

Publishing and Linking Scientific Data

—Experience on Applying Linked Data to Scientific Database Project

SHEN Zhihong

Computer Network Information Center, CAS

2013/9/10

Outline

- Background
- Publishing scientific data
- Linking scientific data
- Summary
- Next work

Background - Scientific DataBase project

□ Scientific DataBase (SDB) project

- A long-term mission funded by CAS(Chinese Academy of Sciences) which started in 1986

- data from research, for research

- Collecting multi-discipline research data and promoting data sharing (2006-2010, during the period of the eleventh Five-Year-Plan of CAS)

- about **61** CAS institutes involved
 - Over **200TB** data available for open access and download



Background - Scientific Databases

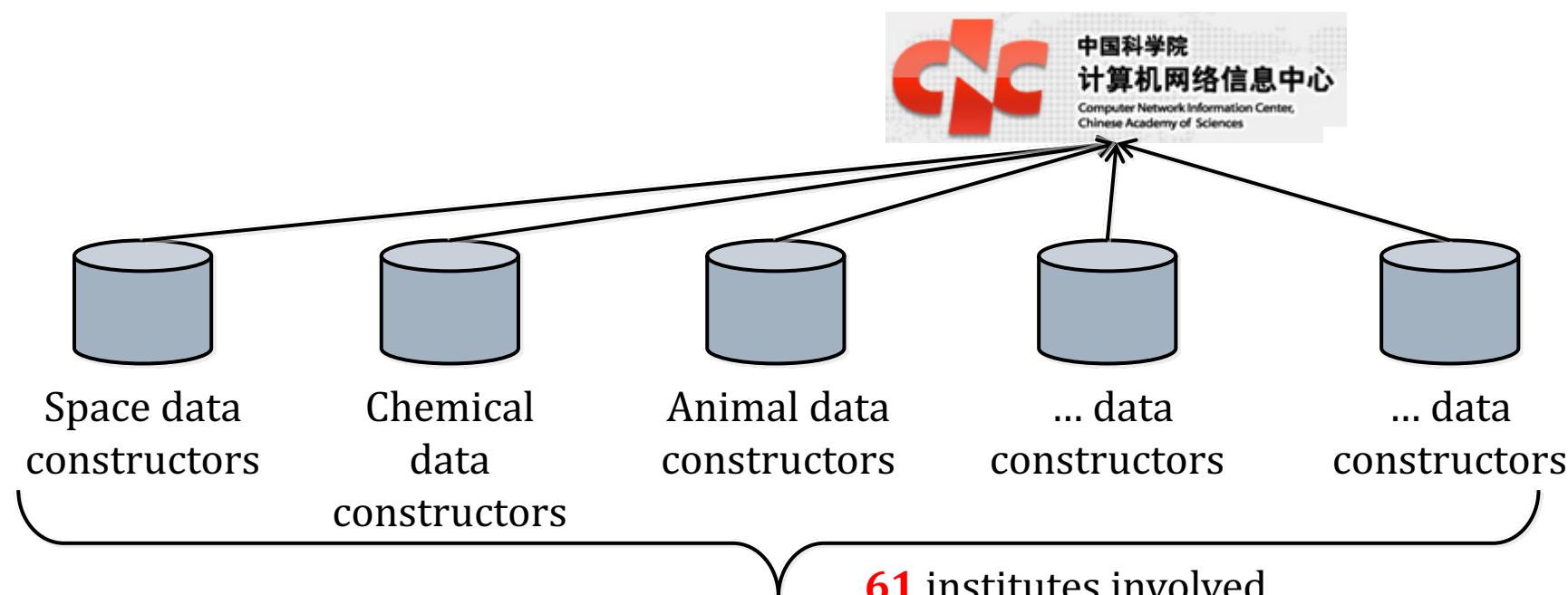
- SDB consists of 51 databases, including ...
- 8 Resource databases
 - Geo-Science
 - Biodiversity
 - Chemistry
 - Astronomy
 - Space Science
 - Micro biology and virus
 - Material science
 - Environment
- 37 institution databases
- 2 Reference databases
 - China Species
 - chemical compound
- 4 Application-Oriented databases
 - High Energy (ITER)
 - Western Environment Research
 - Ecology research
 - Qinghai Lake Research



Background

□ Role of CNIC in SDB project

1. plays the role of the organizer of the project
2. offers the cyberinfrastructure for the project
 - especially data storage environment
3. develops generic tools for building database systems, and develops SDB portal



Background

- Big question to CNIC
 - How to use all data involved and provide unified services for users? Not just a mess
- One possible idea is ...
 - **Put all data together into CNIC?**
 - This idea is very effective, but always means painful.
 - Data ownership: data is owned by different institutes, CNIC has no right to own these data and publish them
 - Inborn heterogeneity: different data, different format, different software environments ...
 - Keep data updated: a difficult task

Background

- **Data linking** makes sense
 - Data ownership: data is published in strict access control by data owners, we just link them!
 - Heterogeneity: data keeps original state, Data linking means no additional works for data owners
 - Keep data updated: data is updated by owners
- So, the best idea for data integration is
 - **Data Linking**
 - Not **moving** all data together, but **linking** them together

Background

□ How we do Data Linking

1. Data owners **publishes** scientific **data** and scientific **metadata** in a good format
2. CNIC Collects all allowed data descriptions and provides **searching** interfaces for users
3. CNIC generates and discovers **links** among data and shows related links according to user query

Background

□ key problems to be solved

1. **Publishing:** A good format? What is a good format?
Some data has little information, how to publish them?
2. **Link generation:** really need automatic methods and runnable frameworks, especially for a large volume of data

Background

Publishing scientific data

- RDF and Linked Data as good formats
- Publishing data records
- Publishing data files

Linking scientific data

Summary

Next Work

Publishing scientific data

□ What is a good format for scientific data?

1. Be able to represent structured data and semi-structured data
 - Relation records, MathML, CML, SMILES, ...
2. It is impossible to invent a new uniform format for all unstructured scientific data, then how about to be a good format for metadata?
 - Metadata of NetCDF, HDF, ...

Data types	examples	requirements
Structured data	Relational data	Be able to represent
Semi-structured data	MathML, CML, SMILES	Be able to represent
Unstructured data	NetCDF, HDF	Be able to represent its metadata

Publishing scientific data

What is a good format for scientific data?

3. able to describe the links between data
4. able to identify each data on the Web
5. machine-readable and understandable
 - May be consumed by programs
6. flexible schema, easy to extend, adding a new property is very simple

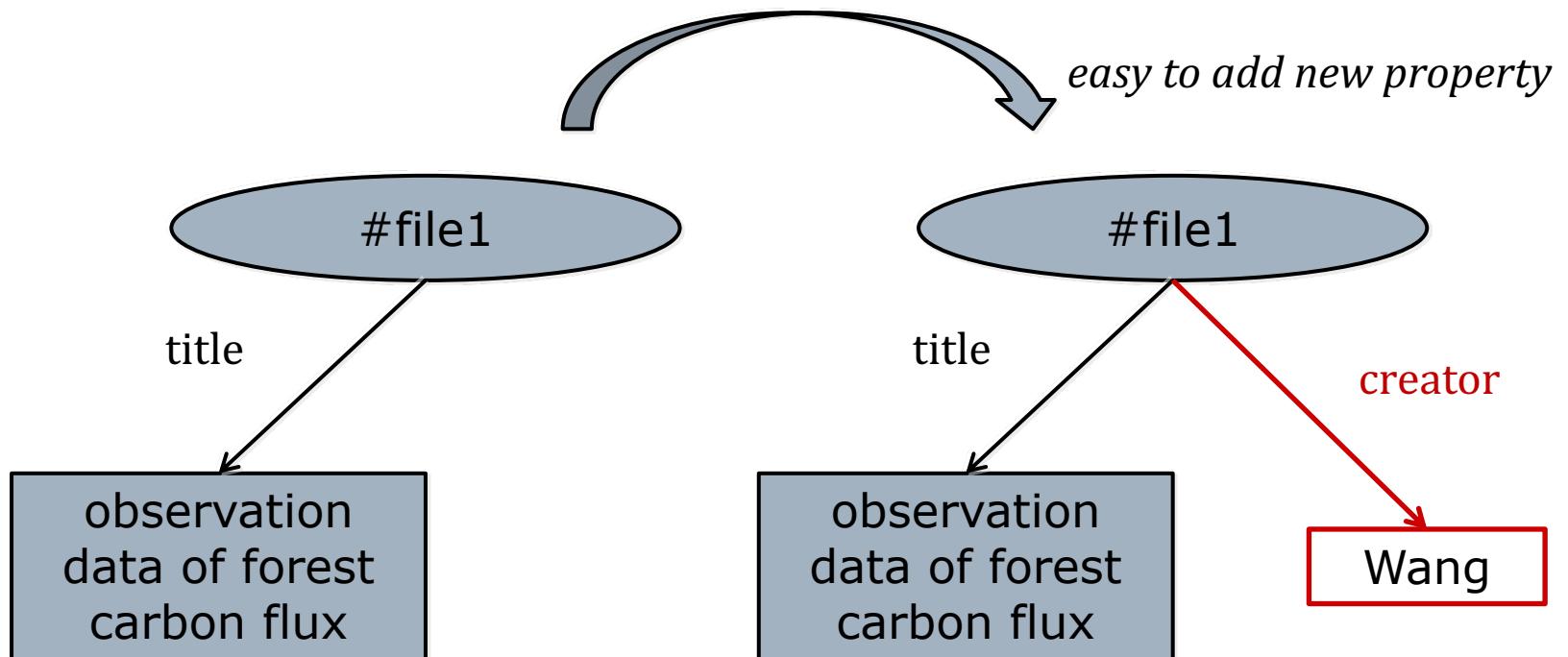
Publishing scientific data

- We chose **RDF** as data representation format
 - RDF: **Resource Description Framework**
 - RDF describes resources in terms of simple **properties** and property **values**
 - it is based upon the idea of *making statements about resources* (in particular web resources) in the form of **subject-predicate-object** expressions
 - e.g

<u><#dataSet1></u>	<u><title></u>	<u>“observation data of forest carbon flux”</u>
subject	predicate	object

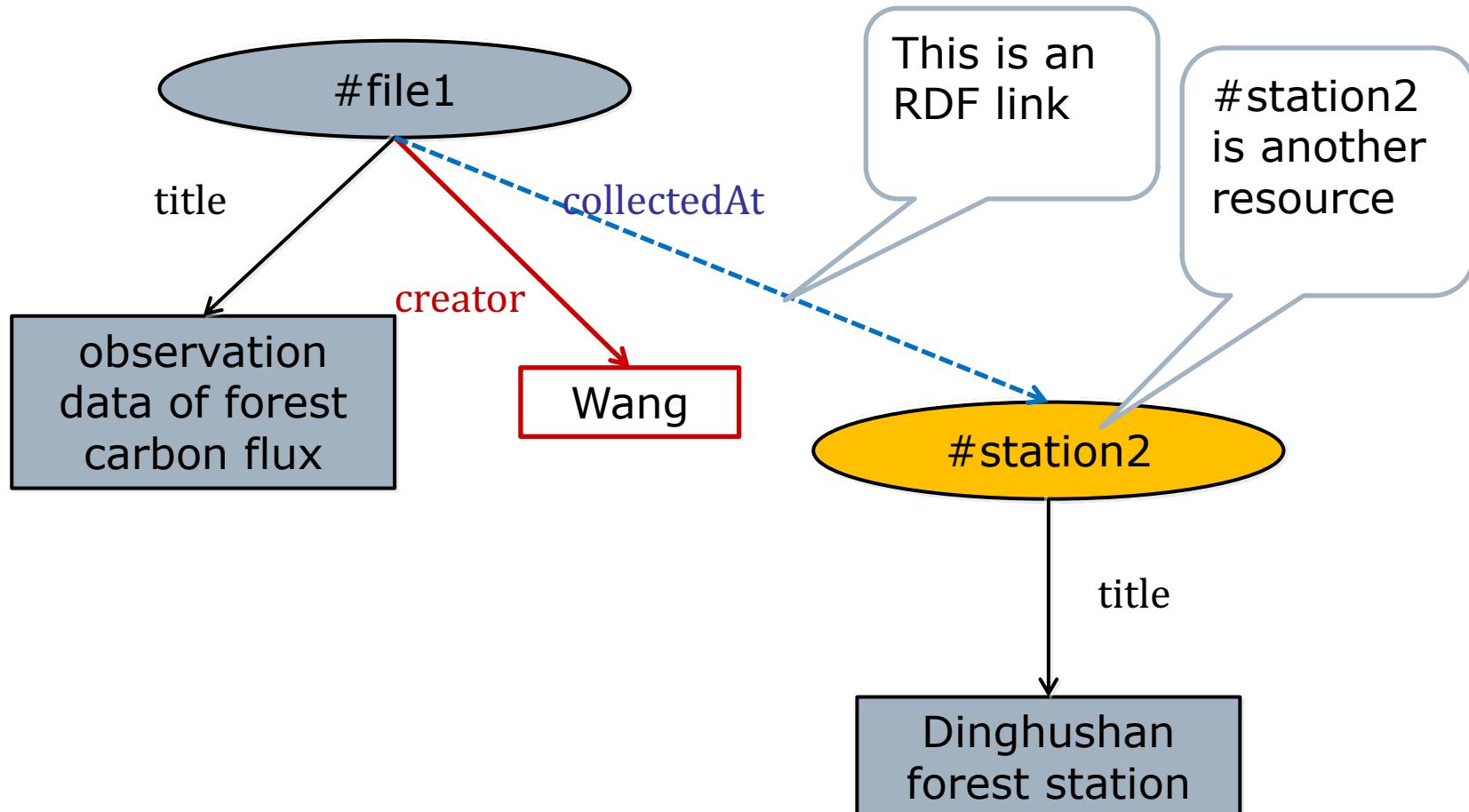
Publishing scientific data

- RDF statements about a resource are often represented as a graph
 - nodes: representing the resources, and their properties values.
 - arcs: representing properties of resources



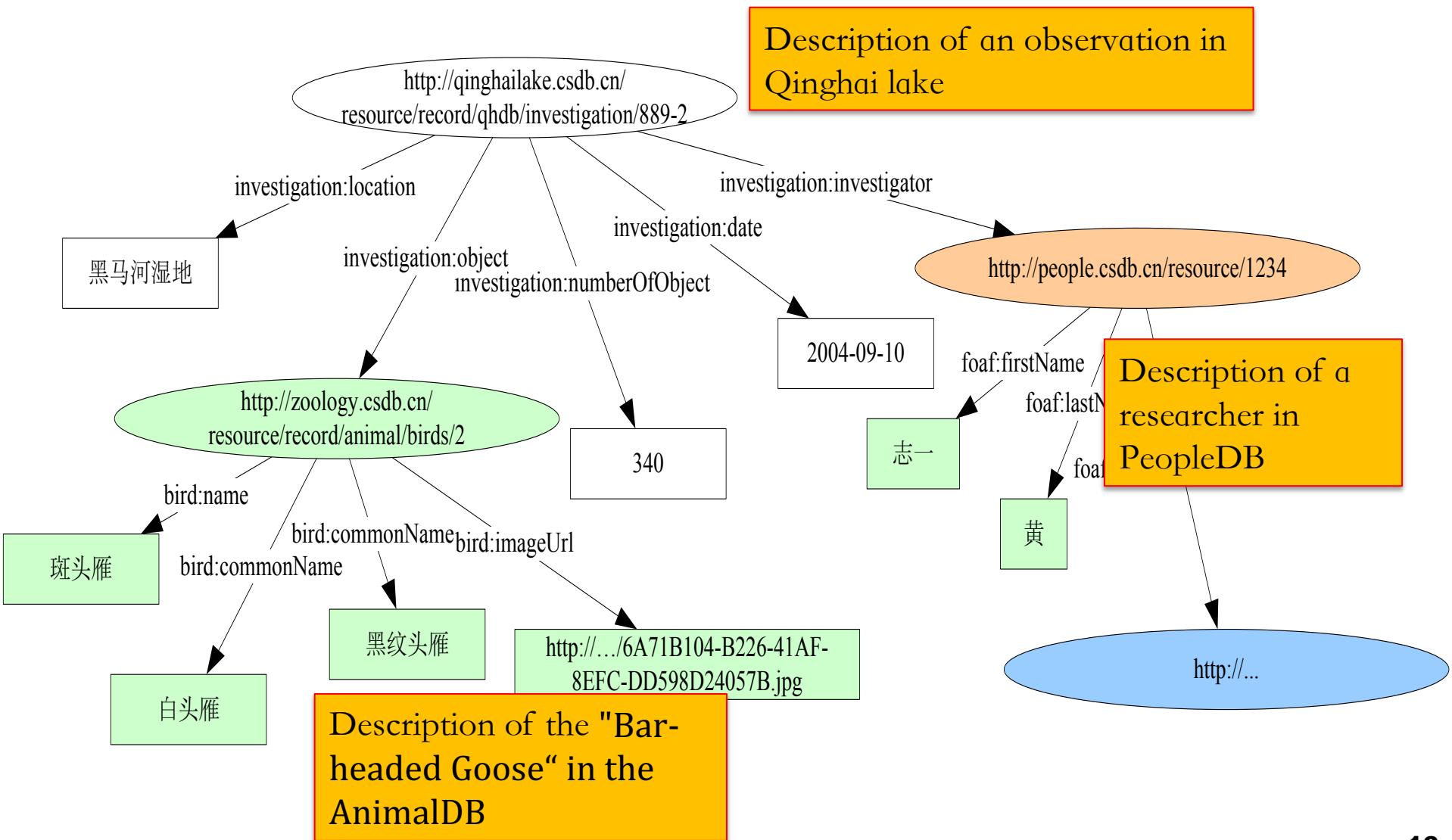
Publishing scientific data

- Most important, RDF data model enables people to set RDF links between data from different sources.



Publishing scientific data

□ More complex RDF data example



Publishing scientific data

- After choosing the data format, we chose **Linked Data** as the interoperation mechanism of data sources
 - **Linked Data**: Tim Berners-Lee coined the term Linked Data in 2006[1].
 - Connect Distributed Data across the Web
 - using the Web to create **typed links** between data from different sources



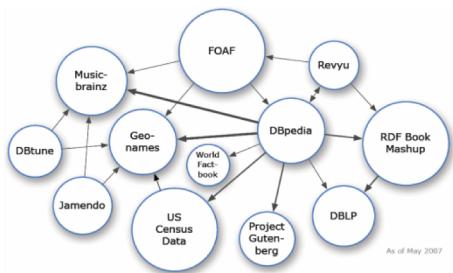
1. T. Berners-Lee, "Design issues: Linked data," Online at <http://www.w3.org/DesignIssues/LinkedData.html>, 2006

Publishing scientific data

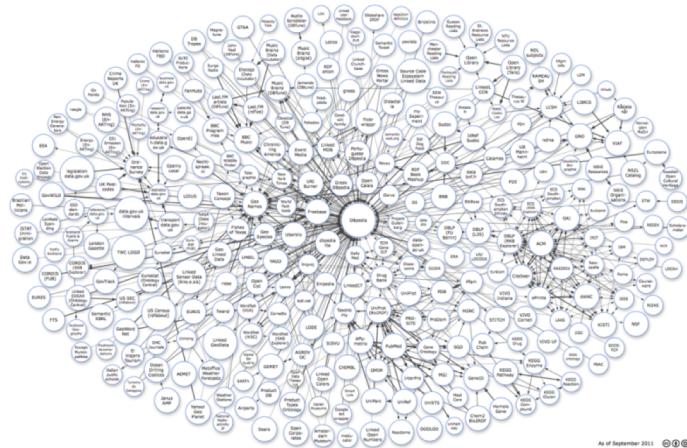
- 4 basic Linked Data Principles
 1. Use URIs as names for things
 2. Use HTTP URIs so that people can look up (dereference) those names.
 3. When someone looks up a URI, provide useful information.
 4. Include links to other URIs so that they can discover more things.

□ Linking Open Data

- Based on the concept of linked data, W3C initiated the **Linking Open Data** movement.
- It has driven many data sets published as Linked Data.
- By September 2011, LOD had covered about 295 datasets with 25 billion RDF triples and about 395 million RDF links.



1. “Linking Open Data cloud diagram, by Richard Cyganiak and Anja Jentzsch. <http://lod-cloud.net/>”



Publishing scientific data

□ Scientific Data samples – Linked Life Data

- Searches and explores over RDF statements from various sources including UniProt, PubMed, EntrezGene and 20 more...
- Performs complex SPARQL(RDF query) queries and retrieves more than one billion RDF resources.



1. Momtchev V, Peychev D, Primov T, et al. Expanding the pathway and interaction knowledge in linked life data[C]. In Proceedings of International Semantic Web Challenge, 2009.

Publishing scientific data

□ Scientific Data samples – LinkedGeoData

- uses the information collected by the OpenStreetMap project and makes it available as an RDF knowledge base according to the Linked Data principles.



The screenshot shows a map of Lukaskirche in Berlin. Red markers indicate specific locations, likely points of interest or data entities. A sidebar on the right lists RDF triples for the location:

- 1. Lukaskirche
amenity: place_of_worship
religion: christian
denomination: lutheran
- 2. Club 11
amenity: nightclub
- 3. amenity: recycling
- 4. amenity: parking
- 5. amenity: parking
- 6. amenity: recycling
- 7. amenity: parking
- 8. Lukas-Apotheke
amenity: pharmacy
- 9. amenity: school
- 10. amenity: parking
- 11. Johanneskirche
denomination: christian_community
religion: christian
amenity:

This facetted Linked Geo Data browser was developed by [AKSW research group](#).

1. Auer, Sören, Jens Lehmann, and Sebastian Hellmann. "Linkedgeodata: Adding a spatial dimension to the web of data." The Semantic Web-ISWC 2009. Springer Berlin Heidelberg, 2009. 731-746.

Publishing scientific data

□ Diseasesome

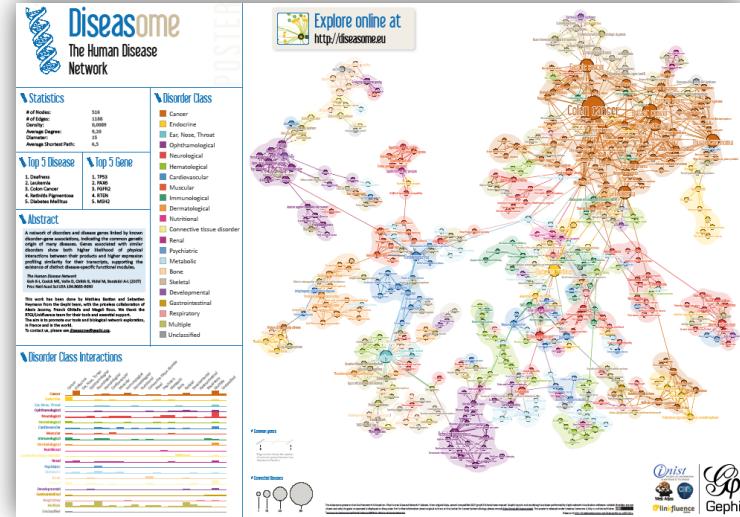
- publishes Linked Data of 4,300 disorders and **disease genes** linked by known disorder-gene associations for exploring all known phenotype and disease gene associations, indicating the common genetic origin of many diseases.

□ Linked Sensor Data

- the first open datasets for sensors and sensor observations, created at Knoesis Center, and converted from weather data at Mesowest.
- Contains descriptions of 20 thousand weather stations and 160 million observations.

□ GeoSpecies Knowledge Base

- Publishing information on Biological Orders, Families, Species as well as species occurrence records and related data, links to geonames, bio2rdf, dbpedia, freebase, umbel.

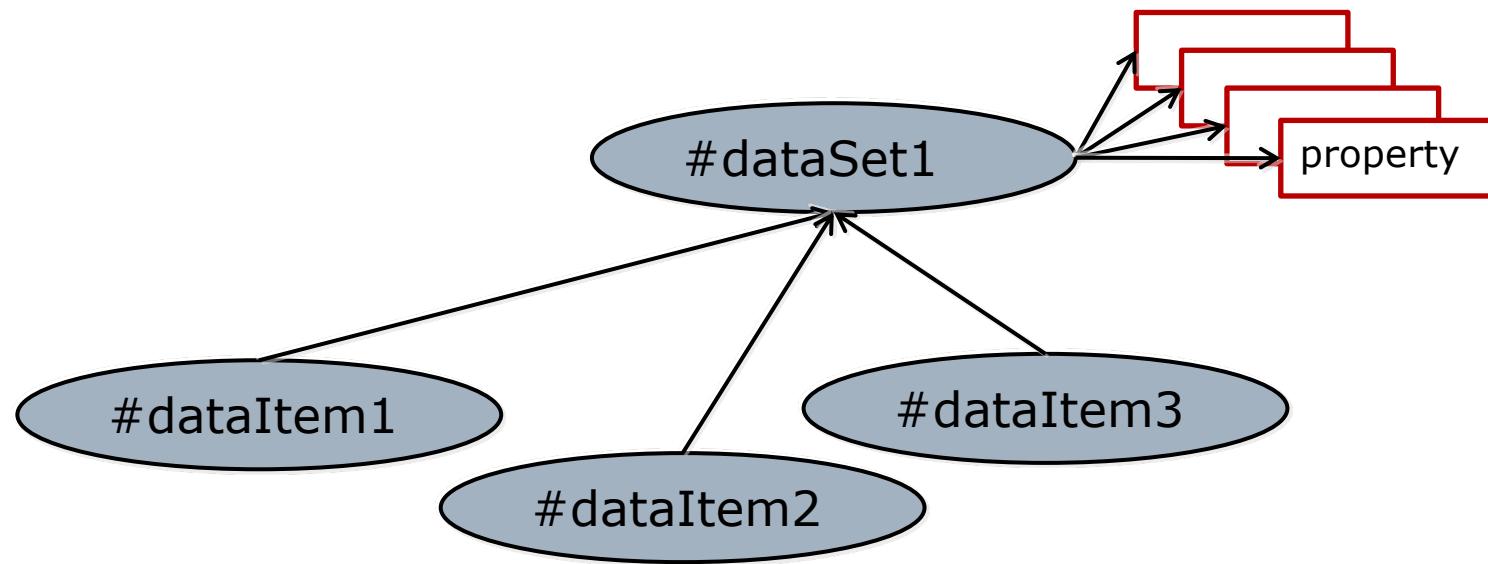


1. Diseasesome | Map: explore the human disease network. Dataset, interactive map and printable poster of gene-disease relationships. <http://diseasome.eu/map.html>
2. http://wiki.knoesis.org/index.php/SSW_Datasets

Publishing scientific data

□ We define data in 2 levels

1. Data item
 2. Data set: a set of data items
- Some data items have little information, but the information of data set can help users to find them



Publishing scientific data

- Dataset metadata is described with mixed RDF vocabularies
 - 1. Dublin-Core
 - <http://purl.org/dc/elements/1.1/>
 - 2. DC-TERMS
 - DCMI Metadata Terms
 - <http://purl.org/dc/terms#>
 - 3. PRISM
 - Publishing Requirements for Industry Standard Metadata
 - <http://prismstandard.org/namespaces/basic/2.0/>

Publishing scientific data

□ An example of dataset metadata:

```
<csdb:Database rdf:about="http://semweb.csdb.cn/csdb/resource/database/12053084">
    <dc:title>蒋家沟降水观测资料</dc:title>
    <prism:publicationDate rdf:datatype="http://www.w3.org/2001/XMLSchema#dateTime">2010-12-
03T09:51:59</prism:publicationDate>
    <dc:subject>地球科学</dc:subject>
    <rdfs:label>蒋家沟降水观测资料</rdfs:label>
    <prism:url>http://ns1.imde.ac.cn;http://www.mountain.csdb.cn/page/showEntity.vpage?uri=mouhazards.cata
HazardObsdata</prism:url>
    <dc:source>中国科学院东川泥石流观测研究站</dc:source>
    <dcterms:accrualPeriodicity>year</dcterms:accrualPeriodicity>
    <dc:contributor>中国科学院东川泥石流观测研究站</dc:contributor>
    <dc:type>004.01</dc:type>
    <dcterms:rights>数据使用者必须与数据提供者签订数据使用共享协议，使用后必须注明数据来源。
</dcterms:rights>
    <csdb:dqInfo rdf:resource="http://semweb.csdb.cn/csdb/resource/dqinfo/31022770"/>
    <prism:creationDate rdf:datatype="http://www.w3.org/2001/XMLSchema#date">2009-08-
08</prism:creationDate>
    <dcterms:rightsHolder rdf:resource="http://semweb.csdb.cn/csdb/resource/contact/30941721"/>
    <prism:keyword>泥石流;降水</prism:keyword>
    <dc:description>本数据集收录中国科学院东川泥石流观测研究站在云南蒋家沟观测到的降水资料。降水资料
包括5个观测站点的长期观测资料。</dc:description>
    <dc:creator>中国科学院水利部成都山地灾害与环境研究所</dc:creator>
    <csdb:sharePolicy>本数据保密期5年，解密后免费使用。</csdb:sharePolicy>
    <dc:coverage rdf:resource="http://semweb.csdb.cn/csdb/resource/coverage/17137914"/>
    <csdb:purpose>本数据集的降水观测资料可以与蒋家沟泥石流暴发资料配合使用，是研究泥石流形成和泥石
流预报的珍贵资料。</csdb:purpose>
</csdb:Database>
```

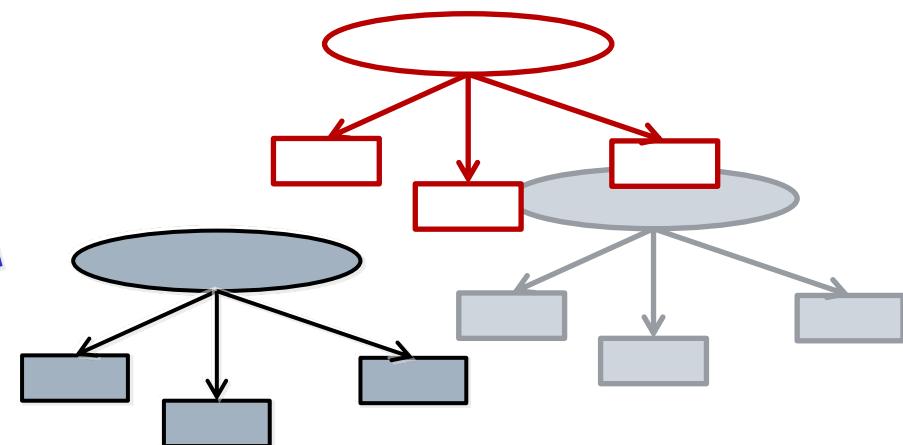
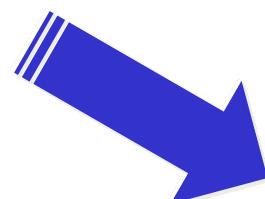
Publishing scientific data

- We distinguish data items into 2 classes
 - Data records:
 - records stored in databases
 - structured
 - Data files
 - file content and its metadata
 - File content is unstructured
 - Metadata is semi-structured

Publishing scientific data

- D2R mapper: Publishing data records
 - A table is mapped to a dataset
 - A record is mapped to an RDF resource
 - Each field is mapped to an RDF property
 - Field values are mapped to RDF property values

	r_regionkey integer	r_name character(25)	r_comment character varving(152)
1	0	AFRICA	lar deposits. blithely final packages ca:
2	1	AMERICA	hs use ironic, even requests. s
3	2	ASIA	ges. thinly even pinto beans ca
4	3	EUROPE	ly final courts cajole furiously final e:
5	4	MIDDLE EAST	uickly special accounts cajole carefully



Publishing scientific data

- F2R mapper: Publishing data files
 - File Content: binary-stream over HTTP (Non-RDF)
 - File Metadata
 - Physical information of a file
 - filename, size, creation time
 - Auto metadata extraction from scientific data files
 - FITS
 - HDF4
 - JPG
 - NetCDF
 - PowerPoint
 - Visio
 - Word

Publishing scientific data

```
<rdf:RDF
  xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
  xmlns:j_0="http://www.semanticdesktop.org/ontologies/2007/05/10/nexif#"
  xmlns:j_1="http://www.csdb.cn/ns/"
  xmlns:rdfs="http://www.w3.org/2000/01/rdf-schema#" >
<rdf:Description rdf:nodeID="AO">
  <j_1:filepath>/test/图片1.JPG</j_1:filepath>
  <j_0:make>NIKON CORPORATION</j_0:make>
  <j_0:exposureBiasValue>5/3</j_0:exposureBiasValue>
  <rdf:type rdf:resource="http://www.semanticdesktop.org/ontologies/2007/05/10/nexif#Photo"/>
  <j_0:width>1680</j_0:width>
  <j_1:filename>图片1.JPG</j_1:filename>
  <rdfs:belongsTo>default</rdfs:belongsTo>
  <j_0:flashpixVersion>48 49 48 48</j_0:flashpixVersion>
  <rdfs:type>file</rdfs:type>
  <j_0:exposureProgram>1</j_0:exposureProgram>
  <j_1:filelastmodified>1317004200000</j_1:filelastmodified>
  <j_0:flash>0</j_0:flash>
  <j_0:height>2520</j_0:height>
  <j_1:filesize>2456762</j_1:filesize>
  <j_0:exposureMode>1</j_0:exposureMode>
  <j_1:filetype>file</j_1:filetype>
  <j_0:exposureTime>1/80</j_0:exposureTime>
</rdf:Description>
</rdf:RDF>
```

extracted metadata of a FITS file

```
Problems @ Javadoc SVN 资源库 Search Console Tasks Debug History SVN 属性 SVN
<terminated> FileMetaDataExtractionJobTest [JUnit] C:\Genuitec\Common\binary\com.sun.java.jdk.win32.x86_1.6.0.013\bin
<rdf:RDF
  xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
  xmlns:j_0="http://www.csdb.cn/ns/"
  xmlns:j_1="http://www.semanticdesktop.org/ontologies/2007/03/22/nco#"
  xmlns:rdfs="http://www.w3.org/2000/01/rdf-schema#" >
<rdf:Description rdf:nodeID="AO">
  <j_1:CRVAL1>53.122751</j_1:CRVAL1>
  <j_1:PI-COI>C. Cesarsky</j_1:PI-COI>
  <j_0:filelastmodified>1317005340000</j_0:filelastmodified>
  <j_1:TEXPTIME>18000</j_1:TEXPTIME>
  <j_1:RA>53.161696</j_1:RA>
  <j_1:EXTEND>true</j_1:EXTEND>
  <j_0:filename>GOODS_ISAAC_03_H_V1.5.fits</j_0:filename>
  <j_1:MJDEND>53397.158255</j_1:MJDEND>
  <j_1:PHOTSYS>AB</j_1:PHOTSYS>
  <j_1:SEEINERR>0.04</j_1:SEEINERR>
  <j_1:NAXIS>2</j_1:NAXIS>
  <j_1:ORIGIN>GOODS</j_1:ORIGIN>
  <j_1:VORELDATE>30 Sep 2005</j_1:VORELDATE>
  <j_1:VERSION>1.5.</j_1:VERSION>
  <j_1:GAIN>80999.999</j_1:GAIN>
  <j_1:CRPIX1>1396.5</j_1:CRPIX1>
  <j_1:CTYPE1>RA---TAN</j_1:CTYPE1>
  <j_1:OBSEND>2005-01-27T03:47:53</j_1:OBSEND>
  <j_1:ASTRERR>0.1</j_1:ASTRERR>
  <j_1:DEC>-27.706750</j_1:DEC>
  <j_1:ASTRCAT>GSC2</j_1:ASTRCAT>
  <j_1:FILTER>H</j_1:FILTER>
  <j_1:PHOTZPER>0.04</j_1:PHOTZPER>
  <j_1:PHOTZP>26.00</j_1:PHOTZP>
  <j_1:OBJECT>GOODS_03</j_1:OBJECT>
  <j_1:PROCSOFT>ESO/MVM</j_1:PROCSOFT>
  <rdfs:belongsTo>default</rdfs:belongsTo>
```

extracted metadata of a JPEG file

Publishing scientific data

- In SDB project, we embed D2R and F2R mappers in VisualDB
- VisualDB is a widely used tool developed by CNIC for helping data owners to manage and publish scientific data on the Web

The screenshot displays the VisualDB platform, which integrates data modeling, storage management, and publishing capabilities.

左侧工具栏 (左侧工具): 包含“建模工具”、“数据库”、“存储位置管理”、“pdc”、“school”、“testLocalFileRepository”、“逻辑建模”和“导航设置”。

中心工作区 (中心区域): 显示了多个数据表的结构视图，通过拖拽操作建立表与表之间的关系。显示的表包括：

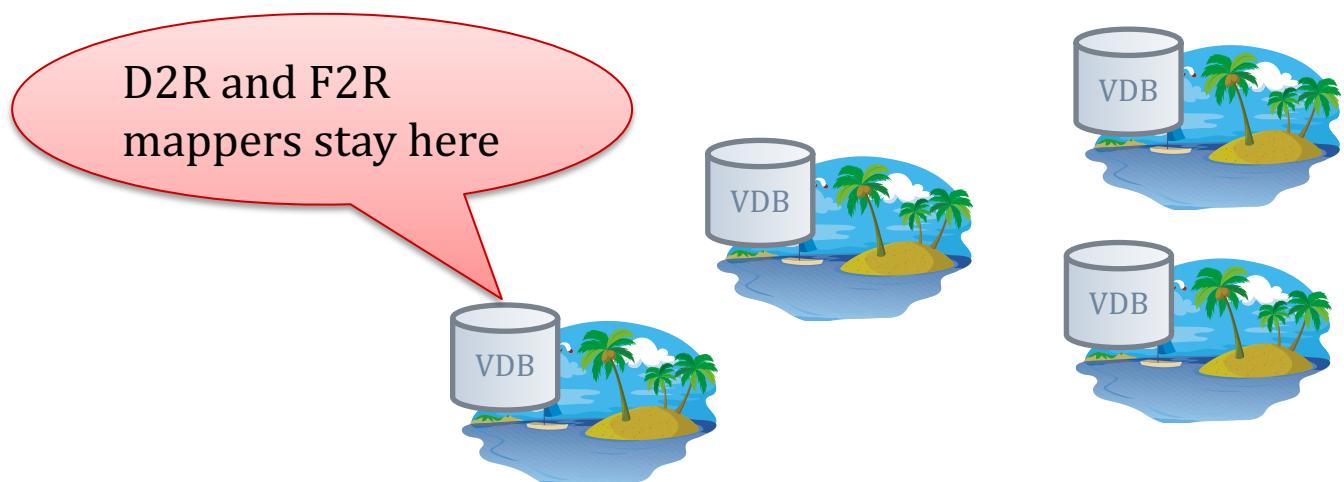
- contact_phone_number (联系人电话号码)
- contact_basic (联系人基础信息)
- users (用户)
- student (学生)
- course (课程)
- conference (会议)
- publisher_salary (出版商薪资)
- publisher (出版商)
- news (新闻)
- 作业元数据 (Assignment Metadata)
- 学生 (Student) 表：包含 编号、姓名、出生年月、班级、性别、年龄、个人照片、选课情况。
- 班级 (Class) 表：包含 编号、名称、描述。
- 作业元数据 (Assignment Metadata) 表：包含 编号、描述。

右侧功能区 (右侧区域): 包含“批量描述”、“旧版xsl升级”、“导出”、“收藏”、“定义实体”、“修改实体”等按钮。下方是“VisualDB DataForge”窗口，显示了“指标管理”下的“新闻”、“新闻类型”、“超链接”、“数据应用”以及“全球经济监测项目”下的“指标基本信息”、“指标数据整合信息”、“指标数据类型”、“指标组”、“季度指标”、“数据来源”和“地区”。右侧还显示了一个“从表信息”的子窗口。

底部状态栏 (底部): 显示了当前操作的记录数（共4283条记录）和一个上下文菜单，包含“编辑数据”、“关联信息”、“复制”、“粘贴”、“删除”、“导入导出”、“按列分组统计”、“随机列”、“显示例”等选项。

Publishing scientific data

- VisualDB has been deployed in more than 30 institutes
 - 4.1 billion records
 - 26 million files
 - are published as Linked Data



Background

Publishing scientific data

Linking scientific data

- General methods & frameworks
- ARIF in SDB
- ARIF & voovle

Summary

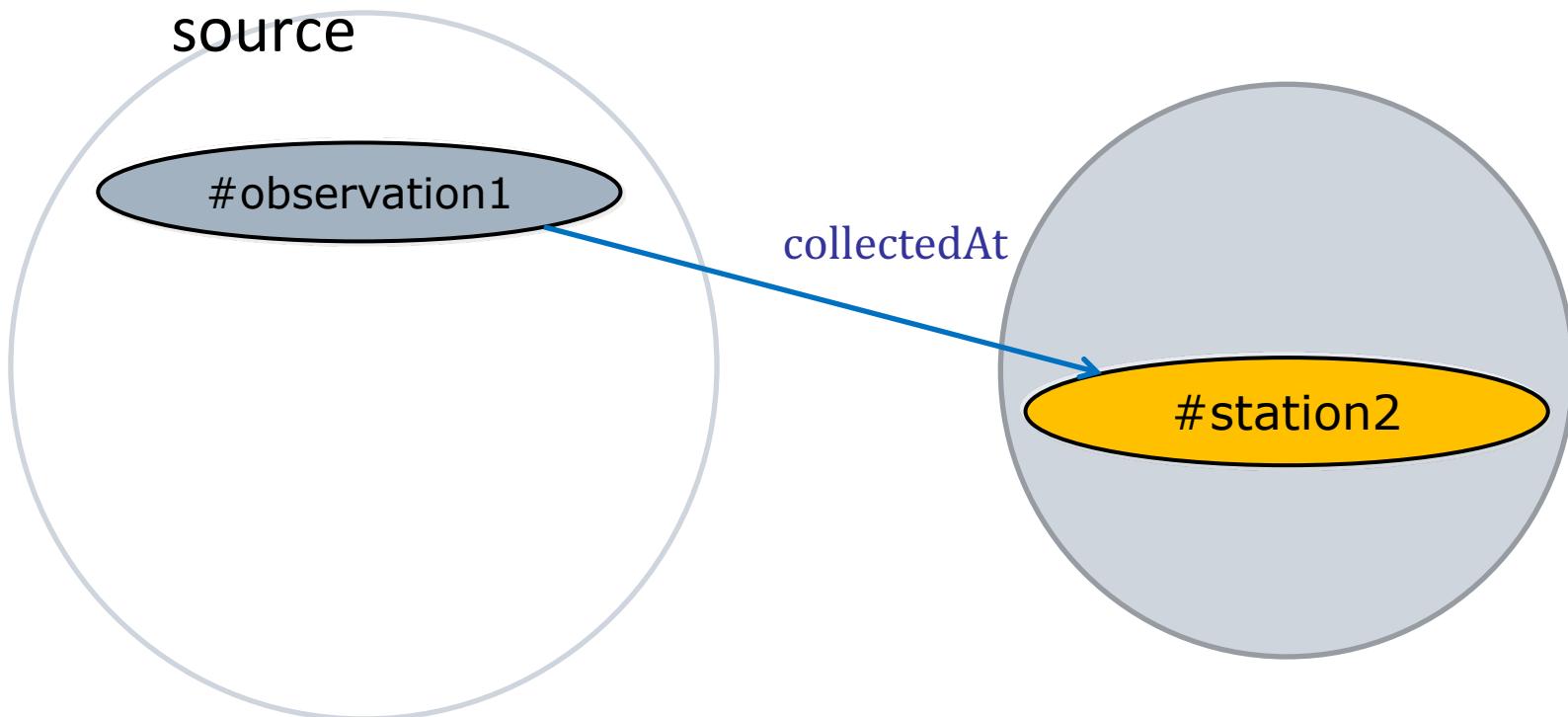
Next Work

Linking scientific data

□ Goal

- **technically speaking**, linking scientific data means set **RDF Links** among data resources, which perhaps are in the same data source or not in the same data

source



Linking scientific data

□ General methods for generating links

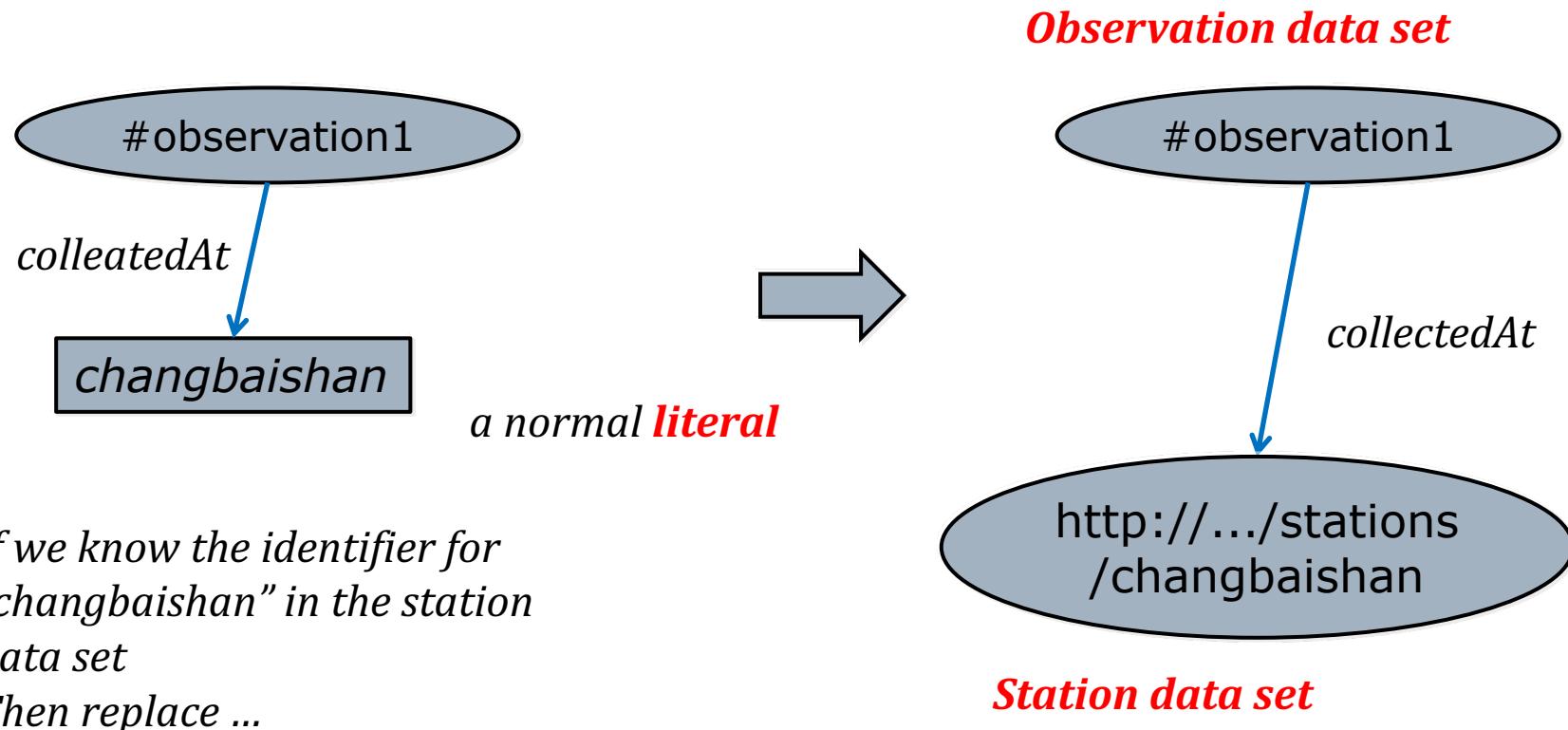
- Setting RDF Links Manually
- Auto-generating RDF Links
 - 1. **identification schemata based**
 - 2. **similarity of entities based**

1. Bizer, Christian, Tom Heath, and Tim Berners-Lee. "Linked data-the story so far." *International Journal on Semantic Web and Information Systems (IJSWIS)* 5.3 (2009): 1-22.

Linking scientific data

□ Method 1: identification schemata based method

- If the source and the target data sets already both support one identification schema, the implicit relationship between entities in both data sets can easily be made explicit as RDF links.



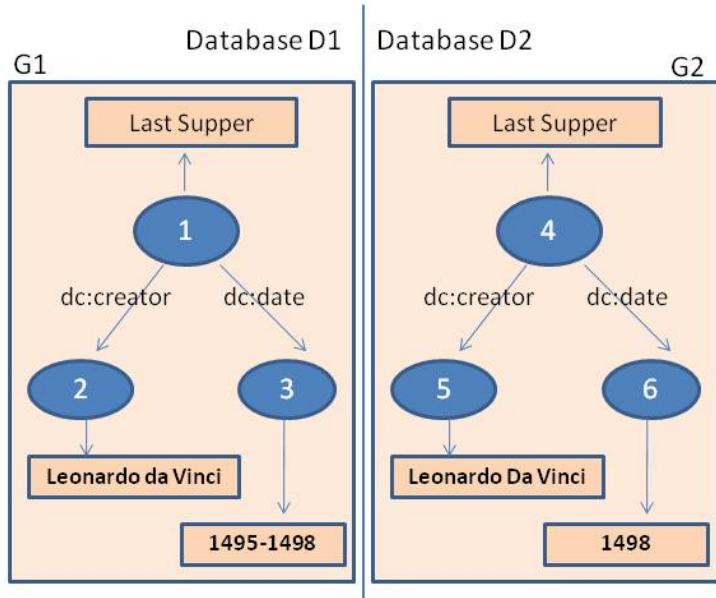
Linking scientific data

- Method 1: identification schemata based method
 - Simple and popular
 - in the publication domain there are ISBN and ISSN numbers
 - in life science, various accepted identification schemata exist for genes, molecules, and chemical substances

Linking scientific data

- **Method 2: similarity of entities based method**
 - complex
 - related work
 - Database community: record linkage and duplicate
 - Semantic Web Community: ontology matching

Linking scientific data



similarity computing

nodes	Sim(n_i, n_j)	
n1	n4	1
n2	n5	0.94
n3	n6	0.44

Generating links between two graphs:

1. retrieve all nodes of two graphs;
2. compute similarity of nodes of shared property one by one;
3. Make a combination of all $\text{sim}(\text{node})$, get $\text{sim}(g_1, g_2)$;
4. If $\text{sim}(g_1, g_2) > T$, then create links;

$$\begin{aligned}
 \text{sim}(g_1, g_2) \\
 = & \text{sim}(n_1, n_4) + \text{sim}(n_2, n_5) + \text{sim}(n_3, n_6) / 3 \\
 = & 0.79
 \end{aligned}$$

Linking scientific data

□ Link Discovery Frameworks

- SILK
 - a link discovery framework using a declarative language for searching relationships between various datasets.
- LIMES
 - LIMES implements novel time-efficient approaches for link discovery in metric spaces.
- RDF-AI
 - RDF-AI focuses on the integration of RDF datasets. Providing modules for pre-processing, matching, fusing, interlinking and postprocessing
- LinQuer
 - LinQuer is a system for generating SQL queries for semantic link discovery over relational data.
- LDIF
 - Intergrates tools including LDSpider、R2R、SILK and Sieve to help application developers with these tasks.

Linking scientific data

□ *SILK as an example*

- *Silk LSL*
 - Flexible, declarative language for specifying linkage rules
- *Silk Link Discovery Engine*
 - responsible for **loading** the instances from the data sources as well as **generating** the links based on the user-provided Link Specifications.
- *Silk Workbench*
 - a web application which guides the user through the process of interlinking different data sources

Save Export as Silk-LS Help

Property Paths

Source: sider

Restriction: ?a rdfs:type sider:drugs .

(custom path)

?a/rdfs:label

?a/sider:sideEffect

?a/rdfs:seeAlso

?a/owl:sameAs

Targeted drugbank

Restriction: ?b rdfs:type drugbank:drugs .

(custom path)

?b/rdfs:label

?b/drugbank:synonym

?b/drugbank:brandName

?b/drugbank:catreet

Transformations

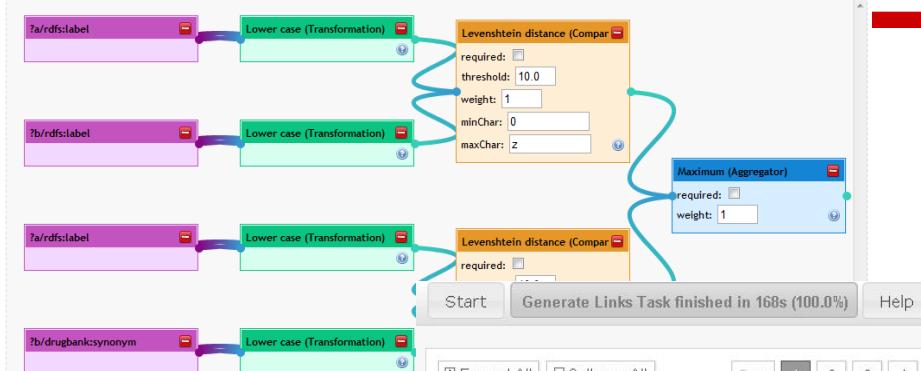
Alpha reduce

Concatenate

Convert Charset

Logarithm

Lower case

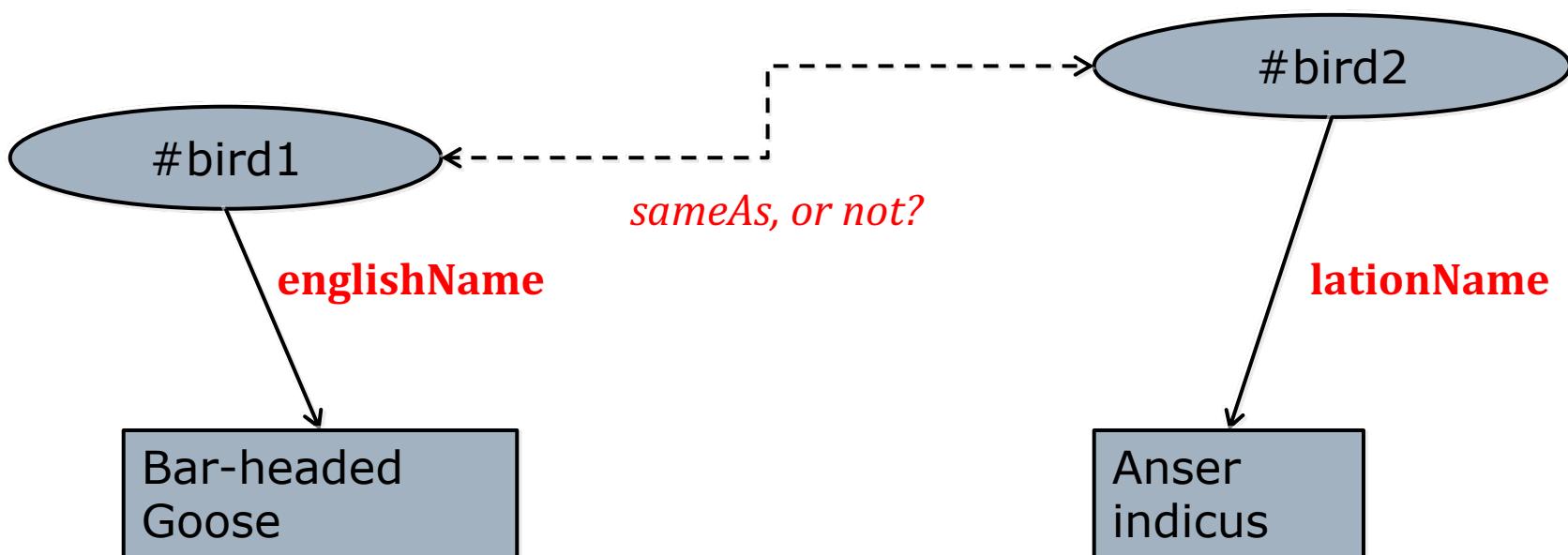


		Source: <input type="text" value="dongwu"/>	Target: <input type="text" value="dongwu"/>	Score	Correct?
				100.0%	
▶	Comparison: jaroWinkler (unnamed_8)	<div style="background-color: #e0f2e0; padding: 2px;">100.0%</div>			
	Input: ?a/<http://purl.org/dc/elements/1.1#title> (unnamed_1)		长尾贼鸥		
	Input: ?b/<http://purl.org/dc/elements/1.1#title> (unnamed_2)		长尾贼鸥		
▶	Comparison: jaroWinkler (unnamed_8)	<div style="background-color: #e0f2e0; padding: 2px;">100.0%</div>			
	Input: ?a/<http://purl.org/dc/elements/1.1#title> (unnamed_1)		日本树莺		
	Input: ?b/<http://purl.org/dc/elements/1.1#title> (unnamed_2)		日本树莺		
▶	Comparison: jaroWinkler (unnamed_8)	<div style="background-color: #e0f2e0; padding: 2px;">100.0%</div>			
	Input: ?a/<http://purl.org/dc/elements/1.1#title> (unnamed_1)		环纹噪鹛		
	Input: ?b/<http://purl.org/dc/elements/1.1#title> (unnamed_2)		环纹噪鹛		
▶	Comparison: jaroWinkler (unnamed_8)	<div style="background-color: #e0f2e0; padding: 2px;">100.0%</div>			
	Input: ?a/<http://purl.org/dc/elements/1.1#title> (unnamed_1)		光滑隐赫杜父鱼		
	Input: ?b/<http://purl.org/dc/elements/1.1#title> (unnamed_2)		光滑隐赫杜父鱼		
▶	Comparison: jaroWinkler (unnamed_8)	<div style="background-color: #e0f2e0; padding: 2px;">100.0%</div>			
	Input: ?a/<http://purl.org/dc/elements/1.1#title> (unnamed_1)		灰头麦鸡		
	Input: ?b/<http://purl.org/dc/elements/1.1#title> (unnamed_2)		灰头麦鸡		
▶	specieslist/specieslist/4d025fe1-e017-4dc6-bd43-a87b4f79b7ad		:oology.csdb.cn/wod/resource/VertebrataCode/code/020490046	<div style="background-color: #e0f2e0; padding: 2px;">100.0%</div>	
▶	specieslist/specieslist/155fb28c-ce05-4bff-8339-48cf2b98237a		:oology.csdb.cn/wod/resource/VertebrataCode/code/022110082	<div style="background-color: #e0f2e0; padding: 2px;">100.0%</div>	
▶	specieslist/specieslist/99361c64-c414-4b06-8ce5-acf9933aac0		:oology.csdb.cn/wod/resource/VertebrataCode/code/052993871	<div style="background-color: #e0f2e0; padding: 2px;">100.0%</div>	
▶	specieslist/specieslist/ac6a3382-c906-4fad-8b2d-b62262d6d912		:oology.csdb.cn/wod/resource/VertebrataCode/code/021430001	<div style="background-color: #e0f2e0; padding: 2px;">100.0%</div>	
▶	specieslist/specieslist/73e0e382-beab-47a4-953d-b4513a98bf08		:oology.csdb.cn/wod/resource/VertebrataCode/code/040210009	<div style="background-color: #e0f2e0; padding: 2px;">100.0%</div>	
▶	specieslist/specieslist/dff57hfr-713f-4275-h36d-f33h775n0321		:oology.csdb.cn/wod/resource/VertebrataCode/code/052690038	<div style="background-color: #e0f2e0; padding: 2px;">100.0%</div>	

Silk Workbench editor and output

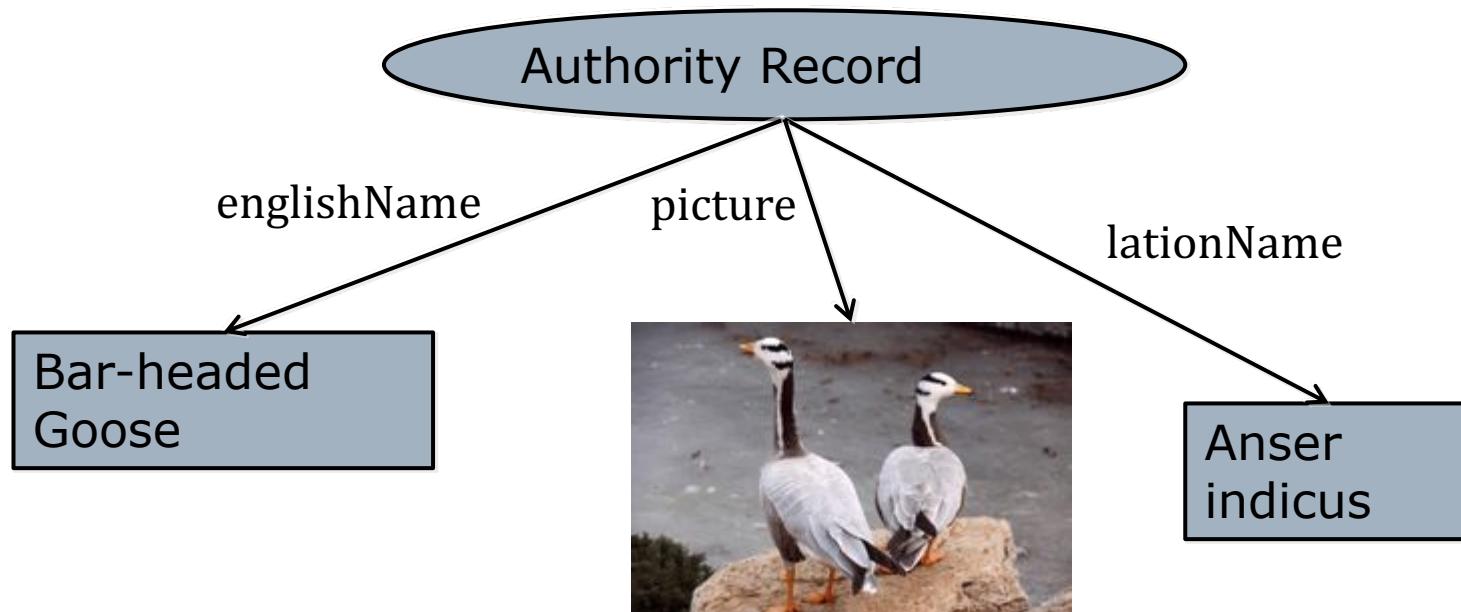
Linking scientific data

- Why these frameworks are not enough for SDB?
- Reason#1
 - Current methods compare two data entities by their **shared** properties. However, it is not common that two distributed data entities have shared properties



Linking scientific data

- An **authority record** with complete properties can help to find meaningful links



- **LORD**: A linked open dataset which contains authority records is called LORD (Linked Open Reference Database)
- Some existing LOD datasets (GeoNames) can be used as LORD
- Frameworks like SILK do not support this method

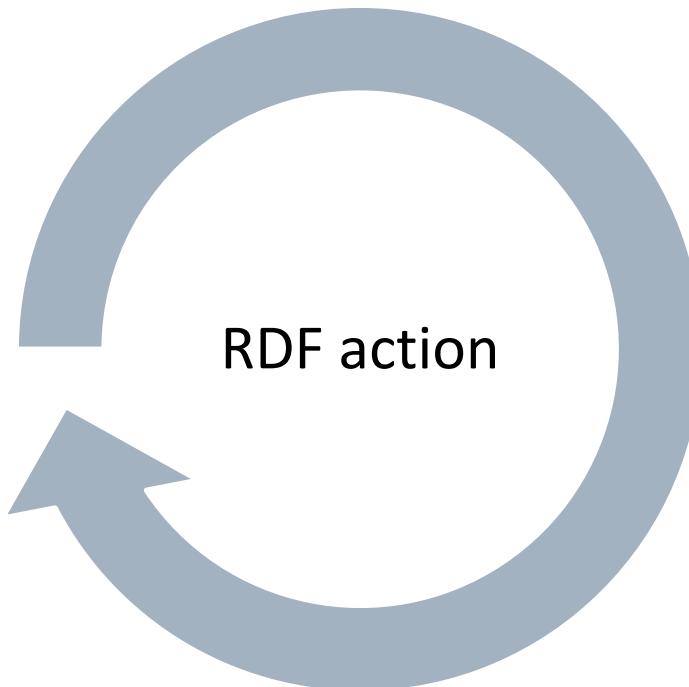
Linking scientific data

- Why these frameworks are not enough for SDB?
- Reason#2
 - The link discovery process is really a long and complex process
 - multiple tasks, paths, steps
 - The process should be able to be planned and a good pipeline mechanism is required
- All frameworks aim at one-step-matching, users have to do some additional work between two steps
 - for example, copy or merge RDF resources

Linking scientific data

- We developed ARIF (**A**nother **R**esource **I**nterlinking **F**ramework)

- All tasks are represented as RDF actions
- ARIF defines 5 kinds of RDF actions

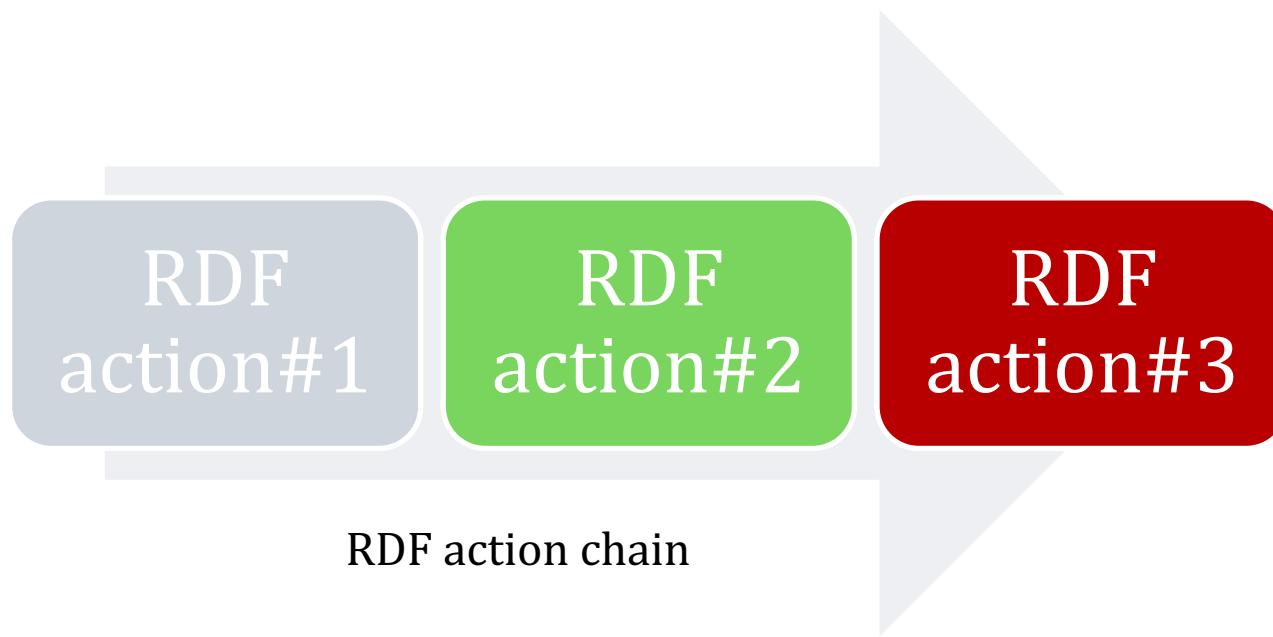


- *RDF dumping*
- *RDF construction*
- *RDF matching*
- *RDF identification*
- *RDF substitution*

LORD based
identification

Linking scientific data

- ARIF (Another Resource Interlinking Framework)
 - RDF actions can be chained with I/O parameters
 - RADL : An RDF action can be represented in RADL language
 - ARIFEngine: Task loader and runner



Linking scientific data

```
<#step1> a :ConstructionStep;
    rdfs:comment "从分类文本中提取分类代码";
    #输入数据源
    :input
    (
        [
            a :QuerySelection;
            :from <#rdfdb>;
            :where """
                ?database
                dc:subject ?subjectText.
                ?database a
                csdb:Database.
                """
            ]
            <#dispatcher>
        );
    #重构规则设置
    :construction
    [
        :output <#extractedSubjects>;
        :template """
            ?extracted a skos:Concept.
            ?extracted dc:identifier ?id.
            ?extracted skos:prefLabel ?title.
        """;
    ];
}
```

An example of a construction task written in RADL

Searching scientific data

- In SDB, ARIF is implemented and deployed as a part of Voovle
- Voovle is a search engine for SDB
- Voovle offers
 - Keywords-based query service
 - SPARQL (RDF Query) query service

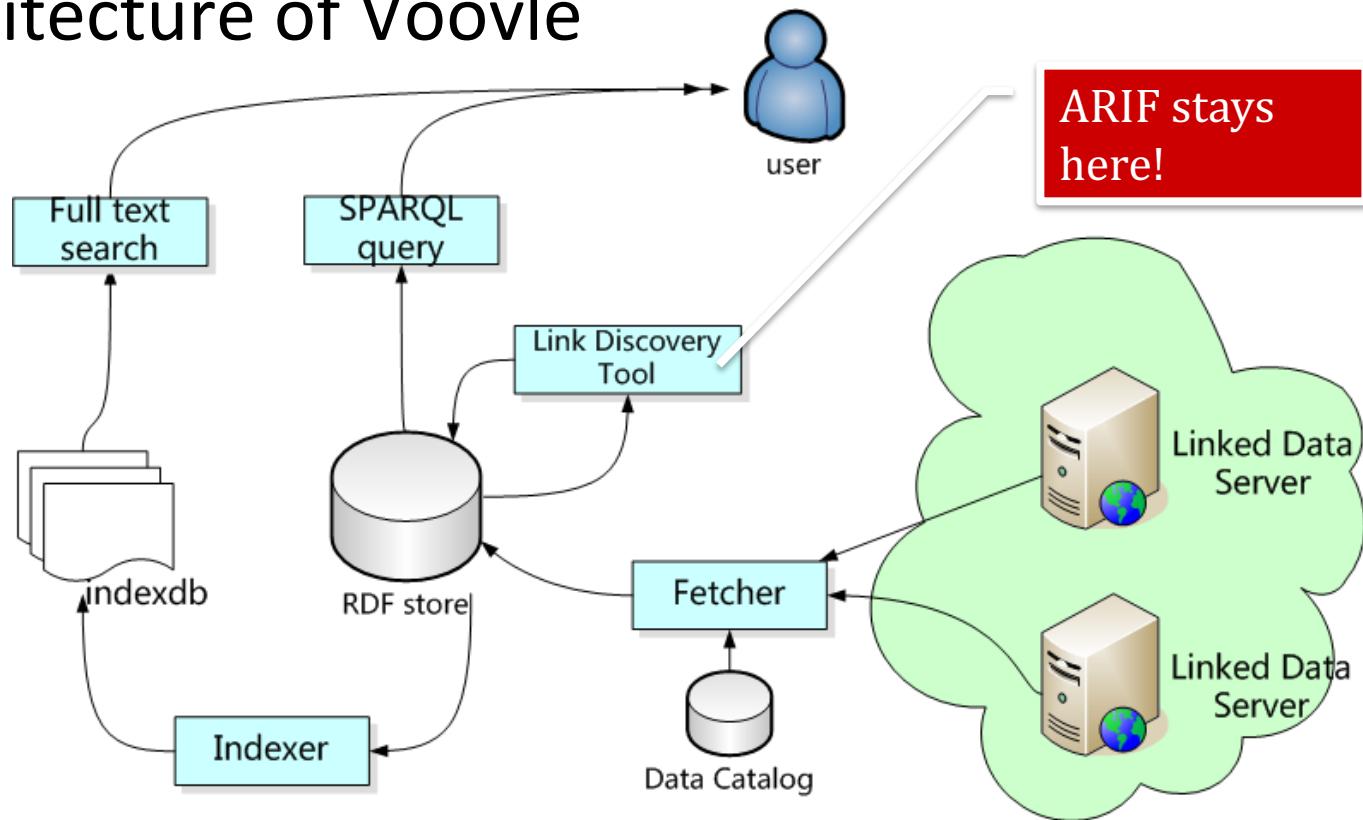


The image displays three separate screenshots of the Voovle search results page for the query '斑头雁' (Bar-headed Goose). Each screenshot shows a list of search results with detailed information about each entry, including titles, descriptions, and URLs. The results are categorized under '资源标题: 斑头雁' (Resource Title: Bar-headed Goose).

- Resource Title: 斑头雁**
 - 【中文科名 : 鸭科】 【中文属名 : 雁属】 【标致 : 斑头雁】 【中文名 : 斑头雁】 【学种名 : 斑头雁】 【中文旧名 : 斑形态种名 : 斑头雁】 【拉丁属名 : Anser】 【拉丁种名 : indicus】 【俗名 : (Latman)】 【拉丁科名 : Anatidae】 【俗名 : 斑头雁】 【物种名 : Anser indicus】 【拉丁名 : (Latman)】
 - 数据来源 [中国动物物种数据库-中国动物物种分类数据库](#) 数据集ID: 2003070019
 - 元数据收集时间 2011-01-22 09:29:30
- 资源标题: 斑头雁**
 - 【中文科名 : 鸭科】 【中文属名 : 雁属】 【标致 : 斑头雁】 【中文名 : 斑头雁】 【学种名 : 斑形态种名 : 斑头雁】 【拉丁属名 : Anser】 【拉丁种名 : indicus】 【俗名 : (Latman)】 【拉丁科名 : Anatidae】 【物种代码 : 0200370019】 【物种旧名 : ANSERIFORMES】 【俗名 : 斑头雁】
 - 数据来源 [中国动物物种数据库-中国脊椎动物分类数据库](#) 数据集ID: 0200370019
 - 元数据收集时间 2011-01-22 07:46:00
- 资源标题: 斑头雁**
 - 【物种识别码 : 10212300000000000000000000000000】 【物种背景生物多样性数据】 【俗名代码 : 】 【中文科名 : 鸭科】 【中文属名 : 斑头雁】 【学种名 : 斑形态种名 : 斑头雁】 【拉丁属名 : Anser】 【拉丁种名 : indicus】 【俗名 : 斑头雁】 【物种旧名 : 】 【物种代码 : 】 【物种旧名 : 斑头雁】
 - 数据来源 [中国动物物种数据库-中国脊椎动物分类数据库](#) 数据集ID: 10212300000000000000000000000000
 - 元数据收集时间 2011-01-20 11:13:12

Searching scientific data

□ Architecture of Voodle



- Voodle collects all RDF data from distributed data sources into a whole RDF store.

Searching scientific data

- In SDB, ARIF has
 - discovered 1552 links between two plant collection databases
 - discovered 4934 links between two animal databases
 - discovered 244 links between Qinghailake ecological database and animal database
 - Need more experiences...

Searching scientific data

- When displaying a record in Vovole, discovered linked data will be listed

The screenshot shows a search result for an observation station. At the top, there are two tabs: '元数据' (Metadata) and '原始数据' (Raw Data). The '原始数据' tab is selected, displaying a table of data with columns for schema labels and their corresponding values.

schema#label	value
rdf-schema#label	长白山站(CBS)
title	长白山站
tower	http://semweb.csdb.cn/flux/resource/tower/1
code	CBS
geoSouthBoundary	41° 41' 49"
climate	温带大陆性气候，具有显著的中纬度山地气候特征
observationStartTime	2002年8月
geoNorthBoundary	42° 25' 18"
precipitation	713mm
geoEastBoundary	128° 16' 48"
temperature	3.6°C
geoWestBoundary	127° 42' 55"
ecosystems	F
soilType	山地暗棕色森林土
vegetation	以红松为主的红松阔叶混交林
canopyHeight	26m
dominantTrees	主要建群树种有红松、椴树、蒙古栎、水曲柳、色木等

Below the table, a section titled '关联数据' (Linked Data) lists 8 entities:

- [1] [flux:entity]: . [EB/OL] http://semweb.csdb.cn/flux/resource/entity/4. []
- [2] [flux:entity]: . [EB/OL] http://semweb.csdb.cn/flux/resource/entity/8. []
- [3] [flux:entity]: . [EB/OL] http://semweb.csdb.cn/flux/resource/entity/6. []
- [4] [flux:entity]: . [EB/OL] http://semweb.csdb.cn/flux/resource/entity/15. []
- [5] [flux:entity]: . [EB/OL] http://semweb.csdb.cn/flux/resource/entity/14. []
- [6] [flux:entity]: . [EB/OL] http://semweb.csdb.cn/flux/resource/entity/11. []
- [7] [flux:entity]: . [EB/OL] http://semweb.csdb.cn/flux/resource/entity/13. []
- [8] [flux:entity]: . [EB/OL] http://semweb.csdb.cn/flux/resource/entity/1. []

Description of an observation station.

observation towers in the station

Background

Publishing scientific data

Linking scientific data

Summary

Next Work

Summary(1)

- Publishing scientific data
 - We chose RDF and Linked Data as publishing standard
 - We publish datasets, files, and records via VisualDB

Summary(2)

- Linking scientific data
 - We think that the popular link discovery frameworks lack supports of third-party authority databases, and lack pipeline mechanisms.
 - We present ARIF (another resource interlinking framework) and integrate it in VooVle

Next Work

□ New challenges

- The stability of data publishing tool, especially on dealing with a large amount of data
- How to evaluate the effectiveness of link discovery framework?
- Improve ARIF to deal with large sized RDF databases

Thanks for your attention!

SHEN Zhihong

bluejoe@cnic.cn



As of September 2011

