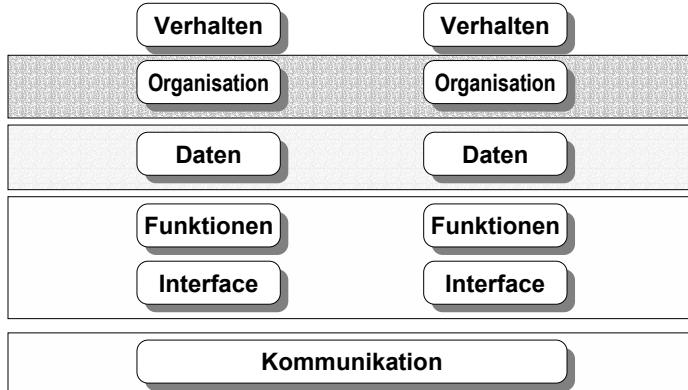


Eine Anleihe bei Prof. Reisig (Vorangegangene RingVL)





Gegenstand heute



Ursprüngliche Vortragsankündigung:

R. Eckstein, S. Heymann

XML , (Modellierung) & Datenintegration - Fortschritt im Life Science Bereich

(Schwerpunkt "Große Datenräume in Web-basierten Umgebungen")

XML hat auch in der Bioinformatik Einzug gehalten - zum Datenaustausch, aber auch zur Repräsentation der komplexen Informationen. Im ersten Teil werden Aspekte von XML sowie weitergehende Entwicklungen vorgestellt, die für den Life Science Bereich von besonderem Interesse sind. Dazu gehören die konzeptionelle Modellierung von Dokumentschemata sowie für semantische Informationen über die Biowerte. Im zweiten Teil wird ein Überblick für XML-Anwendungen im Life Science Bereich gegeben und ein typisches Anwendungsbeispiel aus dem Forschungsbereich Biodiversität & Ökologie erläutert: Die Erfassung und Darstellung von Taxonomie-Daten, den Übergang von klassischen Morphologie-basierten botanischen Schalen zu genbasierter Kladistik. Es wird demonstriert, wie die Konflikte behandelt und die Daten in ein navigierbares Graphenformat transponiert werden. Dabei kommt der GeneViator zum Einsatz.



Herkömmliche Dokumentkonventionen

Die Vortragsankündigung hätte man auch *formularbasiert* schreiben können:

Referenten:

R. Eckstein, S. Heymann

Titel:

XML , (Modellierung) & Datenintegration - Fortschritt im Life Science Bereich

Untertitel:

(Schwerpunkt "Große Datenräume in Web-basierten Umgebungen")

Zusammenfassung:

... Die Erfassung und Darstellung von Taxonomie-Daten, den Übergang von klassischen Morphologie-basierten botanischen Schulen zu genbasierter Kladistik. Es wird ...

//

Lehrstuhl für Datenbanken und Informationssysteme
Ringvortrag Informatik - Prof. J.C. Freytag, Ph.D.

5



Herkömmliche Dokumentkonventionen

Die Vortragsankündigung hätte man auch *formularbasiert* schreiben können:

Referenten:

R. Eckstein, S. Heymann

Titel:

XML , (Modellierung) & Datenintegration - Fortschritt im Life Science Bereich

Untertitel:

(Schwerpunkt "Große Datenräume in Web-basierten Umgebungen")

Zusammenfassung:

... Die Erfassung und Darstellung von Taxonomie-Daten, den Übergang von klassischen Morphologie-basierten botanischen Schulen zu genbasierter Kladistik. Es wird ...

//

(Karteikartenprinzip)



Lehrstuhl für Datenbanken und Informationssysteme
Ringvortrag Informatik - Prof. J.C. Freytag, Ph.D.

6

Gensequenzen: Beispiel SARS

 ACTGATGCGATCGCTAATATCTACCGC 

LOCUS AY274119 29736 bp RNA linear VRL 14-APR-2003

DEFINITION SARS coronavirus TOR2, complete genome.

ACCESSION AY274119

VERSION AY274119.1 GI:29826276

KEYWORDS

ORGANISM SARS Coronavirus Tor2

...

BASE COUNT 8475 a 5940 c 6186 g 9135 t

ORIGIN

1 ctacccagga aaaggccaacc aacctcgatc tcttttagat ctgttctcta aacgaaacttt

61 aaaatctgtg tagctgtcgat tcggctgcatt gccttagtgca cttacgcagt ataaacaata

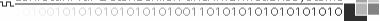
...

29641 agcccttaatg tgtaaaatata atttttagtag tgctatcccc atgtgatccc aatagcttct

29701 taggagaatg aaaaaaaaaaaaaaaa aaaaaaaaaaaaaaaa

//

Lehrstuhl für Datenbanken und Informationssysteme



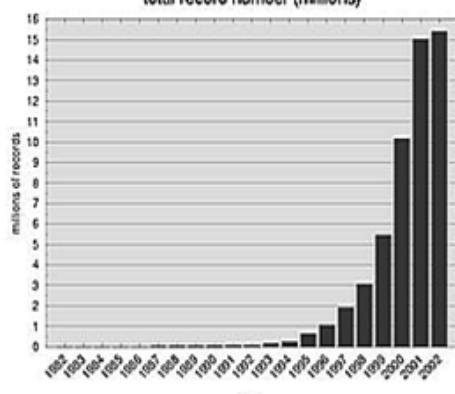
7

Ringvoleiung Informatik - Prof. J.C. Freytag, Ph.D.

Grösse der Datensammlungen

 ACTGATGCGATCGCTAATATCTACCGC 

EMBL Database Growth
total record number (millions)



- EMBL July 2002
 - > 150 Gbytes
- Microarray
 - 1 Petabyte p.A.
- Sanger Centre
 - 20 TB an Daten
- Genome Sequenzen wachsen p.A. um das Vierfache

Lehrstuhl für Datenbanken und Informationssysteme



8

Ringvoleiung Informatik - Prof. J.C. Freytag, Ph.D.

Struktur – Verarbeitbarkeit

ACTGATGCCGATCGCTAATACTACCC

```
<vortrag>
    <autor> R. Eckstein </autor>
    <autor> S. Heymann </autor>
    <titel> XML , (Modellierung) & Datenintegration - Fortschritt im Life Science Bereich</titel>
    <untertitel> Schwerpunkt "Große Datenräume in Web-basierten Umgebungen,</untertitel>
    <zusammenfassung>
        ... Die Erfassung und Darstellung von Taxonomie-Daten, den Übergang von klassischen Morphologie-
        basierten botanischen Schulen zu genbasierter <term ref="http://...> Kladistik </term>. Es wird ...
    </zusammenfassung>
</vortrag>
```

Lehrstuhl für Datenbanken und Informationssysteme

Ringvortrag Informatik - Prof. J.C. Freytag, Ph.D.

9

Struktur – Verarbeitbarkeit

ACTGATGCCGATCGCTAATACTACCC

```
<vortrag>
    <autor> R. Eckstein </autor>
    <autor> S. Heymann </autor>
    <titel> XML , (Modellierung) & Datenintegration - Fortschritt im Life Science Bereich</titel>
    <untertitel> Schwerpunkt "Große Datenräume in Web-basierten Umgebungen,</untertitel>
    <zusammenfassung>
        ... Die Erfassung und Darstellung von Taxonomie-Daten, den Übergang von klassischen Morphologie-
        basierten botanischen Schulen zu genbasierter <term ref="http://...> Kladistik </term>. Es wird ...
    </zusammenfassung>
</vortrag>
```

Lesbarkeit ↔ Struktur ↔ Verarbeitbarkeit

Lehrstuhl für Datenbanken und Informationssysteme

Ringvortrag Informatik - Prof. J.C. Freytag, Ph.D.

10

Strukturierter Text

ACTGATGCGATCGCTAATACTACCC

```
<vortrag>
  <autor> R. Eckstein </autor>
  <autor> S. Heymann </autor>
  <titel> XML , (Modellierung) & Datenintegration - Fortschritt im Life Science Bereich</titel>
  <untertitel> Schwerpunkt "Große Datenräume in Web-basierten Umgebungen,</untertitel>
  <zusammenfassung>
    ... Die Erfassung und Darstellung von Taxonomie-Daten, den Übergang von klassischen Morphologie-
    basierten botanischen Schulen zu genbasierter <term ref= http://www.visualgenomics.ca/gordonp/xml/>
    Kladistik</term>. Es wird ...
  </zusammenfassung>
</vortrag>
```

→ EXTENSIBLE MARKUP LANGUAGE

Lehrstuhl für Datenbanken und Informationssysteme
Ringvortrag Informatik - Prof. J.C. Freytag, Ph.D.

11

EXTENSIBLE MARKUP LANGUAGE

ACTGATGCGATCGCTAATACTACCC

```
<vortrag>
  <autor> R. Eckstein </autor>
  <autor> S. Heymann </autor>
  <titel> XML , (Modellierung) & Datenintegration - Fortschritt im Life Science Bereich</titel>
  <untertitel> Schwerpunkt "Große Datenräume in Web-basierten Umgebungen,</untertitel>
  <zusammenfassung>
    ... Die Erfassung und Darstellung von Taxonomie-Daten, den Übergang von klassischen Morphologie-
    basierten botanischen Schulen zu genbasierter <term ref= http://www.visualgenomics.ca/gordonp/xml/>
    Kladistik</term>. Es wird ...
  </zusammenfassung>
</vortrag>
```

Beschreibungssprachen im Bio-Bereich:
<http://www.visualgenomics.ca/gordonp/xml/>

Lehrstuhl für Datenbanken und Informationssysteme
Ringvortrag Informatik - Prof. J.C. Freytag, Ph.D.

12

EXTENSIBLE MARKUP LANGUAGE

ACTGATGCCGATCGCTAATATCTACCGT

```
<vortrag>
  <autor> R. Eckstein </autor>
  <autor> S. Heymann </autor>
  <titel> XML , (Modellierung) & Datenintegration - Fortschritt im Life Science Bereich</titel>
  <untertitel> Schwerpunkt "Große Datenräume in Web-basierten Umgebungen,</untertitel>
  <zusammenfassung>
    ... Die Erfassung und Darstellung von Taxonomie-Daten, den Übergang von klassischen Morphologie-
    basierter botanischen Schulen zu genbasierter <term ref= http://www.visualgenomics.ca/gordonp/xml/>
    Kladistik</term>. Es wird ...
  </zusammenfassung>
</vortrag>
```

Beschreibungssprachen im Bio-Bereich:
<http://www.visualgenomics.ca/gordonp/xml/>

Handbücher

Dokumenttypdefinitionen

Lehrstuhl für Datenbanken und Informationssysteme
Ringvolese Informatik - Prof. J.C. Freytag, Ph.D.

13

Kompendium aller Erkenntnisse

ACTGATGCCGATCGCTAATATCTACCGT

DTD:

```
<!ENTITY % local_aa_type.value "">
<!ENTITY % aa_type "type (A|B|C|D|E|F|G|H|I|K|L|M
                           |N|P|Q|R|S|T|V|W|X|Y|Z
                           %local_aa_type.value;
                           #REQUIRED)">
```

Lehrstuhl für Datenbanken und Informationssysteme
Ringvolese Informatik - Prof. J.C. Freytag, Ph.D.

14

DTD's: Kompendium neuer Erkenntnisse

ACTGATGCCATCGCTAAATCTACCC

```
<!ENTITY % local_aa_type_value ""> (bis 2002)
```

```
<!ENTITY % aa_type "type (A|B|C|D|E|F|G|H|I|K|L|M
```

```
|N|P|Q|R|S|T|V|W|X|Y|Z
```

```
%local_aa_type_value;) #REQUIRED">
```

```
<!ENTITY % local_aa_type_value ""> (seit Feb. 2003)
```

```
<!ENTITY % aa_type "type (A|B|C|D|E|F|G|H|I|K|L|M
```

```
|N|P|Q|R|S|T|V|U|W|X|Y|Z
```

```
%local_aa_type_value;) #REQUIRED">
```

Lehrstuhl für Datenbanken und Informationssysteme
Ringvoleiung Informatik - Prof. J.C. Freytag, Ph.D.

15

Beispiele aus der Abstammungslehre

ACTGATGCCATCGCTAAATCTACCC



Lehrstuhl für Datenbanken und Informationssysteme
Ringvoleiung Informatik - Prof. J.C. Freytag, Ph.D.
<http://www.visualgenomics.ca/gordonp/xml/>

16

Beispiele aus der Abstammungslehre

ACTGATGCCGATCGCTAATACTACCC



Homo sapiens:

Other names:

man[common name]

Lineage(full)

cellular organisms; Eukaryota; Fungi/Metazoa group; Metazoa; Eumetazoa; Bilateria; Coelomata; Deuterostomia; Chordata; Craniata; Vertebrata; Gnathostomata; Teleostomi; Euteleostomi; Sarcopterygii; Tetrapoda; Amniota; Mammalia; Theria; Eutheria; Primates; Catarrhini; Hominidae; Homo/Pan/Gorilla group; Homo

Ringvoleseung Informatik - Prof. J.C. Freytag, Ph.D.

Lehrstuhl für Datenbanken und Informationssysteme
 http://www.visualgenomics.ca/gordonp/xml/

17

Beispiele aus der Abstammungslehre

ACTGATGCCGATCGCTAATACTACCC



```

<?xml version="1.0" encoding="UTF-8"?>
<?xml-stylesheet type="text/xsl" href="HTMLOutput.xsl"?>

<!DOCTYPE MultipleClassifications SYSTEM "XMLOutput.dtd">

<MultipleClassifications Order="Alpha">
<!--#NEXUS [18-Dec-2001 16:59:36]-->
<ranks>
  <rank RankID="Family"><rankName>Family</rankName><rankValue>10</rankValue></rank>
  <rank RankID="Sub-Family"><rankName>Sub-Family</rankName><rankValue>15</rankValue></rank>
  <rank RankID="Genus"><rankName>Genus</rankName><rankValue>20</rankValue></rank>
  <rank RankID="Species"><rankName>Species</rankName><rankValue>25</rankValue></rank>
</ranks>
<taxa>
  <taxon RankIDREF="Family" TaxonID="TApiaceae">Apiaceae</taxon>
  <taxon RankIDREF="Sub-Family" TaxonID="TPauciugatae">Pauciugatae</taxon>
  ...

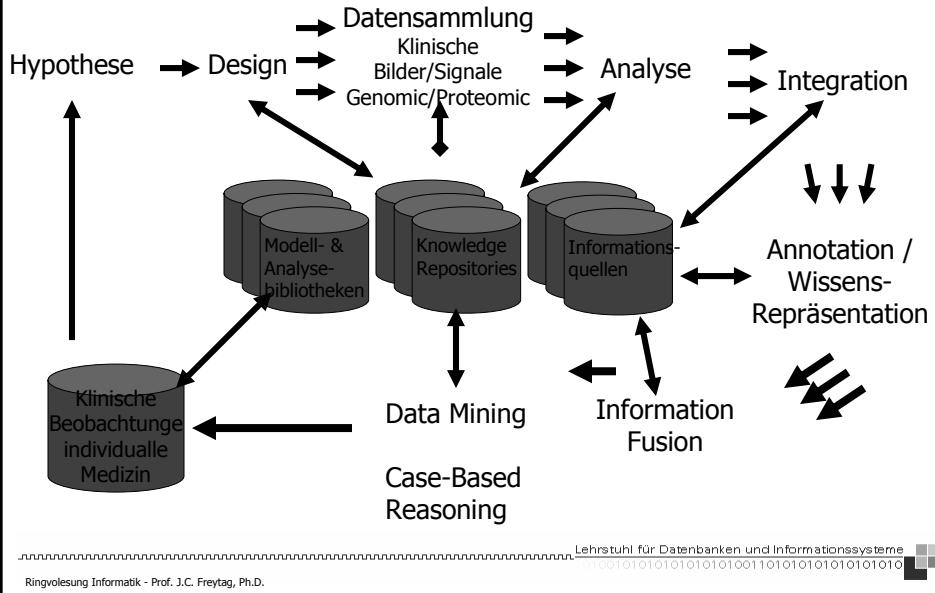
```

Ringvoleseung Informatik - Prof. J.C. Freytag, Ph.D.

Lehrstuhl für Datenbanken und Informationssysteme
 http://www.visualgenomics.ca/gordonp/xml/

18

Beispiel: Genetik & Medizin



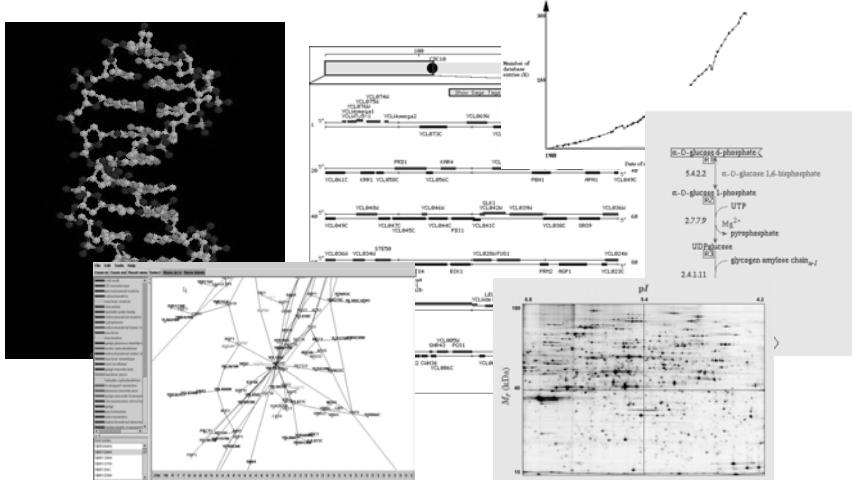
Herausforderungen...

- Formatheterogenität
- Datenheterogenität / Anzahl der Datenquellen
- Umfang der Daten / Grösse der Datensammlungen
- Zugriffsheterogenität

Formatheterogenität

ACTGATGCGATCGCTAATATCTACCC

- Multimedia: Bilder & Video (e.g. microarrays, 3D, ...)



Ringvoles Informatik - Prof. J.C. Freytag, Ph.D.

Lehrstuhl für Datenbanken und Informationssysteme

21

Formatheterogenität (cont.)

ACTGATGCGATCGCTAATATCTACCC

- Text “Annotationen” & Literatur
 - strukturiert vs. semistrukturiert vs. unstrukturiert
- Unterschiedliche Formate, Strukturen, Schemata, Umfänge, ...
- Web-Schnittstellen, Verteilung als Dateien, Datenbank-Dumps, XML-Dokumente, ...



```

ID TRBG361 standard; RNA; PLN; 1859 BP.
XX X56734; S46826;
XX
SV X56734.1
XX
DT 12-SEP-1991 (Rel. 29, Created)
DF 15-MAR-1999 (Rel. 59, Last updated, Version 9)
XX
DE Trifolium repens mRNA for non-cyanogenic beta-glucosidase
XX
KW beta-glucosidase.
XX
OS Trifolium repens (white clover)
OC Eukaryota; Viriplantae; Streptophyta; Embryophyta.
OC Spermatophyta; Magnoliophyta; eudicots; core eudicots.
OC eurosids I; Fabales; Fabaceae; Papilionoideae; Trifoliaceae.
XX

```

Lehrstuhl für Datenbanken und Informationssysteme

22

Ringvoles Informatik - Prof. J.C. Freytag, Ph.D.

Daten-/Inhaltsheterogenität

▪ Genomische, proteomische, transcriptomische, metabolomische, Protein-Protein Interactionen, regulatorische Bio-Netzwerke, Alinierungen, Krankheiten, Patterns & Motifs, Proteine Structuren, Proteinklassifikationen und -familien, spezielle Proteine (Enzyme, Rezeptoren), ...

Lehrstuhl für Datenbanken und Informationssysteme

Ringvoleis Informatik - Prof. J.C. Freytag, Ph.D.

23

Zugriffsheterogenität

NCBI National

EMBL European Bioinformatics Institute

The UniProt Protein Data Bank, Version 16.0

Kyoto Encyclopedia of Genes and Genomes

Lehrstuhl für Datenbanken und Informationssysteme

Ringvoleis Informatik - Prof. J.C. Freytag, Ph.D.

24

Genomisch relevante Bereiche

ACTGATGCCATCGCTAATATCTACGG

Environment

Diseases

Experiments

Pathways

Life

Evolution

DNA
Genome

RNA
Transcriptome

Amino Acids
Proteome

Lehrstuhl für Datenbanken und Informationssysteme
Ringvortrag Informatik - Prof. J.C. Freytag, Ph.D.

25

Inhalt von Datenquellen

Environment

Diseases
OMIM

Experiments
Express

Pathways
Brenda
KEGG

Life
Gene Ontology

Evolution
Taxonomy

DNA
RefSeq
LocusLink

RNA
EMBL (EST)

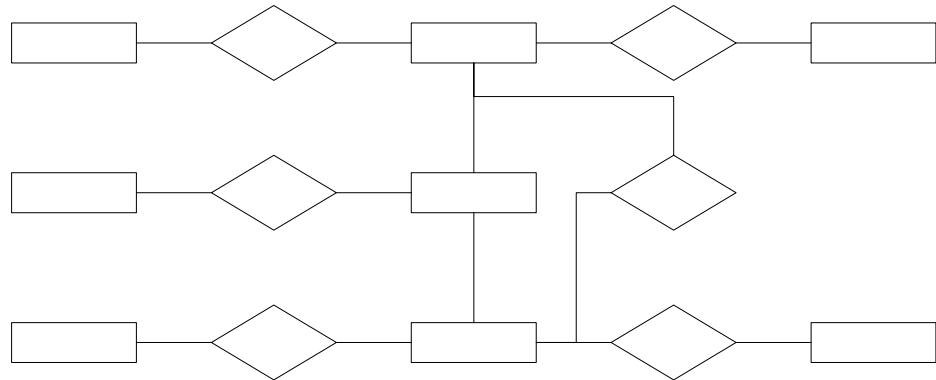
Amino Acids
SWISS-PROT
Proteome
Interpro

Lehrstuhl für Datenbanken und Informationssysteme
Ringvortrag Informatik - Prof. J.C. Freytag, Ph.D.

26

Datenmodellierung

ACTGATGCCGATCGCTAATATCTACCC



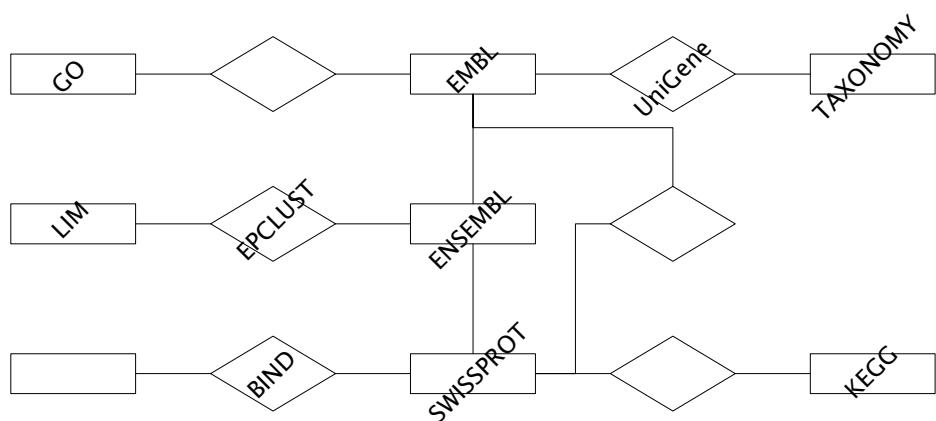
Lehrstuhl für Datenbanken und Informationssysteme

Ringvortrag Informatik - Prof. J.C. Freytag, Ph.D.

27

Datenmodellierung

ACTGATGCCGATCGCTAATATCTACCC

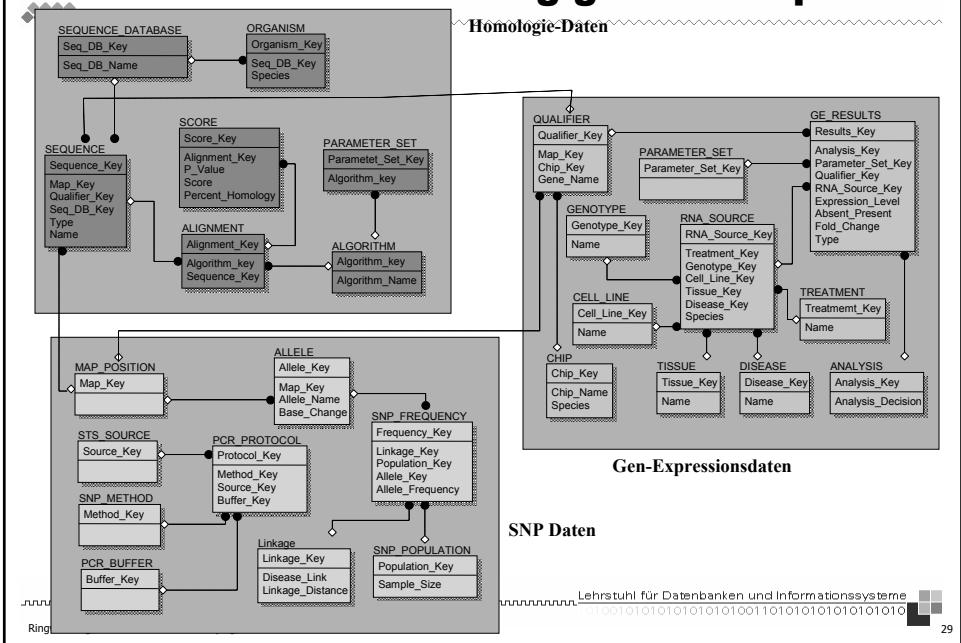


Lehrstuhl für Datenbanken und Informationssysteme

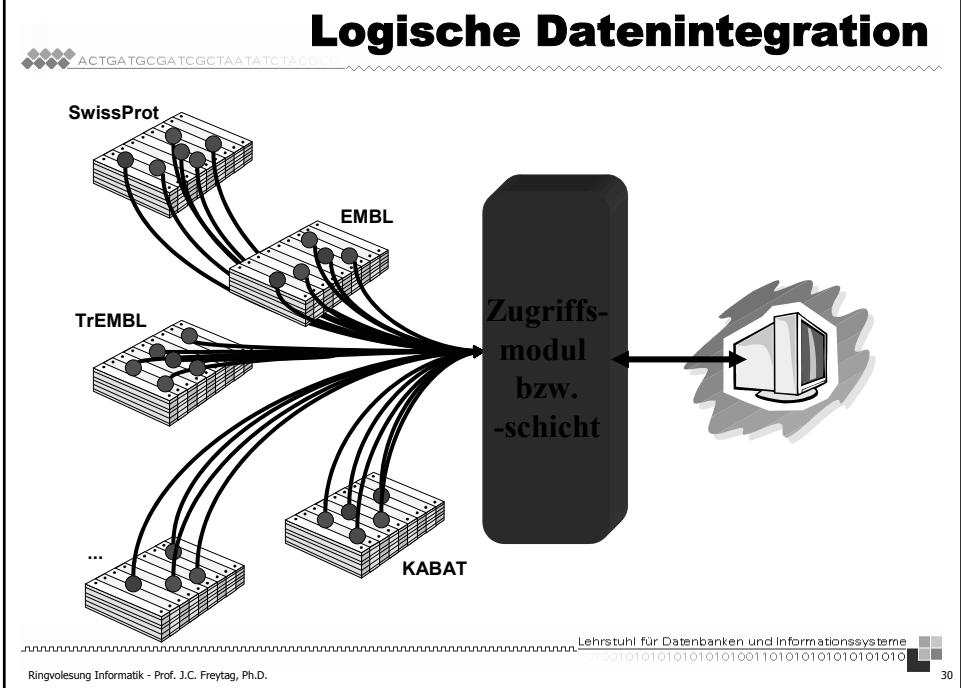
Ringvortrag Informatik - Prof. J.C. Freytag, Ph.D.

28

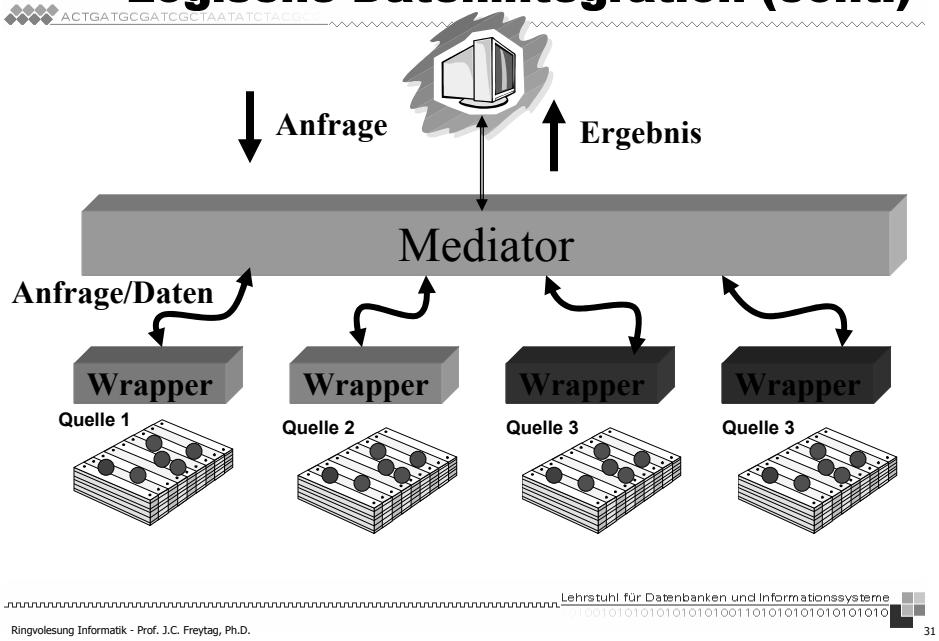
Daten aus unabhängigen Datenquellen



Logische Datenintegration



Logische Datenintegration (cont.)

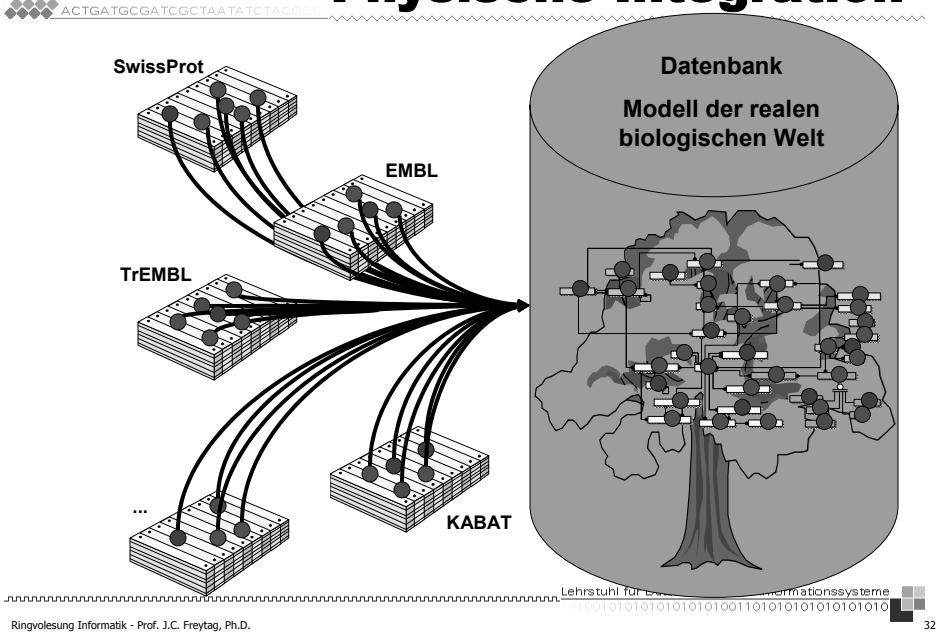


Lehrstuhl für Datenbanken und Informationssysteme

Ringvortrag Informatik - Prof. J.C. Freytag, Ph.D.

31

Physische Integration



Lehrstuhl für Datenbanken und Informationssysteme

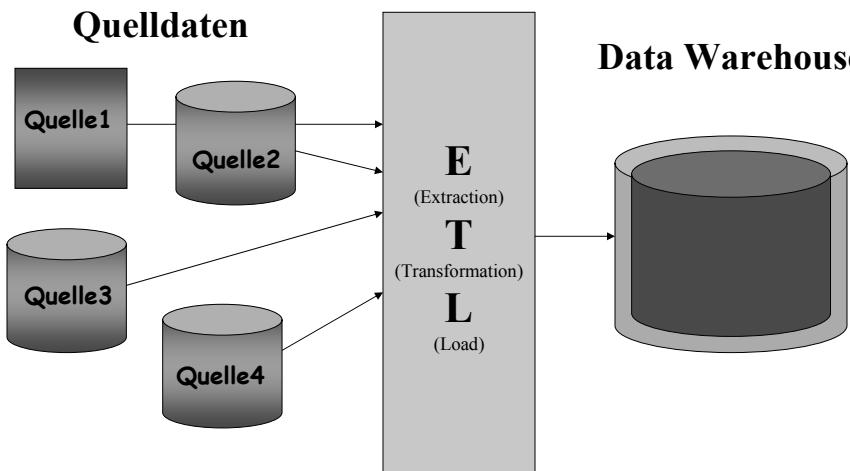
Ringvortrag Informatik - Prof. J.C. Freytag, Ph.D.

32

Data Warehousing

ACTGATGCCGATCGCTAATATCTACCGC

Quelldaten



Data Warehouse

Lehrstuhl für Datenbanken und Informationssysteme
Ringvortrag Informatik - Prof. J.C. Freytag, Ph.D.

33

Gene-EYe Integrationsplattform

ACTGATGCCGATCGCTAATATCTACCGC

Genome Data Warehouse Layer (GDW Schema)

Wissen

Biologische Entitäten -> Biologische Konzepte (e.g. Lebenszyklus)

Genome DataBase Layer (GDB Schema)

Inhalt

Relationale Entitäten -> Biologische Entitäten (e.g. Gene)

Genome Data Store Layer (GDS Schema)

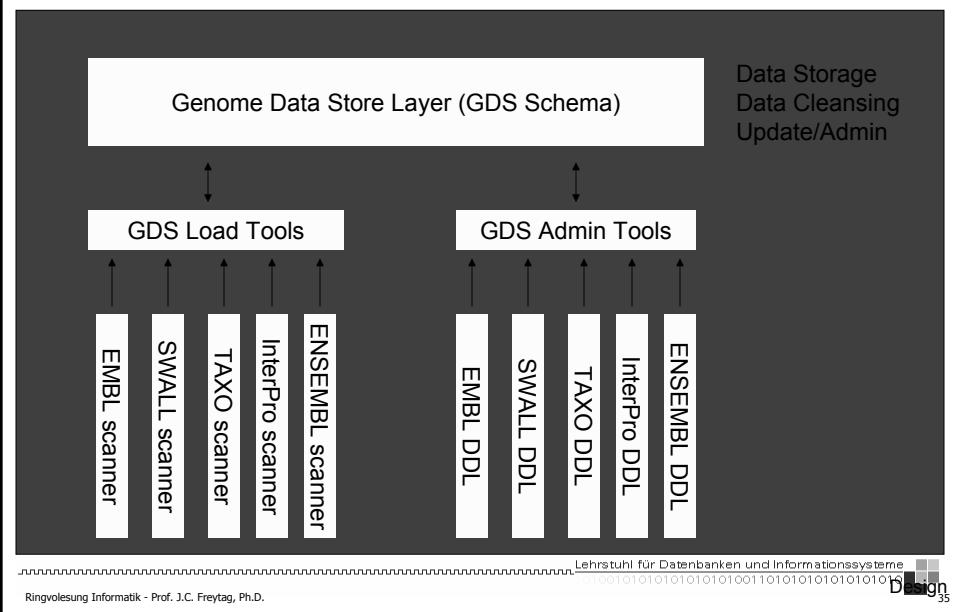
Daten

"Datei"-Daten -> Relationale Entitäten (e.g. EMBL)

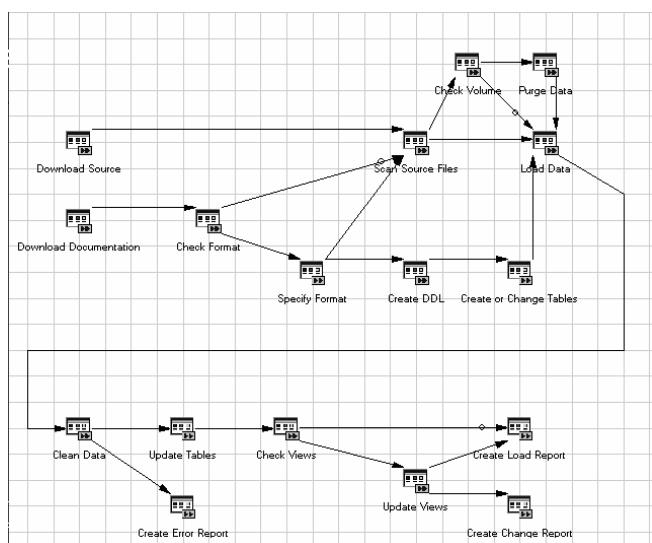
Lehrstuhl für Datenbanken und Informationssysteme
Ringvortrag Informatik - Prof. J.C. Freytag, Ph.D.

34

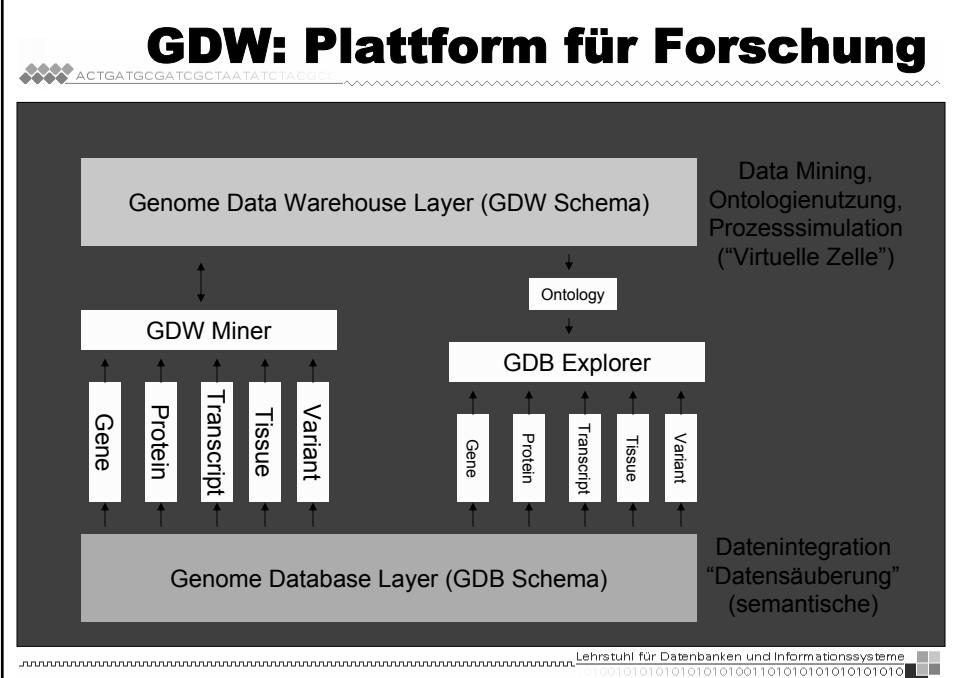
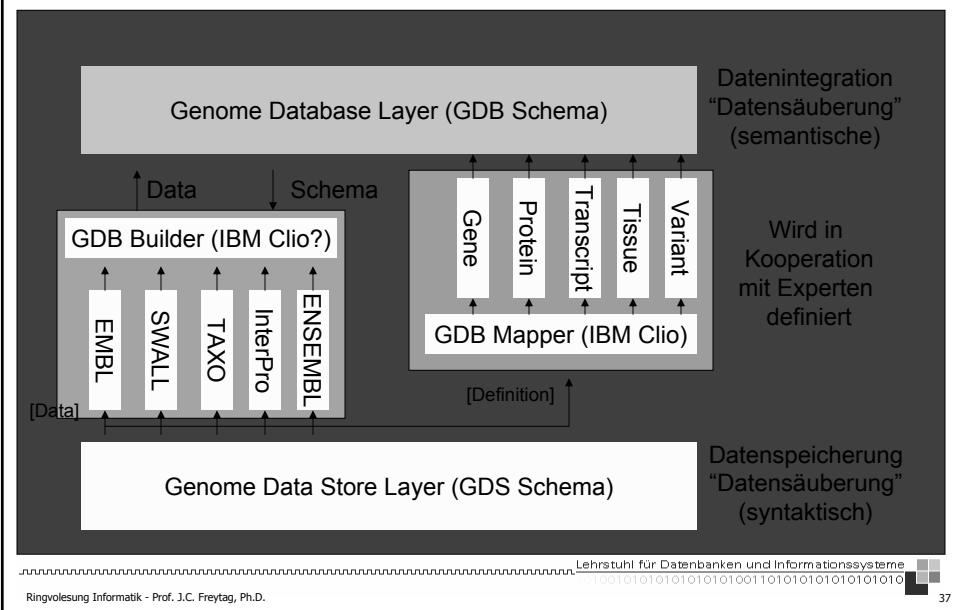
GDS: Von der Datei zur Datenbank



Modellierung des “Wartungsprozesses”

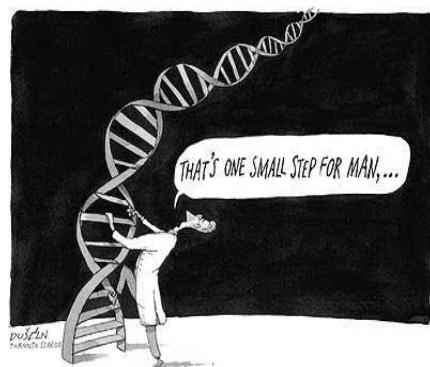


GDB: von den Daten zur Biologie



ACTGATGCCATCGCTAATATCTACCC

What are the goals?



Source: Dusan Petricic, Toronto, Ontario – The Toronto Star
<http://cagle.slate.msn.com/news/gene6.asp>

Fragen??

Ringvortrag Informatik - Prof. J.C. Freytag, Ph.D.

Lehrstuhl für Datenbanken und Informationssysteme



39