



# Science Drivers for Big Data

Joseph Lazio

SKA Program Development Office &  
Jet Propulsion Laboratory, California  
Institute of Technology

# Data Intensive Astronomy



- “There is nothing new under the Sun ....”
- Case studies
  1. Processing: Cosmology and imaging surveys
  2. Data volume: Fundamental physics from pulsar observations and pulsar surveys
  3. Data rates: Observing!
  4. Data visualization: Identifying interference
  5. Data curation: Astronomy for the future

# Data Intensive Astronomy



## Data Volumes



Ἰππάρχος (Hipparchus)

- ca. 135 BCE
- Stellar catalog with 850 entries

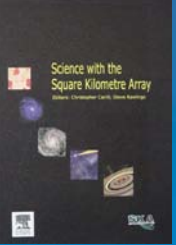
## Computational Limitations



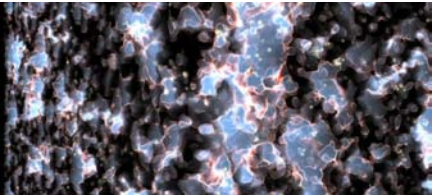
Harvard computers

- Production of stellar plates and spectra (“data rate”) was increasing enormously
- Examined and classified telescope output
- Forerunners of human mathematical computers

Exploring the Universe with the world's largest radio telescope



# Key Science for the SKA (a.k.a. m- and cm- $\lambda$ astronomy)

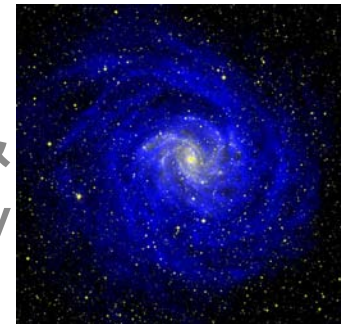


Emerging from the Dark Ages  
& the Epoch of Reionization

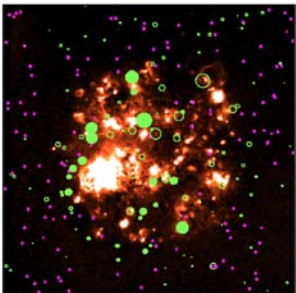
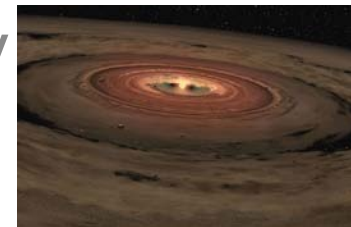


Strong-field Tests of Gravity  
with Pulsars and Black Holes

Galaxy Evolution, Cosmology, &  
Dark Energy



The Cradle of Life & Astrobiology



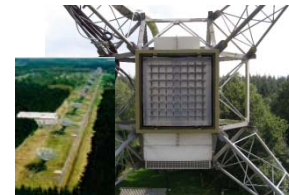
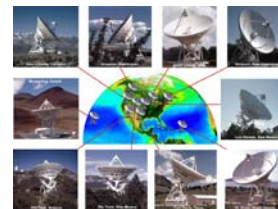
Origin & Evolution of  
Cosmic Magnetism



# SKA Pathfinder



- SKA is ultimate goal, though long-term program
  - Precursors and many pathfinders in existence or under construction
- Data challenges before SKA comes on-line
- Scalability could be an issue

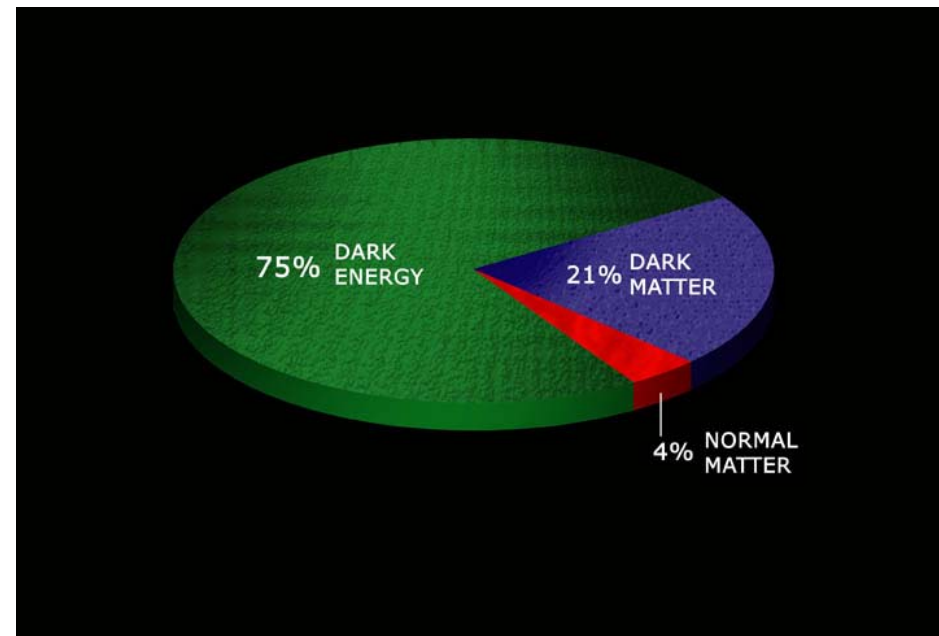


# Case Study 1: Cosmology



## Origin and Fate of the Universe

- Era of “precision cosmology”  
... or precision ignorance
- Need to sample a substantial volume of the Universe



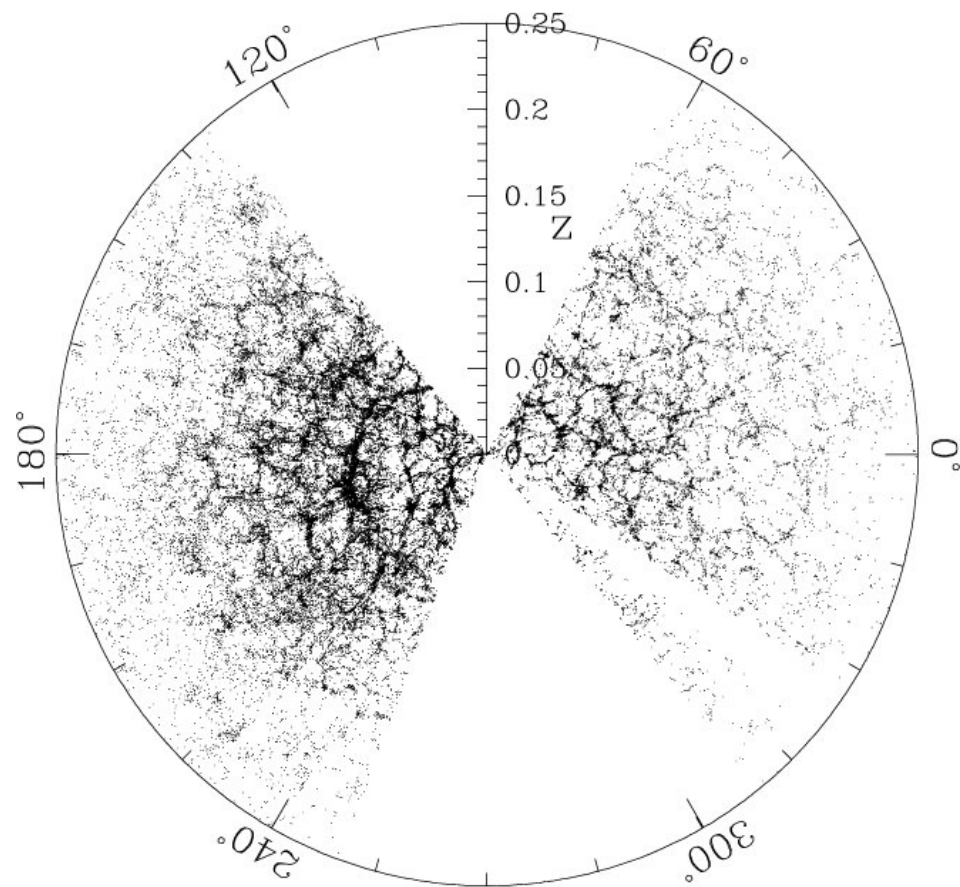
Composition of the Universe

# Cosmology and Sky Surveys



Volume  $\sim D^2 \Delta D \Omega$

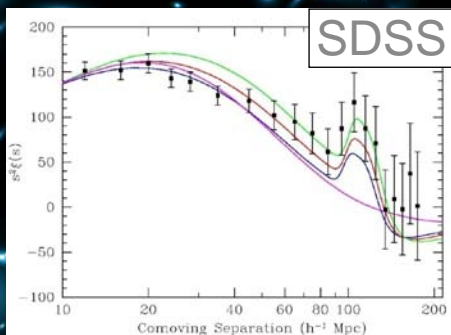
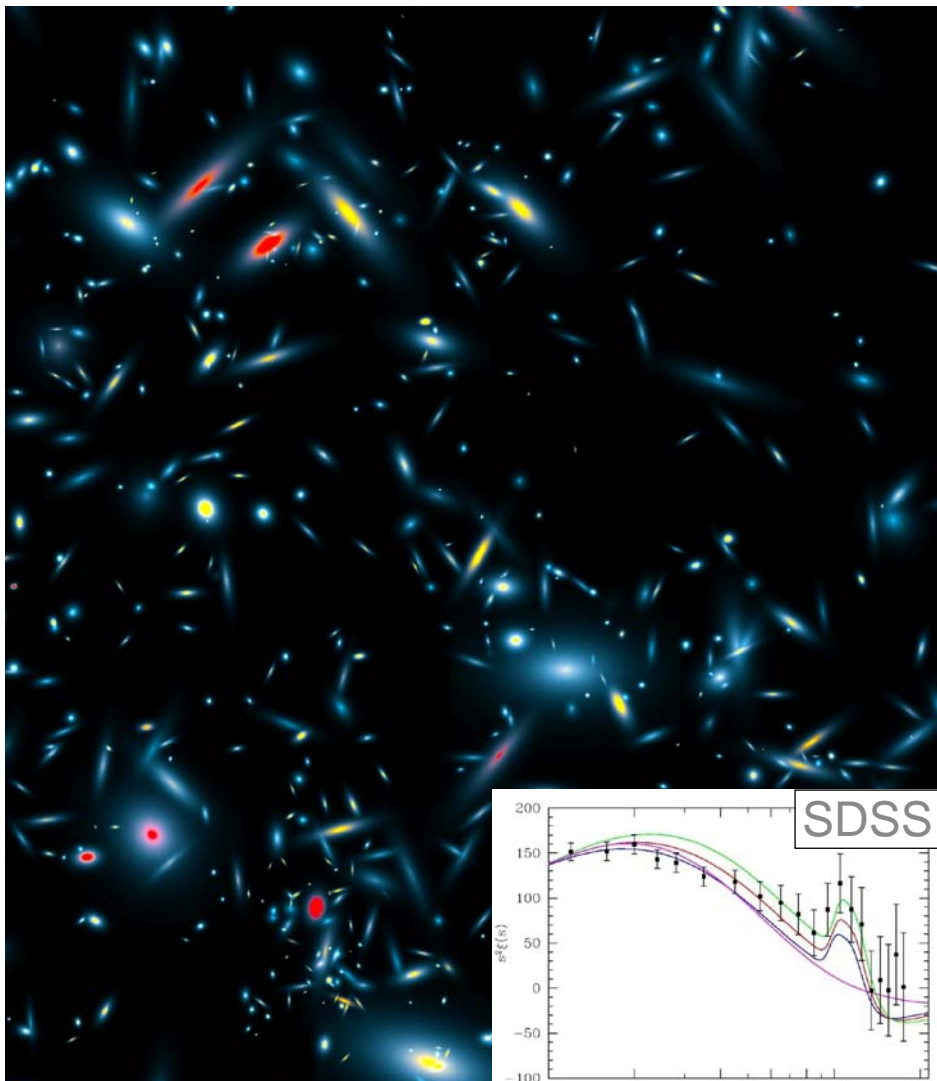
- $D$  – distance;  $\Omega$  – solid angle
- Surveying to larger  $D$  is difficult  $\rightarrow$  need larger telescopes  
“square kilometre” of SKA
- Surveying larger sky areas  $\Omega$  just requires more observing time



Sloan Digital Sky Survey volume



# Cosmology and Sky Surveys



- Image the sky, locating galaxies  
Analysis of locations compared with cosmological models to constrain parameters
- Two broad classes of surveys
  - Continuum: e.g., NVSS, FIRST, *ASKAP/EMU*, *WSRT/APERTIF/WODAN*
  - Spectroscopic: SDSS, Arecibo ALFALFA, *ASKAP/WALLABY*, *SKA H I survey*  
Spectroscopic surveys locate in **3-D space!** very powerful
- Ultimate goal: spectroscopic survey of 1 billion galaxies

Exploring the Universe with the world's largest radio telescope

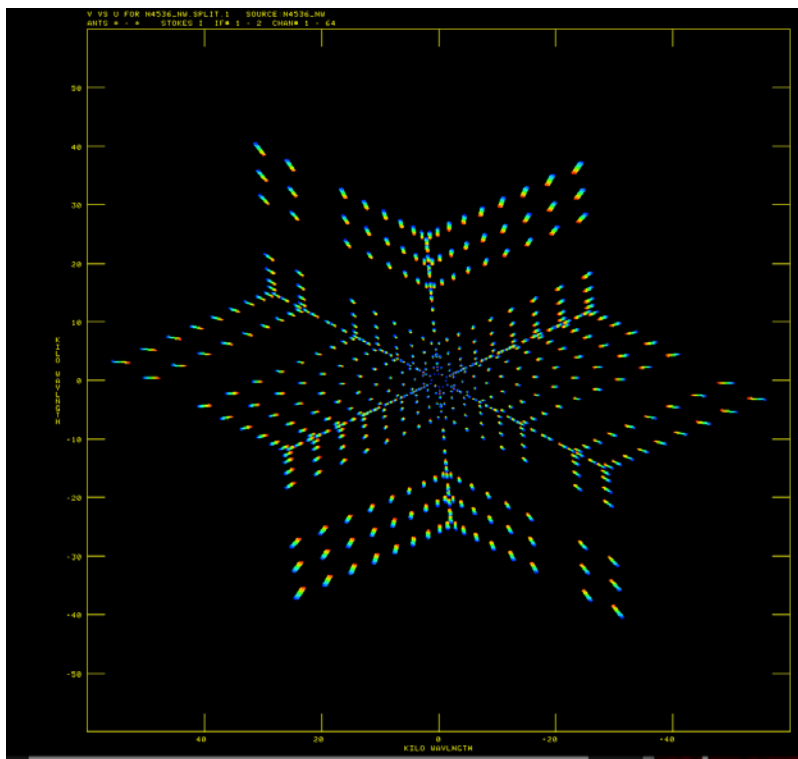




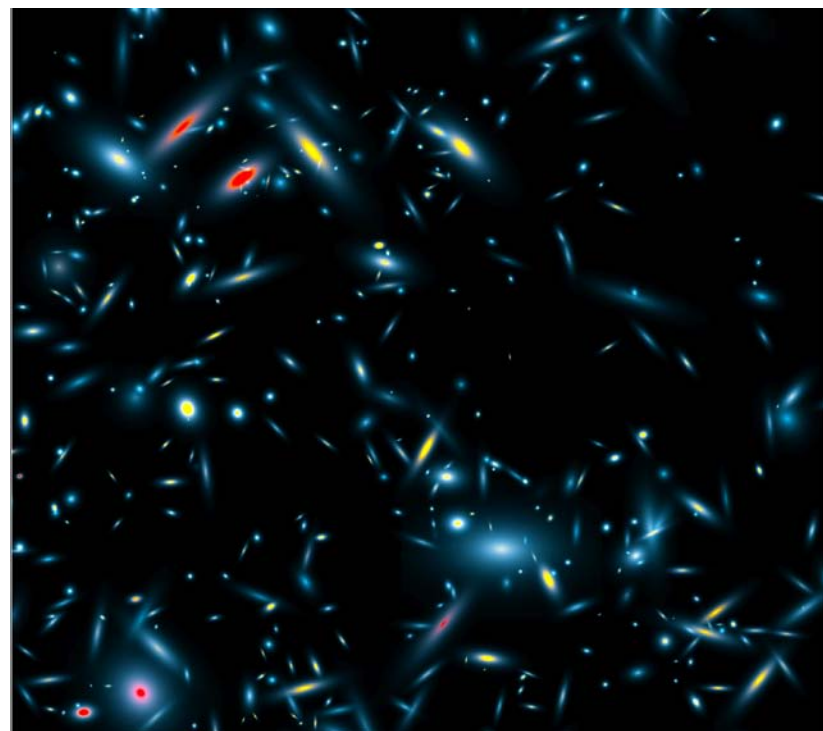
# Imaging with Arrays



## Fourier transform plane



## Image plane



Fourier Transform  
↔

$$N_{\text{data}} \sim N_{\text{antenna}}^2 N_{\text{frequency}} N_{\text{time}}$$

# Imaging Surveys



## Requirements

- Many antennas in order to provide sensitivity and image quality  
large  $N_{\text{antenna}}$
- Spectral resolution because of wide-field effects, line emission from galaxies, or both  
large  $N_{\text{frequency}}$
- Long integrations in order to obtain adequate signal-to-noise ratio  
large  $N_{\text{time}}$ , e.g., 500 hr at 1 s sampling?
- $N_{\text{data}} \sim N_{\text{antenna}}^2 N_{\text{frequency}} N_{\text{beams}} N_{\text{time}}$

| ASKAP                                  | SKA Phase 1                             | SKA Phase 2                             |
|--|---|---|
| $N_{\text{antenna}} = 30$              | $N_{\text{antenna}} \sim 250$           | $N_{\text{antenna}} \sim 1000$          |
| $N_{\text{beams}} = 30$                | $N_{\text{beams}} = 1$                  | $N_{\text{beams}} = 1?$                 |
| $N_{\text{frequency}} \sim 16\text{k}$ | $N_{\text{frequency}} \sim 16\text{k}?$ | $N_{\text{frequency}} \sim 16\text{k}?$ |

# Imaging Surveys II



| ASKAP  | SKA Phase 1  | SKA Phase 2  |
|--|--|--|
| $N_{\text{antenna}} = 30$                                | $N_{\text{antenna}} \sim 250$                            | $N_{\text{antenna}} \sim 1000$                           |
| $N_{\text{beams}} = 30$                                  | $N_{\text{beams}} = 1$                                   | $N_{\text{beams}} = 1?$                                  |
| $N_{\text{frequency}} \sim 16\text{k}$                   | $N_{\text{frequency}} \sim 16\text{k}?$                  | $N_{\text{frequency}} \sim 16\text{k}?$                  |
| $N_{\text{time}} \sim 1.8\text{M}$                       |  |  |
| $N_{\text{data}} \sim 8 \times 10^{14}$                  | $N_{\text{data}} \sim 2 \times 10^{15}$                  | $N_{\text{data}} \sim 3 \times 10^{16}$                  |
| <b><math>N_{\text{OPS}} \sim 8 \times 10^{18}</math></b> | <b><math>N_{\text{OPS}} \sim 2 \times 10^{19}</math></b> | <b><math>N_{\text{OPS}} \sim 3 \times 10^{20}</math></b> |

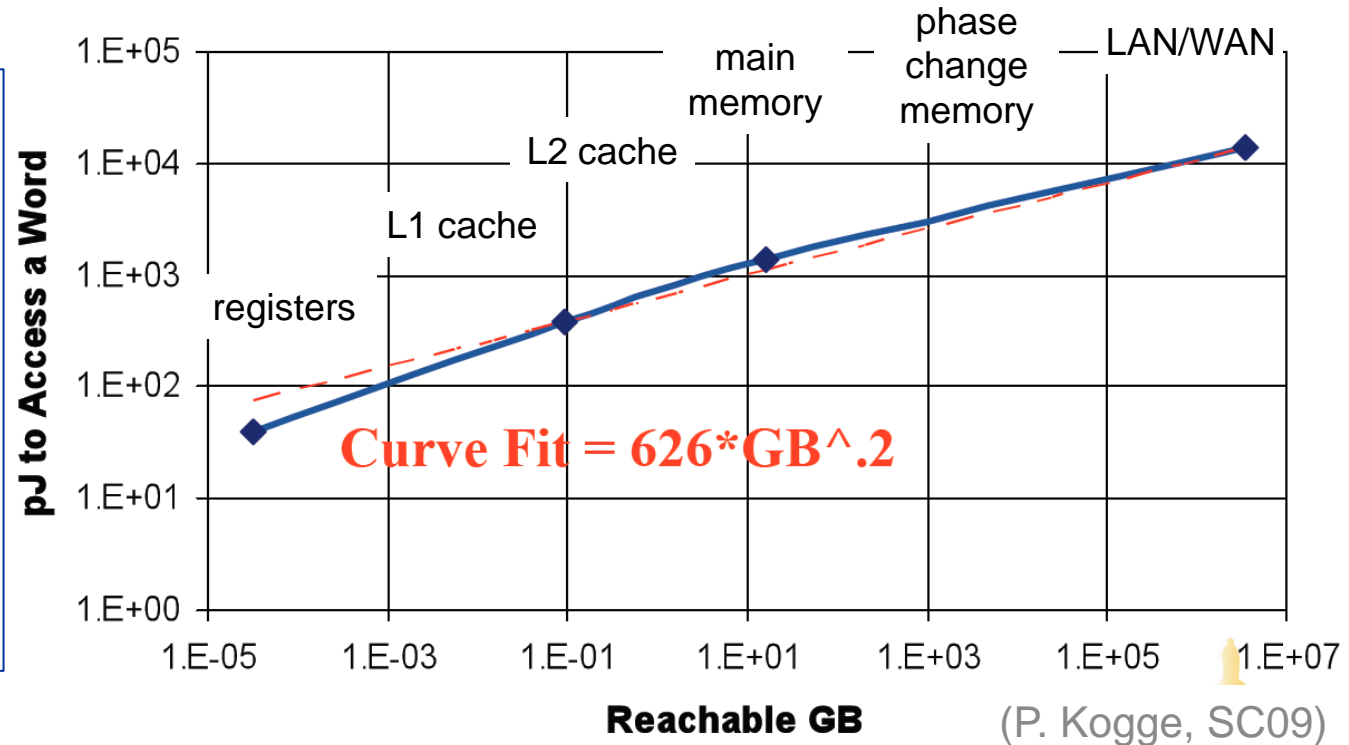
- Imaging is more than “just” an FFT.
  - Gridding, deconvolution, wide-field corrections, self-calibration, ...
- Community estimates are  $10^4$  to  $10^5$  ops per visibility datum(!)

# Energy Consumption



- Really scary could be power consumption

Early estimates for SKA power were 100 MW



- Problem linked to moving data around chips
- See also D'Addario (SKA Memo #130)

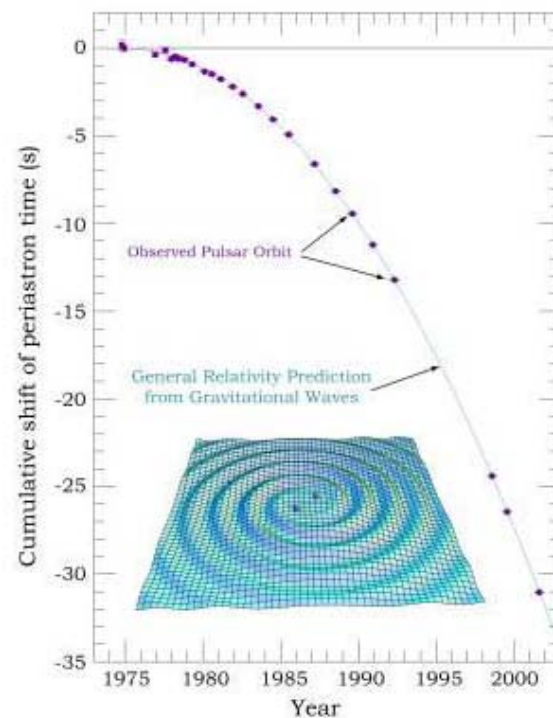


# Case Study 2: Fundamental Physics with Radio Pulsars

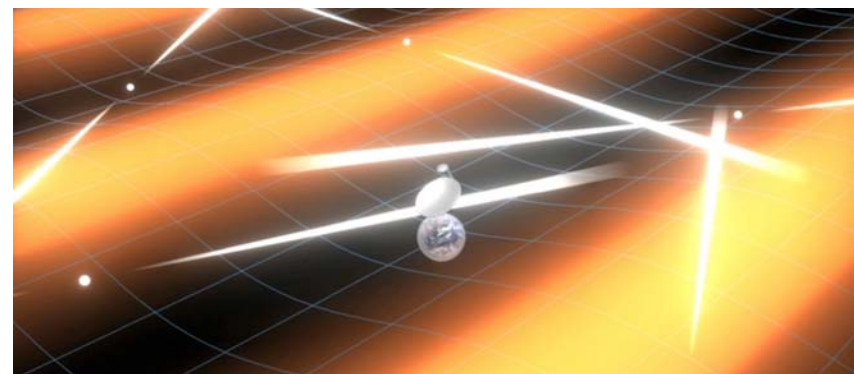


Arrival times of pulses from radio pulsars can be measured with phenomenal accuracy

- Better than 100 ns precision in best cases
- Enables high precision tests of fundamental physics
  - Theories of gravity, gravitational waves, nuclear equation of state
  - 1993 Nobel Prize in Physics
- Problem: Not all pulsars are equal!
- Good “timers” < 10% of total population
- Need to find **many** more!
- All-sky survey



← Ultra-relativistic binaries & gravitational wave studies



# Pulsar Surveys I



## Requirements

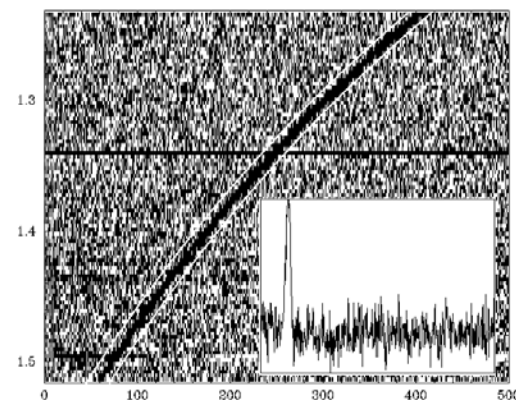
- Large bandwidths because pulsars are faint
- Long integration times because pulsars are faint
- Rapid time sampling in order to resolve pulse profile
- Narrow frequency channelization in order to mitigate interstellar scattering
- For a “pixel” on the sky, accumulate data for a time  $\Delta t$  over a bandwidth  $\Delta \nu$

Suppose  $\Delta t = 20$  min.,  $\Delta \nu = 800$  MHz

- Time sampling  $\delta t$  with frequency channelization  $\delta \nu$

For GBT GUPPI,  $\delta t = 81.92 \mu\text{s}$ ,  $\delta \nu = 24$  kHz

- 60 GB data sets per pixel ...



# Pulsar Surveys II



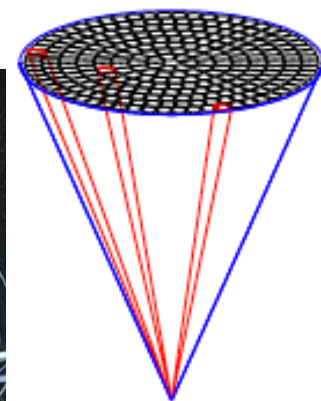
## For GBT

- At 800 MHz, “pixel”  $\sim 16' = 0.3^\circ$
- About 350 kpixels in the sky
- 20 PB data set



## For SKA

- At 800 MHz, “pixel” =  $1.2'$
- About 76 Mpixels in the sky
- 4.6 EB data set

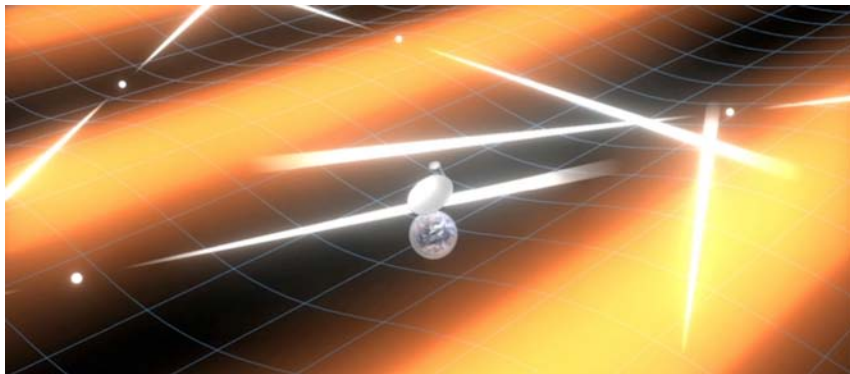


# Case Study 3: Observing!

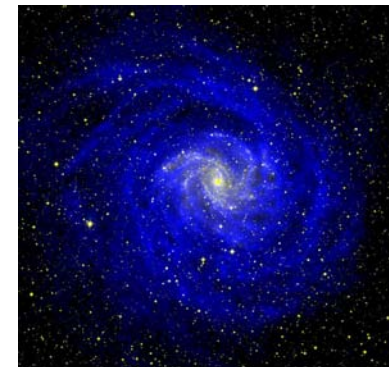


- Gravitational wave studies via a pulsar timing array

Precise measurements of arrival times of pulsars can track changes in pulsar-Earth distance



- Interferometric imaging
  - H I observations or continuum
  - Galaxy evolution through cosmic time, cosmology, magnetic field studies, ...



Exploring the Universe with the world's largest radio telescope



# Case 3a: Pulsar Timing



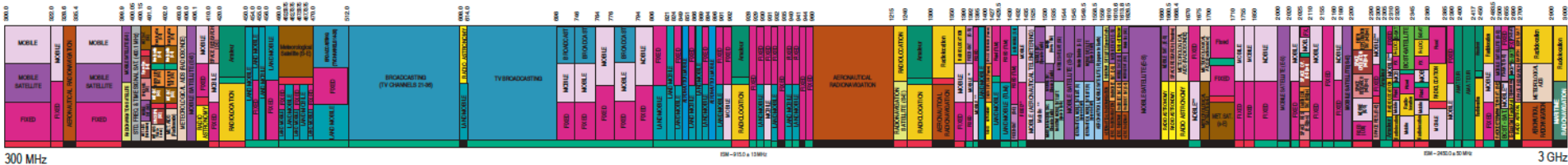
- Nyquist sampling of a 1 GHz bandwidth = 2 Gsamples/s
- 4 polarizations
- 1-byte representation of data?
- 15 minutes per pulsar
- Weekly observations of 100 pulsars (for 5 to 10 years)
- **= 3 PB/month**

# Case 3b: Imaging



- Nyquist sampling of a 1 GHz bandwidth = 2 Gsamples/s
- 4 polarizations
- 1000-element array (499,500 visibilities)
- 1-byte representation of data?
- **= 4 PB/s**

# Case 4: Interference Mitigation



- Most of the radio frequency (RF) spectrum is not reserved for use by radio astronomy  
In fact, **very little** is! ☹
- Other passive users are fine
- Active users can be troublesome!  
GPS, microwave ovens, cell phones, car ignitions, electric fences, ...

# Interference Mitigation II



- Radio flux density measured in Janskys  
 $1 \text{ Jy} = 10^{-26} \text{ W/m}^2/\text{Hz}$
- $10 \text{ } \mu\text{Jy}$  = EVLA, GBT, ASKAP, MeerKAT, ...  
sensitivity
- $10 \text{ nJy}$  = SKA aim
- $100 \text{ Jy}$  ~ cell phone on Moon

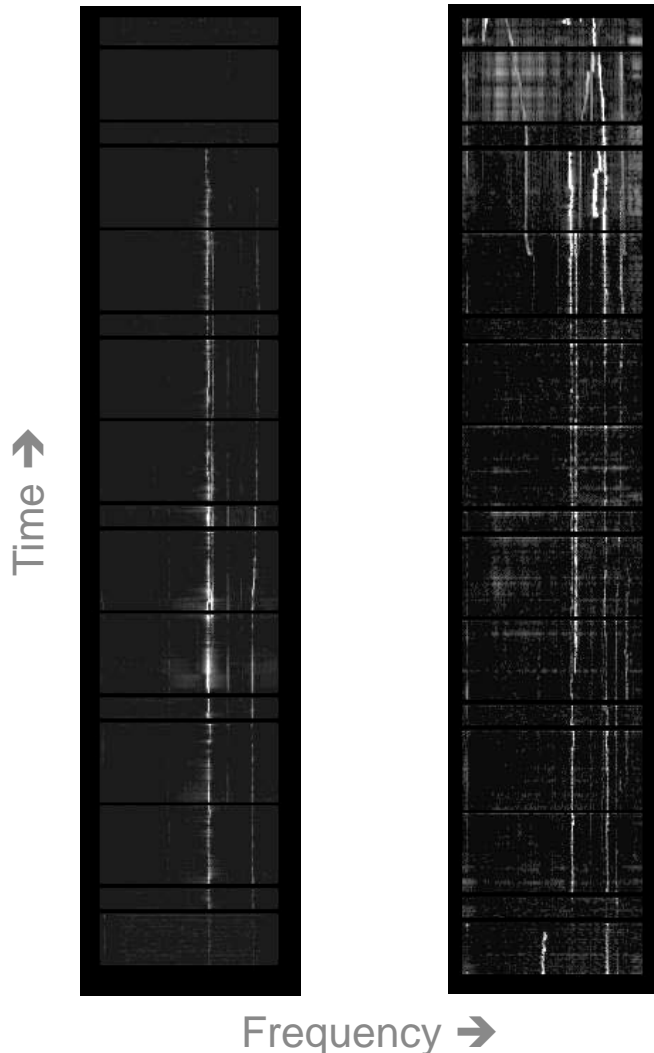


# Data Visualization



Vis. #1

Vis. #2



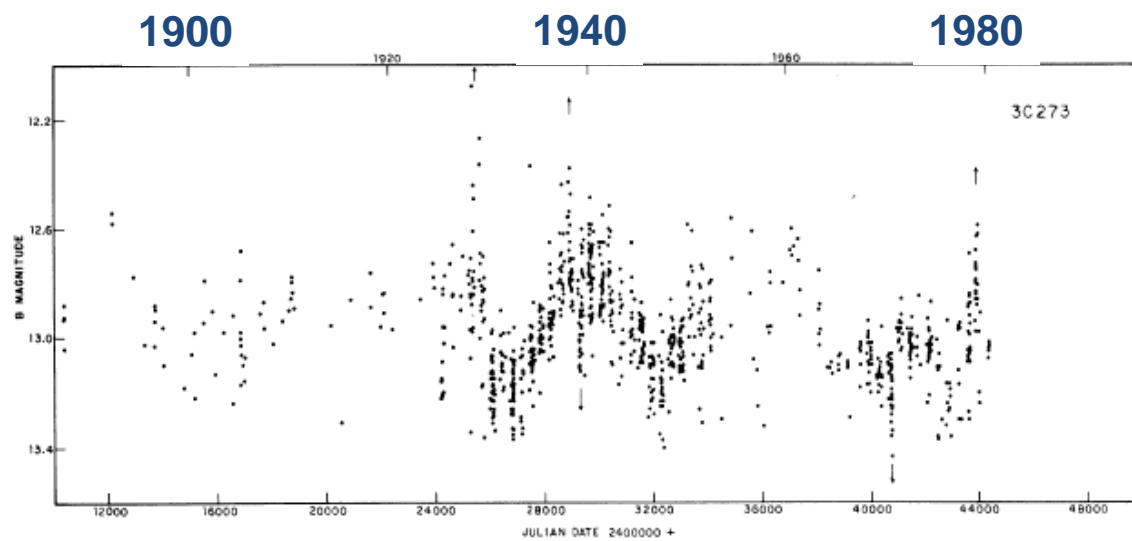
- Data acquired from an array are at least 4-D
  - Visibility (=  $\text{antenna}_i \times \text{antenna}_j$ )
  - Frequency
  - Time
  - Polarization
  - (Beams? for a multi-beam system)
- Best results still obtained by hand

# Case 5: Astronomy for the Future



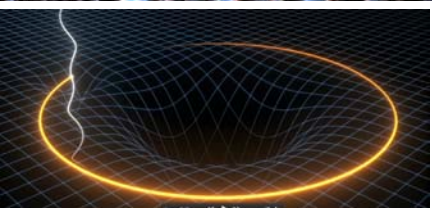
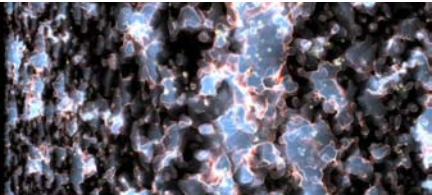
- 100+ year light curves
- BL Lac
  - Variable star #90 in constellation Lacertae?
  - Active galactic nucleus?

- What will our “academic grand-children” want to know?
- What will they be able to know?

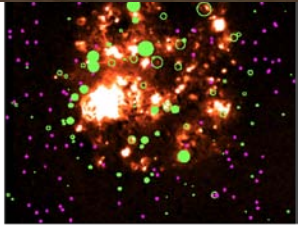
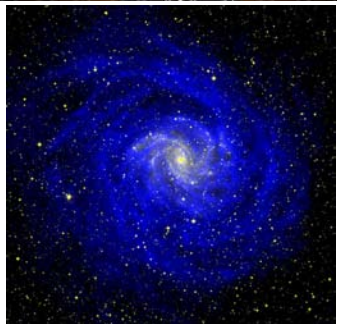


3C 273 (Angione & Smith 1985)

# Summary – Data-Intensive Astronomy



- Exciting science!
- Leads to exciting data challenges
  - Data volume (Exabyte)
  - Processing requirements (Exo-flop)
  - Data rates (PB/s)
  - Data visualization (HMI)
  - Long time scales (100 years?)
- Answers on Thursday



# Pulsar Surveys III



- Processing load