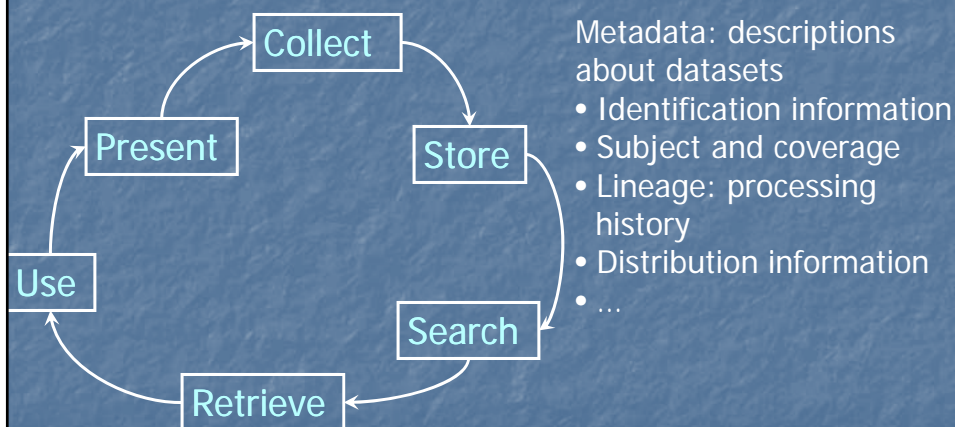


Developing data management project (3): Data organization and preservation

IST400/600

Jian Qin

Review: Science data life cycle



Organization of data

- Organizing datasets
 - Data files have formats
 - Data files have multiple components
 - Storage media
 - Physical locations (file directories)
- Organizing information about datasets
 - Metadata for datasets
 - Metadata schemas
 - Metadata records
 - Storage of these records
 - Retrieval of datasets by using metadata

Datasets need to be associated with metadata so that they can be found and used

Metadata is also data and needs to be stored in databases for retrieval

Storage technology

- Disk systems
- Storage area network



HIGH SPEED

IBM announces first 8 Gb SAN fabric backbone
→ IBM System Storage SAN768B



- Tape systems



Metadata for datasets

- Determine the metadata schema
 - Adopt a standard or develop an in-house schema?
 - Community-based metadata description conventions
 - Understand how scientists search and use data
- Create metadata records
 - Who will create the records? Scientists? Technicians? Data management staff?
 - When will metadata records be created? At the time when data are collected or after?
 - How will metadata be created? Automatic? Manual? Computer-aided, semi-automated?
 - If manual or semi-automated, what tools will be used?

An example

- Institutionalize Metadata Before It Institutionalizes You
 - http://www.fgdc.gov/metadata/documents/InstitutionalizeMeta_Nov2005.doc

A 5-minute quiz (1)

- What are the benefits of metadata for institutions?
- What are the benefits of metadata for individuals?
- _____ is a good way for metadata creation.
- What metadata information is involved in data planning?
- Give two examples of metadata tools.

A 5-minute quiz (2)

- What are the major responsibilities of managers in developing metadata procedures and policies?
- Give an example of metadata production obstacles and the corresponding recommendation

Once metadata schema, procedures, and policies are defined, what next?

IMPLEMENTATION OF METADATA APPLICATIONS

IST400/600 Scientific Data Management

9

Conceptual model vs. implementation model

Identify entities, attributes, and relationships

Conceptual model

Specify constraints, rules, operations

Define data fields, data types, data entry rules and controls

Implementation model

Design queries, triggers, and stored procedures, and user interfaces

IST400/600 Scientific Data Management

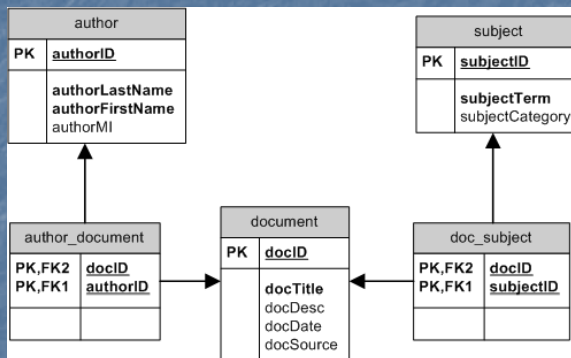
10

Conceptual modeling tools

- E-R Diagram
- Unified Modeling Language (UML)
- XML/RDF (Resource Description Framework)
- Web Ontology Language (OWL)

E-R modeling

- Entities
- Attributes
- Relationships

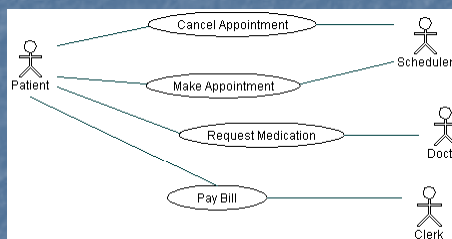


UML modeling (1)



Use case diagrams

- Actor who initiates the events involved in the task
- Communication: connection between actor and use case



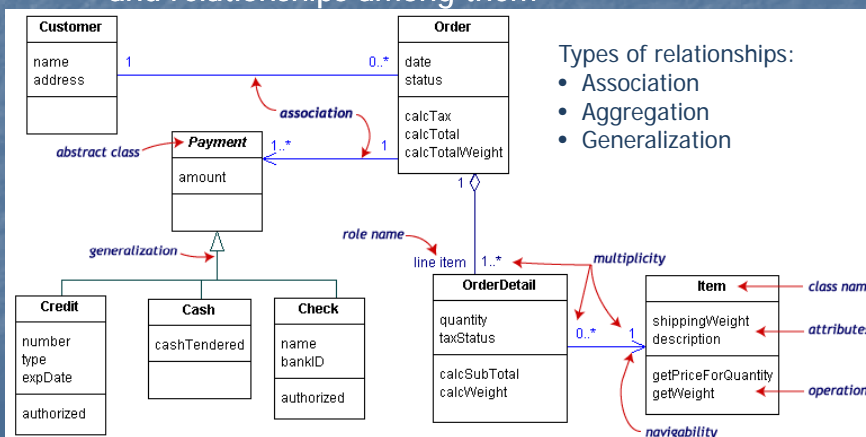
Source: Miller, Randy. (2003). Practical UML: A hands-on introduction for developers. <http://dn.codegear.com/article/31863>

UML modeling (2)



Class diagrams:

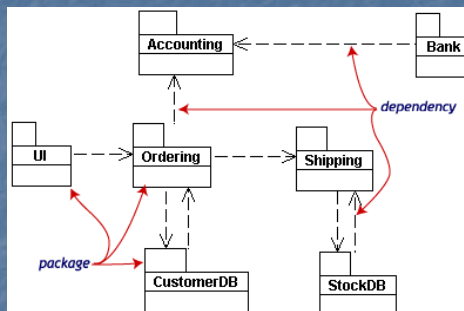
- Provides an overview of a system by showing its classes and relationships among them



UML modeling (3)



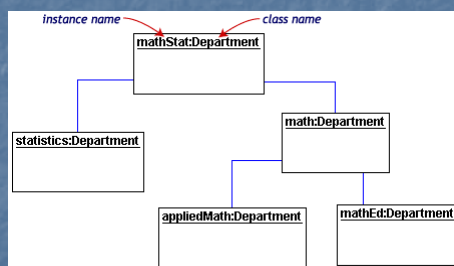
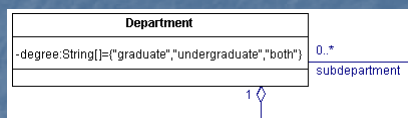
- Packages
 - Package: a collection of logically related UML elements
 - Dependencies: one package depends on another if changes in the other could possibly force changes in the first



UML modeling (4)



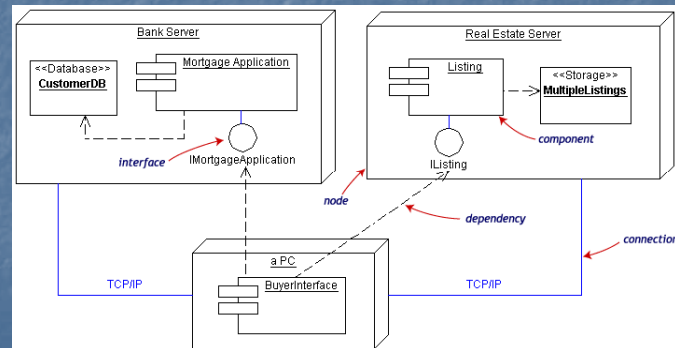
- Object diagrams
 - Contains instances instead of classes



UML modeling (5)



- Component and deployment diagrams
 - Component: a code module
 - Deployment diagram: shows physical configurations of software and hardware.



IST400/600 Scientific Data Management

17

UML modeling (6)



- Other diagrams in UML:
 - Sequence diagrams
 - Collaboration diagrams
 - Statechart diagrams
 - Activity diagrams
 - Component diagrams

IST400/600 Scientific Data Management

18

Example: Metadata fields mapped to data flow

Data Development Stage	Metadata Information
Data Planning	Identification Information title, originator, abstract, purpose, keywords, time period Data Organization point, raster, vector Spatial Referencing coordinate system and datum Entity and Attributes (planned)
Data Processing	Data Quality completeness, positional accuracy, geoprocessing steps
Data Analysis	Data Quality attribute accuracy, analysis steps Entity and Attributes (results) Metadata Reference

IST400/600 Scientific Data Management 19

UML modeling for the example

- What would a use case diagram look like?
- What would a package diagram look like?
- What would an object diagram look like?