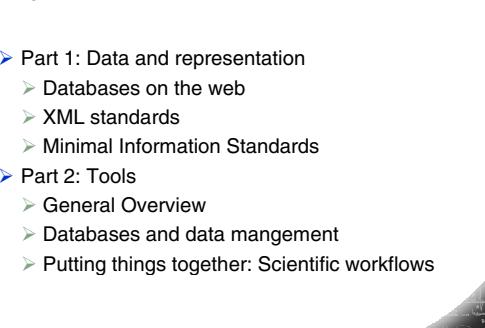
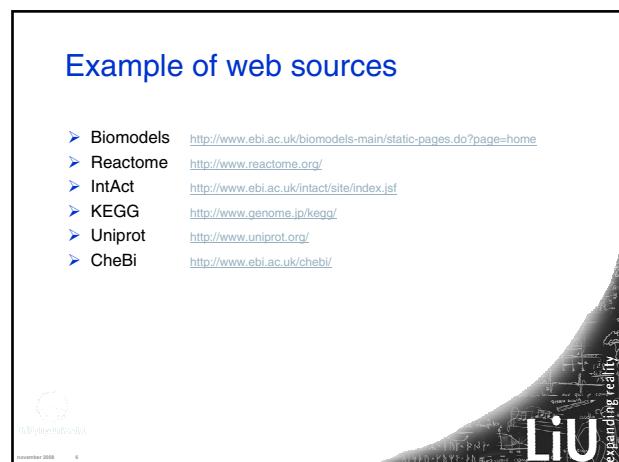
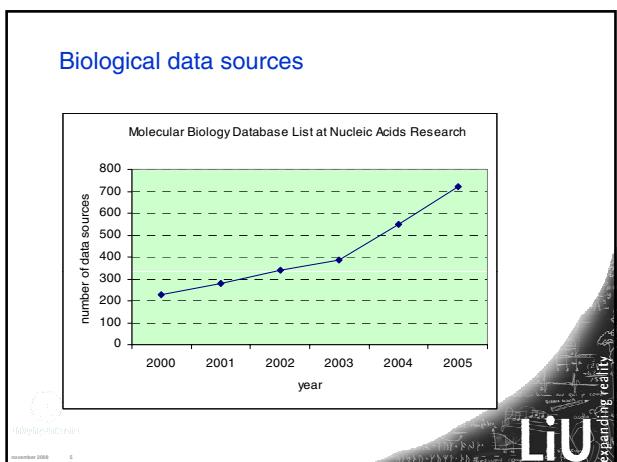


Topics

- Part 1: Data and representation
 - Databases on the web
 - XML standards
 - Minimal Information Standards
- Part 2: Tools
 - General Overview
 - Databases and data management
 - Putting things together: Scientific workflows

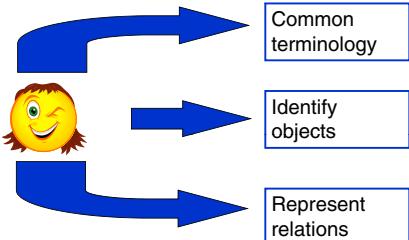




So ...

- Availability of web sources within bioinformatics is unique
- How can we make best use of them?

How to make use of this data?



How to represent data ...

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

SBML

- Defined by: System Biology Development workbench group.
- Aim: Standard for information exchange.
- Mathematical formulas can be defined, discrete events.

```
<model>
  <listOfCompartments>
    <listOfSpecies>
      <listOfReactions>
        <reaction>
          <listOfReactants>
            <listOfProducts>
            <listOfModifiers>
          </reaction>
        </listOfReactions>
      </listOfSpecies>
    </listOfCompartments>
  </model>
```

XML model example

```
<model id="Tyson1991CellModel_6"
  name="Tyson1991_CellCycle_6var">
  <listOfSpecies>
    +<species id="C22" 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>
    <rdflib>
      <rdflib:resource>http://www.reactome.org#REACT_6308</rdflib>
    </annotation>
    <listOfReactants>
      <speciesReference species="M"/>
    </listOfReactants>
    <listOfProducts>
      <speciesReference species="C22"/>
      <speciesReference species="YP"/>
    </listOfProducts>
    <listOfParameters>
      <math xmlns="http://www.w3.org/1998/Math/MathML">
        <apply> <times> </apply> k6 <ci> M </ci> </apply></math>
      </listOfParameters>
    <KineticLaw>
      <reaction id="Reaction2" name="cdc2k phosphorylation">
        ... more reactions
      </reaction>
    </KineticLaw>
  </reaction>
  <listOfReactions>
    ... more reactions
  </listOfReactions>
</model>
</sbml>
```

PSI MI

- Defined by: Proteomics Standards Initiative.
- Aim: Standard for representation of molecular interaction.
- Thorough description of experiments and protein structure.

```
<entry>
  <experimentList>
    <interactorList>
      <interactionList>
        <interaction>
          <experimentList>
            <participList>
              </interaction>
            </participList>
          </experimentList>
        </interactionList>
      </interactorList>
    </experimentList>
  </entry>
```

Available XML standards

Name	Vet.	Year	Defined by	Purpose	Tools	Data
SBML	2	2003	Systems Biology Working Group	A computer-readable format for representing mathematical models of biological systems.	Many tools available.	Data available from many databases, e.g. KEGG and BioCyc.
PSI-MI	2.5	2003	Proteomics Standards Initiative	Standardized data representations for protein-protein interactions to facilitate data comparison, reuse and analysis.	Tools for viewing and analysis.	Data available from many sources, e.g. UniProt, DIP and MINT.
BioPAX	2	2005	The BioPAX group	A collaborative effort to create a data exchange format for biological pathways data.	Existing tools for OWL, such as Protégé.	Datasets available from Reactome.
CeMML	1.1	2002	University of Auckland and Physics Sciences	Standardization of models of cellular and subcellular processes.	Tools for publication, visualization, creation and simulation.	CeMML Model Repository (~240 models).
CML	2.2	2003	Peter Murray-Rust, Henry S. Rzepa	Interchange of chemical structures for the Internet and other.	Molecule browsers, editors.	BioCyc.
EMBL-ML	1.0	2005	EBI	More mobility and flexibility of exchange of nucleotide sequence information.	API support in Bio2RDF.	EMBL.
INSDSeq	1.9	2003	International Nucleotide Sequence Database Collaboration	INSDSeq uses ASN.1 for the storage and retrieval of data such as nucleotide and protein sequences. Data encoded in ASN.1 can be converted to XML.	API support in Bio2RDF.	EMBL, DBI and Genbank.
Sequencher	n/a	n/a	NCBI	Sequencher uses ASN.1 for the storage and retrieval of data such as nucleotide and protein sequences.	SGD BioWarehouse and ProteinStructureFactory's ORFeus.	Emaze.
BSML	3.1	2002	LabBook.com	Facilitating the interchange of data for scientific experiments within the same or different software environments.	Labbook's Generic Experiment Viewers.	Previously supported by EMBL.
HUPO-ML	0.9	2003	JHUPO	A protein-centered markup language for exchange of proteomic data between researchers.	HUPO-ML Editor.	
MAGE-ML	1.1	2003	MGED	To facilitate the exchange of gene expression data.	Converters.	

November 2008 12

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



November 2008 13

LiU expanding reality

BioPAX

- ▼ ● entity
 - physicalEntity
 - dna
 - rna
 - protein
 - complex
 - smallMolecule
- ▼ ● interaction
 - physicalInteraction
 - conversion
 - complexAssembly
 - biochemicalReaction
 - transportWithBiochemicalReaction
 - transport
 - transportWithBiochemicalReaction
 - control
 - catalysis
 - modulation
 - pathway

November 2008 15

LiU expanding reality

- Defined by: The BioPAX working groups.
- Aim: General format for interactions.
- Ontology, implemented in OWL
- Pathway and molecular interactions.

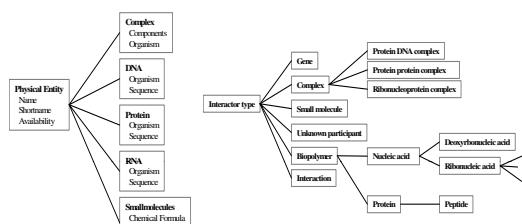
Summary of standards

Name	Substances			Interactions	Pathways	Compounds	Organism	Experiments
	DNA, RNA	Protein	Other					
SBML	UL	UL	UL	SOL	SOL	SL		
PSI-MI	SOL	SOL	SOL	SOL	L	SL	S	
BioPAX	SOL	SOL	SOL	SOL	S	L	L	
CellML	L	L	L	S	S	U	U	
CML				S	S			
ENB-XML	VL	SL					L	
INSDSeq	SL	SL					L	
Sequencher	SL	SL					L	
BSML	SL	SL					L	S
HUPO-ML	SL	SL					L	S
MAGE-ML	L	L					L	S
mzXML								SO
mzData								S
AGML							U	S

November 2008 16

LiU expanding reality

Type of objects



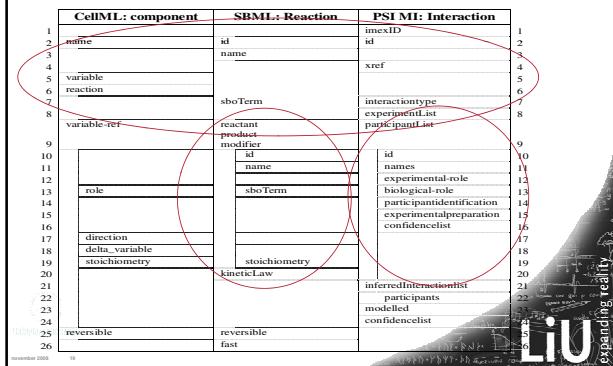
Representation of objects

SBML: Species	PSI MI: Interactor	CellML: component
1 id	1 id	1 name
2 name	2 names	2 datatype
3 speciesType	3 ref	3 alternative
4 speciesType	4 interactome	4 emtabio_entry
5 organism	5 organism	5 emtabio_species
6 ncbiTaxId	6 names	6 emtabio_species
7 celltype	7 celltype	7 emtabio_species
8 compartment	8 compartment	8 compartment
9 mouse	9 mouse	9 mouse
10 sequence	10 sequence	10 emtabio_acs
11 sequence	11 sequence	11 variable
12 initialAmount	12 initialAmount	12 initialvalue
13 initialConcentration	13 initialConcentration	13 units
14 substanceUnits	14 substanceUnits	14
15 spatialUnits	15 spatialUnits	15
16 boundarySubstanceUnits	16 boundarySubstanceUnits	16
17 boundaryCondition	17 boundaryCondition	17
18 charge	18 charge	18
19 constant	19 constant	19
20		20
21		21

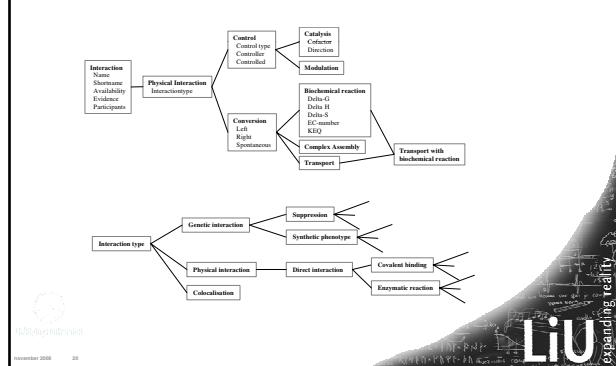
November 2008 47

LiU expanding reality

Representation of interactions



Interaction types



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
-

Part 2: Tools

- Many of the standards are supported by software tools:
 - SBML Tools: http://sbml.org/SBML_Software_Guide/SBML_Software_Matrix
 - CellML Tools: <http://www.cellml.org/tools>
- Data management
- Scientific workflows

Data management of XML data

- How can the data be efficiently stored and accessed?
- Native XML databases
- Hybrid solutions

XML as a data model

- 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

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]

```

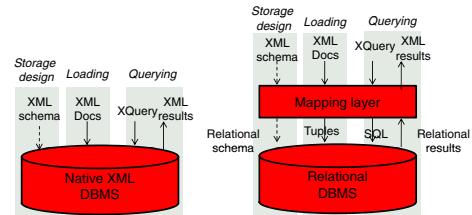


XML as a data model

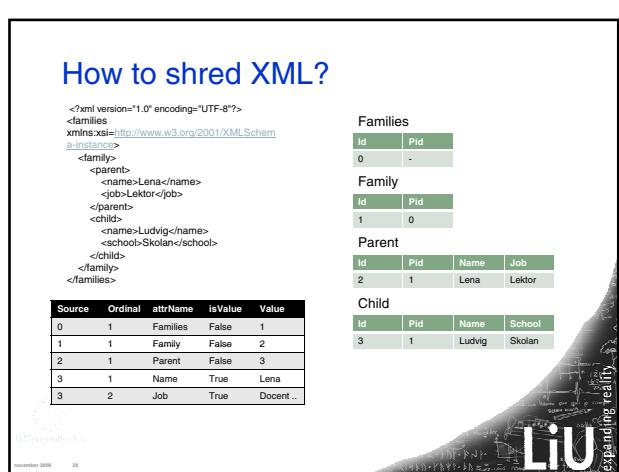
- XML is richer than the relational model
 - Tree structure,
 - Order
 - ...
 - Vary from highly structured to unstructured
 - Database export
 - ...
 - Annotated text documents
 - Can contain links to other type of entities



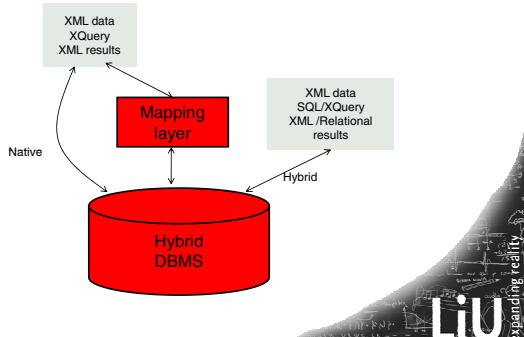
Storage possibilities for XML



Shredding



Hybrid XML Storage



New possibilities...

```

<model id="Tyros1991_Cdc2Model_6"
  name="Tyros1991_Cdc2Model_6var">
  <listOfSpecies>
    <!-- Species -->
    <species id="C2" name="cdc2k" compartment="cell">
      <!-- Species M -->
      <species id="M" name="p-cyclin" compartment="cell">
        <!-- Species YP -->
        <species id="YP" name="p-cyclin" compartment="cell"> ...
    </listOfSpecies>
    <listOfCompartment>
      <compartment id="cell" name="Cell/Cycle" ...>
    </listOfCompartment>
  </listOfSpecies>
  <annotation id="Reaction1" name="cyclin_cdc2k_dissociation">
    <reaction id="R1">
      <listOfReactants>
        <!-- Annotations -->
        <listOfAnnotations>
          <annotation id="r1d1" resource="http://www.reactome.org/REACT_63087"/>
          <annotation id="r1d2" resource="http://www.genontology.org/I GO:0000079"/>
        </listOfAnnotations>
        <listOfReactants>
          <speciesReference species="M"/>
        </listOfReactants>
        <listOfProducts>
          <speciesReference species="C2">
            <speciesReference species="YP"/>
          </listOfProducts>
        <listOfKinetics>
          <kineticLaw id="K1" ...>
            <math ...>http://www.w3.org/1998/Math/MathML</math>
            <param id="K1s" value="1" />
          </kineticLaw>
        <listOfParameters>
          <parameter id="K2" value="1" />
        </listOfParameters>
        <listOfKinetics>
          <kineticLaw id="K3" ...>
            <math ...>http://www.w3.org/1998/Math/MathML</math>
            <param id="K3s" value="1" />
          </kineticLaw>
        <listOfParameters>
          <parameter id="K4" value="1" />
        </listOfParameters>
      </listOfReactants>
      <listOfReactions>
        <!-- Reaction id="Reaction1" name="cdc2k phosphorylation" -->
        <reaction id="R2" ...>
          <listOfReactions>
            <!-- More reactions -->
          </listOfReactions>
        </reaction>
      </listOfReactions>
    </listOfKinetics>
  </model>
  <listOfSBMLs>
    <sbml id="1" version="2" ...>
  </listOfSBMLs>
</model>

```

Species:		
ID	Name	Compartment
C2	cdc2k	cell
M	p-cyclin_cdc2	cell
YP	p-cyclin	cell

Reaction:

ID	Name	Annotations	Formula
Reaction1	cyclin_cdc2k dissociation	-annotation >>	<kinetic_law>
Reaction2	cdc2k phosphorylation	-annotation >>	<kinetic_law>

Reactants

<u>Reaction</u>	<u>Species</u>
Reaction1	M
Reaction2	...
....

Reaction



liu

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:fnxmlcolumn('SBML_DATA.SBML_DOC')
  /model/listOfReactions/reaction/listOfModifiers/modifierSpeciesReference,
  $z in db2:fnxmlcolumn('SBML_DATA.SBML_DOC')/model/listOfSpecies/species[@id = $y/@species]
return <products>{$y/../@name} {$z/@name}</products>

```

```

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

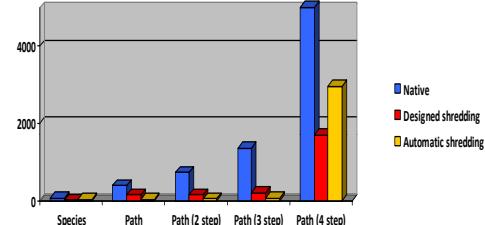
```



November 2008

31

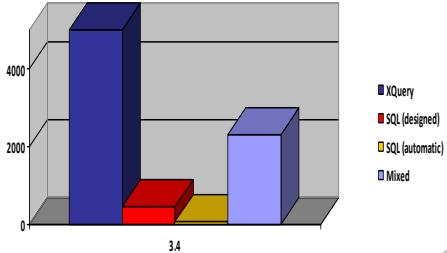
Efficiency: Increasing query complexity



November 2008

32

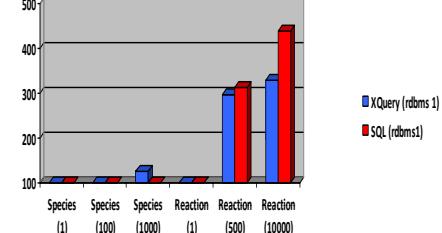
Efficiency: Combining representations



November 2008

33

Efficiency: Return the result as XML



November 2008

34

So...

- Native databases –
 - Easy to use but sometimes too inefficient
- Automatic shredding –
 - Can give datamodels that are hard to work with
- Manual shredding /Hybrid solutions –
 - Requires time consuming design

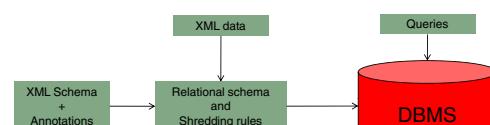


November 2008

35

Tool development: ShreX

- Need a tool to speed up the process



- Extension of an old tool to allow hybrid storage



November 2008

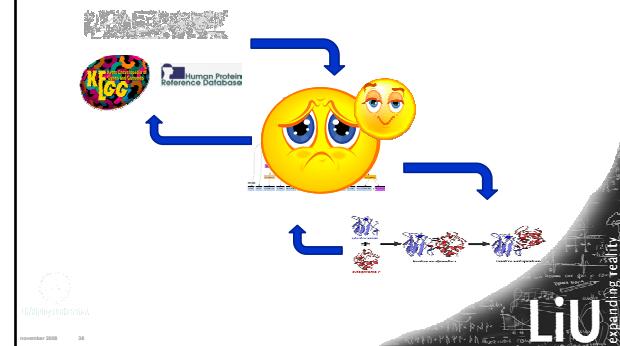
36

Demo ShreX



LiU
expanding reality

Workflows for data exploration



Capturing provenance

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

Dictionary

prove•nance (prāvənəns)
verb
the place of origin or earliest known history of something : *an orange rug of Indian 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' + *venir* 'come.'



LiU
expanding reality

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
 - >2,000 downloads since beta release in Jan 2007
 - 100% Python--runs on Windows, Linux and Mac

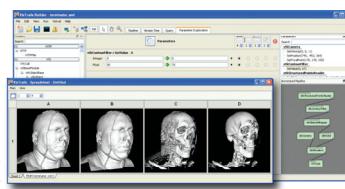
Demo Vistrails



LiU
expanding reality

Enhanced functionality

- Parameter exploration
- Query to find workflows of interest
- Compute difference between workflows
- Create an analogy based on computed differences



Interesting issues:

- Reuse of other peoples efforts
 - Provenance server
 - Co-work
- Bio-specific version of Vistrails
 - Handling SBML - Libsbml
 - Annotations – use of ontologies
 - Easy to use module library
 - Combining webdata, own results and visualization



Innsbruck 2008

40



expanding reality

Questions?

Thanks!



Working with ShreX:

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

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



Innsbruck 2008

41



expanding reality

Working with ShreX:

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

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



Innsbruck 2008

41



expanding reality

Working with ShreX:

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

Families	
Id	Pid
0	-

Families_family	
Id	Pid
1	0

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

shreX:tablename='person'

shreX:complexType name='parentType' shreX:withParent='true'



Innsbruck 2008

41



expanding reality

Working with ShreX:

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

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

shreX:complexType name='parentType' shreX:withParent='true'

shreX:complexType name='childType'



Innsbruck 2008

41



expanding reality