

Andreas Holzinger
VO 709.049 Medical Informatics
04.11.2015 11:15-12:45

Lecture 04

Biomedical Databases: Data Acquisition, Storage, Information Retrieval and Use

a.holzinger@tugraz.at
Tutor: markus.plass@student.tugraz.at
<http://hci-kdd.org/biomedical-informatics-big-data>



A. Holzinger 709.049 1/92 Med Informatics L04



Schedule

- 1. Intro: Computer Science meets Life Sciences, challenges, future directions
- 2. Back to the future: Fundamentals of Data, Information and Knowledge
- 3. Structured Data: Coding, Classification (ICD, SNOMED, MeSH, UMLS)
- 4. Biomedical Databases: Acquisition, Storage, Information Retrieval and Use
- 5. Semi structured and weakly structured data (structural homologies)
- 6. Multimedia Data Mining and Knowledge Discovery
- 7. Knowledge and Decision: Cognitive Science & Human-Computer Interaction
- 8. Biomedical Decision Making: Reasoning and Decision Support
- 9. Intelligent Information Visualization and Visual Analytics
- 10. Biomedical Information Systems and Medical Knowledge Management
- 11. Biomedical Data: Privacy, Safety and Security
- 12. Methodology for Info Systems: System Design, Usability & Evaluation

A. Holzinger 709.049 2/92 Med Informatics L04



Keywords of the 4th Lecture

- Bayes' Rule
- Biomedical data warehouse
- Business hospital information system
- Clinical workflow
- Data integration
- Enterprise data modeling
- Information retrieval (IR)
- Probabilistic Model
- Quality of information retrieval
- Set theoretic model
- Vector Space Model (VSM)

A. Holzinger 709.049 3/92 Med Informatics L04



Advance Organizer (1/4)

- **Business intelligence (BI)** = a type of application software designed to report, analyze, and present information on real-time management dashboards, i.e., integrated displays of metrics that measure the performance of a system;
- **Cassandra** = an open source and free database management system designed to handle huge amounts of data on a distributed system. This system was developed at Facebook and is now managed as a project of the Apache Software foundation.
- **Cladogram** = a phylogenetic tree to show evolutionary relationships with species represented by nodes and lines of descent represented by links (unrooted or rooted);
- **Classification system** = arbitrary in nature, there is no standard measure of difference that defines a species, genus, family, or order;
- **Cloud computing** = a computing paradigm in which highly scalable computing resources, often configured as a distributed system, are provided as a service
- **CPOE (Computerized physician order entry)** = a process of electronic entry of medical practitioner instructions for the treatment of patients (particularly hospitalized patients) under his or her care;
- **Data Mart (DM)** = access layer of a data warehouse environment that is used to get data to the users. The DM is a subset of the DW, usually oriented to a specific business line or team to provide data to users usually through business intelligence tools;
- **DBGET** = a data retrieval tool (simpler than ENTREZ) from the Kyoto University, which covers more than 20 databases, related to the Kyoto Encyclopedia of Genes and Genomes
- **Distance matrix method** = work by two most closely related taxa in a distance matrix and clustering them;

A. Holzinger 709.049 4/92 Med Informatics L04



Advance Organizer (2/4)

- **Ensembl** = database format;
- **ENTREZ** = a dedicated data retrieval tool;
- **Extract, transform, and load (ETL)** = Software tools used to extract data from outside sources, transform them to fit operational needs, and load them into a database or data warehouse;
- **Federated data base system** = type of meta-database management system, which integrates multiple autonomous database systems into a single federated database;
- **Genetic algorithm** = a technique used for optimization inspired by the process of natural evolution or "survival of the fittest." Often described as a type of "evolutionary algorithm," these algorithms are well-suited for solving nonlinear problems;
- **Genomes OnLine Databases (GOLD)** = a general genomics gateway;
- **Hadoop** = An open source (free) software framework for processing huge datasets on certain kinds of problems on a distributed system. Its development was inspired by Google's MapReduce and Google File System.
- **Hbase** = An open source (free), distributed, non-relational database modeled on Google's Big Table. It was originally developed by Powerset and is now managed as a project of the Apache Software foundation as part of the Hadoop.
- **Information Extraction (IE)** = automatic assignment of meaning to elementary textual entities and complex structured information objects;
- **Information Retrieval (IR)** = indexing and retrieval of information in documents;
- **KEGG** = Kyoto Encyclopedia of Genes & Genomes, a combined database containing information on types of proteins (receptors, signal transduction components, enzymes)

A. Holzinger 709.049 5/92 Med Informatics L04



Advance Organizer (3/4)

- **MapReduce** = A software framework introduced by Google for processing huge datasets on certain kinds of problems on a distributed system.³² Also implemented in Hadoop;
- **Mashup** = An application that uses and combines data presentation or functionality from two or more sources to create new services. These applications are often made available on the Web, and frequently use data accessed through open application programming interfaces or from open data sources;
- **MEDLINE** = Literature data bank;
- **Metadata** = Data that describes the content and context of data files, e.g., means of creation, purpose, time and date of creation, and author;
- **MMDB** = Molecular Modeling Database, can be accessed at the NCBI (National Center for Biotechnology Information) using ENTREZ;
- **Natural language processing (NLP)** = a set of machine learning techniques from computer science and linguistics that uses computer algorithms to analyze human (natural) language;
- **Neural networks** = computational models, inspired by the structure and workings of biological neural networks (i.e., the cells and connections within a brain), that find non linear patterns in data;
- **Non-relational database** = a database that does not store data in tables (rows and columns). (In contrast to relational database);
- **Online Mendelian Inheritance in Man (OMIM)** = a database as resource for the study of human genetics and human molecular medicine;
- **PDB** = Protein Data Bank contains data derived from X-ray crystallography and NMR (nuclear magnetic resonance) studies;

A. Holzinger 709.049 6/92 Med Informatics L04

Advance Organizer (4/4)

- **Phylogenetics** = similarities and differences among species can be used to infer evolutionary relationships (=phylogenies); Examples for phylogenetic software: PAUP, PHYLIP;
- **PROSITE** = database containing sequence patterns associated with protein family membership, specific protein functions and post-translational modifications;
- **R** = An open source (free) programming language and software environment for statistical computing and graphics;
- **Relational database** = a database made up of a collection of tables (relations), i.e., data are stored in rows and columns. Relational database management systems (RDBMS) store a type of structured data. SQL is the most widely used language for managing relational databases (see there);
- **Semi-structured data** = Data that do not conform to fixed fields but contain tags and other markers to separate data elements. Examples of semi-structured data include XML or HTML-tagged text. Contrast with structured data and unstructured data.
- **Similarity table** = distance table;
- **SQL** = Originally an acronym for structured query language, SQL is a computer language designed for managing data in relational databases. This technique includes the ability to insert, query, update, and delete data, as well as manage data schema (database structures) and control access to data in the database;
- **SRS** = Sequence Retrieval System, a data retrieval tool based on open source software
- **SWISS-PROT** = is a databank containing a collection of confirmed protein sequences with annotations relating to structure, function and protein family assignment;
- **UniGene** = experimental facility for the clustering of GenBank sequences and is related to EST (expressed sequence tag) data;

A. Holzinger 709.049

7/92

Med Informatics L04

Glossary

- ACeDB = A C elegans Data Base
- ADE = adverse drug events
- CDSS = clinical decision support system
- CPOE = computerized physician order entry
- DBMS = Data Base Management System
- EMAC = electronic medication administration chart
- EO = electronic order
- ERT = error registration table
- GFR = glomerular filtration rate
- HIS = Hospital Information System (DE: KIS = Krankenhaus Informations System)
- HWO = handwritten order
- NICU = neonatal intensive care unit
- NOE = nurse order entry (followed by physician's verification and countersignature)
- PBMAC = paper-based medication administration chart
- POE = physician order entry
- RR = rate ratio
- UniProt = Universal Protein Ressource

A. Holzinger 709.049

8/92

Med Informatics L04

Learning Goals ... at the end of the 4th lecture you ...

- ... have an overview about the general **architecture of an Hospital Information System**
(Details in lecture 10: Medical Information Systems and Biomedical Knowledge Management!);
- ... know some principles of **hospital databases**;
- ... have an overview on some important **biomedical databases**;
- ... are familiar with some basic methods of **information retrieval**;

A. Holzinger 709.049

9/92

Med Informatics L04

Key Challenges

- Increasingly large and *complex data sets* due to **data intensive biomedicine** [1]
- Increasing amounts of **non-standardized** and **un-structured information** (e.g. "free text")
- **Data quality, data integration, universal access**
- **Privacy, security, safety, data protection, data ownership, fair use of data** (see →Lecture 11) [2]
- **Time aspects in databases** [3]

[1] Shah, N. H. & Tennenbaum, J. D. 2012. The coming age of data-driven medicine: translational bioinformatics' next frontier. Journal of the American Medical Informatics Association, 19, (E1), E2-E4.

[2] Kieseberg, P., Hobel, H., Schrittwieser, S., Weippl, E. & Holzinger, A. 2014. Protecting Privacy in Data-Driven Biomedical Science. In: LNCS 8401. Berlin Heidelberg: Springer pp. 301-316.

[3] Gschwandtner, T., Gärtner, J., Aigner, W. & Miksch, S. 2012. A taxonomy of dirty time-oriented data. In: LNCS 7465. Heidelberg, Berlin: Springer, pp. 58-72.

A. Holzinger 709.049

10/92

Med Informatics L04

Let us start with a look into the Hospital ...

A. Holzinger 709.049

11/92

Med Informatics L04

Slide 4-1 Hospital Information System: Typical Scenario



G'sund Net, Ausgabe 45, März 2005

A. Holzinger 709.049

12/92

Med Informatics L04

Slide 4-2 HIS: Typical View on the Clinical Workplace

Ra.	OP	Patient	OP R leg. AM-OF	Diagnose
OP_1	GYNOP	08.09 (M, 53)	✓ GEMOC	Üb-Schmerzen bei Adenomensis uteri
		10.17 (M, 43)	✓ GEMOC	Cyst ov.
		11.28 (M, 39)	✓ GEMOC	Plazenzarrest
		12.52 (M, 57)	✓ GEMOC	BPM
		12.57 (M, 41)	✓ GEMOC	Bilirubin-Hemmungsdosis
		12.58 (M, 52)	✓ GEMOC	Uterusmyome
OP_2	UNFOP	08.51 (M, 79)	✓ GEMIB	Vakuumextraktions
		10.51 (M, 71)	✓ GEMIB	Koagulations
		14.25 (M, 39)	✓ GEMIB	Ste. Weber C Frakur, op. 2.12.2010
		17.02 (W, 77)	✓ GEMIB	Schenkelhalsfraktur medial garden V/re b. liegend
SEC	GYNOP	09.01 (W, 40)	✓ GEMIA	Sekto primär Einling (Rehruung Mutter)
		10.23 (W, 38)	✓ GEMIB	Befindlichkeit
		12.39 (W, 34)	✓ GEMIA	Sekto V. a. von Präzessatze
				Sekto primär Einling (Rehruung Mutter)

G'sund Net, Ausgabe 70, Juni 2011

A. Holzinger 709.049 13/92 Med Informatics L04

Slide 4-3: Much of hospital work is teamwork ...

... and requires a lot of communication and information exchange ...

Holzinger, A., Geierhofer, R., Ackerl, S. & Searle, G. (2005). CARDIAC@VIEW: The User Centered Development of a new Medical Image Viewer. Central European Multimedia and Virtual Reality Conference, Prague, Czech Technical University (CTU), 63-68.

A. Holzinger 709.049 14/92 Med Informatics L04

Slide 4-4: The medical report is the most important output

Radiologischer Befund

angelegt am 06.05.2006/20:26
geändert von:
gedruckt am 17.11.2006/08:34
Arzt: NCHHN

Kurzamnese: St.p. SHT
Fragestellung: -
Untersuchung: Thorax eine Ebene liegend

SB
Bewegungsartefakte. Zustand nach Schädelhirntrauma.
Das Cor in der Größenorm, keine akuten Stauungszeichen.
Fragliches Infiltrat parahilär li. im UF, RW-Erguss li.
Zustand nach Anlage eines ET, die Spitze ca. 5cm cranial der Bifurkation, liegt MS. orthotop positioniert. ZVK über re., die Spitze in Proj. auf die VCS. Kein Hinweis auf Pneumothorax.
Der re. Rezessus frei.

Mit kollegialen Grüßen

*** Elektronische Freigabe durch am 09.05.2006 ***

Holzinger, A., Geierhofer, R. & Errath, M. 2007. Semantic Information Extraction in medical Informationssystemen. *Informatik Spektrum*, 30, (2), 69-78.

A. Holzinger 709.049 15/92 Med Informatics L04

Slide 4-5 Excuse: Chess Game versus Natural Language

<http://stanford.edu/~cziech/cs221/apps/deepBlue.html>

A. Holzinger 709.049 16/92 Med Informatics L04

German Example: Synonymy and Ambiguity

„die Antrumschleimhaut ist durch Lymphozyten infiltriert“
„lymphozytäre Infiltration der Antrummukosa“
„Lymphoyteninfiltration der Magenschleimhaut im Antrumbereich“

HWI = Harnwegsinfekt, Hinterwandinfarkt, Hakenwurminfektion, Halswirbelimmobilisation, Hinterwandischämie, Hip Waist Index, Height-Width Index, Häufig wechselnder Intimpartner, Hepatitis weight index ...

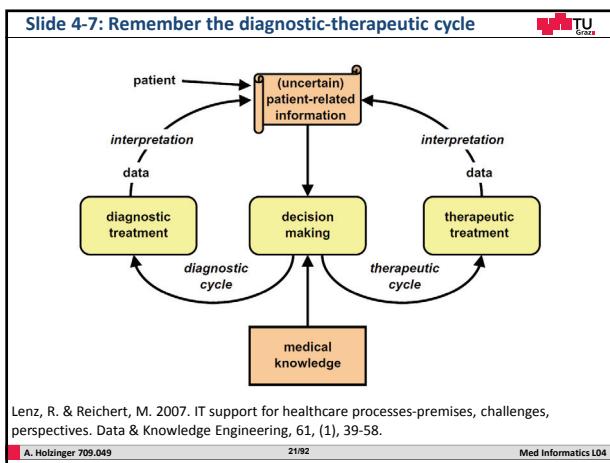
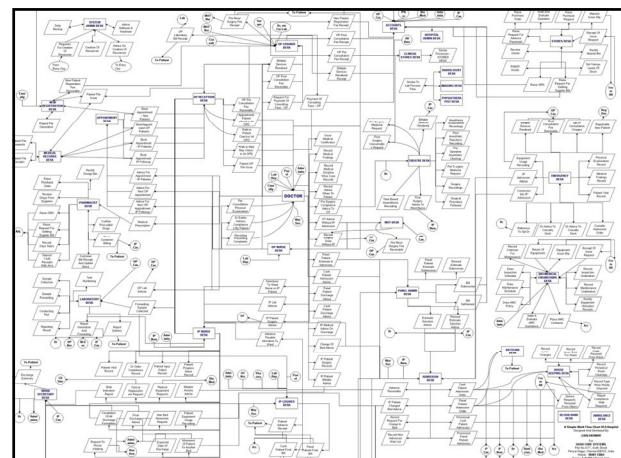
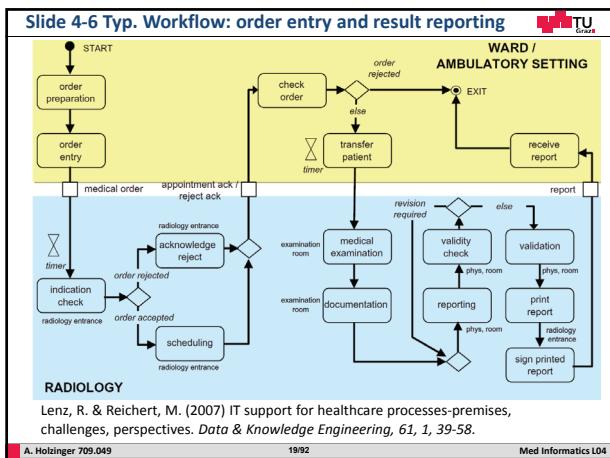
Leitung = Nervenleitung, Abteilungsleitung, Stromleitung, Wasserleitung, Harnleitung, ...

<http://www.medizinische-abkuerzungen.de/suche.html>

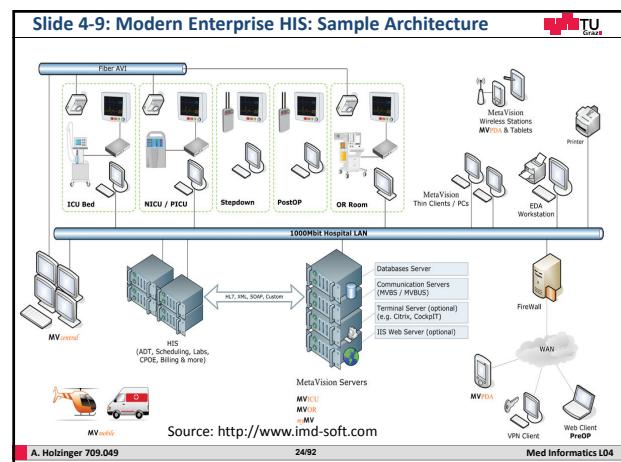
A. Holzinger 709.049 17/92 Med Informatics L04

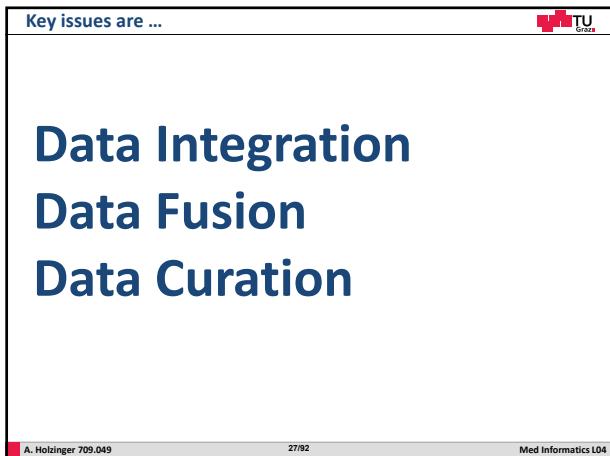
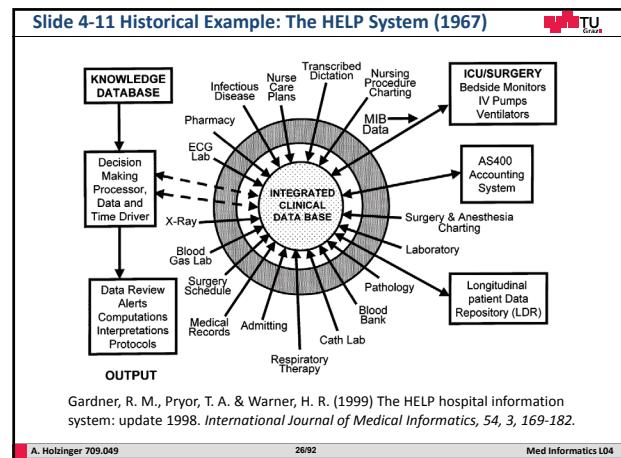
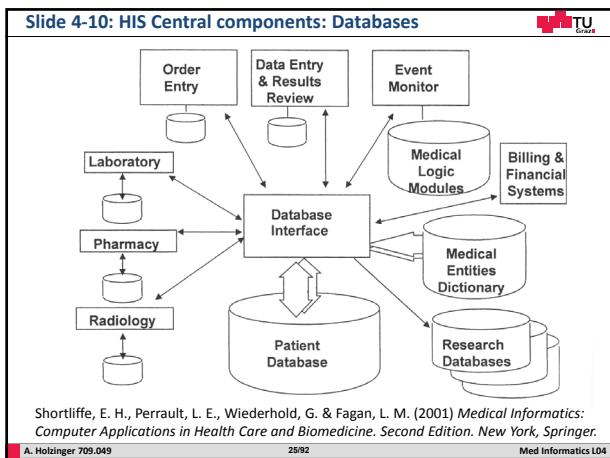
Hospital workflows are also complex ...

A. Holzinger 709.049 18/92 Med Informatics L04

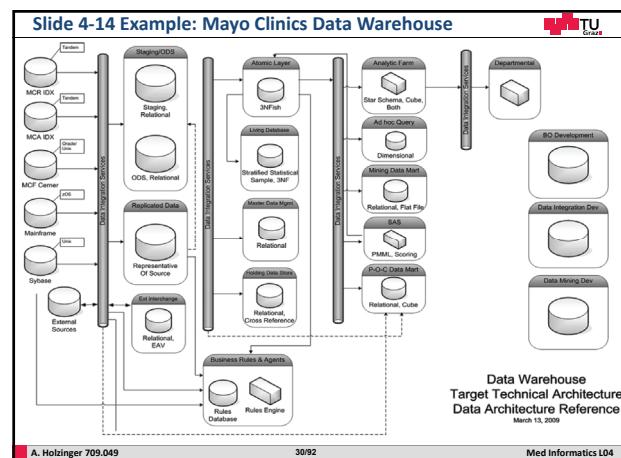
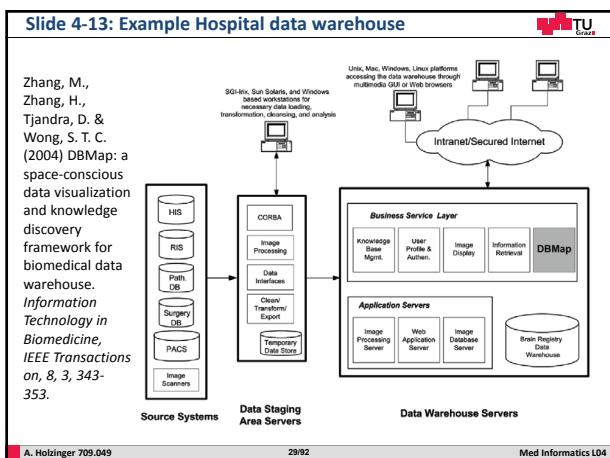


What is the architecture of an hospital information system?





- Slide 4-12 Database – fundamental terms and definitions**
- **Database (DB)** is the organized collection of data through a certain data structure (e.g. hash-table, adjacency matrix, graph structure, etc.).
 - **Database management system (DBMS)** is software which operates the DB. Well known DBMSs include: Oracle, IBM DB2, Microsoft SQL Server, Microsoft Access, MySQL, SQLite. Examples for Graph Databases include InfoGrid, Neo4j, or BrightstarDB.
 - The used DB is not generally portable, but different DBMSs can inter-operate by using standards such as SQL and ODBC.
 - **Database system (DBS)** = DB + DBMS. The term database system emphasizes that data is managed in terms of accuracy, availability, resilience, and usability.
 - **Data warehouse (DWH)** is an integrated repository used for reporting and long term storage of analysis data.
 - **Data Marts (DM)** are access layers of a DWH and are used as temporary repositories for data analysis.
- A. Holzinger 709.049 28/92 Med Informatics L04



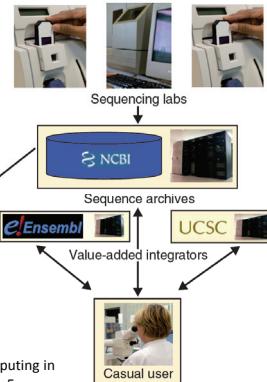
What about cloud-based Information Systems?

A. Holzinger 709.049

31/92

Med Informatics L04

Slide 4-15: Traditional Genome Information System



Stein, L. D. (2010) The case for cloud computing in genome informatics. *Genome Biology*, 11, 5.

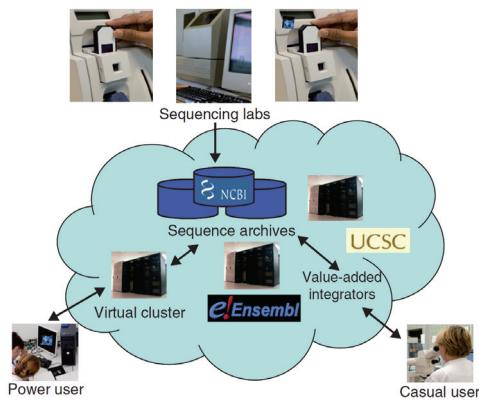
A. Holzinger 709.049

32/92

Med Informatics L04

Slide 4-16: Genome Info Ecosystem cloud computing

Stein, L. D.
(2010) The
case for cloud
computing in
genome
informatics.
*Genome
Biology*, 11, 5.



A. Holzinger 709.049

33/92

Med Informatics L04

Slide 4-17: Example Clinical Cloud Computing: PACS Cloud

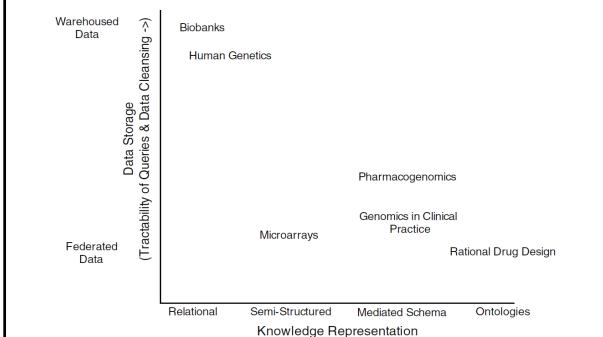
Bastiao-Silva, L. A., Costa, C., Silva, A. & Oliveira, J. L. (2011) A PACS Gateway to the Cloud. 6th Iberian Conference on Information Systems and Technologies (CISTI). 1-6.

A. Holzinger 709.049

34/92

Med Informatics L04

Slide 4-18: Federated Data vs. Warehoused Data



Louie, B., Mork, P., Martin-Sanchez, F., Halevy, A. & Tarczy-Hornoch, P. (2007) Data integration and genomic medicine. *Journal of Biomedical Informatics*, 40, 1, 5-16.

A. Holzinger 709.049

35/92

Med Informatics L04

What is the difference between hospital databases and Biomedical databases?

Slide 4-19: Biomedical databases ...

- ... are libraries of life science data, collected from scientific experiments and computational analyses.
- ... contain (clinical, biological, ...) data from clinical work, genomics, proteomics, metabolomics, microarray gene expression, phylogenetics, etc.
- Examples:
 - Text: e.g. PubMed, OMIM (Online Mendelian Inheritance in Man);
 - Sequence data: e.g. Entrez, GenBank (DNA), UniProt (protein).
 - Protein structures: e.g. PDB, Structural Classification of Proteins (SCOP), CATH (Protein Structure Classification);

A. Holzinger 709.049

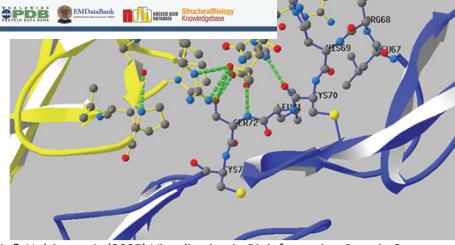
37/92

Med Informatics L04

Slide 4-20 Example Database: PDB

RCSB PDB Deposit Search Visualize Analyze Download

PDB An Information Portal to 113331 Biological Macromolecular Structures



Wiltgen, M. & Holzinger, A. (2005) Visualization in Bioinformatics: Protein Structures with Physicochemical and Biological Annotations. In: Central European Multimedia and Virtual Reality Conference, Prague, Czech Technical University (CTU), 69-74

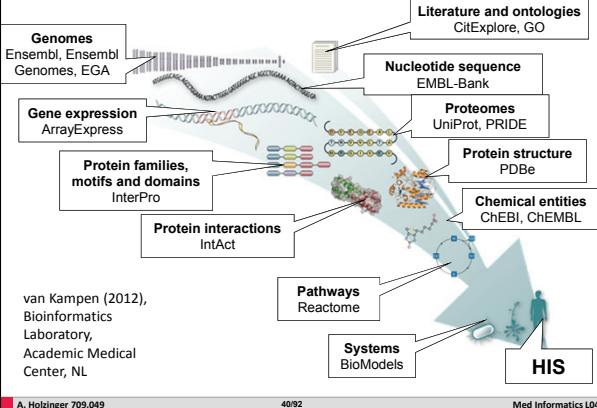
A. Holzinger 709.049

38/92

TU Graz

Med Informatics L04

Slide 4-21 Databases: From Molecules to Systems



A. Holzinger 709.049

40/92

Med Informatics L04

Slide 4-22: Example Genome Database: Ensembl

TU Graz

A. Holzinger 709.049

41/92

Med Informatics L04

Slide 4-23 Ex. Gene Expression Database: ArrayExpress

TU Graz

<http://www.ebi.ac.uk/arrayexpress/>

A. Holzinger 709.049

42/92

Med Informatics L04

Slide 4-24: Example Protein Interaction Database: IntAct

IntAct provides a freely available, open source database system and analysis tools for molecular interaction data. All interactions are derived from literature curation or direct user submissions and are freely available. To perform a search in the IntAct database use the search box above.

Publications: 12151 Experiments: 31292 Interactions: 434941 Interactors: 79905

Citing IntAct

<http://www.ebi.ac.uk/intact/>

A. Holzinger 709.049 43/92 Med Informatics L04

Slide 4-25: Example for Systems Database: BioModels

BioModels Database

McAuley et al., (2012). A whole-body mathematical model of cholesterol metabolism and its age-associated dysregulation.

October 2013 model of the month by Nick Judy

Original model: BIOMODELSDB/BIOMODEL/2012/01

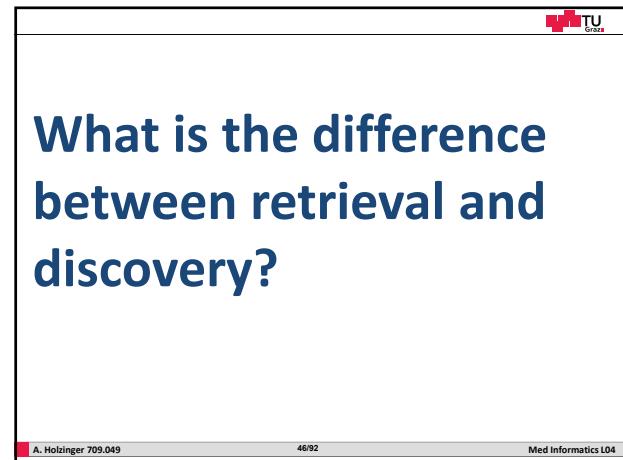
The hepatic system is central in cholesterol metabolism, with the liver able to synthesize VLDs (very low density lipoproteins), which are converted into IDLs (intermediate density lipoproteins) (IDLs) through the addition of triglycerides. Both IDLs and VLDs are taken up by the liver via the LDL receptor (LDLR). The intestine also takes up VLDs and IDLs, and releases VLDLs into the blood. The hepatic receptor is transcriptionally regulated by intracellular cholesterol levels.

It has been demonstrated that: a) There is age-associated decline in the clearance rate of LDL-C from the blood, which may contribute to the number of hepatic LDLs; b) Intestinal cholesterol absorption increases with age in some species.

In this paper, the authors take a mechanistic approach to construct a model, with these observations in mind, to predict the effect of age on cholesterol metabolism. The model is able to predict the effect of age on cholesterol levels in the blood over a period of 10 years. The model incorporates dietary cholesterol absorption in the intestine, and hepatic LDL-C clearance from the plasma [1]. BIOMODELSDB/BIOMODEL/2012/01 consists of 6 compartments (Figure 1), and is composed of a series of coupled ODEs.

<http://www.ebi.ac.uk/biomodels-main/>

A. Holzinger 709.049 44/92 Med Informatics L04

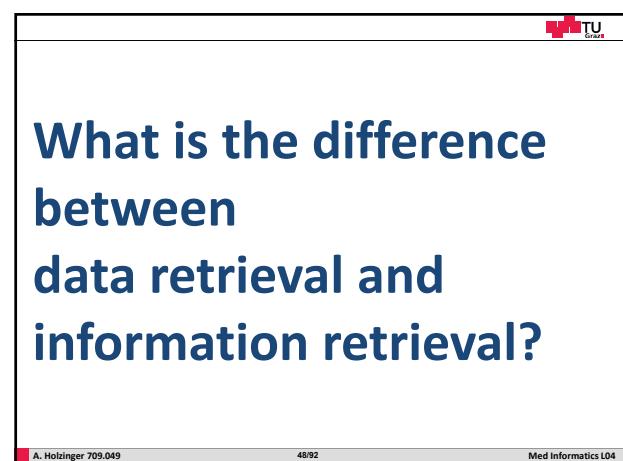


Slide 4-26: Data Mining/KDD versus Information Retrieval

Fayyad, U., Piatetsky-Shapiro, G. & Smyth, P. 1996. The KDD process for extracting useful knowledge from volumes of data. *Communications of the ACM*, 39, (11), 27-34.

Baeza-Yates, R. & Ribeiro-Neto, B. 2011. *Modern Information Retrieval: The Concepts and Technology behind Search*, Harlow et al., Pearson.

A. Holzinger 709.049 47/92 Med Informatics L04



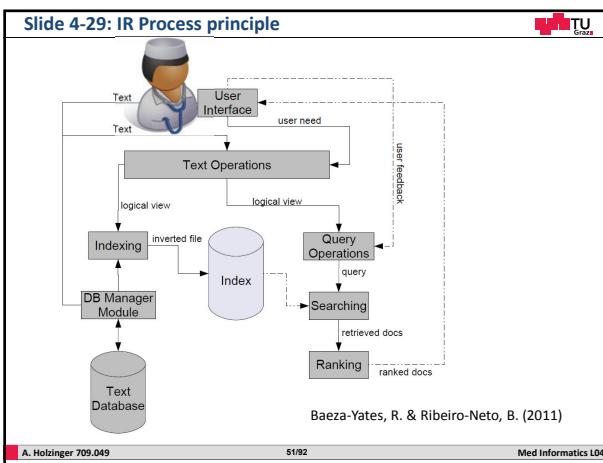
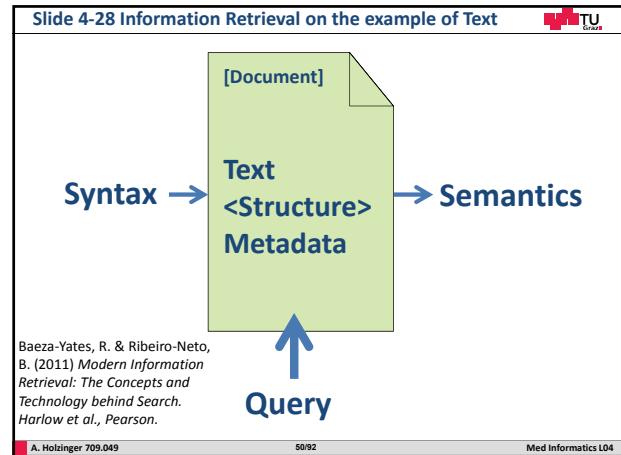
Slide 4-27: Data retrieval (DR) vs. Information retrieval (IR)

- IR is used to satisfy the end-users' information needs.
- Def.: IR deals with the representation, storage, organization of and access to information objects.

Factor	Data Retrieval (DR)	Information Retrieval (IR)
Model	Deterministic	Probabilistic
Matching	Exact match	Partial (best match)
Inference	Deduction	Induction
Classification	Monothetic*	Polythetic**
Query language	Artificial (abstract)	Natural
Query specification	Must be complete	Can be incomplete
Items wanted	matching	relevant
Error response	sensitive	insensitive

*Monothetic = type in which all members are identical on all characteristics;
**Polythetic = type in which all members are similar, but not identical;
Van Rijsbergen, C. J. (1979) *Information Retrieval (Second Edition)*. London, Butterworths.

A. Holzinger 709.049 49/92 Med Informatics L04



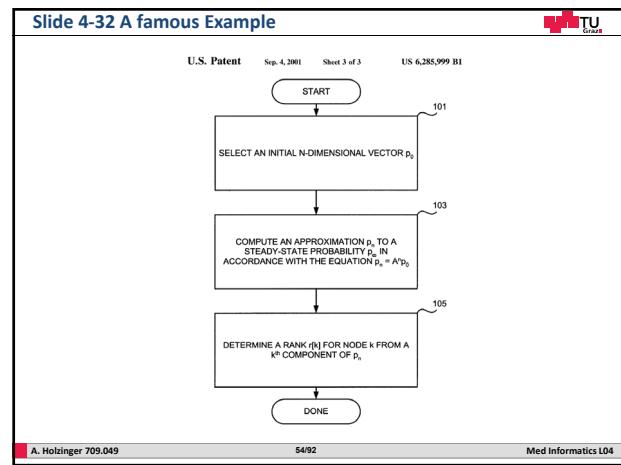
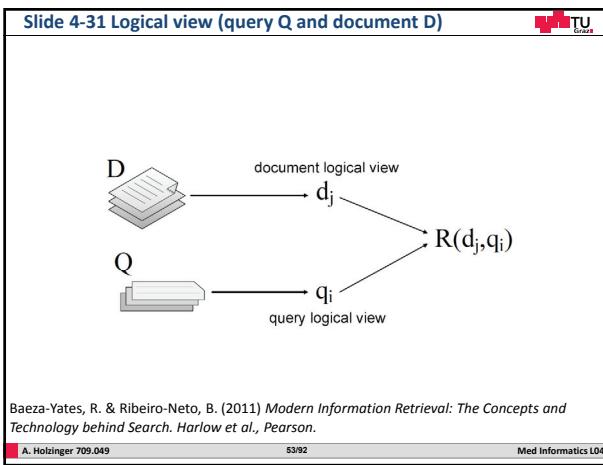
Slide 4-30: Formal Description of IR Models

Definition: Let the IR Model be a quadruple $\{D, Q, \mathcal{F}, R(q_i, d_j)\}$

- D** is a set composed of logical views (representation component) of the **documents** within a collection;
- Q** is a set of logical views (representation component) of the user information needs (these are called "**queries**");
- F** is a framework for modeling document representations, queries and their relationships (reasoning component);
This includes sets and Boolean relations, vectors and linear algebra operations, sample spaces and probability distributions;
- R** (q_i, d_j) is a ranking function that associates a real number with a query representation $q_i \in Q$ and a document representation $d_j \in D$.
Such ranking defines an ordering among the docs with regard to the query q_i

Baeza-Yates, R. & Ribeiro-Neto, B. (2011) *Modern Information Retrieval: The Concepts and Technology behind Search*. Harlow et al., Pearson.

A. Holzinger 709.049 52/92 Med Informatics L04



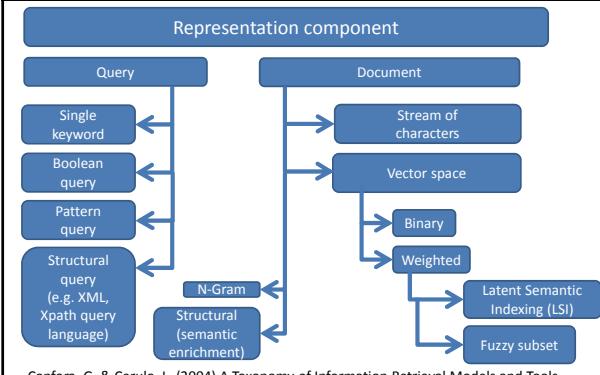
Remember: We have two components: Representation and Reasoning component

A. Holzinger 709.049

55/92

Med Informatics L04

Slide 4-33: Taxonomy of Information Retrieval Models 1/3

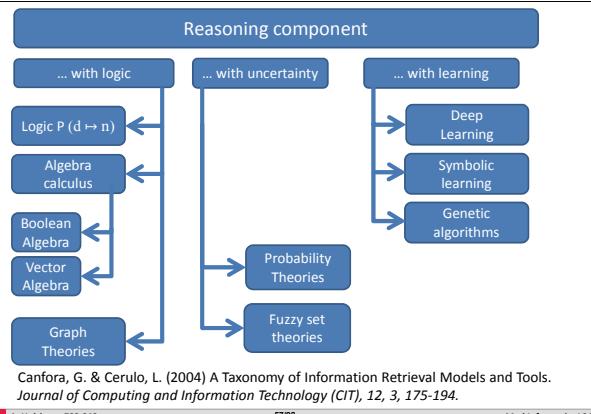


A. Holzinger 709.049

56/92

Med Informatics L04

Slide 4-34: Taxonomy of Information Retrieval Models 2/3

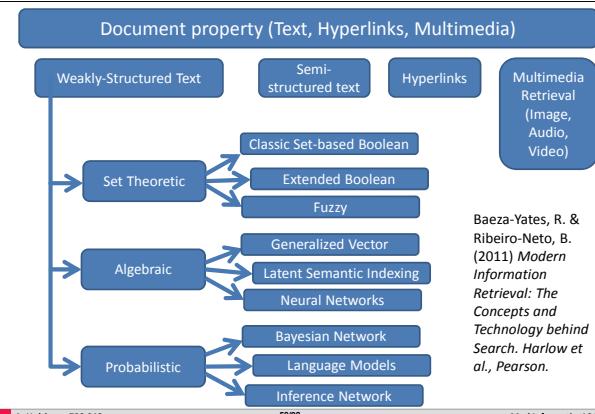


A. Holzinger 709.049

57/92

Med Informatics L04

Slide 4-35: Taxonomy of Information Retrieval Models 3/3



A. Holzinger 709.049

58/92

Med Informatics L04

Slide 4-36: Set Theoretic Example: Boolean Model

- Documents and queries are represented as a set of index terms; the queries are Boolean expressions (AND, OR, NOT);

"For the Boolean model, the index term weight variables are all binary i.e., $\omega_{i,j} \in \{0, 1\}$. A query q is a conventional Boolean expression. Let \vec{q}_{dnf} be the disjunctive normal form for the query q . Further, let \vec{q}_{cc} be any of the conjunctive components of \vec{q}_{dnf} . The similarity of a document d_j to the query q is defined as

$$sim(d_j, q) = \begin{cases} 1 & \text{if } \exists \vec{q}_{cc} | (\vec{q}_{cc} \in \vec{q}_{dnf}) \wedge (\forall k_i, g_i(\vec{q}_j) = g_i(\vec{q}_{cc})) \\ 0 & \text{otherwise.} \end{cases}$$

If $sim(d_j, q) = 1$ then the Boolean model predicts that the document d_j is relevant to the query q (it might not be). Otherwise, the prediction is that the document is not relevant."

Baeza-Yates, R. & Ribeiro-Neto, B. (2011)

A. Holzinger 709.049

59/92

Med Informatics L04

Slide 4-37: Set Theor. Model - Boolean Model Pros & Cons

Advantages	Disadvantages
Easy to understand	No partial matches
Exact formalism	The "bag-of-words" representation does not accurately consider the semantics of documents *
Query language is expressive	Query language is complicated
	Retrieved documents cannot be ranked

*) refer to: Vallet, D., Fernández, M. & Castells, P. (2005) An Ontology-Based Information Retrieval Model. In: Gómez-Pérez, A. & Euzenat, J. (Eds.) *The Semantic Web: Research and Applications*. Berlin, Heidelberg, Springer, 103-110.

A. Holzinger 709.049

60/92

Med Informatics L04

Slide 4-38 Example Algebraic Model: Vector Space Model

$D = \langle d_1, d_2, \dots, d_n \rangle$ (collection of medical docs)
 $d_i = t_1, t_2, \dots, t_k$ (every document consists of terms)
 Now we carry out a document transformation and get vectors:

$$w_{i,j} = \begin{cases} 1, & t_i \in d_j \\ 0, & t_i \notin d_j \end{cases} \rightarrow d_j = (0, 1, 1, 0, 1, \dots, 1)^T$$

Now we count the frequency of the terms and get:

$$w_{i,j} = \begin{cases} (1 + \log f_{i,j}) * \log \frac{N}{n_i}, & \text{if } f_{i,j} > 0 \\ 0 & \text{otherwise} \end{cases}$$

A. Holzinger 709.049

61/92

Med Informatics L04

Slide 4-39: As a result we get a matrix ...

$$D_{m \times n} = \begin{pmatrix} w_{1,1} & w_{1,2} & \cdots & w_{1,n-1} & w_{1,n} \\ w_{2,1} & w_{2,2} & & w_{2,n-1} & w_{2,n} \\ \vdots & & \ddots & & \vdots \\ w_{m-1,1} & w_{m-1,2} & & w_{m-1,n-1} & w_{m-1,n} \\ w_{m,1} & w_{m,2} & \cdots & w_{m,n-1} & w_{m,n} \end{pmatrix}$$

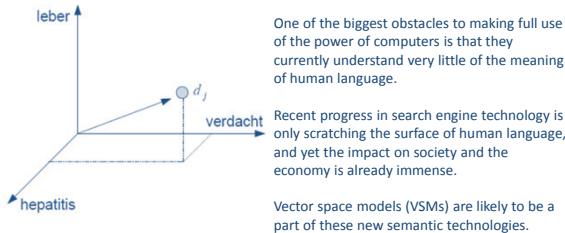
Salton, G., Wong, A. & Yang, C. S. 1975. Vector-Space Model for automatic indexing. *Communications of the ACM*, 18, (11), 613-620.

A. Holzinger 709.049

62/92

Med Informatics L04

Slide 4-40: d_j can thus be seen as a point in n-dim space



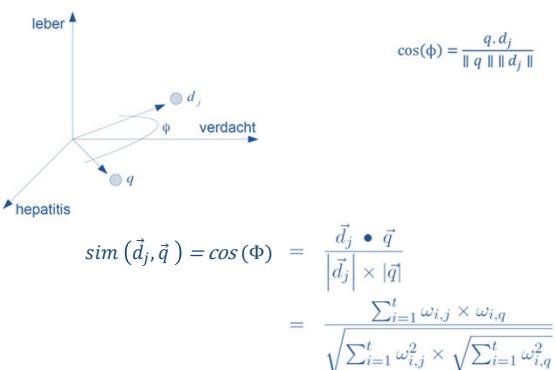
Turney, P. D. & Pantel, P. 2010. From frequency to meaning: Vector space models of semantics. *Journal of artificial intelligence research*, 37, (1), 141-188. Survey article 922 citations yet ...

A. Holzinger 709.049

63/92

Med Informatics L04

Slide 4-41: Use the cos-similarity for ranking similar docs



A. Holzinger 709.049

64/92

Med Informatics L04

Slide 4-42: Algebraic Vector Space Model: Pros & Cons

Advantages	Disadvantages
Easy to understand	Higher effort to calculate similarity
Partial matches possible	The "bag-of-words" representation does not accurately consider the semantics of documents *
Sorting of documents by rank	
Using term weighting schemes	

* refer to: Vallet, D., Fernández, M. & Castells, P. (2005) An Ontology-Based Information Retrieval Model. In: Gómez-Pérez, A. & Euzenat, J. (Eds.) *The Semantic Web: Research and Applications*. Berlin, Heidelberg, Springer, 103-110.

A. Holzinger 709.049

65/92

Med Informatics L04

Slide 4-43: Example: Probabilistic Model (Bayes' rule)

"For the probabilistic model, the index weight variables are all binary i.e., $\omega_{i,j} \in [0, 1]$, $\omega_{i,j} \in [0, 1]$. A query q is a subset of index terms. Let R be the set of documents known (or initially guessed) to be relevant. Let \bar{R} be the complement of R (i.e., the set of non-relevant documents). Let $P(R|\vec{d}_j)$ be the probability that the document d_j is relevant to the query q and $P(\bar{R}|\vec{d}_j)$ be the probability that d_j is non-relevant to q . The similarity $\text{sim}(d_j, q)$ of the document d_j to the query q is defined as the ratio

$$\text{sim}(d_j, q) = \frac{P(R|\vec{d}_j)}{P(\bar{R}|\vec{d}_j)}$$

$$\text{sim}(d_j, q) = \frac{P(\vec{d}_j|R) \times P(R)}{P(\vec{d}_j|\bar{R}) \times P(\bar{R})}$$



$$\text{sim}(d_j, q) \sim \frac{(\prod_{g_i(\vec{d}_j)=1} P(k_i|R)) \times (\prod_{g_i(\vec{d}_j)=0} P(\bar{k}_i|R))}{(\prod_{g_i(\vec{d}_j)=1} P(k_i|\bar{R})) \times (\prod_{g_i(\vec{d}_j)=0} P(\bar{k}_i|\bar{R}))}$$

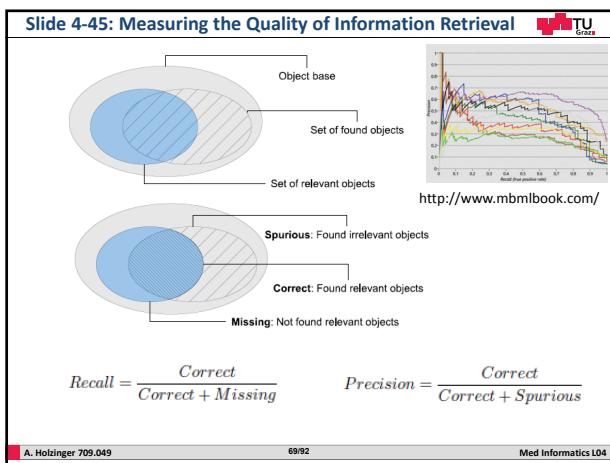
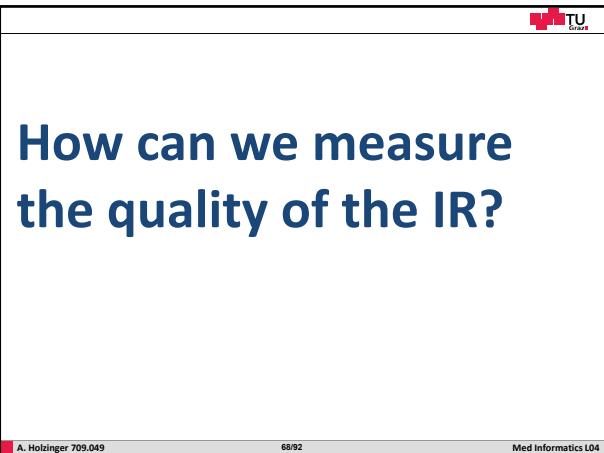
A. Holzinger 709.049

66/92

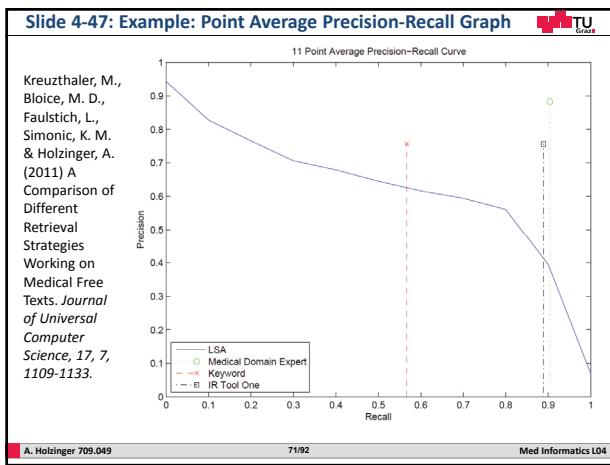
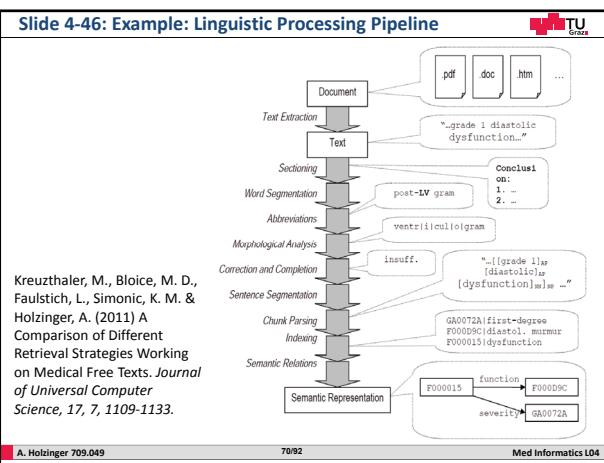
Med Informatics L04

Slide 4-44: Probabilistic Model: Pros & Cons	
Advantages	Disadvantages
Documents can be ranked by relevance	It is a binary model (\rightarrow binary weights)
	The index terms are assumed to be independent and a lack of document normalization
	There is a need to guess the initial separation of documents into relevant and non-relevant sets

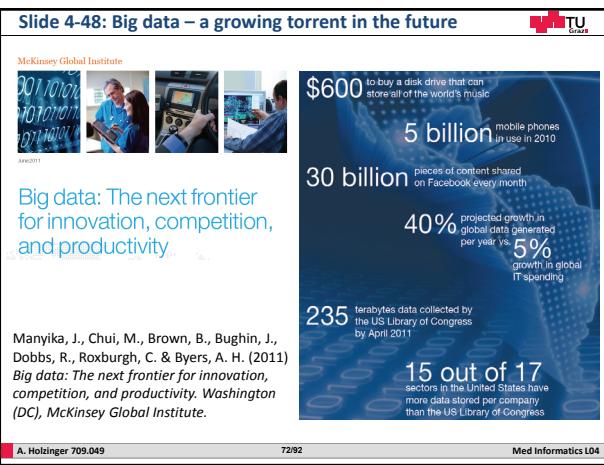
A. Holzinger 709.049 67/92 Med Informatics L04



A. Holzinger 709.049 68/92 Med Informatics L04



A. Holzinger 709.049 71/92 Med Informatics L04



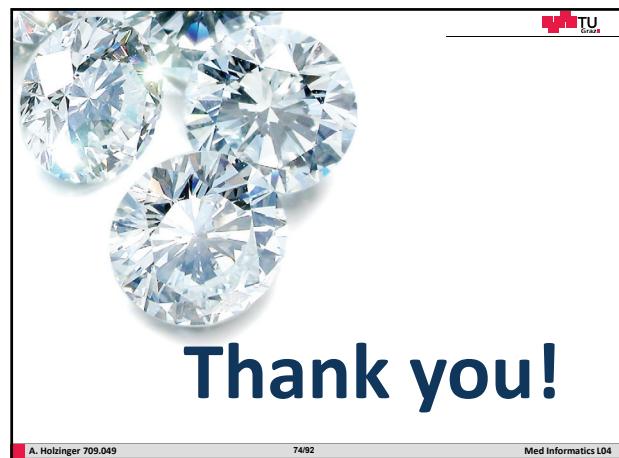
A. Holzinger 709.049 72/92 Med Informatics L04

Not big data is the real challenge ... but

Complex Data!

**What is interesting?
What is relevant?**

A. Holzinger 709.049 73/92 Med Informatics L04



Sample Questions (1)

- What is typical for medical workflows?
- How is the workflow in the clinical control loop?
- What does each shell in the Hospital Activity Shell model express?
- Of which main parts does the classic conceptual model of a Hospital Information System consist?
- What is a data mart?
- Why is the physician order entry a critical process?
- What is business intelligence in the context of a HIS?
- What is the difference between Information Extraction and Information Retrieval?
- Which differences exist between Data Retrieval and Information Retrieval?
- What advantages/disadvantages does cloud computing in health care have?
- What is a PACS cloud?

A. Holzinger 709.049 75/92 Med Informatics L04

Sample Questions (2)

- What is the purpose of the Protein Structure Database (PDB)?
- What advantages does a integrated HIS offer?
- What is the difference between monothetic data types and polythetic data types?
- What is the purpose of medical documentation?
- How does a typical medical document look?
- What are the big difficulties in medical documents?
- How can an Information Retrieval Model be formally described?
- What is the difference between a representation component and a reasoning component?
- What advantages/disadvantages does the Boolean model have?
- Describe the principles of the Vector space model!
- Which advantage does the Probabilistic model offer?
- What is the big disadvantage of an Ontology-Based Model?
- How can you determine the quality of information retrieval?

A. Holzinger 709.049 76/92 Med Informatics L04

Backup Slide: Example: Physician Order Entry (Paper)

PO	date	sig
REUPUSID SUSP 1MG/ML 10ML	10 mg [08:00]	
clazipride (1:1 water)	15 mg [18:00]	
N	START: 28-07-0410:22	
	STOP: 28-07-0415:15	
621C	verb. Short Medicator *33725	
LET OP - STOPDATUM IS INGEVULD		
41537		

Handwritten Physician Order:

Tramadol caps 50 mg 3 x per dag Test A/B

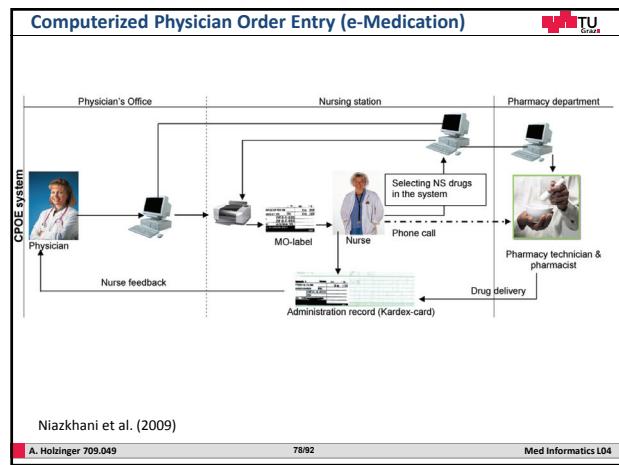
Handwritten Prescription:

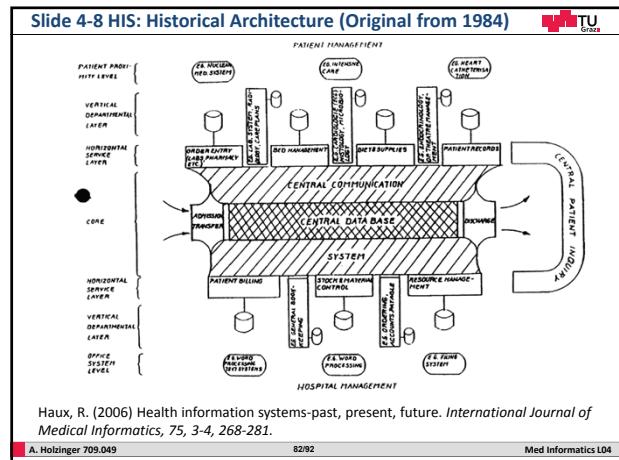
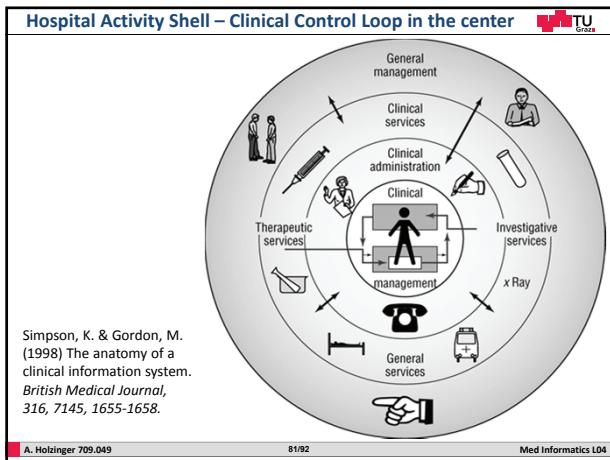
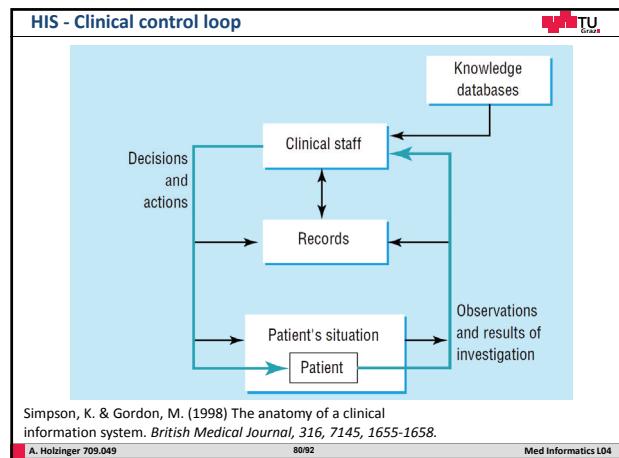
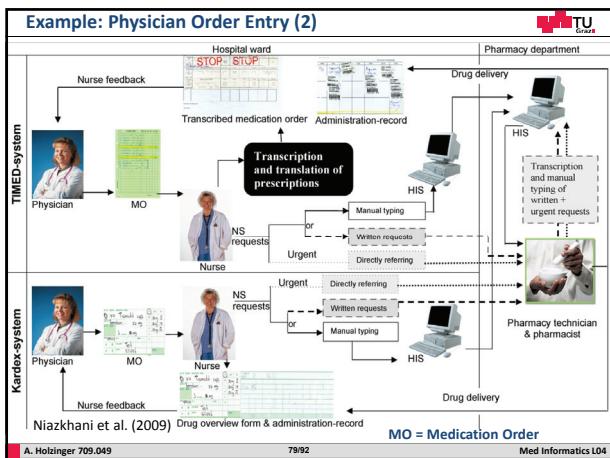
Tramadol 50 mg 30 mg 14 20 mg 22

Notes:

Niazkhani, Z., van der Sijs, H., Pirnejad, H., Redekop, W. K. & Aarts, J. (2009) Same system, different outcomes: Comparing the transitions from two paper-based systems to the same computerized physician order entry system. *Int. Journal of Medical Informatics*, 78, 3, 170-181.

A. Holzinger 709.049 77/92 Med Informatics L04





Example: Enterprise Data Modeling (EDM) at Mayo Clinic

Subjects = the highest level areas that define the activities of the enterprise (e.g. Individual)
Concepts = the collections of data that are contained in one or more subject areas (e.g., Patient, Provider, Employee, Referrer, Volunteer, etc.)
Business Information Models = the organization of the data that support the processes and workflows of the enterprise's defined Concepts.

Chute, C. G., Beck, S. A., Fisk, T. B. & Mohr, D. N. (2010) The Enterprise Data Trust at Mayo Clinic: a semantically integrated warehouse of biomedical data. *Journal of the American Medical Informatics Association*, 17, 2, 131-135.

A. Holzinger 709.049 83/92 Med Informatics L04

Backup: For your own experiments: www.care2x.org

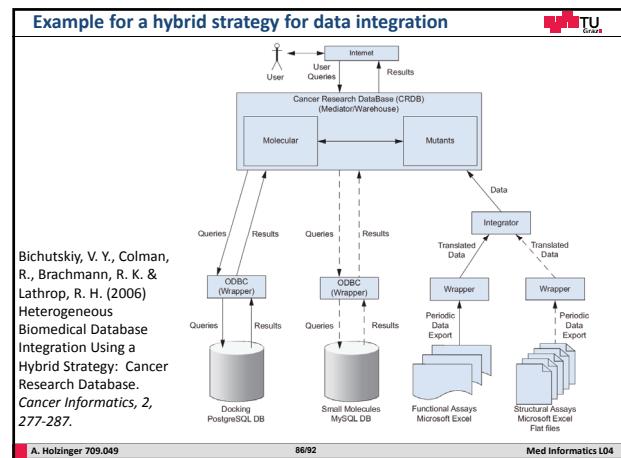
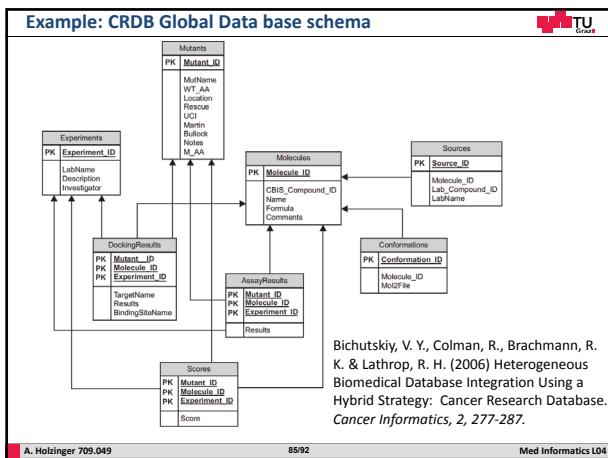
The screenshot shows the Care2x.org software interface for patient admission data. The main window displays a patient profile with the following details:

- Admission number: 2008010003
- Title: Mr.
- Family name: Bakaria
- Given name: Parvez
- Date of birth: 05/06/2003
- Address: Mahabir 1224
- City: Doha
- Country: Qatar
- Diagnosis: Ill fits fit
- Therapy: I.V. nnm
- Billing Type: Health Fund
- Insurance: Advance Bank
- Admitted by: Elsyda Latifah

Checkboxes at the bottom include:

- I need to admit a patient
- I am looking for a patient
- I need to research in the archive

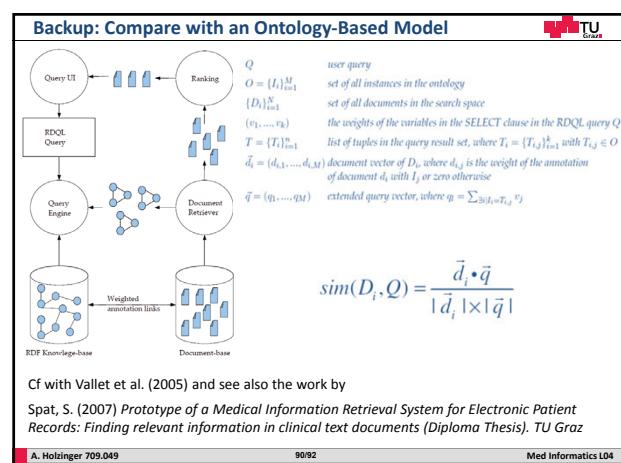
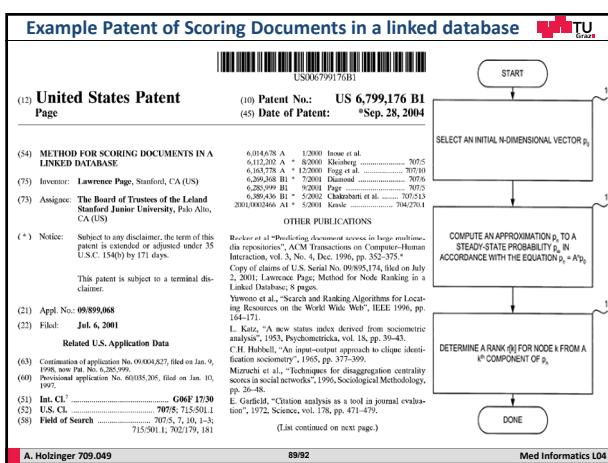
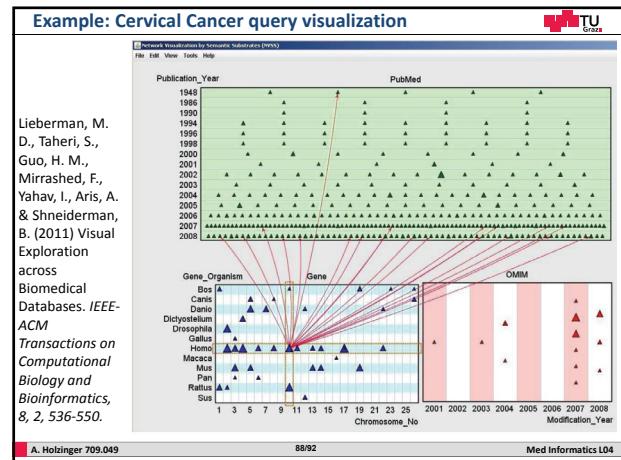
Care2x 2003 beta 1.0.07 - License - Contact - Our Privacy Policy - Legal - Credits
Page generation time: 2.779150327768
A. Holzinger 709.049 84/92 Med Informatics L04



Slide 4.22 Example Database: Protein Structure Data Bank

The screenshot shows the PDB website interface. The main search bar has the query "Biological Macromolecular Resource". Below it, a search result for "Biological Macromolecular Resource" is displayed, including a thumbnail of a protein structure, the title "Biological Macromolecular Resource", and a link to "Full Article | Archive | PDB Structural Biology Knowledgebase".

A. Holzinger 709.049 87/92 Med Informatics L04



Ontology Based Model: Pros & Cons	
Advantages	Disadvantages
Documents can be ranked by relevance	Works only if adequate knowledge base is available
Semantics of the documents can be considered	Only usable for already known facts – completely useless to discover new items
Model outperforms classic IR models	Big effort to build and maintain a adequate knowledge base

A. Holzinger 709.049

91/92

Med Informatics L04

Some Useful Links	
▪ http://www.library.tufts.edu/hsl/resources/dbases.html	
▪ http://www.ncbi.nlm.nih.gov/omim	
▪ http://lucene.apache.org/java/docs/	
▪ http://www.dcs.gla.ac.uk/Keith/Preface.html	
▪ http://hive.apache.org/	
▪ http://www.cs.waikato.ac.nz/ml/weka/	
▪ http://scikit-learn.sourceforge.net/stable/	
▪ http://www.eecs.wsu.edu/mgd/gdb.html	

A. Holzinger 709.049

92/92

Med Informatics L04