



# **Scientific Data Formats**

### **Data Diversity**

- Scientific data come from a variety of sources (remote-sensing instruments, sensors, experiments, simulations...), often in proprietary data formats.
- Typical "raw data" and data products include multi-dimensional arrays, images, spectra, vector fields
- Different scientific communities have different data management and data processing needs, and thus use different data formats.

### The purpose of specialized scientific data formats is to

- store, manage, exchange, share, and archive data from scientific applications
- provide the base for data management and analysis tools, such as data integration, sharing, visualization, and archiving
- Provide complete software solutions with support of a variety of programming languages, ranging from Fortran 77 to Java

SDM Tutorial, EDBT'06, Gertz, Ludäsche

## Scientific Data Formats (2)

### There are several key data formats Earth & Space Sciences

- HDF (Hierarchical Data Format)
- CDF/netCDF (Common Data Format)
- FITS (Flexible Image Transport System)
- XSIL (eXtensible Scientific Interchange Language)
- GML (Geography Markup Language)
- GRIB (Grids in Binary)
- CCM (Community Climate Model History Tape Format)

### ...and many XML-based languages in the Life Sciences...

Chemical Markup Language (*CML*) – Molecular Dynamics [Markup] Language (*MoDL*) – MicroArray and Gene Expression Markup Language (*MAGE-ML*) – Genome Annotation Markup Elements (*GAME*) – BIOpolymer Markup Language (*BIOML*) – Numerical Data Markup Language (*NDML*) – Protein Extensible Markup Language (*PROXIML*) – Systems Biology Markup Language (SBML) ...

(see also http://xml.coverpages.org/)









# **Common Data Format (CDF)**

Designed and developed in 1985 by the National Space Science Data Center at NASA Goddard Space Flight Center.

- Conceptual data abstraction for storing, manipulating, and accessing multidimensional data sets.
- Basic component: software programming interface that is a device independent view of the CDF data model.
- CDF files created on any given platform can be transported to any other platform onto which CDF is ported and used with any CDF tools or layered applications.
- CDF software package is used by hundreds of government agencies, universities, and private and commercial organizations as well as independent researchers on both national and international levels.
- Adopted by the International Solar-Terrestrial Physics project as well as the Central Data Handling Facilities (CDHF) as their format of choice for storing and distributing data.

nssdc.gsfc.nasa.gov/cdf/cdf\_home.html

SDM Tutorial, EDBT'06, Gertz, Ludäscher

# Network Common Data Format (netCDF)

- Developed by University Corporation for Atmospheric Research around 1990.
- netCDF library defines a machine-independent format for representing scientific data.
- Together, the APIs for array-oriented data access, library, and format support the creation, access, and sharing of scientific data.



### **Objectives:**

- Self-describing. A netCDF file includes information about the data it contains.
  Portable. A netCDF file can be accessed by computers with different ways of storing integers, characters, and floating-point numbers.
- Direct-access. A small subset of a large dataset may be accessed efficiently,
- Appendable. Data may be appended to a properly structured netCDF file
- Sharable. One writer and multiple readers may simultaneously access file
- Archivable. Access to all earlier forms of netCDF data will be supported by current and future versions of the software.

http://www.unidata.ucar.edu/





# Flexible Image Transport System (FITS)

- Developed 1981 by NASA
- Data format most widely used within astronomy for transporting, analyzing, and archiving scientific data files.
- More than just another image format (such as JPG or GIF)
- Primarily designed to store scientific data sets consisting of multi-dimensional arrays (images) and 2-D tables organized into rows and columns of data.



- FITS file is comprised of segments called Header/Data Units (HDUs)
- First HDU is called the "Primary Array", which can contain a 1-999 dimensional array of integers or floating point numbers.
- A typical primary array could contain a 1-D spectrum, a 2-D image, or a 3-D data cube.
- Any number of additional HDUs may follow the primary array; these are referred to as FITS "extensions" (Images, ASCII tables, binary tables)















# Trends and Developments

There are many efforts two wrap scientific data represented in the data format XYZ into XML. For example,

- XML-based markup language called CDF Markup Language to describe CDF data and metadata and created the following two utilities in Java:
  - CDF2CDFML dumps the contents of a CDF file into a XML file that conforms to the CDF DTD or CDF schema.
  - CDFML2CDF creates a CDF file from an XML file that conforms to the CDF DTD or CDF schema.
- To promote data exchanges among space scientists, the CDF office has developed a Web service called Data Translation Web Service (DTWS) and a client to talk to DTWS via a web browser. The DTWS is a web service based on SOAP. Supported translations include
  - CDF-to-netCDF, CDF-to-FITS, CDF-to-CDFML, netCDF-to-CDF
  - FITS-to-CDF, HDF4-to-CDF

http://translators.gsfc.nasa.gov

- XSLT is a helper in doing such transformations!

SDM Tutorial, EDBT'06, Gertz, Ludäscher

## **Processing Scientific Data**

### In Earth Sciences and Astrophysics

- Almost exclusively based on a file-processing approach:
  - 1. obtain files from instrument, experiment, or simulation (ftp, gridFTP,...)
  - 2. run program on files, typically local
  - 3. record data product as file and store it either locally or remotely
- Data and metadata go through numerous complex processing steps formulated as (image) pipelines and workflows (→ Module 4)
- Pipelines are assembled manually, often through scripting (Perl, Python, ...)

There are many applications settings in which incoming data is delivered in a continuous, streaming fashion but the processing occurs in a file-based approach.

- Weather forecasting, wildfire detection, hurricane tracking...
- In general, there are many types of remote-sensing applications in which real-time processing is crucial.









# **Data Stream Processing**

### **Objectives:**

- Exploit models, techniques, and concepts developed for traditional (relational and XML) data stream management systems
- Develop stream management system against which users can formulate queries and analytical operations
- Move query processing and computation to the data

### How to go about this?

- Get an understanding of the science domain, talk to scientists
- Develop formal data model for scientific data of concern (e.g., multidimensional arrays in HDF), streams of such data, and operations on data streams
- Start with the design and implement a stream simulator that takes science data files and converts them into stream



There are some scientific data models out there (developed by the DB community) that look at non-streaming data, e.g.,

- Multi-dimensional arrays (Marathe and Salem, VLDB Journal 2002; Libkin, SIGMOD 1996)
- Array algebra (Baumann, VLDB Journal 1994, RasDaMan)
- Gridded data sets (Howe and Maier, VLDB 2004)

Can't we just take "conventional" data stream processing models, techniques, and architectures?

- It is not a good idea to map hundreds of millions of pixels to relations
- Operations on multi-dimensional data (e.g., raster images) are often much more complex than selection, project, join (e.g., re-projections or affine image transforms)
- Data sets exhibit spatio-temporal characteristics
- Queries are typically formulated using graphical user interfaces

SDM Tutorial, EDBT'06, Gertz, Ludäsche

### **Streaming Geospatial Image Data**

### A case study:

- GeoStreams: A query processing architecture for streaming remotely-sensed image data <u>http://geostreams.ucdavis.edu</u>
- GOES West, geostationary weather satellite operated by NOAA
- 5 imager bands, 19 sounder channels
- Delivers data in GVAR format at a rate of 2.1Mbits/sec (~22GB/day)



### How does the stream data model look like?

- Infinite *point set*, with points of the form (x,y,t), each point has *point value* x,y = spatial location of point (pixel), t = timestamp
- Point set is a topological space, value set is a homogeneous algebra
- Spatial component of point set corresponds to a regularly space lattice and is *geo-referenced* (we are talking about field-based data...)









# **References Module III**

- Hierarchical Data Format (HDF). <u>http://hdf.ncsa.uiuc.edu/</u>
- HDF: the hierarchical data format. Brand Fortner, DR DOBB'S J SOFTWARE TOOLS
  PROF PROGRAM. Vol. 23, no. 5, pp. 42, 44-48. May 1998
- Common Data Format (CDF). <u>http://cdf.gsfc.nasa.gov/</u>
- A software package for the data-independent management of multidimensional data L.Treinish, M. Gough, EOS. Vol. 68, pp. 633-635. 14 July 1987
- Network Common Data Format (netCDF). <u>http://www.unidata.ucar.edu/software/netcdf/</u>
- NetCDF: an interface for scientific data access, R. Rew, G. Davis, Computer Graphics and Applications, Jul 1990, 10(4): 76-82
- Flexible Image Transport System (FITS). <u>http://fits.gsfc.nasa.gov/</u>
- FITS a Flexible Image Transport System, D.C. Wells, E.W. Greisen, International Workshop on Image Processing in Astronomy. Proceedings of the 5th. Colloquium on Astrophysics, held in Trieste, Italy, June 4-8, 1979. Editors, G. Sedmak, M. Capaccioli, R.J. Allen.; Publisher, Osservatorio Astronomico di Trieste, Trieste, Italy, 1979
- eXtensible Scientific Interchange Language (XSIL). http://www.cacr.caltech.edu/SDA/xsil/
- A High-Performance Active Digital Library, Roy Williams, Bruce Sears, Proceedings of HPCN98, Amsterdam, April 1998, eds. L. O. Herzberger and P. M. A. Sloot, Lect. Notes Comp. Sci. (Springer)
- Geography Markup Language. <u>http://www.opengeospatial.org/</u>
- Deep Lens Survey. http://dls.bell-labs.com/
- GeoStreams project. <u>http://geostreams.ucdavis.edu</u>