

Big Data Initiatives at MAST

A 5-Year Plan to Enhance MAST Capabilities for the Community

Marc Postman

Space Telescope Science Institute

June 29, 2016



Archive for Space Telescopes



AST is a NASA astrophysics data archive center

Archive established with HST
launch in 1990

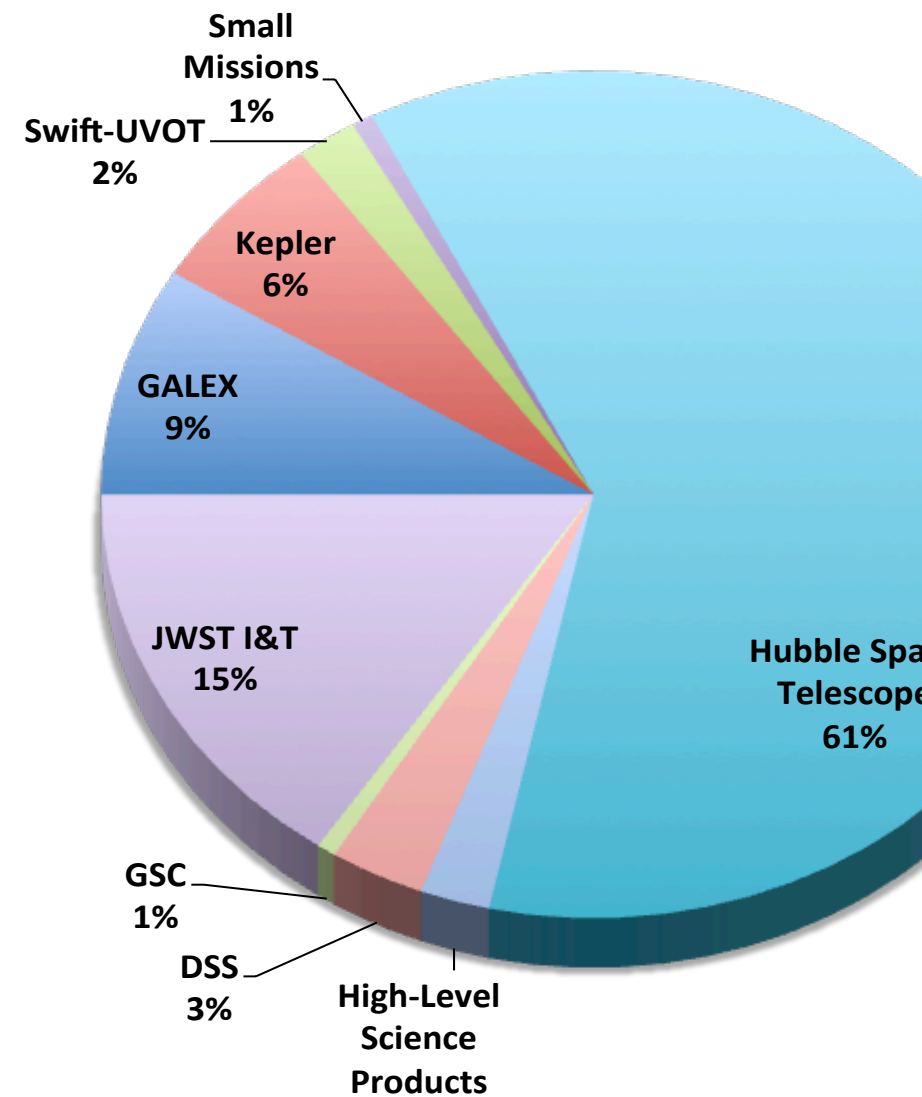
Multi-mission since addition of IUE
in 1998

active missions including
Hubble, Kepler

Many legacy missions: GALEX,
IUE, FUSE, ...

future: TESS, JWST, WFIRST

includes some relevant ground-
based archives: GSC, PanSTARRS*



MAST data are diverse

ST data are diverse:

data types: images, spectra, light curves, catalogs, models

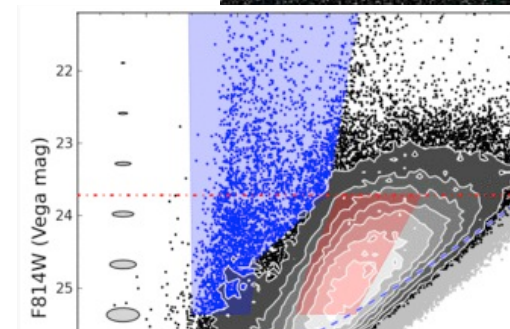
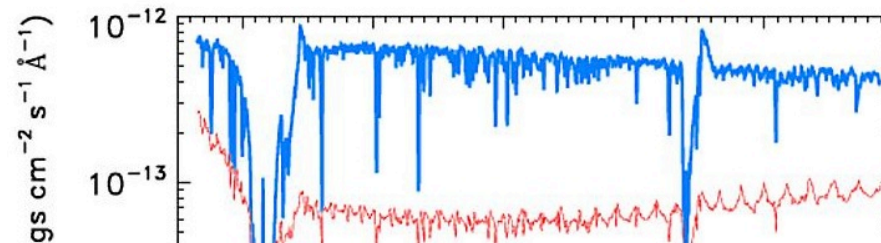
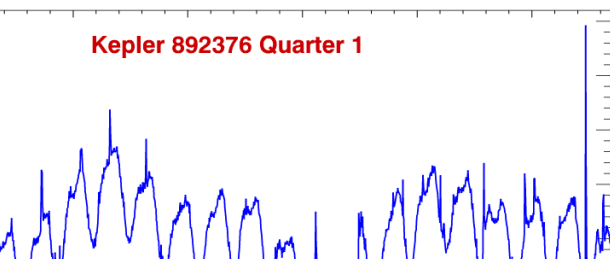
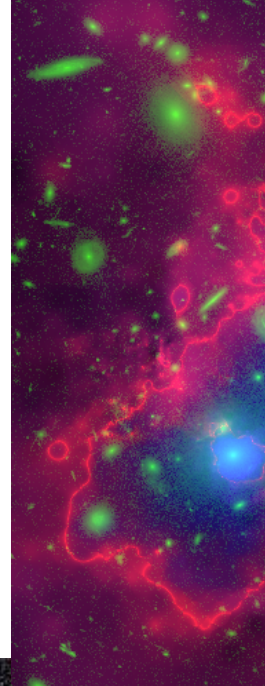
scale: from Pan-STARRS (2 PB images + 100 TB database) to small shuttle-based missions

community-contributed projects with just a few files

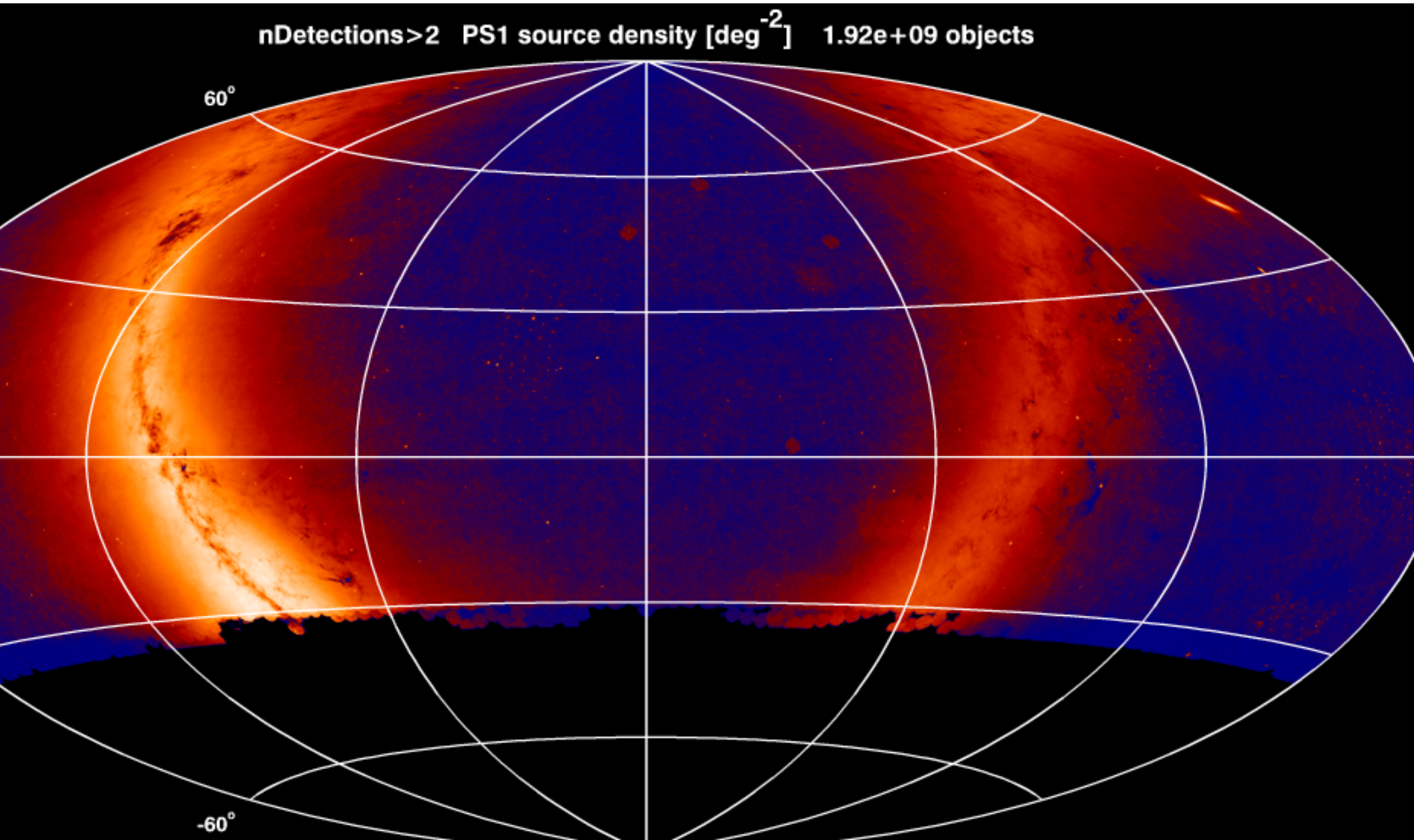
processing level: from raw data packets delivered directly from spacecraft to science-ready, high-level science products

many different missions and instruments

Hubble alone: 12 different instruments, 17 varieties of detector, hundreds of instrument modes/filters/etc.



...ive users are getting more sophisticated every year, and their queries (and analy...
...easingly crossing over archive and wavelength boundaries. By the time WF...
...ched this will be even more so, thus **the MAST archive must evolve wi**
...iderations and capabilities in mind.



What we mean by Big Data

“Big Data” is a term borrowed from industry, which collects data on a massive scale from users and sensors, and is typically used to serve targeted advertisements with a higher return on investment than would otherwise be possible.

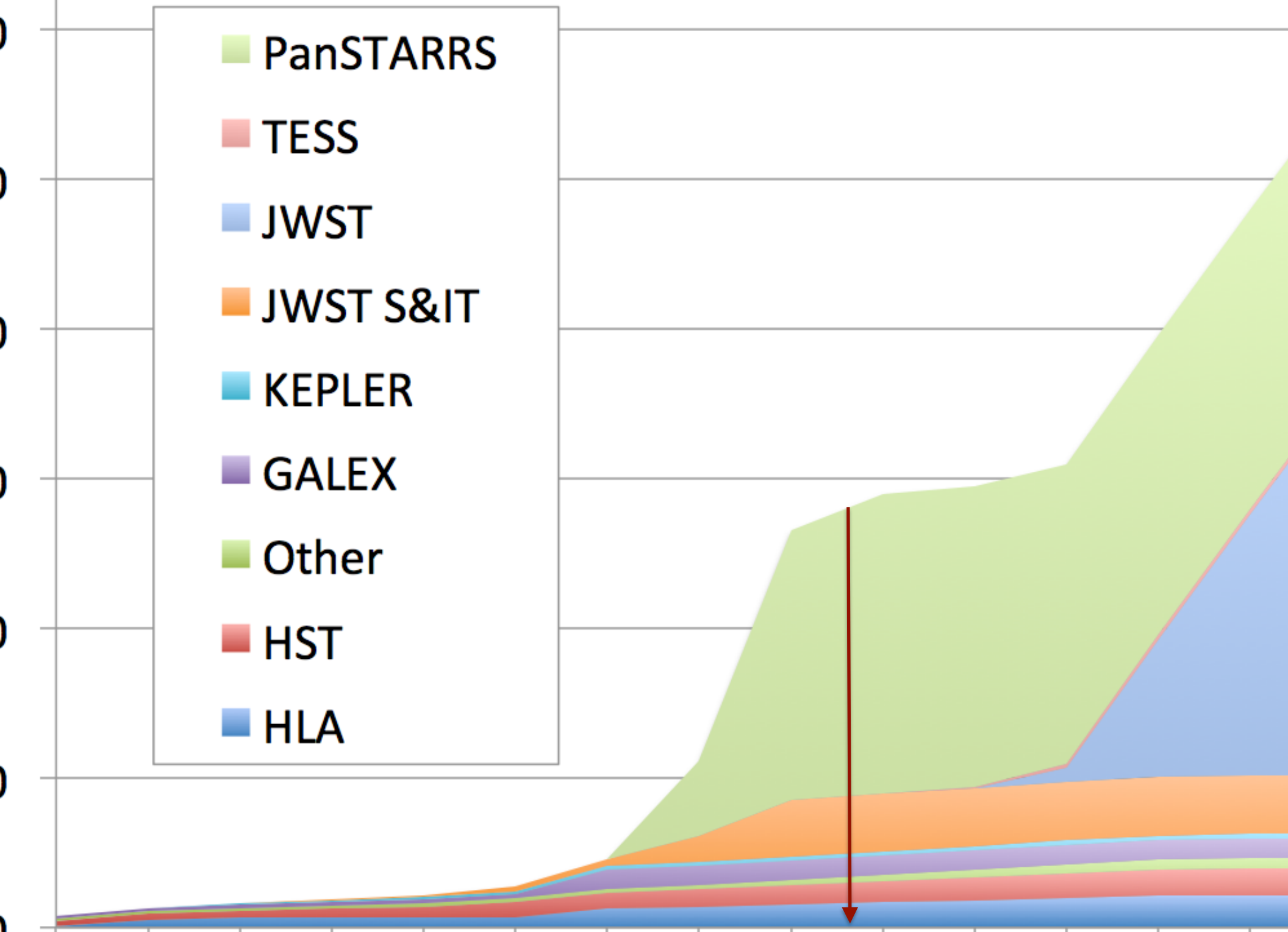
At MAST, we redefine it for astrophysics applications as data that meets one or more of the following criteria:

Data whose raw form is so large that we must qualitatively change the way in which we reduce, store, and access it.

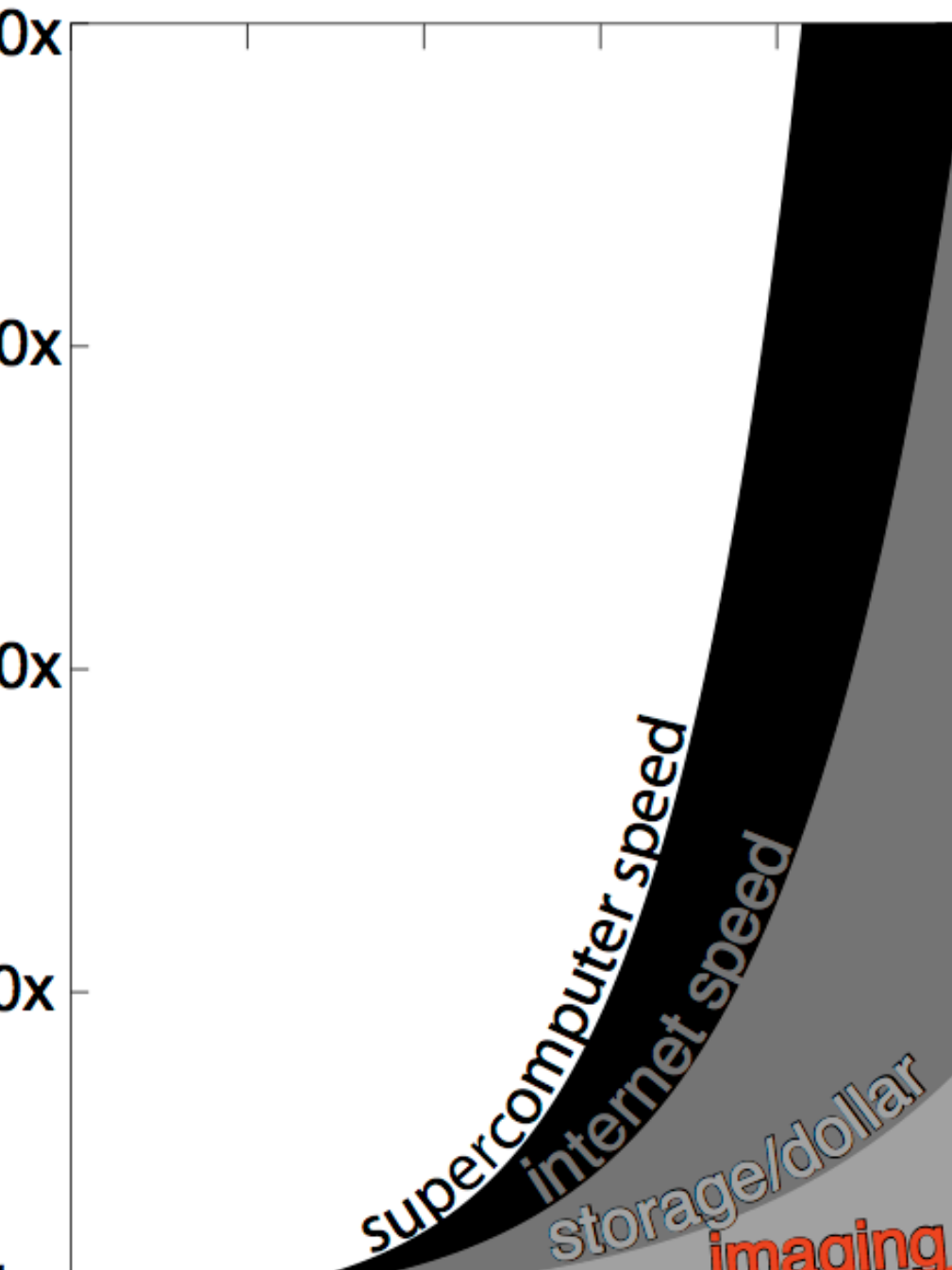
Data whose reduced form is so large that we must qualitatively change the way in which we interact with and explore it.

Data whose structure is so complex that our current tools cannot efficiently extract the scientific information we seek.

Past and Projected Data Volume



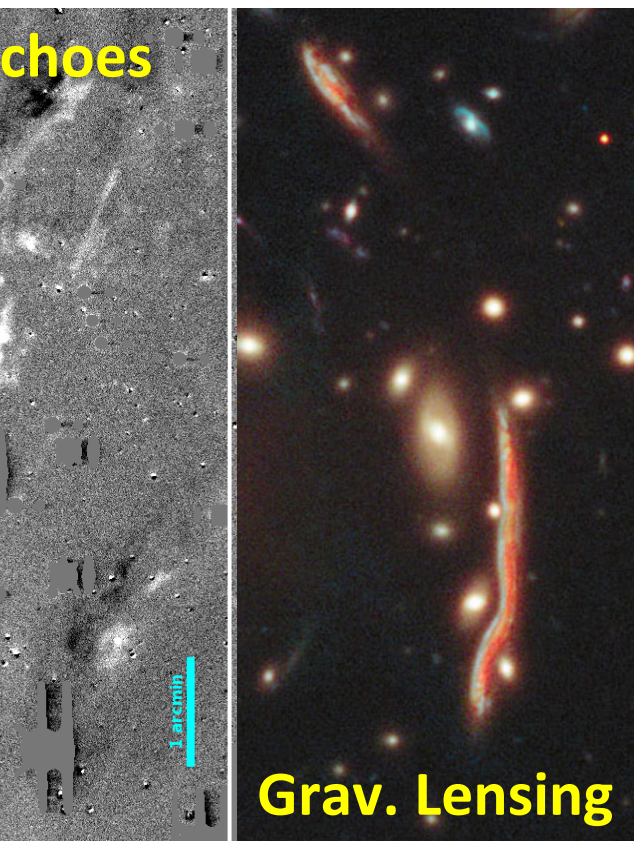
er 10 years, the rate of UVOIR imaging typically grows by a factor of 10x.
In contrast, the amount of storage one can buy, internet speed, and supercomputer speed grow by factors of dozens to hundreds.



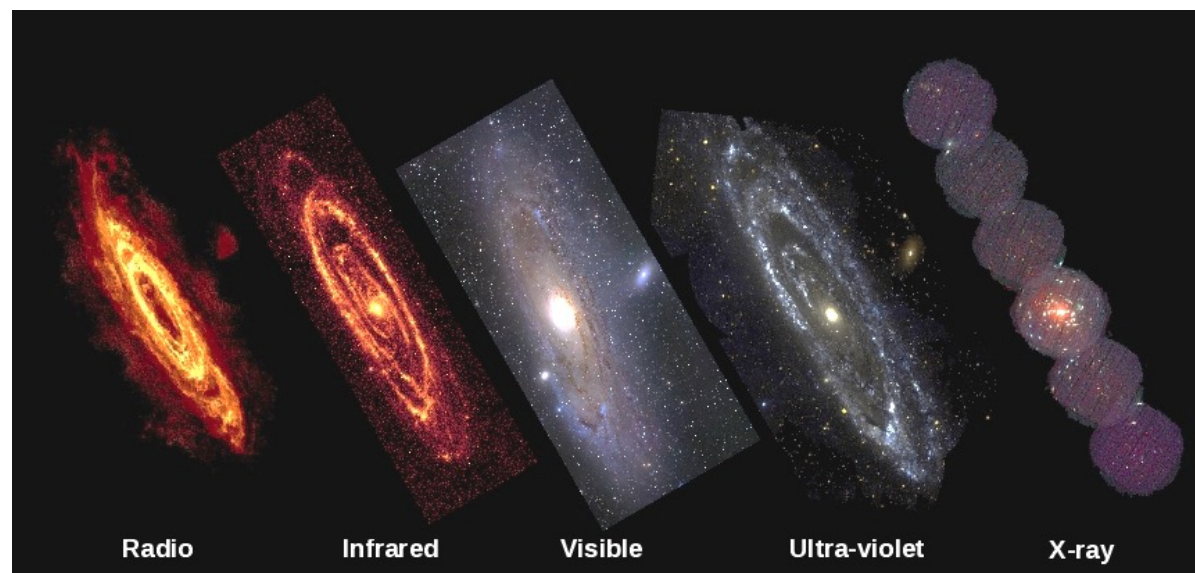
Data source: Barentsen & Peek,
<https://github.com/barentsen/tech-progress-data>

It is not our computers that keep up, but our brains: our ability to explore, comprehend, assimilate, and infer won't keep up with the flow of data unless we build sophisticated tools for interacting with and communicating the ideas hidden in these enormous datasets.

Automated source classification

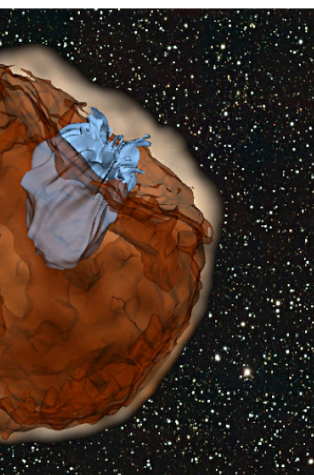


Multi-wavelength data collection and cross-correlation

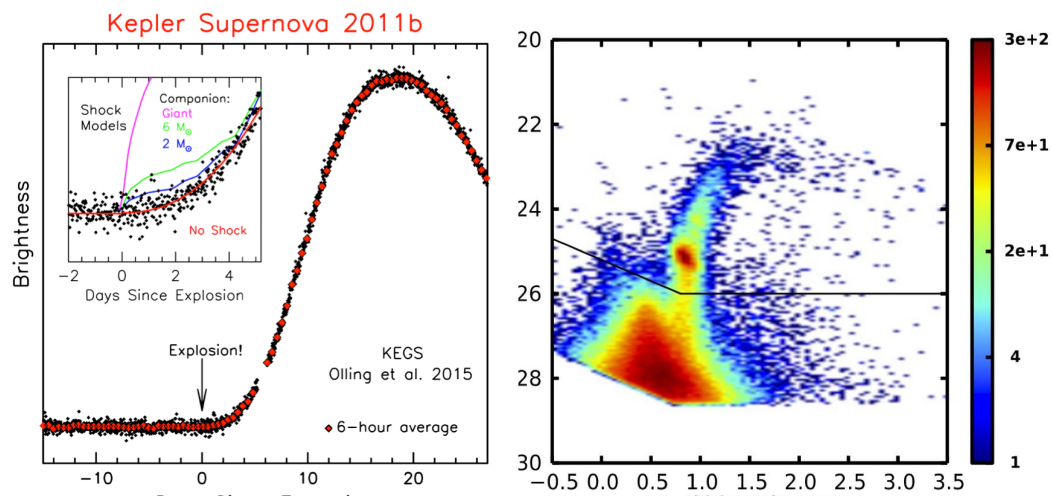


Stellar Populations, Galaxies, Active Galactic Nuclei, Galaxy Photometric Redshifts

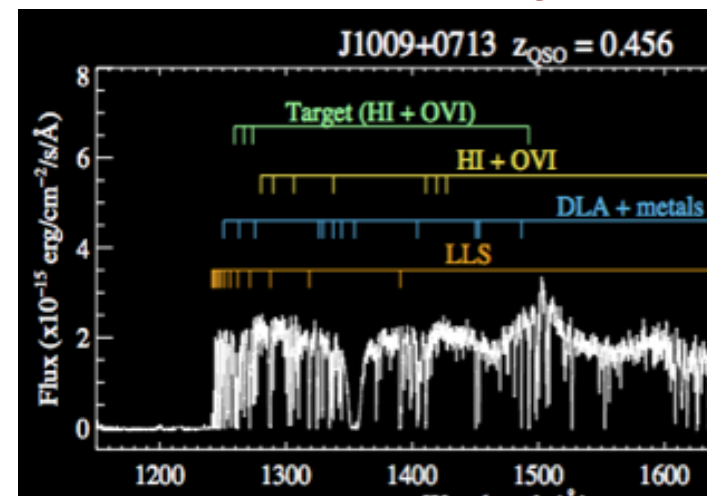
Main analyses



Model fitting to large or complex datasets



Disentangling datasets into constituent systems



Advances

Derive Star Formation Histories for 100 Nearest Galaxies (photometry required for $10^9 - 10^{10}$ stars).

Galaxy – Black Hole Co-Evolution Study in Population of 10^7 Galaxies

Automated detection and classification of light echoes (requires time domain data). Light echo structure is complex and sources are rare.

Automated detection and classification of amorphous astronomical structures (lensed galaxies, gaseous filaments, etc.). Source structure is complex.

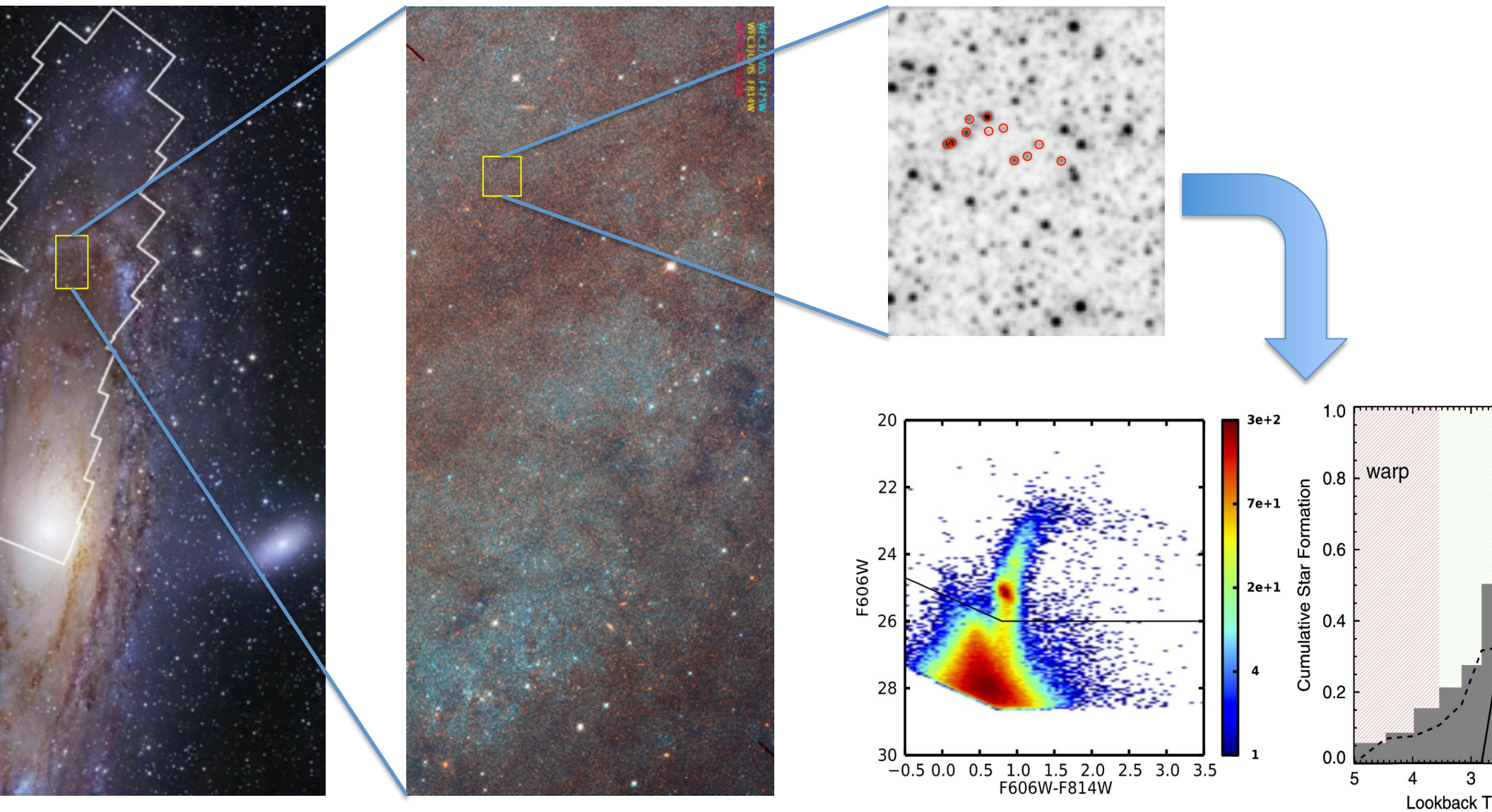
Determine 3D distribution of 10^9 galaxies based on photometric redshifts derived from multi-wavelength data.

Automated Classification of light-curve data from Kepler/K2/TESS (hundreds of data points for 10^5 stars).

Probabilistic Classification of Transient sources via cross-matching with large multi-wavelength catalogs.

Identification of spectroscopically rare sources via PCA template matching

Tracking Star Formation Histories



quality becomes better, larger fraction of users can use automated path for their science.

data

ive – processed

pipeline or HLA)

ive – raw data

3. Do photom (& completeness)

- Use HLA/HSC (DAOphot, SExtractor)
- Use other existing catalog (e.g., HLSP)
- Use SExtractor yourself
- Use DAOphot yourself
- Use other photometry package yourself
- **Completeness - not included in HLA/HSC**
- **Completeness - limited public tools exist**

5. SED fitting of single objects

- **Limited public tools exist**
- TADA
- CHORIZOS (limited support)
- BEAST (not public)
- Use own software

7. Analysis/V

- Use Discovery
- CasJobs - **limi**
- IDL
- Pyraf (IRAF)
- DS9
- Glue (3D visu
- Use own soft

Typical project timeline

e images

strodizzle)

rizzle yourself

software

4. Crossmatching data (if needed)

- Use HSC
- Use Discovery Portal to do crossmatch
- Use other tools (e.g., IRAF/TMATCH, DAOPHOT/match, ...)
- Use own software

6. SFH fitting of populations

- **Limited public tools exist**
- StarFISH (public – outdated)
- MATCH (not public)
- Use own software

8. Simulation

- **Limited pub**
- DAOphot/ac
- TFIT - Match
- resolution
- Sunrise (n-b
- Use own so

Compute Servers for Science:

Science servers, ranging from 8 to 32 cores per server. RAM: 32 to 400 GB.

Utilization of science servers is managed via informal communications. No checking for resource conflicts (e.g., running cpu-intensive jobs simultaneously).

Science servers cannot be easily reconfigured to match needs of different kinds of applications.

TScl has aggregate of ~2,000 cores but they cannot be utilized for large-scale computing due to lack of virtualized environment.

Current Large Mission systems:

HST mission systems highly physical (non-virtualized) and fixed capacity though use of High Throughput Computing (HTC) via HTCondor does provide job scheduling and resource allocation for the Hubble Data Management System.

WST mission systems are being implemented as virtualized from ground up. Some in-house applications (like DMS) that are primarily CPU-bound may stay non-virtualized.

bandwidth:

Mix of 1 Gb/s and 10 Gb/s internal networks.

Upgrades to 10 Gb/s are made during normal refresh cycles or earlier, if needed for specific engineering requirements (e.g. Connection to storage systems)

Connection to internet: 500 Mb/s

Connection to internet2: 1 Gb/s (as of Dec. 2015)

Plans being made to upgrade standard internet connection to 1 Gb/s in short term (<12 months?) and 10 Gb/s in longer term.

Storage:

~8 PB of storage using mix of NAS and SAN systems

Largest single dedicated system is for PanSTARRS
(~5 PB) followed by MAST (1.2 PB).

Advanced Software and Analysis Tools

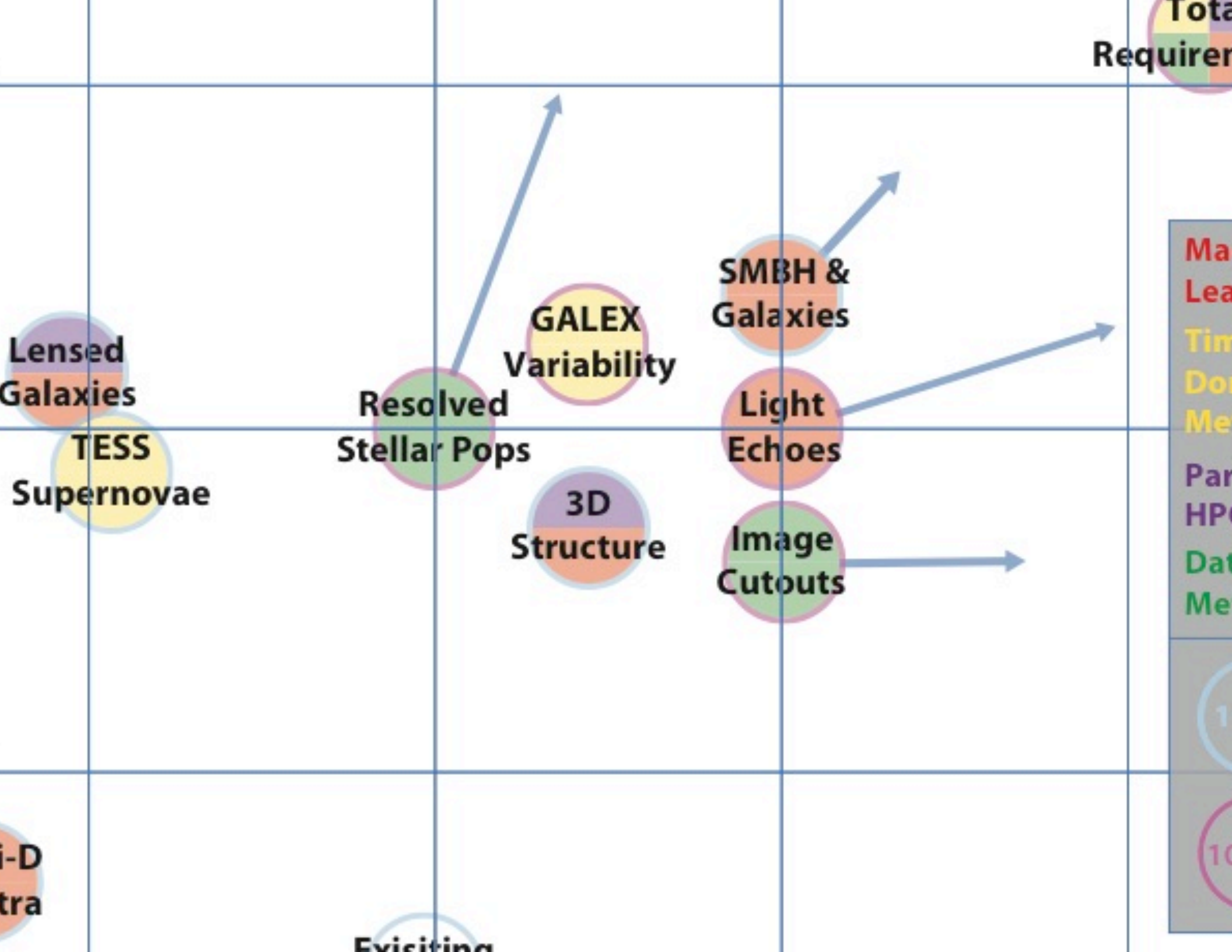
Archive users are getting more sophisticated every year, and their queries (and analyses) are increasingly crossing over archive and wavelength boundaries.

Archives of the near future will be more than a simple massive file or database servers. Users will rather interact with them algorithmically through well-defined and well-designed Applications Programming Interfaces (APIs).

Cloud computing (as well as dedicated local clouds) should be ubiquitous in the next 5 to 10 years.

Given the amount of astronomical data in MAST in near future, we expect that server-side analyses will be commonplace for the users, thus an advanced scripting capability must be supported.

Science Case	Computing	Storage	Bandwidth	Software
Exoplanet detection classification (PS1, etc.)	>1000 cpu core	~1 to 10 PB	~10 Gb/s needed for citizen science/ data transfer.	Machine Learning (ML) classification, feature v (FV) mapping.
Galaxy redshift (HST, JWST, etc.)	~1000 cpu core	~10 TB (if image cutout service avail.)	~1 Gb/s	Parallelizable code management, ML code
Galactic stellar population studies (HST, WFIRST)	>1000 cpu core (10K core better)	Few hundred TB	~10 Gb/s (uses non-local data)	Automated pipeline, ef database query tools
AGN / Galaxy evolution studies	>2000 cpu core	>1 PB	>1 Gb/s	ML classification / FV mapping / efficient x-co
Structure in	~500 cpu core	>200 TB	>1 Gb/s	Parallelizable code; hig connectivity to other ar
Supernovae	~1000 cpu core	~10 TB for raw data	~1 Gb/s (if data local)	Transient detection cod highly undersampled d
Inter-visit photometry study	Several hundred cpu core	Few hundred TB	High bw to local db (>10 Gb/s)	Transient detection cod capable of running on 1 curves.
High-dimensional cosmological datasets	Not a driver	~100 GB	Current OK	FV mapping, ML codes, Efficient x-corr, PCA and



Recommendations: Infrastructure

Significant improvements in internal and external network bandwidths are needed. External bw: 10 Gb/s asap; 100 Gb/s by 2018-2019. Internal bw: 100 Gb/s by 2018, 1000 Gbps by 2021.

Multiple virtualized, dynamically configurable nodes with upwards of 1000 cores each are needed in next 2 years. Multiple 10K core systems will likely be needed in 5 years.

Storage: Increase storage capacity by 6 PB by end of 2018; increase to 30 – 50 PB by 2021.

Recommendations: System Design

Data workspaces needed to bring “scientists to the data.”

Workspace must support a variety of tools, scripting languages and software.

Must provide adequate computing resources (cpu, memory, storage)

Must have good internal bandwidth between workspace, data sources, and cpu cores.

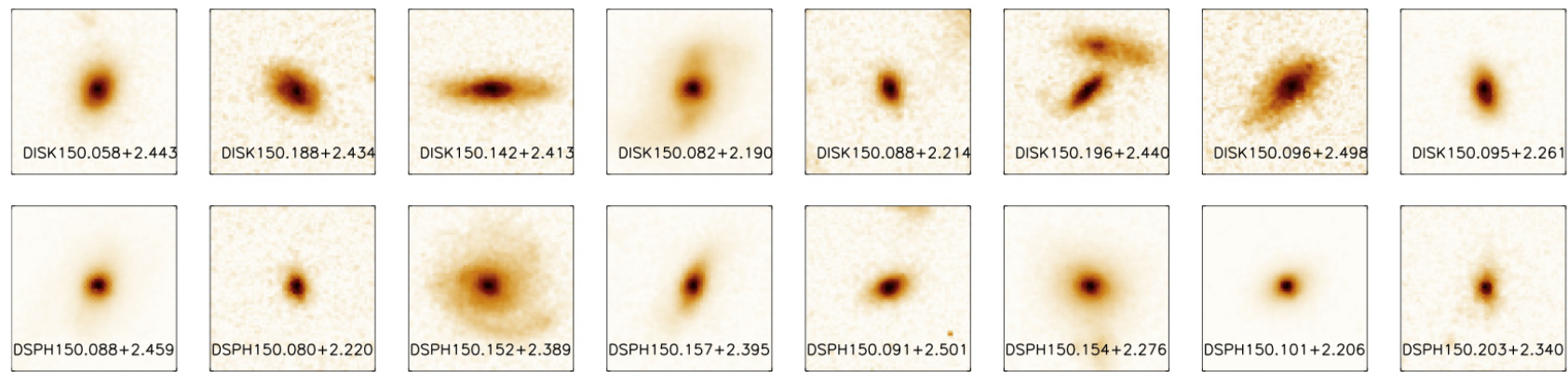
In 5 years timeframe, STScI will need to invest in hybrid cloud (Flexible datacenter/private + public cloud/partner institution [e.g., JHU/MARCC = Mid-Atlantic Regional Computing Center])

Recommendations: Science Software

Develop integrated data visualization tools that run on the server-side to enable researchers to efficiently explore all our data holdings (e.g., *GLUE*).

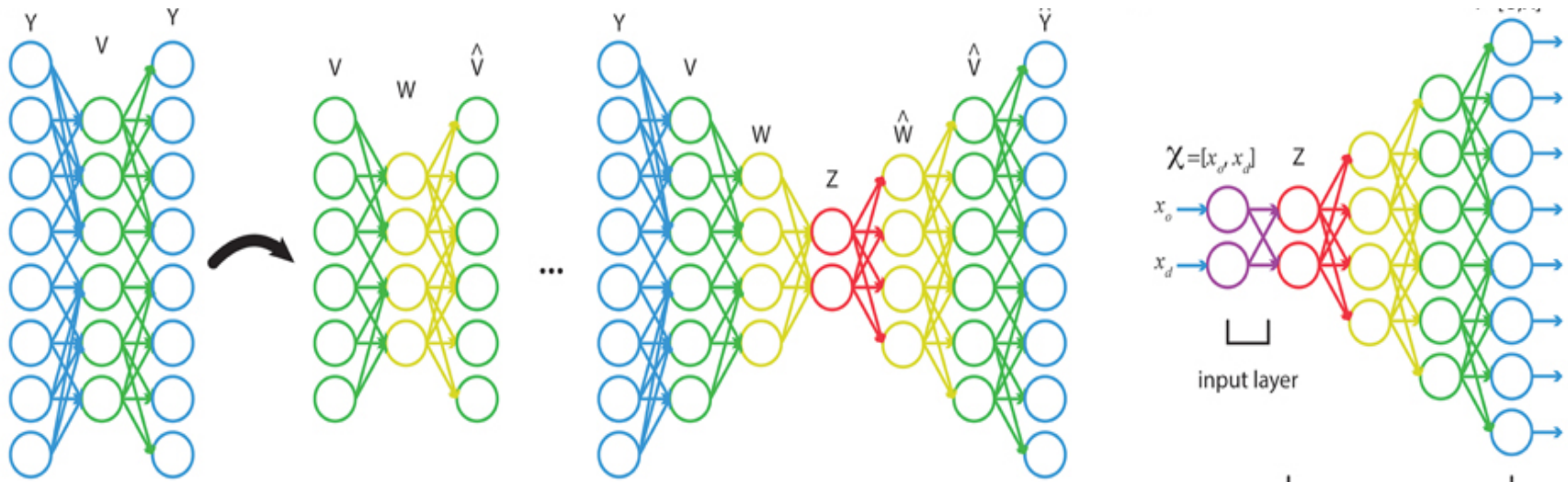
Explore, select and deploy machine-learning architectures to support archive researchers who will, more frequently, require advanced classification and regression analyses.

Initiate development of automated spectral feature classification tools. The software should be adaptable to data from all current and future spectroscopic data in MAST.



Server-Side Scripting Machine Learning (ML) & visualization

- Allows users to refine (in real-time) their queries and generate new hypotheses without having to download, query, investigate, and repeat.
- On-the-fly ML clustering and dimensionality reduction of search results can generate easy to investigate data visualizations

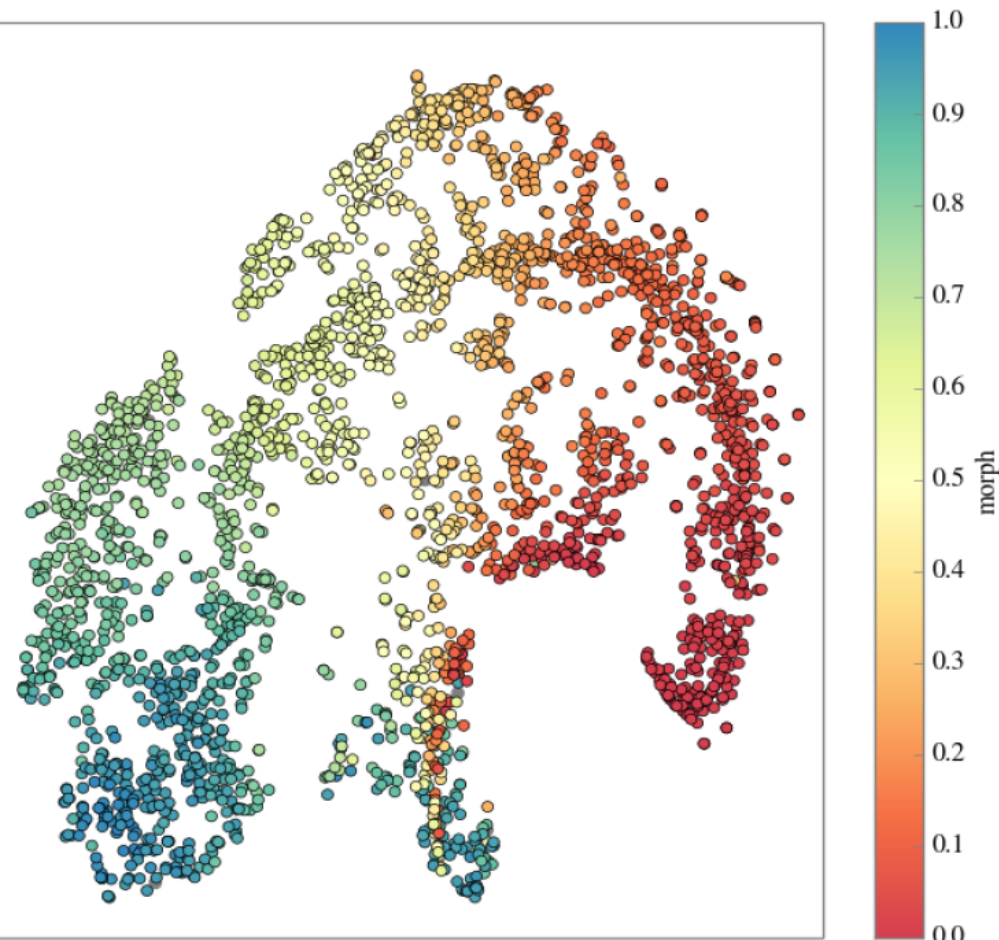


Deep auto-encoder dimensionality reduction: unsupervised learning methods take very high dimensional data and reduce them to 2 or 3 key dimensions

Server-Side Scripting Machine Learning and visualization

- t-SNE attempts to cluster objects into like groups. If we can enable users to perform this clustering and dimensionality reduction “on the fly” on search results, it will allow users to further refine their searches and do more advanced hypothesis generation before downloading any large data sets.

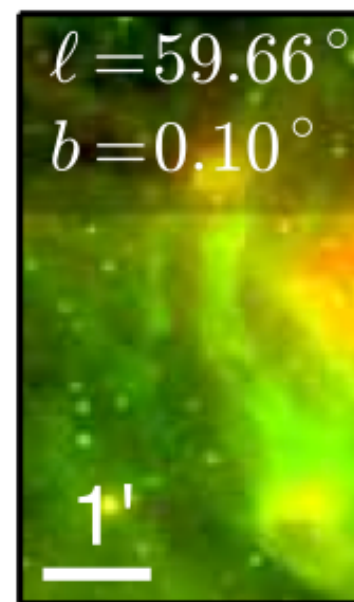
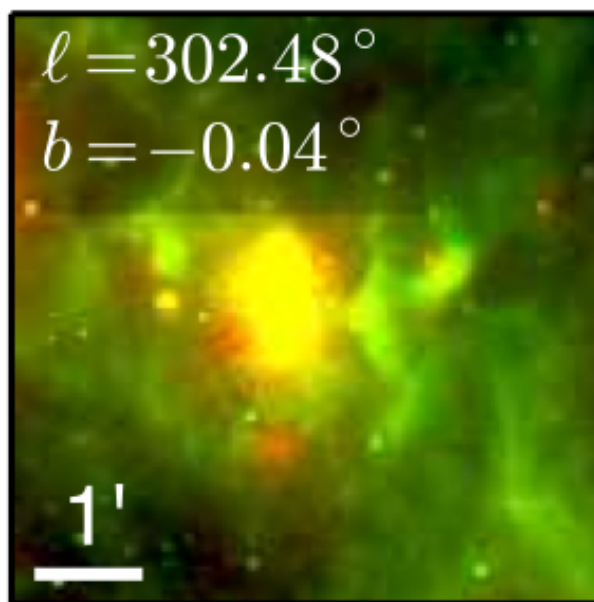
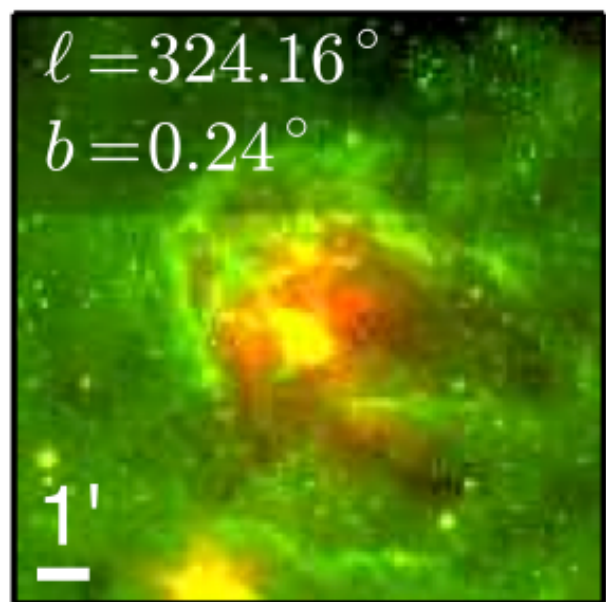
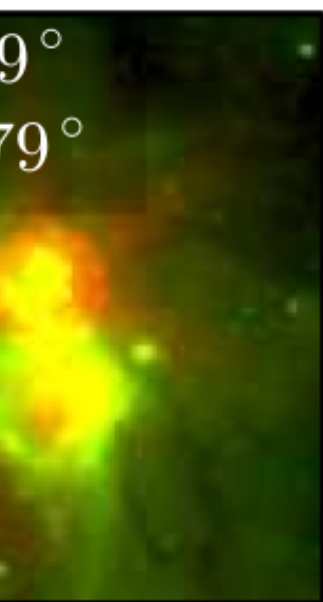
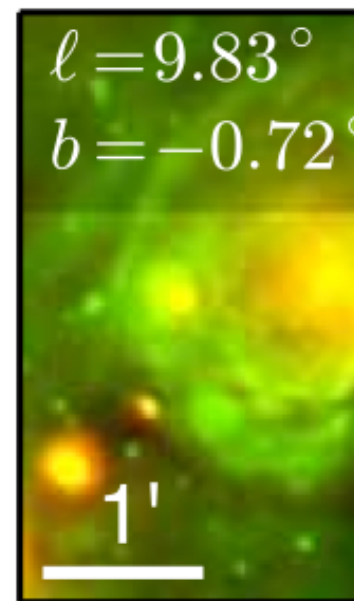
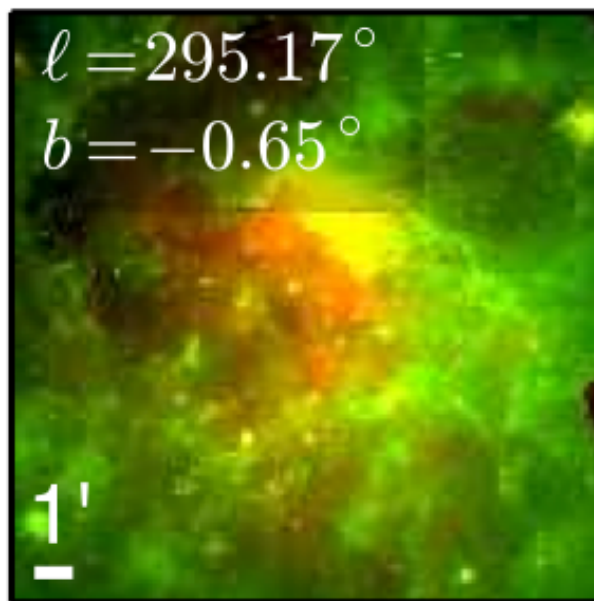
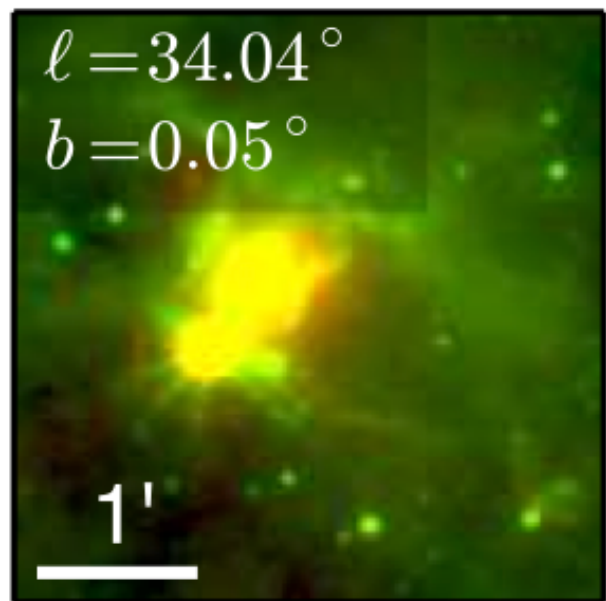
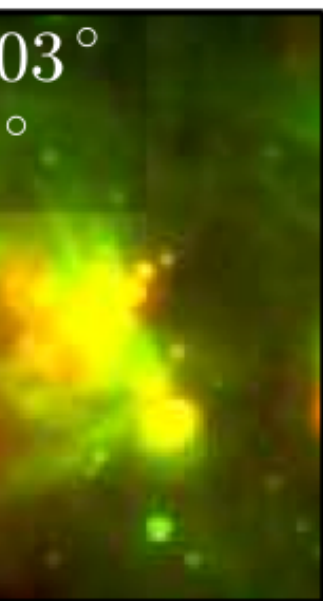
<http://keplerebs.villanova.edu/tsne>



t-SNE clustering applied to
Kepler eclipsing binary star data.
(Matijevic et al. 2012)*

**t-Distributed Stochastic Neighbor Embedding*

Ionizing bubbles in interstellar gas (due to star forming regions)



The ML algorithm “Brut” was trained on citizen science classifications

Recommendations: Skill sets

initiate the ***Barry M. Lasker Data Science Postdoctoral Fellowship***: focused specifically on early-career researcher whose interests and expertise are centered on big data science.

- Program has been initiated. First fellow selected (41 applied).

Recommendations: Skill sets

Expand staff expertise with server-side API' and scripting environments.

Expand staff expertise in the areas of machine-learning and automated classification techniques.

Expand staff expertise in the development and application of code to reduce the dimensionality of highly complex datasets.

Recommendations: Organizational Structure

consolidate leadership and oversight of our current and upcoming science archive initiatives (HST, Kepler, TESS, JWST, MAST, PanSTARRS, WFIRST) with goal of maximizing science potential, supported with state-of-the-art tools and interfaces.

Establish a Data Science Mission Office (DSMO) at STScI.

Archive project scientist and Archive project technologist will be resident in DSMO.

- Project scientist responsible for developing and maximizing scientific impact of STScI archives.
- Project technologist responsible for developing and prioritizing utilization of appropriate technologies to support archival research and big data science.

Staff members from other missions and divisions will be matrixed with DSMO to support its activities.



Big Data @ STScI

Enhancing STScI's Astronomical Data Science
Capabilities over the Next Five Years

$$\begin{aligned} & + \gamma(\mu + \beta) \frac{\partial I_v}{\partial r} \\ & \frac{\partial}{\partial \mu} \left\{ \gamma(1 - \mu^2) \left[\frac{1 + \beta\mu}{r} - \gamma^2(\mu + \beta) \frac{\partial \beta}{\partial r} \right. \right. \\ & \left. \left. - \gamma^2(1 + \beta\mu) \frac{\partial \beta}{\partial t} I_v \right] \right\} - \frac{\partial}{\partial v} \left\{ \gamma v \left[\frac{\beta(1 - \mu^2)}{r} \right. \right. \\ & \left. \left. - \gamma^2\mu(\mu + \beta) \frac{\partial \beta}{\partial r} + \gamma^2\mu(1 + \beta\mu) \frac{\partial \beta}{\partial t} I_v \right] \right\} \\ & \left\{ \frac{2\mu + \beta(3 - \mu^2)}{r} + \gamma^2(1 + \mu^2 + 2\beta\mu) \frac{\partial \beta}{\partial r} \right. \\ & \left. - \gamma^2[2\mu + \beta(1 + \mu^2)] \frac{\partial \beta}{\partial t} \right\} I_v = v_{\text{sc}} - \gamma_{\text{sc}} I_{\text{sc}} \quad (1) \end{aligned}$$

Science Definition Team Report

February 2016

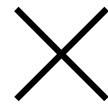
Discovery in Astronomically Big Data

An AURA (STScI – NOAO – LSST) workshop and hands-on exploration

February 27 - March 2, 2017 at STScI

Astro Domains

Galaxy Beyond 3D	Big Spectroscopy
Galaxies & the Time Domain	Morphology



Data Methods

Citizen Science	Data Int
Machine Learning	Visu

Goals

Learn about new scientific discoveries in each astronomy domain

Acquire skills for discovery methods in current datasets

Discuss and reduce cultural barriers to work in the discovery space

Develop partnerships for common visualization and exploration tools