

# 以关联数据发布科学数据

## *Publishing Scientific Data as Linked Data*

沈志宏, 何少鹏, 侯艳飞

Zhihong SHEN, Shaopeng HE, Yanfei HOU

中国科学院计算机网络信息中心

Computer Network Information Center, Chinese Academy of Sciences

2014/03/13



中国科学院  
计算机网络信息中心

Computer Network Information Center,  
Chinese Academy of Sciences

# *Outline*

---

- 背景  
background
- 为什么选择关联数据  
why Linked Data?
- 关键问题与发布方案  
publishing methods and steps
- 发布工具与应用系统  
publishing tool and Linked Data applications
- 总结与展望  
future work on data publication

# *background - Scientific DataBase project*

- data sharing project SDB
- 中国科学院信息化专项 “科学数据库项目 (SDB)”
- A long-term mission funded by CAS(Chinese Academy of Sciences)
- started in 1986
- mission
  - Collected multi-discipline research data ( 收集不同学科领域的科研数据 )
  - promoted data sharing ( 促进数据共享 )
- *data from research, data for research*



<http://www.csdb.cn>

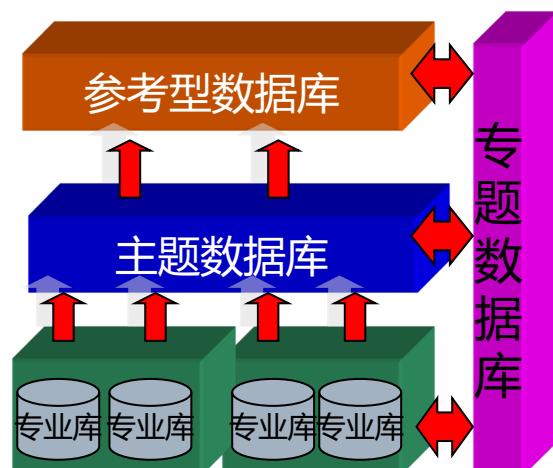
# *background - Scientific DataBase project*

- During the period of the eleventh Five-Year-Plan of CAS (2006-2010)
- “十五”期间
  - about **61** CAS institutes involved
  - **有61个中国科学院研究所参与**
  - Over **200TB** data available for open access and download
  - **超过 200TB 的数据已经开放使用和下载**



# *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*



- 化合物和植物物种2个参考型库
- 化学、材料、空间、天文、遥感、人地系统、动物、微生物等8个主题库
- 聚变、青海湖、冰雪冻土、生态功能区划等4个专题库
- 土壤、海洋、地球化学、指纹、语料、光学、基因组、蛋白组等37个专业库

# *background - Scientific Databases*

- 51 distributed web sites were set up for databases, following a set of construction standards ( 数据规范、建站规范、服务规范 )



<http://www.space.csdb.cn/>

<http://www.chemdb.csdb.cn>

<http://www.cryosphere.csdb.cn>

<http://www.fusion.csdb.cn/>

<http://www.chemcpd.csdb.cn/>

<http://www.data.ac.cn>

<http://www.qinghailake.csdb.cn>

<http://www.lakesci.csdb.cn>

<http://www.ocdb.csdb.cn>

<http://www.soil.csdb.cn>

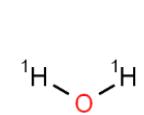
# *background - Scientific Databases*

## □ Examples

$\text{H}_2\text{O}$ 的性质数据  
attributes of  $\text{H}_2\text{O} \rightarrow$

from <http://www.chemdb.csdb.cn>

基本性质



化合物名称：water  
 化学式：H<sub>2</sub>O  
 分子量：17.9994  
 CAS号：7732-18-5  
 SRN：149729862  
 Std.InChIKey : XLYOFNOQVPJJNP-UHFFFAOYSA-N

3d cosmo

热化学数据

基本热化学性质					
性质	数值	单位	温度(K)	状态	
标准焓 $S^\circ_{298.15}$	-241.83	kJ/mol	298.15	g	Chase, 1998
	-241.826	kJ/mol	298.15	g	Cox, Wagman, et al., 1984
	-285.83	kJ/mol	298.15	l	Cox, Wagman, et al., 1984
	-285.83	kJ/mol	298.15	l	Chase, 1998
标准熵 $S^\circ_{298.15}$	69.95	J/(mol*K)	298.15	l	Cox, Wagman, et al., 1984

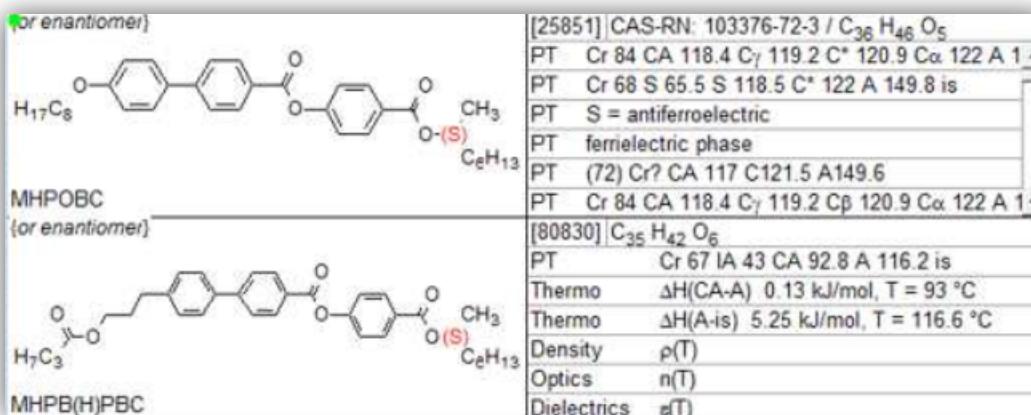
鸟的GPS信息  
GPS locations of a bird →

from <http://qinghailake.csdb.cn>

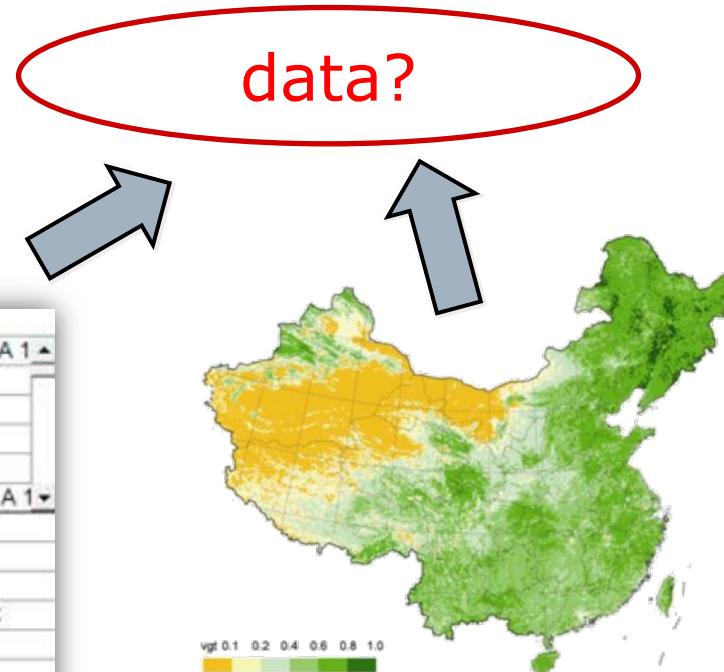
Id	465347
Obs	11
编号	BH07_67580
Ptt	67580
Date	2007-06-23
Record Id	LATEST ARGOS LOCS
Time	15:16:04
Days Ago	652
Hrs Ago	15644.87
Latitude	36.132
Longitude	98.805
Lc94	LB
Nmess	2
Voltage	
Days Dply	89

# *background – problems*

- The big problem is that ...
- Scientific data can be **browsed** easily on the Web, but can **not** be easily **accessed** and **understood** by computer programs
  
- 数据的展示很精彩，但...  
应用程序无法获取到数据！



HTML presentation



HTML presentation

# *background – motivation*

---

- Thus we need a standard method to make distributed scientific data
  - *both Accessible*
  - *and Understandable*
  
- Such a standard method includes
  - *an unified data representation format* (统一的数据描述格式)
  - *an unified open access mechanism* (统一的数据访问机制)

# *Outline*

---

- 背景  
background
- 为什么选择关联数据  
why Linked Data?
- 关键问题与发布方案  
publishing methods and steps
- 发布工具与应用系统  
publishing tool and Linked Data applications
- 总结与展望  
future work on data publication

# *why Linked Data?*

---

- main features of scientific data ( 科学数据的特征 )
  - **variety** of data formats ( 学科领域差异带来的格式多样性 )
    - Storage forms ( 存储格式的多样性 )
      - structured, semi-structured, non-structured
    - Property value types ( 属性值类型的多样性 )
      - text, date, GPS location, images...
  - in-born **interlinked** ( 普遍存在的关联性 )
    - gene and disease ( 基因与疾病 )
    - Materials and chemical compounds ( 材料与化学 )
    - researchers, locations, objects of ecological observation records
    - 生态考察的观测人员、地理位置、对象信息

# *why Linked Data?*

---

- requirements on data representation format ( 科学数据发布格式的要求 )
  1. able to represent data in **different forms and formats** ( 包容格式多样性 )
  2. able to describe the **links** among data ( 能表达关联 )
  3. ensure consistent understanding of **names** ( class names, property names... ) 保证领域内名字含义的一致性
  4. **flexible** schema, easy to extend, adding a new property is very simple ( 无模式 : 属性易扩展 )
  5. makes full use of **existing description tools**, esp. XML ( 对现有描述格式的继承 , 特别是 XML 格式 )
    - CML、MathML、SMILES, ...
    - DarwinCore

# why Linked Data?

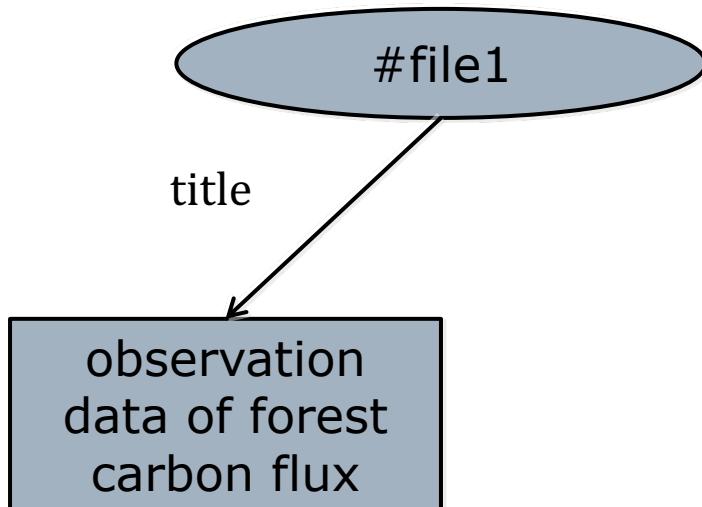
---

- We chose **RDF** as data representation format
  - **RDF: Resource Description Framework**  
<http://www.w3.org/RDF>
  - RDF describes resources in terms of simple **properties** and property **values** ( 属性 : 属性值 )
  - it is based upon the idea of *making statements about resources* (in particular web resources) in the form of **subject-predicate-object** expressions
  - 声明 : 主-谓-宾 三元组
  - e.g.

<u>&lt;#file1&gt;</u>	<u>&lt;#title&gt;</u>	<u>“observation data of forest carbon flux”</u>
subject	predicate	object
resource	property	value

# *why Linked Data?*

- RDF statements about a resource are often represented as a **graph** ( RDF图 : 节点、弧 )
  - **nodes**: representing the resources, and their properties values.
  - **arcs**: representing properties of resources

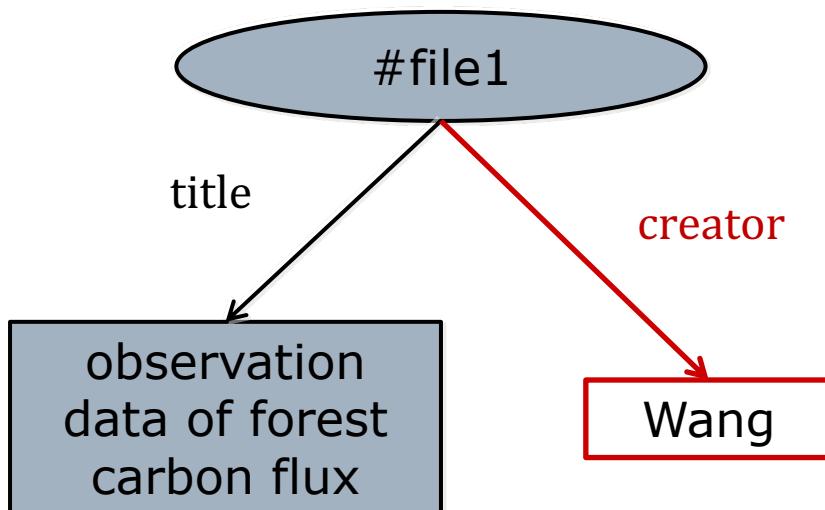


<#file1>    <#title>

“observation data of forest carbon flux”

# *why Linked Data?*

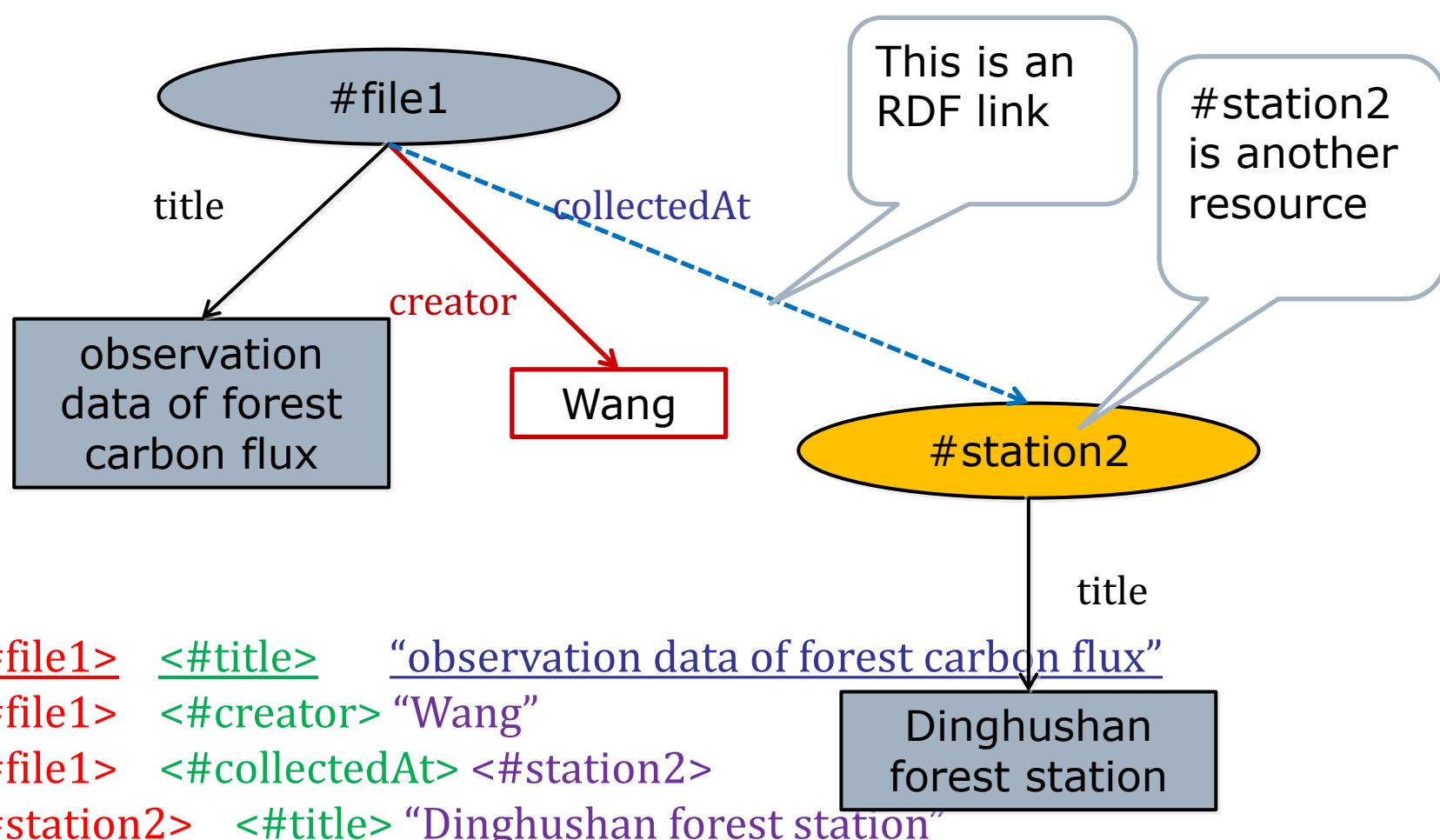
- RDF statements about a resource are often represented as a **graph** ( RDF图 : 节点、弧 )
  - **nodes**: representing the resources, and their properties values.
  - **arcs**: representing properties of resources



[`<#file1> <#title> “observation data of forest carbon flux”`](#)  
[`<#file1> <#creator> “Wang”`](#)

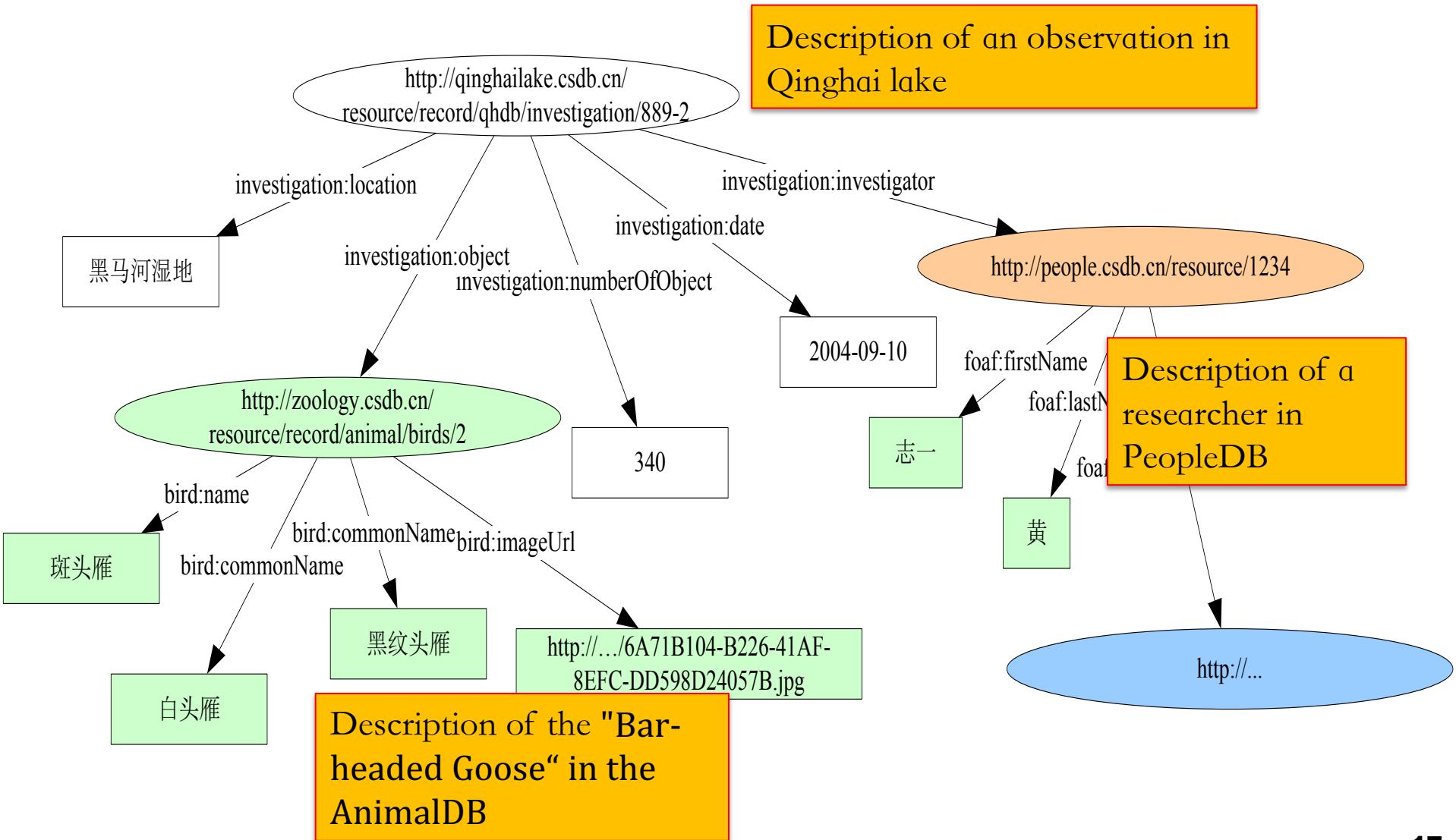
# why Linked Data?

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



# why Linked Data?

## □ More complex RDF data example



# *why Linked Data?*

- After choosing the data format, we chose **Linked Data** as the open access mechanism of data sources
  - Tim Berners-Lee, the father of World Wide Web, coined the term **Linked Data** in 2006[1].
  - 万维网之父Tim Berners-Lee在2006年提出Linked Data
  - 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

# *why Linked Data?*

---

- Four basic principles of Linked Data

Linked Data 的 “四项基本原则”

1. *Use URIs as names for things*

采用URIs 来命名事物

2. *Use HTTP URIs so that people can look up (dereference) those names.*

采用HTTP URIs 以便用户可以查询这些名字

3. *When someone looks up a URI, provide useful information.*

当用户查询URI时，提供有用的信息

4. *Include links to other URIs so that they can discover more things.*

包含到其他资源的链接，以便发现更多的信息

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

# *why Linked Data?*

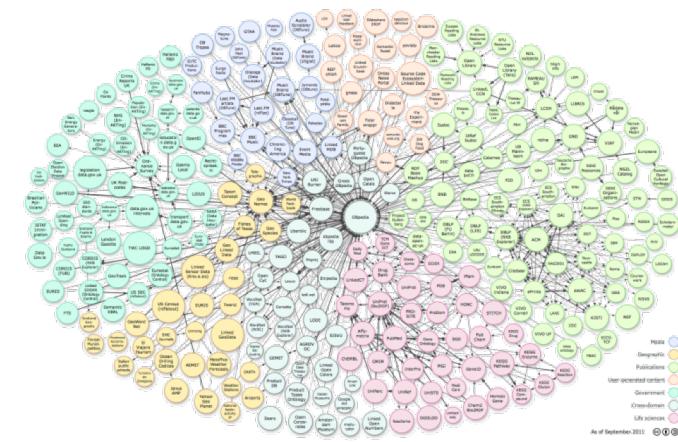
## □ 关联数据

- 完全架构于目前的Web体系之上，这对科学数据库来说几乎意味着零升级成本。
- Linked Data is a standard solution which is totally built on current Web architecture, that means very **little cost** for upgrading SDB project.



Web of Document

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



Web of Data

# *why Linked Data? - scientific data examples*

## □ Linked Life Data

- Searches and explores over RDF statements from various sources including UniProt, PubMed, EntrezGene and 20 more...
- 包括UniProt, PubMed, EntrezGene 等20多个数据源，10亿条记录
- Performs complex SPARQL(RDF query) queries and retrieves more than one billion RDF resources.



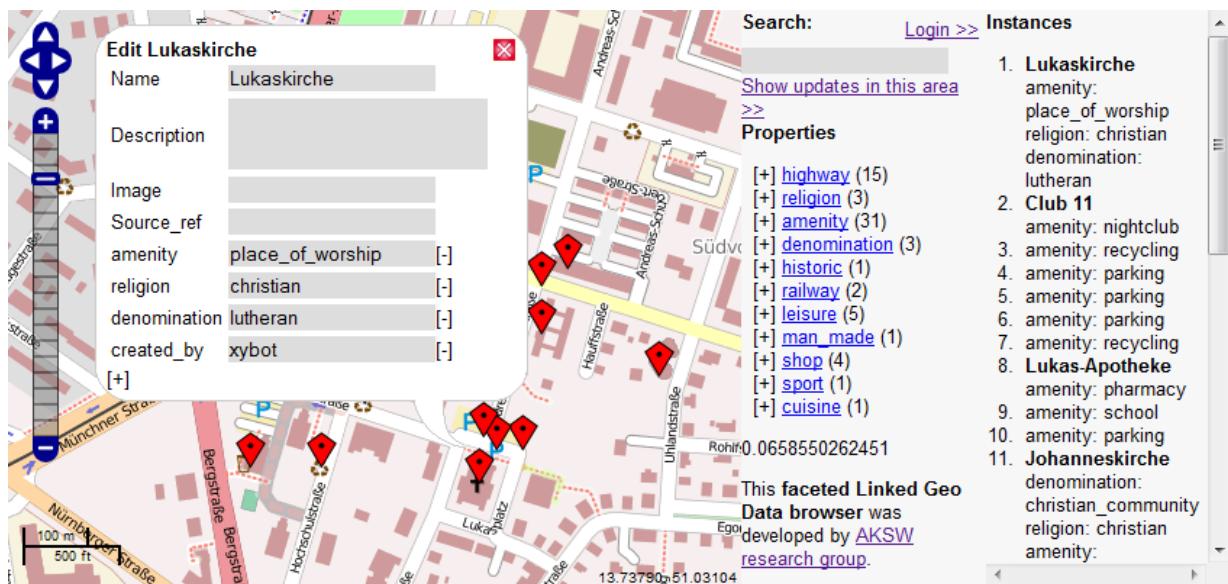
a semantic data integration platform for the biomedical domain

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.

# *why Linked Data? - scientific data examples*

## □ LinkedGeoData

- uses the information collected by the OpenStreetMap project and makes it available as an RDF knowledge base according to the Linked Data principles.
- 采用OpenStreetMap 收集的信息并依据Linked Data原则将其发布为可用的RDF知识集



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.

# *why Linked Data? - scientific data examples*

## □ Diseasesome

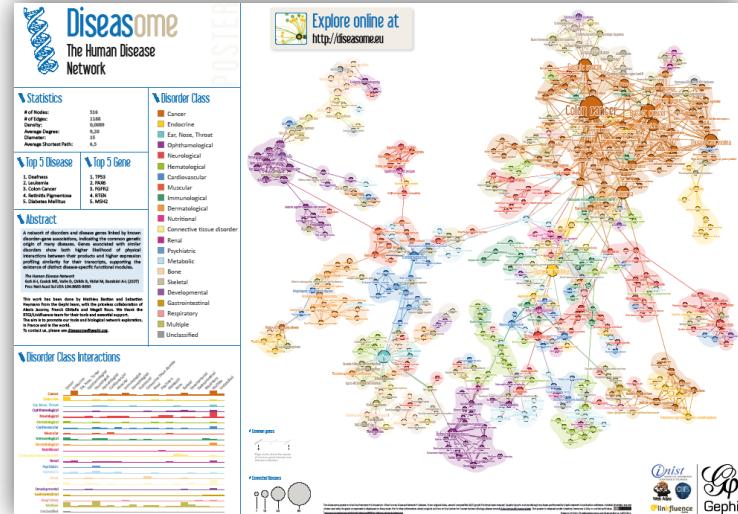
- publishes Linked Data of 4,300 disorders and **disease genes** linked by known disorder-gene associations for exploring all known phenotype and disease gene associations, indicating the common genetic origin of many diseases.
- 发布了4,300疾病和**疾病基因**的关联数据，探索疾病表现和基因之间的关联关系，并绘制血多疾病常见的基因起源

## □ 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.
- 发布关于生物序列、家族、物种、物种出险记录、相关数据、以及到geonames, bio2rdf, dbpedia, freebase, umbel的链接



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](http://wiki.knoesis.org/index.php/SSW_Datasets)

# Outline

---

- 背景  
background
- 为什么选择关联数据  
why Linked Data?
- 关键问题与发布方案  
**publishing methods and steps**
- 发布工具与应用系统  
publishing tool and Linked Data applications
- 总结与展望  
future work on data publication

# *publishing methods and steps*

---

- two issues should be considered in the publishing process

- Issue #1: **data items own little information**

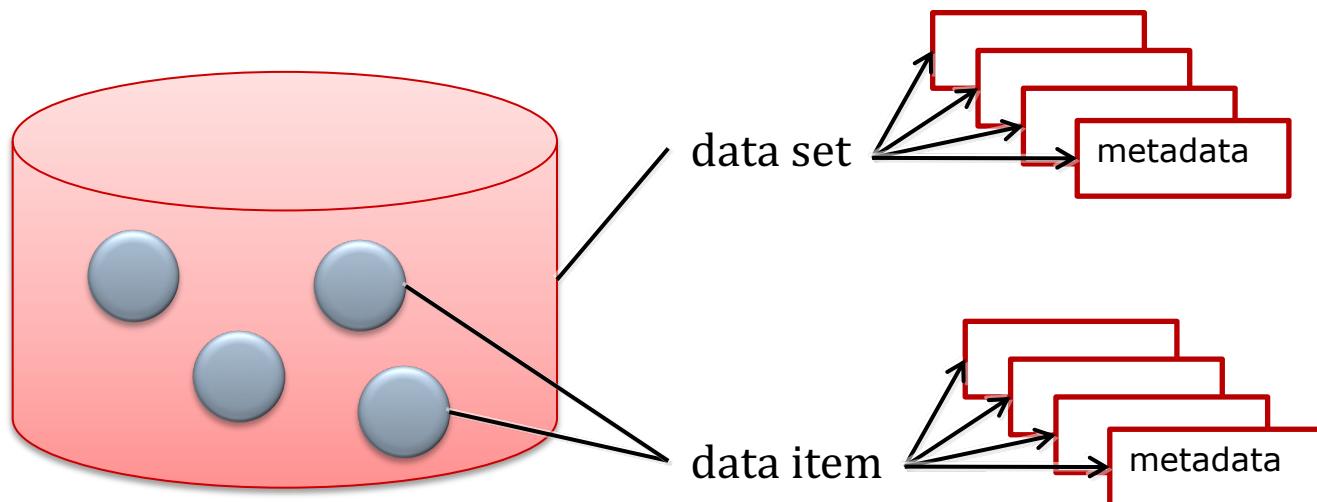
Scientific data contains many files/records which are not easy to be found only by their own property values ( 科学数据的记录非常多，但是单凭它的属性还很难被发现 )

- Issue #2: **RDF data are structured, but most scientific data are unstructured.**

unstructured scientific data (usually exist as data files) are much more than structured scientific data ( 除了结构化数据，科学数据更多的是文件 )

# *publishing methods and steps*

- Issue #1: data items own little information
- We define data in two 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 methods and steps*

---

- Data set 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/>

Elements:

*<dc:title>*  
*<dc:subject>*  
*<dc:source>*  
*<dc:description>*  
*<dc:coverage>*  
*<dc:creator>*  
*<dc:contributor>*  
*<dcterms:accrualPeriodicity>*  
*<dcterms:rightsHolder>*  
*<prism:creationDate>*  
*<prism:keyword>*  
*<dc:type>*  
*<dcterms:rights>*  
*<prism:publicationDate>*  
*<rdfs:label>*  
*<prism:url>*  
*<csdb:dqInfo>*  
*<csdb:sharePolicy>*  
*<csdb:purpose>*

# *publishing methods and steps*

## □ An example of data set 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 methods and steps*

- Issue #2: RDF data are structured, but most scientific data are unstructured
- We distinguish data items into two classes

## Data records (数据记录)

- records usually stored in relational databases
- structured

	A	B	C		
1	College Enrollment 2005 - 2006				
2					
3	Student ID	Last Name	Initial	Age	Program
4	ST348-245	Smith	B.	21	Drafting
5	ST348-246	Wilson	C.	19	Science
6	ST348-247	Thompson	S.	18	Business
7	ST348-248	James	D.	23	Nursing
8	ST348-249	Ramirez	A.	37	Science
9	ST348-250	Graham	T.	20	Arts
10	ST348-251	Rosen	O.	26	Business
11	ST348-252	Hirsch	W.	22	Arts
12	ST348-253	Russell	E.	20	Nursing
13	ST348-254	Robitaille	K.	19	Drafting

Each row in the table is a data record

## Data files (数据文件)

- File content is unstructured
- Metadata is structured



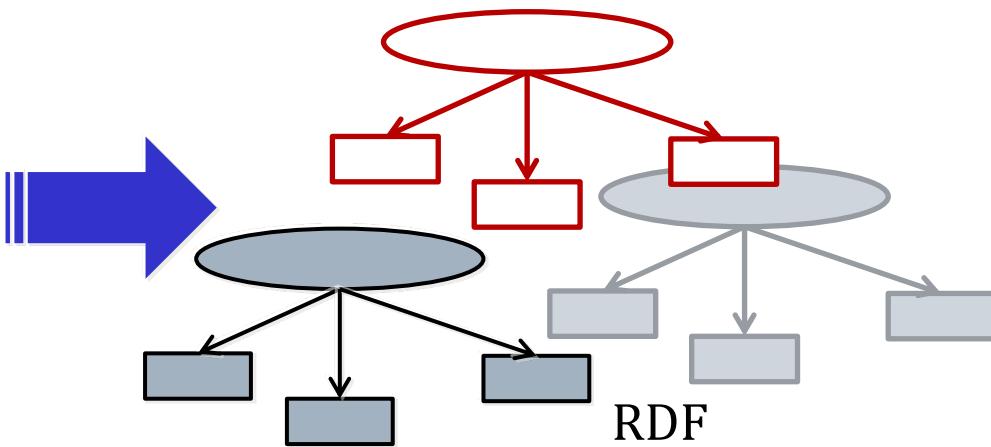
# *publishing methods and steps*

## □ Publishing data records - rules:

1. *each table is mapped to a data set*  
一个表对应一个数据集
2. *each record is mapped to an RDF resource*  
一条记录对应一个RDF资源
3. *each field is mapped to an RDF property*  
每个字段对应RDF的一个属性
4. *field values are mapped to RDF property values*  
字段值对应RDF的属性值
5. *foreign key relationships are mapped to RDF links*  
外键关联映射成RDF链接

id	filepath	owner	visibility	size	description	lastmodified	filename
1	/44444.jpg	0140e4d634d00002 +		165208	(Null)	1379628477599	44444.jpg
4	/QQ图片20130140e4d634d00002 +			173269	(Null)	1379628477631	QQ图片20130140e4d634d00002 +.jpg
5	/333.jpg	0140e4d634d00002 +		282042	(Null)	1379628477642	333.jpg
6	/EAC2.png	0140e4d634d00002 +		295745	(Null)	1379628477650	EAC2.png
7	/2.jpg	0140e4d634d00002 +		320109	(Null)	1379628477660	2.jpg
8	/QQ图片20130140e4d634d00002 +			330207	(Null)	1379628477666	QQ图片20130140e4d634d00002 +.jpg
9	/dm.zip	0140e4d634d00002 +		409135	(Null)	1379628477681	dm.zip
10	/untitled file.txt	0140e4d634d00002 +		19	(Null)	1379628558050	untitled file.txt
11	/public/2.jpg	0140e4d634d00002 +		320109	(Null)	1379643370380	2.jpg
12	/public/333.jpg	0140e4d634d00002 +		282042	(Null)	1379643370480	333.jpg
13	/public/44444.jpg	0140e4d634d00002 +		165208	(Null)	1379643370562	44444.jpg
14	/public/dm.zip	0140e4d634d00002 +		409135	(Null)	1379643370604	dm.zip
15	/public/EAC2.rtf	0140e4d634d00002 +		295745	(Null)	1379643370646	EAC2.rtf
16	/public/QQ图片20130140e4d634d00002 +			330207	(Null)	1379643370683	QQ图片20130140e4d634d00002 +.jpg

data records



# *publishing methods and steps*

---

## □ Publishing data files – rules:

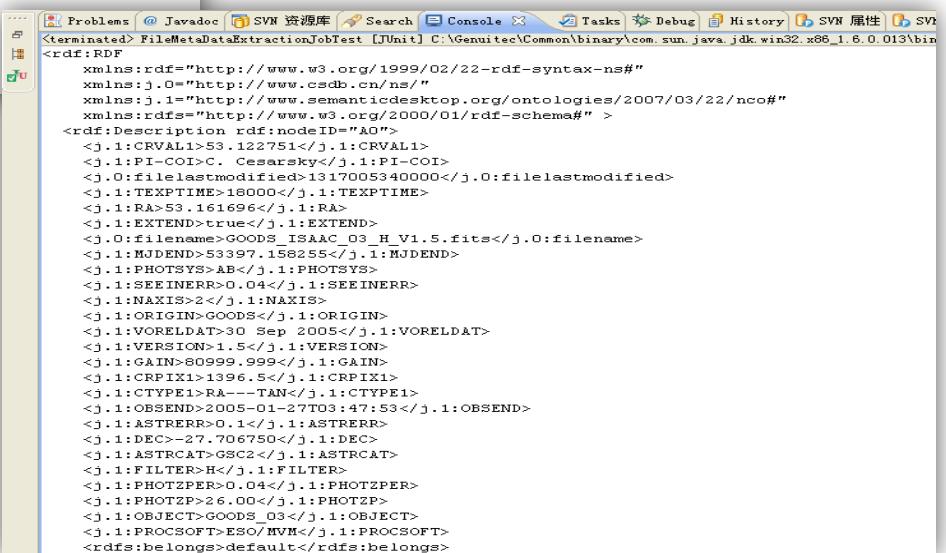
1. *Each file is assigned two URLs, one for file download, the other for file metadata*  
两个URL : 文件下载+文件描述
2. *When users access download URL, provides file contents in binary-stream (Non-RDF)*  
发布文件的二进制内容
3. *When users access file description URL, provides mixed description in RDF format*  
发布文件的描述信息
4. *file metadata includes file attributes and extracted metadata*
  1. *file attributes (文件的属性信息)*
    - *filename, size, creation time*
    - *often retrieved from file systems*
  2. *extracted metadata from content (自动抽取的元数据)*
    - *FITS、HDF4、JPG、NetCDF、PowerPoint、Visio、Word*

# *publishing methods and steps*

## □ Examples of extracted metadata

```
<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:belongs>default</rdfs:belongs>
  <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



The screenshot shows an IDE interface with multiple tabs open. The main tab displays the following RDF code:

```
... Problems @ Javadoc SVN 资源库 Search Console Tasks Debug History SVN 属性 SVN
...
<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>16000</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/HVM</j_1:PROCSOFT>
  <rdfs:belongs>default</rdfs:belongs>
```

extracted metadata of a JPEG file

# *Outline*

---

- 背景  
background
- 为什么选择关联数据  
why Linked Data?
- 关键问题与发布方案  
publishing methods and steps
- 发布工具与应用系统  
publishing tool and Linked Data applications
- 总结与展望  
future work on data publication

# *publishing tool - VisualDB*

---

- 科学数据发布工具VisualDB
- In SDB project, we developed **VisualDB** help data owners to publish scientific data both as HTML pages and Linked Data on the Web
- 完全配置化发布，无需编程
- VisualDB allows users **define publishing rules** in Web interfaces, without programming



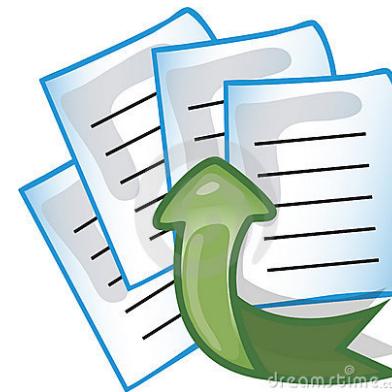
*Visual Scientific Data Management and Publishing Tool*

# *publishing tool - VisualDB*

- In SDB, VisualDB has been deployed in 37 institutes
  - 4.1 billion records ( 41亿条记录 )
  - 26 million files ( 2600万个文件 )
  - are published as Linked Data

	A	B	C		
1	College Enrollment 2005 - 2006				
2					
3	Student ID	Last Name	Initial	Age	Program
4	ST348-245	Smith	B.	21	Drafting
5	ST348-246	Wilson	C.	19	Science
6	ST348-247	Thompson	S.	18	Business
7	ST348-248	James	D.	23	Nursing
8	ST348-249	Ramirez	A.	37	Science
9	ST348-250	Graham	T.	20	Arts
10	ST348-251	Rosen	O.	26	Business
11	ST348-252	Hirsch	W.	22	Arts
12	ST348-253	Russell	E.	20	Nursing
13	ST348-254	Robitaille	K.	19	Drafting

Each row in the  
table is a data  
record



4.1 billion records

26 million files

# *publishing tool - VisualDB*

## □ Sample URIs

*request for  
a html page*

`http://www.zoology.csdb.cn/wod  
/page/VertebrataCode/code/020  
370019`

*request for an  
RDF graph*

`http://www.zoology.csdb.cn/wod  
/data/VertebrataCode/code/020  
370019`

The screenshot shows the VisualDB publishing tool interface. On the left, a URL box contains the URI `http://www.zoology.csdb.cn/wod/page/VertebrataCode/code/020370019`. An arrow points from this URL box to the main content area. The main content area displays a web browser window showing a page from the '中国脊椎动物分类代码数据库' (Chinese Vertebrate Classification Code Database). The page features a navigation menu on the left and a detailed table of species information on the right. Below the table, a large block of RDF code is shown, representing the data from the selected URI. Another arrow points from the URL box containing the RDF URI to this block of RDF code.

http://www.zoology.csdb.cn/page/showItem.vpage

首页 数据检索 . 查找数据库 . 数据服务 . 服务指南 . 服务案例 服务公告 关于本库 .

数据库导航

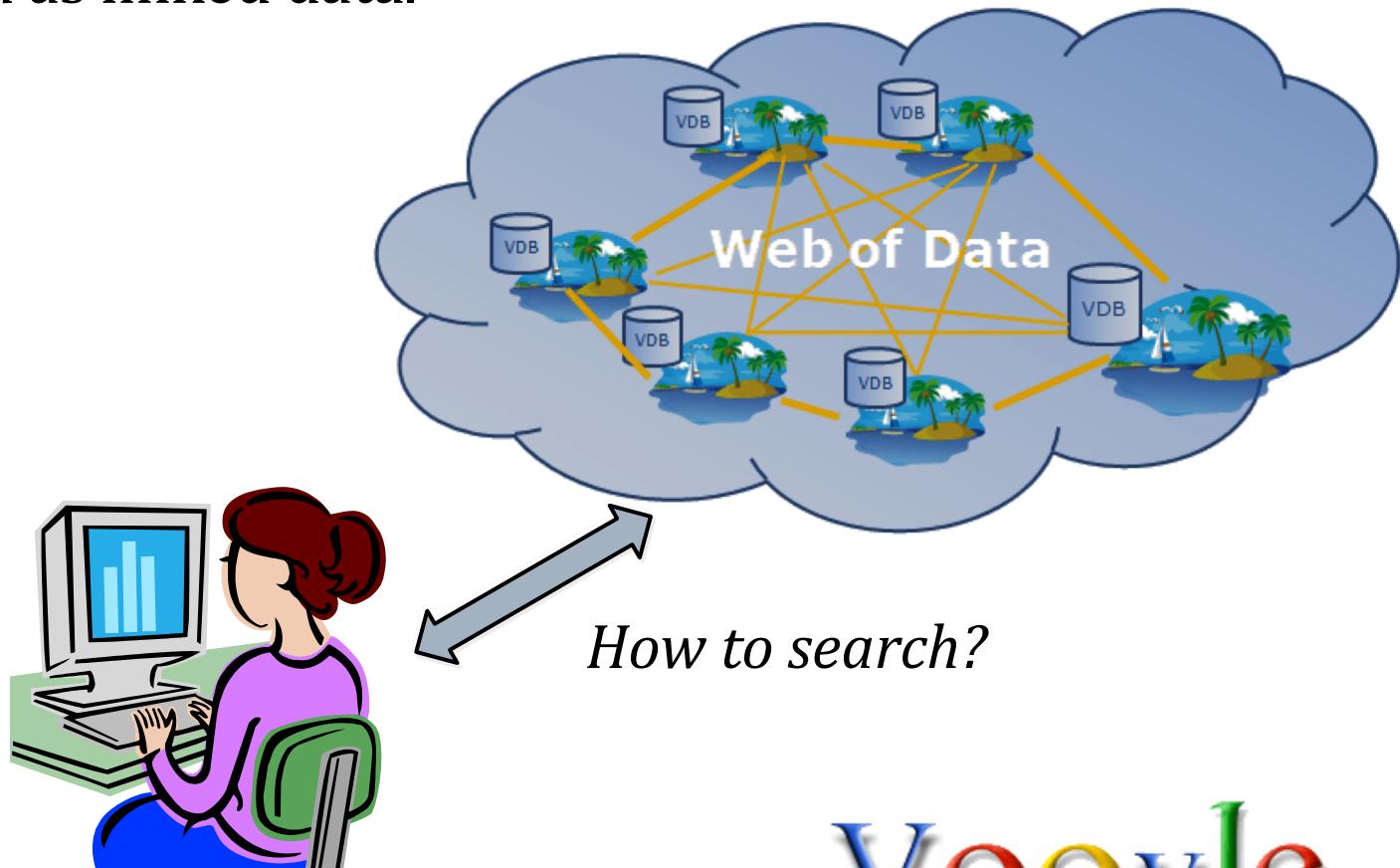
- 中国动物志数据库
- 中国动物图谱数据库
- 中国鸟类数据库
- 中国蜜蜂数据库
- 中国内陆水体鱼类数据库
- 中国两栖爬行动物数据库
- 中国直翅目昆虫数据库
- 中国跳虫目昆虫数据库
- 中国动物物种编目数据库
- 中国漫长类物种文献及地名
- 西南丘陵脊椎动物分布名册

物种代码 020370019  
中文目名 雁形目  
拉丁目名 ANSERIFORMES  
中文科名 鸭科  
拉丁科名 Anatidae  
中文属名 雁属  
拉丁属名 Anser

<rdf:RDF  
xmlns:dct="http://purl.org/dc/elements/1.1#"  
xmlns:code="http://www.zoology.csdb.cn/wod/resource/vocab/VertebrataCode/code/"  
xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"  
xmlns:j:0="http://www.zoology.csdb.cn/wod/resource/vocab/VertebrataCode/"  
xmlns:dataset="http://voovle.csdb.cn/ns/dataset#">  
<rdf:Description  
rdf:about="http://www.zoology.csdb.cn/wod/resource/VertebrataCode/code/020370019">  
<code:family>Anatidae</code:family>  
<code:genus>Anser</code:genus>  
<code:cfamily>鸭科</code:cfamily>  
<code:cgenus>雁属</code:cgenus>  
<code:corder>雁形目</code:corder>  
<code:csubfamily>null</code:csubfamily>  
<dataset:belongsEntity  
rdf:resource="http://www.zoology.csdb.cn/wod/resource/vocab/VertebrataCode/code"/>  
<rdf:type  
rdf:resource="http://www.zoology.csdb.cn/wod/resource/vocab/VertebrataCode/code"/>  
<code:species>indicus</code:species>  
<code:order>ANSERIFORMES</code:order>  
<code:subfamily>null</code:subfamily>  
<code:subspecies>null</code:subspecies>  
<code:author>(Latham, 1790)</code:author>  
<code:code>020370019</code:code>  
<code:cname>斑头雁</code:cname>  
<dataset:belongs  
rdf:resource="http://www.zoology.csdb.cn/wod/resource/vocab/VertebrataCode"/>  
<dc:title>斑头雁</dc:title>  
</rdf:Description>  
</rdf:RDF>

# *Linked Data application - Vooyle*

- 科学数据搜索引擎Vooyle
- We developed Vooyle as a **search engine** for scientific data published as linked data.



*How to search?*

Vooyle

# Linked Data application - Voovle

- Voovle offers free text search by inputting keywords ( 文本检索 )
- Voovle also offers semantic query by submitting a SPARQL query ( 支持 SPARQL查询 )

```
SELECT distinct (?s) WHERE
{
?s ?p '斑头雁'
}
```

No	time is:1400 ms
1	<a href="http://www.zoology.csdb.cn/wod/resource/Animalnames/name/AE8132F8-4FF5-40BC-BCBF-326D15796110">http://www.zoology.csdb.cn/wod/resource/Animalnames/name/AE8132F8-4FF5-40BC-BCBF-326D15796110</a>
2	<a href="http://www.zoology.csdb.cn/wod/resource/cnAtlas/tableTaxa/8A40049F-D2C0-455F-8409-A54F5B419DED">http://www.zoology.csdb.cn/wod/resource/cnAtlas/tableTaxa/8A40049F-D2C0-455F-8409-A54F5B419DED</a>
3	<a href="http://www.zoology.csdb.cn/wod/resource/cnfauna/tableTaxa/1478B56D-6609-4A58-8305-CFB75621AB0A">http://www.zoology.csdb.cn/wod/resource/cnfauna/tableTaxa/1478B56D-6609-4A58-8305-CFB75621AB0A</a>
4	<a href="http://www.zoology.csdb.cn/wod/resource/specieslist/specieslist/27db1a96-b1e4-40b2-bd50-bf5ace063bb1">http://www.zoology.csdb.cn/wod/resource/specieslist/specieslist/27db1a96-b1e4-40b2-bd50-bf5ace063bb1</a>
5	<a href="http://www.zoology.csdb.cn/wod/resource/VertebrataCode/code/020370019">http://www.zoology.csdb.cn/wod/resource/VertebrataCode/code/020370019</a>
6	<a href="http://www.qinghailake.csdb.cn/wod/resource/cn.csdb.qhnew/circlelakeVirussample/1813">http://www.qinghailake.csdb.cn/wod/resource/cn.csdb.qhnew/circlelakeVirussample/1813</a>
7	<a href="http://www.qinghailake.csdb.cn/wod/resource/cn.csdb.qhnew/circlelakeVirussample/1812">http://www.qinghailake.csdb.cn/wod/resource/cn.csdb.qhnew/circlelakeVirussample/1812</a>
8	<a href="http://www.qinghailake.csdb.cn/wod/resource/cn.csdb.qhnew/circlelakeVirussample/1826">http://www.qinghailake.csdb.cn/wod/resource/cn.csdb.qhnew/circlelakeVirussample/1826</a>
9	<a href="http://www.qinghailake.csdb.cn/wod/resource/cn.csdb.qhnew/birds/2">http://www.qinghailake.csdb.cn/wod/resource/cn.csdb.qhnew/birds/2</a>

[1] Apache Lucene - Welcome to Apache Lucene. Online at <http://lucene.apache.org/>

The screenshot shows the Voovle search interface. At the top, there is a search bar with the text '斑头雁' and a '搜索' (Search) button. Below the search bar, the text '检索关键字示例： 斑头雁 棉花 红土' is displayed. The main area shows search results for '斑头雁'. Each result includes a thumbnail image, the URL, and some descriptive text. The results are paginated with '1/1' at the bottom. The footer of the page contains the text '版权所有：中国科学院计算机网络信息中心·科学数据中心'.

# Linked Data application - Vooiple

# Vooiple



search engine for CERN observation data  
<http://cern.vooiple.csdb.cn>

search engine for SDB  
<http://vooiple.csdb.cn>

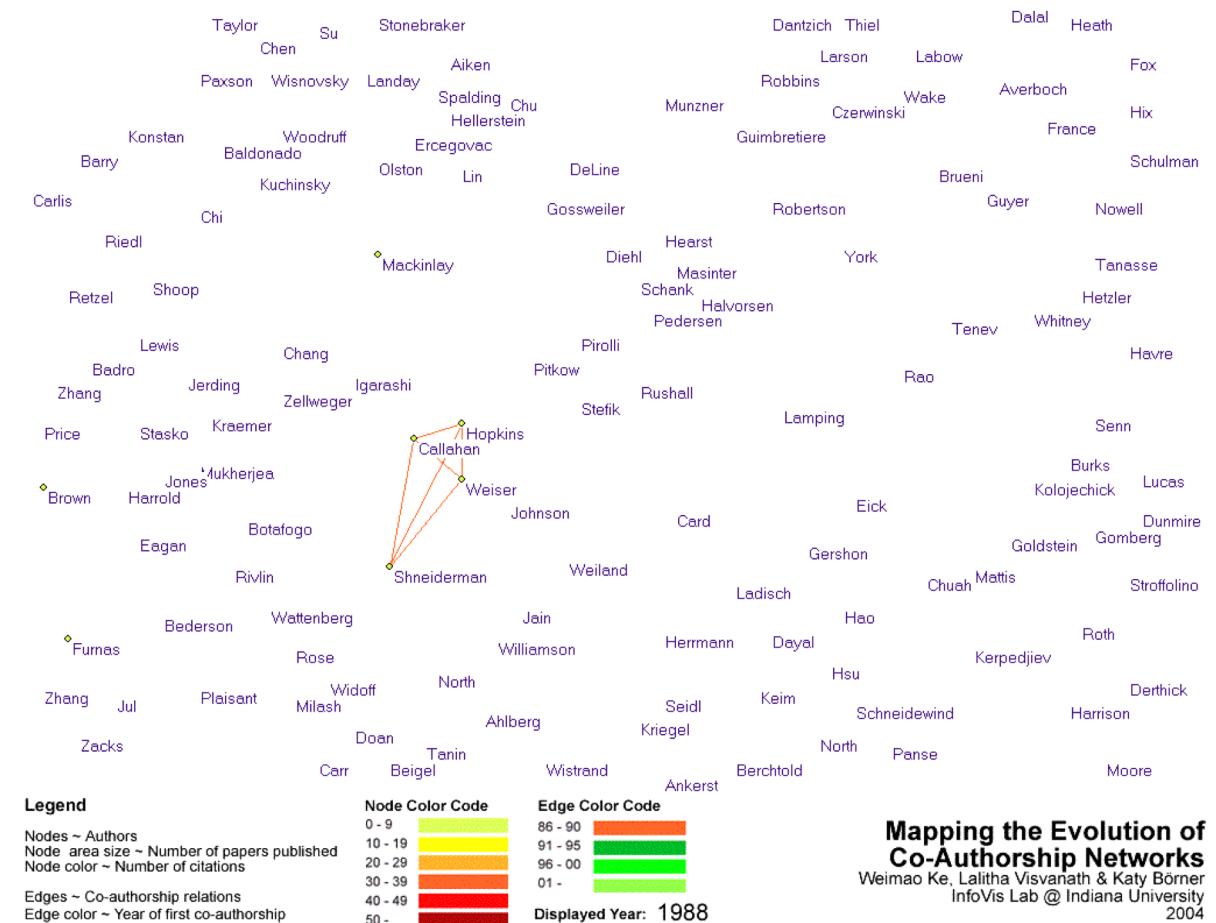
124 datasets, 5.64 million records

search engine for Scientific Data Grid

<http://datagrid.vooiple.csdb.cn>

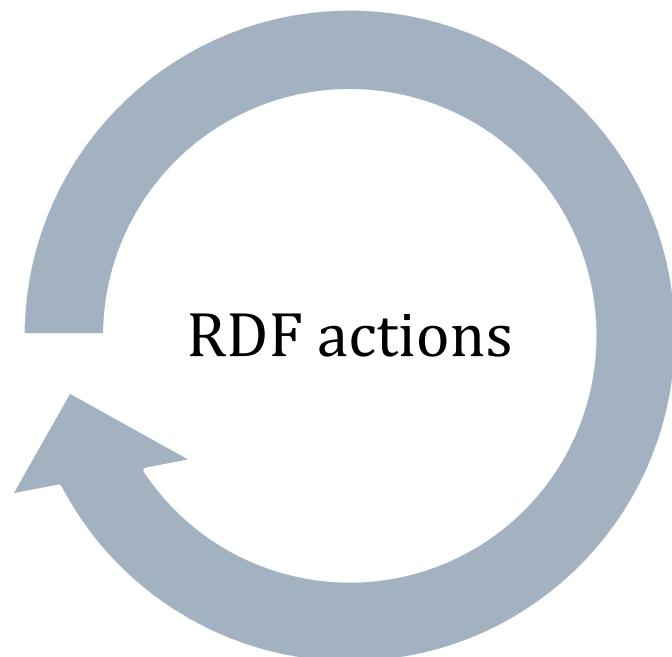
# Linked Data application – link discovery tool

- Upon RDF store built by Voovle, we developed tools to find useful links among scientific data.
- 跨数据集的关联发现



# *Linked Data application – link discovery tool*

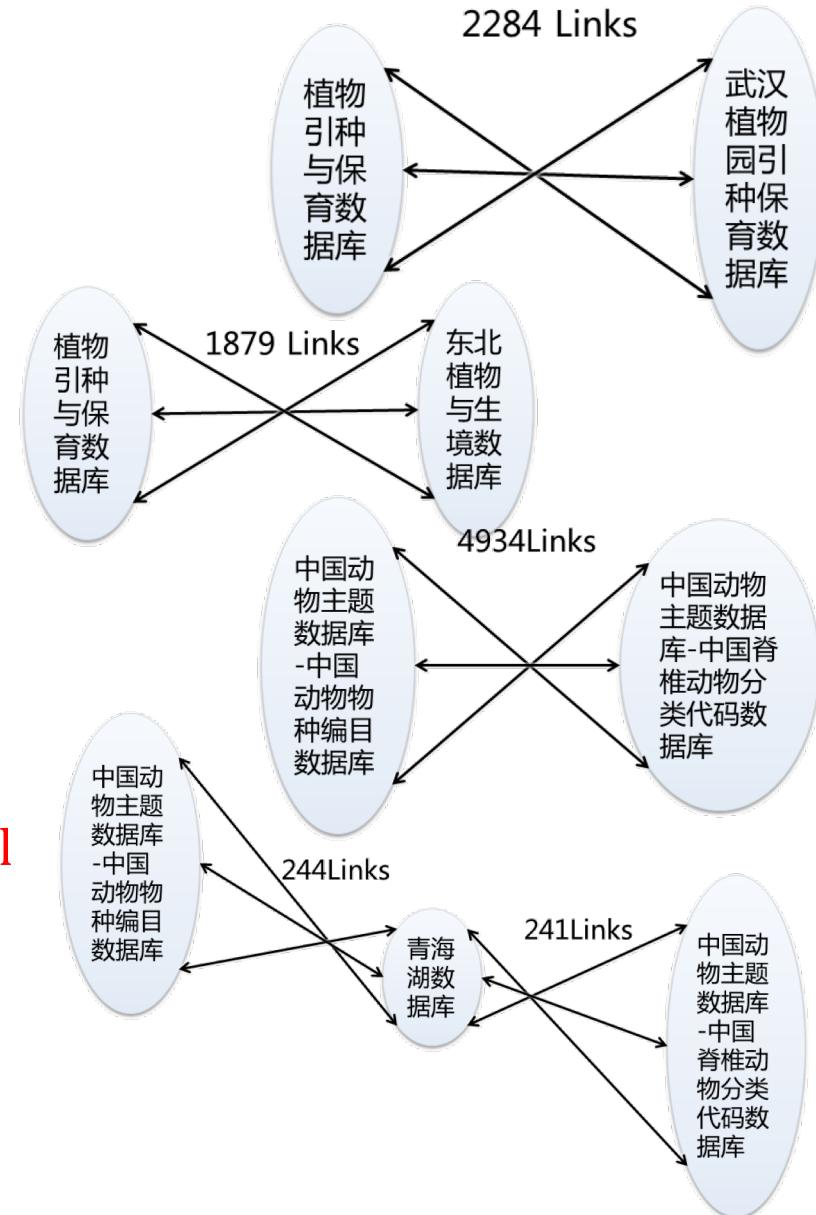
- ARIF is such a tool based on **graph similarity** ( 图形相似度 )
- In ARIF, all RDF graph operations involved in the link discovery process are defined as different “**RDF actions**”.
- ARIF归纳了关联发现过程中的所有RDF图形操作，并定义为RDF action。



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

# *Linked Data application – link discovery tool*

- 部分关联发现效果
- experiments and results with ARIF
  - discovered **2284 links** between two **plant collection databases**
  - discovered **4934 links** between two **animal databases**
  - discovered **244 links** between Qinghai lake **ecological observation database** and **animal database**
  - trying more experiments...



# *Outline*

---

- 背景  
background
- 为什么选择关联数据  
why Linked Data?
- 关键问题与发布方案  
publishing methods and steps
- 发布工具与应用系统  
publishing tool and Linked Data applications
- 总结与展望  
future work on data publication

# *conclusions*

---

- 我们发布科学数据以促进数据共享
  - We published scientific data to **promote data sharing** in Chinese Academy of Sciences
- 我们选择RDF作为科学数据的发布格式
  - We chose **RDF** as data representation format for scientific data
- 我们选择关联数据作为科学数据开放访问机制
  - We chose **Linked Data** as open access mechanism of scientific data

## *new problems*

---

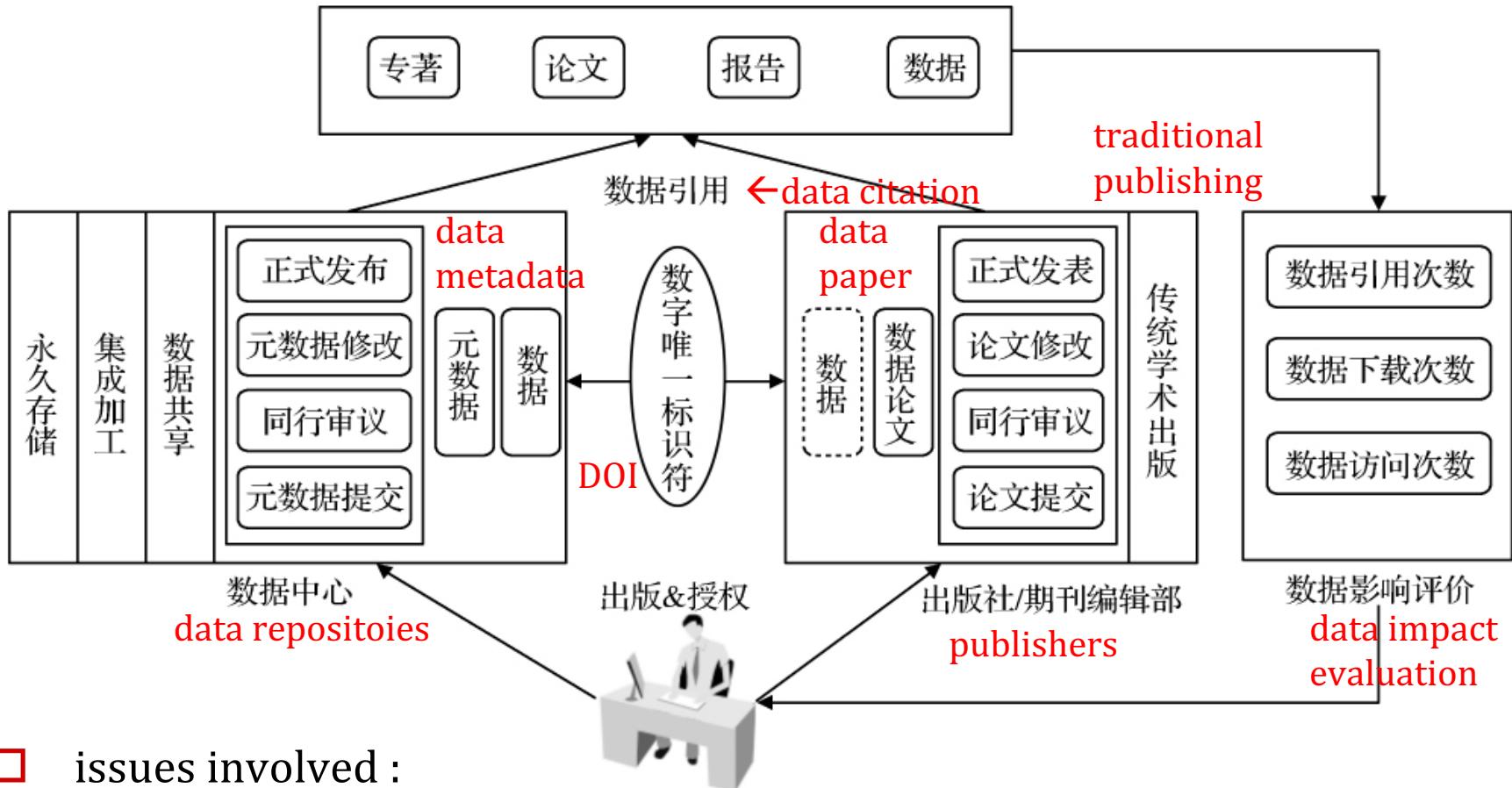
- no guarantee of quality  
数据质量没有保证
- unstable data access URLs  
数据URL不稳定
- unable to evaluate the values of data  
数据的学术价值无法评估

# *possible answers?*

- no guarantee of quality  
数据质量没有保证
  - peer reviews on data?*
- unstable data access URLs  
数据URL不稳定
  - data repositories provide data storage and preservation?*
- unable to evaluate the values of data.  
数据的学术价值无法评估
  - data DOIs enable data citation?*



# *big picture of Data Publishing(or Publication)*



□ issues involved :

- metadata management / data management / identifier / peer review
- long term storage / data citation / data impact evaluation

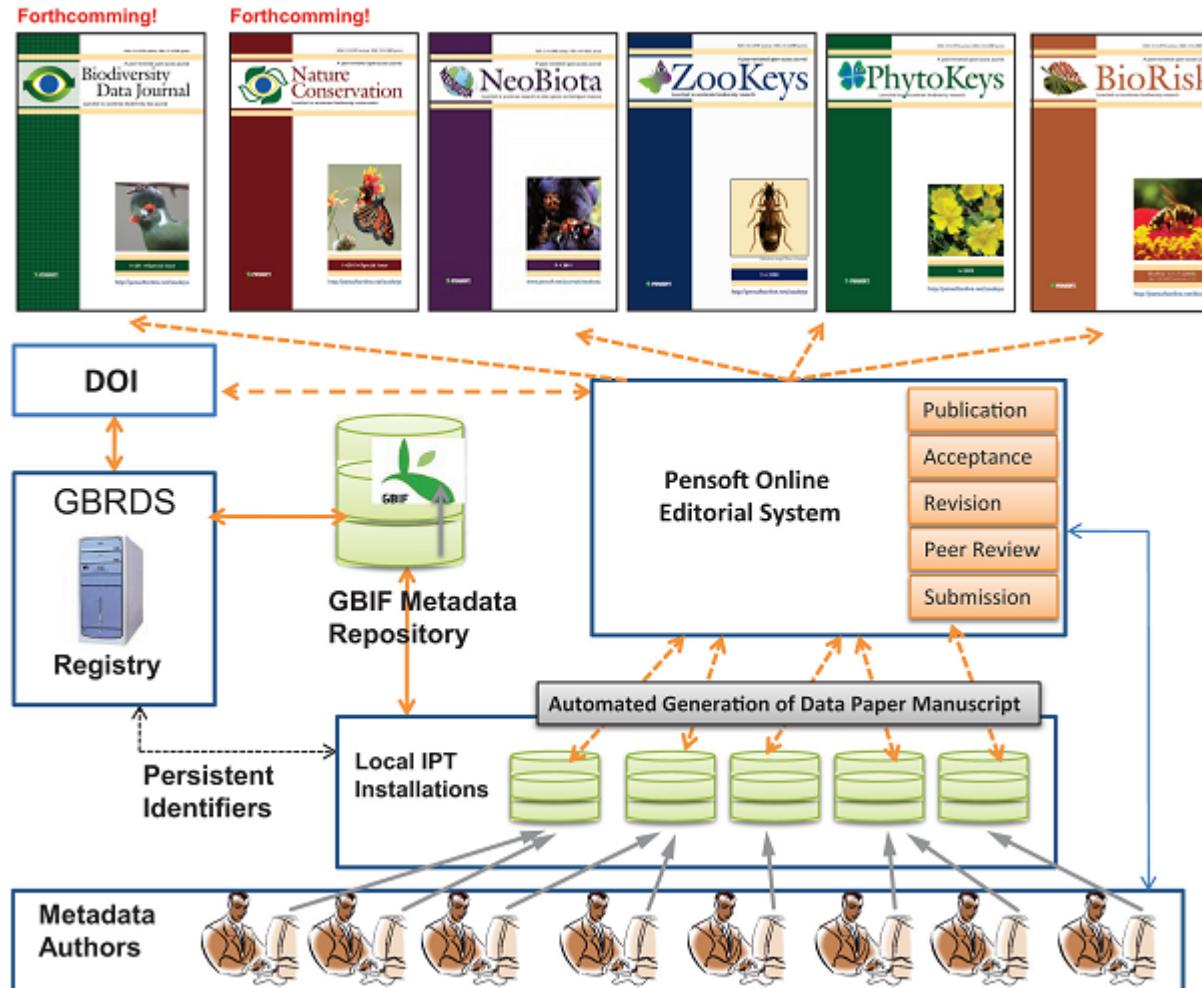
□ roles involved: data repository / publisher / journal / author

1. 吴立宗等, "科学数据出版现状及其体系框架," 遥感技术与应用, 2013

# *similar framework for Data Publishing*

the GBIF/Pensoft workflow of data publishing and automated generation of Data Paper

- important issues involved in data publication:
  - Metadata
  - Identifier
  - markup
  - link
  - Archiving
  - migration
  - Exchange
  - rights
  - ...

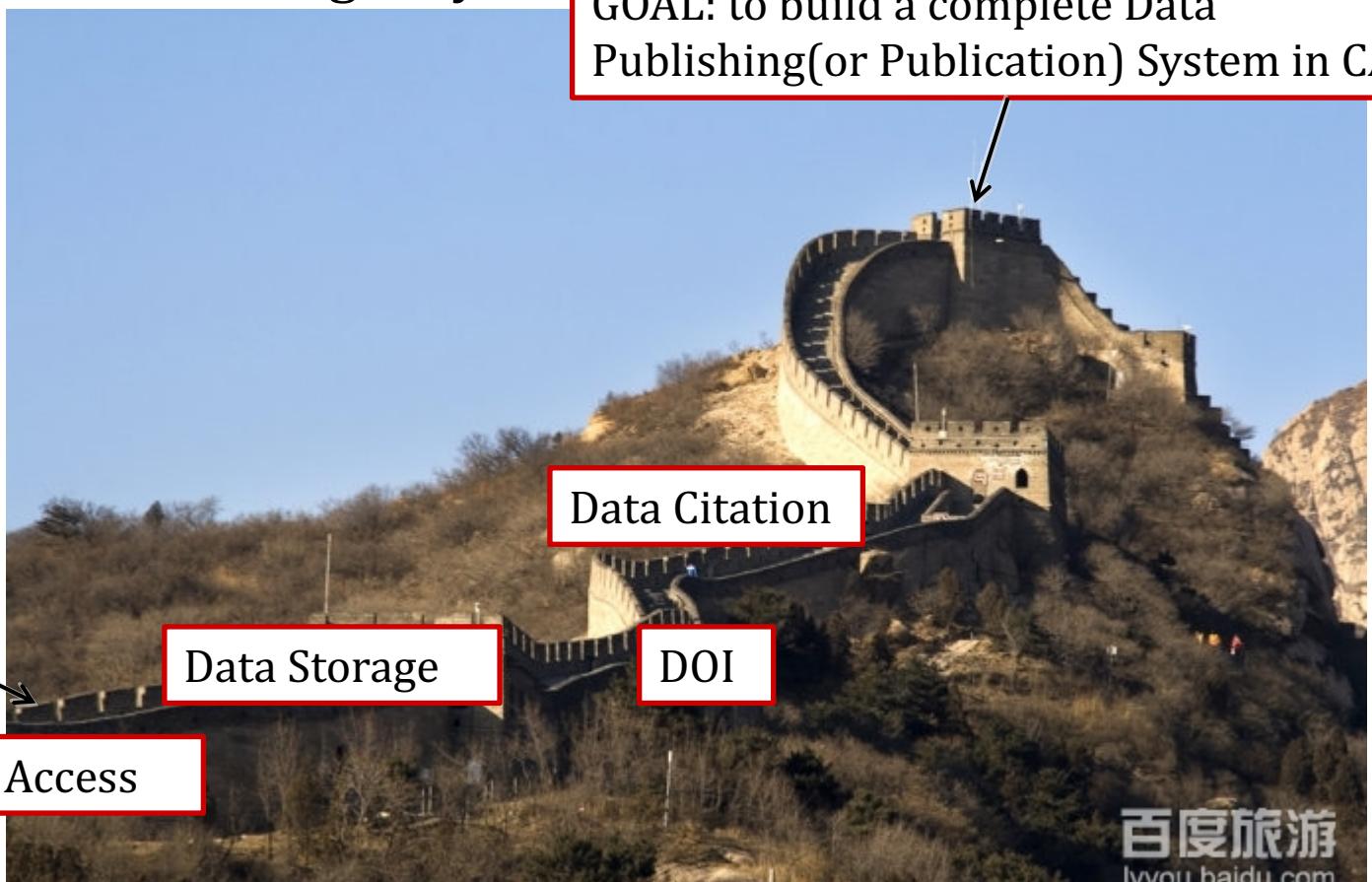


1. Lyubomir Penev et al. Pensoft Data Publishing Policies and Guidelines for Biodiversity Data. Implemented by Pensoft Publishers, 26th of May 2011

# Keep walking

- 千里之行，始于足下
- Our goal is to build a complete data publishing(or publication) system based on CNIC cyber infrastructure, but we still have a long way to go.

GOAL: to build a complete Data Publishing(or Publication) System in CAS



---

*Thank you for your attention!*  
请提宝贵意见！

报告人：沈志宏  
Zhihong SHEN  
中国科学院 计算机网络信息中心

[bluejoe@cnic.cn](mailto:bluejoe@cnic.cn)  
<http://www.csdb.cn>