

## NASA's Long-Term Astrophysics Data Archives

L. M. Rebull, IRSA 17 Oct 2016



## NASA

## With Input From:

- Vandana Desai (IRSA)
- Harry Teplitz (IRSA)
- Steve Groom (IRSA)
- Rachel Akeson (NExScl)
- Bruce Berriman (NExScl)
- George Helou (IPAC)
- David Imel (IPAC)
- Joe Mazzarella (NED)

- Alberto Accomazzi (ADS)
- Tom McGlynn (GSFC)
- Alan Smale (GSFC)
- Rick White (MAST)



# Having said that...

 I'm an astronomer at IRSA, so that is what I know best, and most of my examples are IRSAfocused by necessity.



## NASA's Commitment to Astrophysics Data Archives

- "NASA has regarded data handling and archiving as an integral part of space missions."
  - "This support now provides the major return on the considerable investment the agency made... over the past 20 years."



## "A Sustainable Archive"

- Provides data discovery and analysis tools. Facilitates new science.
- Contains high-quality, reliable data.
- Provides simple and useful tools to a broad community.
- Provides user support to the novice as well as to the power user.
- Adapts and evolves in response to community input.



## An Archive's Job

- Ingest new data (and reprocessing of old data).
- Maintain/serve vital repository of irreplaceable data:
  - Support for <u>observation</u> planning and <u>mission</u> planning.
  - Resource for original science.
  - High level science products.
- Enable cutting-edge research:
  - API and Virtual Observatory.
  - User support by experts.
  - New/enhanced services.
  - Multi-wavelength projects.







## Archival Papers Outnumber Non-Archival (GO/PI) Papers



### Archives double an observatory's output!





## **IRSA Science Highlights**



WISE+Spitzer discover the coldest brown dwarf (Luhman 2014, ApJL 786, L18)





WISE morphological study of Wolf-Rayet nebulae, Toala et al. (2015, A&A 578, 66)



WISE+2MASS+PanSTARRS data may reveal super-void in CMB cold spot seen by Planck; (Szapudi et al. 2015, MNRAS, 450, 288)

**Combination of Surveys** 

## **MAST Science Highlight:** R8799 b,c,d imaged by HST in 1998

#### Exoplanet HR 8799 System



NASA, ESA, and R. Soummer (STScI)

### planet b: planet c: **NASA's Astrophy** 83,000x fainter than star 36,000x fainter than star at 1.72 arcsec at 0.96 arcsec

planet d: 33,000x fainter than star at 0.60 arcsec

Post-processing speckle subtraction, >an order of magnitude contrast improvement over "state of the art" when data taken in 1998.

Soummer et al. 2011, Pueyo et al. 2015



## **HEASARC Science Highlight**



Tarantula Nebula: **Combining 6** years of Fermi data to discover the first extragalactic gamma-ray pulsar!

Credit: NASA/DOE/Fermi LAT Collaboration, 2015, Science, 350, 801

## Some NASA Archives by Center

- IPAC:
  - IRSA IR, sub mm
  - NED Extragalactic
  - NASA Exoplanet Archive
  - KOA (w/ WMKO) Keck Observatory
- STScI : MAST UV, optical, IR
- GSFC : HEASARC high energy, CMB
- SAO/CfA:
  - ADS literature
  - Chandra











IRSA

http://irsa.ipac.caltech.edu/

Caltech

- IRSA = NASA/IPAC Infrared Science Archive, @ Caltech, @ IPAC.
- Charter is to provide interface to all NASA infrared and sub-mm data sets. (~1  $\mu m$  -> ~1 cm).
- Founded 1993, original home to IRAS data (1983).
- IRSA ensures the legacy of NASA's "golden age" of IR:
  - Enable research that has not yet been envisioned.
  - Priorities set by missions and the community.
  - Support future flight missions.
- IRSA datasets are cited in about 10% of astronomical refereed journal articles.
- Total holdings over a petabyte (>1000 TB); >120 billion rows in catalogs as of 9/2016.
- Total number queries: Over 33.7 million queries, 255 TB downloaded 1-9/2016.









- NED = NASA/IPAC Extragalactic Database (@ Caltech, @ IPAC).
- Primary hub for multi-wavelength research on extragalactic science.
- Merges data from catalogs and literature.
- 1000s of extragalactic papers per year, with unique measurements for millions of objects.
- 215 million objects with 256 million cross-IDs (from >102,000 articles/catalogs)!
  - 2 billion photometric data points joined into spectral energy distributions.
- Myriad cross-links, notes, etc.
- Updates every few months.



http://ned.ipac.caltech.edu/

Caltech



The NASA/IPAC Extragalactic Database (NED) is operated by the Jet Propulsion Laboratory, California Institute of Technology, under contract with the National Aeronautics and Space Administration.







## Caltech

## **NASA Exoplanet Archive**



http://exoplanetarchive.ipac.caltech.edu/

- (Also @ Caltech, @ IPAC)
- Focused on confirmed and candidate exoplanets and on data searching for exoplanets
- Includes Kepler data, and US portal to CoRoT data.
- Online tools to work with these data, like the periodogram service.
- Place for observers to upload/ share data (Exo-FOP).



(ipac)

#### NASA EXOPLANET ARCHIVE A SERVICE OF NASA EXOPLANET SCIENCE INSTITUTE

For the Public PLANETQUEST













KOA



https://koa.ipac.caltech.edu/

- A collaboration between NExScI and the W.
  M. Keck Observatory.
- Public data for **all ten Keck instruments** since the Observatory saw first-light in 1994.
  - Browse-quality images of raw data.
  - Browse-quality, reduced data for HIRES, NIRC2, OSIRIS, and LWS, created by automating pipelines.
- Contributed data: Keck Observatory
  Database of Ionized Absorption toward
  Quasars (KODIAQ; N. Lehner, PI).
- Coming soon: NIRSPEC extracted spectra; moving target services.
- See Luca Rizzi's poster for more details.







- + KI Search Form
- + Contributed Datasets
- + KOA User Guide
- + Program Interface
- + Reducing Keck Data
- + FAQ
- + KOA Helpdesk
- + KOA News
- + KOA Bibliography
- + Related Archives
- + Login / Logout

#### KOA Data Access Service - v11.8.3

STATE PLANE

Real Property lies

Currently not logged in. [Log in]

+ Publicly Available Data The W. M. Keck Observatory Archive (KOA) ingests and curates data from the following instruments: DEIMOS. ESI, HIRES, KI, LRIS, LWS, MOSFIRE, NIRC, NIRC2, NIRSPEC, and OSIRIS. The schedule for public access to non-proprietary data for these and future instrument releases can be found in the description of User Access and Proprietary Periods, KOA serves data released by the Keck Observatory Database of Ionized Absorbers towards Quasars (KODIAQ), created by combining archival observations of quasars made with public HIRES data.



OSIRIS SPEC Upgrade - OSIRIS is back online for science after its SPEC detector upgrade. Interactive Visualizers are now available for calibrated NIRC2 images and reduced OSIRIS data cubes. A FITS header viewer is available from the quicklook column on the results pages.

The KOA web pages are supported under most common hardware platforms and operating systems, but on tablets such as iPads and Androids, features using mouse click-and-drag functionality are not available. Internet Explorer is not supported - we recommend Firefox or Chrome for Windows users.

Search

Reset Form

1. Choose Instrument (all modes):									
HIRES	NIRC2	NIRSPEC							
DEIMOS	ESI	LRIS							
🛛 LWS	MOSFIRE	NIRC							
OSIRIS									
	Check All	Clear All							
(To retrieve public)	Keck Interferometer data	, use the dedicated KI Search Form							



MAST



- Mikulski Archive for Space Telescopes @ STScl.
- Archive established with HST launch in 1990.
- Multi-mission since addition of IUE in 1998.
- Optical, UV, IR.
- Includes Hubble, Kepler, GALEX, IUE, FUSE, TESS, JWST, Pan-STARRS, DSS, GSC2, ...
- >700 TB of data (soon to jump to 2.5 PB with Pan-STARRS release), 2 million searches per month, 1200 refereed papers per year.



#### MAST: Barbara A. Mikulski Archive for Space Telescopes

The MAST Portal lets you search multiple collections of astronomical datasets all in one place. Use this tool to find astronomical data, publications, and images.

 Select a collection and enter a new search target OR upload an existing list

Use the filters and analysis tools to find

the exact data you're looking for. 3. Add files to the download basket to

control your download options. See the <u>User's Guide</u> for more detailed documentation and <u>video tutorials</u>.

Quick Start:

of targets.

#### Currently available data collections:

- MAST Observations: Millions of observations from Hubble, Kepler, GALEX, IUE, FUSE, and more.
- Virtual Observatory: Search thousands of astronomical data archives from around the world for images, spectra, and catalogs.
- Hubble Source Catalog: A master catalog with a hundred million measurements of objects in Hubble images.

#### Featured tutorial: Conducting a search



See all video tutorials on our YouTube channel.

MAST is managed by <u>Space Telescope Science</u> Institute. For more collections, visit <u>archive.stsci.edu</u>. Information about acknowledging the use of this resource may be found <u>here</u>.



# IACA? A atmomberaina A mahirra



## HEASARC



http://heasarc.gsfc.nasa.gov/

- High Energy Astrophysics Science Archive Research Center, @GSFC, since 1990.
- Extremely energetic cosmic phenomena ranging from black holes to the Big Bang.
- Chandra, XMM-Newton, Fermi, Suzaku, NuSTAR, INTEGRAL, ROSAT, Swift, & more than 20 others.
- Merged with Legacy Archive for Microwave Background Data Analysis (LAMBDA) in 2008 (CMBR): WMAP, COBE, ACT, etc.



**SPDF** 

SSC

data (from Xue et al. 2016, ApJS,

224, 15) is now available in Browse and Xamin. VMM Newton AO.46



## ) **astrophysics** data system



- ADS = Astrophysics Data System (@SAO/CfA).
- Indexes 12 million publications in astronomy, physics, arXiv.
- Complete coverage of astronomy and refereed physics literature.
- Tracks citations, institutional and telescope bibliographies, links to data products (back to the other archives).
- New interface and API integrating ORCID, full-text search, analytics.







## astrophysics data system

Classic Form

Modern Form

Paper Form

QUICK FIELD:	Author	First Author	Abstract	Year	Fulltext	All Search Terms	*	
Advanced -								۹

j") 🔞
") Ø

Use a classic ADS-style form

Learn more about searching the ADS Access ADS data with our API

</>

**NASA's Astrophysics Archives** 

25

## More ...

- There are other archives (based at these centers), not necessarily NASA-funded, that follow this model.
  - Ex: Pan-STARRS, VLA-FIRST @STScl
  - Ex: Palomar Oschin wide-field survey @ IRSA:
    - Zwicky Transient Facility (2017+)
    - intermediate Palomar Transient Factory (iPTF; 2013-2016)
    - Palomar Transient Factory (2009-2012)
  - There are (of course) many other non-NASA archives (SDSS, NRAO,...) and non-US archives (Simbad, ESA, ESO,...).
  - Also, observers can deliver data back to these centers for distribution (which may include data beyond original program).









## Lessons Learned ...



## Ease of Access

- Researchers at all levels (team members, emeriti, summer students) need to be able to get and use data.
- Intuitive, web-based interface.
  - No extra software installation.
  - Visualization of resources, data, tools...
  - Easy choices to "just give me the table", etc.
- Help needs to be there when users need it, easily found or promptly answered.



**NASA's Astrophysics Archives** 

## Support



- Need to have knowledgeable staff, who have done science with the data products, who can (a) find problems; (b) pass on valuable experience to new users.
- Helpdesk:
  - Speed and accuracy matters!
  - Questions can be complex.
- Documentation:
  - Tools/data releases.
  - Documentation updates in response to tickets.
  - Demos:
    - Live (AAS, ADASS, DPS, etc.).
    - Video tutorials (IRSA has >60 videos; >4500 views total).
  - The complexity of Science User needs increases with time.



# Vis

## Visualization

- Data, catalogs, plots
- What do you have?
- What do I need?
  - What I know I need...
  - What did I just find?

BEC GGGGG SPICEMBITE





## Ease of Use: High Level Science Products

- *Greatly* enhance the science return of the archives.
- Hubble Legacy high-level science products (HLSP) are used 10x as much as typical pipeline products.
- Make complex data sets accessible to a wider audience of researchers.
- Expand the use of large, coherent projects:
  - Hubble Treasury
  - Spitzer Legacy, Exploration Sci
- Generated by the community or by the center.



Spitzer/GLIMPSE

# NASA's Astrophysics Archives



# Ease of Use: Multi-Survey

- Combining information across wavelengths, surveys, missions...
- Source lists from entire missions:
  - Spitzer, Hubble, Chandra, Herschel, WISE ...





chives

## NED Science Example

- In context of assessing NED completeness, looking at fusion of GALEX, SDSS, 2MASS, WISE, ...
- Found super-luminous spiral galaxies!
- Ogle et al., 2016, ApJ, 817, 109





# **Changing Mission**

- This is a result that came from looking at what was in the archive already.
- As data get bigger and bigger, won't be able to pull data out of the archive to work with it.
- Mission evolving from "search-and-retrieve" to "do [some] analysis in situ."
- Science discoveries waiting in the archives that were never imagined or expected by the mission or even program PIs.



# Community Feedback

- Need to talk to community to find out what is needed, wanted, wished for.
- Mission members, user committees, surveys, helpdesk, talking at conferences, and review cycles all feed into setting priorities.

## **Commitment to Archival Research**

- NASA as a whole (+ sometimes missions) explicitly funds archives, and archival research: NASA ADAP (Astrophysics Data Analysis Program).
- Well-designed archive & products can greatly enhance research value of the dataset.
  - Reducing barriers: finding data, making data accessible (reliability, units, file format, artifacts, documentation).
  - NASA enables new ideas of things to do with older data.
  - NASA has strong tradition of active collaboration between missions and archives.
    - Thinking about archive during the mission!



# Ease of Use: VO

- VO= Virtual Observatory, IVOA=: International Virtual Observatory Alliance
  - Standardized *protocols* for interoperability between archives (i.e., NOT the applications that use the protocols).
  - Data discovery.
- Use interface you know, to get to data elsewhere.
- Interoperability of tools, within archives and across archives.
- (People) communication and collaboration across archives:
  - Astronomy Data Centers Executive Committee (ADEC).
  - US Virtual Observatory Alliance (USVOA).
  - NASA Astronomical Virtual Observatories (NAVO).
- NAVO:
  - Provide comprehensive and consistent access to all NASA data through VO protocols.
  - Coordinate NASA interactions with international and national VO communities.



## **IRSA VO Weekly Queries**



**NASA's Astrophysics Archives** 





NASA



# Ease of Use: API

- Application Program Interface.
  - (e.g., call by programs, scripts, command line)
- Allows scripted access to archive data.
- Enables complex projects.
- Enables rapid queries ...
- (also inadvertent DoS so need to watch and throttle!)





# **Operational Issues**

- Keep archive running while making it better.
  - Assembling the plane while in flight!
- Growing audience, usage within existing resources.
  - Be efficient in how use resources.
  - Use the same software across multiple data sets (X. Wu, earlier).

## NASA

## Innovations

- Interactive UI, using pieces developed by others.
- Machine learning.
  - NED: data embedded in free-form text, tables not standard (e.g., RA, Ra, ra, R.A., ...); pilot project to apply ML to classify data and facilitate extraction.
- Improving scalability, extensibility, data prospecting...
- Greater integration of functionality and content across systems.
  - ADS: ORCID claiming; search by object via SIMBAD TAP service; embedding of publisher images via APIs.



# Technical: Data Ingest

- Spitzer Legacy programs changed culture by requesting products be delivered back to the community.
- Now a common feature of Spitzer proposals.
- Brings these data to larger audience via central archive.
- IRSA has to have resources to ingest these products.



- Expense is not necessarily in TB but in education of the people delivering the products.
  - Need to have delivery <u>well-organized and documented</u>.
  - People who have done it a lot: easy.
  - People new at this: not necessarily easy.
  - Tools to help people new at this.
  - Complexity is not just about size!
- Can end up with, e.g., optical and UV data available through the Spitzer archive (SINGS, LVL).



# What's Next: Big Data

- "Big data" in some missions, certainly "big data" across all NASA Archives.
- To some extent, have already been working with "big data"!
- IRSA has already invested in data visualization services to help people identify and experiment with data quickly.
- Planning: identifying the most critical needs of users, including increased analysis at the archive facilitated by user workspaces.
- Richer services for in situ analysis.
- All archives thinking about this in some way.



## Summary

- Long-term, stable archives greatly increase the return on observatory investment. (Doubles papers!)
- Robust support for both expert and novice users pays off.
- User support by instrument experts is crucial.
- Standardization of tools within an archive increases efficiency.
- Interoperability between archives benefits everyone.
- High level data products can expand the reach of large data sets.
- Shift in approach from "search and retrieve" to "analyze in situ".