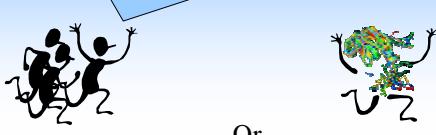


**GET THAT PROTEIN!**



Or

Implementation of and access to biological databanks

Patrick Lambrix  
Department of Computer and Information Science  
Linköpings universitet

1

## Outline

- Biological databanks
- Storing and accessing textual information
- Accessing multiple biological databanks  
(link-driven federations,  
view integration)

2

## Beyond the scope

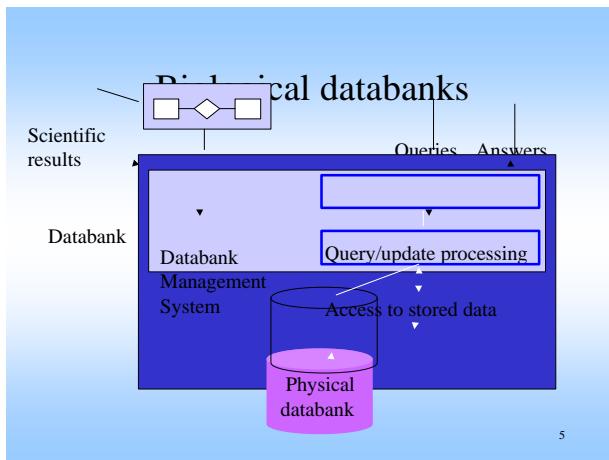
- Non-textual data
- Alignments (e.g. BLAST)
- Visualization
- Data mining

3

## Biological databanks

- Biological data in electronic form
- well-known sources:  
e.g. SWISS-PROT, EMBL, DDBJ, PDB,  
GENBANK, KEGG, ACEDB
- Used in every day research

4



## Outline

- Biological databanks
- **Storing and accessing textual information**
- Accessing multiple biological databanks
  - link-driven federations
  - view integration

6

## Storing and accessing textual information

- What information is stored?
- How is the information stored?
  - high level
- How is the information retrieved?

7

## What information is stored?

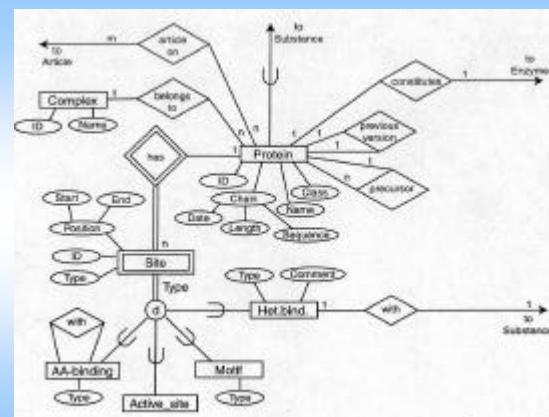
- Model the information
  - Entity-Relationship model (ER)
  - Unified Modeling Language (UML)

8

## What information is stored? - ER

- entities and attributes
- entity types
- key attributes
- relationships
- cardinality constraints
- EER: sub-types

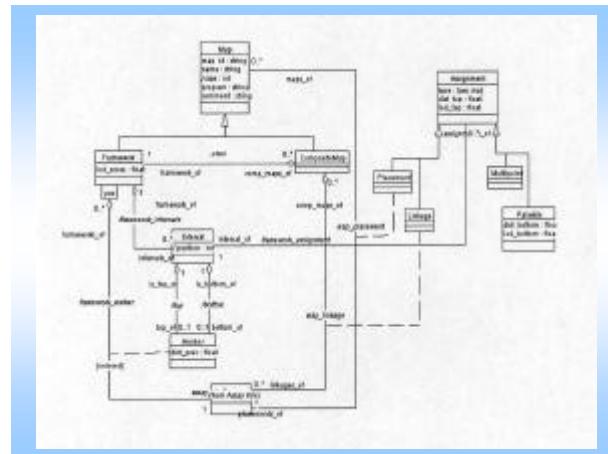
9



## What information is stored? - UML

- classes, objects
- attributes
- operations
- associations, links

11



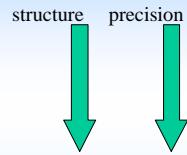
## Storing and accessing textual information

- What information is stored?
- How is the information stored?
  - high level
- How is the information retrieved?

13

## Storing textual information

- Text (IR)
- Semi-structured data
- Data models (DB)
- Rules + Facts (KB)



14

```

1 tgcacccg gccccggc tgcgtttt ccccaaccac gccccggc tgcaccc
61 ccggcccg tgcgtttt tgcgtttt gccccggc tgcgtttt gccccggc
121 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
181 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
241 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
301 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
361 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
421 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
481 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
541 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
601 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
661 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
721 accaaatg tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
781 ctttttgtt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
841 agaaagttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
901 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
961 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
1021 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
1081 ttcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
1141 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
1201 tcaaccat tttttttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
1261 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
1321 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
1381 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
1441 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
1501 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
1561 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
1621 tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt
1681 tttttttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt tgcgtttt

```

15

DEFINITION	Homo sapiens adrenergic, beta-1-, receptor
ACCESSION	NM_000684
SOURCE ORGANISM	human
REFERENCE	1
AUTHORS	Frielle, Collins, Daniel, Caron, Lefkowitz, Kobilka
TITLE	Cloning of the cDNA for the human beta 1-adrenergic receptor
REFERENCE	2
AUTHORS	Frielle, Kobilka, Lefkowitz, Caron
TITLE	Human beta 1- and beta 2-adrenergic receptors: structurally and functionally related receptors derived from distinct genes

Homo sapiens adrenergic, beta-1-, receptor  
NM\_000684  
human  
1  
Frielle, Collins, Daniel, Caron, Lefkowitz, Kobilka  
Cloning of the cDNA for the human beta 1-adrenergic receptor  
2  
Frielle, Kobilka, Lefkowitz, Caron  
Human beta 1- and beta 2-adrenergic receptors: structurally and functionally related receptors derived from distinct genes

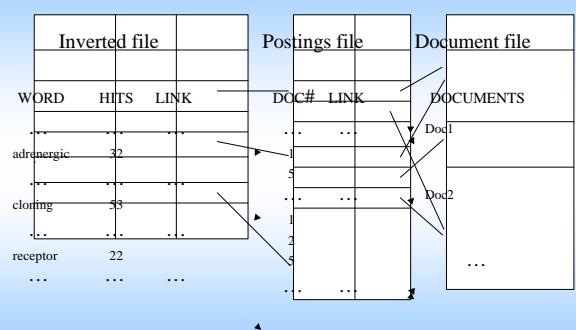
16

## Storing textual information - Text - Information Retrieval

- search using words
- conceptual models:
  - boolean, vector, probabilistic, ...
- file model:
  - flat file, inverted file, ...

17

## IR - File model: inverted files



18

## IR - formal characterization

Information retrieval model: (D,Q,F,R)

- D is a set of document representations
- Q is a set of queries
- F is a framework for modeling document representations, queries and their relationships
- R associates a real number to document-query-pairs (ranking)

19

## IR - conceptual models

Classic information retrieval

- Boolean model
- Vector model
- Probabilistic model

20

## Boolean model

	adrenergic	cloning	receptor	
Doc1	yes	yes	no	--> (1 1 0)
Doc2	no	yes	no	--> (0 1 0)

Q1: cloning and (adrenergic or receptor)  
--> (1 1 0) or (1 1 1) or (0 1 1)      Result: Doc1

Q2: cloning and not adrenergic  
--> (0 1 0) or (0 1 1)      Result: Doc2

21

## Boolean model

### Advantages

- based on intuitive and simple formal model (set theory and boolean algebra)

### Disadvantages

- binary decisions
  - words are relevant or not
  - document is relevant or not, no notion of partial match

22

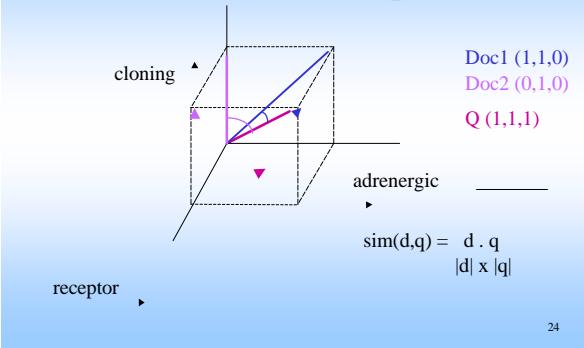
## Boolean model

	adrenergic	cloning	receptor	
Doc1	yes	yes	no	--> (1 1 0)
Doc2	no	yes	no	--> (0 1 0)

Q3: adrenergic and receptor  
--> (1 0 1) or (1 1 1)      Result: empty

23

## Vector model (simplified)



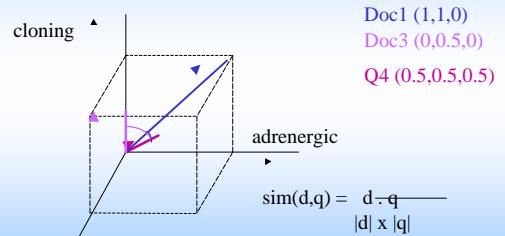
24

## Vector model

- Introduce weights in document vectors (e.g. Doc3 (0, 0.5, 0))
- Weights represent importance of the term for describing the document contents
- Weights are positive real numbers
- Term does not occur -> weight = 0

25

## Vector model



26

## Vector model

- How to define weights? tf-idf

$$d_j (w_{1,j}, \dots, w_{t,j})$$

$w_{i,j}$  = weight for term  $k_i$  in document  $d_j$   
 $= f_{i,j} \times idf_i$

27

## Vector model

- How to define weights? tf-idf

term frequency  $\text{freq}_{i,j}$ : how often does term  $k_i$  occur in document  $d_j$ ?

normalized term frequency:

$$f_{i,j} = \text{freq}_{i,j} / \max_i \text{freq}_{i,j}$$

28

## Vector model

- How to define weights? tf-idf
- document frequency : in how many documents does term  $k_i$  occur?

$N$  = total number of documents

$n_i$  = number of documents in which  $k_i$  occurs  
inverse document frequency  $idf_i$ :  $\log(N / n_i)$

29

## Vector model

- How to define weights for query?

recommendation:

$$q = (w_{1,q}, \dots, w_{t,q})$$

$$w_{i,q} = \text{weight for term } k_i \text{ in } q \\ = (0.5 + 0.5 f_{i,q}) \times idf_i$$

30

## Vector model

- Advantages
  - term weighting improves retrieval performance
  - partial matching
  - ranking according to similarity
- Disadvantage
- assumption of mutually independent terms?

31

## Probabilistic model

weights are binary ( $w_{i,j} = 0$  or  $w_{i,j} = 1$ )  
 R: the set of relevant documents for query q  
 Rc: the set of non-relevant documents for q  
 $P(R|d_j)$ : probability that  $d_j$  is relevant to q  
 $P(Rc|d_j)$ : probability that  $d_j$  is not relevant to q

$$\text{sim}(d_j, q) = P(R|d_j) / P(Rc|d_j)$$

32

## Probabilistic model

$$\text{sim}(d_j, q) = P(R|d_j) / P(Rc|d_j)$$

(Bayes' rule, independence of index terms,  
take logarithms,  $P(k_i|R) + P(\text{not } k_i|R) = 1$ )

$\Rightarrow \text{SIM}(d_j, q) =$

$$\sum_{i=1}^t w_{i,q} \times w_{i,j} \times$$

$$(\log(P(k_i|R) / (1 - P(k_i|R))) +$$

$$\log(P(k_i|Rc) / (1 - P(k_i|Rc))))$$

33

## Probabilistic model

- How to compute  $P(k_i|R)$  and  $P(k_i|Rc)$ ?
    - initially:  $P(k_i|R) = 0.5$  and  $P(k_i|Rc) = n_i/N$
    - Repeat: retrieve documents and rank them
- V: subset of documents (e.g. r best ranked)  
 Vi: subset of V, elements contain  $k_i$
- $$P(k_i|R) = |V_i| / |V|$$
- and  $P(k_i|Rc) = (n_i - |V_i|) / (N - |V|)$

34

## Probabilistic model

- Advantages:
  - ranking of documents with respect to probability of being relevant
- Disadvantages:
  - initial guess about relevance
  - all weights are binary
  - independence assumption?

35

## IR - measures

---

Precision = 
$$\frac{\text{number of found relevant documents}}{\text{total number of found documents}}$$

---

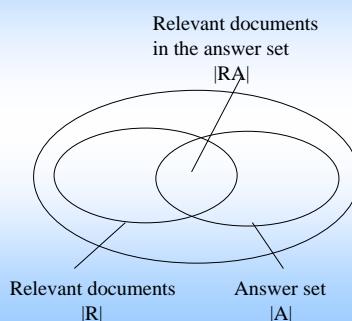
Recall = 
$$\frac{\text{number of found relevant documents}}{\text{total number of relevant documents}}$$

36

## IR - measures

$$\text{Precision} = \frac{|RA|}{|A|}$$

$$\text{Recall} = \frac{|RA|}{|R|}$$



37

## PERL

### Practical Extraction and Reporting Language

```
while ($line=<INPUTFILE>) {
    if ($line=~ /^DEFINITION/) {
        ($head, $rest) = split("DEFINITION", $line);
        print OUTPUTFILE $rest;
        print OUTPUTFILE "\n";
    }
}
```

38

## Storing textual information - Databases

- Relational databases:
  - model: tables + relational algebra
  - query language (SQL)
- Object-oriented databases
  - model: persistent objects, messages, encapsulation, inheritance
  - query language (e.g. OQL)
- Systems: GDB (R), ACEDB (OO)

39

### Relational databases

PROTEIN					DESCRIBED-IN	
PROTEIN-ID	ACCESSION	DEFINITION	SOURCE		PROTEIN-ID	REFERENCE-ID
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human		1	1
ARTICLE-AUTHORS			ARTICLE-TITLE			
REFERENCE-ID	AUTHOR	REFERENCE-ID			TITLE	
1	Fribille		1		Cloning of the cDNA for the human beta 1-adrenergic receptor	
1	Collins				Human beta 1- and beta 2- adrenergic receptors: structurally and functionally related receptors derived from distinct genes	
1	Daniel					
1	Caron			2		
1	Lefkowitz					
1	Kobilka					
2	Fribille					
2	Kobilka					
2	Lefkowitz					
2	Caron					

40

## SQL queries

```
select source
from protein
where accession = NM_000684;
```

PROTEIN	ACCESSION	DEFINITION	SOURCE
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human

41

## SQL queries

```
select title
from protein, article-title, described-in
where protein.accession = NM_000684
and protein.protein-id
= described-in.protein-id
and described-in.reference-id
= article-title.reference-id;
```

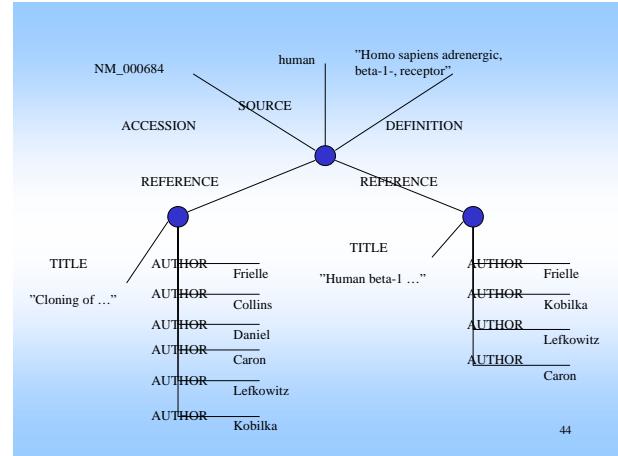
PROTEIN				ARTICLE-TITLE	
PROTEIN-ID	ACCESSION	DEFINITION	SOURCE	REFERENCE-ID	TITLE
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human	1	Cloning of the ...

42

## Storing textual information - Semi-structured data

- Connected network of nodes
- query: path search
- systems: Genbank, OMIM

43

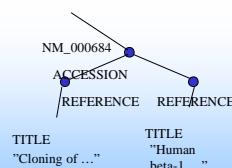


44

## Queries

```
select source
from protein
where fills accession NM_000684;
```

NM\_000684 → human  
ACCESSION SOURCE



45

## Storing textual information - Knowledge bases

- Often based on a logic
- Query answering based on inference mechanism
- Knowledge bases often fit in main memory
- Useful for ontologies

46

## Storing textual information - Knowledge bases

(F) source(NM\_000684, Human)  
(R) source(P?,Human) => source(P?,Mammal)  
(R) source(P?,Mammal) => source(P?,Vertebrate)

**Q:** ?- source(NM\_000684, Vertebrate)

**A:** yes

**Q:** ?- source(x?, Mammal)  
**A:** x? = NM\_000684

47

## Outline

- Biological databanks
- Storing and accessing textual information
- Accessing multiple biological databanks
  - link-driven federations
  - view integration

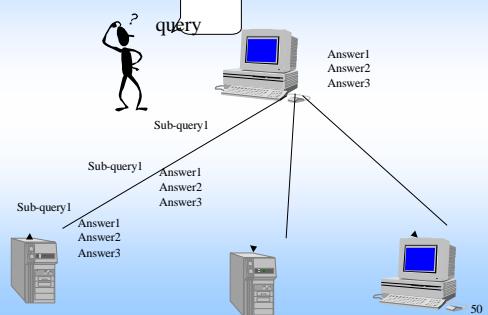
48

## Access to multiple databanks - Problems

- Users of biological databanks need to have good knowledge about where to find information and how to retrieve the information from the sources.
- Representations in different databanks may not conform or not be consistent with each other.

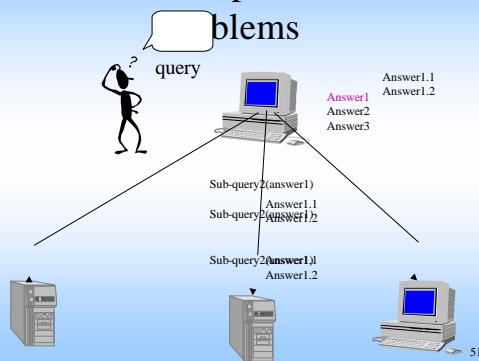
49

## Access to multiple databanks - Problems



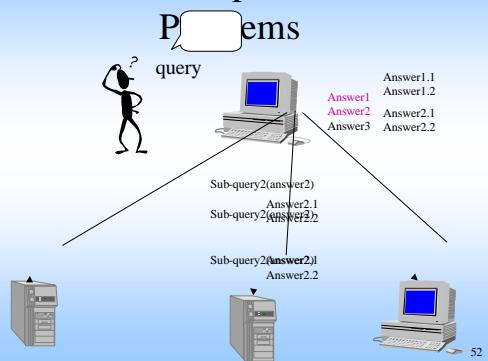
50

## Access to multiple databanks - Problems



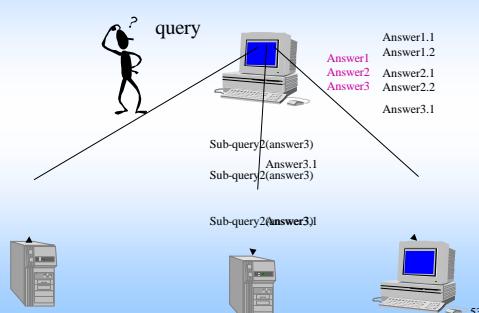
51

## Access to multiple databanks - Problems



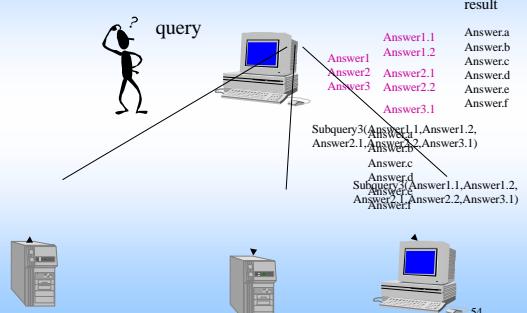
52

## Access to multiple databanks - Problems



53

## Access to multiple databanks - Problems



54

## Current approaches

- Link-driven federations
- View integration

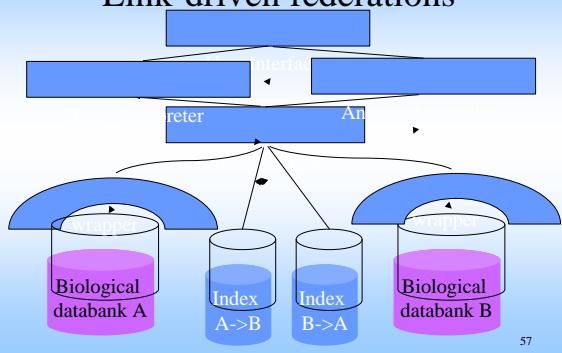
55

## Current approaches

- Link-driven federations
  - create explicit links between databases
  - query: retrieve interesting results and use web links to go to other related data sources
  - SRS, Entrez

56

## Link-driven federations



57

## Link-driven federations

- Advantages
  - complex queries
  - fast
- Disadvantages
  - good knowledge required
  - syntax-based
  - discrepancies in terminology not solved

58

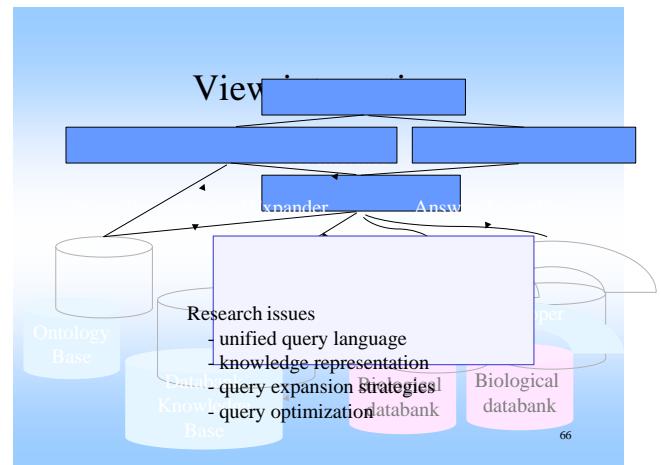
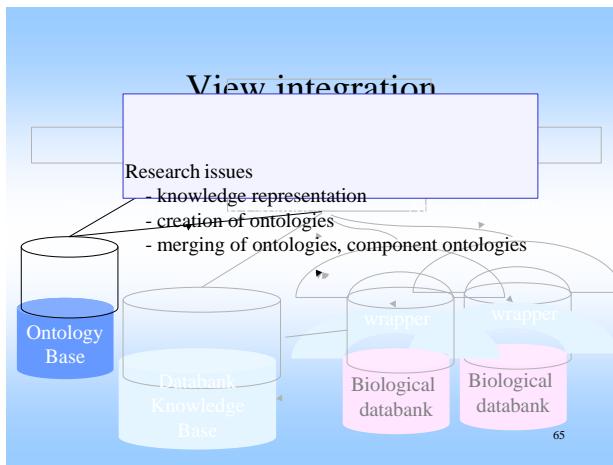
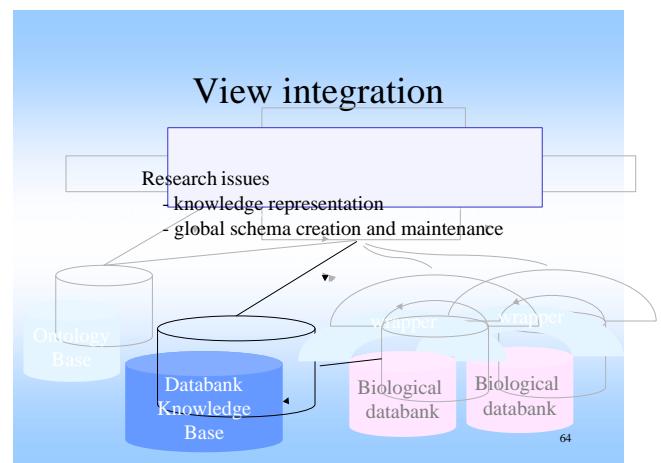
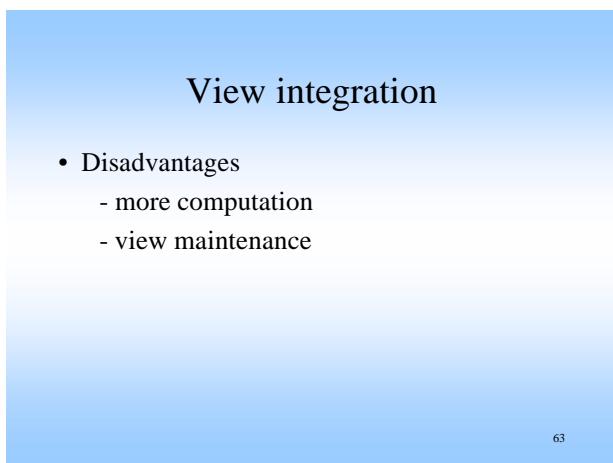
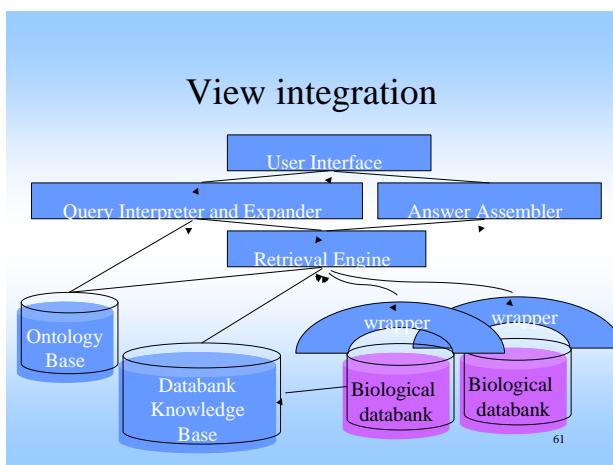


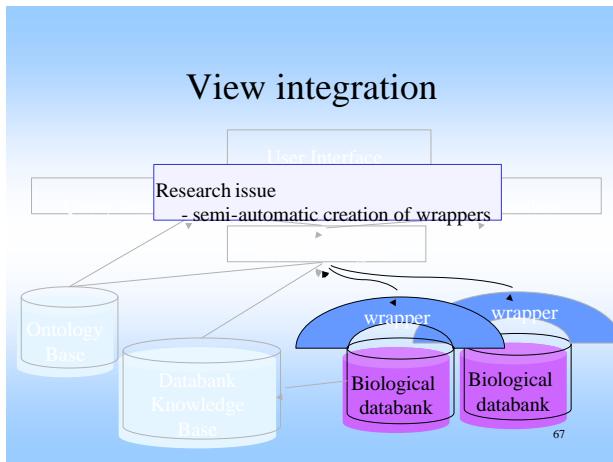
59

## Current approaches

- View integration
  - define global schema over the databases
  - query: using a high-level query language
  - BioKleisli, K2, TAMBIS, BioTRIFU

60

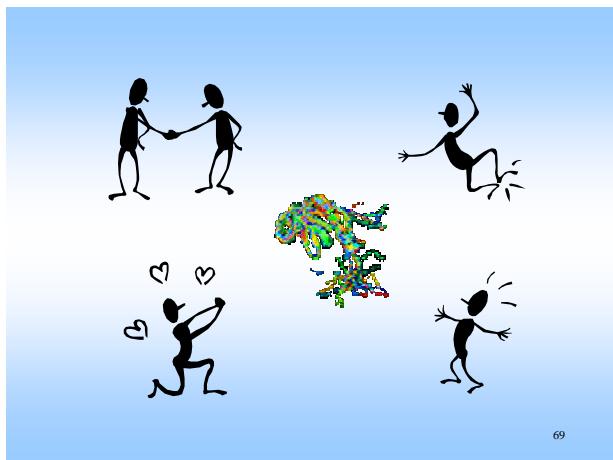




## Related work at IISLAB Laboratory for Intelligent Information Systems

- Transparent access to multiple biological databanks
- Live help systems (web, ML)
- Extraction of information from text (IE,ML)
- TRIFU (KBIR)
- Description logics (KR)

68



## Literature

### • Information Retrieval

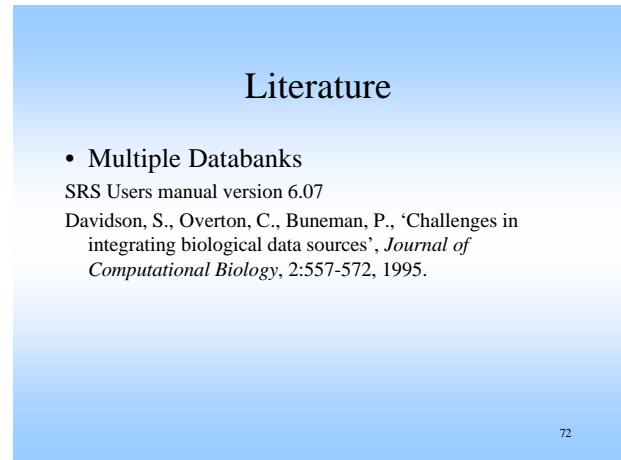
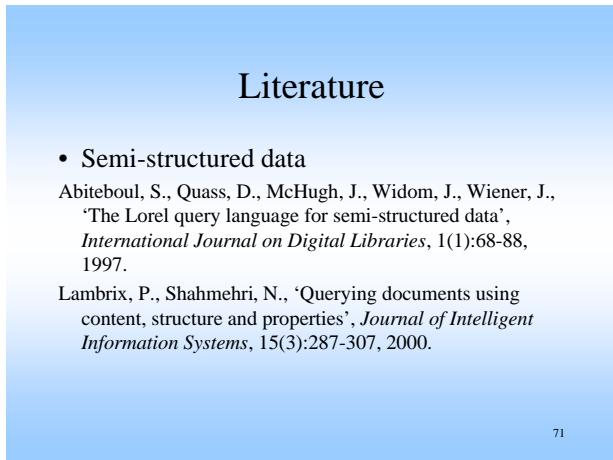
Baeza-Yates, R., Ribeiro-Neto, B., *Modern Information Retrieval*, Addison-Wesley, 1999.

Schwartz, R., Wall, L., *Learning Perl*, O'Reilly & Associates, 1994.

### • Databases

Elmasri, R., Navathe, S., *Fundamentals of database systems*, 3rd edition, Addison Wesley, 2000.

70



## Literature

- Multiple Databanks

- Davidson, S., Crabtree, J., Brunk, B., Schug, J., Tannen, V., Overton, C., Stoeckert, C., 'K2/Kleisli and GUS: Experiments in integrated access to genomic data sources', *IBM Systems journal*, 40(2):512-531, 2001.
- Goble, C., Stevens, R., Ng, G., Bechhofer, S., Paton, N., Baker, P., Peim, M., Brass, A., 'Transparent access to multiple bioinformatics information sources', *IBM Systems Journal*, 40(2):532-551, 2001.

73

## Assignment

Choose one of:

- Challenges in integrating biological databanks
- Compare integration systems
- own topic - literature
- own topic - implementation

74