

# Big Data in astronomy: an overview

SOMACHINE 2020

Federica B. Bianco

University of Delaware

Physics and Astronomy

Biden School of Public Policy and Administration

Data Science Institute

NYU Center for Urban Science and Progress

Rubin Observatory LSST Science  
Collaborations Coordinator

Rubin LSST Transients and Variable Stars  
Science Collaborations Chair

this slide deck:

<https://slides.com/federicabianco/bdastro>

1. historical perspective: BD context
2. BD from astronomical surveys
  - optical
  - radio
  - gravitational waves
3. space-based astronomy BD problem
4. crowdsourcing approach
5. time domain astronomy BD problems
6. platforms
  - computational platforms in astro
  - computational platforms in other fields
  - VO

# 1/6

## *Historical perspective*

# Historical perspective

"Data larger than can be analyzed with typical tool"

"Data that stresses the infrastructure"

"Data that does not fit in memory"

# Historical perspective

"Data larger than can be analyzed with typical tool"

*John R. Mashey Chief Scientist, SGI, mid-1990s*



# Historical perspective

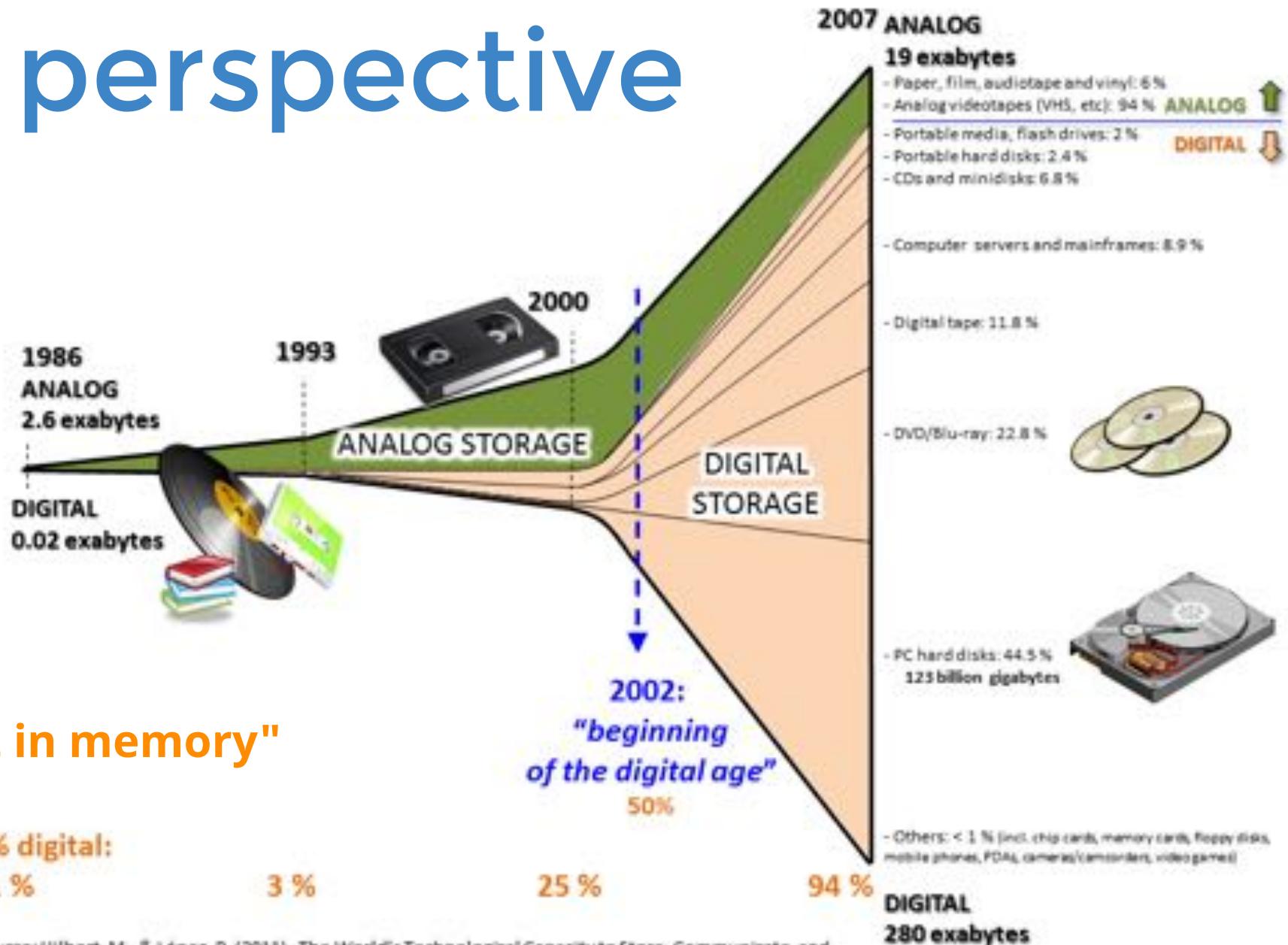
"Data larger than can be analyzed with typical tool"

*John R. Mashey Chief Scientist, SGI, mid-1990s*

"Data that stresses the infrastructure"



# Historical perspective



"Data that does not fit in memory"

% digital:

1 %

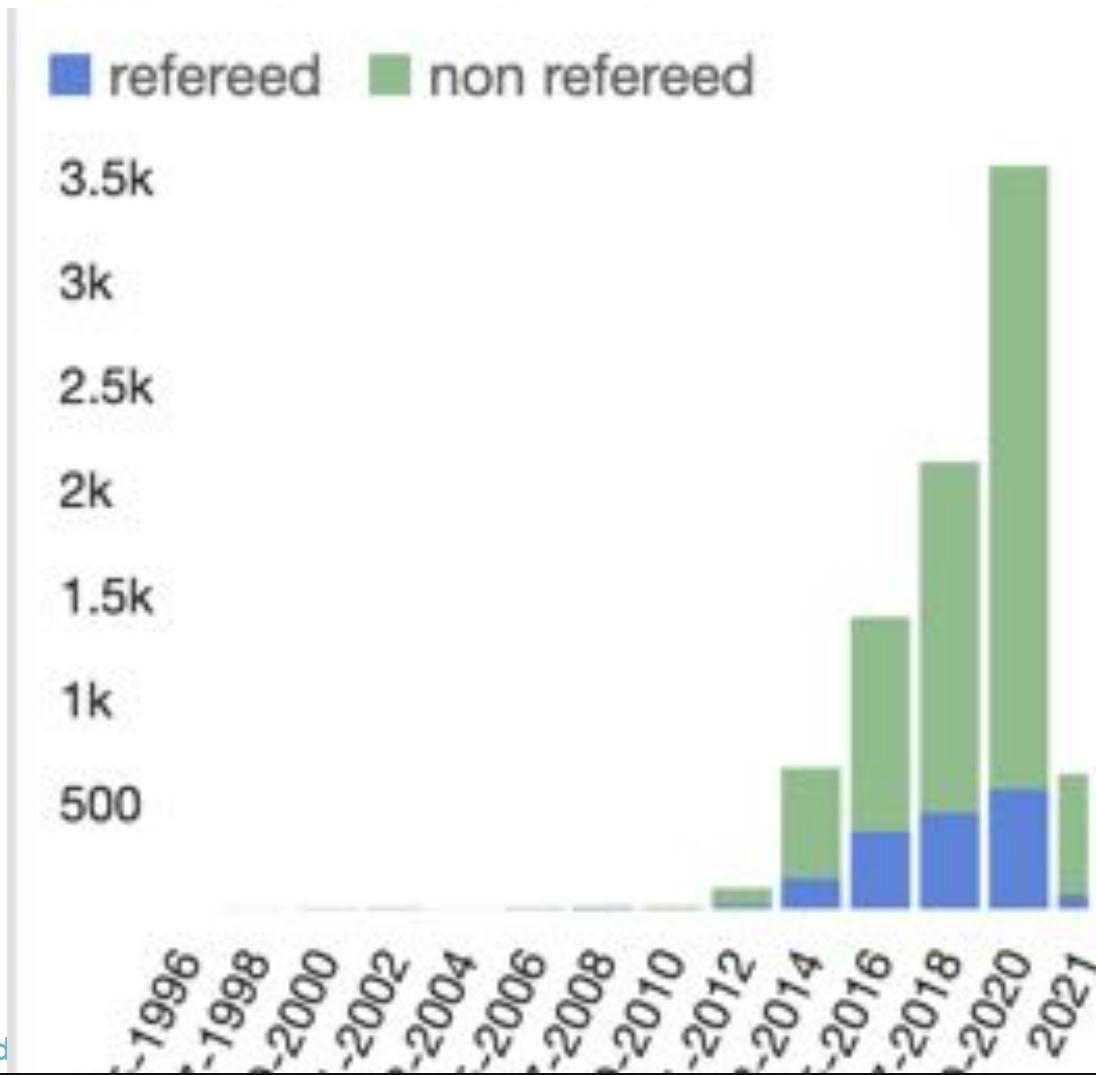
3 %

25 %

94 %

Source: Hilbert, M., & López, P. (2011). The World's Technological Capacity to Store, Communicate, and Compute Information. *Science*, 332(6025), 60–65. <http://www.martinhilbert.net/WorldInfoCapacity.html>

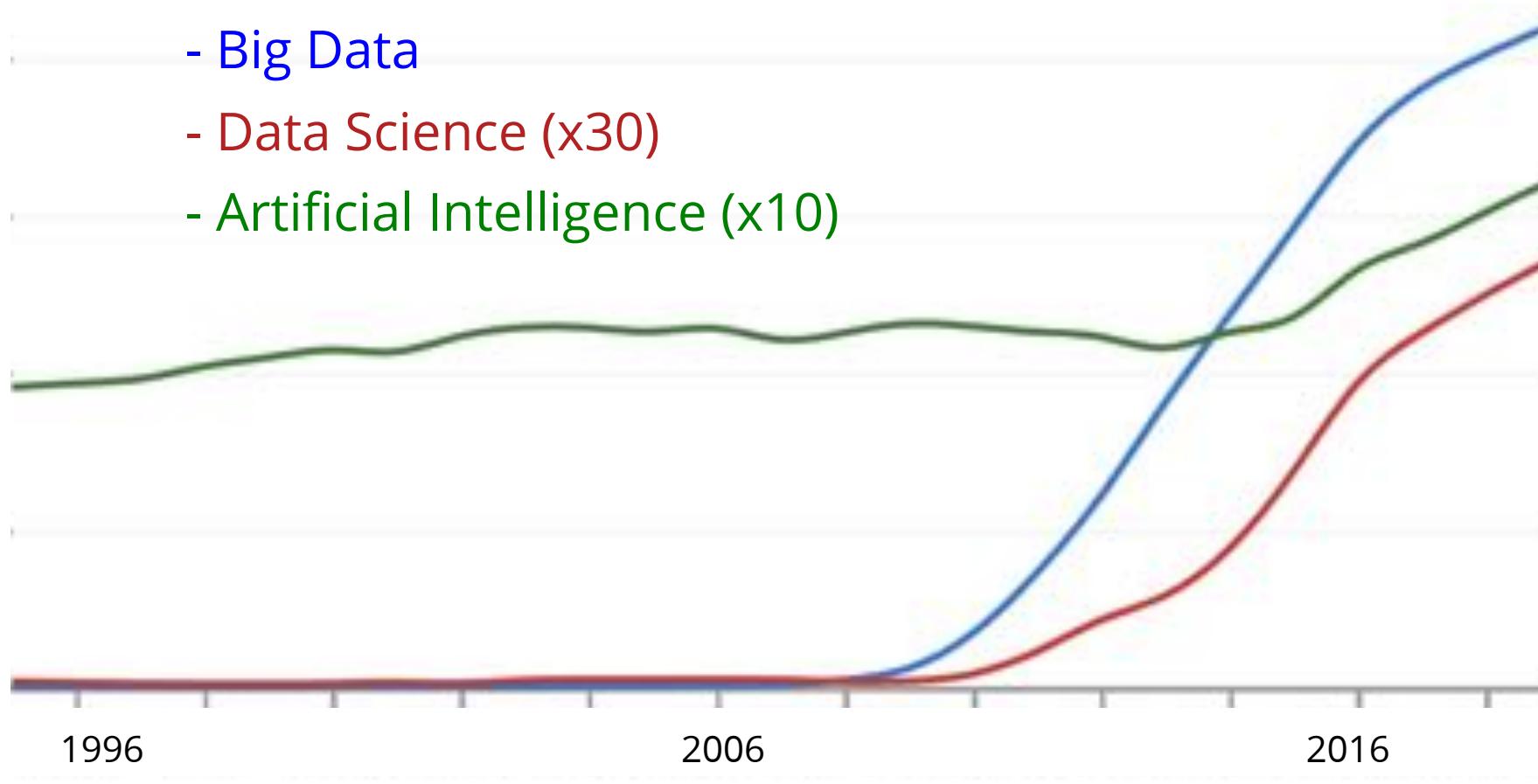
# Historical perspective



Big Data in astronomy papers  
(source: ADS)

# Historical perspective

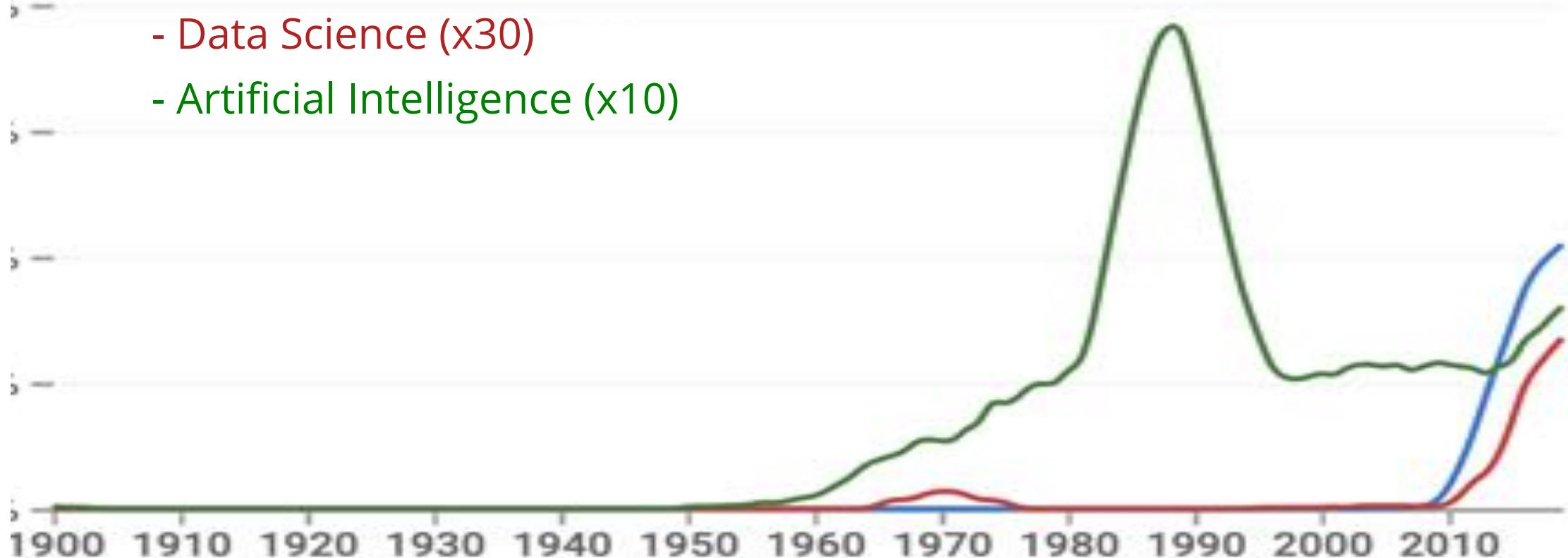
- Big Data
- Data Science (x30)
- Artificial Intelligence (x10)



occurrence of term in Google-books corpus <https://books.google.com/ngrams>

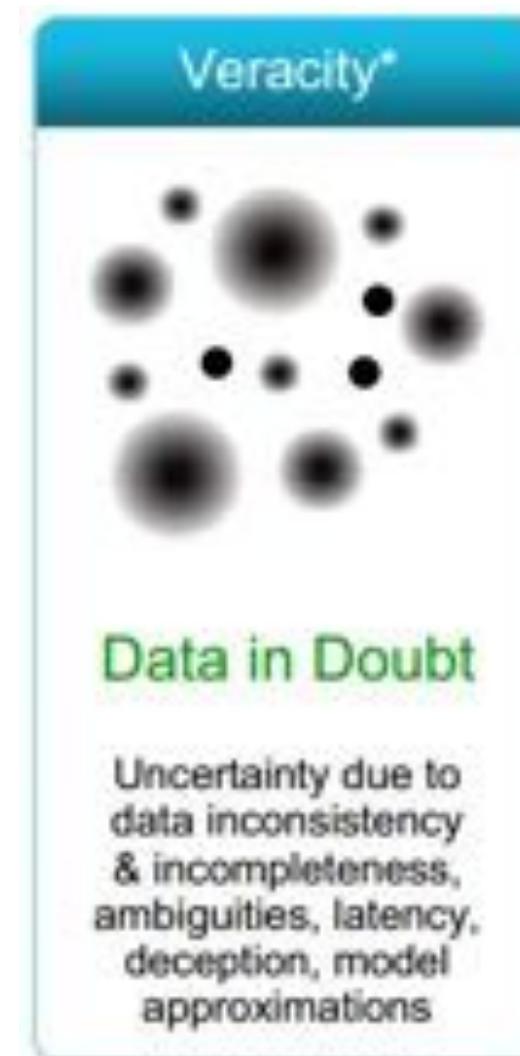
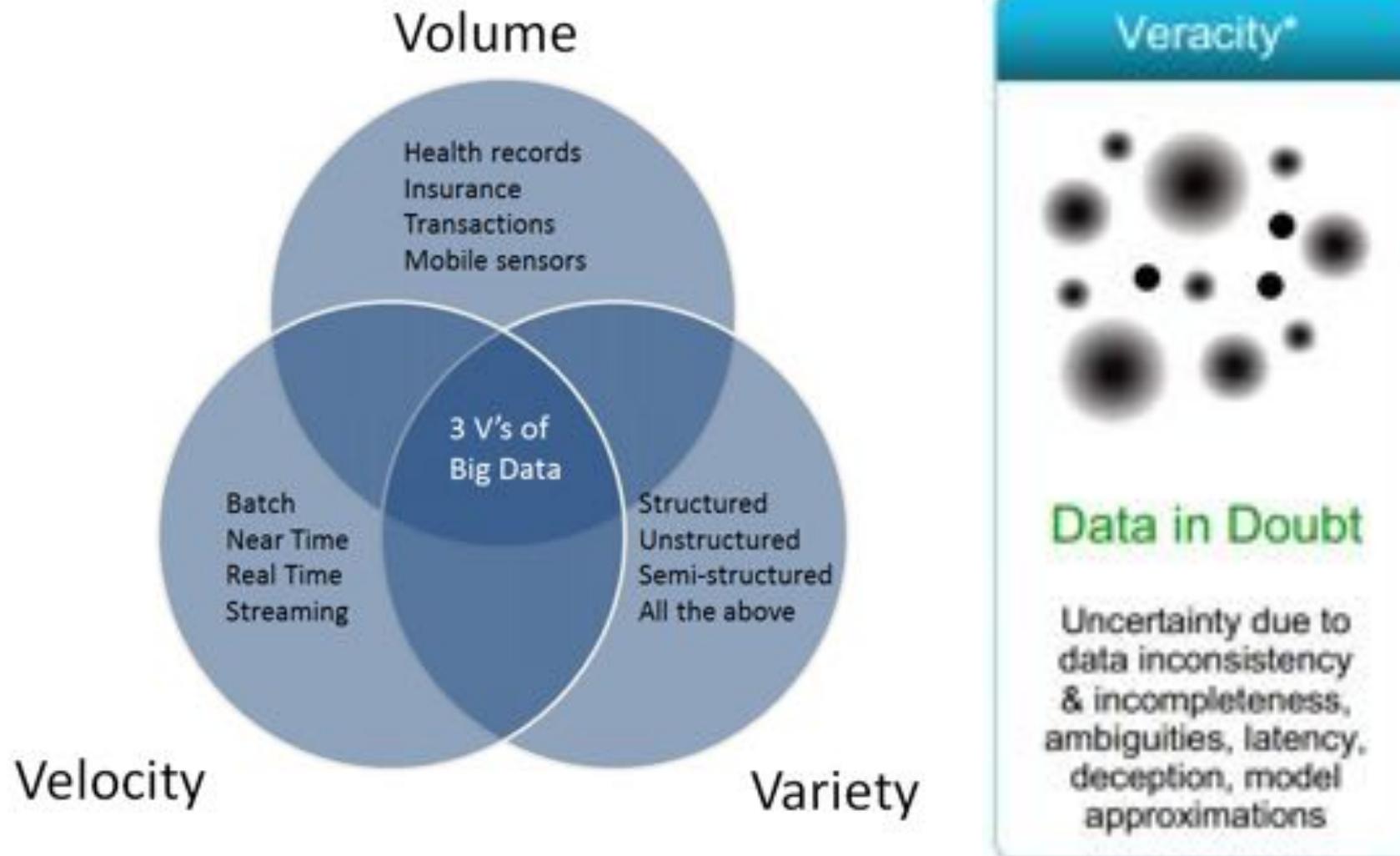
# Historical perspective

- Big Data
- Data Science (x30)
- Artificial Intelligence (x10)

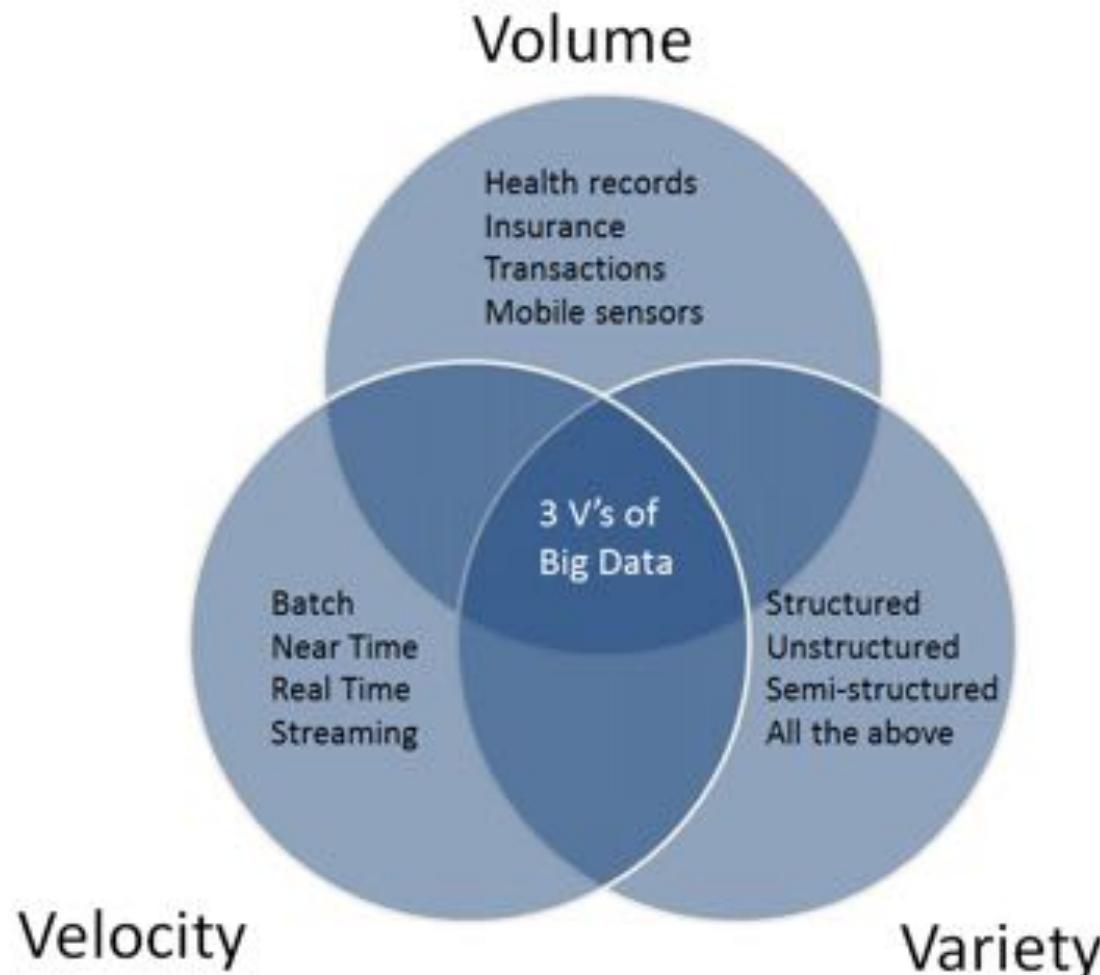


occurrence of term in Google-books corpus <https://books.google.com/ngrams>

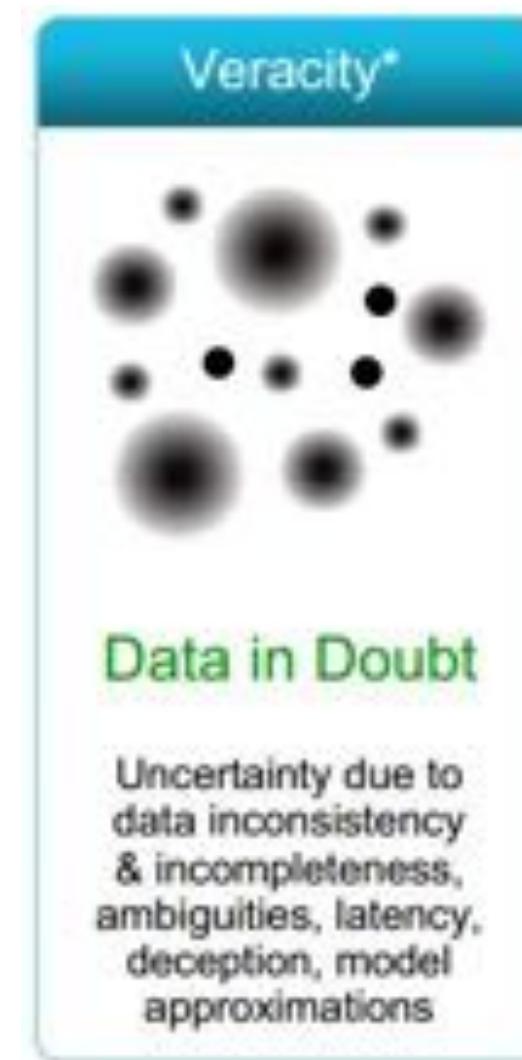
# Historical perspective



# Historical perspective



*Gartner report 2001*



Michael Walker on 28 November 2012

# 4-V of Big Data

## V1: Volume

Number of bites

Number of pixels

Number of rows in a  
data table x number  
of columns for  
catalogs

## V2: Variety

Diverse science return  
from the same dataset.

Multiwavelength  
Multimessenger

Images and spectra

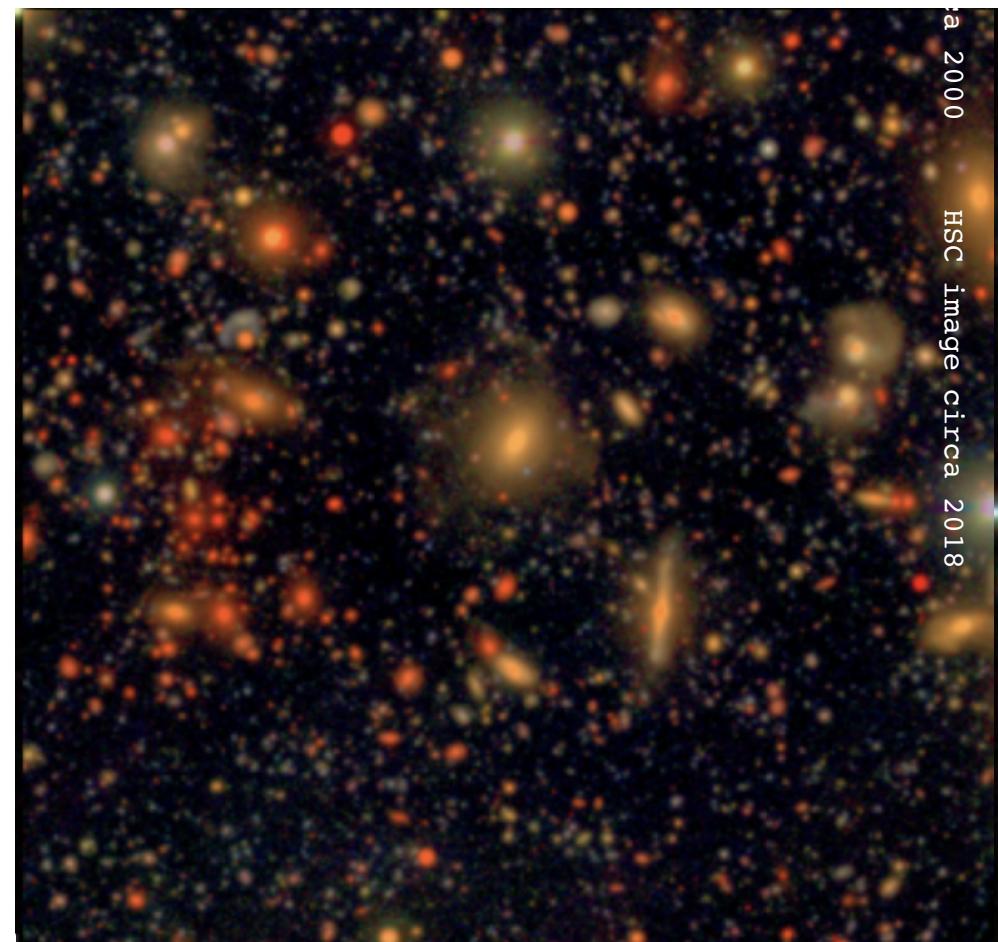
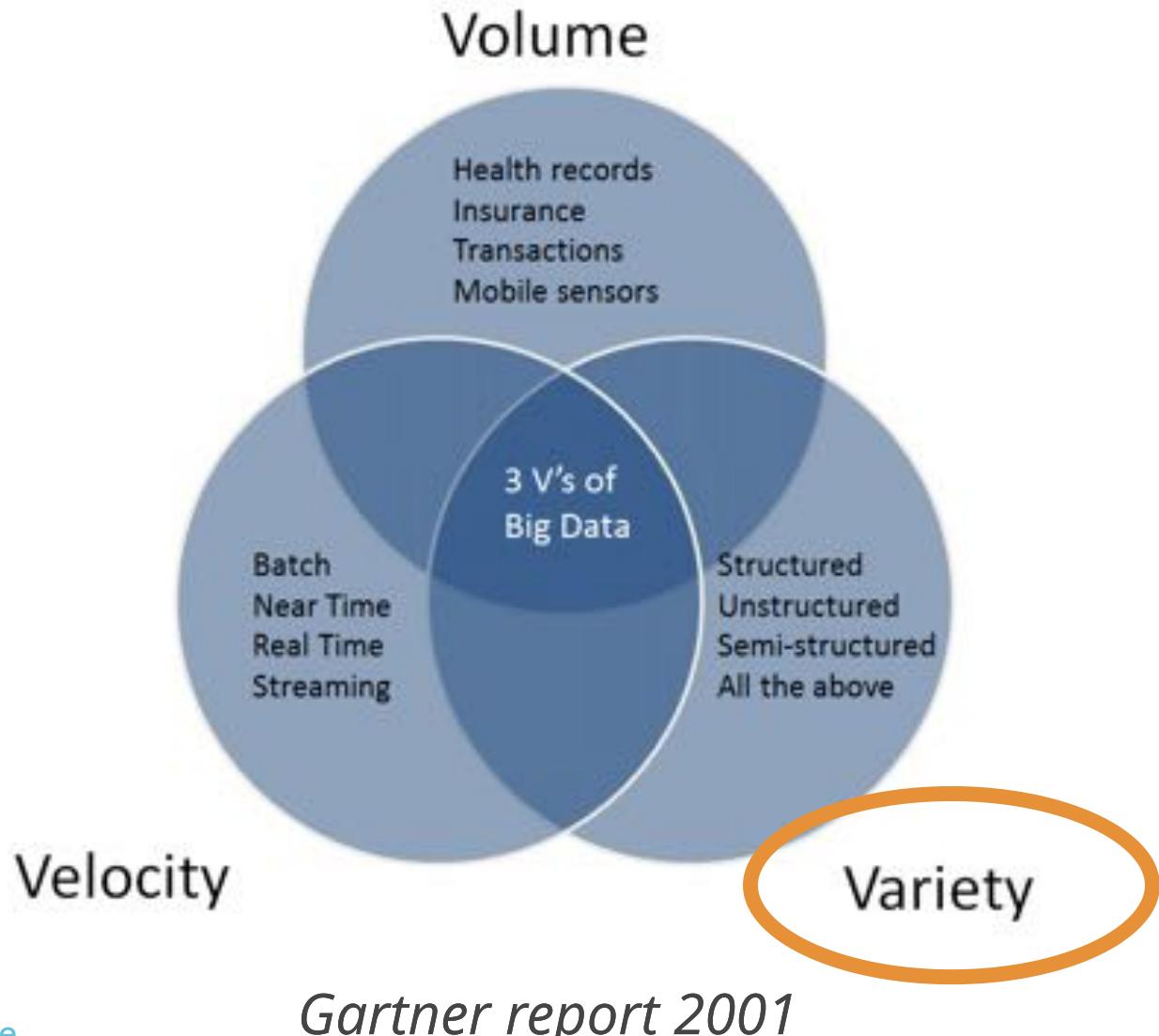
## V3: Velocity

real time analysis,  
edge computing,  
data transfer

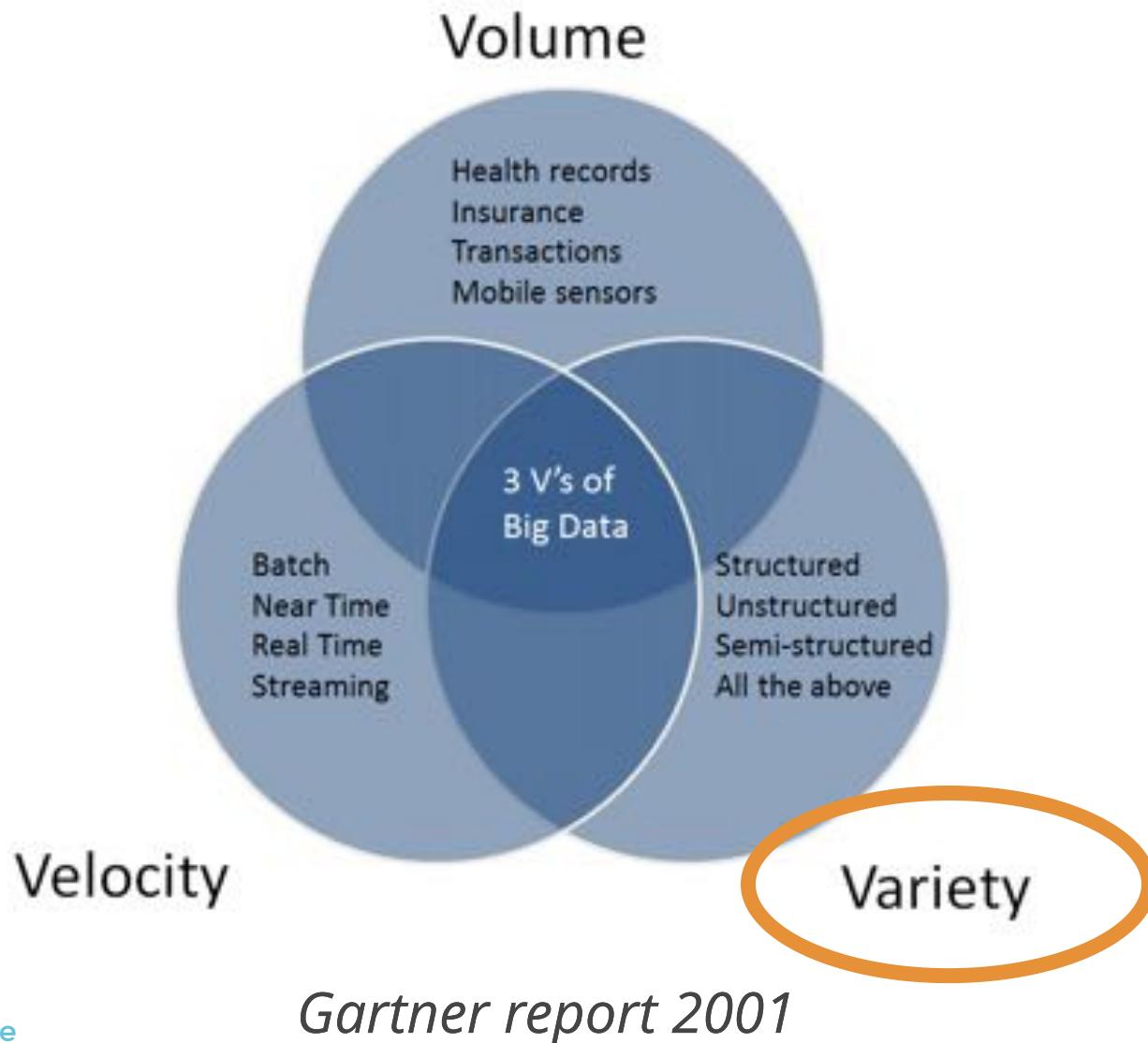
## V4: Veracity

This V will refer to  
both data quality  
and availability  
(added in 2012)

# Historical perspective



# Historical perspective



complexity

Astronomical data mainly include ***images spectra, time-series data, and simulation data.***

Most of the data are saved in catalogues or databases. The ***data from different telescopes or projects have their own formats***, which causes difficulty with integrating data from various sources in the analysis phase. In general, ***each data item has a thousand or more features***; this causes a large dimensionality problem. Moreover, data have many data types: structured, semi-structured, unstructured, and mixed.



# Rubin Observatory LSST

next talk by M Rawls!

each night is 20TB data,  
coming in at the rate of  
30GB per minute



# Big Data: Astronomical or Genomic?

Zachary D. Stephens, Skylar Y. Lee, Faraz Faghri, Roy H. Campbell, Chengxiang Zhai, Miles J. Efron, Ravishankar Iyer, Michael C. Schatz , Saurabh Sinha , Gene E. Robinson 

Published: July 7, 2015 • <https://doi.org/10.1371/journal.pbio.1002195>

Data Phase	Astronomy	Twitter	YouTube	Genomics
Acquisition	25 zetta-bytes/year	0.5–15 billion tweets/year	500–900 million hours/year	1 zetta-bases/year
Storage	1 EB/year	1–17 PB/year	1–2 EB/year	2–40 EB/year
Analysis	In situ data reduction	Topic and sentiment mining	Limited requirements	Heterogeneous data and analysis
	Real-time processing	Metadata analysis		Variant calling, ~2 trillion central processing unit (CPU) hours
	Massive volumes			All-pairs genome alignments, ~10,000 trillion CPU hours
Distribution	Dedicated lines from antennae to server (600 TB/s)	Small units of distribution	Major component of modern user's bandwidth (10 MB/s)	Many small (10 MB/s) and fewer massive (10 TB/s) data movement

[doi:10.1371/journal.pbio.1002195.t001](https://doi.org/10.1371/journal.pbio.1002195.t001)

# what drives scientific discovery in astronomy



# what drives astronomy

**Experiment driven**

Following: Djorgovski

<https://events.asiaa.sinica.edu.tw/science/20170904/talk/djorgovski1.pdf>

Observations Details				
2. Mar. 1610	2. P. Gassendi	March H. 12	O	**
3. Mar.	30. March		** O	*
2. Apr.	2. April		O	*** *
3. Apr.	3. April		O	* *
3. Apr.	3. April	H. 13.	* O	*
9. Apr.	9. April		* O	**
6. May	6. May		** O	*
8. May	8. May	H. 13.	* * * O	
10. May	10. May		*	* * O *
11.			*	* O *
12. May	12. May	H. 4 night	*	O *
13. May	13. May		*	* O *
14. May	14. May		*	* * O *

Galileo Galilei 1610

# what drives astronomy

*Experiment driven*

*Theory driven | Falsifiability*



# what drives astronomy

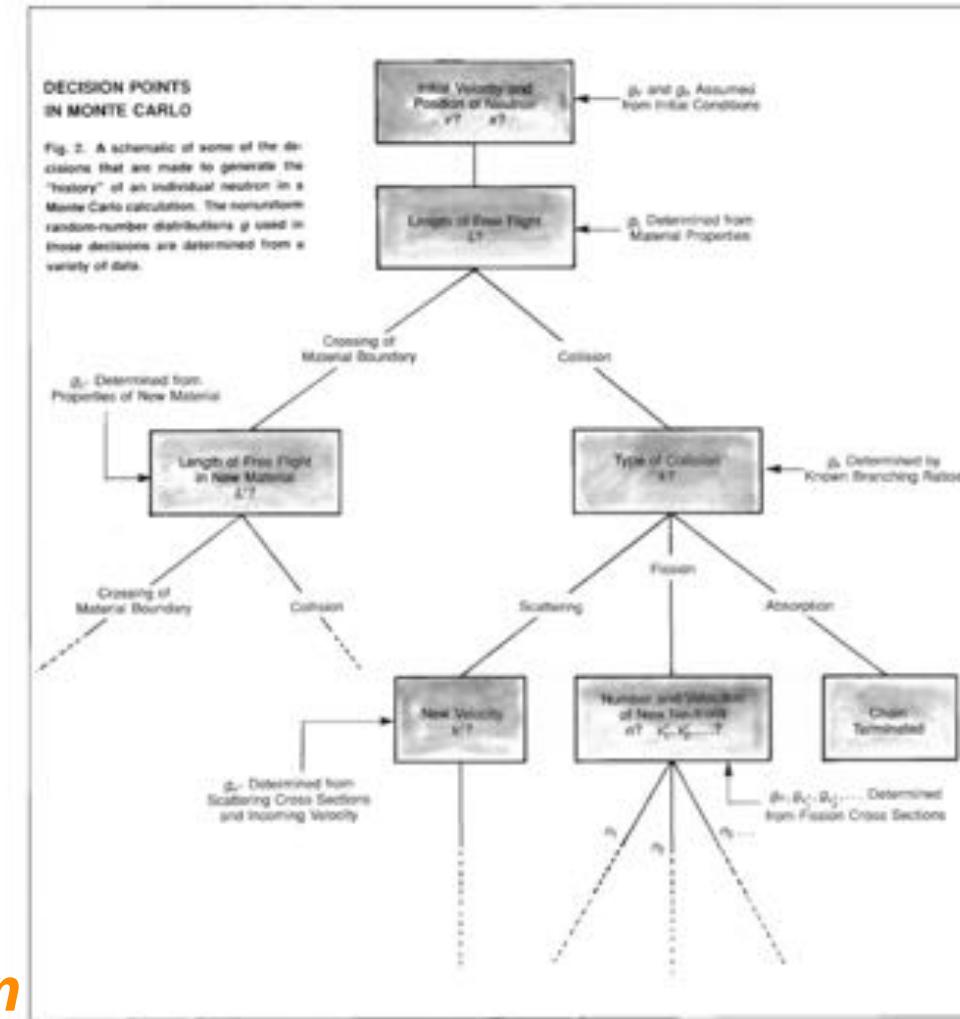


Stanislav Ulam

Experiment driven

Theory driven | Falsifiability

Simulations | Probabilistic inference | Computation



[http://www-star.st-and.ac.uk/~kw25/teaching/mcrt/MC\\_history\\_3.pdf](http://www-star.st-and.ac.uk/~kw25/teaching/mcrt/MC_history_3.pdf)

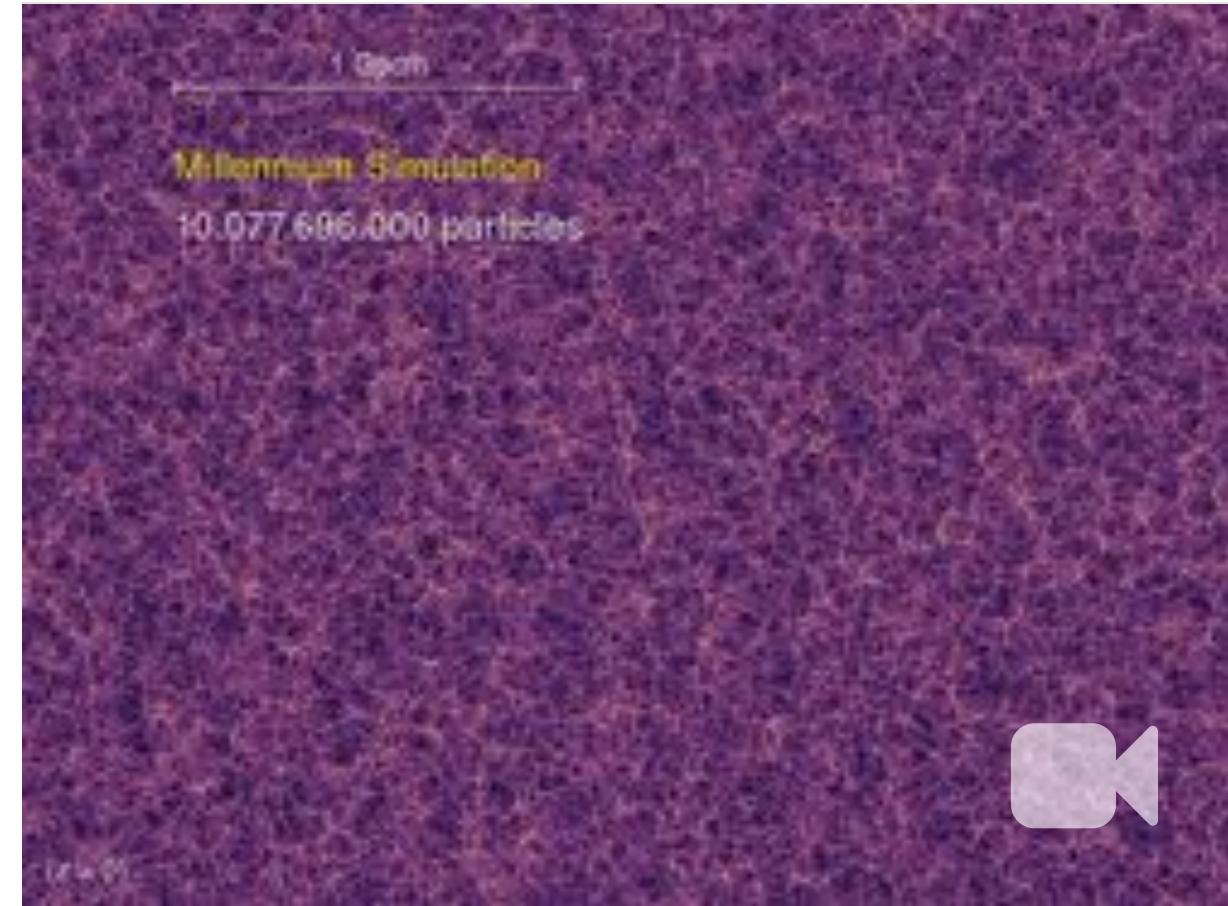
Ulam 1947

# what drives astronomy

*Experiment driven*

*Theory driven | Falsifiability*

*Simulations | Probabilistic inference | Computation*



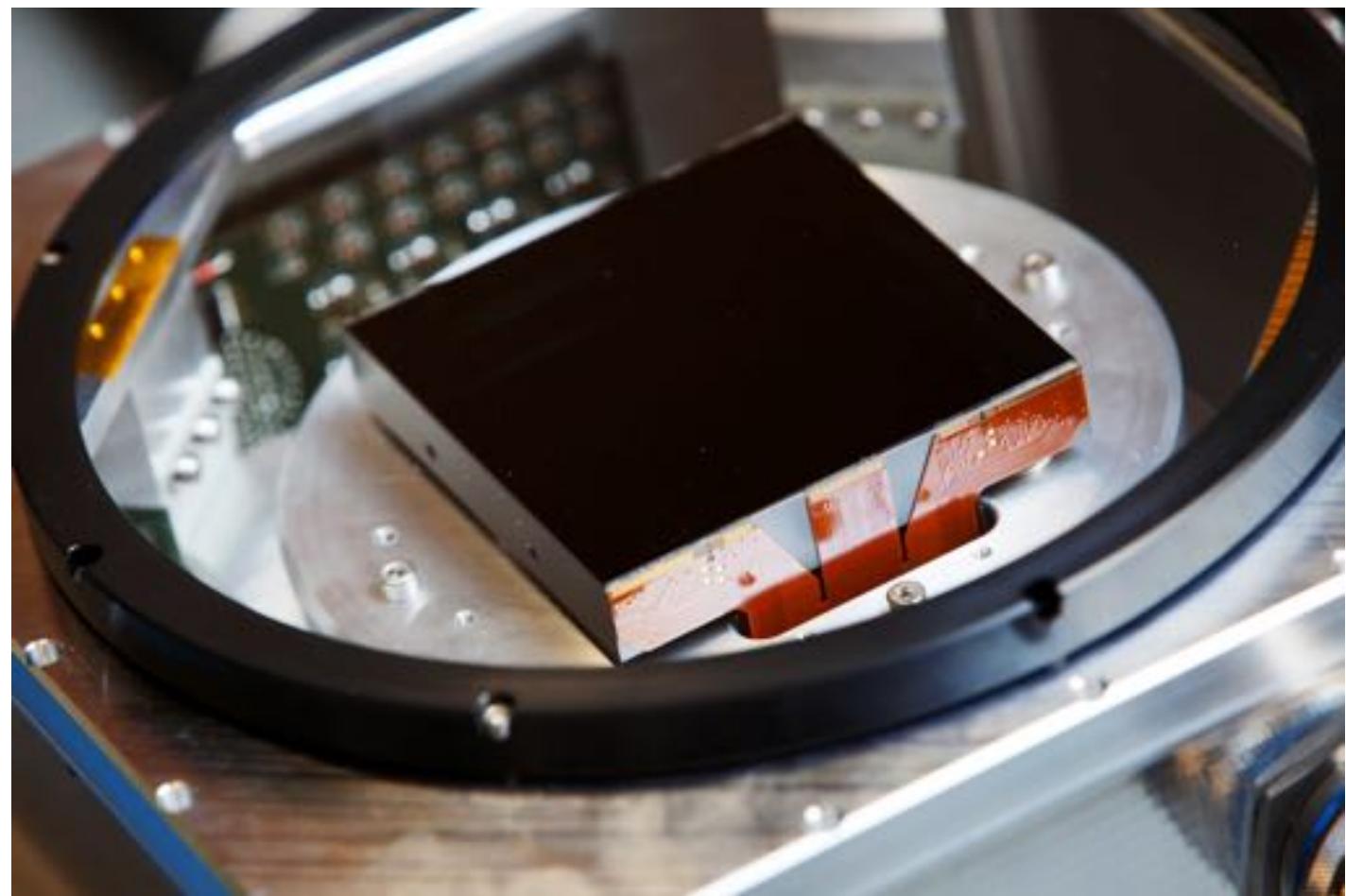
# what drives astronomy

*Experiment driven*

*Theory driven | Falsifiability*

*Simulations | Probabilistic inference | Computation*

***Data | Survey astronomy | Computation | pattern discovery***



# what drives astronomy

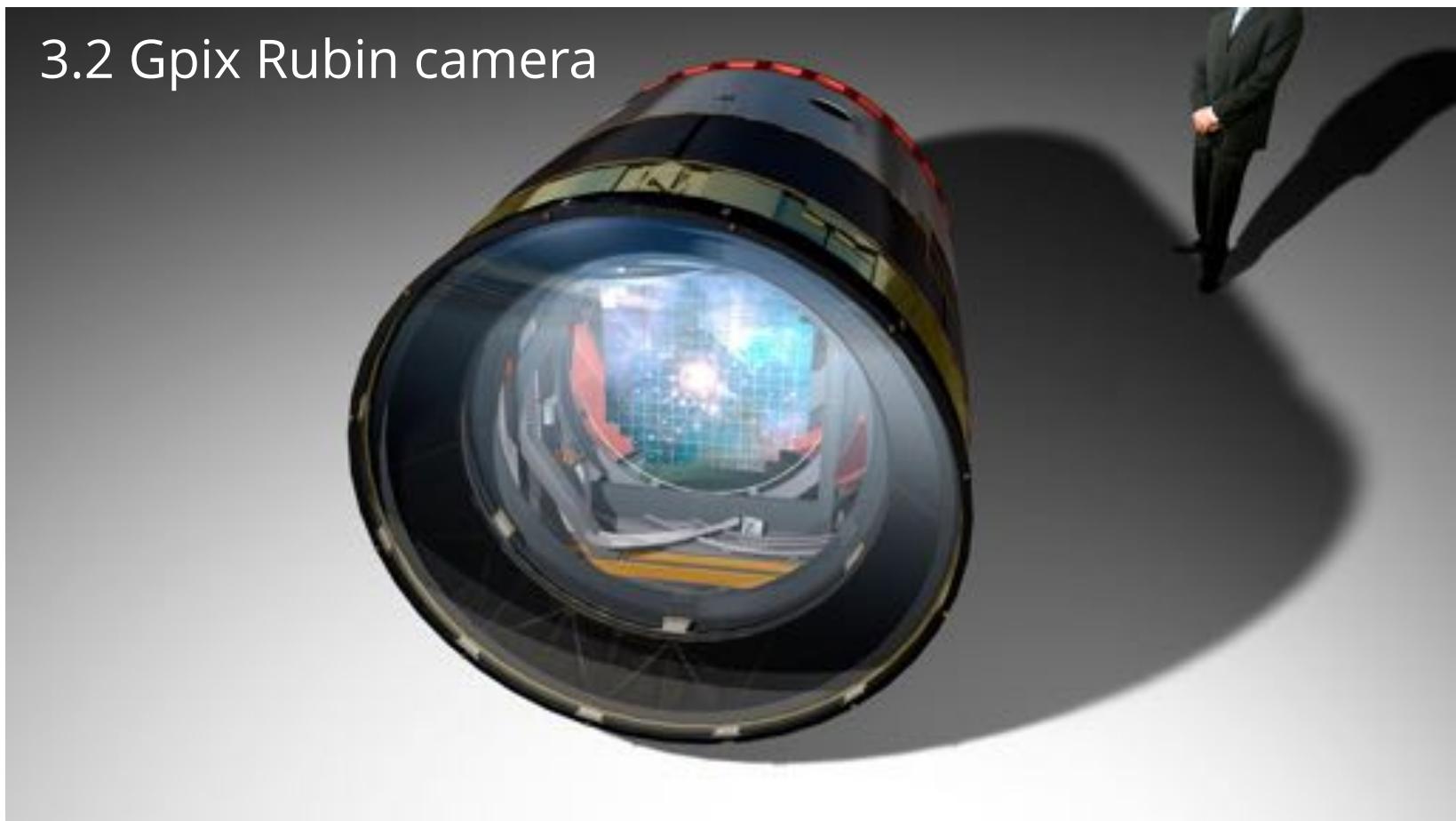
*Experiment driven*

*Theory driven | Falsifiability*

*Simulations | Probabilistic inference | Computation*

***Data | Survey astronomy | Computation | pattern discovery***

3.2 Gpix Rubin camera



# what drives astronomy

*Experiment driven*

*Theory driven | Falsifiability*

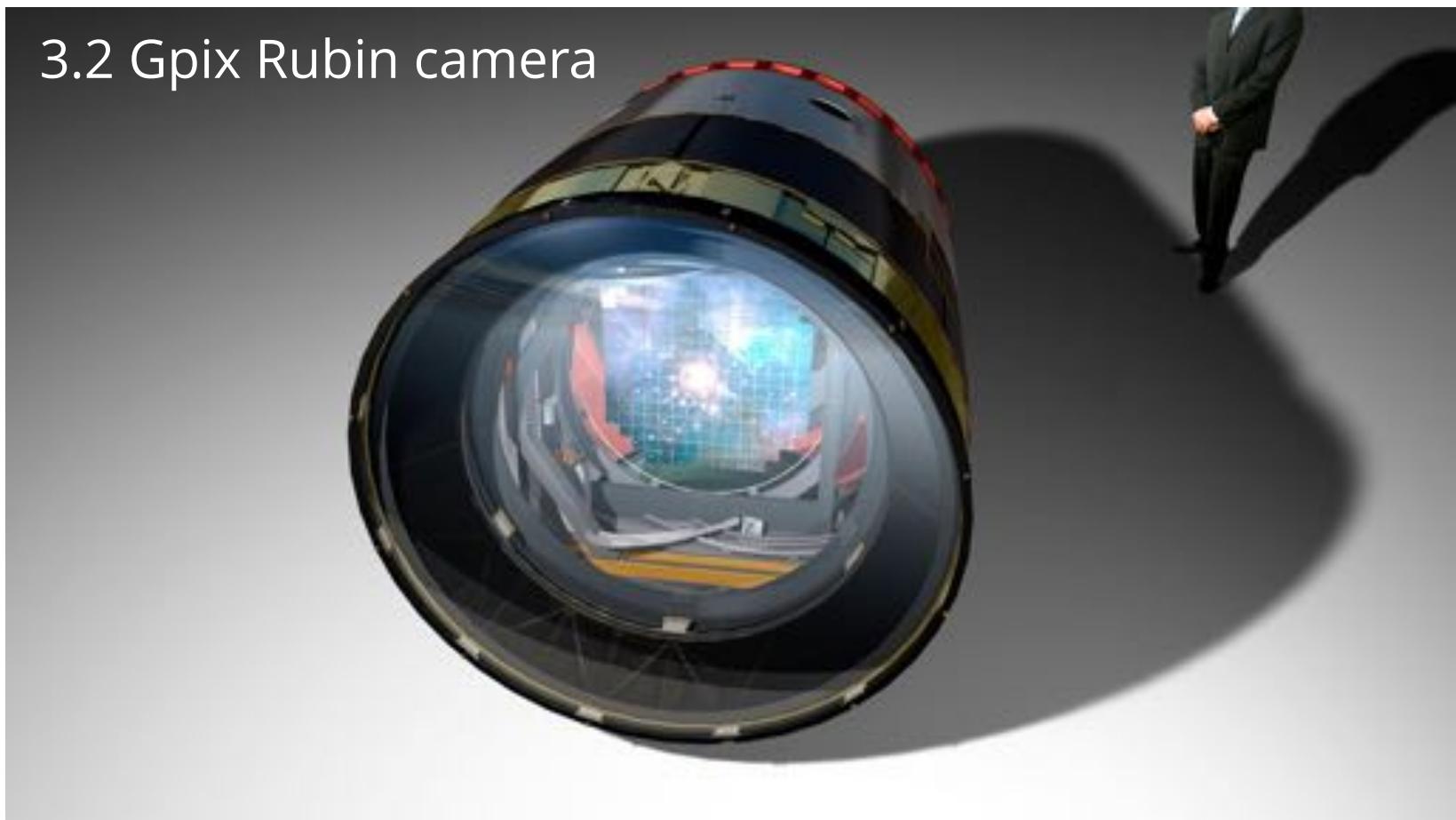
*Simulations | Probabilistic inference | Computation*

***Data | Survey astronomy | Computation | pattern discovery***

lazy learning

learning by example  
(supervised learning)

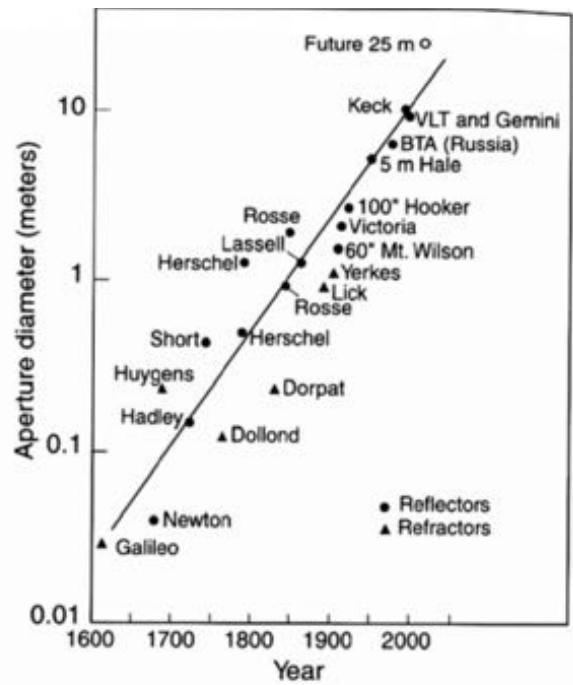
pattern discovery  
(unsupervised learning)



3.2 Gpix Rubin camera

26  
BD from  
astronomical  
surveys

# astronomical data production



**Fig. 1.** Evolution of telescope aperture diameter over the last four centuries. According to the trend line shown, the diameter of the largest telescopes doubles about every 40 years. The 20- to 30-meter class telescopes planned for the 2015 time frame display a somewhat faster growth rate than the historical trend.

# astronomical data production

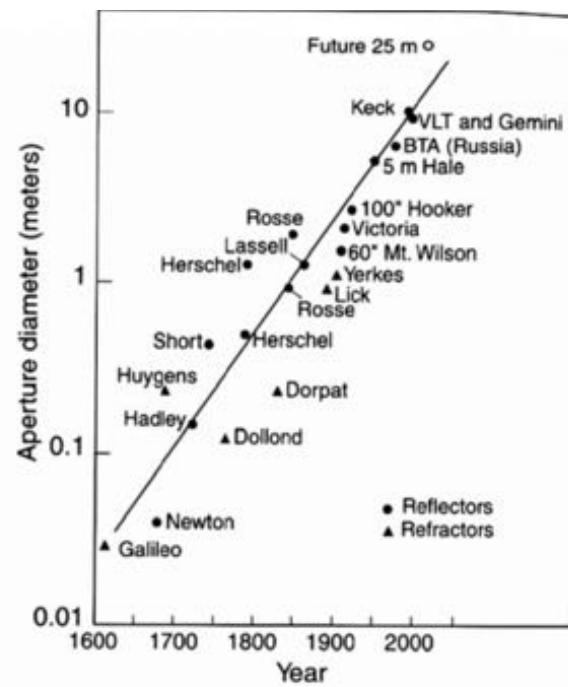
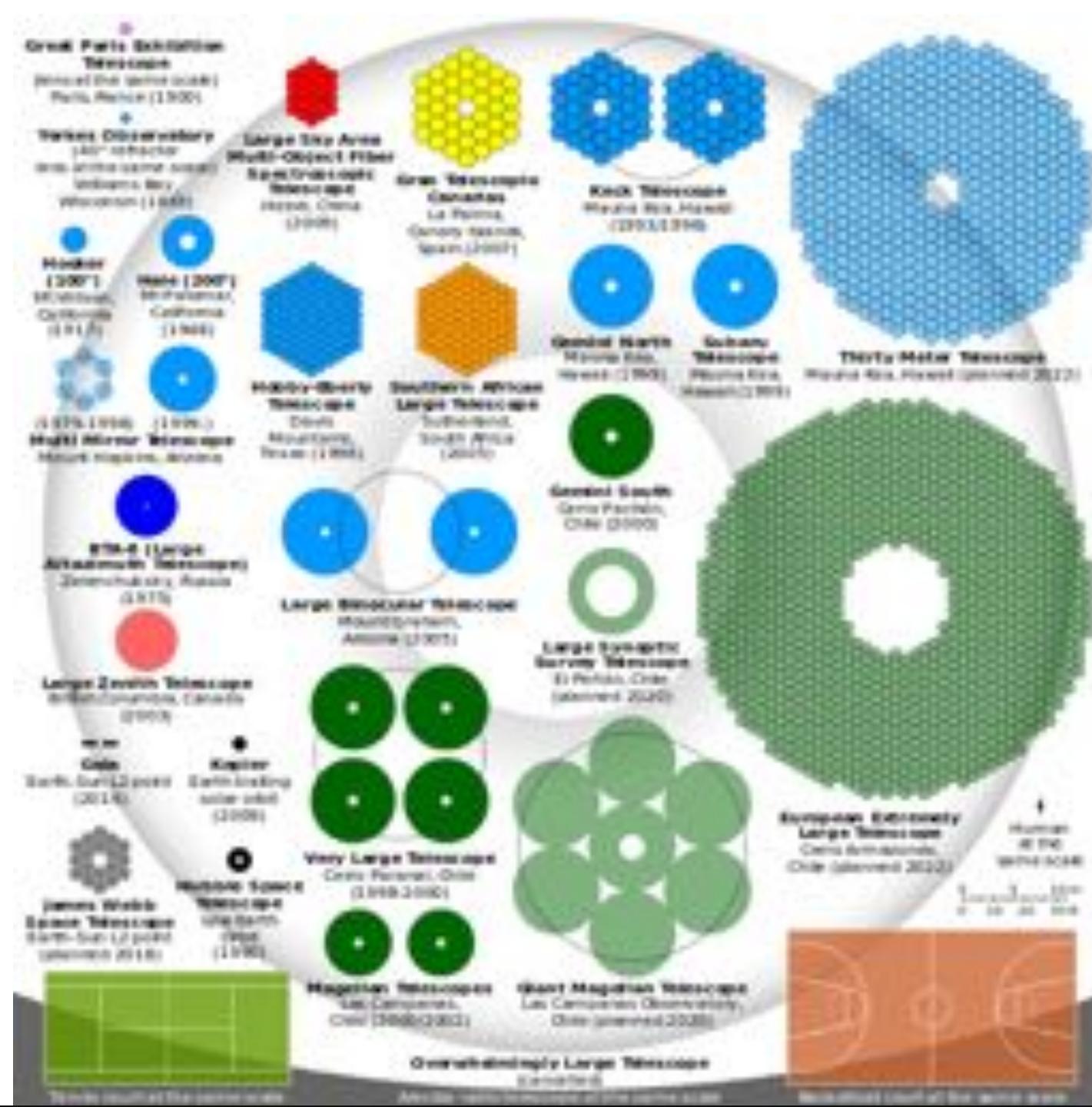


Fig. 1. Evolution of telescope aperture diameter over the last four centuries. According to the trend line shown, the diameter of the largest telescopes doubles about every 40 years. The 20- to 30-meter class telescopes planned for the 2015 time frame display a somewhat faster growth rate than the historical trend.

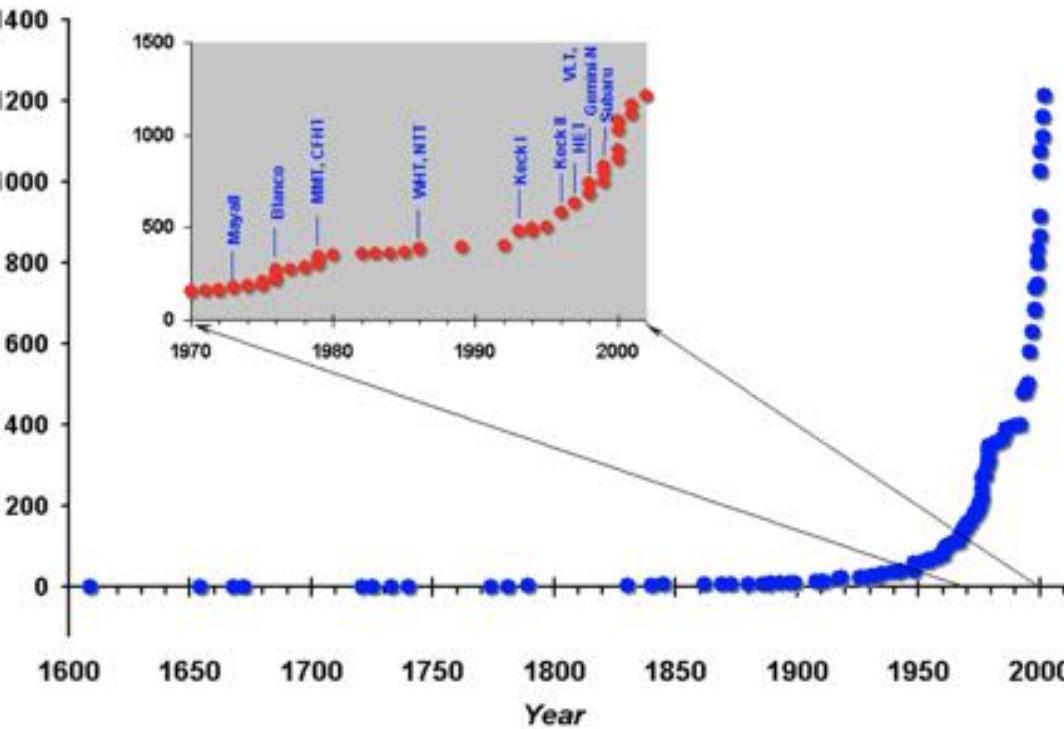


Bely, The Design and Construction of Large Telescopes



@fedhere

# astronomical data production



Mountain & Gillett, "The Revolution in Telescope Aperture",



<http://hat.astro.princeton.edu/>

# astronomical data production

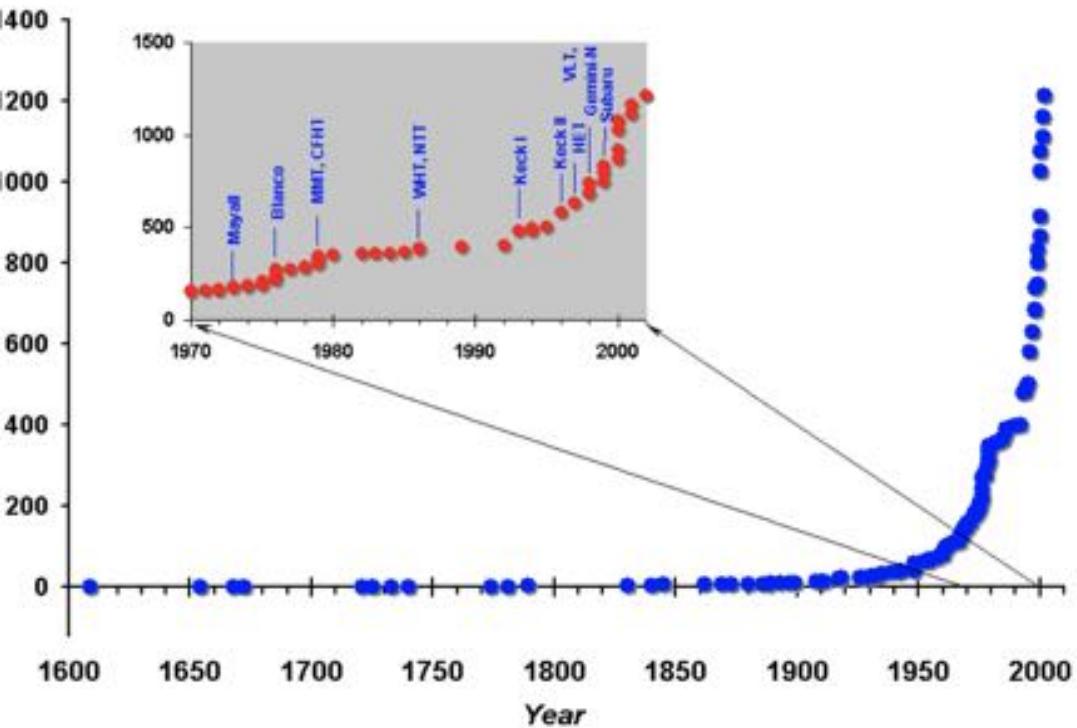
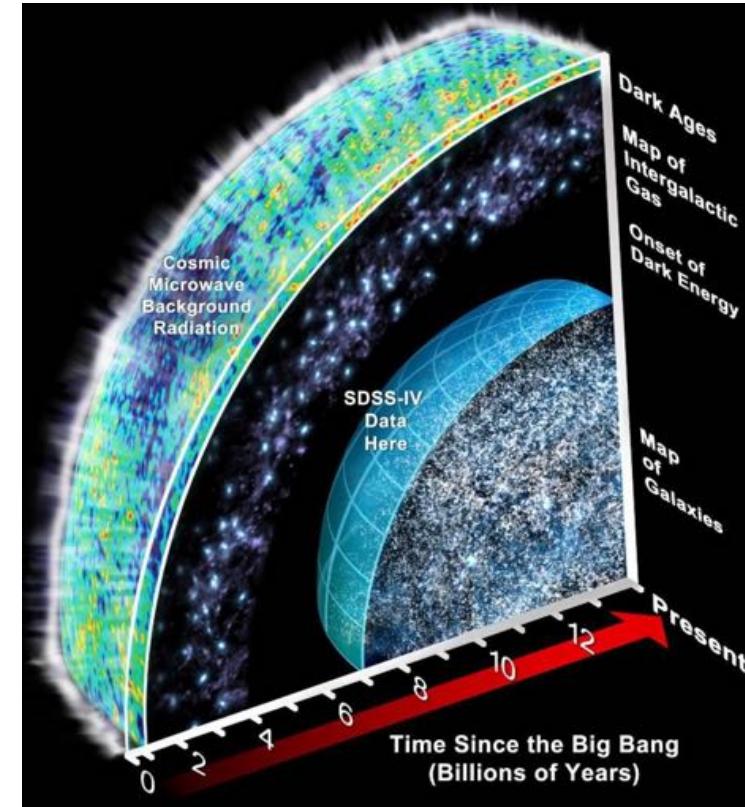


Figure 1 The growth in cumulative telescope collecting area over the past 400 years, with each point representing a completed ground-based telescope. The combination of the ability to manufacture and support large mirrors combined with adaptive optics has given the new generation of large telescopes tremendous scientific gains over the previous 4-m telescopes. For example, an 8-m telescope delivering images of 0.1 arcsec can observe point-like objects at least 20 times fainter than a conventional 4-m telescope delivering 1.0 arcsec images. Will we see such gains in the next generation of telescopes? MMT, Multiple Mirror Telescope; CFHT, Canada-France-Hawaii Telescope; WHT, William Herschel Telescope; NTT, New Technology Telescope; HET, Hobby-Eberly Telescope; VLT, Very Large Telescope; Gem-N, Gemini North Telescope.

## Area vs Volume



Both data volumes and data rates grow exponentially, with a doubling time  $\sim 1.5$  years

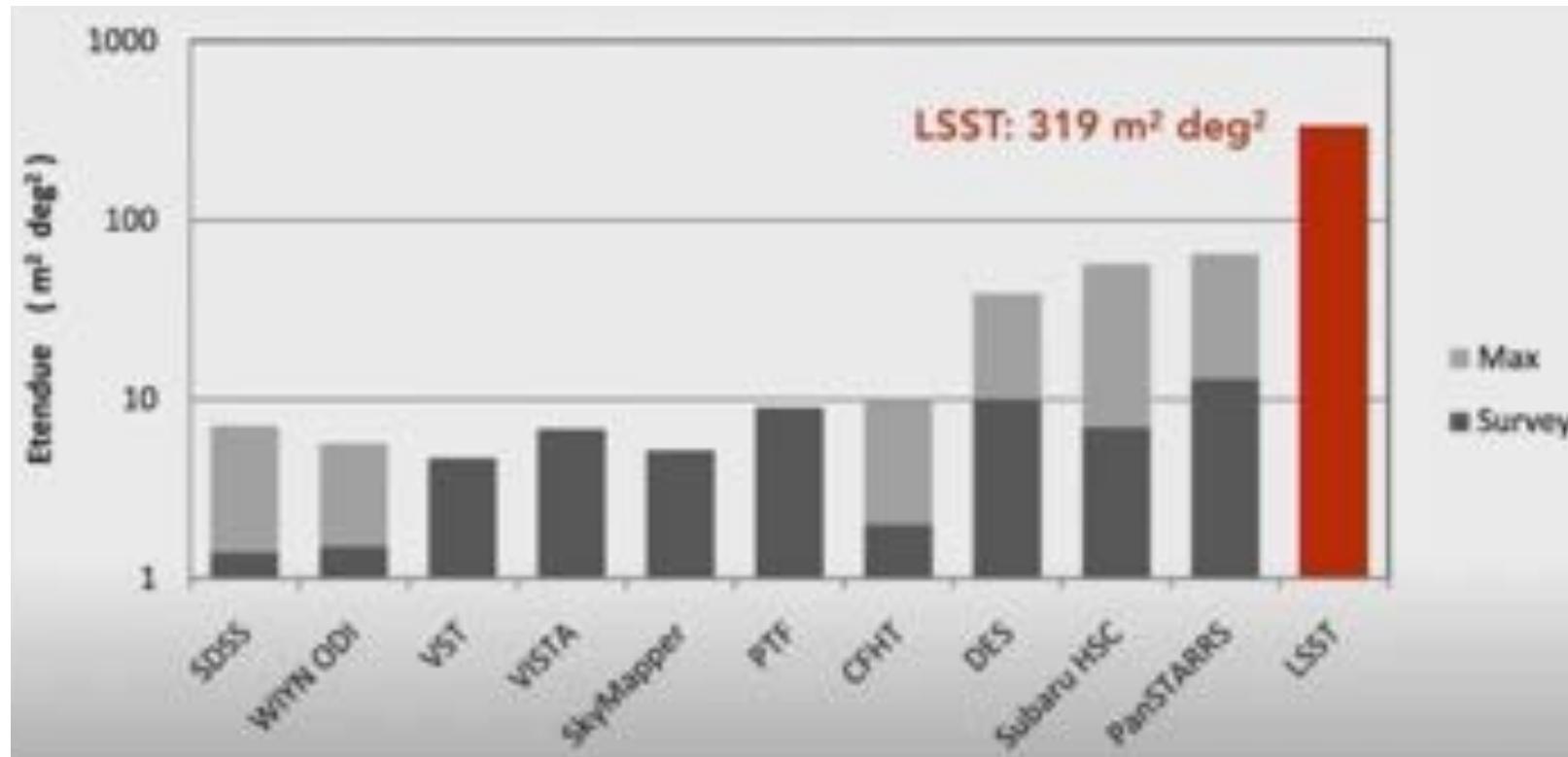
It is also estimated that everyone has access to 50% of the existing data!



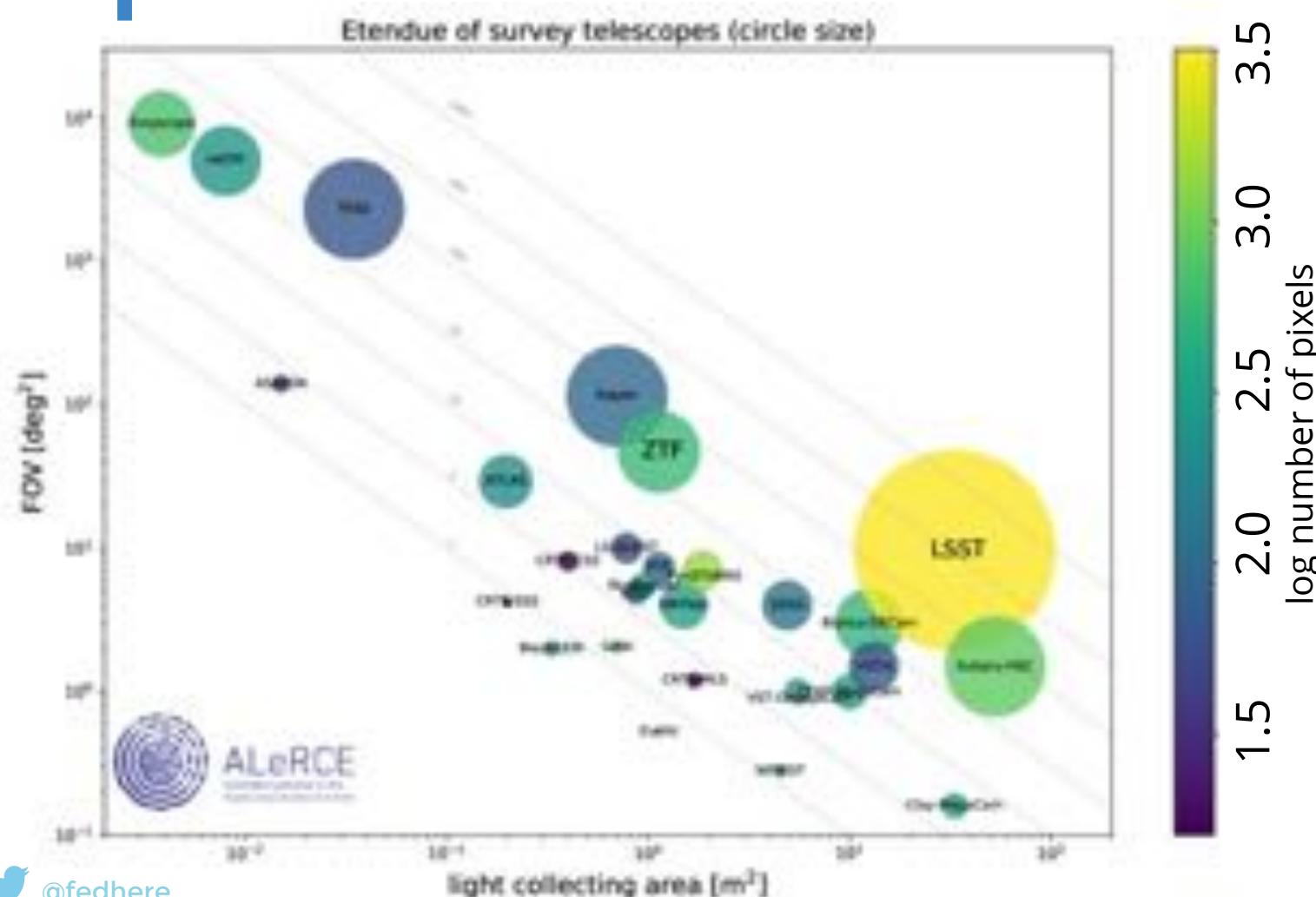
@fedhere

# astronomical data production

**Etendue:**  $area \times FoV$



# astronomical data production



# **Etendue: $area \times FoV$**

**data volume :**  
*area* x *FoV*  
x  
*resolution*  
x  
*sensitivity*

# 4-V of Big Data in astronomy

## V1: Volume

Number of bites

Number of pixels

Number of rows in a  
data table x number  
of columns for  
catalogs

# astronomical data volume

# number of sources

**Table 1.** Main data for the most important all-sky and large-area astronomical surveys providing multi-wavelength photometric data. Catalogues are given in the order of increasing wavelengths.

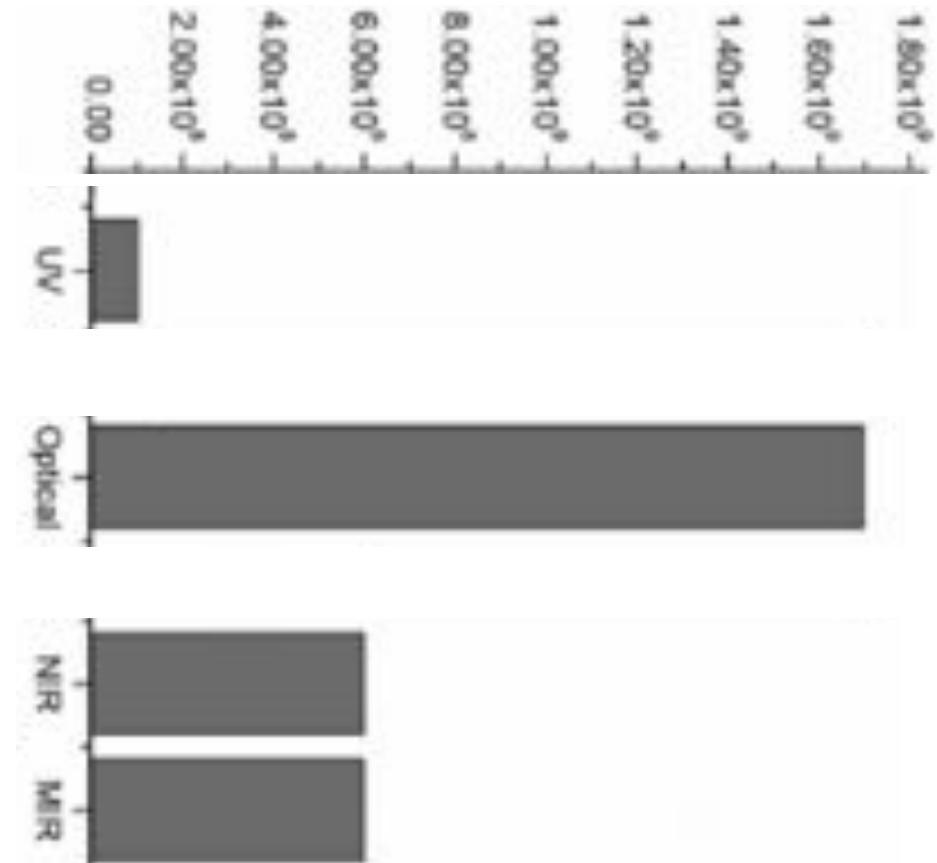
Survey, catalogue	Years	Spectral range	Sky area (deg <sup>2</sup> )	Sensitivity (mag/mJy)	Number of sources	Density (obj/deg <sup>2</sup> )
Fermi-GLAST	2008–2014	10 MeV–100 GeV	All-sky		3033	0.07
CGRO	1991–1999	20 keV–30 GeV	All-sky		1300	0.03
INTEGRAL	2002–2014	15 keV–10 MeV	All-sky		1126	0.03
ROSAT BSC	1990–1999	0.07–2.4 keV	All-sky		18,806	0.46
ROSAT FSC	1990–1999	0.07–2.4 keV	All-sky		105,924	2.57
GALEX AIS	2003–2012	1344–2831 Å	21,435	20.8 mag	65,266,291	3044.85
APM	2000	opt <i>b</i> , <i>r</i>	20,964	21.0 mag	166,466,987	7940.61
MAPS	2003	opt O, E	20,964	21.0 mag	89,234,404	4256.55
USNO-A2.0	1998	opt <i>B</i> , <i>R</i>	All-sky	21.0 mag	526,280,881	12,757.40
USNO-B1.0	2003	opt <i>B</i> , <i>R</i> , <i>I</i>	All-sky	22.5 mag	1,045,913,669	25,353.64
GSC 2.3.2	2008	opt <i>j</i> , <i>V</i> , <i>F</i> , <i>N</i>	All-sky	22.5 mag	945,592,683	22,921.79
Tycho-2	1989–1993	opt <i>BT</i> , <i>VT</i>	All-sky	16.3 mag	2,539,913	61.57
SDSS DR12	2000–2014	opt <i>u</i> , <i>g</i> , <i>r</i> , <i>i</i> , <i>z</i>	14,555	22.2 mag	932,891,133	64,094.20
DENIS	1996–2001	0.8–2.4 μm	16,700	18.5 mag	355,220,325	21,270.68
2MASS PSC	1997–2001	1.1–2.4 μm	All-sky	17.1 mag	470,992,970	11,417.46
2MASS ESC	1997–2001	1.1–2.4 μm	All-sky	17.1 mag	1,647,599	39.94
WISE	2009–2013	3–22 μm	All-sky	15.6 mag	563,921,584	13,669.83
AKARI IRC	2006–2008	7–26 μm	38,778	50 mJy	870,973	22.46
IRAS PSC	1983	8–120 μm	39,603	400 mJy	245,889	6.21
IRAS FSC	1983	8–120 μm	34,090	400 mJy	173,044	5.08
IRAS SSSC	1983	8–120 μm	39,603	400 mJy	16,740	0.42
AKARI FIS	2006–2008	50–180 μm	40,428	550 mJy	427,071	10.56
Planck	2009–2011	0.35–10 mm	All-sky	183 mJy	33,566	0.81
WMAP	2001–2011	3–14 mm	All-sky	500 mJy	471	0.01
GB6	1986–1987	6 cm	20,320	18 mJy	75,162	3.70
NVSS	1998	21 cm	33,827	2.5 mJy	1,773,484	52.43
FIRST	1999–2015	21 cm	10,000	1 mJy	946,432	94.64
SUMSS	2003–2012	36 cm	8,000	1 mJy	211,050	26.38
WENSS	1998	49/92 cm	9,950	18 mJy	229,420	23.06
7C	2007	198 cm	2,388	40 mJy	43,683	18.29

# astronomical data volume

# number of sources

**Table 1.** Main data for the most important all-sky and large-area astronomical surveys providing multi-wavelength photometric data. Catalogues are given in the order of increasing wavelengths.

Survey, catalogue	Years	Spectral range	Sky area (deg <sup>2</sup> )	Sensitivity (mag/mJy)	Number of sources	Density (obj/deg <sup>2</sup> )
Fermi-GLAST	2008–2014	10 MeV–100 GeV	All-sky		3033	0.07
CGRO	1991–1999	20 keV–30 GeV	All-sky		1300	0.03
INTEGRAL	2002–2014	15 keV–10 MeV	All-sky		1126	0.03
ROSAT BSC	1990–1999	0.07–2.4 keV	All-sky		18,806	0.46
ROSAT FSC	1990–1999	0.07–2.4 keV	All-sky		105,924	2.57
GALEX AIS	2003–2012	1344–2831Å	21,435	20.8 mag	65,266,291	3044.85
APM	2000	opt <i>b</i> , <i>r</i>	20,964	21.0 mag	166,466,987	7940.61
MAPS	2003	opt O, E	20,964	21.0 mag	89,234,404	4256.55
USNO-A2.0	1998	opt <i>B</i> , <i>R</i>	All-sky	21.0 mag	526,280,881	12,757.40
USNO-B1.0	2003	opt <i>B</i> , <i>R</i> , <i>I</i>	All-sky	22.5 mag	1,045,913,669	25,353.64
GSC 2.3.2	2008	opt <i>j</i> , <i>V</i> , <i>F</i> , <i>N</i>	All-sky	22.5 mag	945,592,683	22,921.79
Tycho-2	1989–1993	opt <i>BT</i> , <i>VT</i>	All-sky	16.3 mag	2,539,913	61.57
SDSS DR12	2000–2014	opt <i>u</i> , <i>g</i> , <i>r</i> , <i>i</i> , <i>z</i>	14,555	22.2 mag	932,891,133	64,094.20
DENIS	1996–2001	0.8–2.4 μm	16,700	18.5 mag	355,220,325	21,270.68
2MASS PSC	1997–2001	1.1–2.4 μm	All-sky	17.1 mag	470,992,970	11,417.46
2MASS ESC	1997–2001	1.1–2.4 μm	All-sky	17.1 mag	1,647,599	39.94
WISE	2009–2013	3–22 μm	All-sky	15.6 mag	563,921,584	13,669.83
AKARI IRC	2006–2008	7–26 μm	38,778	50 mJy	870,973	22.46
IRAS PSC	1983	8–120 μm	39,603	400 mJy	245,889	6.21
IRAS FSC	1983	8–120 μm	34,090	400 mJy	173,044	5.08
IRAS SSSC	1983	8–120 μm	39,603	400 mJy	16,740	0.42
AKARI FIS	2006–2008	50–180 μm	40,428	550 mJy	427,071	10.56
Planck	2009–2011	0.35–10 mm	All-sky	183 mJy	33,566	0.81
WMAP	2001–2011	3–14 mm	All-sky	500 mJy	471	0.01
GB6	1986–1987	6 cm	20,320	18 mJy	75,162	3.70
NVSS	1998	21 cm	33,827	2.5 mJy	1,773,484	52.43
FIRST	1999–2015	21 cm	10,000	1 mJy	946,432	94.64
SUMSS	2003–2012	36 cm	8,000	1 mJy	211,050	26.38
WENSS	1998	49/92 cm	9,950	18 mJy	229,420	23.06
7C	2007	198 cm	2,388	40 mJy	43,683	18.29



# 4-V of Big Data in astronomy

## V1: Volume

Number of bites

Number of pixels

Number of rows in a  
data table x number  
of columns for  
catalogs

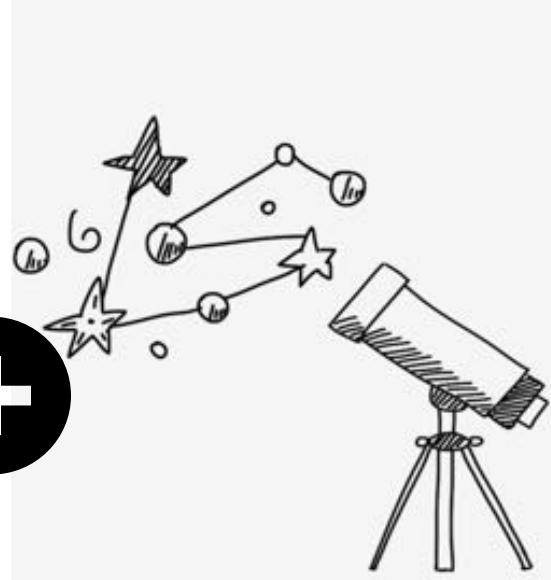
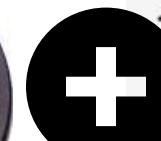
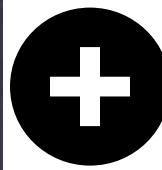
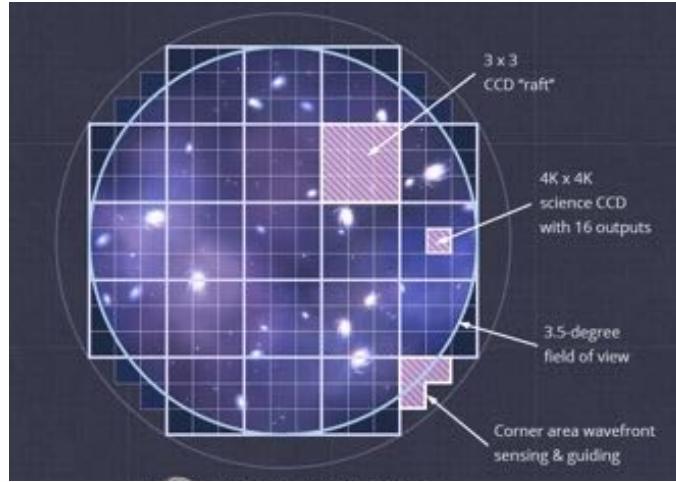
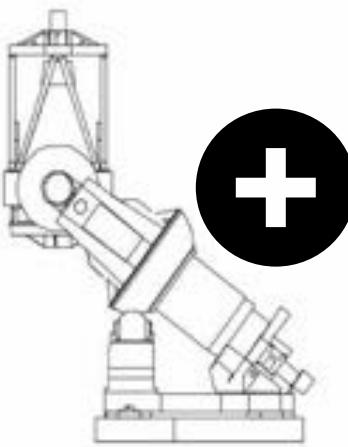
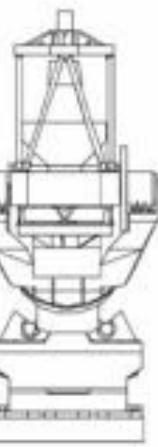
## V2: Variety

Diverse science return  
from the same dataset.

Multiwavelength  
Multimessenger

Images and spectra

# ground based how do the data get big?



filters

→ variety (complexity)

telescope size

→ fainter, more distant

FoV

→ more sky area at once

camera size

→ more data units

resolution

→ more objects/details

# optical

*"The Sloan Digital Sky Survey has created the most detailed three-dimensional maps of the Universe ever made, with deep multi-color images of one third of the sky, and spectra for more than three million astronomical objects. Learn and explore all phases and surveys—past, present, and future—of the SDSS."*

5 bands

2.5m

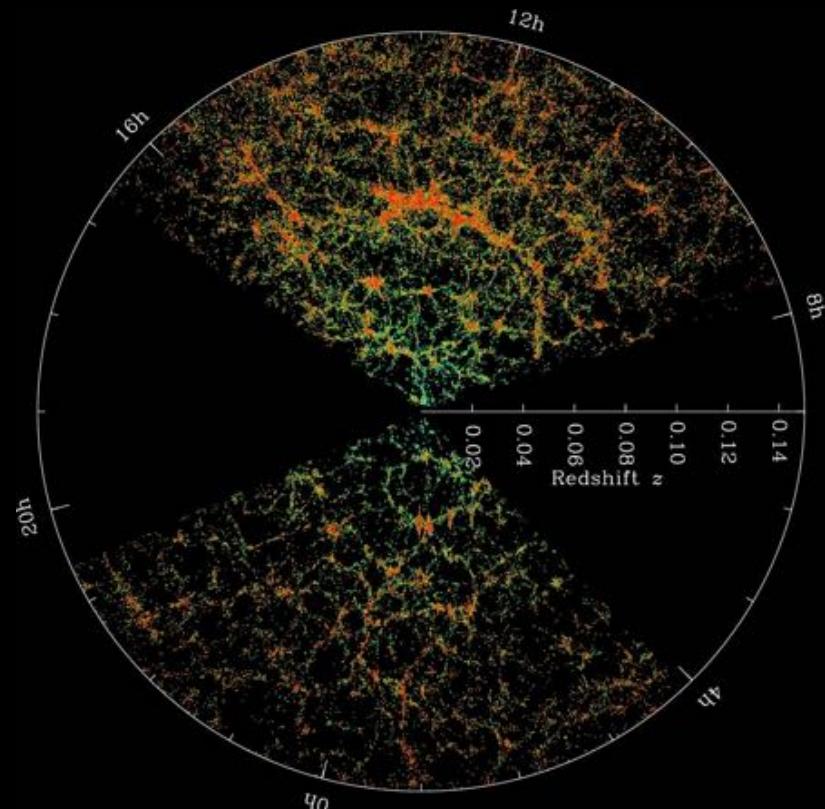
6 sq degree

4Mpix

1"/pix



# SDSS



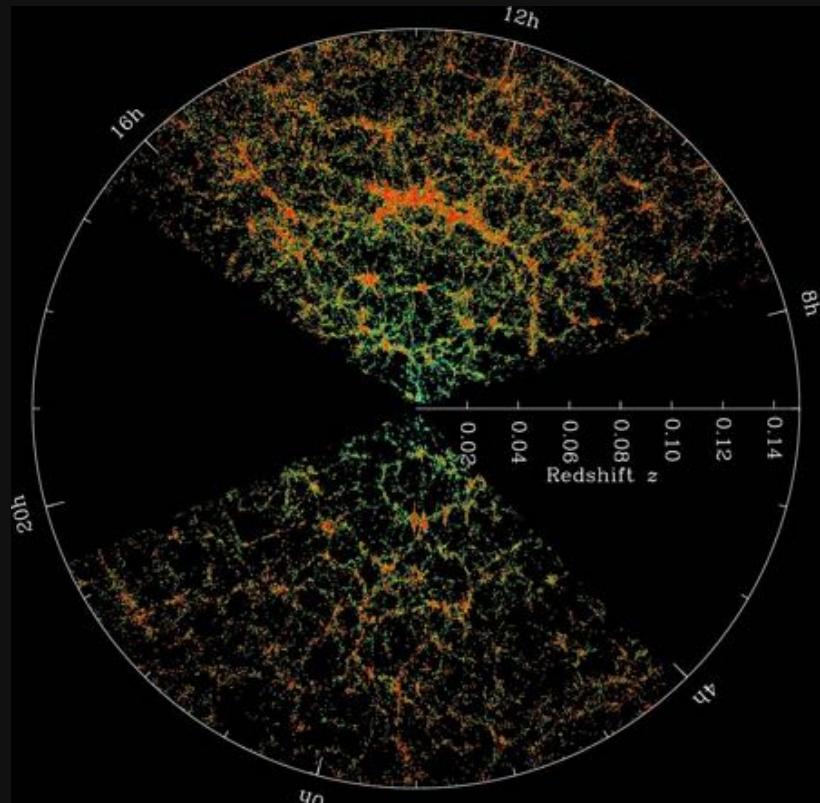
# optical: data releases



# SDSS

photometric parameters for 53 million unique objects.

SDSS DR	images	catalog	1D+2D spectra
2003 DR1	2.3Tb	0.5Tb	
2003 DR2	5Tb	0.7Tb	
2004 DR3	6Tb	1.2Tb	
...			
2009 DR7	15.7Tb	18Tb	3.5 TB
...			
<b>2019 DR16</b>			<b>273 TB</b>



The SDSS map of the Universe. Each dot is a galaxy; the color bar shows the local density.

# optical: data releases

SLOAN DIGITAL SKY SURVEY  
**SkyServer DR16 plus** 

SciServer 

Not logged in [Help](#) [Login](#)

Home Data Schema Education Astronomy SDSS Contact Us Download Site Search Help

**DR16 Tools**

Getting Started  
Famous places  
Get Images  
Scrolling sky  
Visual Tools  
Search  
- Radial  
- Rectangular  
- Search Form  
- **SQL** NEW!  
- Imaging Query  
- Spectro Query  
- IR Spec Query

CrossID  
Skyquery CrossMatch  
CasJobs

**SQL Search**

This page allows you to directly submit a SQL (Structured Query Language) query to the SDSS database. You can modify the default query as you wish, or cut and paste a query from the [SDSS Sample Queries](#) page.

**Please note:** To be fair to other users, queries run from SkyServer search tools are restricted in how long they can run and how much output they return, by **timeouts** and **row limits**. Please see the [Query Limits](#) help page for more information. If your query is not restricted by a timeout or number of rows returned, please use the [CasJobs](#) batch query service.

```
-- This query does a table JOIN between the imaging (PhotoObj) and spectra
-- (SpecObj) tables and includes the necessary columns in the SELECT to upload
-- the results to the SAS (Science Archive Server) for FITS file retrieval.
SELECT TOP 10
    p.objid, p.ra, p.dec, p.u, p.g, p.r, p.i, p.z,
    p.run, p.rerun, p.camcol, p.field,
    s.specobjid, s.class, s.z as redshift,
    s.plate, s.mjd, s.fiberid
FROM PhotoObj AS p
    JOIN SpecObj AS s ON s.bestobjid = p.objid
WHERE
    p.u BETWEEN 0 AND 19.6
    AND g BETWEEN 0 AND 20
```

[Check syntax](#) **Output Format**  HTML  XML  CSV  JSON  VOTable  FITS  MyDB NEW!

**Submit query** **Table name**

To find out more about the database schema use the [Schema Browser](#).



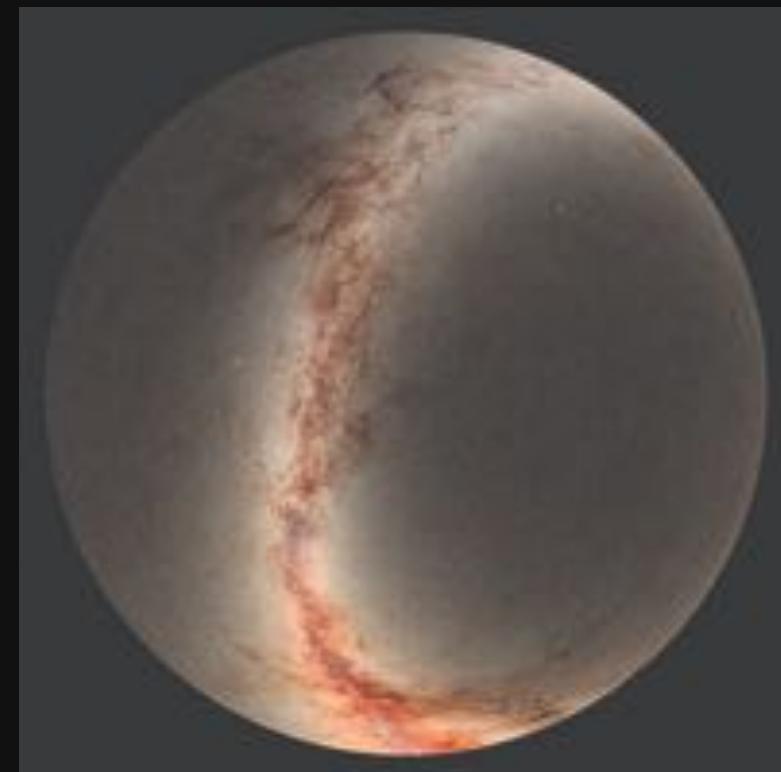
@fedhere

# optical PS1

2019 DR2 PanSTARRS

1.6 Pb

The amount of imaging data is equivalent to two billion selfies, or 30,000 times the total text content of Wikipedia. The catalog data is 15 times the volume of the Library of Congress.



# optical: data releases

<b>Number of raw on-sky camera exposures ingested:</b>	<b>40,000</b>	<b>23 TB</b>
<b>Volume (with ancillary files):</b>		<b>195 TB</b>
<b>Number of lightcurves</b>	<b>319</b>	<b>110 GB</b>



# optical: DES



THE DARK ENERGY SURVEY

570 MP camera with a 3 deg<sup>2</sup> field of view  
installed at the prime focus of the Blanco 4 m

5 bands

4m

3 sq degree

0.5Gpix

0.2"/pix

# optical: ZTF



2 band

1.2m

47 sq degree

0.5Gpix

# optical: Rubin LSST

6 bands

8m (6.5 effective)

9 sq degree

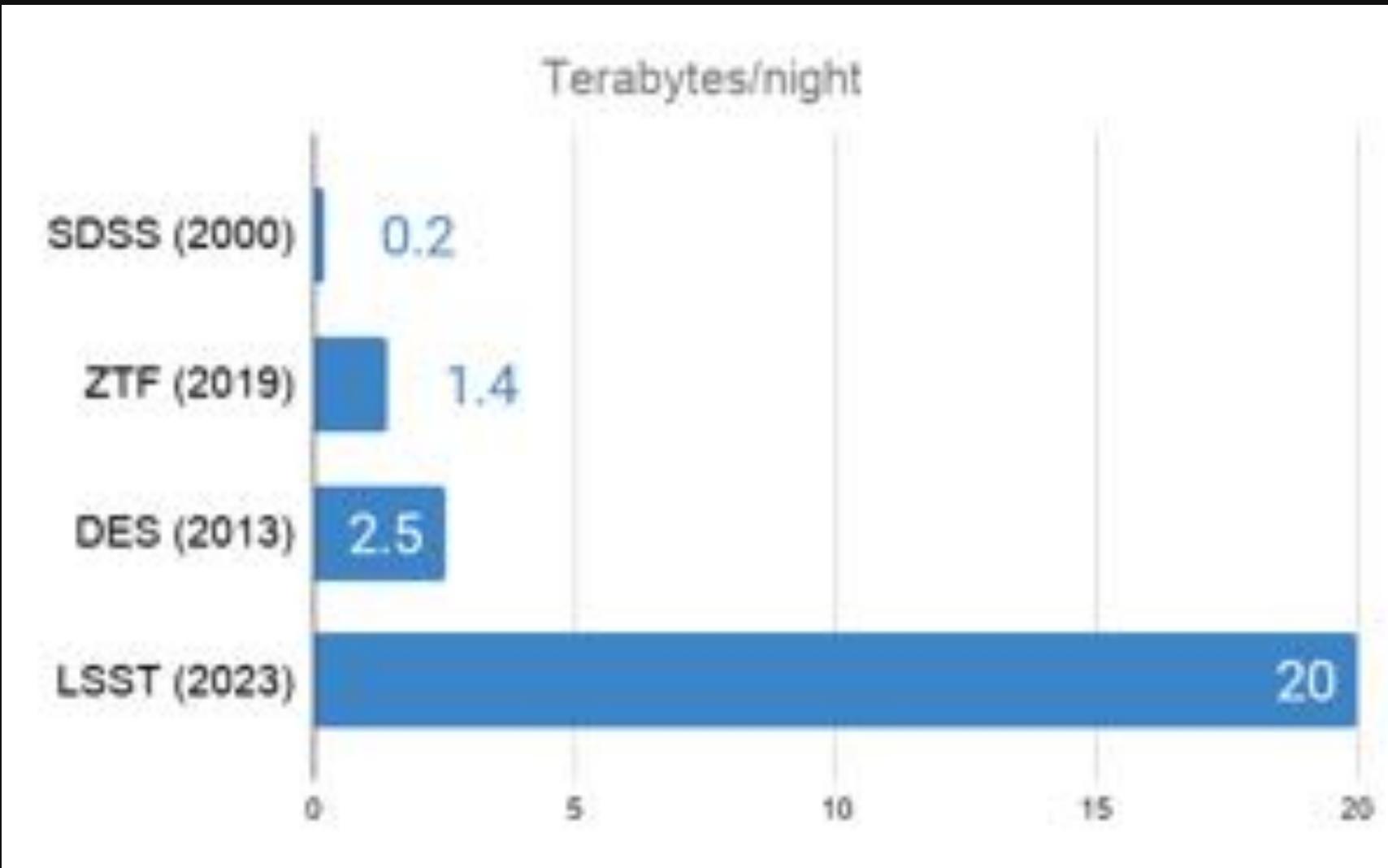
0.2"/pix

3.2Gpix

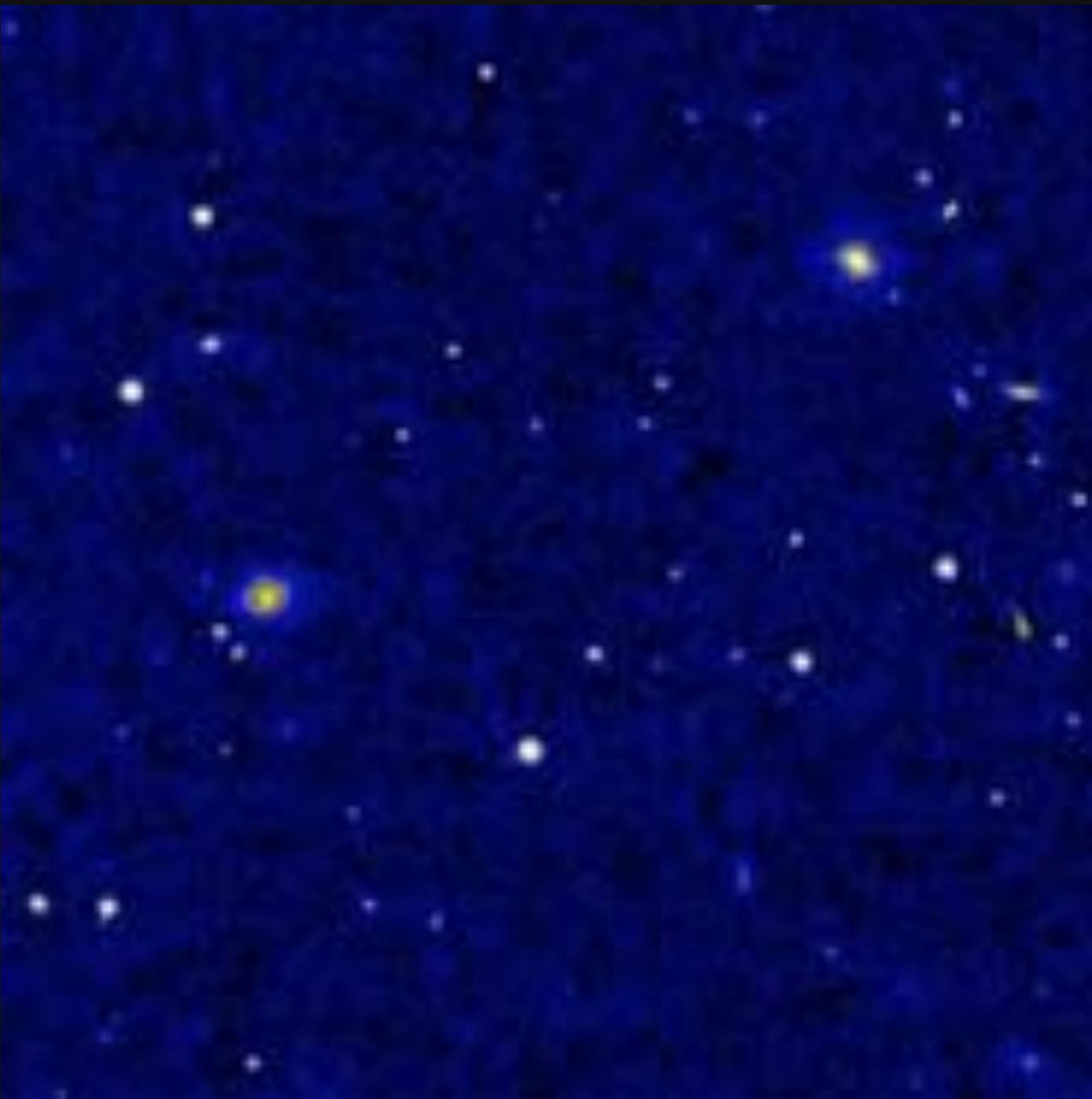


**VERA C. RUBIN  
OBSERVATORY**

# optical: Rubin LSST



DSS:  
digitized  
photographic  
plates



One quarter the diameter of the moon

# SDSS



One quarter the diameter of the moon

# DLS

20 sq deg ultra-deep multi-band sky survey.



One quarter the diameter of the moon

# Rubin LSST (simulated)



One quarter the diameter of the moon

# optical: other surveys

MACHO (Microlensing)

Catalina Sky Survey

Digitized Sky Survey

SNLS (Supernovae)

OGLE (Microlensing)

DLA (weak lensing)

VIMOS-VLT Deep Survey

Palomar Distant Solar System Survey

DESI Legacy Imaging Surveys

# 5Pb data

*Extreme imaging via physical model inversion: seeing around corners and imaging black holes,*  
Dr. Katie Bouman



# 5Pb data



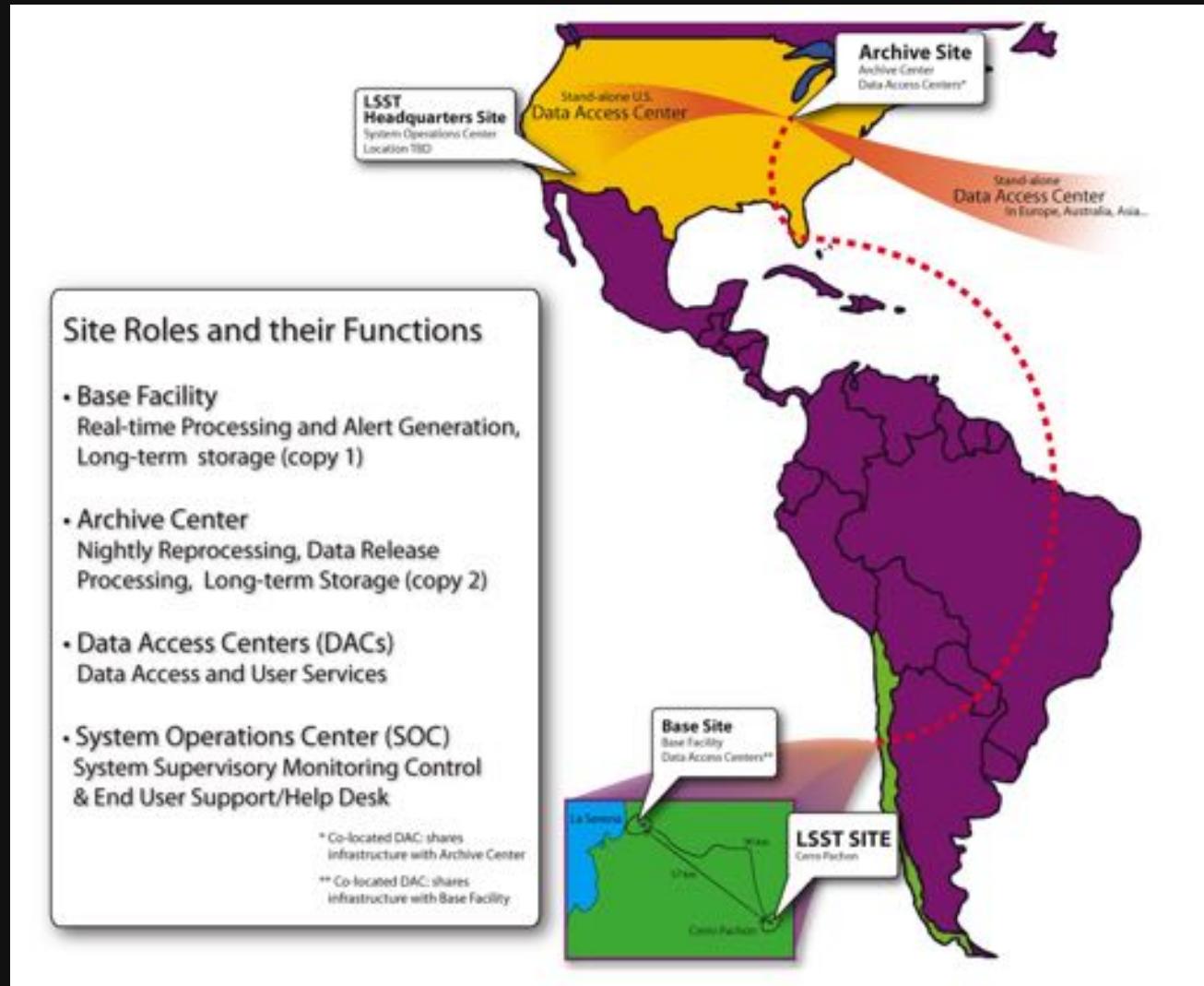
*Extreme imaging via physical model inversion: seeing around corners and imaging black holes,*  
Dr. Katie Bouman

reconstruct images and video from a sparse telescope array distributed around the globe. Additionally, it presents a number of evaluation techniques developed to rigorously evaluate imaging methods in order to establish confidence in reconstructions done with real scientific data.

# 5Pb data

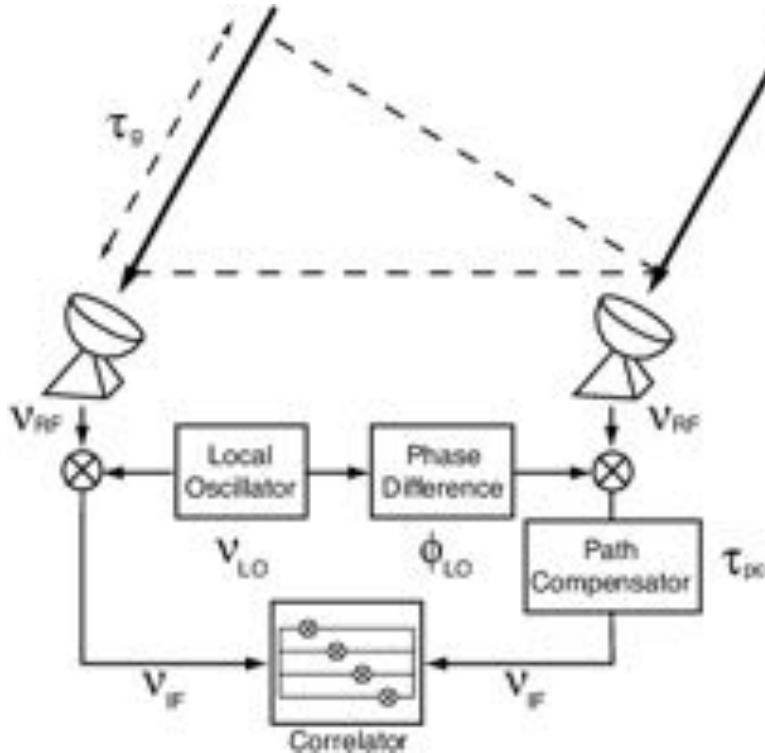
*Extreme imaging via physical model inversion: seeing around corners and imaging black holes,*  
Dr. Katie Bouman





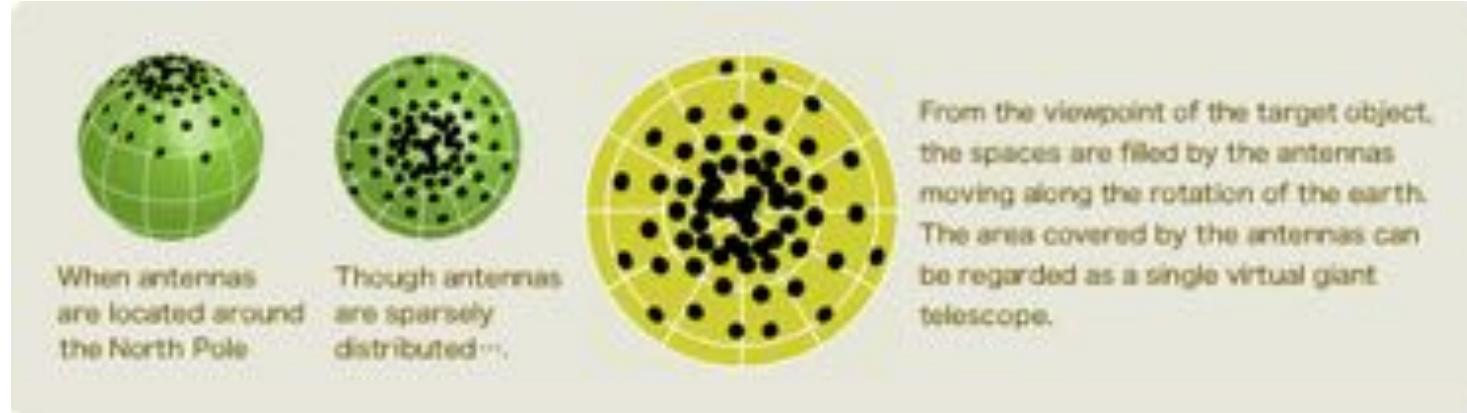
# radio

Radio interferometers do not image the sky directly. Instead they measure the amount of power on different angular scales,



interferometry:

Create a virtual telescope by combining multiple antennae.



\*The actual ALMA antenna location differs from the figure above. The figure is a conceptual illustration to explain the principle of the "aperture synthesis" technique (interferometric imaging method) in a very simple way.

Cross correlating the signal from each antenna produces coherent interpretable radio images of the sky

The directed use of the Earth's rotation for filling the Fourier plane is known as *Earth rotation aperture synthesis* and was the subject of the 1974 Nobel Prize in Physics

... if you thought LSST was  
Big Data...SKA



# Square Kilometer Array

<https://www.skatelescope.org/ska-community-briefing-18jan2017/>

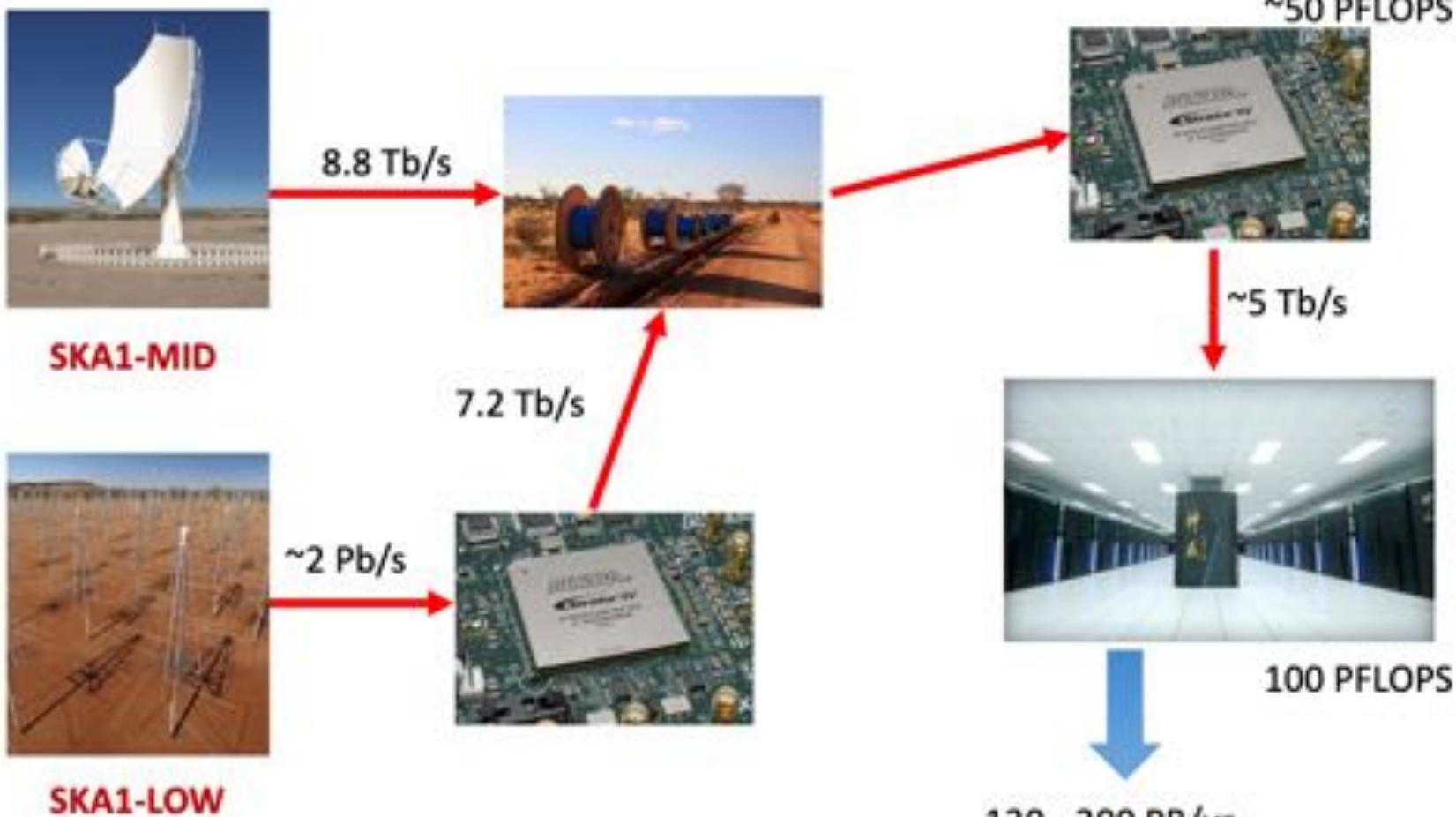
## SKA– Key Science Drivers: The history of the Universe



Extremely broad range of science!

# Square Kilometer Array

## Data Flow through the SKA

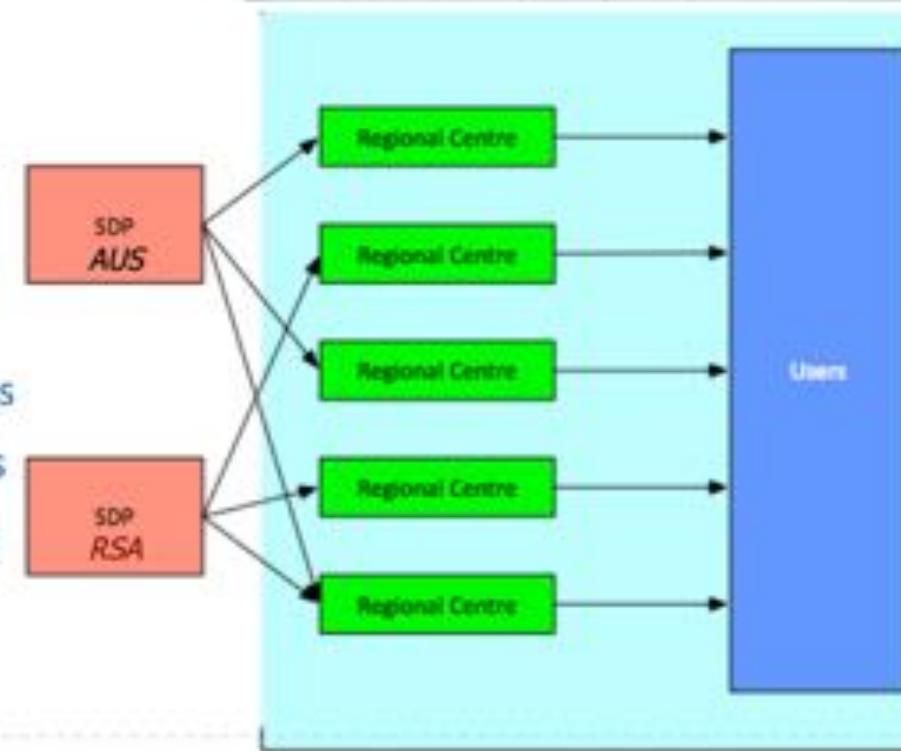


<https://www.skatelescope.org/ska-community-briefing-18jan2017/>

# Square Kilometer Array

## SKA Regional Centres – outside SKAO scope

- Required
  - capacity for reprocessing data and their analysis
  - storage for a long-term archive
  - local user support
- Intent
  - SKA partner countries planning SKA regional Centres
  - National super-computing centres
  - Provide local support to scientists
  - Development of new techniques, new algorithms
  - Deliver SKA science



## **Data rate: 0.5–1 TB s<sup>-1</sup>.**

Data from the individual antennas of the SKA1-MID, or the individual stations of the SKA1-LOW, are transported to the central signal processing facility, where the data from each pair of antennas/stations are correlated to produce the visibility data

$$\text{300PB/telescope/year} = \\ \text{8.5EB}$$

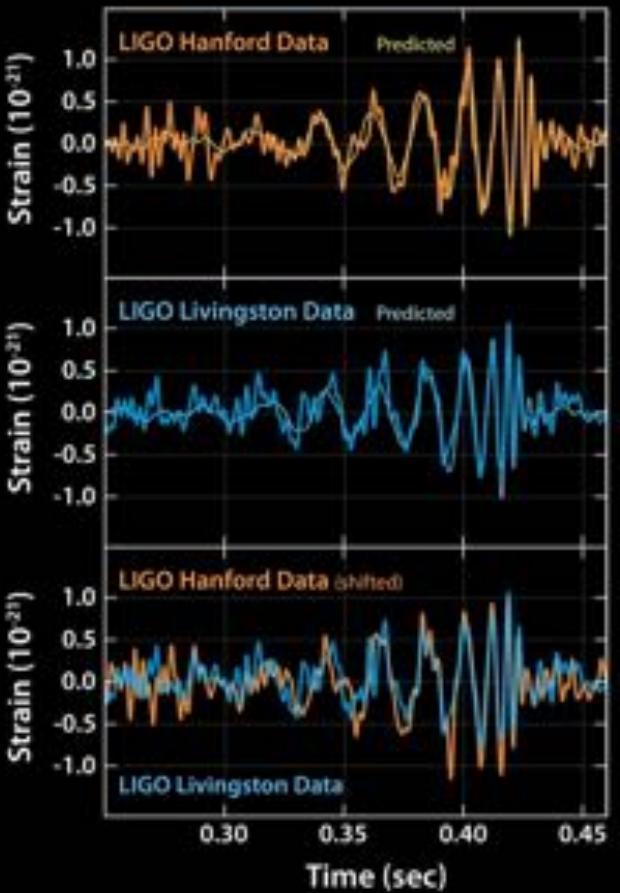
Typical images will have spatial axes with 2e15 x 2e15 pixels up to 2e16 frequency channels. This dimensionality results in petabyte scale volumes for individual data products,

Performing multiple Fourier transforms on data with image sizes as large as 2e15 or even 2e16 pixels on a side is computationally prohibitive.

[...] In spite of these challenges, the computing for the SKA over the coming decade is certainly achievable, and it is likely that its implementation will be instrumental in driving the next generation of global e-infrastructure.

*The computers must be able to make decisions on objects of interest, and remove data which is of no scientific benefit, such as radio interference from things like mobile phones or similar devices, even with the remote locations which will host the SKA.*

# Gravitational Wave and MMA



## Computation and Data Collection

Computers are required both to run the LIGO instruments and to process the data that it collects.

When it is in 'observing' mode, LIGO generates terabytes (1000's of gigabytes) of data *every day*. All of this information must be transferred to a network of supercomputers for storage and archiving. Such supercomputers are located at each of the observatories, at Caltech, at MIT, and at various other institutions. Once the data is secured, scientists can use customized computer programs to scour the data for gravitational waves.

The amount of data LIGO collects is as incomprehensively large as gravitational wave signals are small. LIGO's archive already holds the equivalent over 1-million DVDs of data and will add the equivalent of about 178-thousand DVDs each year to its archive. In actual numerical terms, the data archive at Caltech holds over 4.5 Petabytes (Pb) of data, and will grow at a rate of about 0.8 Pb (800 terabytes) per year. What's a petabyte? If you wanted to count up to a petabyte by counting one byte per second, it would take you 35.7 million **years** to reach one petabyte!

Storing information is one thing; processing it is another. Processing and analyzing all of LIGO's data requires a vast computing infrastructure. For LIGO's first observing run in 2015, the LIGO Lab will provide 35 MSU (**million service units**) worth of computing cycles/time. This is equivalent to running a modern 4-core laptop computer for 1,000 years! The amount of computing time is expected to grow by a factor of 10 to around 400 MSU by the time LIGO has completed its third observing run.

If you'd like to learn even more about all of LIGO's remarkable technology and engineering, visit [Look Deeper](#).

# Gravitational Wave and MMA

To search for binaries with components more massive than  $m_{\min}=0.2M_{\odot}$  while losing no more than 10% of events the initial LIGO interferometers will require  $1 \times 10^{11}$  flops for data analysis to keep up with data acquisition.

Advanced LIGO will require  $7.8 \times 10^{11}$  flops, and VIRGO  $4.8 \times 10^{12}$  flops.

If the templates are stored rather than generated as needed, storage requirements range

$1.5 \times 10^{11}$  (TAMA) -  $6.2 \times 10^{14}$  (VIRGO)

<https://journals.aps.org/prd/abstract/10.1103/PhysRevD.60.022002>

Owen+1999

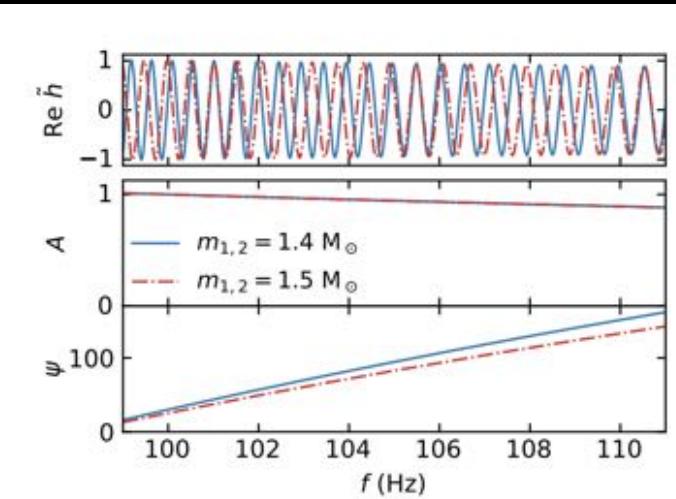


FIG. 1: An example of two waveforms that look very different in the frequency domain (top panel) but have very similar amplitude and phase profiles (middle and bottom panels). The amplitude and phase profiles can be well captured by a low-dimensional linear space spanned by a few basis functions. Waveform amplitudes are shown in arbitrary units.

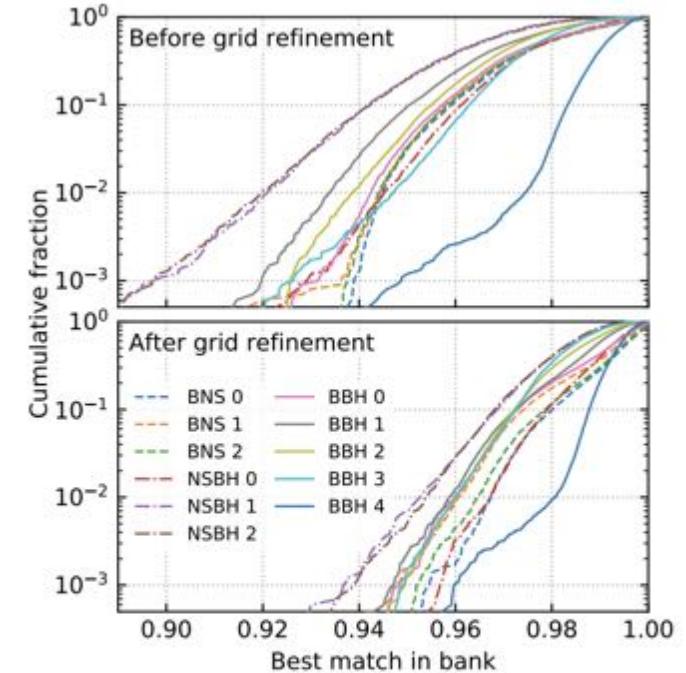
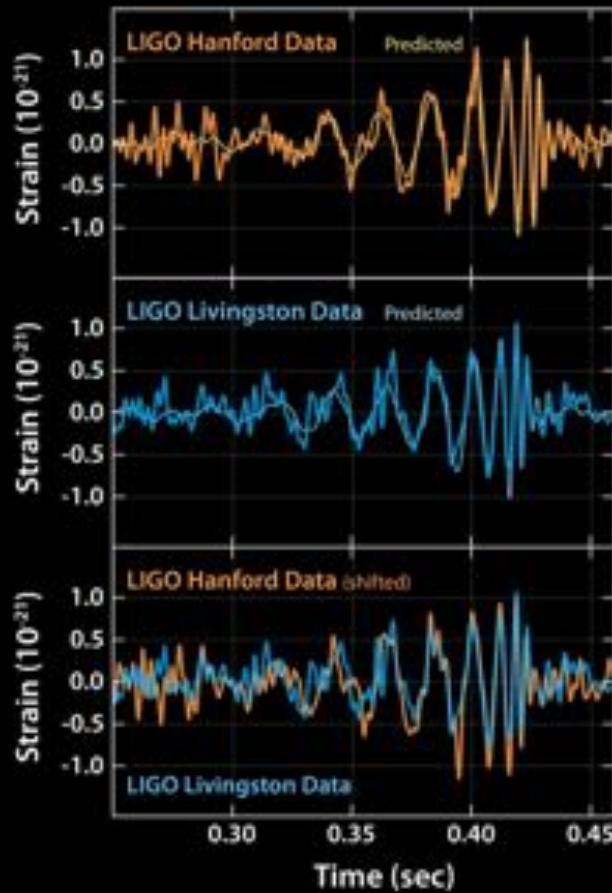


FIG. 5: Effectualness of our template banks, tested on random waveforms drawn from a distribution uniform in individual masses and aligned spins. The vertical axis shows the fraction of the random trials that do not achieve a given match in the bank.

<https://arxiv.org/pdf/1904.01683.pdf>  
Roulet+2019

# Gravitational Wave and MMA



## Historical data releases

All data from Advanced LIGO can be found at [GWOSC](#). GWOSC also hosts data releases from initial LIGO (S5 and S6 observing runs; constraints on GRB051103, and data from a blind injection study).

Some additional historical data (which are not formally LSC data products) are listed below:

Nov 18,  
2015

Localization of Short Duration Gravitational-wave Transients with the Early Advanced LIGO and Virgo Detectors

Apr 23,  
2014

The First Two Years of Electromagnetic Follow-Up with Advanced LIGO and Virgo



## Gravitational Wave Open Science Center

[https://www.gw-openscience.org/eventapi/html/O3\\_Discovery\\_Papers/GW190521/v2/](https://www.gwopenscience.org/eventapi/html/O3_Discovery_Papers/GW190521/v2/)

**3/6**

*space-based  
astronomy  
problems*

# 4-V of Big Data in astronomy

## V1: Volume

Number of bites

Number of pixels

Number of rows in a  
data table x number  
of columns for  
catalogs

## V2: Variety

Diverse science return  
from the same dataset.

Multiwavelength  
Multimessenger

Images and spectra

## V3: Velocity

real time analysis,  
edge computing,  
data transfer

# space based Hubble deep field:

When in 1993 Dr. Robert "Bob" Williams became the new director of the Space Telescope Science Institute

The new director realized that a facility so powerful could maximize its scientific return only if its unique data were made available to the whole world community immediately after acquisition, a dramatic paradigm shift in the way astronomical observatories had been operated until then. In a bold move, in December 1995 Dr. Williams used his "Director Discretionary Time", the fraction of a telescope's observing time that the director can use for special projects, for an unprecedented experiment: stare on the same spot on the sky for ten consecutive days to image the faintest and most distant astronomical sources ever unveiled by humans: the Hubble Deep Field was born. He made the images available to anyone who had the curiosity to look and study them and the findings turned out to be nothing short of transformative.

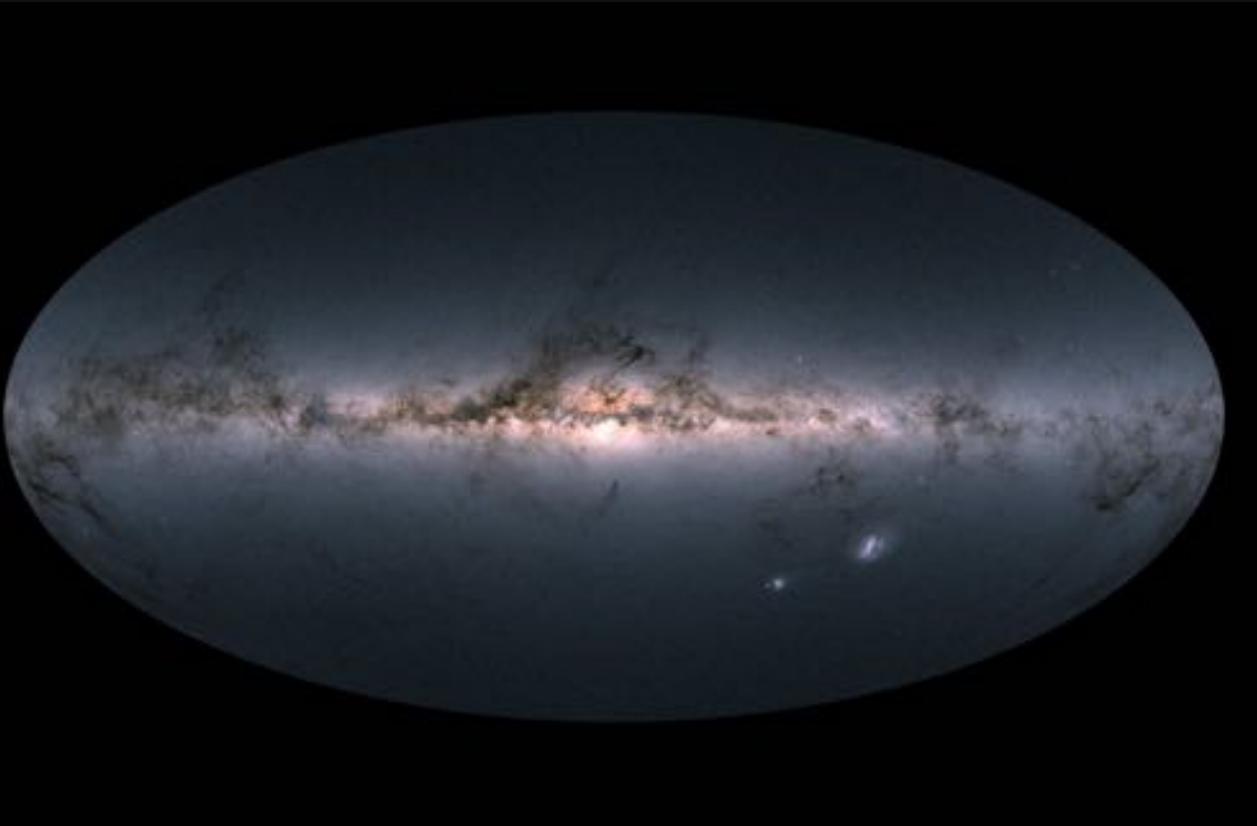
Bob decided that this data was so overwhelmingly powerful, in terms of what it was telling us about the universe, that it was worth it for the community to be able to get their hands on the data immediately. And so the original deep field team processed the data, found the objects in it, and then catalogued each of them, so that every object in the deep field had a description in terms of size, distance, color, brightness and so forth. And that catalogue was available to researchers from the very start--it started a whole new model, where the archive does all the work.

# γ- and X-ray

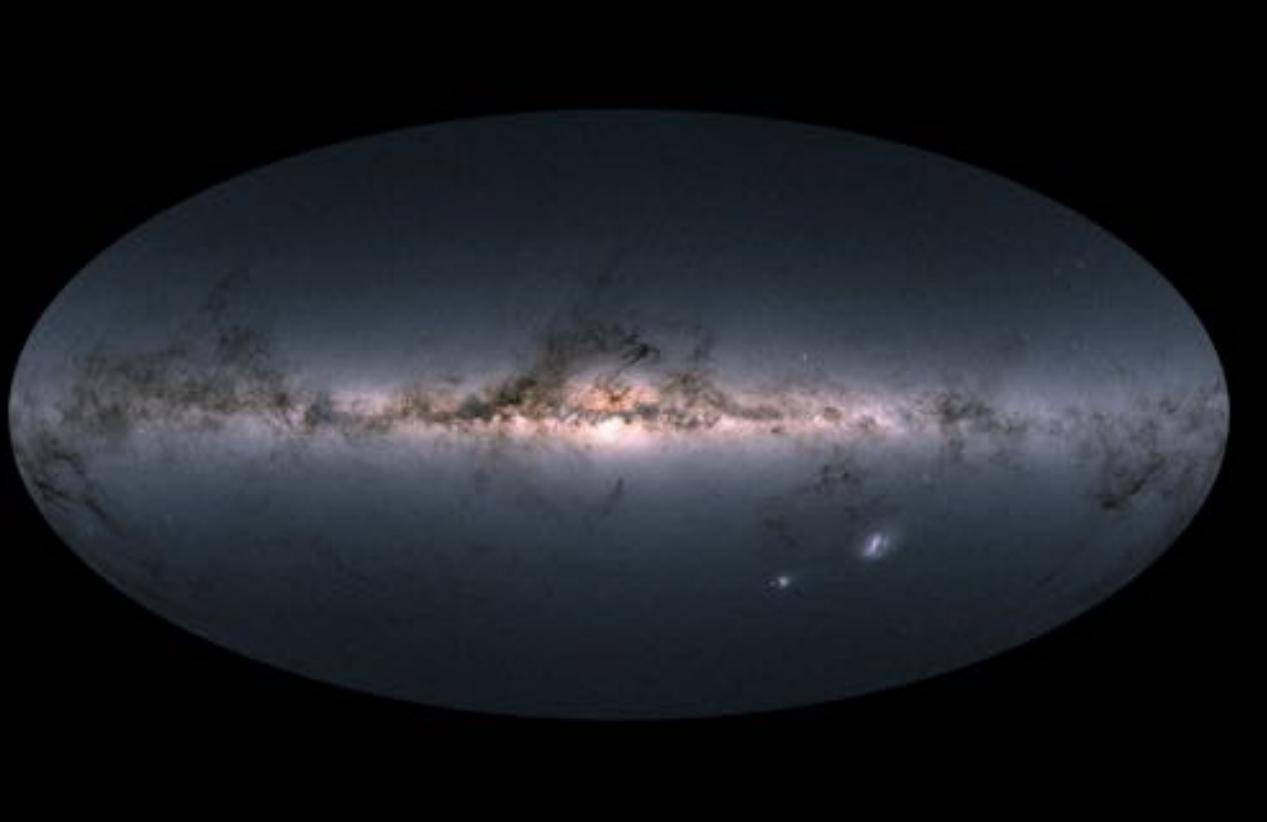
**Table 1.** Main data for the most important all-sky and large-area astronomical surveys providing multi-wavelength photometric data. Catalogues are given in the order of increasing wavelengths.

Survey, catalogue	Years	Spectral range	Sky area (deg <sup>2</sup> )	Sensitivity (mag/mJy)	Number of sources	Density (obj/deg <sup>2</sup> )
Fermi-GLAST	2008–2014	10 MeV–100 GeV	All-sky		3033	0.07
CGRO	1991–1999	20 keV–30 GeV	All-sky		1300	0.03
INTEGRAL	2002–2014	15 keV–10 MeV	All-sky		1126	0.03
ROSAT BSC	1990–1999	0.07–2.4 keV	All-sky		18,806	0.46
ROSAT FSC	1990–1999	0.07–2.4 keV	All-sky		105,924	2.57
GALEX AIS	2003–2012	1344–2831Å	21,435	20.8 mag	65,266,291	3044.85

# Gaia: big data and telemetry do not get along...

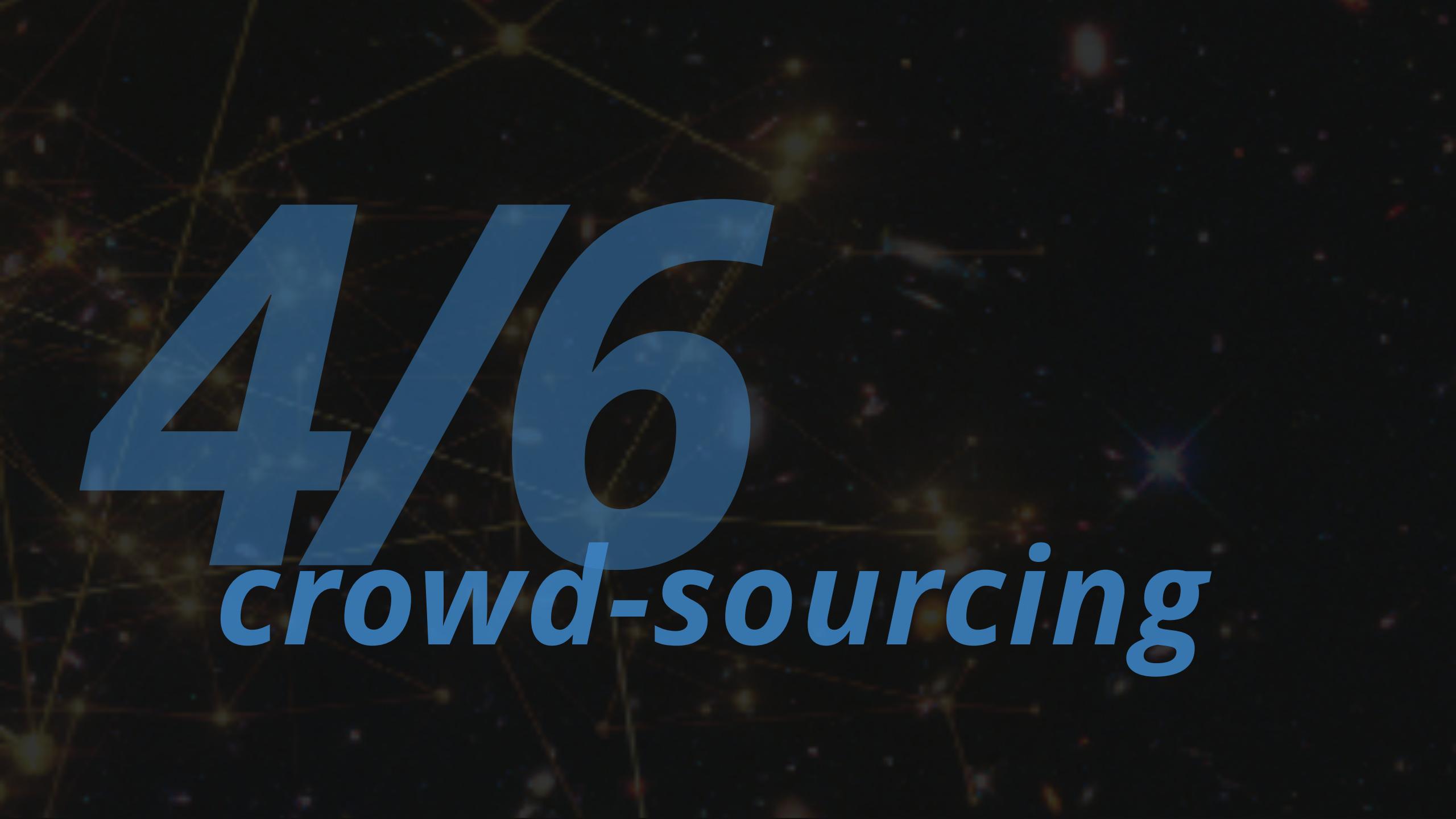


# Gaia: big data and telemetry do not get along...



Besides, the whole content of all CCDs cannot be downloaded because of the telemetry bottleneck: the content of the Astrometric focal plane CCDs only would already amount to about 6000 Mbps, while the actual bandwidth reaches a few Mbps only.

Let us assume that about 55 stars/s will be observed on the average in each of the CCDs. *Detecting* then *windows* each star with say  $6 \times 12$  pixels of 16 bits each, imply that 10 Mbps would be needed, already gaining a factor 600. Then, because one-dimensional measurements allow for achieving the astrometric performances, a factor 12 can be gained by a  $1 \times 12$  binning<sup>1</sup>. Besides, this *sampling* greatly improves the signal to noise ratio. Finally, a factor  $> 2$  can perhaps be gained with a lossless compression scheme Portell et al. (2005).



4/6

*crowd-sourcing*

# 4-V of Big Data in astronomy

## V1: Volume

Number of bites

Number of pixels

Number of rows in a  
data table x number  
of columns for  
catalogs

## V2: Variety

Diverse science return  
from the same dataset.

Multiwavelength  
Multimessenger

Images and spectra

## V3: Velocity

real time analysis,  
edge computing,  
data transfer

## V4: Veracity

This V will refer to  
both data quality  
and availability  
(added in 2012)

# crowd sourcing



[home](#) | [project summary](#) | [people](#) | [gallery](#) | [news](#) | [related links](#) | [bibliography](#) | [data](#) | [use](#) | [download](#) | [forum](#)

## Gallery of Solved Images

In the images below, the red circles are stars our algorithm automatically detects in the image, and the green circles are stars from our master index which appear in the query image. Nebulae, constellations and other objects can be automatically overlayed on the image after it has been solved.

A shot of the Great Nebula, by Jerry Lodriguss (c.2006), from [astropix.com](#)



# crowd sourcing



Galaxy Zoo

Language English

ABOUT

CLASSIFY

TALK

COLLECT

EAGLE galaxies have landed! Read the blog to find out more about them and what to do if some of them appear clumpy.



LSST: 2x3.2 GPix images/minute for 10 years:

scaling the Galaxy Zoo the entire population of the Earth would be insufficient to study the full dataset.

# crowd sourcing

An all-sky search for continuous wave signals in the frequency range 50-1190 Hz + with frequency derivative range from -20e-10 Hz/s to 1.1e-10 Hz/s collected in 2005-2007 during the fifth LIGO science run.

Hundreds of thousands of host machines, that contributed a total of approximately 25000 CPU

## What is Einstein@Home?

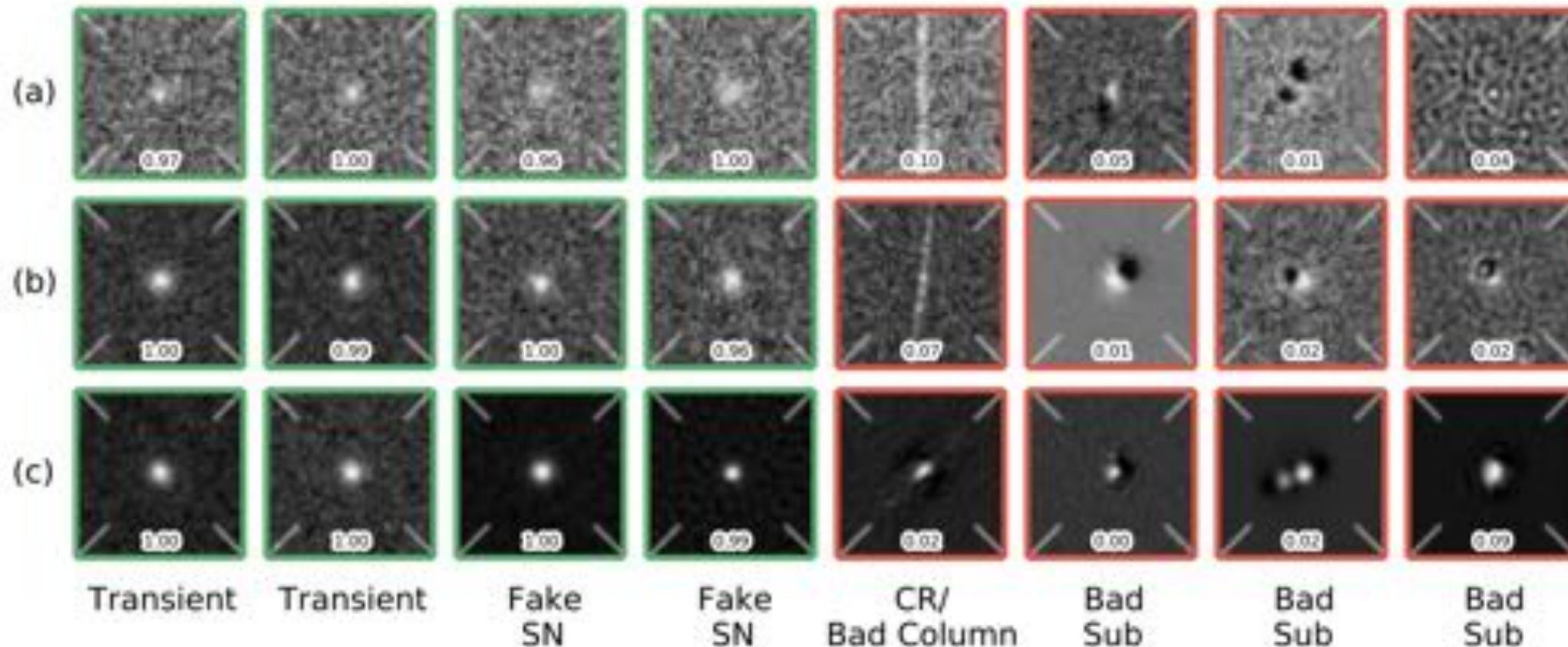
Einstein@Home uses your computer's idle time to search for weak astrophysical signals from spinning neutron stars (often called pulsars) using data from the LIGO gravitational-wave detectors, the Arecibo radio telescope, and the Fermi gamma-ray satellite.

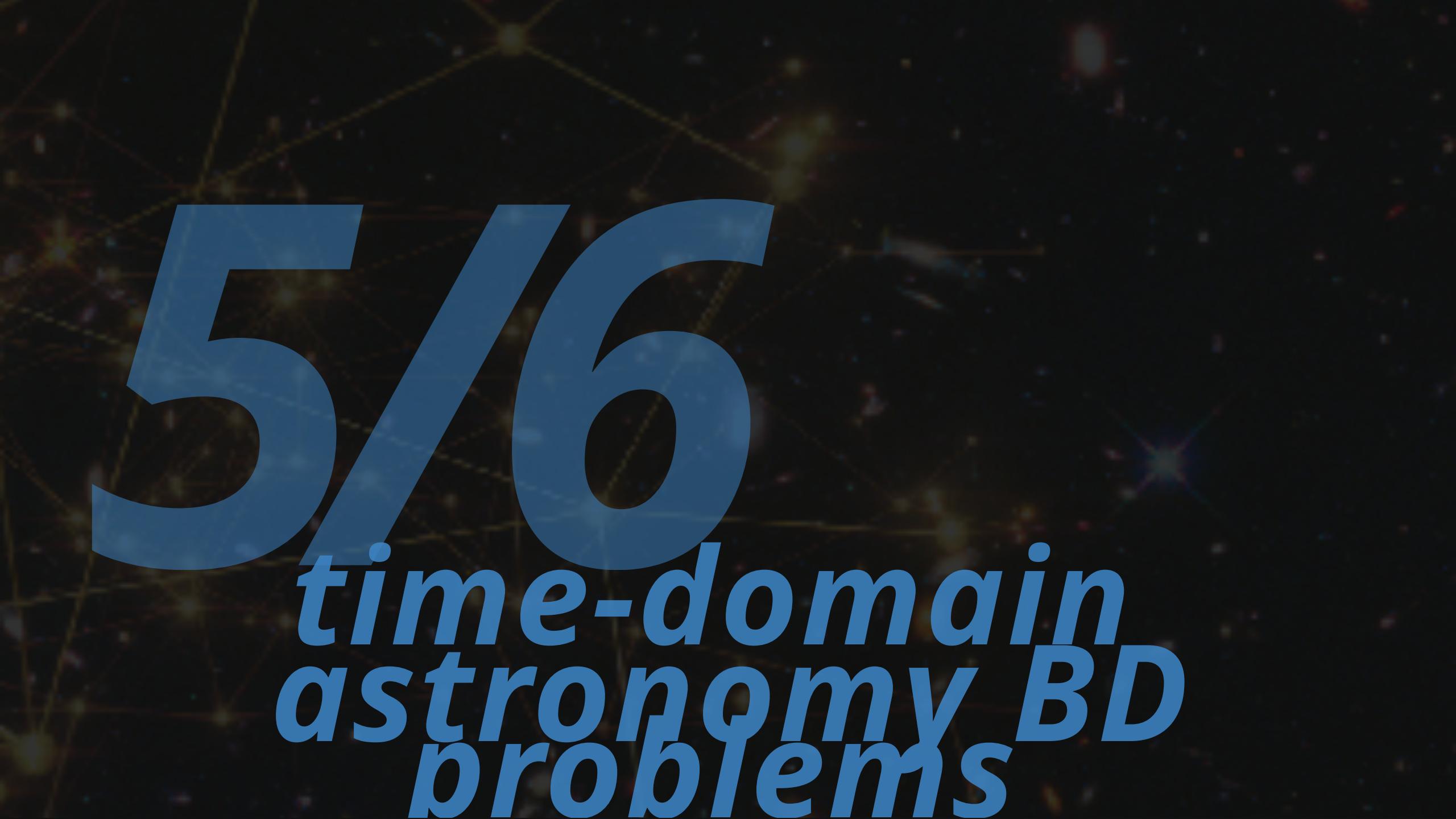
[Learn more](#)

[JOIN NOW](#)

# crowd sourcing

But some of the most crucial problems of veracity  
are ... boring





**5/6**

*time-domain  
astronomy BD  
problems*

# Big data and time domain astrophysics



Rubin  
Observatory

**LSST**  
Legacy Survey of Space and Time

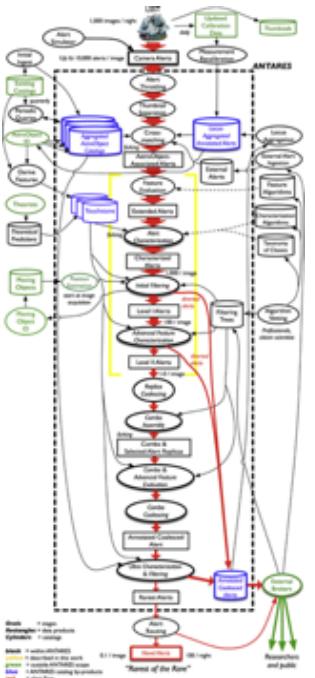


ALeRCE  
Automatic Learning for the  
Rapid Classification of Events



**ANTARES**  
beta

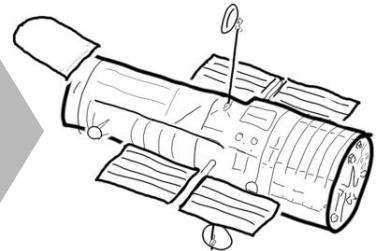
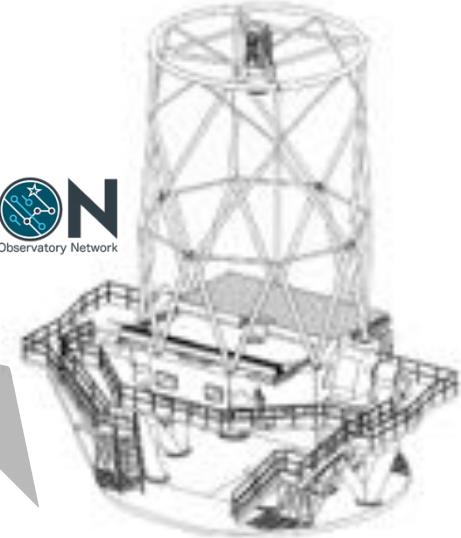
**Lasair**



**TOM**  
TOM TOOLKIT  
Las Cumbres Observatory LCO

the astronomy discovery chain

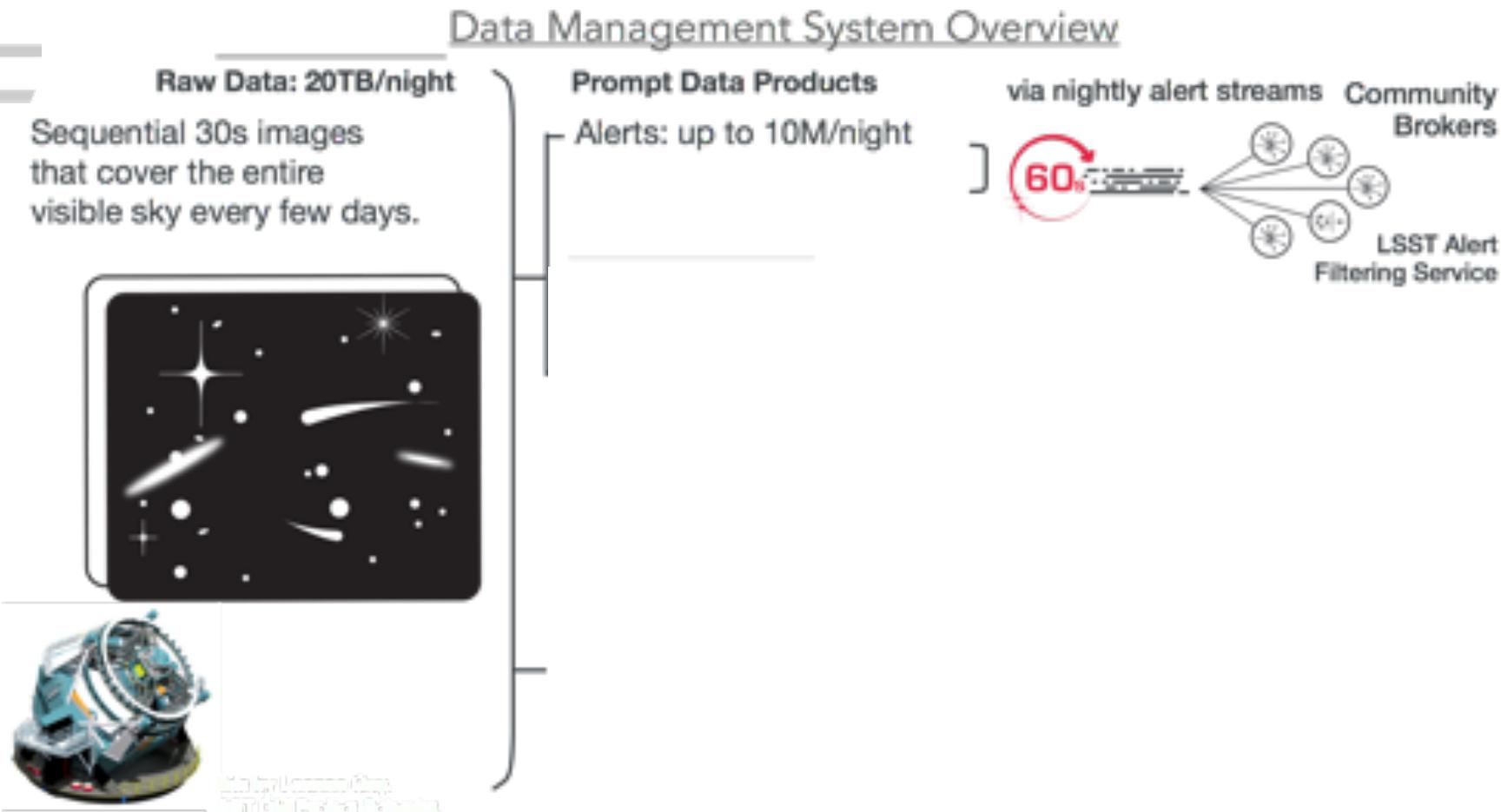
**AEON**  
Astronomical Event Observatory Network



# Big data and time domain astrophysics Discovery



Rubin  
Observatory



<https://www.youtube.com/embed/ZdvEGPt4s0Y?enablejsapi=1>

# Big data and time domain astrophysics

# Discovery



Rubin  
Observatory



ZTF realtime pipeline

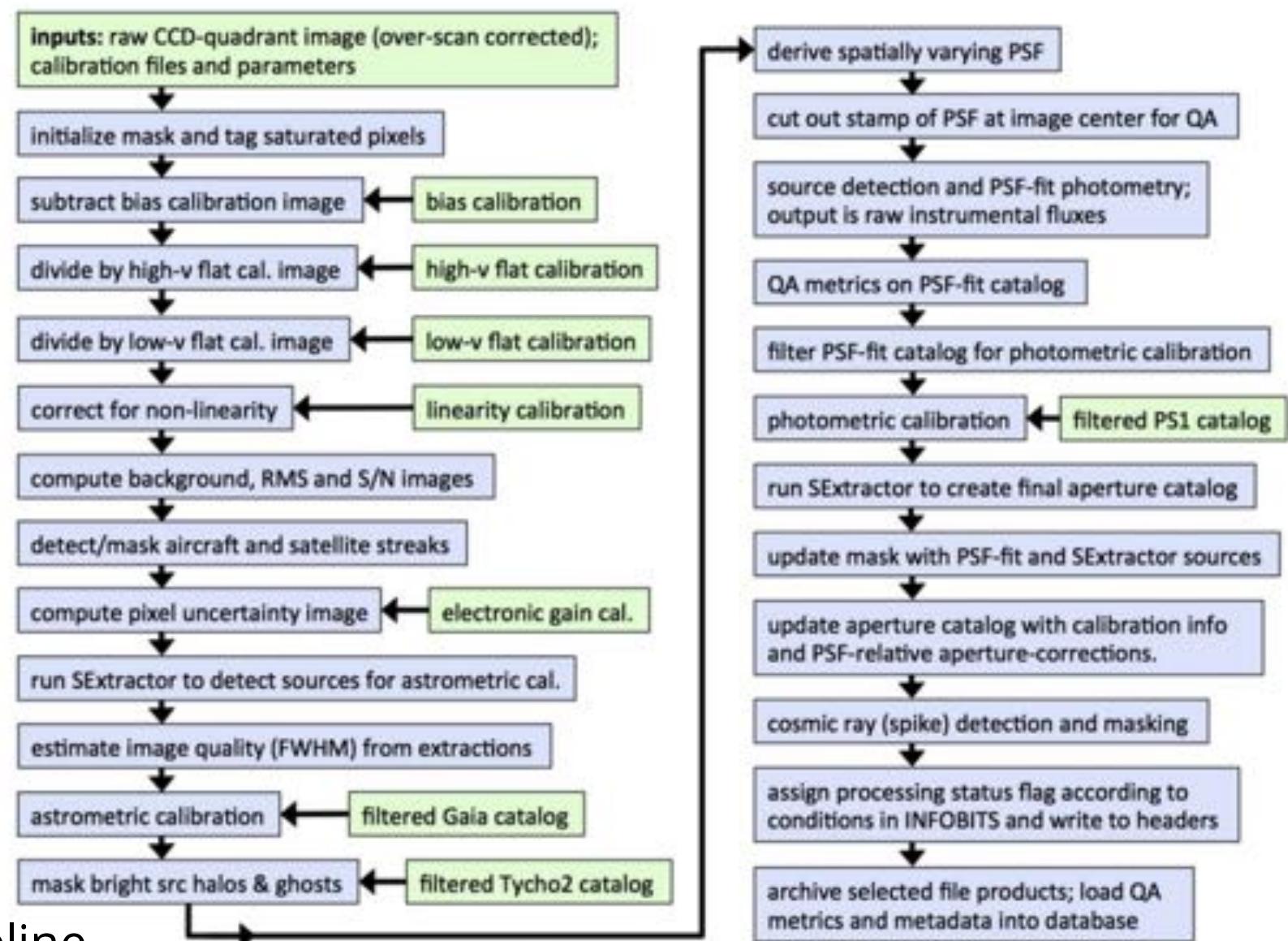


Figure 3. Processing flow in the instrumental calibration pipeline. This represents the first phase of the real-time pipeline.

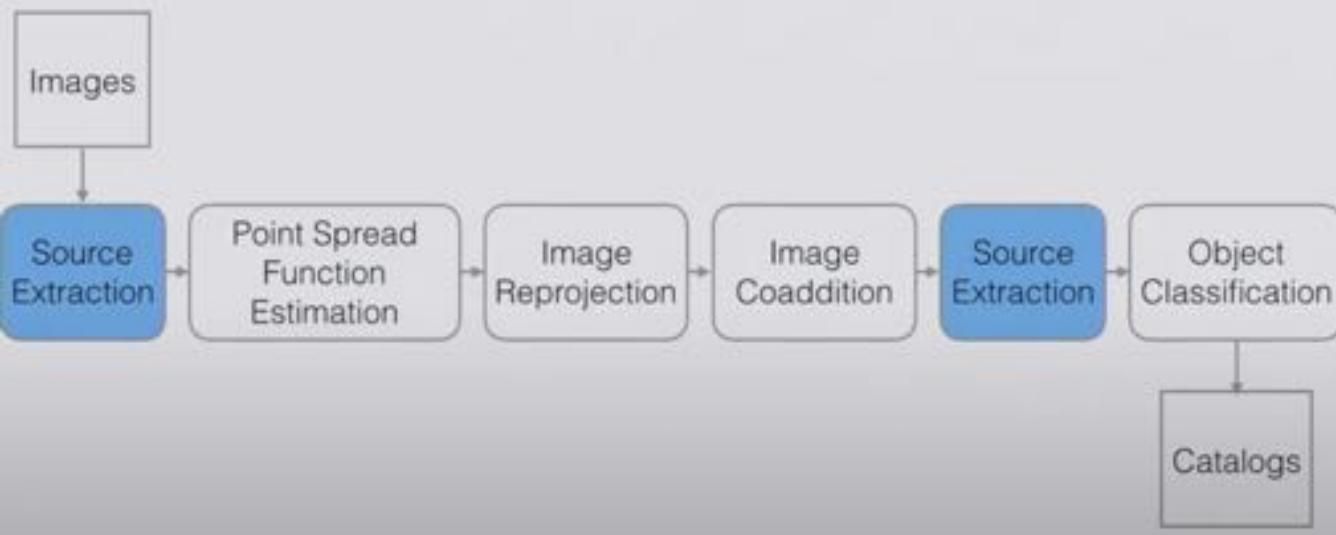
# Demonstrated computational gain by using Big Data platforms

## Kira: Processing Astronomy Imagery Using Big Data Technology

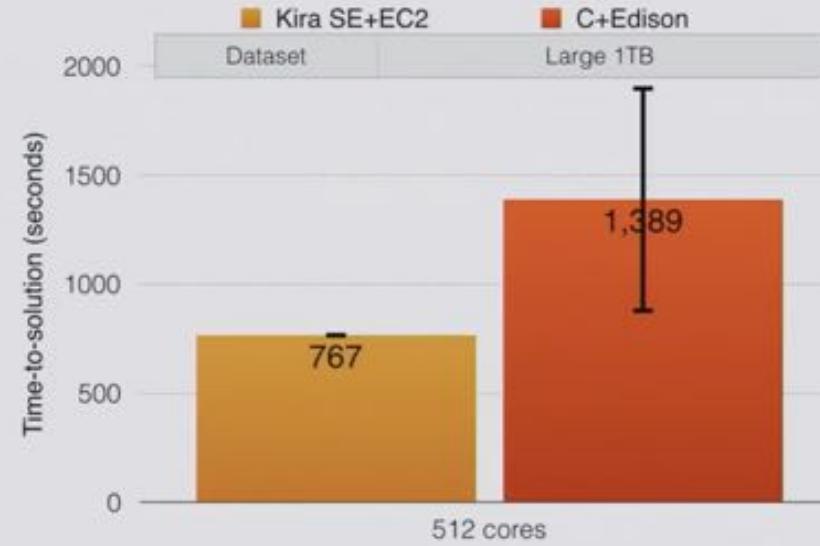
BIDS Spring 2016 Data Science Faire | May 3, 2016 UC Berkeley Zhao Zhang

Particularly because the data (IO) intensive applications are

### A Typical Supernovae Detection Pipeline



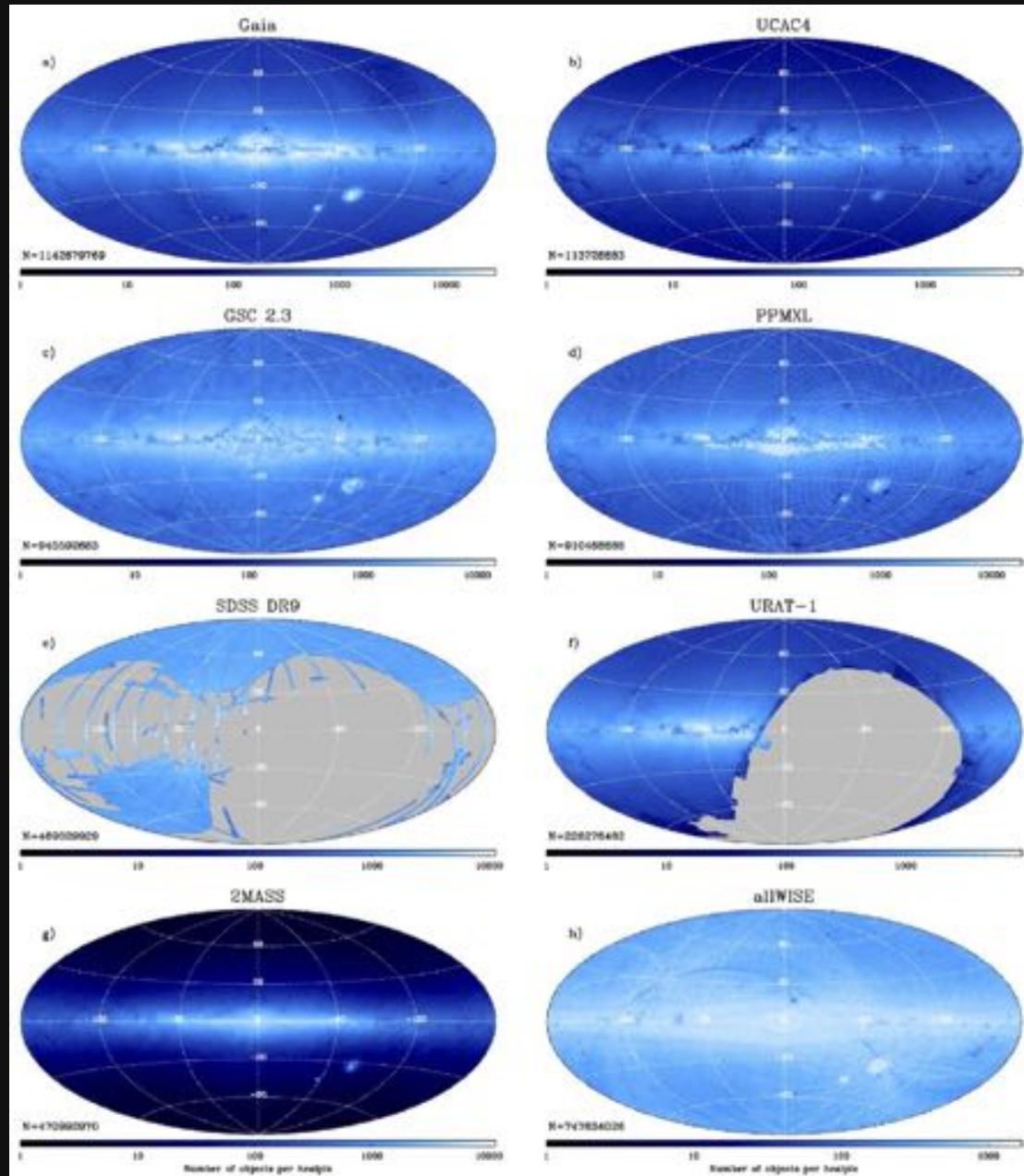
### Kira SE VS. C 1TB Dataset Performance on Supercomputer



# Big data and time domain astrophysics

## Cross-matching

Cross matching catalogs is vital for physical inference. The complexity of the data scales  $\sim$ with the square of the data sources



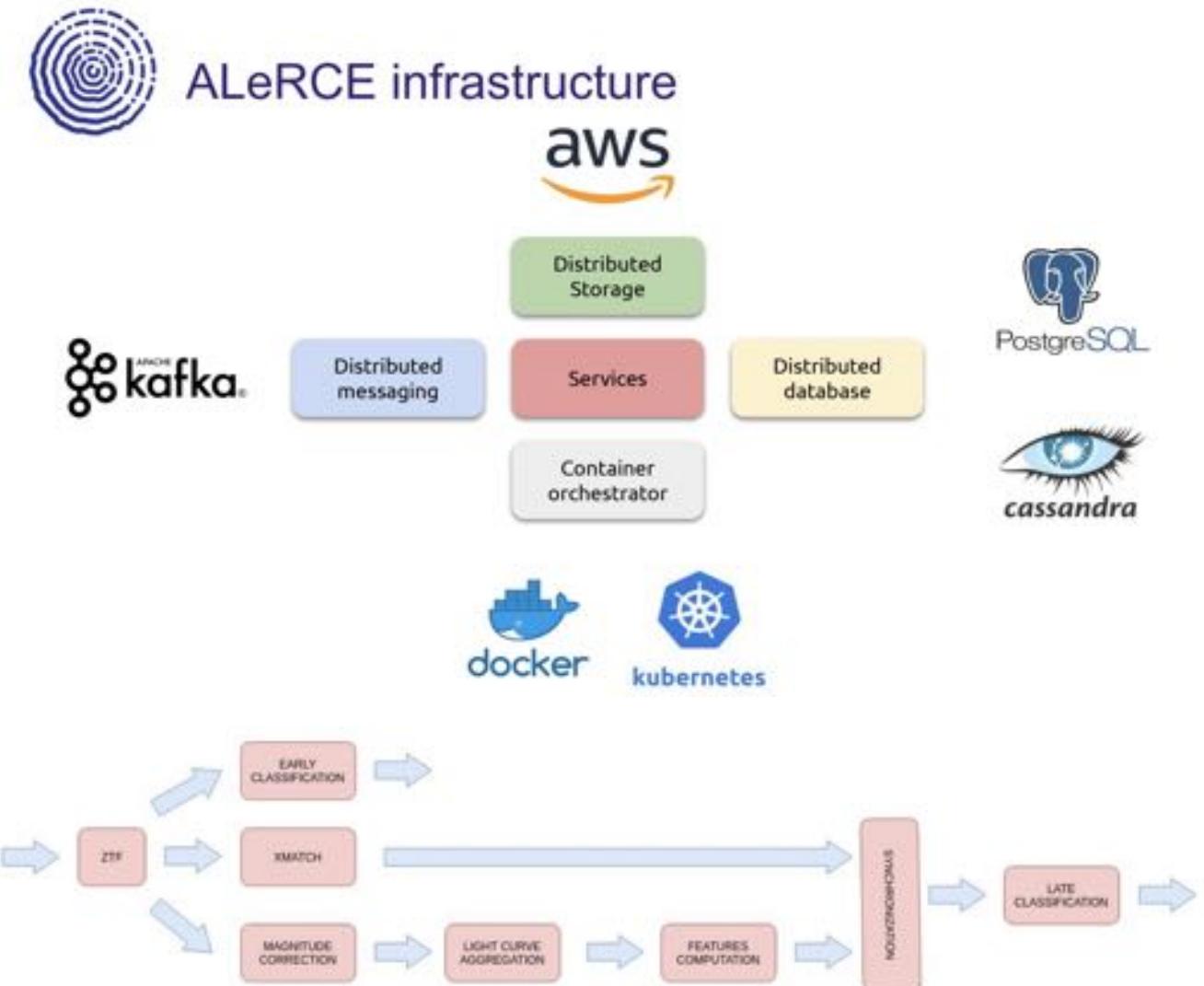
# Big data and time domain astrophysics

## Alert brokers augmentation and classification



# Lasair

<https://cpb-us-e1.wpmucdn.com/sites.northwestern.edu/dist/a/2770/files/2019/08/CASTILLO.pdf>



# Big data and time domain astrophysics



Rubin  
Observatory

**LSST**  
Legacy Survey of Space and Time

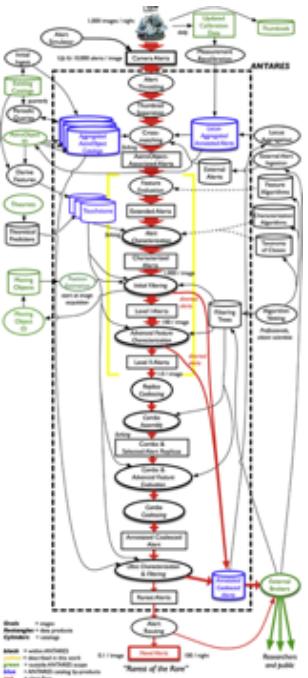


ALeRCE  
Automatic Learning for the  
Rapid Classification of Events



**ANTARES**  
beta

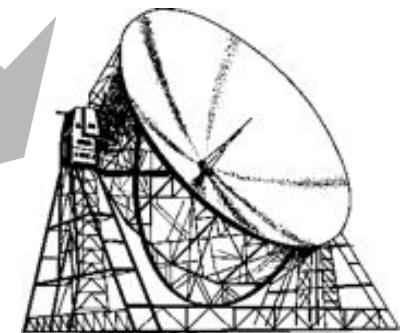
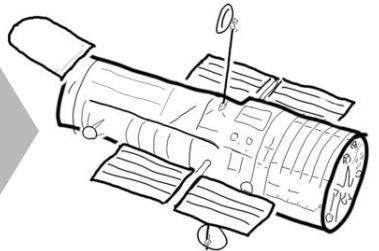
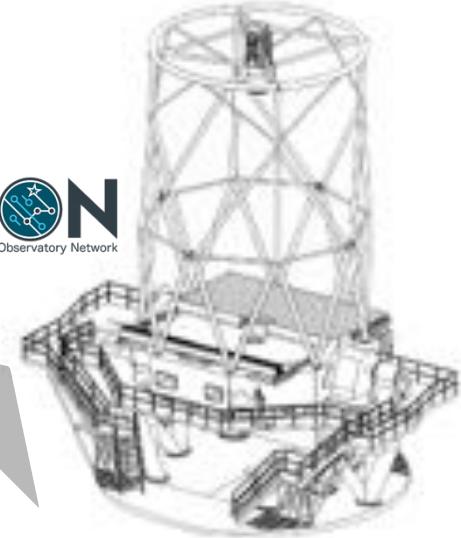
**Lasair**



**TOM**  
TOM TOOLKIT  
Las Cumbres Observatory LCO

the astronomy discovery chain

**AEON**  
Astronomical Event Observatory Network



**6/6**  
*platforms*

# data platforms

The screenshot shows the MAST Portal search interface. At the top, there's a navigation bar with links like "MAST Home", "About MAST", "Contact Us", and "Log In". Below the navigation is a search bar with placeholder text "Search MAST Portal". To the right of the search bar is a "Search" button. The main area features a large orange banner with the text "MAST Portal" in white. Below the banner, a descriptive paragraph reads: "The MAST Portal lets you search multiple collections of astronomical data-sets from one place. Use this tool to find astronomical data, including images, spectra, catalogs, timeseries, publication records and more." On the left side, there are two expandable sidebar menus: "Mission" and "Instrument". The "Mission" menu lists several missions with their counts: Hubble (1794), Spitzer (1774), WISE (177), GALEX (177), SDSS (177), 2MASS (177), and others. The "Instrument" menu lists instruments with their counts: HRC (177), WFC3 (177), WFC3-IR (177), WFC3-GSA (177), WFC3-UVIS (177), and WFC3-IR (177). The central part of the page displays a grid of search results, each row containing a thumbnail image, instrument name, filter, and magnitude. A vertical sidebar on the right contains a "MAST Data" section with a "Data Catalogs" link and a "MAST Images" section with a "Search" link.

# data platforms

<https://www.slideshare.net/databricks/astronomical-data-processing-on-the-lsst-scale-with-apache-spark>



## **AXS - Astronomical Data Processing on the LSST Scale with Apache Spark**

Petar Zečević, SV Group, University of Zagreb  
Mario Jurić, DIRAC Institute, University of Washington

#UnifiedDataAnalytics #SparkAISummit

# data platforms

<https://www.slideshare.net/databricks/astromical-data-processing-on-the-lsst-scale-with-apache-spark>



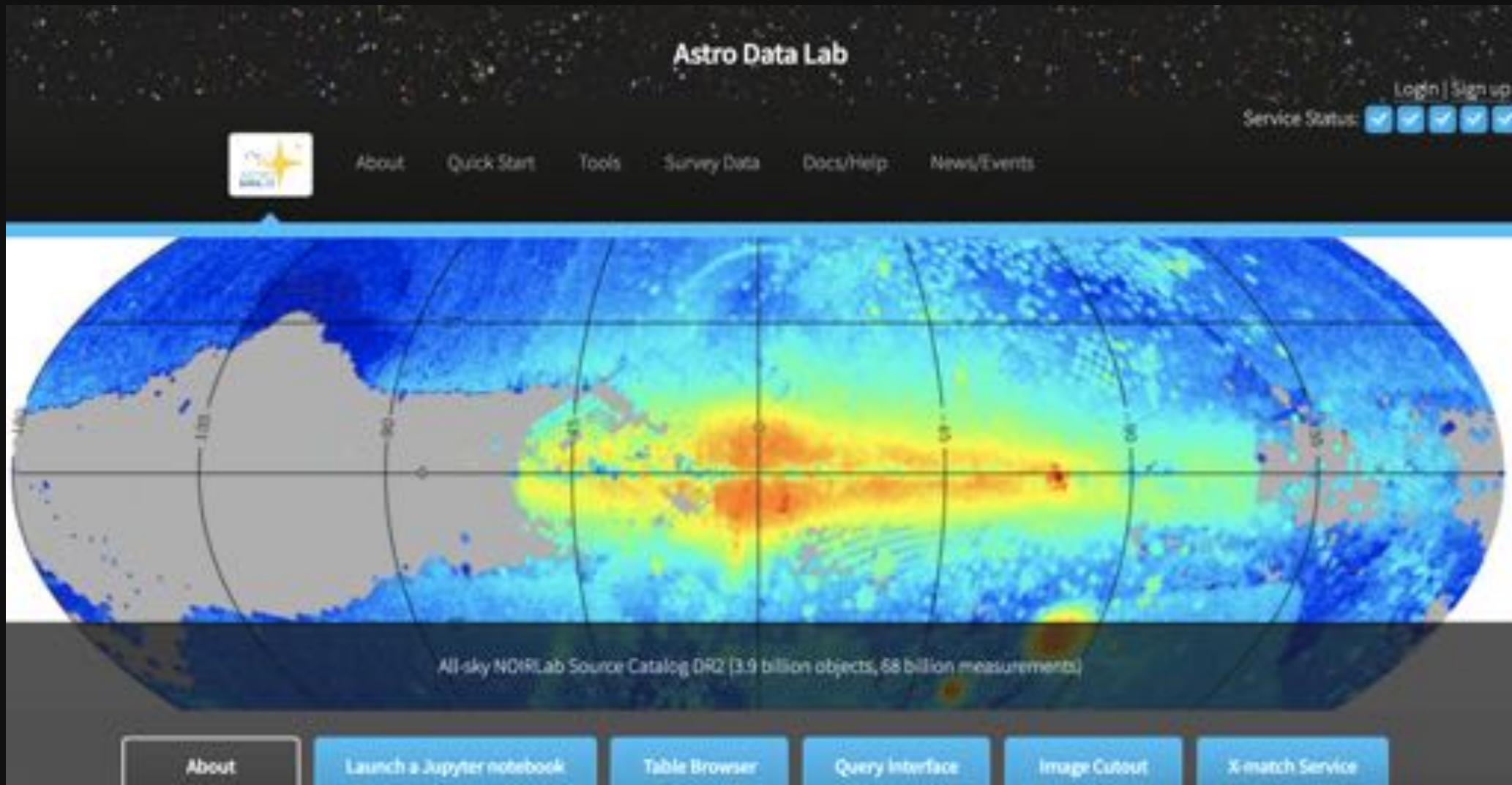
## Enter Spark, AXS

- AXS: Astronomy eXtensions for Spark
- The main idea:
  - Spark is a proven, scalable, cloud-ready and widely-supported analytics framework with full SQL support (legacy support).
  - Extend it to exploratory data analysis.
  - Add a scalable positional cross-match operator
  - Add a domain-specific Python API layer to PySpark
  - Couple to S3 API for storage, Kubernetes for orchestration...
- ... A scalable platform supporting an arbitrarily sized dataset and a large number of users, deployable on either public or private cloud.

# data platforms

<https://www.sciserver.org/>

# data platforms



# data platforms

The screenshot shows a YouTube video player with the following details:

- Title:** Rubin Science Platform: Three Aspects
- Thumbnail:** A portrait photo of a woman with long dark hair.
- Description:** A diagram illustrating the three aspects of the Rubin Science Platform: Portal Aspect, Notebook Aspect, and API Aspect.
- Portal Aspect:** Exploratory analysis and visualization of the Rubin archive.
- Notebook Aspect:** In-depth 'next-to-data' analysis and creation of added-value data products.
- API Aspect:** Remote access to the Rubin archive via industry-standard APIs.
- Platform Interface:** A screenshot of the Rubin Science Platform (RSP) web interface. It features a top navigation bar with icons for PORTAL, NOTEBOOKS, and WEB APIs. Below this is a main menu with icons for DATA RELEASES, ALERT FILTERING SERVICE, USER DATABASES, USER FILES, USER COMPUTING, and SOFTWARE TOOLS.
- Video Player Controls:** Standard YouTube controls for play, volume, and progress.
- Video Statistics:** 79 views, posted on Jan 11, 2021.
- Interaction Buttons:** Like (2), Dislike (0), Share, Save, and more.

[https://www.youtube.com/watch?  
v=4irmRLrNGeE&t=151s](https://www.youtube.com/watch?v=4irmRLrNGeE&t=151s)

VO

Provide and federate content (data, metadata) services, standards, and analysis/compute services – Develop and provide data exploration and discovery tools

[https://ivoa.netastronomers/getting\\_started.html](https://ivoa.netastronomers/getting_started.html)

Provide and federate content (data, metadata) services, standards, and analysis/compute services – Develop and provide data exploration and discovery tools

## **Virtual Observatory Eulogy**

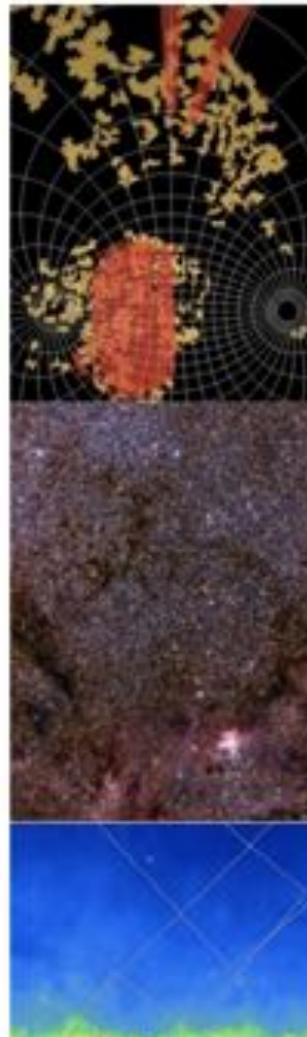
*(de mortibus nil nisi bonum)*

**The Good:**

- Progress on interoperability, standards, etc.
- A global ***data grid of astronomy***
- Empowering a broad community
- Some useful web services
- Community training, outreach
- Better than most other fields (yes!)

**The Not So Good:**

- Data exploration and mining tools
  - That is where the science comes from!
  - Thus, little VO-enabled science
  - Thus, a slow community buy-in



# Are VO's successful in other fields...

*Google Earth Engine*

<https://www.youtube.com/embed/MnCf9Gjz720?enablejsapi=1>

<https://events.asiaa.sinica.edu.tw/school/20170904/talk/djorgovski1.pdf>

# OPINION: what astro did right about BD

FITS files: universal data storage

Strong pressure on making data public

Strong tradition of collaboration

# OPINION: what astro did right about BD

Still lack of trust in cloud services

sparse collaboration between institutes generating  
solutions, a ton of platforms that work differently

slow integration of methods

# Shameless plug: Rubin LSST Science Collaborations

## Four Science Goals

Rubin  
Observatory

### Dark Matter, Dark Energy

- Weak Lensing
- Baryon acoustic oscillations
- Supernovae, Quasars



### Cataloging the Solar System

- Potentially Hazardous Asteroids
- Near Earth Objects
- Object inventory of the Solar System



### Milky Way Structure & Formation

- Structure and evolutionary history
- Spatial maps of stellar characteristics
- Reach well into the halo



### Exploring the Transient sky

- Variable stars, Supernovae
- Fill in the variability phase-space
- Discovery of new classes of transients



# Rubin LSST Science Collaborations



No federally funded LSST science



No science is reserved for any one group



**Since there is no science team,  
science preparation is done by the  
Science Collaborations... for free!**



**1500+ members**

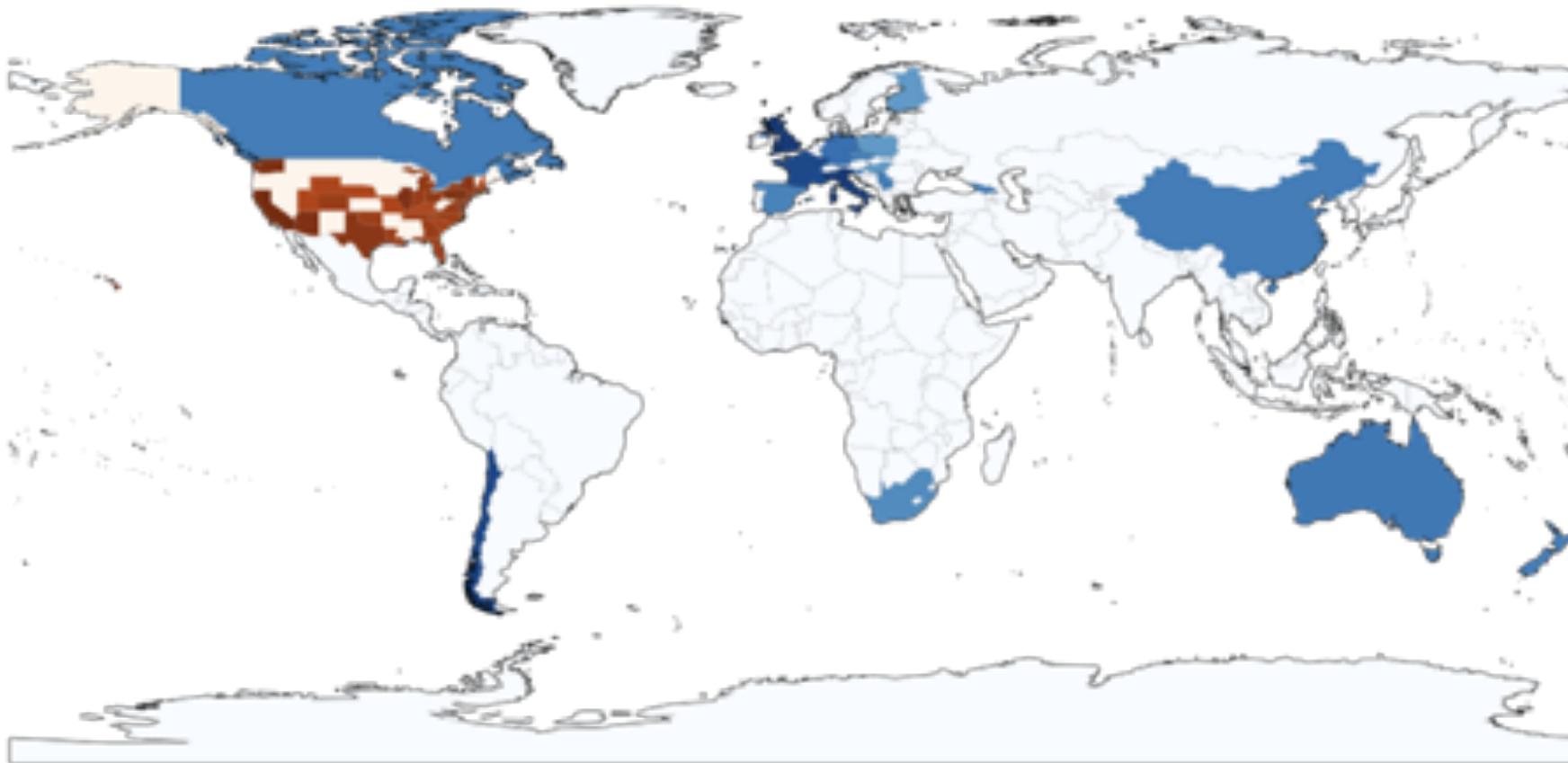
*Science Collaborations*

federica bianco [fbianco@udel.edu](mailto:fbianco@udel.edu)

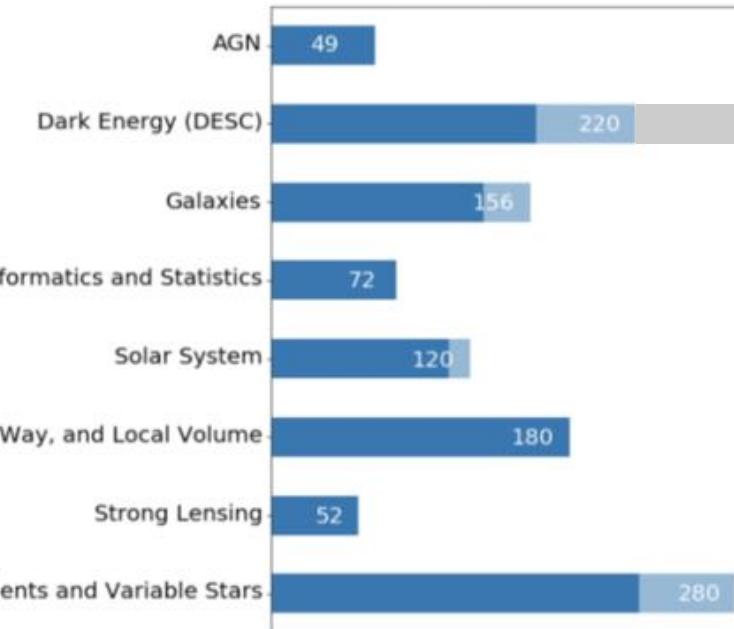
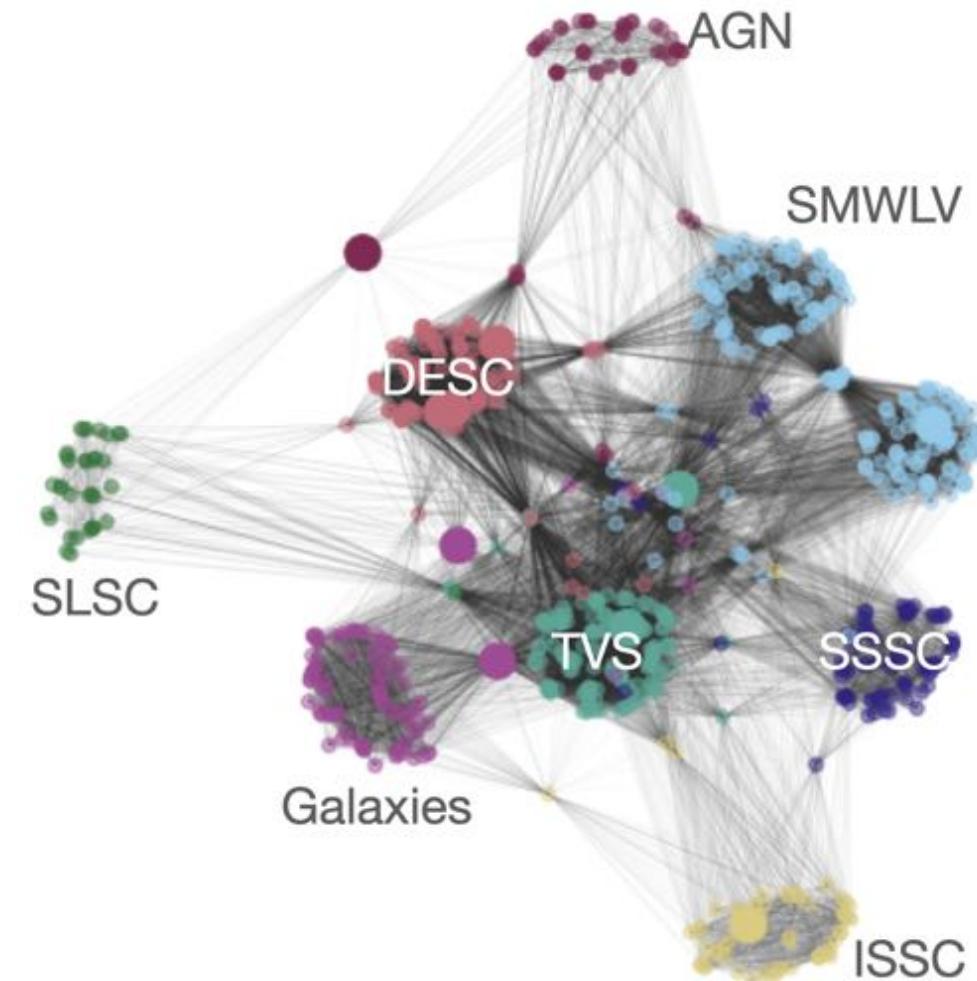
@fedhere



# Rubin Observatory LSST SCs



# Rubin Observatory LSST SCs



# Rubin LSST Science Collaborations



*We aspire to be an inclusive, equitable, and ultimately just group and we are working with renewed vigor in the wake of the recent event that exposed inequity and racism in our society to turning this aspiration into action.*



## #desc-for-black-lives

@heather999 created this channel on June 9th. This is the very beginning of the #desc-for-black-lives channel. Description: Dialogues about how each of us as individual DESC members and our collaboration as a whole can help eradicate anti-Black racism. ([edit](#))

<https://lsst-tvssc.github.io/calltoaction.html>



Diversity Equity and Inclusion council of the SCs

# Thank you!

Federica B. Bianco  
University of Delaware  
Physics and Astronomy  
Biden School of Public Policy and Administration  
Data Science Institute  
NYU Center for Urban Science and Progress

Rubin Observatory LSST Science  
Collaborations Coordinator

Rubin LSST Transients and Variable  
Stars Science Collaborations Chair

please email me if you have questions!  
[fbianco@udel.edu](mailto:fbianco@udel.edu)