

Andreas Holzinger
VO 709.049 Medical Informatics
11.11.2015 11:15-12:45

Lecture 05

Semi structured, weakly structured data Graphs, Networks and Homologies

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



A. Holzinger 709.049 1/78 Med Informatics L05



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 709.049 278 Med Informatics L05

Keywords of the 5th Lecture

- Big data pools
- Complex networks
- Computational graph representation
- Electronic patient record (EPR)
- Homology modeling
- Macroscopic structures
- Medical documentation
- Metabolic network
- Microscopic structures
- Network metrics
- Structural data dimension
- Topological structures

A. Holzinger 709.049 3/78 Med Informatics L05



Advance Organizer (1/3) A-G

- **Adjacency matrix** = simplest form of computational graph representation, in which 0 or 1 denotes whether or not there is a directed edge from one node to another (in graph theory adjacent nodes in a graph are linked by an edge);
- **Artifacts** = not only a noise disturbance, which is contaminating and influencing the signal (surrogates) but also data which is wrong, however interpreted as to be reliable, consequently may lead to a wrong decision;
- **Computational graph representation** = e.g. by adjacency matrices
- **Data fusion** = data integration techniques that analyze data from multiple sources in order to develop insights in ways that are more efficient and potentially more accurate than if they were developed by analyzing a single source of data. Signal processing techniques can be used to implement some types of data fusion (e.g. combined sensor data in Ambient Assisted Living);
- **Global Distance Test (GDT)** = a measure of similarity between two protein structures with identical amino acid sequences but different tertiary structures. It is most commonly used to compare the results of protein structure prediction to the experimentally determined structure as measured by X-ray crystallography or protein NMR;
- **Graph theory** = study of mathematical structures to model relations between objects from a certain collection;
- **Graphs** = a hypothetical structure consisting of a series of nodes connected by weighted edges (graphs can be directed/undirected and stoichiometric/non-stoichiometric regarding interaction classes);

A. Holzinger 709.049 4/78 Med Informatics L05

Advance Organizer (2/3) H-P

- **Homology** = in mathematics (especially algebraic topology and abstract algebra), it is (ὅμοιος homos = "identical") a certain general procedure to associate a sequence of Abelian groups (i.e. does not depend on their order) or modules with a given mathematical object such as a topological space or a group;
- **Homology modeling** = comparative modeling of protein, refers to constructing an atomic-resolution model of the "target" protein from its amino acid sequence and an experimental three-dimensional structure of a related homologous protein (the "template"); in Bioinformatics, homology modeling is a technique that can be used in molecular medicine.
- **In silico** = via computer simulation, in contrast to *in vivo* (within the living) or *in vitro* (within the glass);
- **Multi-scale representation** = in a graph, nodes do not have to represent biological objects on the same scale, one node (e.g. a molecule) may have an edge connecting it to a node representing a cell or tissue (the edge indicates that the molecule exerts an effect on the cell/tissue);
- **Network** = graphs containing cycles or alternative paths;
- **Network analysis** = a set of techniques used to characterize relationships among discrete nodes in a graph or a network;
- **Network topology** = the shape or structure of a network;
- **Petri-Net** = a special class of graph, consisting of two general classes or node: place and transition nodes;
- **Predictive modeling** = a set of techniques in which a mathematical model is created or chosen to best predict the probability of an outcome (e.g. regression);
- **P-System** = addresses the slowness of Petri-nets

A. Holzinger 709.049 5/78 Med Informatics L05



Advance Organizer (3/3) R-V

- **Radius of a graph** = average minimum path length (biological networks are not arranged in a regular or symmetrical pattern);
- **Scale-free Topology** = ensures that there are very short paths between any given pair of nodes, allowing rapid communication between otherwise distant parts of the network (e.g. the Web has such a topology);
- **Semi-structured data** = does not conform with the formal structure of tables/data models assoc. with relational databases, but at least contains tags/markers to separate semantic elements and enforce hierarchies of records and fields within the data; aka schemaless or self-describing structure; the entities belonging to the same class may have different attributes even though they are grouped together;
- **Spatial analysis** = a set of techniques, applied from statistics, which analyze the topological, geometric, or geographic properties encoded in a data set;
- **Structural homology** = similar structure but different function;
- **Supervised learning** = machine learning techniques that infer a function or relationship from a set of training data (e.g. classification and support vector machines);
- **Time series analysis** = set of techniques from both statistics and signal processing for analyzing sequences of data points, representing values at successive times, to extract meaningful characteristics from the data;
- **Time series forecasting** = use of a model to predict future values of a time series based on known past values of the same or other series (e.g. structural modeling); decomposition of a series into trend, seasonal, and residual components, which can be useful for identifying cyclical patterns in the data;
- **Unstructured data** = complete randomness, noise; (wrongly, text is called unstructured, but there is some structure, too, so text data is a kind of weakly structured data);
- **Vertex degree** = within a topology, the numbers of edges connecting to a node;

A. Holzinger 709.049 6/78 Med Informatics L05

TU
Graz

Glossary

- ANSI = American National Standards Institute
- CD = cardiac development
- CDA = Clinical Document Architecture
- CHD = congenital heart disease
- CMM = Correlated motif mining
- DPI = Dossier Patient Integre' = integrated patient record
- E = Edge
- EPR = Electronic Patient Record
- G(V,E) = Graph
- GI = gastrointestinal
- HER = Electronic Health Record
- HL7 = Health Level 7
- KEGG = Kyoto Encyclopedia of Genes and Genomes
- NP = nondeterministic polynomial time
- OWL = Web Ontology Language
- PPI = Protein-Protein Interaction
- SGML = Standard Generalized Markup Language
- TF= Transcription factor
- TG = Target Gene
- V = Vertex
- XML = Extensible Markup Language

A. Holzinger 709.049

7/78

Med Informatics L05

Learning Goals ... at the end of the 5th lecture you ...

TU Ceilidh

- ... have an idea of the **complexity of data** in biomedical informatics
- ... are aware of the various **contents** of Electronic Patient Records
- ... have seen some application examples of **network structures** from both macro-cosmos and micro-cosmos and are fascinated about it;
- ... have a rough overview about some basics of how to **get point clouds** out of data sets
- ... have an understanding of the challenges of **network science**

Key Problems



- Automated Machine Learning algorithms need much training data – focus is on adjusting model parameters without fully **understanding the data** that the learning algorithm is modeling [1]
- Curse of dimensionality [2] – need for privacy and **anonymization** [3] (see lecture 11)
- **Weakly structured data** [4]

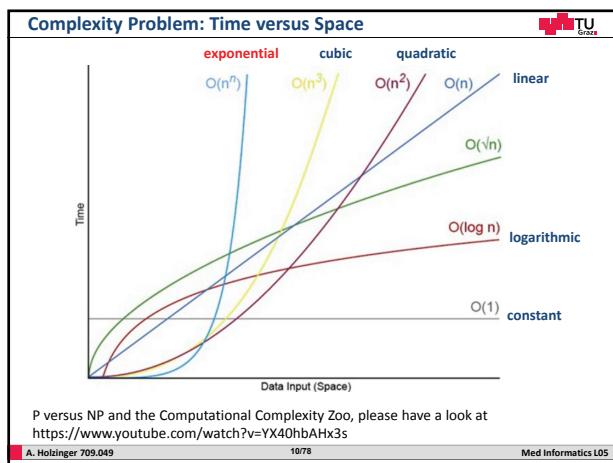
[1] Smith, M. R., Martinez, T. & Giraud-Carrier, C. 2014. An instance level analysis of data complexity. *Machine learning*, 95, (2), 225-256.

[2] Friedman, J. H. 1997. On bias, variance, 0/1-loss, and the curse-of-dimensionality. Data mining and knowledge discovery, 1, (1), 55-77.

[3] Aggarwal, C. C. On k-anonymity and the curse of dimensionality. Proceedings of the 31st international conference on Very large data bases VLDB, 2005. 901-909

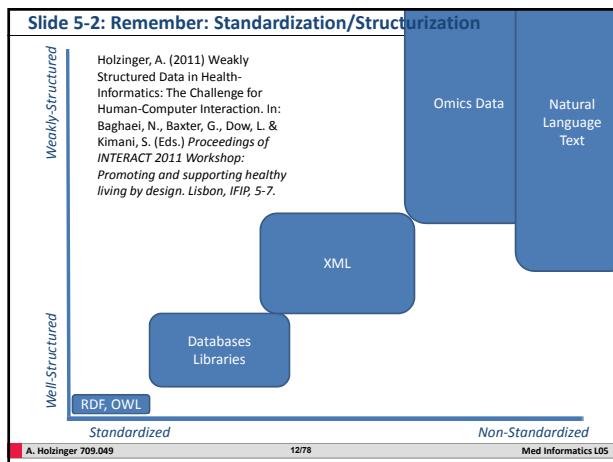
[4] Holzinger, A., Stocker, C. & Dehmer, M. 2014. Big Complex Biomedical Data: Towards a Taxonomy of Data. In: CCIS 455. Berlin Heidelberg: Springer pp. 3-18.

A. Holzinger 709.049 9/78 Med Informatics L05



Complex and High dimensional

Geschwind, D. H. & Konopka, G. 2009. Neuroscience in the era of functional genomics and systems biology. *Nature*, 461, (7266), 908-915.



Slide 5-3: Example: Well-Structured Data

A. Holzinger 709.049 13/78 Med Informatics L05

Slide 5-4: Example: Semi-structured Data: XML

```
<?xml version="1.0"?>
<patient>
  <patient-id>11111</patient-id>
  <Name>Chen</Name>
  <Date of Birth>1.1.1900</Date of Birth>
  <diagnosis>
    <code>123</code>
    <diagnosistext>Myocardinfarct</diagnosistext>
  </diagnosis>
</patient>
```

Holzinger, A. (2003) Basiswissen IT/Informatik. Band 2: Informatik. Das Basiswissen für die Informationsgesellschaft des 21. Jahrhunderts. Wuerzburg, Vogel Buchverlag.

A. Holzinger 709.049

14/78

Med Informatics L05

Slide 5-5 Example: Generic XML template for a med. report

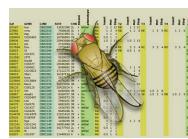
Rassinoux, A.-M., Lovis, C., Baud, R. & Geissbuhler, A. (2003) XML as standard for communicating in a document-based electronic patient record: a 3 years experiment. *International Journal of Medical Informatics*, 70, 2-3, 109-115.

A. Holzinger 709.049 15/78 Med Informatics L05

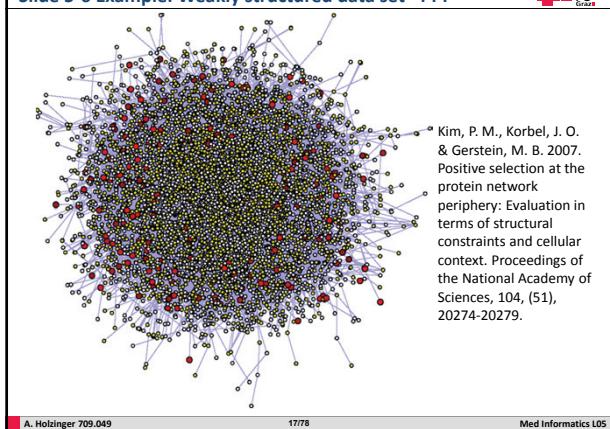
Slide 5-6 Comparison of XML - RDF/OWL in Bioinformatics

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 16/78 Med Informatics L05

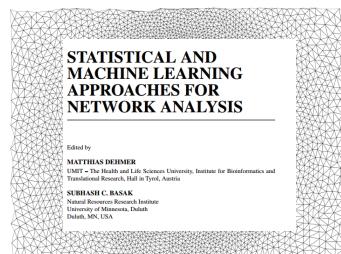


Slide 5-6 Example: Weakly structured data set - PPI



Network Science – Graph Theory

Networks = Graphs



<http://www.wired.com/tag/network-science/>

A. Holzinger 709.049

18/78

Med Informatics L05

Slide 5-7: Complex Biological Systems key concepts

- In order to understand complex biological systems, the three following key concepts need to be considered:
- (i) **emergence**, the discovery of links between elements of a system because the study of individual elements such as genes, proteins and metabolites is insufficient to explain the behavior of whole systems;
- (ii) **robustness**, biological systems maintain their main functions even under perturbations imposed by the environment; and
- (iii) **modularity**, vertices sharing similar functions are highly connected.

Network theory can largely be applied for biomedical informatics, because many tools are already available

A. Holzinger 709.049 19/78 Med Informatics L05

Slide 5-8: Networks on the Example of Bioinformatics

G(V, E) Graph
 $V \dots \text{vertex}$
 $E \dots \text{edge } \{a, b\}$
 $a, b \in V; a \neq b$

Hodgman, C. T., French, A. & Westhead, D. R. (2010) *Bioinformatics*. Second Edition. New York, Taylor & Francis.

A. Holzinger 709.049 20/78 Med Informatics L05

Slide 5-9: Computational Graph Representation

Adjacency ($\text{a}'\text{j}\text{a-s}^{\text{o}}\text{n(t)-s}^{\text{e}}$) Matrix $A = (a_{jk})$ $a_{jk} = \begin{cases} 1, & \text{if } (j, k) \in E \\ 0, & \text{otherwise} \end{cases}$

$a_{jk} = \begin{bmatrix} 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \end{bmatrix}$

Simple graph, symmetric, binary

Directed and weighted

For more information: Diestel, R. (2010) *Graph Theory*, 4th Edition. Berlin, Heidelberg, Springer.

A. Holzinger 709.049 21/78 Med Informatics L05

Slide 5-10: Example: Tool for Node-Link Visualization

Jean-Daniel Fekete http://wiki.cytoscape.org/InfoVis_Toolkit

Fekete, J.-D. The infovis toolkit. Information Visualization, INFOVIS 2004, 2004. IEEE, 167-174.

Excusus: Do not mix up Image Processing with Visualization (see L 09)

<p>Image Formation object in → image out</p>	<p>Image Processing image in → image out</p>
<p>Image Analysis image in → features out</p>	<p>Computer Graphics numbers in → image out</p>
<p>Computer Vision image in → interpretation out</p>	<p>Visualization image in → representation out</p>

Meijering, Erik & Cappellen, Gert (2006) *Biological Image Analysis* Primer, available via <http://www.imagescience.org/meijering/publications/1009/> Erasmus University Medical Center

Slide 5-11: Some Network Metrics (1/2)

Order = total number of nodes n ; Size = total number of links (a):

$$\sum_i \sum_j a_{ij}$$

Clustering Coefficient (b) = the degree of concentration of the connections of the node's neighbors in a graph and gives a measure of local inhomogeneity of the link density:

$$C_i = \frac{2t_i}{k(k_i - 1)}$$

$$C = \frac{1}{n} \sum_i C_i$$

Path length (c) = is the arithmetical mean of all the distances:

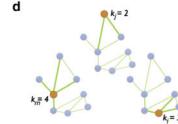
$$l = \frac{1}{n(n-1)} \sum_{i,j} d_{ij}$$

Costa, L. F., Rodrigues, F. A., Travieso, G. & Boas, P. R. V. (2007) Characterization of complex networks: A survey of measurements. *Advances in Physics*, 56, 1, 167-242.

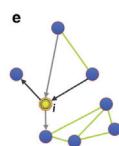
A. Holzinger 709.049 24/78 Med Informatics L05

Slide 5-12: Some Network Metrics (2/2)

- Centrality (d)** = the level of "betweenness-centrality" of a node i ("hub-node" in Slide 28);



- Nodal degree (e)** = number of links connecting i to its neighbors: $k_i = \sum_i a_{ij}$

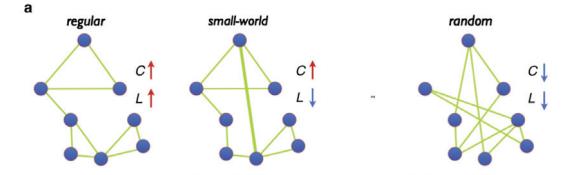


Modularity (f) = describes the possible formation of communities in the network, indicating how strong groups of nodes form relative isolated sub-networks within the full network (refer also to Slide 5-8).

A. Holzinger 709.049

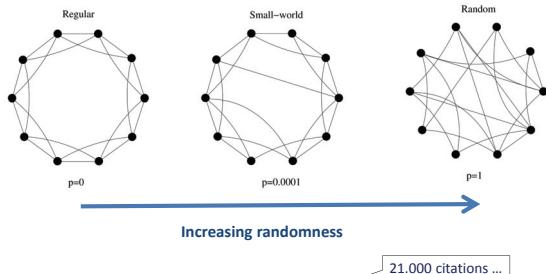
25/78

Med Informatics L05

Slide 5-13: Network Topologies

b Scale-free network
Modular network

A. Holzinger 709.049
26/78
Van Heuvel & Hulshoff (2010)
Med Informatics L05

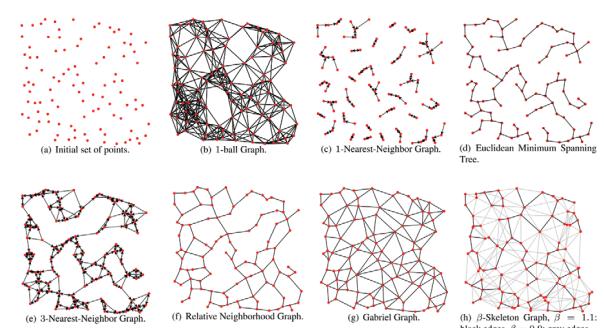
Slide 5-14: Small-World Networks

Watts, D. J. & Strogatz, S. (1998) Collective dynamics of small-world networks. *Nature*, 393, 6684, 440-442.
Milgram, S. 1967. The small world problem. *Psychology today*, 2, (1), 60-67.

A. Holzinger 709.049

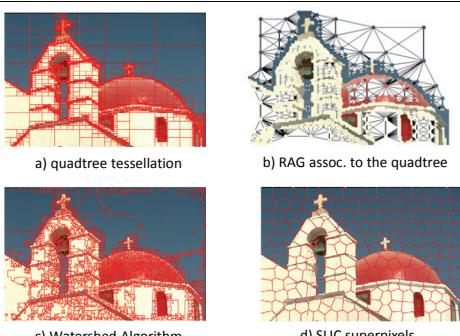
27/78

Med Informatics L05

Slide 5-15 Graphs from Point Cloud Data Sets

Lézoray, O. & Grady, L. 2012. Graph theory concepts and definitions used in image processing and analysis. In: Lézoray, O. & Grady, L. (eds.) *Image Processing and Analysing With Graphs: Theory and Practice*. Boca Raton (FL): CRC Press, pp. 1-24.

A. Holzinger 709.049
28/78
Med Informatics L05

Slide 5-16 Graphs from Images

Lézoray, O. & Grady, L. 2012. Graph theory concepts and definitions used in image processing and analysis. In: Lézoray, O. & Grady, L. (eds.) *Image Processing and Analysing With Graphs: Theory and Practice*. Boca Raton (FL): CRC Press, pp. 1-24.

A. Holzinger 709.049

29/78

Med Informatics L05

Slide 5-18 Example Watershed Algorithm

Algorithm 4.2 Watershed transform w.r.t. topographical distance based on image integration via the Dijkstra-Moore shortest paths algorithm.

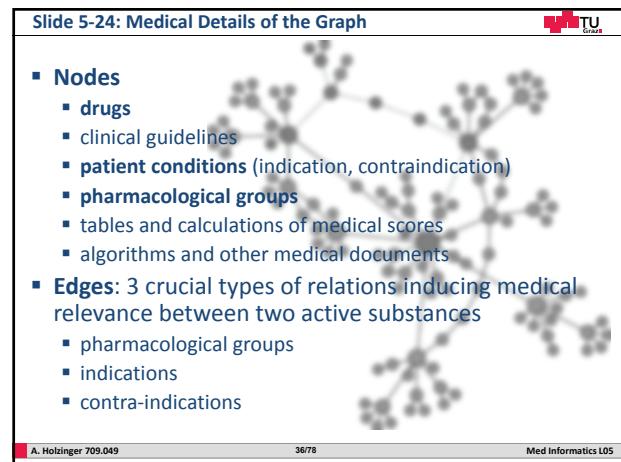
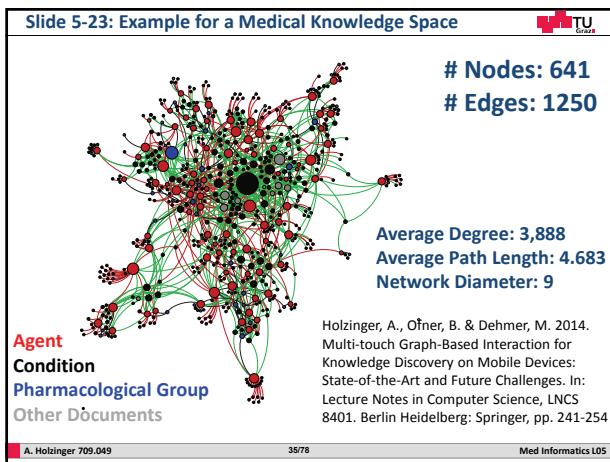
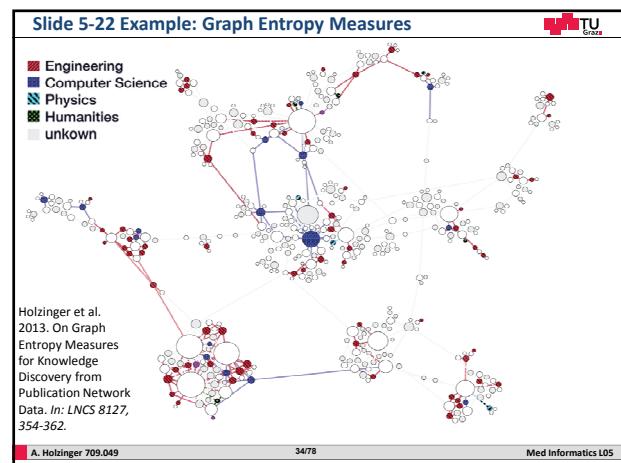
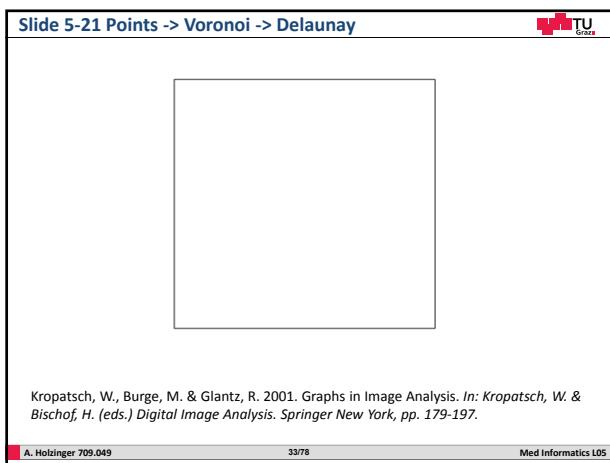
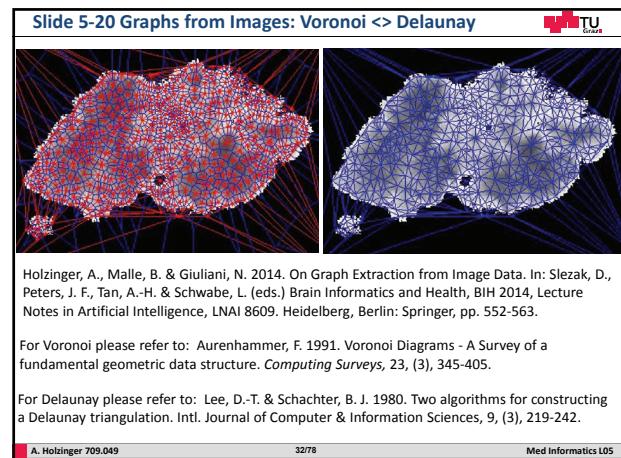
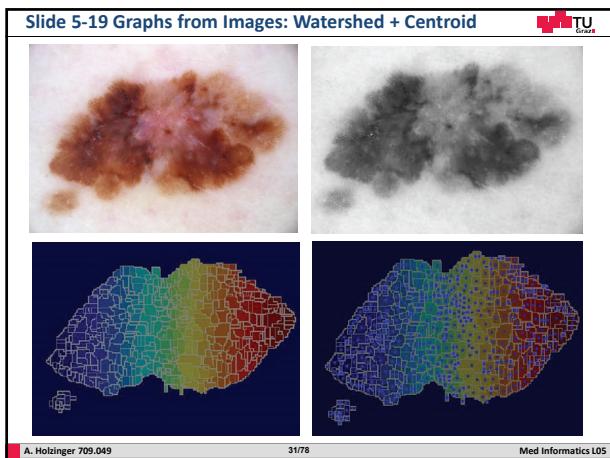
```

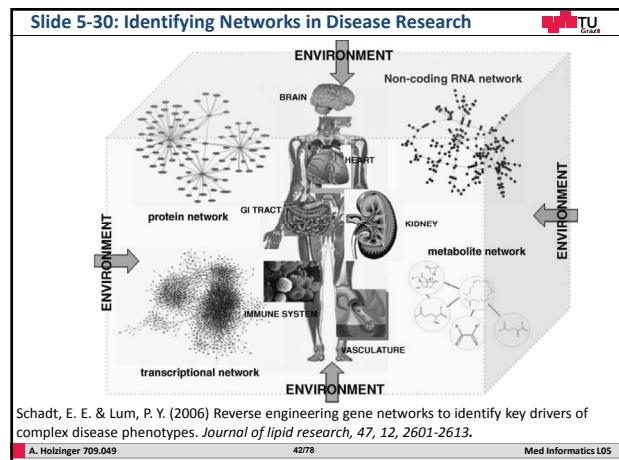
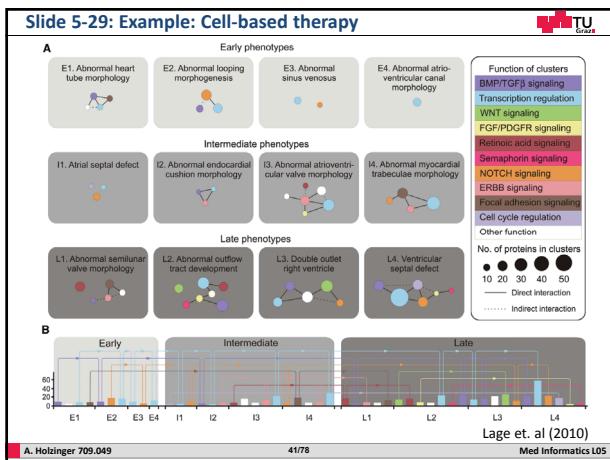
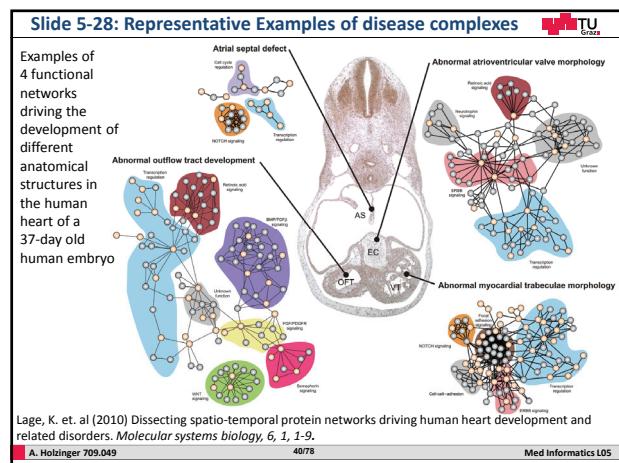
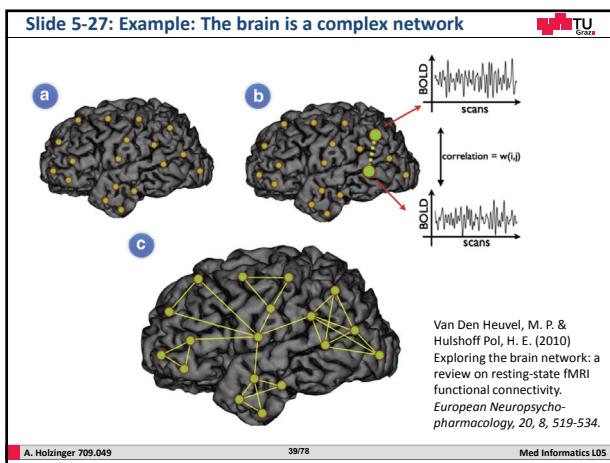
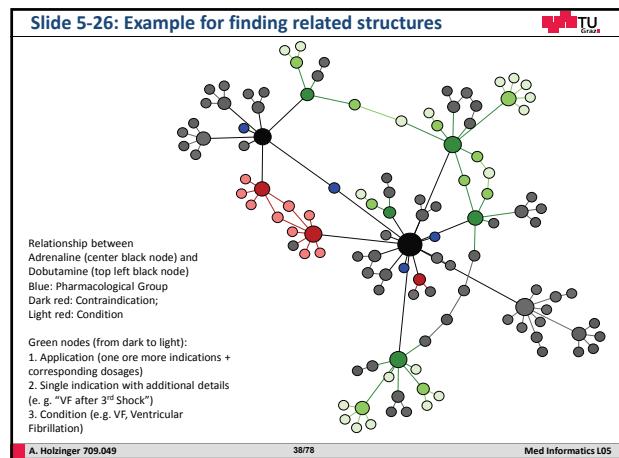
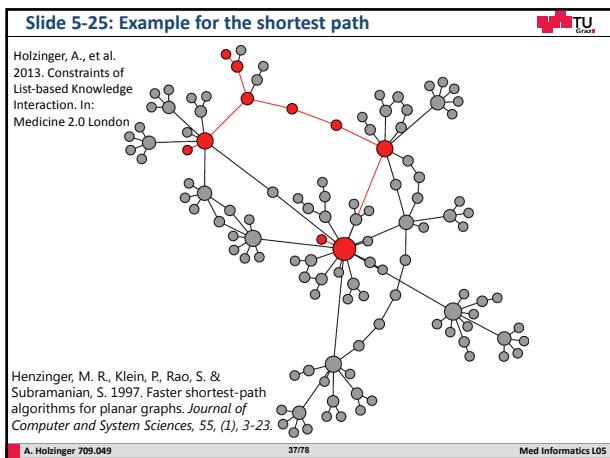
1: procedure ShortestPathWatershed;
2: INPUT: lower complete digital grey scale image  $G = (V, E, im)$  with cost function  $cost$ .
3: OUTPUT: labelled image  $lab$  on  $V$ .
4: #define WSHED 0
5: (* Uses distance image  $dist$ . On output,  $dist[v] = im[v]$ , for all  $v \in V$ . *)
6: 
7: for all  $v \in V$  do (* Initialize *)
8:    $lab[v] \leftarrow 0$ ;  $dist[v] \leftarrow \infty$ 
9: end for
10: for all local minima  $m_i$  do
11:   for all  $v \in m_i$  do
12:      $lab[v] \leftarrow i$ ;  $dist[v] \leftarrow im[v]$  (* initialize distance with values of minima *)
13:   end for
14: end for
15: while  $V \neq \emptyset$  do
16:    $u \leftarrow GetMinDist(V)$  (* find  $u \in V$  with smallest distance value  $dist[u]$  *)
17:    $V = V - \{u\}$ 
18:   for all  $v \in V$  with  $(u, v) \in E$  do
19:     If  $dist[u] + cost[u, v] < dist[v]$  then
20:        $dist[v] \leftarrow dist[u] + cost[u, v]$ 
21:        $lab[v] \leftarrow lab[u]$ 
22:     else if  $lab[v] \neq WSHED$  and  $dist[u] + cost[u, v] = dist[v]$  and  $lab[v] \neq lab[u]$  then
23:        $lab[v] = WSHED$ 
24:     end if
25:   end for
26: end while

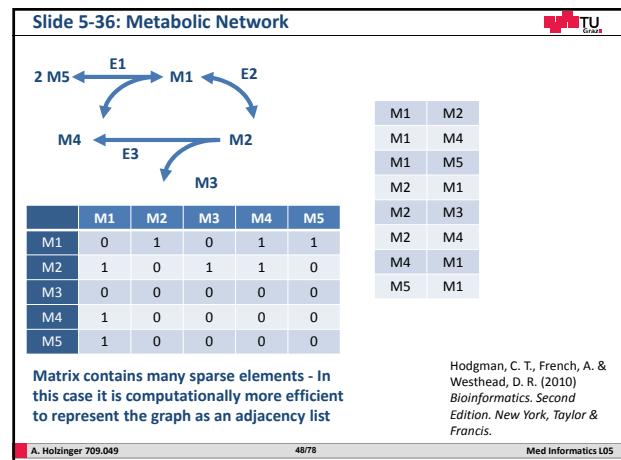
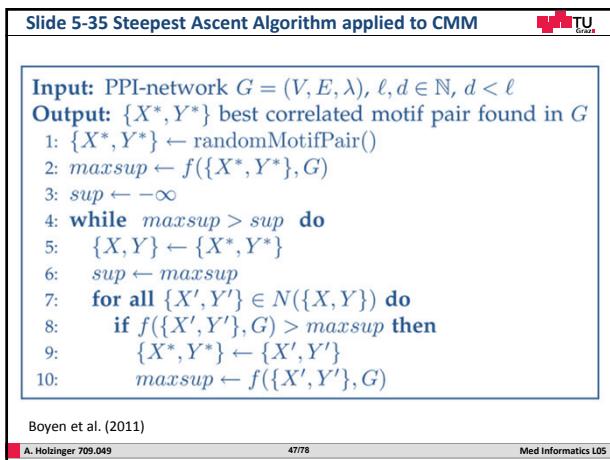
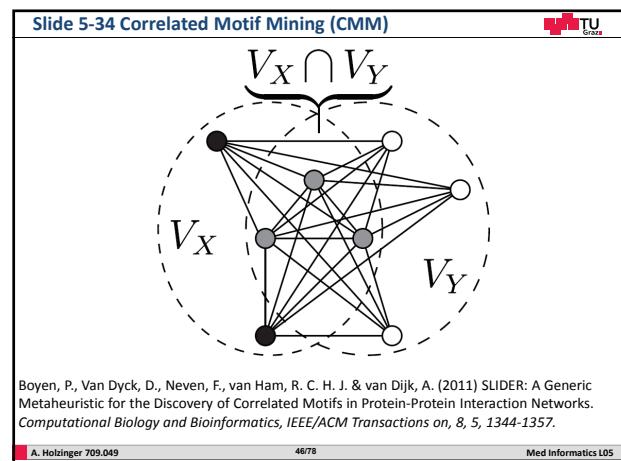
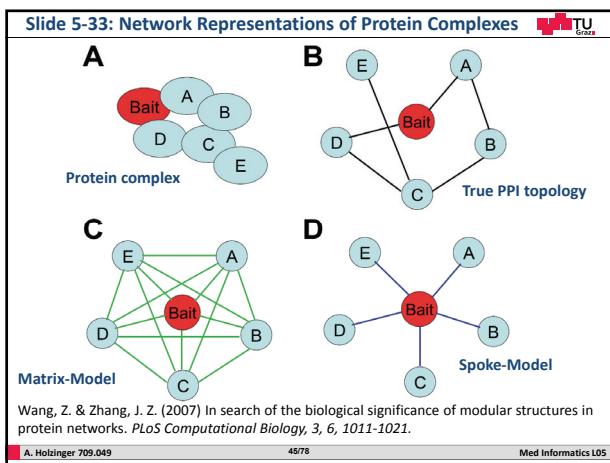
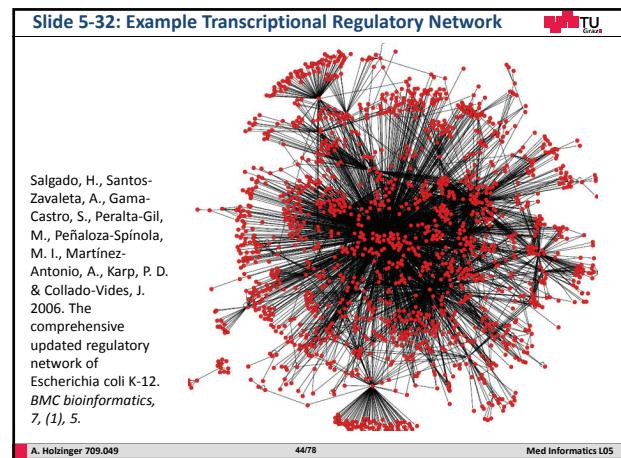
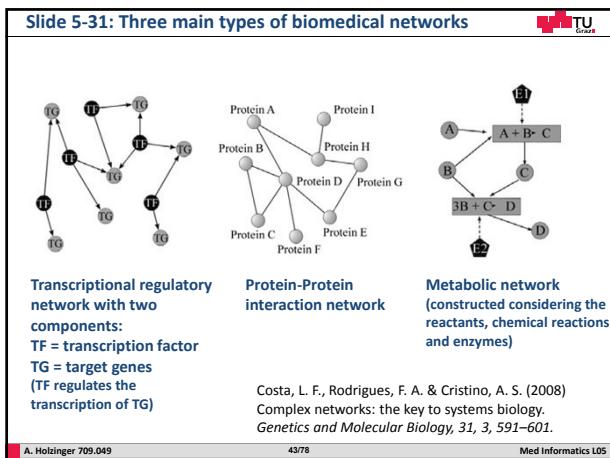
```

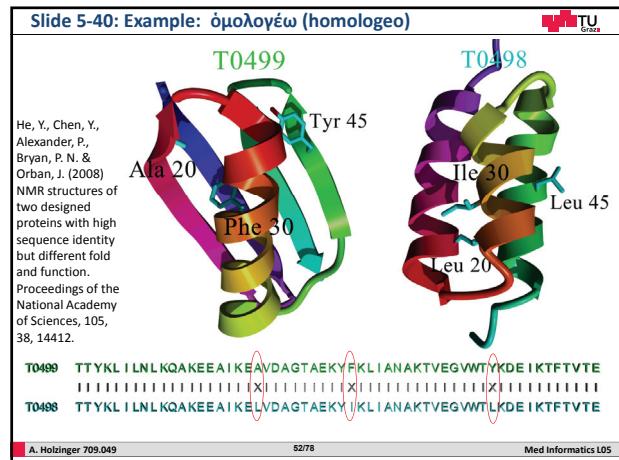
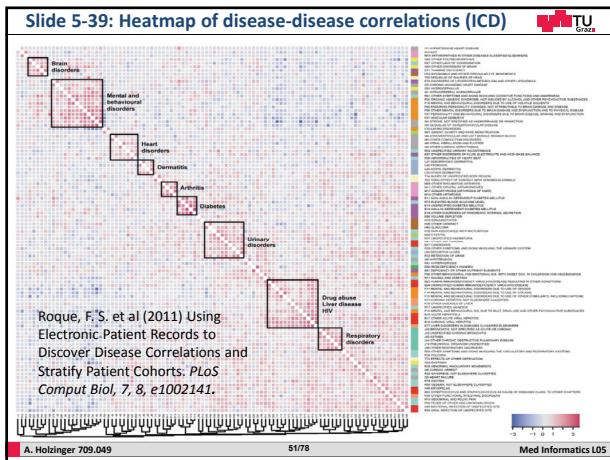
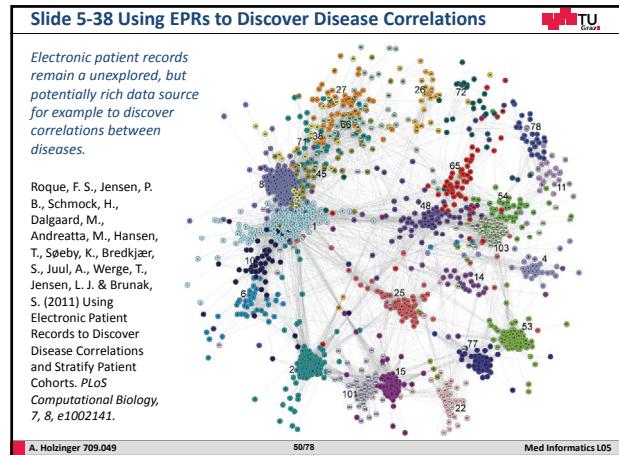
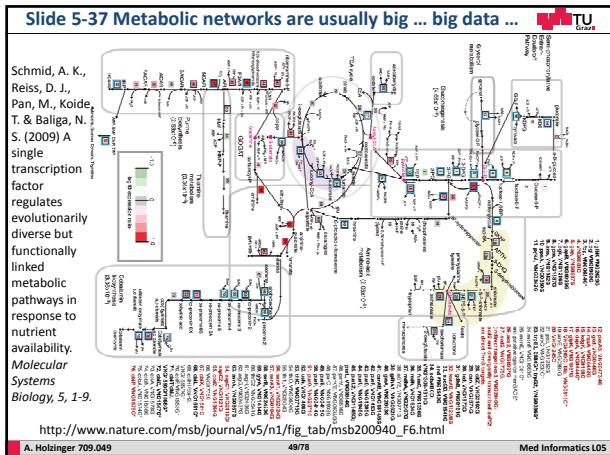
Meijster, A. & Roerdink, J. B. A proposal for the implementation of a parallel watershed algorithm. Computer Analysis of Images and Patterns, 1995. Springer, 790-795.

A. Holzinger 709.049
30/78
Med Informatics L05

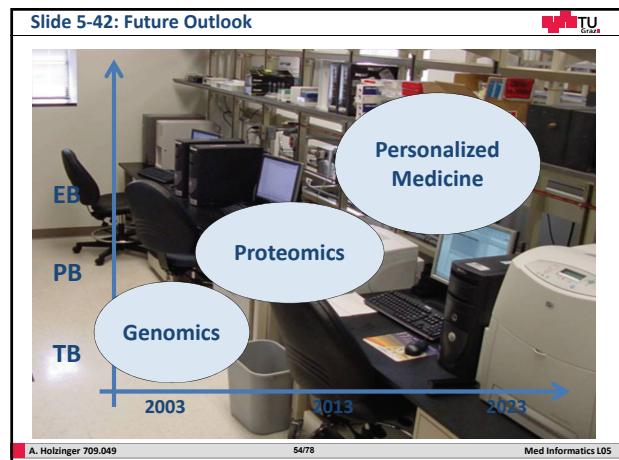








- Slide 5-41 Conclusion**
- Homology modeling is a knowledge-based prediction of protein structures.
 - In homology modeling a protein sequence with an unknown structure (the target) is aligned with one or more protein sequences with known structures (the templates).
 - The method is based on the principle that homologue proteins have similar structures.
 - **Homology modeling will be extremely important to personalized and molecular medicine in the future.**
- A. Holzinger 709.049 53/78 Med Informatics L05





Sample Questions

- Which are the four main “big data” pools in the health care domain and what problems involved?
- What is the main problem in medical documentation?
- What is the advantage of an integrated Patient record?
- What are the advantages/disadvantages of XML/OWL for data in bioinformatics?
- What are the three key concepts in order to understand complex biological systems?
- What are the main symbols describing a network as used in Bioinformatics?
- How can networks represented computationally effectively?
- What are the main network metrics?
- What are the main network topologies used in Biomedical informatics?
- What is the Small-World Theory?
- Why is the study of networks relevant for medical professionals?
- Which are the three main types of biomedical networks?
- What is a Motif?
- What benefits can we gain from Correlated Motif Mining (CMM)?
- What is more efficient if a matrix contains many sparse elements?
- Why are structural homologies interesting for biomedical informatics?

A. Holzinger 709.049

56/78

Med Informatics L05

Some Useful Links

TU Graz

- <http://www.cdisc.org>
- <http://www.w3.org/Math/>
- <http://www.sgpp.org/structures.shtml>
- <http://salilab.org/modeller>
- <http://swissmodel.expasy.org>
- <http://www.expasy.org/tools>
- <http://www.geneticseducation.nhs.uk>

A. Holzinger 709.049 57/78 Med Informatics L05

Appendix: clustering network motifs in integrated networks

TU Graz

A. Holzinger 709.049 58/78 Med Informatics L05

<http://omics.frias.uni-freiburg.de/>

Example from Immunology: Structural Homology

TU Graz

Nature Reviews | Immunology

Calandra, T. & Roger, T. 2003. Macrophage migration inhibitory factor: a regulator of innate immunity. *Nat Rev Immunol*, 3, 791–800.

A. Holzinger 709.049 59/78 Med Informatics L05

Klein Bottle

TU Graz

<http://www.maa.org/cvm/1998/01/tprppoh/article/Pictures/KleinBottle.gif>

A. Holzinger 709.049 60/78 Med Informatics L05

Medical Documentation – Patient Record

A. Holzinger 709.049 61/78 Med Informatics L05

Medical Documentation - Electronic Patient Record

A. Holzinger 709.049 62/78 Med Informatics L05

Challenge is in Genomic medicine ...

- ... to integrate and analyze these diverse and voluminous data sources to elucidate both normal and disease physiology.
- XML is suited for describing semi-structured data including a natural modeling of biological entities, because it allows features as e.g. nesting ...

A. Holzinger 709.049 63/78 Med Informatics L05

Example: Comparison of XML and OWL data in bioinformatics

difficulty of modeling many-to-many relationships, such as the relationship between genes and functions

```
<xml version="1.0">
<!DOCTYPE Generalist SYSTEM "Generalist.dtd">
<Generalist>
<!-- symbol <CREB-17A> organism="D. melanogaster">
<Sequence>ATCGACGCGCCGCGCTGCT</Sequence>
<Sequence>ATCGACGCGCCGCGCTGCT</Sequence>
<Function id="0007616" status="confirmed">Term-long-term memory</Function>
<Function id="0007616" status="inferred">Term-long-term memory</Function>
</Generalist>
```

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 64/78 Med Informatics L05

On time and space of data ...

A. Holzinger 709.049 65/78 Med Informatics L05

... to microscopic atomistic structures

ATOM	1	N	GLY	A	1	44.642	51.634	101.504	0.01	27.20
ATOM	2	CA	GLY	A	1	45.240	51.230	100.593	0.01	26.93
ATOM	3	C	GLY	A	1	46.692	49.449	101.308	0.01	26.80
ATOM	4	O	GLY	A	1	46.895	50.222	102.391	0.01	26.91
ATOM	5	N	SER	A	2	47.283	48.516	100.951	1.00	26.26
ATOM	6	CA	SER	A	2	48.277	47.866	101.761	1.00	26.17
ATOM	7	C	SER	A	2	49.212	47.031	100.845	1.00	24.21
ATOM	8	O	SER	A	2	49.068	47.195	99.630	1.00	19.77
ATOM	9	CB	SER	A	2	47.430	47.091	102.800	1.00	26.31
ATOM	10	OG	SER	A	2	46.270	46.330	100.904	1.00	27.95
ATOM	11	H	ND1	HIS	3	50.047	46.886	101.378	1.00	26.53
ATOM	12	CA	HIS	A	3	51.129	45.389	100.609	1.00	21.44
ATOM	13	C	HIS	A	3	50.953	43.905	100.849	1.00	20.32
ATOM	14	O	HIS	A	3	50.530	43.595	101.950	1.00	22.00
ATOM	15	CB	HIS	A	3	52.555	45.674	100.990	1.00	19.69
ATOM	16	CG	HIS	A	3	52.940	47.090	100.611	1.00	21.44
ATOM	17	ND1	HIS	A	3	53.371	47.470	99.422	1.00	20.87
ATOM	18	CD2	HIS	A	3	52.956	48.175	101.433	1.00	21.69
ATOM	19	CE1	HIS	A	3	53.676	48.730	99.476	1.00	20.57

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 66/78 Med Informatics L05

