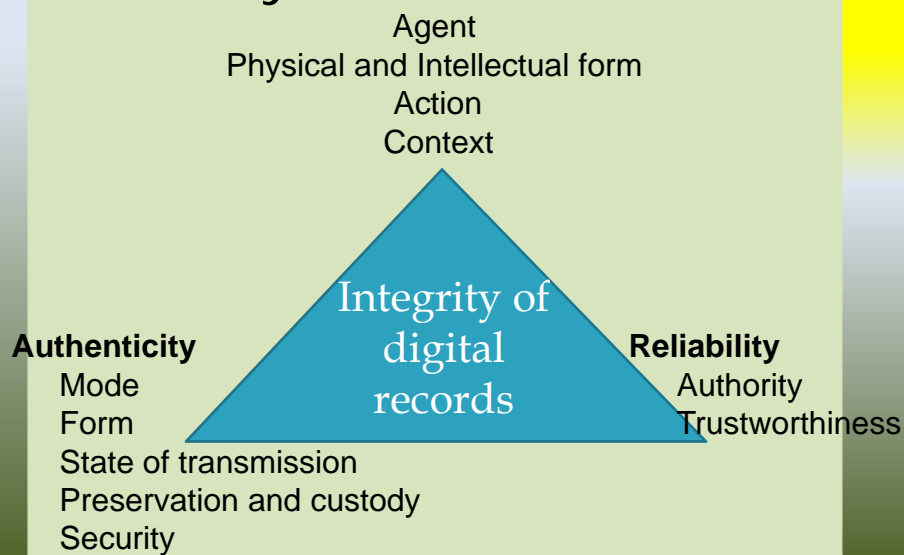


Managing Data: Preservation and Access

IST400/600

Jian Qin

Summary of Duranti's article



Case 1: Submission of genetic sequence data (1)

- Discipline: Biological sciences
- Needs for data:
 - Most journals now expect that DNA and amino acid sequences that appear in articles will be submitted to a sequence database before publication.
 - NCBI relies on data contributions to provide timely and accurate processing and biological review of new entries and updates to existing entries
- How do we translate these needs into data management needs?

Data preservation and access

3

Case 1: Submission of genetic sequence data (2)

- Need
 - Descriptive data elements for genetic sequence data
 - Submission interface for submitting and updating data sets
 - Tool development
 - Search interface for submitted sequence data
 - Backend management
 - Linking between relevant databases
 - Linking between data and publications
 - Data and database administration
 - Long-term preservation of data

Compiled based on <http://www.ncbi.nlm.nih.gov/Genbank/index.html>

Data preservation and access

4

Case 2: querying image data (1)

- Data case: SkyServer (<http://cas.sdss.org/dr6/en/>)
- Discipline: Astronomy
- Data size: 80GB containing 14 million objects and 50 thousand spectra
- Query type:
 - Point-and-click requests for
 - Images of the sky
 - Images of spectra
 - Tabular outputs of the SDSS database

Compiled based on:

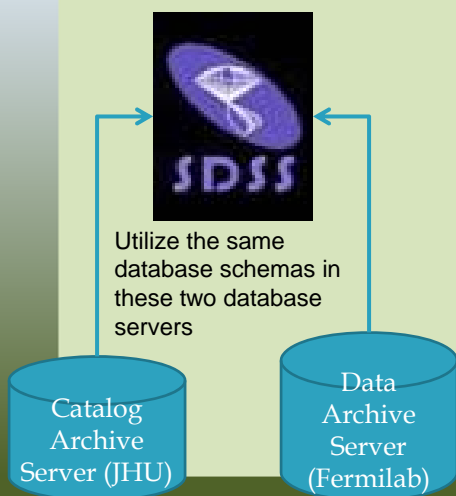
Szalay, A. et al. (2002). <ftp://ftp.research.microsoft.com/pub/tr/tr-2001-104.pdf>

Gray, J. et al. (2002). <http://www.sdss.jhu.edu/ScienceArchive/pubs/msr-tr-2002-01.pdf>

Data preservation and access

7

Case 2: querying image data (2)



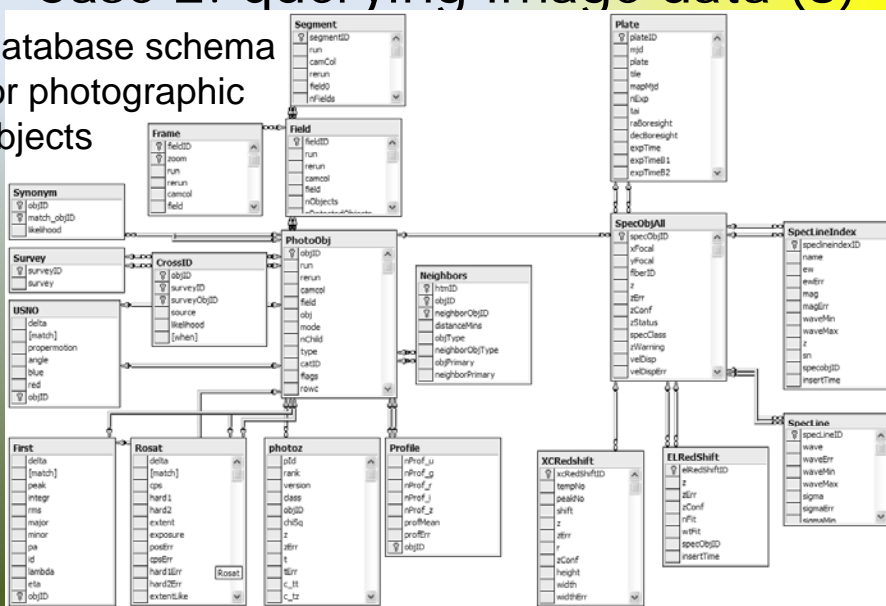
- Data mining:
 - Question answering
 - Defined 20 typical queries
 - Designed the SkyServer database to answer those questions
- Data objects:
 - Photographic objects
 - Spectroscopic objects

Data preservation and access

8

Case 2: querying image data (3)

Database schema for photographic objects



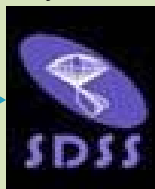
Data preservation and access

9

Case 2: querying image data (4)

Data Transformation Service (DTS) load and convert data as well as check for integrity

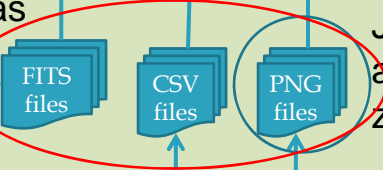
SkyServer



Converted to JPEG images at various zoom levels



The Main 2.5-meter Telescope



Data preservation and access

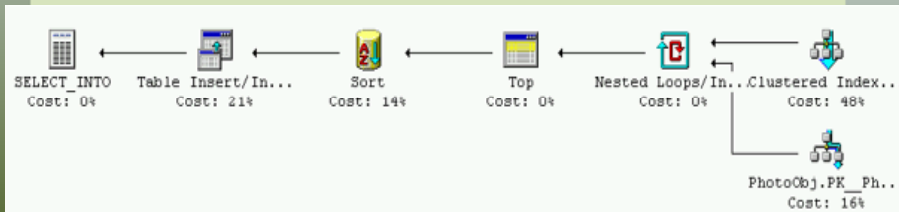
10

Case 2: querying image data (5)

Query: find all galaxies without saturated pixels within 1' of a given point

```

declare @saturated bigint; -- initialized "saturated" flag
set @saturated = dbo.fPhotoFlags('saturated'); -- avoids SQL2K optimizer problem
select G.objID, GN.distance -- return Galaxy Object ID and
into ##results -- angular distance (arc minutes)
from Galaxy as G -- join Galaxies with
join fGetNearbyObjEq(185,-0.5, 1) as GN -- objects within 1' of ra=185 & dec=-.5
on G.objID = GN.objID -- connects G and GN
where (G.flags & @saturated) = 0 -- not saturated
order by distance -- sorted nearest first
    
```



Data preservation and access

11

Case 2: query image data (6)

Query: Provide a list of moving objects consistent with an asteroid

```

select objID, -- return object ID
sqrt( power(rowv,2) + power(colv, 2) ) as velocity, -- velocity
dbo.fGetUrlExpId(objID) as Url -- url of image to examine it.
into ##results
from PhotoObj -- check each object.
where (power(rowv,2) + power(colv, 2)) between 50 and 1000 -- square of velocity
and rowv >= 0 and colv >=0 -- negative values indicate error
    
```



Data preservation and access

12

What do the two cases tell us?

- A lot of data loading and conversion on a daily basis
- Preserving data is an ongoing process
- Preserving data needs metadata (organizing)
 - Who created the data, owns it, what it is about, use policy, etc.
- Short- and long-term access to data (access)
 - what is available, where the available data are located, and how to get them

Data preservation and access

13

Data preservation challenges

- Data formats
 - Vary in data types, e.g. vector and raster data types
 - Format conversions, e.g. from an old version to a newer one
- Data relations
 - e.g. there are data models, annotations, classification schemes, and symbolization files for a digital map
- Semantic issues
 - Naming datasets and attributes

Data preservation and access

14

Data access challenges

- Reliability
- Authenticity
- Leverage technology to make data access easier and more effective
 - Cross-database search
 - Integration applications