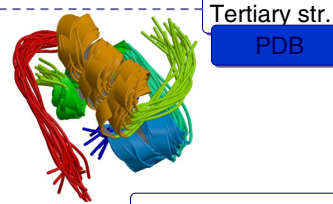


## Biological data

LOCUS HUMINS01 4044 bp DNA  
 DEFINITION Human insulin gene, complete cds.  
 ACCESSION J00265  
 VERSION J00265.1 GI:186429  
 \*\*\*  
 3841 aggggcccag gatgtgggg ccaactgagaa gtgacttttt gtccaatgac tctggactct  
 3901 tggagtcgcc agagaccttg ttcaggaag ggaatgagaa cattccagca attttccccc  
 3961 caccctagccc tcccaggttc tatttttaga gttattctg atggaatccc tgtggaggga  
 4021 ggagactggg ctgagggagg ggggt

DNA seq.

GenBank



### Entry information

Entry name **INS\_HUMAN**  
 Primary accession number **P01308**

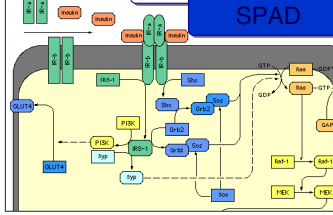
Protein seq.

SWISS-PROT

### Comments

- FUNCTION:** Insulin decreases blood glucose concentration. It increases cell permeability to monosaccharides, amino acids and fatty acids. It accelerates glycolysis, the pentose phosphate cycle, and glycogen synthesis in liver.
  - SUBUNIT:** Heterodimer of a B chain and an A chain linked by two disulfide bonds.
  - SUBCELLULAR LOCATION:** Secreted.
  - DISEASE:** Defects in INS are the cause of familial hyperproinsulinemia [MIM:176720]
- 10 20 30 40 50 60  
 MALWRRLLPL LALLALWQPD PAAAFVNOHL CSHLVEALY LVCGERGFY TPKTRREAE  
 70 80 90 100 110  
 LVGQVELGG QPAGSLQPL ALEGLQKRG IVEQCCTSTG SLVQLNTGN

### Signaling pathway



### General information about the entry

Entry name **INSULIN**  
 Accession number **PS00262**  
 Entry type **PATTERN**

Secondary str.

PROSITE

### Name and characterization of the entry

Description Insulin family signature.  
 Pattern C-C-[P]-[F]-x-C-[STDNEKPI]-x(3)-[LIVMFS]-x(3)-C.

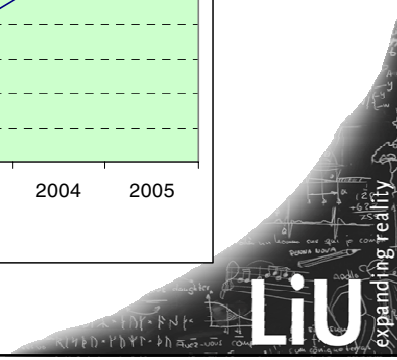
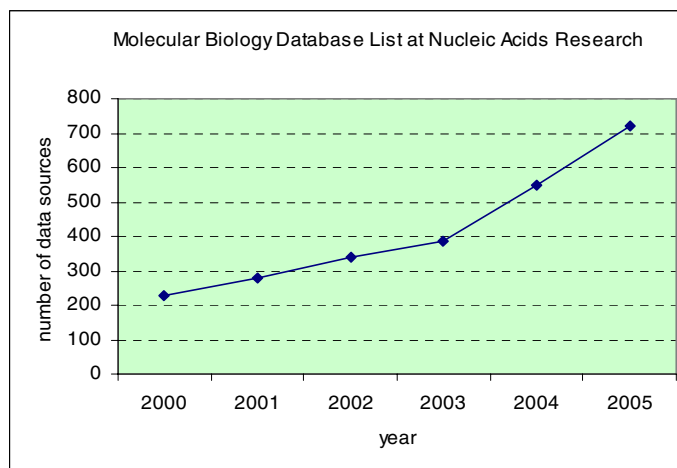
INSULIN

### Taxonomy

AmiGO

- all: all (184297)
- GO:0003674 : molecular function
  - GO:0005488 : binding (33177)
  - GO:0005102 : receptor binding (1884)
  - GO:0005179 : hormone activity (325)
  - GO:0004871 : signal transducer activity (3374)
  - GO:0005102 : receptor binding (1884)
  - GO:0005179 : hormone activity (325)

## Biological data sources

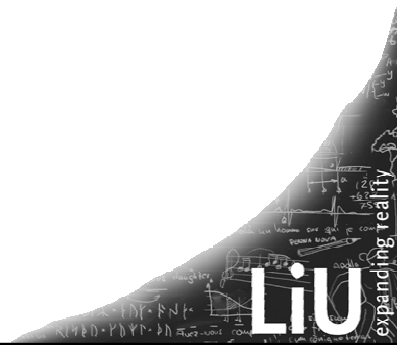


## Outline

- What is the problem?
- Standardization for the Web
- Impact on database technology

mars 2009

7

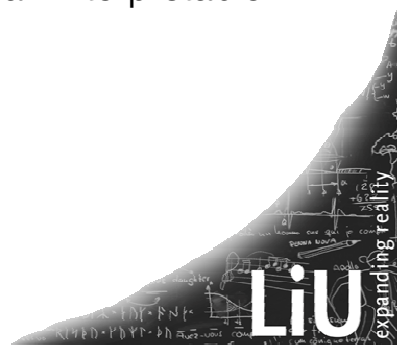


## What is the problem?

- The user's effort is not enough for the task
- The data describes complex real world objects
- The data is not easily human interpretable

mars 2009

8

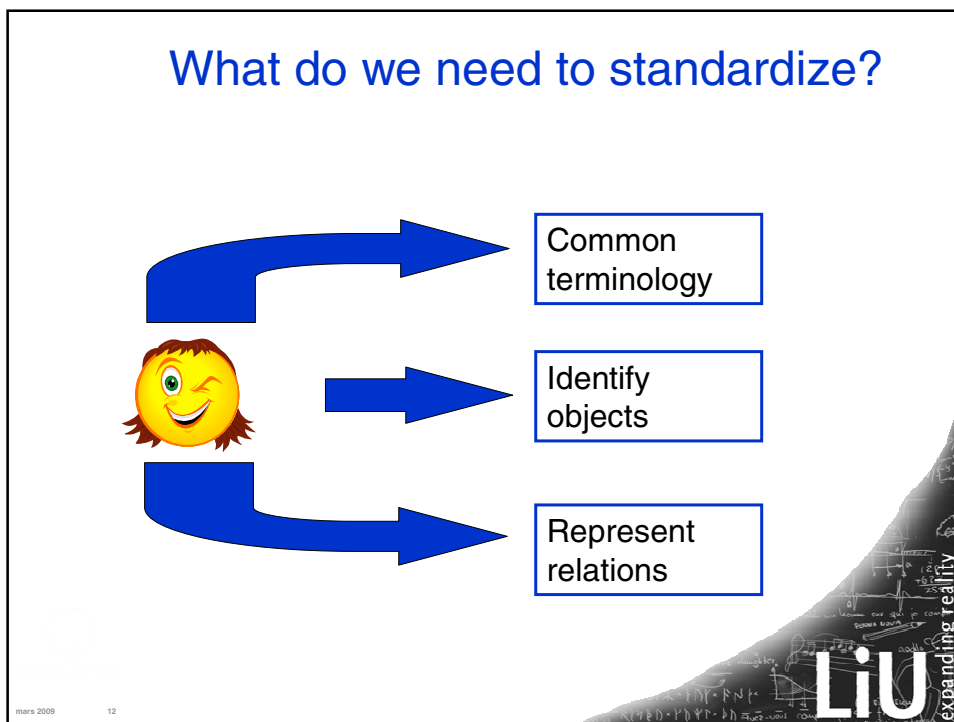
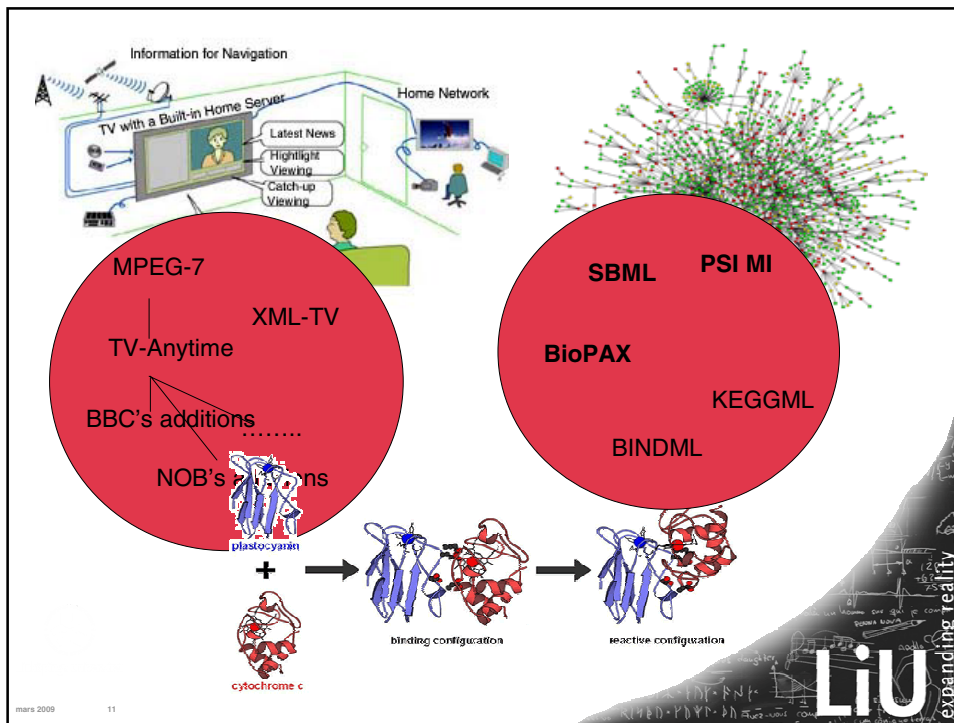


## A contradiction?

- Data integration is a hard problem
- Standardization and agreement of common format is a prerequisite for efficient data integration
- The Web is an ad-hoc platform where new data formats and actors occurs all the time

## Standardization for the Web

- How are standards used?
- What do they contain?



## How to represent data ...

- Web data is often less structured than traditional applications
  - Semi structured
  - Integration
- XML is the most common language
- OWL and RDF are alternatives

mars 2009

13

LiU

expanding reality

## The TV domain

### TV-Anytime

```
<TVAMain>
  <ProgramDescription>
  <ProgramLocationTable>
  <Schedule serviceIDRef='SVT1'
    start='...' end='...'>
    <ScheduleEvent>
      <Program crid='crid:...'>
    </ScheduleEvent>
    </Schedule>
  </ProgramLocationTable>
  <ProgramInformationTable>
  <ProgramInformation
    programId='crid:...'>
  <BasicDescription>
    <Title>SVT News</Title>
    <Genre href='...'>
    <Name>News</Name>
    </Genre> ...
  </BasicDescription>
  </ProgramInformation>
  </ProgramInformationTable>
</ProgramDescription>
</TVAMain>
```

### XMLTV

```
<tv>
  <channel id="C1">
    <display-name lang="se">
      SVT1</display-name>
  </channel>
  <programme
    start="200006031633"
    channel="C1">
    <title lang="sv">Nyheter</title>
    <title lang="en">News</title>
    <desc lang="sv">...</desc>
    <category>News</category>
    <country>SE</country>
  </programme>
</tv>
```

mars 2009

14

LiU

expanding reality

## Within molecular interactions ....

Name	Ver.	Year	Defined by	Purpose	Tools	Data
SBML	2	2003	Systems Biology Workbench development group.	A computer-readable format for representing models of biochemical reaction networks.	Many tools available.	Data available from many databases, for instance, KEGG and Reactome.
PSIMI	2.5	2005	Proteomics Standards Initiative.	A standard for data representation for protein-protein interaction to facilitate data comparison, exchange and verification.	Tools for viewing and analysis.	Datasets available from many sources, for instance IntAct, DIP and MINT.
BioPAX	2	2005	The BioPAX group.	A collaborative effort to create a data exchange format for biological pathway data.	Existing tools for OWL such as Protégé.	Datasets available from Reactome.
CellML	1.1	2002	University of Auckland and Physiome Sciences, Inc.	Support the definition of models of cellular and subcellular processes.	Tools for publication, visualization, creation and simulation.	CellML Model Repository (~240 models).
CML	2.2	2003	Peter Murray-Rust, Henry S. Rzepa.	Interchange of chemical information over the Internet and other networks.	Molecular browsers, editors.	BioCYC.
EMBLxml	1.0	2005	EBL	More stability and fine-grained modelling of nucleotide sequence information.	API support in BioJavaX.	EMBL.
INSDseq	1.4	2005	International Nucleotide Sequence Database Collaboration.	The purpose of INSDSeq is to provide a near-uniform representation for sequence records.	API support in BioJavaX.	EMBL, DBJ and GenBank.
Seqentry	n/a	n/a	NCBI	NCBI uses ASN.1 for the storage and retrieval of data such as nucleotide and protein sequences. Data encoded in ASN.1 can be transferred to XML.	SR's BioWarehouse and ProteinStructureFactory's ORFer.	Entrez.
BSML	3.1	2002	Labbook.com.	Facilitate the interchange of data for more efficient communication within the life sciences community.	Labbook's Genomic Browser and Sequence Viewer. Converters.	Previously provided by EMBL.
HUP-ML	0.8	2003	JHUPO.	A proteomics-oriented markup language for exchanging proteome data between researchers.	HUP-ML Editor.	
MAGE	1.1	2003	MGED.	To facilitate the exchange	Converters.	

mars 2009

15

LiU

expanding reality

## What do the standards contain?

- Information about objects:
  - Proteins/Complexes
  - Genes/DNA
  - Other molecules
- Interaction information
- Information about experiments
  - Kind of experiment
  - Evidence of the experiment
- More ....

mars 2009

16

LiU

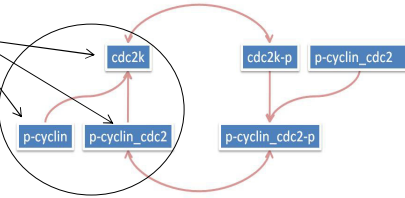
expanding reality



## SBML model example

```

<model id="Tyson1991CellModel_6"
  name="Tyson1991_CellCycle_6var">
  <listOfSpecies>
    + <species id="C2" name="cdc2k" compartment="cell">
    + <species id="M" name="p-cyclin_cdc2" compartment="cell">
    + <species id="YP" name="p-cyclin" compartment="cell"> ...
  </listOfSpecies>
  <listOfReactions>
    <reaction id="Reaction1" name="cyclin_cdc2k dissociation">
      <annotation>
        <rdf:li rdf:resource="http://www.reactome.org/#REACT_6308"/>
        <rdf:li
          rdf:resource="http://www.geneontology.org/#GO:0000079"/>
        </annotation>
      <listOfReactants>
        <speciesReference species="M"/>
      <listOfReactants>
      <listOfProducts>
        <speciesReference species="C2"/>
        <speciesReference species="YP"/>
      </listOfProducts>
      <kineticLaw>
        <math xmlns="http://www.w3.org/1998/Math/MathML">
          <apply> <times> <ci> k6 </ci> <ci> M </ci> </apply> </math>
        </listOfParameters> <parameter id="k6" value="1">
        </listOfParameters>
      </kineticLaw>
    </reaction>
    + <reaction id="Reaction2" name="cdc2k phosphorylation">
      ... more reactions
    </listOfReactions>
  </model>
</sbml>
  
```



mars 2009

17

LiU

expanding reality

## Representation of objects

SBML: Species	PSI MI: Interactor	CellML: component
1 id	id	name
2 name	names	dc:title
3	xref	dterms:alternative
4 speciesType	interactorType	cmeta:bio_entity
5	organism	
6	ncbiTaxId	
7	names	cmeta:species
8	celltype	
9 compartment	compartment	(~ group)
10	tissue	
11	sequence	cm:sequence
12		cm:sequence
13		cm:sequence
14 initialAmount		initialvalue
15 initialConcentration		
16 substanceUnits		units
17 spatialSizeUnits		
18 hasOnlySubstanceUnits		
19 boundaryCondition		
20 charge		
21 constant		

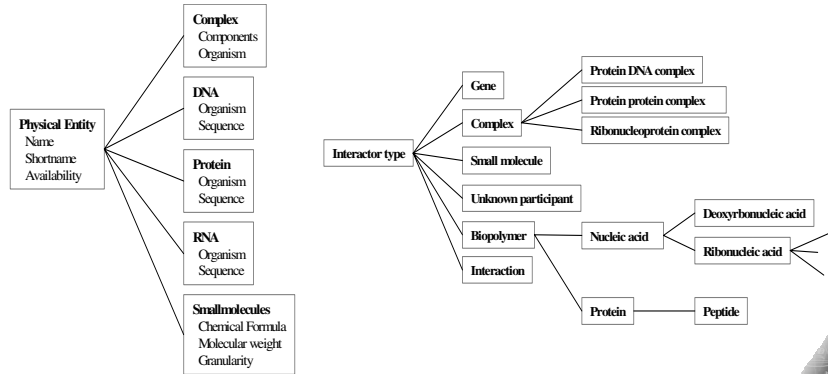
mars 2009

18

LiU

expanding reality

## Type of objects - ontologies



mars 2009

19

LiU

expanding reality

## Representation of interactions

	CellML: component	SBML: Reaction	PSI MI: Interaction	
1			imexID	1
2	name	id	id	2
3		name		3
4	variable		xref	4
5	reaction			5
6		sboTerm	interactiontype	6
7	variable-ref	reactant	experimentList	7
8		product	participantList	8
9		modifier		9
10		id	id	10
11		name	names	11
12	role	sboTerm	experimental-role	12
13			biological-role	13
14			participantidentification	14
15			experimentalpreparation	15
16			confidencelist	16
17	direction			17
18	delta_variable			18
19	stoichiometry	stoichiometry		19
20		kineticLaw		20
21			inferredInteractionList	21
22			participants	22
23			modelled	23
24			confidencelist	24
25	reversible	reversible		25
26		fast		26

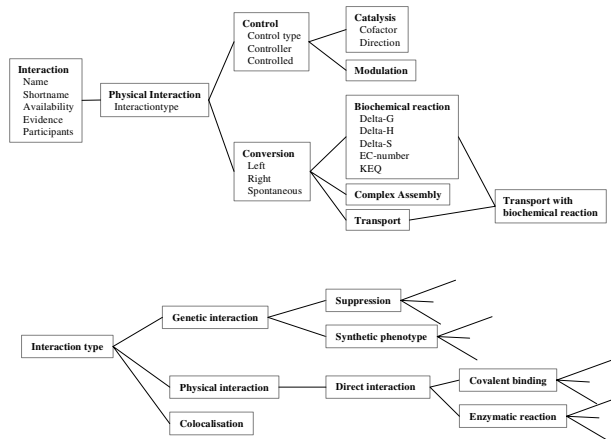
mars 2009

20

LiU

expanding reality

## Interaction types



mars 2009 21

LiU

expanding reality

## Bioinformatics: Minimal Information Standards.

- **MIRIAM : Minimal Information Requested in the Annotation of biochemical Models**
- **MIAPE: The Minimum Information About a Proteomics Experiment**
  - **MIAPE: GE**(Gel Electrophoresis)
  - **MIAPE: MS** (Mass Spectrometry)
  - **MIAPE: CC** (Column Chromatography)
  - **MIAPE: CE** (Capillary Electrophoresis)
  - ....
- **MIMIx: The minimum information required for reporting a molecular interaction experiment ....**
- ....

mars 2009 22

LiU

expanding reality

## XML and storage

- XML provides a data model
- The valid XML data structures can be defined by
  - XML Schema
  - DTD
- XML has its own query languages
  - XPath
  - XQuery
- XML is richer than the relational model
  - Tree structure,
  - Order
  - ...
- Vary from highly structured to unstructured

mars 2009

23

LiU

expanding reality

## Expressing Queries in XQuery

**Find information on a given protein. Protein id is given.**

```
document("rat_small.xml")//proteinInteractor[@id="EBI-77471"]
```

**Find the protein information for the proteins that participate in a given interaction. Interaction id is given.**

```
for $ref in document("rat_small.xml")//interaction  
[names/shortLabel="interaction1"]  
/participantList/proteinParticipant/proteinInteractorRef/@ref  
return document("rat_small.xml")//proteinInteractor[@id=$ref]
```

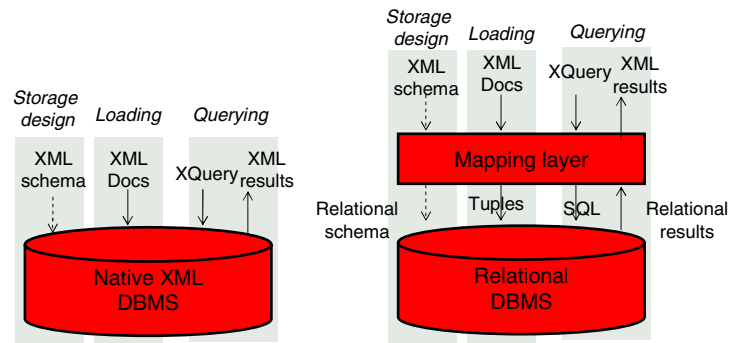
mars 2009

24

LiU

expanding reality

# Storage possibilities for XML



Shredding

LiU

expanding reality

# How to shred XML?

```
<?xml version="1.0" encoding="UTF-8"?>
<families
xmlns:xsi="http://www.w3.org/2001/XMLSchema
a-instance">
  <family>
    <parent>
      <name>Lena</name>
      <job>Lektor</job>
    </parent>
    <child>
      <name>Ludvig</name>
      <school>Skolan</school>
    </child>
  </family>
</families>
```

Source	Ordinal	attrName	isValue	Value
0	1	Families	False	1
1	1	Family	False	2
2	1	Parent	False	3
3	1	Name	True	Lena
3	2	Job	True	Docent ..

## Families

Id	Pid
0	-

## Family

Id	Pid
1	0

## Parent

Id	Pid	Name	Job
2	1	Lena	Lektor

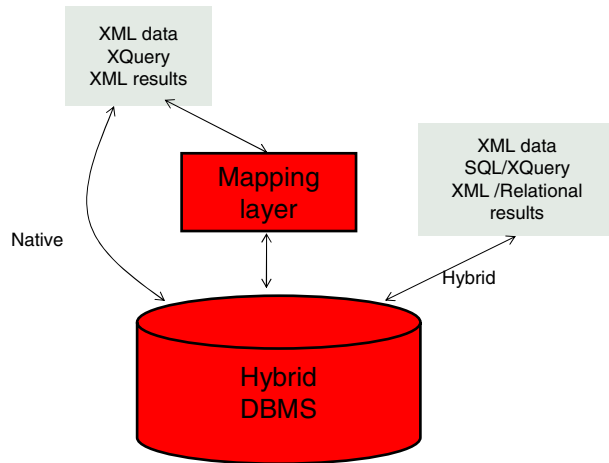
## Child

Id	Pid	Name	School
3	1	Ludvig	Skolan

LiU

expanding reality

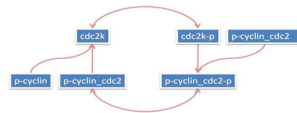
# Hybrid XML Storage



mars 2009 27

LiU expanding reality

# New possibilities...



```

<model id="Tyson1991CellModel_6"
  name="Tyson1991_CellCycle_6var">
<listOfSpecies>
+ <species id="C2" name="cdc2k" compartment="cell">
+ <species id="M" name="p-cyclin_cdc2" compartment="cell">
+ <species id="YP" name="p-cyclin" compartment="cell"> ...
</listOfSpecies>
<listOfReactions>
<reaction id="Reaction1" name="cyclin_cdc2k dissociation">
<annotation>
<rdf:li rdf:resource="http://www.reactome.org/#REACT_6308"/>
</rdf:li>
rdf:resource="http://www.geneontology.org/#GO:000079"/>
</annotation>
<listOfReactants>
<speciesReference species="M"/>
</listOfReactants>
<listOfProducts>
<speciesReference species="C2"/>
<speciesReference species="YP"/>
</listOfProducts>
<kineticLaw>
<math xmlns="http://www.w3.org/1998/Math/MathML">
<apply> <times/> <ci> k6 </ci> <ci> M </ci> </apply></math>
<listOfParameters> <parameter id="k6" value="1">
</listOfParameters>
</kineticLaw>
</reaction>
+ <reaction id="Reaction2" name="cdc2k phosphorylation">
... more reactions
</listOfReactions>
</model>
</sbml>

```

Species:

Id	Name	Compartment
C2	cdc2k	cell
M	p-cyclin_cdc2	cell
YP	p-cyclin	cell
....	....	....

Reaction:

Id	Name	Annotation	Formula
Reaction1	cyclin_cdc2k dissociation	<annotation ...>	<kinetic_law ...>
Reaction2	cdc2k phosphorylation	<annotation ...>	<kinetic_law ...>
....	....	....	....

Reactants:

Id	Species
Reaction1	M
Reaction2	....
....	....

Products:

Id	Species
Reaction1	C2
Reaction1	YP
....	....

mars 2009 28

LiU expanding reality

## SQL and Xpath/XQuery

```
select, r.name, s.name  
from reaction r, products p, species s  
where r.id = p.id and p.species=s.id;
```

```
xquery  
for $y in db2-fn:xmlcolumn('SBML_DATA.SBML_DOC')  
/model/listOfReactions/reaction/listOfModifiers/modifierSpeciesReference,  
$z in db2-fn:xmlcolumn('SBML_DATA.SBML_DOC')/model/listOfSpecies/species[@id = $y/@species]  
return <product> {$y/././@name} {$z/@name} </product>
```

```
SELECT p.reaction, species.name  
from species,  
(SELECT X.*  
FROM reactome_data,  
XMLTABLE ('$d/model/listOfReactions/reaction/listOfModifiers/modifierSpeciesReference' passing  
reactome_doc as "d"  
COLUMNS  
product VARCHAR(200) PATH '@species',  
reaction VARCHAR(200) PATH '././@id') AS X) p  
where p.product=species.id
```

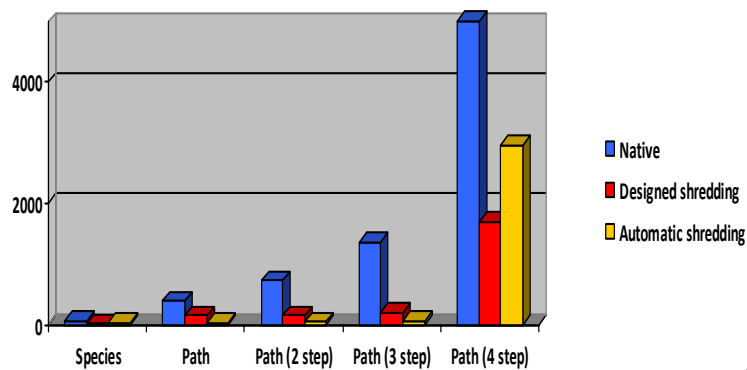
mars 2009

29

LiU

expanding reality

## Efficiency: Increasing query complexity



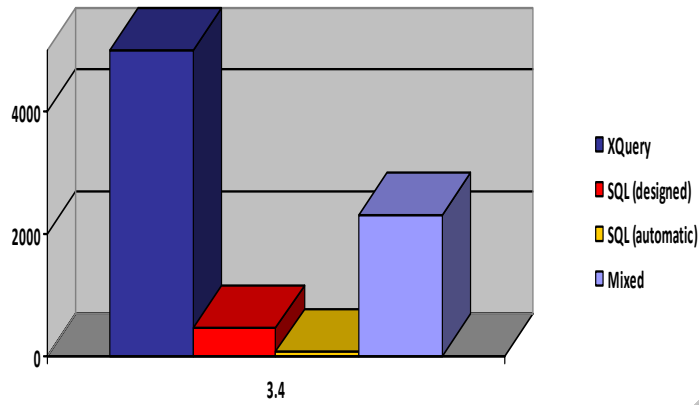
mars 2009

30

LiU

expanding reality

## Efficiency: Combining representations

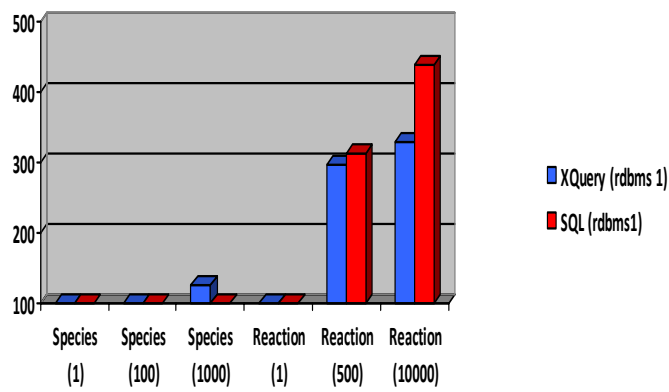


mars 2009 31

LiU

expanding reality

## Efficiency: Return the result as XML



mars 2009 32

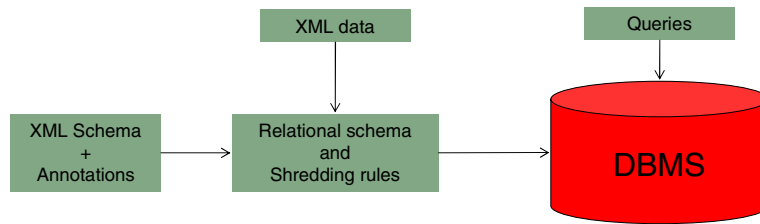
LiU

expanding reality



## Tool development: HShreX

- Need a tool to speed up the process



- (Further development of work by Amer-Yahia et al.)

mars 2009

33

LiU

expanding reality

## Working with HShreX:

```

<?xml version="1.0" encoding="UTF-8"?>
<xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema"
  xmlns:shrex="http://www.cse.ogi.edu/shrex">
  <xs:element name="families">
    <xs:complexType>
      <xs:sequence maxOccurs="unbounded">
        <xs:element name="family" type="familyType"/>
      </xs:sequence>
    </xs:complexType>
  </xs:element>

  <xs:complexType name="familyType">
    <xs:sequence>
      <xs:element name="parent" type="parentType" />
      <xs:element name="child" type="childType" />
    </xs:sequence>
  </xs:complexType>

  <xs:complexType name="parentType">
    <xs:sequence>
      <xs:element name="name" type="xs:string"/>
      <xs:element name="job" type="xs:string"/>
    </xs:sequence>
  </xs:complexType>

  <xs:complexType name="childType">
    <xs:sequence>
      <xs:element name="name" type="xs:string"/>
      <xs:element name="school" type="xs:string"/>
    </xs:sequence>
  </xs:complexType>
</xs:schema>
  
```

### Families

Id	Pid
0	-

### Families\_family

Id	Pid
1	0

### Families\_family\_parent

Id	Pid	Name	Job
2	1	Lena	Lektor

### Families\_family\_child

Id	Pid	Name	School
3	1	Ludvig	Skolan

mars 2009

34

LiU

expanding reality

## Working with HShreX:

```
<?xml version="1.0" encoding="UTF-8"?>
<xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema"
  xmlns:shrex="http://www.cse.ogi.edu/shrex">

  <xs:element name="families">
    <xs:complexType>
      <xs:sequence maxOccurs="unbounded">
        <xs:element name="family" type="familyType"/>
      </xs:sequence>
    </xs:complexType>
  </xs:element>

  <xs:complexType name="familyType">
    <xs:sequence>
      <xs:element name="parent" type="parentType" >
        <xs:element name="child" type="childType"
          shrex:maptoxml="true">
        </xs:sequence>
      </xs:complexType>

  <xs:complexType name="parentType">
    <xs:sequence>
      <xs:element name="name" type="xs:string"/>
      <xs:element name="job" type="xs:string"/>
    </xs:sequence>
  </xs:complexType>

  <xs:complexType name="childType">
    <xs:sequence>
      <xs:element name="name" type="xs:string"/>
      <xs:element name="school" type="xs:string"/>
    </xs:sequence>
  </xs:complexType>
</xs:schema>
```

### Families

Id	Pid
0	-

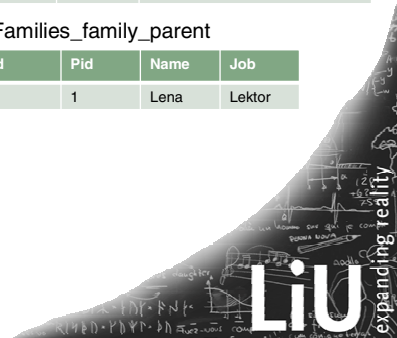
### Families\_family

Id	Pid	Child
1	0	<child>             <name>Ludvig</name>             <school>Skolan/school</school>           </child>

### Families\_family\_parent

Id	Pid	Name	Job
2	1	Lena	Lektor

mars 2009



## Working with HShreX:

```
<?xml version="1.0" encoding="UTF-8"?>
<xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema"
  xmlns:shrex="http://www.cse.ogi.edu/shrex">

  <xs:element name="families">
    <xs:complexType>
      <xs:sequence maxOccurs="unbounded">
        <xs:element name="family" type="familyType"/>
      </xs:sequence>
    </xs:complexType>
  </xs:element>

  <xs:complexType name="familyType">
    <xs:sequence>
      <xs:element name="parent" type="parentType"
        shrex:tablename="person" >
        <xs:element name="child" type="childType"
          shrex:tablename="person" >
        </xs:sequence>
      </xs:complexType>

  <xs:complexType name="parentType">
    <xs:sequence>
      <xs:element name="name" type="xs:string"/>
      <xs:element name="job" type="xs:string"/>
    </xs:sequence>
  </xs:complexType>

  <xs:complexType name="childType">
    <xs:sequence>
      <xs:element name="name" type="xs:string"/>
      <xs:element name="school" type="xs:string"/>
    </xs:sequence>
  </xs:complexType>
</xs:schema>
```

### Families

Id	Pid
0	-

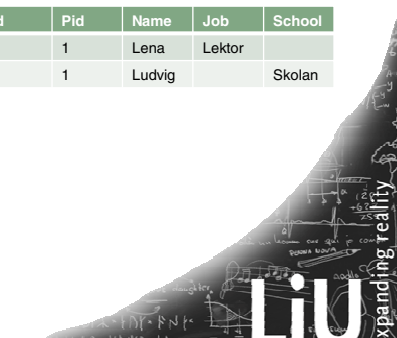
### Families\_family

Id	Pid
1	0

### Person

Id	Pid	Name	Job	School
2	1	Lena	Lektor	
3	1	Ludvig		Skolan

mars 2009



## Working with HShreX:

```
<?xml version="1.0" encoding="UTF-8"?>
<xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema"
  xmlns:shrex="http://www.cse.ogi.edu/shrex">

  <xs:element name="families">
    <xs:complexType>
      <xs:sequence maxOccurs="unbounded">
        <xs:element name="family" type="familyType"/>
      </xs:sequence>
    </xs:complexType>
  </xs:element>

  <xs:complexType name="familyType">
    <xs:sequence>
      <xs:element name="parent" type="parentType"
        shrex:withparent="true">
      <xs:element name="child" type="childType">
    </xs:sequence>
  </xs:complexType>

  <xs:complexType name="parentType">
    <xs:sequence>
      <xs:element name="name" type="xs:string"/>
      <xs:element name="job" type="xs:string"/>
    </xs:sequence>
  </xs:complexType>

  <xs:complexType name="childType">
    <xs:sequence>
      <xs:element name="name" type="xs:string"/>
      <xs:element name="school" type="xs:string"/>
    </xs:sequence>
  </xs:complexType>
</xs:schema>
```

### Families

Id	Pid
0	-

### Families\_family

Id	Pid	Name	Job
1	0	Lena	Lektor

### Families\_family\_child

Id	Pid	Name	School
3	1	Ludvig	Skolan

mars 2009

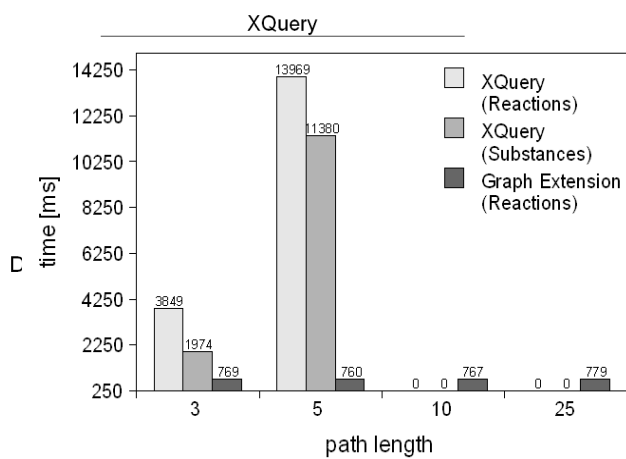
## Future challenges

- Hybrid storage, to be as efficient as possible in each case.
- Special solutions for networks and hierarchies

mars 2009

38

## XML with graph queries:



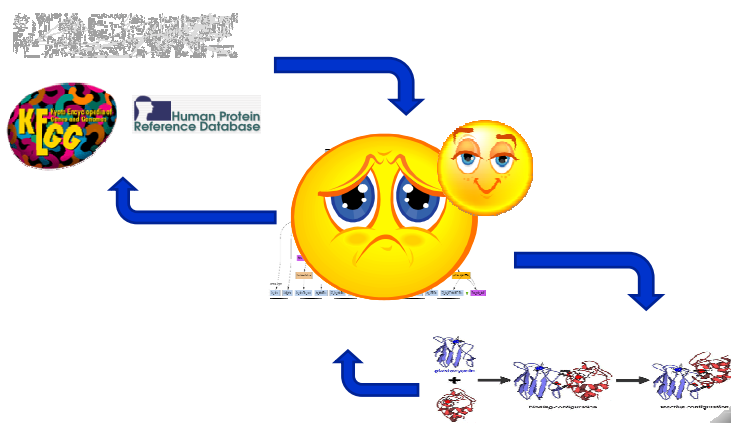
mars 2009

39

LiU

expanding reality

## Workflows for data exploration



mars 2009

40

LiU

expanding reality

# Capturing provenance

- Provenance of scientific artifacts is necessary to reproduce, validate and share scientific results
- Provenance can be as important as the results!

Dictionary

prov•e•nance |ˈprəvənəns|

noun

the place of origin or earliest known history of something : *an orange rug of Iranian provenance.*

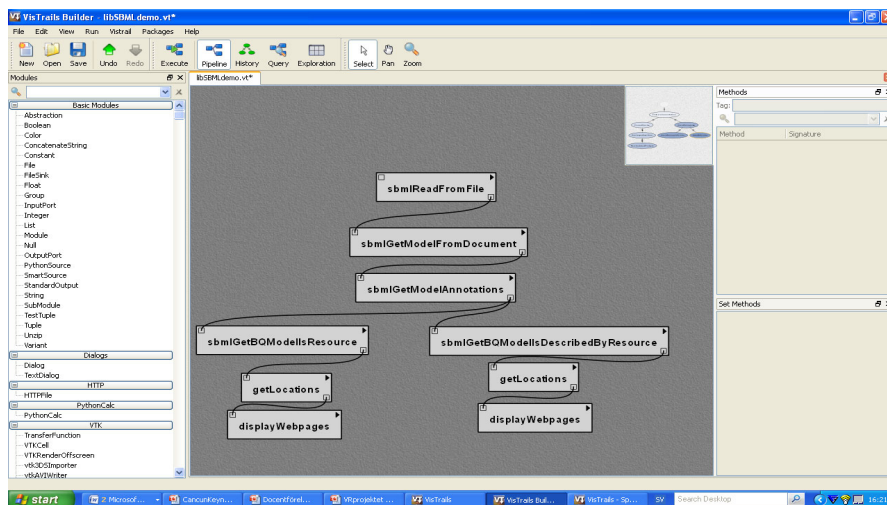
- the beginning of something's existence; something's origin : *they try to understand the whole universe, its provenance and fate.*

See note at **ORIGIN**.

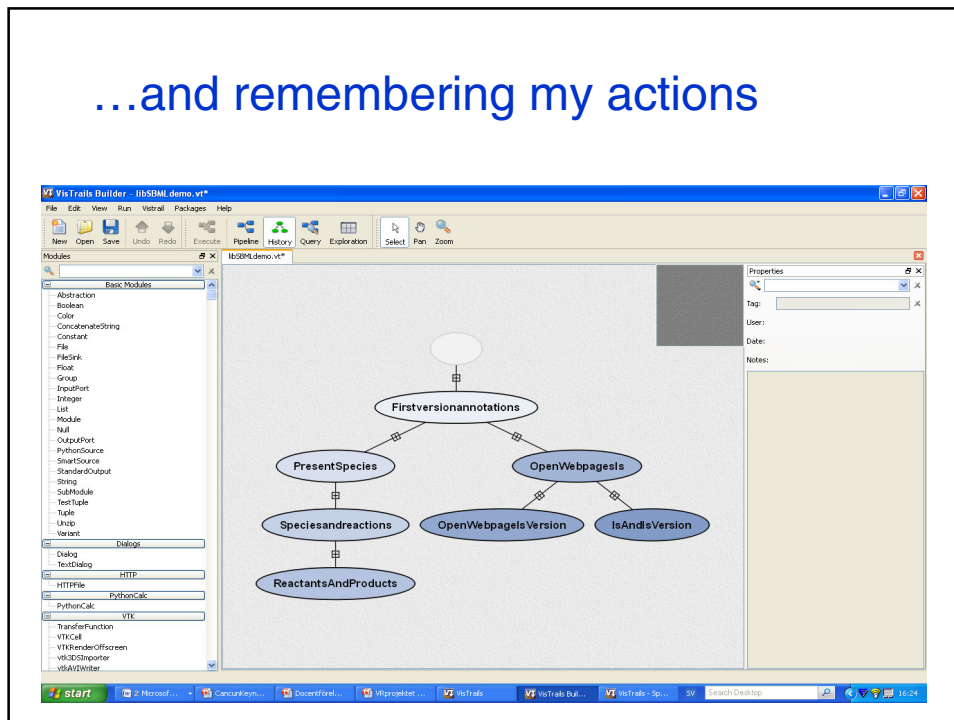
- a record of ownership of a work of art or an antique, used as a guide to authenticity or quality : *the manuscript has a distinguished provenance.*

**ORIGIN** late 18th cent.: from French, from the verb *provenir* 'come or stem from,' from Latin *provenire*, from *pro-* 'forth' + *venire* 'come.'

# Specifying my steps....

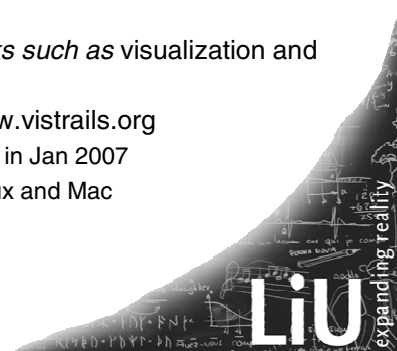


...and remembering my actions



## The Vistrails system (Freire et al. University of Utah)

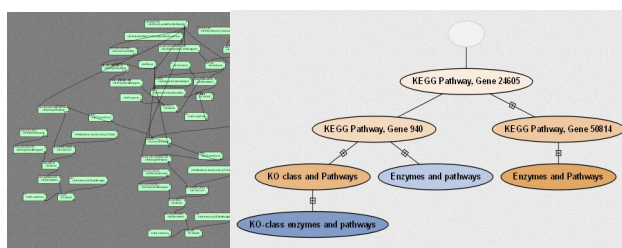
- *Vision: Provenance enable the world*
- Comprehensive provenance infrastructure for computational tasks
  - Captures provenance transparently
  - Provides intuitive query interfaces for exploring provenance data
  - Supports collaboration
- Designed to support *exploratory tasks such as* visualization and data mining
- VisTrails system is open source: [www.vistrails.org](http://www.vistrails.org)
  - >2,000 downloads since beta release in Jan 2007
  - 100% Python--runs on Windows, Linux and Mac



## Interesting issues:

### ➤ Reuse of other peoples efforts

- Provenance server
- Presentation of search results
- Co-work
- Interoperability between systems



mars 2009

45

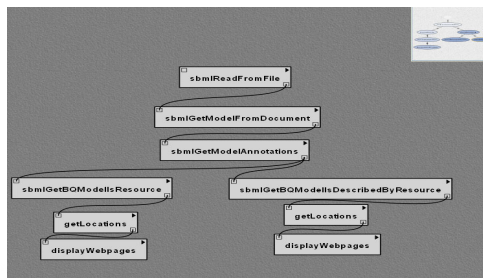
LiU

expanding reality

## Interesting issues:

### ➤ Bio-specific version of Vistrails

- Handling XML – LibSBML, XQuery
- Annotations – use of ontologies
- Easy to use module library for bioinformatics
- Combining webdata, own results and visualization



mars 2009

46

LiU

expanding reality

# Questions?

People that have contributed to this work:

**Linköpings Universitet:** Patrick Lambrix, José M Pena, He Tan,  
Tommy Ellqvist, David Hall, Mikael Åsberg

**University of Utah:** Juliana Freire, David Koop, Lauro Lins,  
Claudio Silva

**University of Rostock:** Dagmar Köhn, Andreas Heuer, Adelinde  
Uhrmacher, Stefan Schmidt

Thanks!

