Title: Data Science in Radio Cosmology

Date: Mar 24, 2016  11:00 AM

URL: http://pirsa.org/16030130

Abstract: <p>In recent decades probing for the subtle indications of new physics in<br>
experimental data has become increasingly difficult. The datasets have gotten<br>
much bigger, the experiments more complex, and the signals ever smaller. Success<br>
stories, like LIGO and Kepler, require a sophisticated combination of statistics<br>
and computation, coupled with an appreciation of both the experimental realities<br>
and the theoretical framework governing the data.<br>
<br>
In this talk I will look broadly at data science in physics, and how and why it<br>
has taken an increasingly central role. I'll highlight specifically my current<br>
area of research, radio cosmology: discussing why it is one of the most<br>
challenging areas for data science, and describing my work developing optimal<br>
and efficient statistical methods for turning terabytes of timestreams into<br>
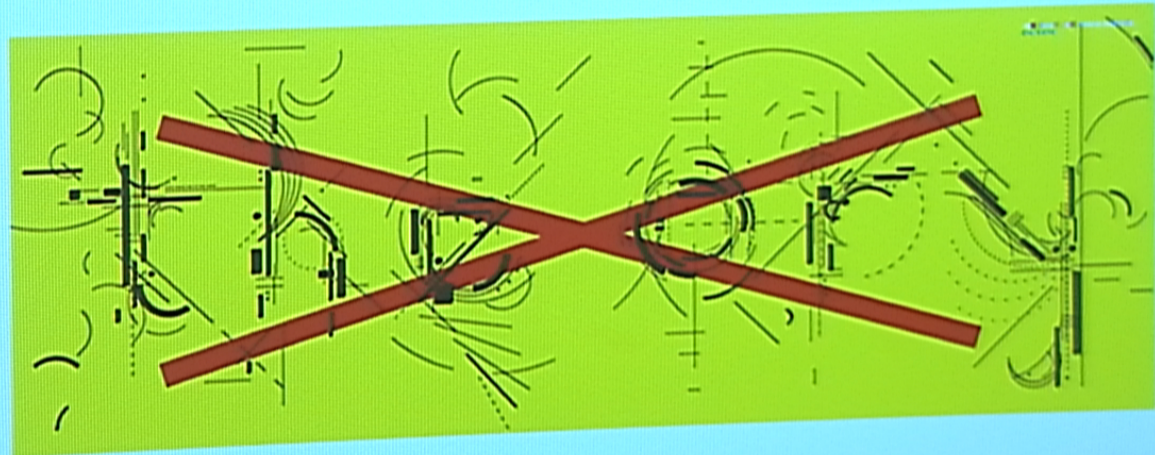cosmology.</p>

# Data Science

- Combination of:
  - ▶ Statistics
  - ▶ Signal processing
  - ▶ Machine learning
  - ▶ High performance computing
  - ▶ **Physics** (theory and experiment)

CHRIS ANDERSON   MAGAZINE   06.23.08   12:00 PM

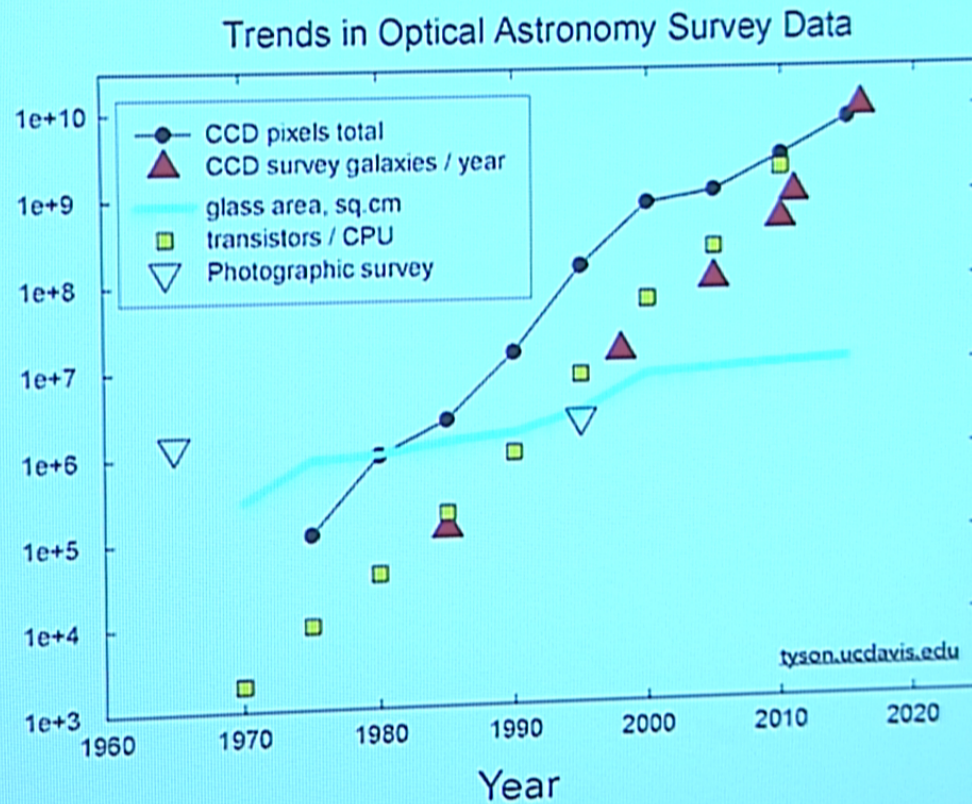# THE END OF THEORY: THE DATA DELUGE MAKES THE SCIENTIFIC METHOD OBSOLETE

# Data Science *in* Science

- For Science
  - ▸ Model fitting
  - ▸ Model selection

- Tools:
  - ▸ *The fashionable stuff*
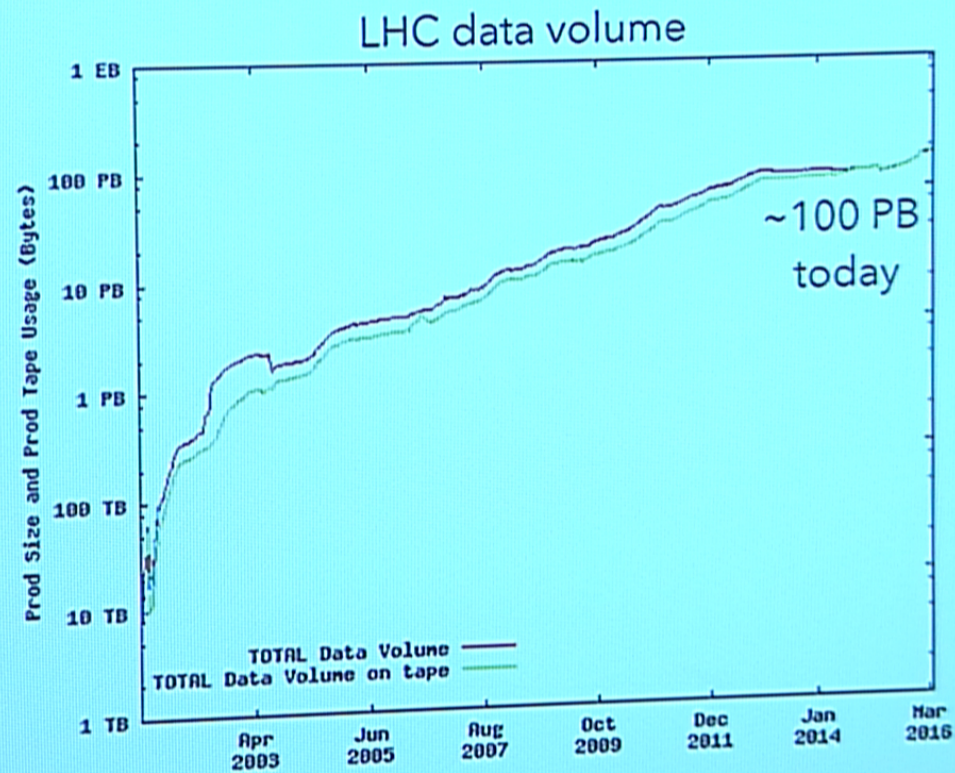  - ▸ Machine learning (classification, regression…)

Why now? Data acquisition...

Trends in Optical Astronomy Survey Data

# Data growth…



LHC data volume

~100 PB today

Exponential growth in data size

# Why now?

- Huge increase in accumulation of data
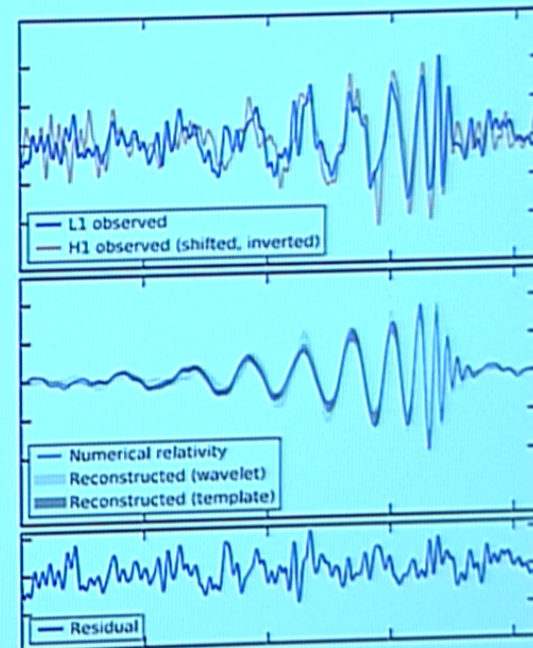- Volume and complexity means it is becoming a specialised just to do anything with it.

# LHC

- Need to identify which collisions are interesting
- Each collision is described by many tens of parameters
- Machine classification problem (Boosted Decision Trees)
  - ▶ Trained on simulated data
  - ▶ Selects between events of interest and background events

- Used in real-time in software triggers (LHCb)

- Used for event selection for Higgs detection (CMS)

Gligorov 2014, CMS Collaboration 2012

# LIGO event detection

- Real-time matched filter

  ‣ Effectively search against 250k template waveforms

  ‣ Look for peaks in likelihood ratio, keep those that exceed some false positive rate

  ‣ Keep only events coincident within 15ms in both detectors (~light travel time)

L1 observed
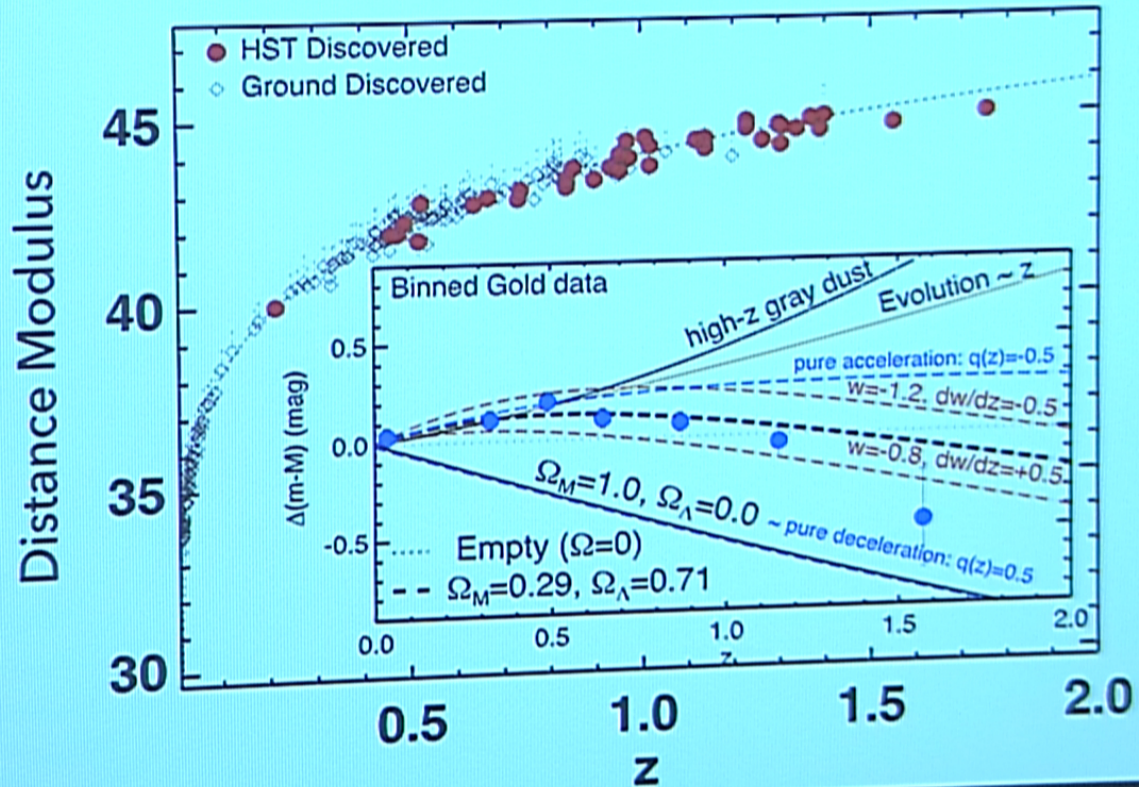H1 observed (shifted, inverted)

Numerical relativity
Reconstructed (wavelet)
Reconstructed (template)

Residual

Abbott et al. 2016

# Data lifecycle

- During acquisition
  - ▶ Real-time phase

- Post-acquisition
  - ▶ Basic results
  - ▶ Processing to likelihood function

- Science
  - ▶ Parameter estimation
  - ▶ Synthesis with other data

# Accelerating Universe

Luminosity distance
$$D_L = (L/4\pi F)^{1/2}$$

# Probing Dark Energy

- Acceleration 'explained' by dark energy (cosmological constant, quintessence ...)

- Expansion is governed by Friedmann equation

$$H(z)^2 \approx \Omega_m (1+z)^3 + \Omega_{DE} \exp\left[\int_0^z (1+w(z))\frac{dz}{1+z}\right]$$

- Fundamental physics is contained within the equation of state $w(\rho) = p/\rho < -1/3$

# Probing Dark Energy

- Construct likelihood function for dataset i.e. Pr(supernovae|expansion)

- Use Markov-Chain Monte-Carlo to infer distribution of relevant parameters

- Propose parameters ($w_0$, $w_a$, $\Omega_m$, …)

- Generate expansion history

- Compare to data (likelihood), accept parameters if likely

- Repeat until enough samples



Supernova Cosmology Project
Suzuki, et al., ApJ (2011)

Union2.1 SN Ia
Compilation
with SN
Systematics

BAO

CMB

SNe

# Baryon Acoustic Oscillations

- Sounds waves propagating in the early Universe. Leave acoustic peaks in the CMB

- Weaker imprint left in the matter distribution

- Gives a standard (statistical) ruler

CMB angular power spectrum



$$r_s = \int_0^{\tau_*} c_s d\tau \sim 100\, h^{-1}\, \mathrm{Mpc}$$

Known from CMB

# Galaxy redshift surveys



Galaxy Correlation Function

$$r_s = \Delta\theta \, d_A(z) \qquad r_s = \frac{c\Delta z}{H(z)}$$

# Cosmological 21cm

Higher energy state  Spin flip

1420 MHz
$\lambda = 21$ cm

- 21cm line is the transition between parallel and anti-parallel spins of neutral Hydrogen

- The ratio between the two occupancies determines the spin temperature $T_S$

$$n_1/n_0 = (g_1/g_0)\exp(-T_*/T_S)$$

- We can observe the contrast relative to the CMB

$$\Delta T = 23.8 \left(\frac{1+z}{10}\right)^{1/2} [1 - \bar{x}(1+\delta_x)](1+\delta_b)(1-\delta_v)\left[\frac{T_S - T_\gamma}{T_S}\right] \text{mK}$$

# Hydrogen in the Universe

z = 1100

z = 20

z = 6

Dark ages

Reionisation

HI in galaxies

Djorgovski et al. (Caltech)

# 21cm Intensity Mapping

- In 21cm the frequency gives the redshift.

- Observe the diffuse emission from many unresolved galaxies

- Changes the game in telescope design:

  - Previously: large field of view, large collecting area, large angular resolution (SKA?)

  - Now: large field of view, large collecting area, modest angular resolution (compact arrays, single dishes).

Chang, Pen, Peterson and McDonald , 2008, http://arxiv.org/pdf/0709.3672

# Foreground Challenges



400 MHz

-0.00045    0.00043

Cosmological 21cm Signal ~ 1mK

# A way out?

# Cross correlation detection

- Cross-correlation with of GBT data with DEEP2 Galaxy survey by Chang et al.(2010) - *avoids foreground problem!*

- Updated using WiggleZ survey (Masui et al. 2012)

**Cross power spectrum**

$$\Omega_{HI} = \left[ 0.62^{+0.25}_{-0.15} \right] \times 10^{-3}$$

- 15 hr
- 1 hr
- $\Omega_{HI}\, b_{HI}\, r = 0.43\ 10^{-3}$

**7.4 sigma detection**

$\Delta^2(k)$ [mK]

k [h / Mpc]

- cross-power
- deep auto-power
- wide auto-power
- combined

Intensity Mapping at Green Bank

# The Future?

- Work at GBT will continue with the aim of measuring the 21cm *autocorrelation*.

- However, observations like this are slow. To survey the whole sky to this depth ~ 20 years

  ▸ Is there a better way to do this?

# Interferometers

$$\Delta\phi = 2\pi\hat{n} \cdot d_{ij}/\lambda$$

- Visibility is instantaneous correlation of 2 antennas

$$V_{ij} = \langle F_i F_j^* \rangle$$

advancing wave crests

baseline

$F_i$ $\qquad$ $F_j$

- Each pair measures a Fourier mode of the sky

- Written explicitly:

$$V_{ij}(t) = \frac{1}{\Omega_{ij}} \int d^2\hat{n}\, A_i(\hat{n};t) A_j^*(\hat{n};t) e^{2\pi i \hat{n} \cdot \mathbf{u}_{ij}(t)} T(\hat{n})$$

# Data rate

- For full $N^2$ correlation
  - ▶ ~5 GB/s
  - ▶ ~400 TB/day
  - ▶ ~140 PB/year

- Need a way to significantly compress the data!

# Highly redundant array

Not so fast! Calibration

- Each feed has an unknown, time-variable, complex gain, from amplifier and cable behaviour
- Must correct for this, or the baselines will average incoherently
- This process is *calibration*, and must be done in **real-time**
  ▶ Nearly optimal solution via eigenvalue decomposition
  ▶ Use injected calibration signal
  ▶ Sky signal pulsars

Newburgh + CHIME, arXiv:1406.2267

# Interferometers

$$V_{ij}(t) = \frac{1}{\Omega_{ij}} \int d^2\hat{\mathbf{n}}\, A_i(\hat{\mathbf{n}}; t) A_j^*(\hat{\mathbf{n}}; t) e^{2\pi i \hat{\mathbf{n}} \cdot \mathbf{u}_{ij}(t)} T(\hat{\mathbf{n}})$$

- Write in terms of a beam transfer function:

$$V_{ij}(t) = \int d^2\hat{\mathbf{n}}\, B_{ij}(\hat{\mathbf{n}}; t) T(\hat{\mathbf{n}}) + n_{ij}(t)$$

# Transit Interferometers

- Timeseries is periodic on the sidereal day $\quad t \to \phi$

  ▸ Apply this restriction and see how the analysis goes.

$$V_{ij}(\phi) = \int d^2\hat{n}\, B_{ij}(\hat{n}; \phi) T(\hat{n}) + n_{ij}(\phi)$$

**Spherical Harmonic Transform**

$$V^{ij}(\phi) = \sum_{lm} B^{ij}_{lm}(\phi) a^{T}_{lm} + n^{ij}(\phi)$$

**Fourier Transform**

$$V^{ij}_{m} = \sum_{l} B^{ij}_{lm} a^{T}_{lm} + n^{ij}_{m}$$

# Transit Interferometers

- Timeseries is periodic on the sidereal day $\quad t \to \phi$

  ▸ Apply this restriction and see how the analysis goes.

$$V_{ij}(\phi) = \int d^2\hat{\mathbf{n}}\, B_{ij}(\hat{\mathbf{n}}; \phi)T(\hat{\mathbf{n}}) + n_{ij}(\phi)$$

**Spherical Harmonic Transform**

$$V^{ij}(\phi) = \sum_{lm} B^{ij}_{lm}(\phi)a^T_{lm} + n^{ij}(\phi)$$

**Fourier Transform**

$$V^{ij}_m = \sum_{l} B^{ij}_{lm} a^T_{lm} + n^{ij}_m$$

# m-mode transform

- Mapping does not mix m's (each is independent)

$$V_m^\alpha = \sum_l B_{lm}^\alpha a_{lm}^T + n_m^\alpha$$

- Write in vector form

$$v = Ba + n.$$

- Simple, linear mapping from the information on the sky, to the measured degrees of freedom

- Discrete relation, with finite number of degrees, can apply all the standard statistical, signal processing techniques.

- Computationally efficient: For 1000 m's an $O(N^3)$ matrix operation becomes $10^6$ times faster

# Interferometric Imaging

- Traditional imaging is based around the 2D Fourier Transform approximation to the interferometry equation (only valid on small patches instantaneously)

- Use a series of steps to relax this approximation and increase field of view (w-projection, mosaicking, A-projection)

  ▶ eg. w-term. From non coplanarity of array and sky. Solve by iteratively deconvolving the effects

$$V = \int dx dy A^2(x,y) e^{2\pi i (ux + vy + w\sqrt{1-x^2-y^2})} I(x,y)$$

# m-mode Imaging

- For our restricted domain (transit telescopes), we can solve the problem exactly.

- Measurement is linear mapping:

$$v = B a + n .$$

- How do we make an image of the sky? Use standard tools of signal processing:

  ▶ Pseudo-inverse to solve and regularize (*Maximum likelihood*)

  ▶ Wiener Filter (*Bayesian expectation*)

- Conceptually straightforward. Deals naturally with all full sky effects, polarisation etc.

# Observed sky *(from time stream)*   $\hat{a} = \left(N^{-\frac{1}{2}}B\right)^{+}N^{-\frac{1}{2}}v$

2x15m wide cylinders, 60 feeds, 0.25m spacing 400-600 MHz

# Foreground Removal

- Spectral smoothness allows separation of 21cm
  - ▶ Measure components and model (Liu, Dillon etc.)
  - ▶ Power spectrum removal (Foreground wedge)
  - ▶ Delay-space filtering (Parsons et al. 2012)
- Most methods have difficulties:
  - ▶ *Mode mixing* of angular and frequency fluctuations by frequency-dependent beams (esp. interferometers)
  - ▶ *Robustness* Biasing introduced if foreground model poorly understood (esp. non-gaussianities)
  - ▶ *Statistical Optimality* Need to keep track of transformations on statistics, for optimal PS estimation
  - ▶ *Polarisation leakage* mixes fluctuations from polarised foreground

# Karhunen-Loeve Transform

- Old CMB idea - E/B mode separation (Bunn et al. 2003)

- An 'optimal' treatment - m-modes makes it feasible.

- Construct the covariances of the signal and foregrounds in the measured basis

$$S = \langle ss^\dagger \rangle = B \langle a_s^* a_s^T \rangle B^\dagger \qquad F = B \langle a_f a_f^\dagger \rangle B^\dagger$$

- Jointly diagonalise both (eigenvalue problem)

$$Sx = \lambda Fx$$

- Gives a new, uncorrelated basis. Corresponding eigenvalue gives the expected signal to foreground power ratio.

Most foreground

smooth in freq

Most signal

oscillates in freq

# Foreground Cleaning



Foregrounds $10^6$ times larger than signal
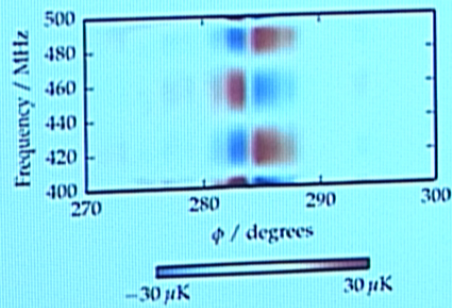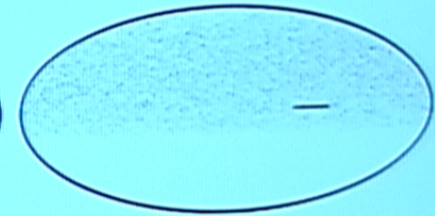
# Foreground Cleaning

Unpolarised Foreground    Polarised Foreground (Q)    21cm Signal

Foreground residuals significantly smaller than signal

# Foreground Cleaning
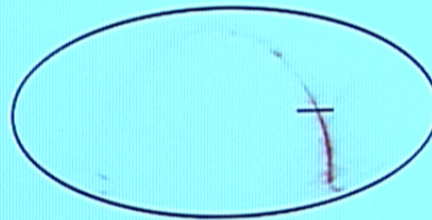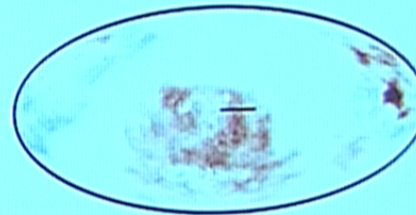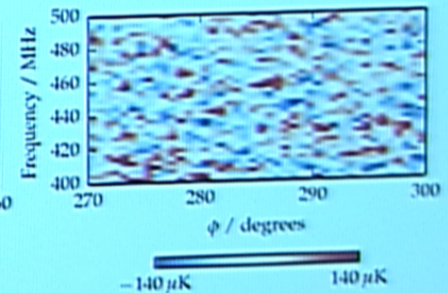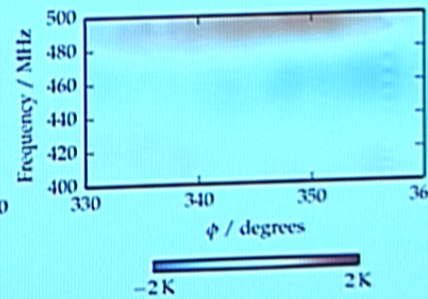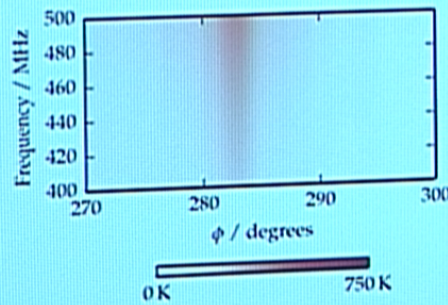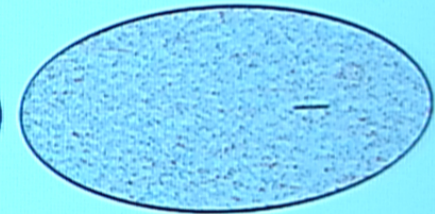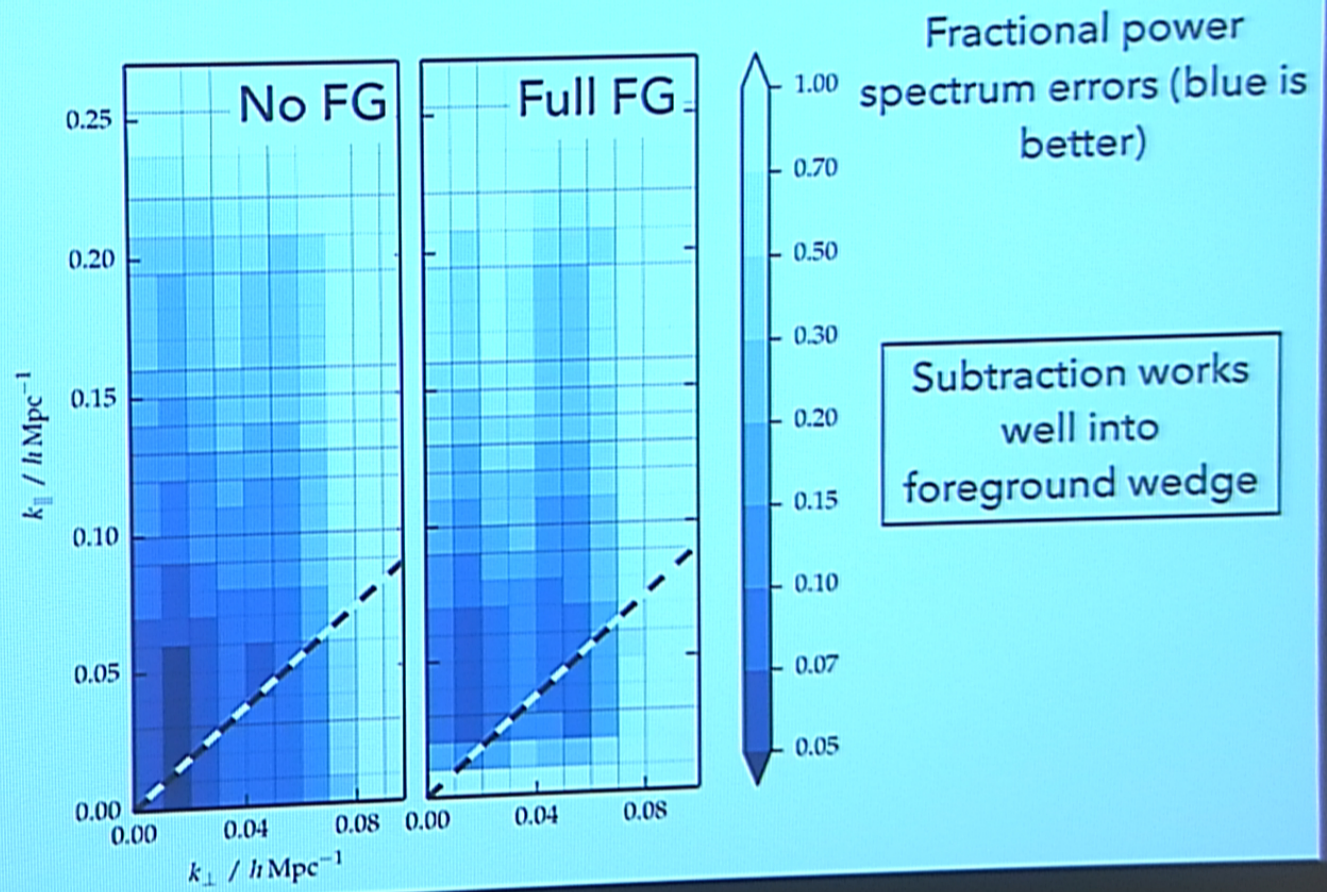


Unpolarised Foreground    Polarised Foreground (Q)    21cm Signal

Foregrounds $10^6$ times larger than signal

2D Power spectrum Estimation

# Summary

- Data Science is becoming an increasingly large and distinct part of physics

- 21cm Intensity Mapping is a promising technique for mapping the Universe and measuring BAOs.

- Data volume and foregrounds are challenging

- New techniques, like the m-mode formalism show promise for surmounting them