



**NRC · CNRC**

*From Discovery to Innovation...*

# *The Virtual Observatory: The Future of Astrophysics Data Handling*

**David Schade**

**Canadian Astronomy Data Centre**

**Herzberg Institute for Astrophysics**

**National Research Council Canada**

*With support from the Canadian Space Agency*



National Research  
Council Canada

Conseil national  
de recherches Canada

Canada

## *Astronomy and Astrophysics*

Does fairly well in information technology

- Has excellent online literature services
  - ADS Abstracts
  - Journals
  - Preprints
- Has a good history of data archiving
- Has reasonable data access policies
- **BUT**
  - As a scientist it is frustrating and time-consuming to locate suitable data and data quality is often sub-standard



# A Brief History of Data Archiving in Astronomy

## *History*

- NASA has been a driving force in data archiving for astronomy
- Canada-France-Hawaii Telescope (CFHT) was a pioneer in archiving data from ground-based observatories
- Digital Revolution in astronomy happened in the 1980's





# The Canadian Astronomy Data Centre

## *History*

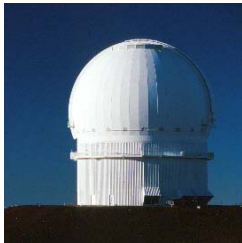
- Canadian Astronomy Data Centre was created in 1986
- Astronomers and Computer Scientists
- supported by the Canadian Space Agency
- original mandate: to serve Hubble Space Telescope

## CADC Firsts

- First web interface in astronomy
- Previews of data
- On-the-fly calibration
- Advanced image processing



## *Current Collection at CADC*



**Canada-France-Hawai'i Telescope  
Hawai'i**



**Dominion Radio Astrophysical Observatory  
British Columbia**

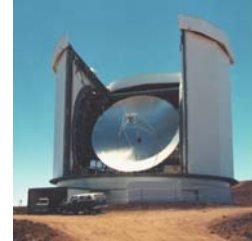


**Gemini South Telescope  
Cerro Pachón, Chile**

**Hubble Space Telescope**



**James Clerk Maxwell Telescope  
Hawai'i**



**Gemini North Telescope  
Hawai'i**



# The Canadian Astronomy Data Centre

 <a href="#">HST Archive</a>	 <a href="#">CFHT Archive</a>	 <a href="#">JCMT Archive</a>
 <a href="#">Digitized Sky Survey</a>	 <a href="#">SIMBAD WWW Access</a>	 <a href="#">IRAS HCON Access</a>
 <a href="#">CD-ROM Access</a>	 <a href="#">CISTI/HIA Library</a>	 <a href="#">Astronomical Meetings</a>
 <a href="#">CGPS Archive</a>	 <a href="#">ESO Archives</a>	 <a href="#">LaPalma Archives</a>
 <a href="#">AAT Archives</a>	 <a href="#">ATNF Archives</a>	 <a href="#">SN1987A Archives</a>
 <a href="#">USNO Guide Star</a>	 <a href="#">UKIRT archive</a>	 <a href="#">MDS Survey</a>
 <a href="#">Hipparcos Catalogue</a>	 <a href="#">GSC Catalogue</a>	 <a href="#">Other Catalogues</a>

- Many Services
  - Digitized Sky Survey
  - Archive Inter-operability
- Meta-Data Catalogues
  - 19 databases
  - 80,000,000 rows
  - 34 gigabytes
- Data Files
  - 12,000,000 files
  - 20 terabytes



# A Brief History of Data Archiving in Astronomy

- Archiving is a word that does not adequately describe the activities, capabilities, and functions of data centres
  - Store, protect, catalogue, facilitate access, lobby for open data policy
  - Lobby for effective handling of data and metadata
  - Develop processing pipelines to add value
  - Execute processing on request





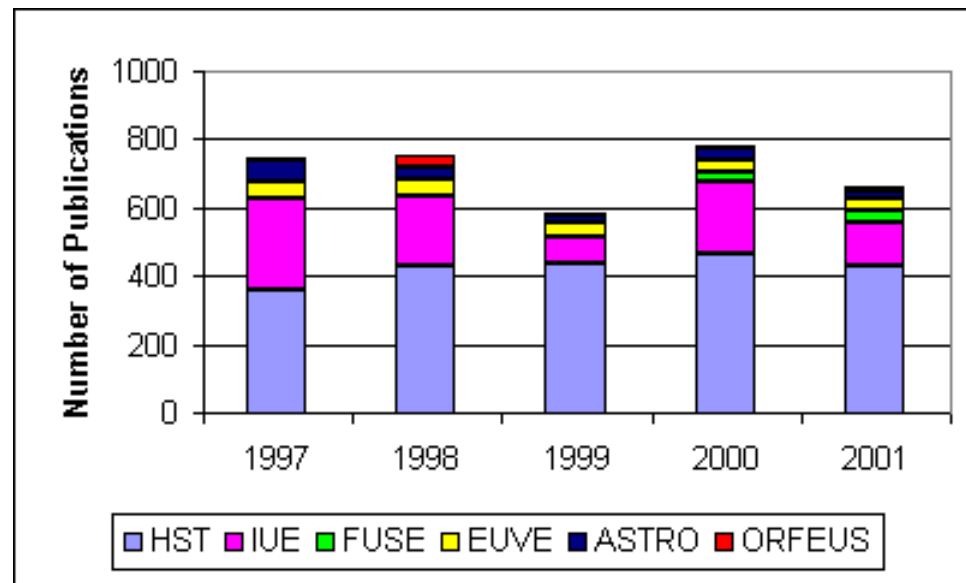
# A Brief History of Data Archiving in Astronomy

Do astronomers publish research based on archival data?



# Scientific Impact of Multi-Mission Archive at Space Telescope Science Institute

- ◆ **~10% of the most-cited papers in the ISI database are based on MAST archival data**
- ◆ Over 600 papers/year with HST and other archives
- ◆ HST Data: Retrieval rate is 4 times the ingestion rate
- ◆ Over 30,000 datasets requested per month (over 8,000 are non-HST data); ~400,000 web hits per month



From Megan Donahue  
STScI



**NRC · CNRC**

## *International VO initiatives*

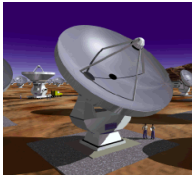
- Massive homogeneous survey datasets are being created
  - Sloan Digital Sky Survey
  - 2MASS infrared survey
  - Canada-France-Hawaii Legacy Survey
- Multi-wavelength survey datasets can be constructed
- Network bandwidth is increasing
- Astronomers have embraced many online services
- Funding agencies are receptive



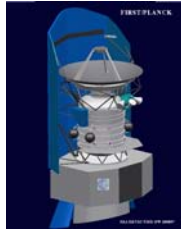
*New types of science will be possible with new modes of data access*



**NRC · CNRC**



**ALMA**



**FIRST**



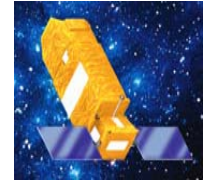
**NGST**



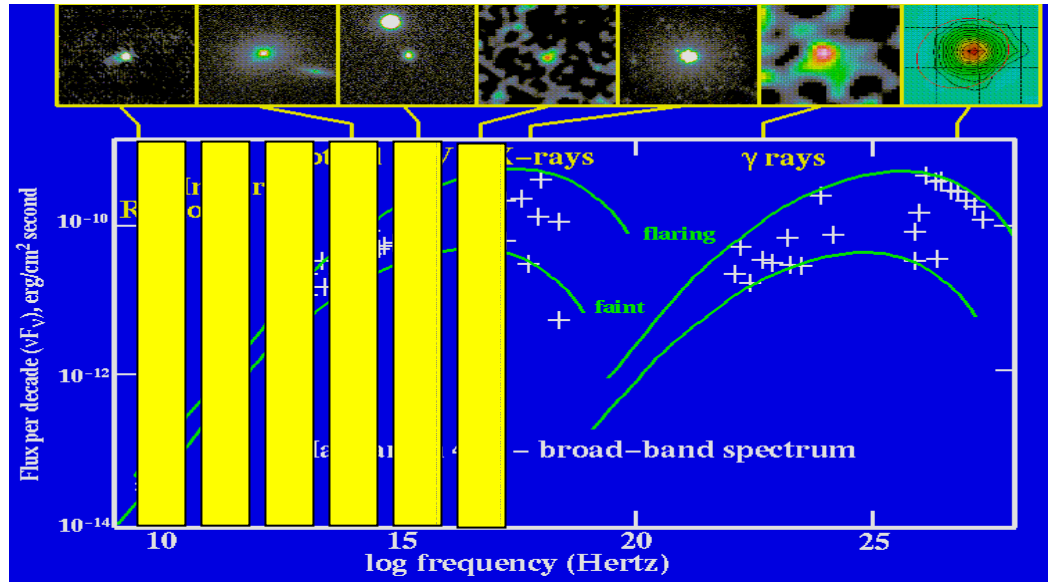
**GEMINI**



**CFHT**



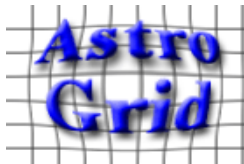
**FUSE**





## *International initiatives: Different strokes for different folks*

- Major initiatives in Canada, the United States, the European community, the United Kingdom. (Australia, India, Russia)
- Each VO group has their own view of what it means to produce a VO and what the priorities should be.
- U.S.: A high-level distributed infrastructure, tools.
- U.K.: Several thrusts: data pipelines, ontology, data mining
- Europe: VO closely associated with operational data centers and other groups
- Canada: VO is within the Canadian Astronomy Data Centre
- **Data-centric versus infrastructure-centric views**



## *Definition*

The Virtual Observatory will be said to exist when astronomers can successfully execute scientific queries that seamlessly cross archive boundaries and wavelength boundaries, can combine the returned datasets in a way that permits their joint processing, and can achieve this without the need to understand engineering-level details of the instrument that produced the returned datasets.

- Discussions of online toolsets, grid computing, distributed datasets, etc. are implementation details.
- “Observatory” implies that the product is pixel data
- Are analysis tools and catalogues legitimate products?
- The Virtual Observatory needs to be defined in terms of capabilities delivered to scientists (the users).



## *Convergence ?*

- Despite the differences in viewpoint at this early stage of the VO game, the approaches will converge as projects become reality.
- Interoperability
- Standards
- Integration
  
- But there need to be new investments in data archiving centres to match the investment in higher level infrastructure.
- POTENTIAL CONTENT CATASTROPHE FOR VO



## Standard Practice

- Proprietary period of 1-2 years during which only the proposer of the observations may access those data
- Some data is calibrated and much is not
- Data quality is an issue
- Metadata completeness is an issue
- Metadata quality is an issue





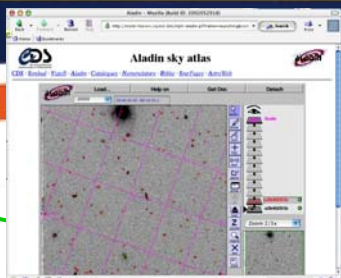
- Past history
  - Canada has benefited enormously from open data access (and facility access) policies of the United States
    - Data access: Largely NASA
    - Facility access: NOAO and many others
- NASA has been very progressive
- Many facilities have had no channels to access data (NOAO) , some do not save and protect data (e.g., Keck telescopes: U. California and California Institute of Technology)
- Europe has been very progressive: **BUT now the archives of the European Southern Observatory are CLOSED to astronomers outside of Europe.**



- Present-day data policies are very mixed:
  - Tension between observatory operations and archiving needs
- Canada has been progressive
  - Canada-France-Hawaii Telescope archives since 1980s
    - Data quality has been fair
- Canada and Chile were the leading forces in creating an archive for the Gemini telescopes (partners U.S., U.K., Canada, Argentina, Chile, Brazil, Australia)
- Canada and France are considering a long (~ 3 years) proprietary period for the CFHT Legacy Survey

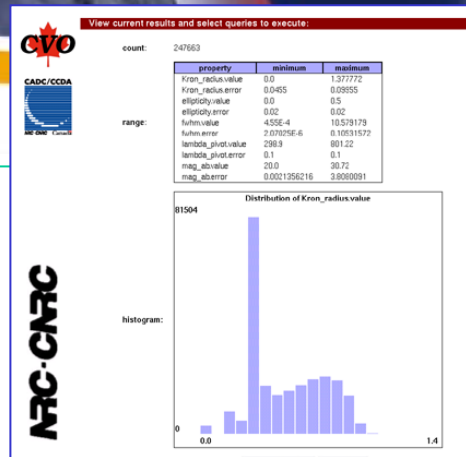


# CVO Architecture



**Archives**

**VoPix**



**NRC-CNRC**

Archives publish to the VO



Web interface to archive

**VoProc**

**VoSrc**



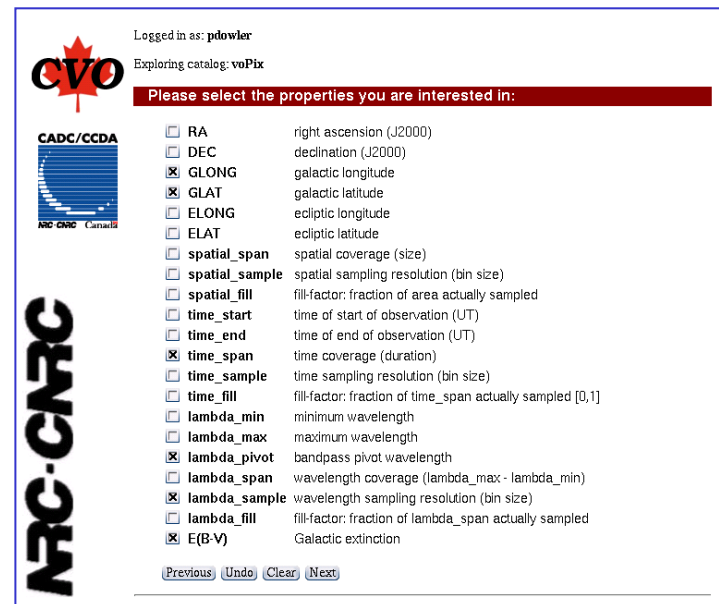
CVO is a software layer above the archive level

**NRC-CNRC**



The CVO system provides a view on archive content:

- High-level
- Scientific descriptors
- Not instrument specific
- Integrates different archive content



Logged in as: pdowler  
Exploring catalog: voPix

**Please select the properties you are interested in:**

<input type="checkbox"/>	RA	right ascension (J2000)
<input type="checkbox"/>	DEC	declination (J2000)
<input checked="" type="checkbox"/>	GLONG	galactic longitude
<input checked="" type="checkbox"/>	GLAT	galactic latitude
<input type="checkbox"/>	ELONG	ecliptic longitude
<input type="checkbox"/>	ELAT	ecliptic latitude
<input type="checkbox"/>	spatial_span	spatial coverage (size)
<input type="checkbox"/>	spatial_sample	spatial sampling resolution (bin size)
<input type="checkbox"/>	spatial_fill	fill-factor: fraction of area actually sampled
<input type="checkbox"/>	time_start	time of start of observation (UT)
<input type="checkbox"/>	time_end	time of end of observation (UT)
<input checked="" type="checkbox"/>	time_span	time coverage (duration)
<input type="checkbox"/>	time_sample	time sampling resolution (bin size)
<input type="checkbox"/>	time_fill	fill-factor: fraction of time_span actually sampled [0,1]
<input type="checkbox"/>	lambda_min	minimum wavelength
<input type="checkbox"/>	lambda_max	maximum wavelength
<input checked="" type="checkbox"/>	lambda_pivot	bandpass pivot wavelength
<input type="checkbox"/>	lambda_span	wavelength coverage (lambda_max - lambda_min)
<input checked="" type="checkbox"/>	lambda_sample	wavelength sampling resolution (bin size)
<input type="checkbox"/>	lambda_fill	fill-factor: fraction of lambda_span actually sampled
<input checked="" type="checkbox"/>	E(B-V)	Galactic extinction

[Previous](#) [Undo](#) [Clear](#) [Next](#)










Multi-wavelength, hierarchical object catalogues are a representation of the state of our understanding of the universe.

Logged in as: pdowler  
Exploring catalog: voPix

**Please select the properties you are interested in:**

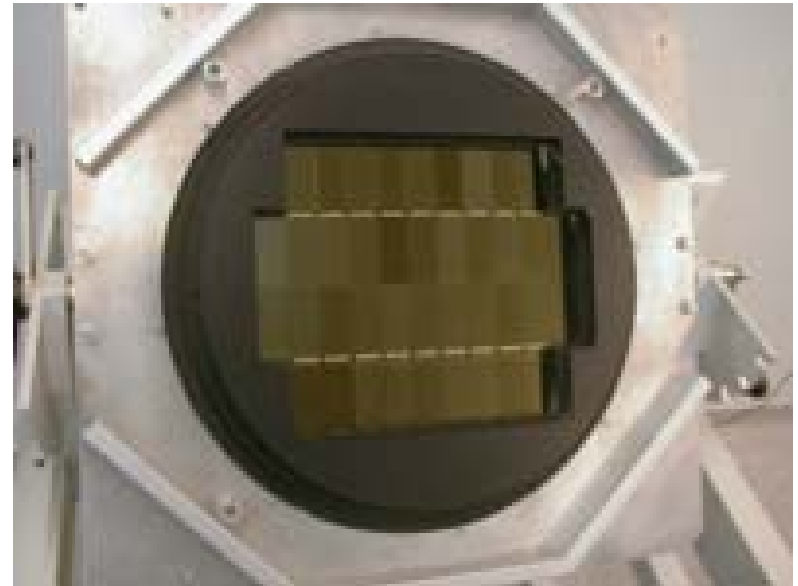
<input type="checkbox"/>	RA	right ascension (J2000)
<input type="checkbox"/>	DEC	declination (J2000)
<input checked="" type="checkbox"/>	GLONG	galactic longitude
<input checked="" type="checkbox"/>	GLAT	galactic latitude
<input type="checkbox"/>	ELONG	ecliptic longitude
<input type="checkbox"/>	ELAT	ecliptic latitude
<input type="checkbox"/>	spatial_span	spatial coverage (size)
<input type="checkbox"/>	spatial_sample	spatial sampling resolution (bin size)
<input type="checkbox"/>	spatial_fill	fill-factor: fraction of area actually sampled
<input type="checkbox"/>	time_start	time of start of observation (UT)
<input type="checkbox"/>	time_end	time of end of observation (UT)
<input checked="" type="checkbox"/>	time_span	time coverage (duration)
<input type="checkbox"/>	time_sample	time sampling resolution (bin size)
<input type="checkbox"/>	time_fill	fill-factor: fraction of time_span actually sampled [0,1]
<input type="checkbox"/>	lambda_min	minimum wavelength
<input type="checkbox"/>	lambda_max	maximum wavelength
<input checked="" type="checkbox"/>	lambda_pivot	bandpass pivot wavelength
<input type="checkbox"/>	lambda_span	wavelength coverage (lambda_max - lambda_min)
<input checked="" type="checkbox"/>	lambda_sample	wavelength sampling resolution (bin size)
<input type="checkbox"/>	lambda_fill	fill-factor: fraction of lambda_span actually sampled
<input checked="" type="checkbox"/>	E(B-V)	Galactic extinction

[Previous](#) [Undo](#) [Clear](#) [Next](#)



## CFHT MegaCam

- A 40 CCD camera
  - 320 Megapixels
  - 1 square degree on the sky
- Raw Data Rate
  - 720 megabytes per image!
  - 100 gigabytes per night!
  - 20 Terabytes per year!



## CFHT Legacy Survey

- SCIENCE
  - *Determine the fate of the universe*
- Data Policy
  - Data are released immediately to the Canadian and French communities and to the world after a proprietary period



## CFHT Legacy Survey

- Partnership between CFHT (Hawaii), CADC (Victoria), TERAPIX (Paris), CDS (Strasbourg)
- Science: Supernovae, Weak Lensing, Kuiper Belt
- 5 years / 500 nights
- 20 Terabytes per year
- 50 million objects with high-quality imaging
- Processed image products and catalogues
- **100 Terabyte project**

## Data Distribution via network

- 150 Mbps **continuously** for 5 years
- CANET/BCNET
- Need Gbit network





- DVD jukeboxes
  - 4.7 Gbytes/disk
  - 16 \$/Gbyte
  - 11.5 Tbytes/m<sup>2</sup>
  - 6 jukeboxes/year
  - 3,900 disks/year
- High overhead
  - Operationally
  - Physical space



- Spinning disks
  - 20 Terabytes in each rack
- Processing
  - 20 1.5 GHz CPUs in each rack
- Cost effective
- Effective use of space
- Reliability ???



## *Astronomy and Astrophysics*

- Virtual Observatory recognizes the value and effectiveness of good information management in astrophysics
- Astronomy has a good IT foundation to build upon
- Funding agencies are receptive
- Data access policies need to be monitored for problems
- Virtual Observatory needs to invest in both infrastructure **and in data**

**THE END**



**NRC · CNRC**





# A Brief History of Data Archiving in Astronomy

## Outline

- History and CADC
  - I will neglect NASA and concentrate on what I know
    - CFHT archived their digital data in the 1980's
    - Plates were taken home but remained the property of the observatory which never recalled them
  - Hubble Space Telescope opened doors in archiving for optical astronomers
  - Archiving is a word that has outlived its usefulness
    - Archive functions: Store,protect,catalogue, facilitate access, lobby for open data policy
    - Non-archive functions: Develop processing pipelines to add value
  - Archive Status: Do astronomers do research with archival data? YES HST examples
    - Deliver the archive over and over/ Megan's publication numbers





## Outline

- Virtual Observatory initiatives: IVOA
  - Definition of the VO
  - Different strokes for different folks
  - High-level infrastructure
  - Where's the data?
  - There need to be data-centric initiatives also
  - THE GOALS ARE WELL-ALIGNED

