



# Source extraction and characterisation I – continuum

Dr. Paul Hancock

 [paul.hancock@curtin.edu.au](mailto:paul.hancock@curtin.edu.au)  
 @drpaulhancock



International  
Centre for  
Radio  
Astronomy  
Research

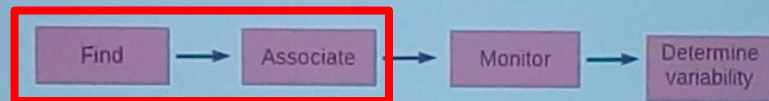
# Intro

Paul will  
talk about  
this later



## Determining variability

General transients pipeline:



Variability statistics:

$$\eta_{\nu} = \frac{1}{N-1} \sum_{i=1}^N \frac{(I_{\nu,i} - \xi_{I_{\nu}})^2}{\sigma_{\nu,i}^2}$$

Chi-square probability of constant flux

$$V_{\nu} = \frac{s}{\bar{I}_{\nu}} = \frac{1}{\bar{I}_{\nu}} \sqrt{\frac{N}{N-1} (\overline{I_{\nu}^2} - \bar{I}_{\nu}^2)}$$

Coefficient of variation (modulation)

Final steps:

... testing

# Everything you do wrong looks like variability

If you care about variability, then you care about all the ways that things can go wrong.

Eg the presence or changes in:

- Observing conditions
- RFI
- Calibration
- Imaging
- Detection of features
- Characterisation of features
- Analysis and methodology
- Work-flows

Variability can be:

1. Astrophysical
  - a. Intrinsic (SNe)
  - b. Extrinsic (scintillation)
2. Environmental (RFI, the ionosphere)
3. Instrumental (gain, bandpass stability)
4. Methodological (dodgy math!)

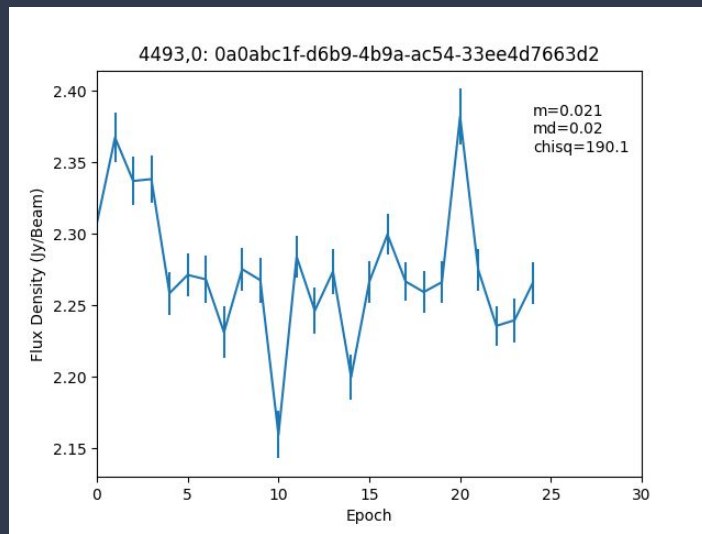
# Measuring Variability

## Problems:

- Masked/missing data points
- Upper/lower limits
- Non-uniform uncertainties
- Inaccurate uncertainties
- Separating significance and degree

## Solutions:

- Better source characterisation
- Better statistical models



# Source Finding Done Right

Assumptions:

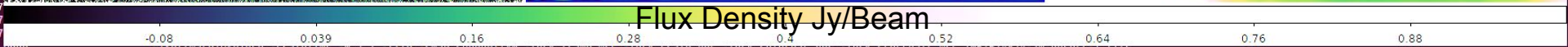
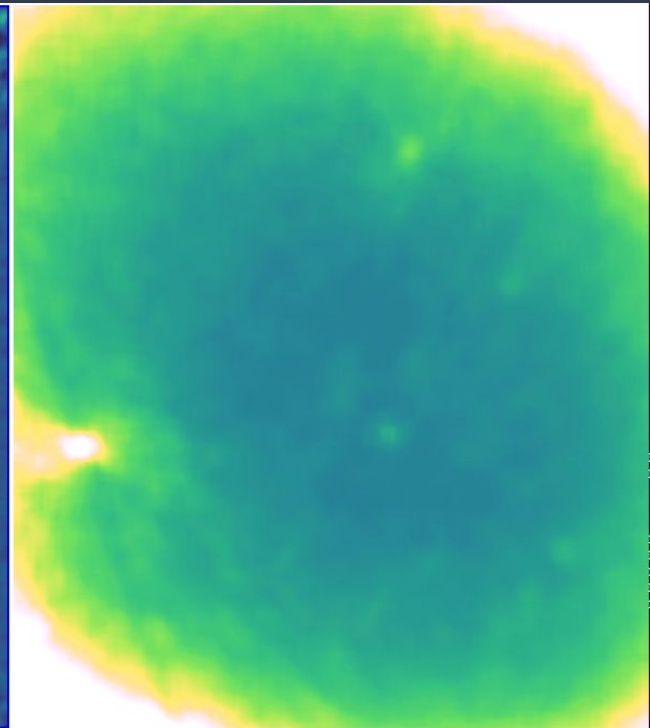
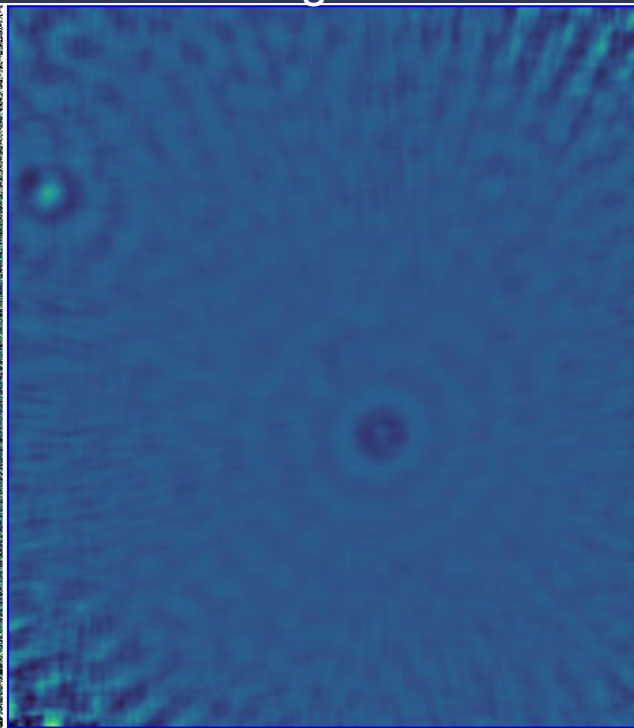
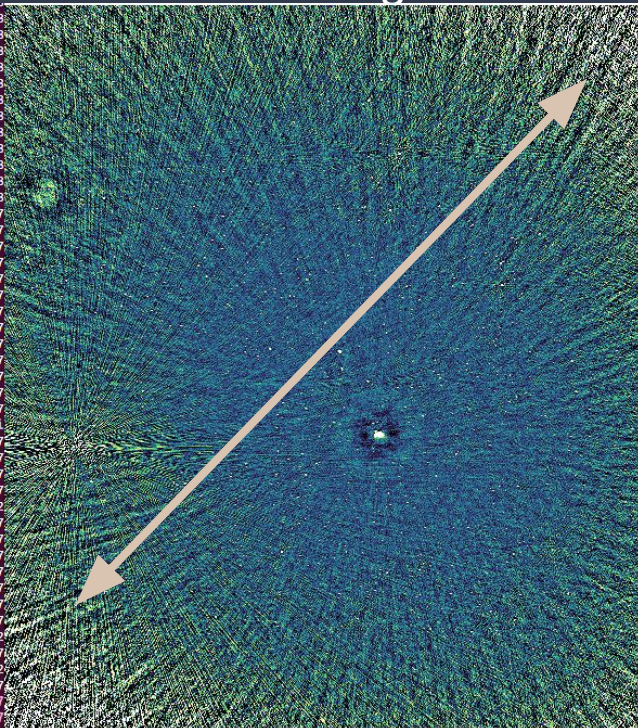
- Compact sources
- Continuum images

# Snapshot Image: Data model

Image

Background

Noise



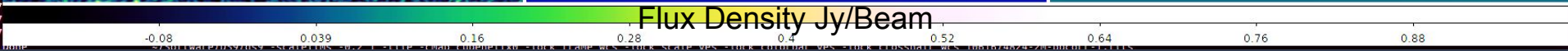
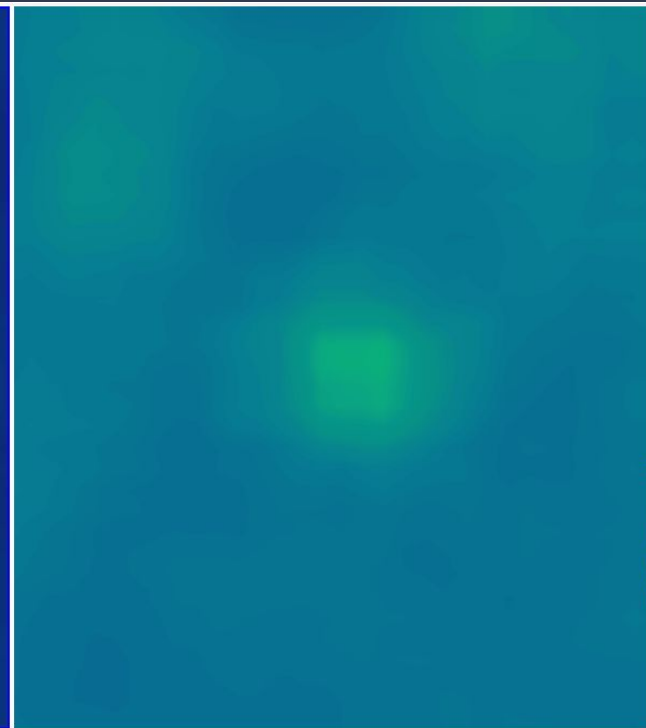
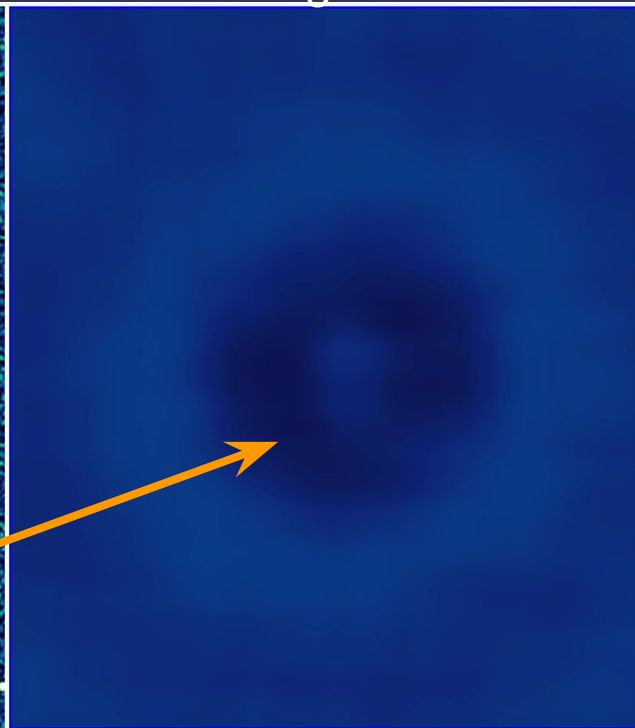
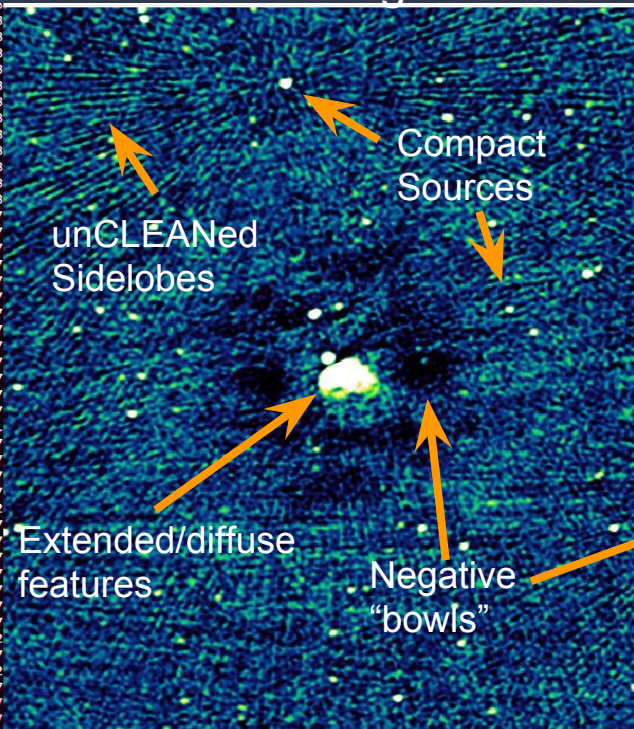


# Zoom

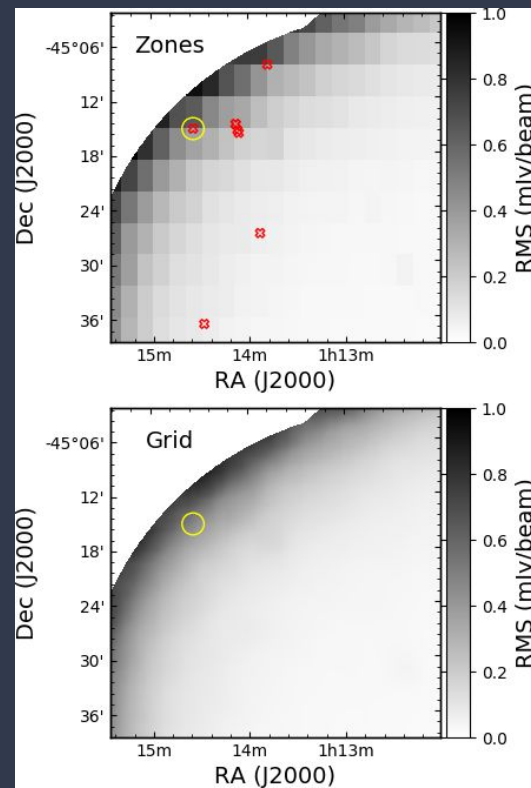
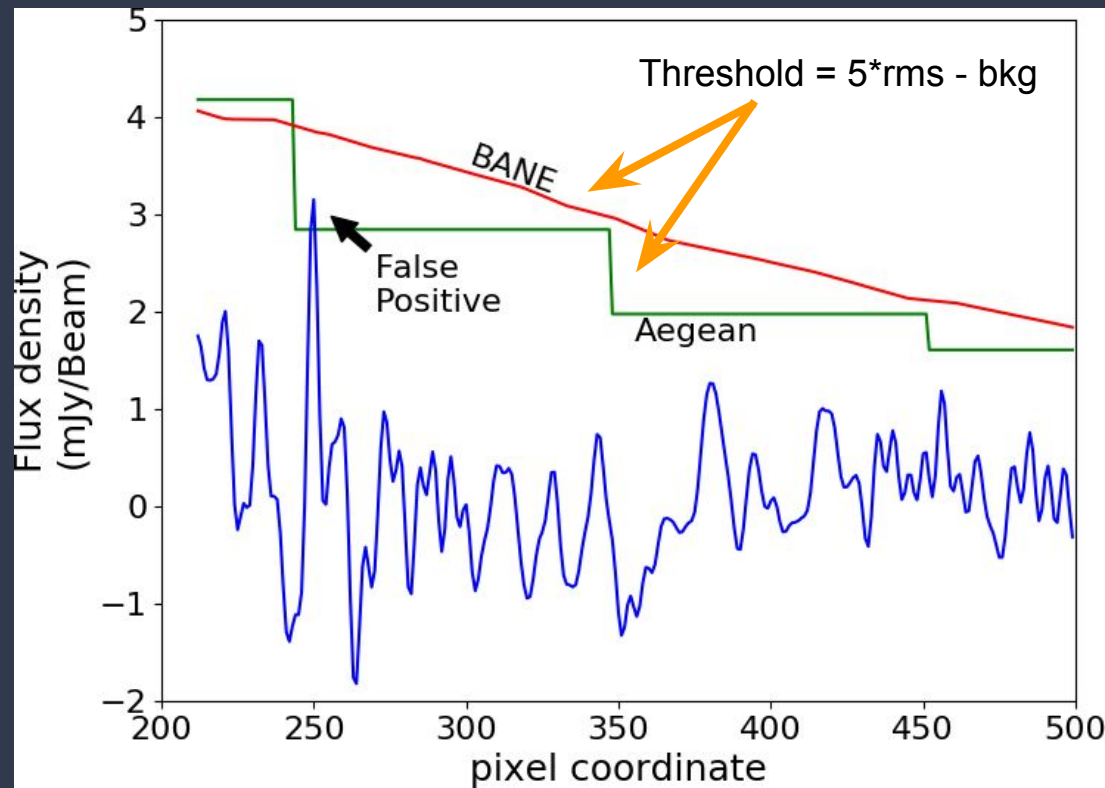
Image

Background

Noise



# Finding sources – thresholding

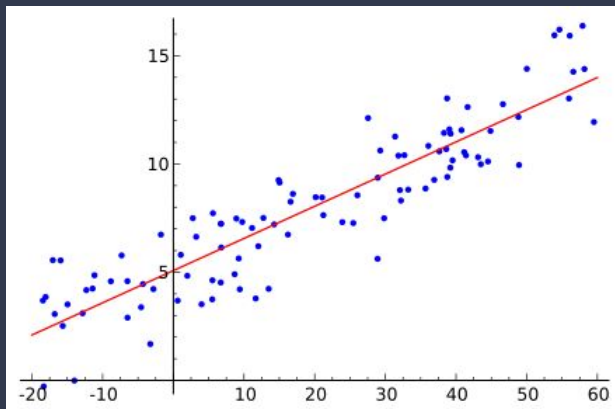




# (linear) Least squares fitting

Given:

- $x$  - data
- $f(\theta; x)$  - model data with parameters  $\theta$



commons.wikimedia.org

Goal:

- Minimise the sum of the square of the residuals

$$\arg \text{Min} \sum (f(\theta; x) - x)^2$$

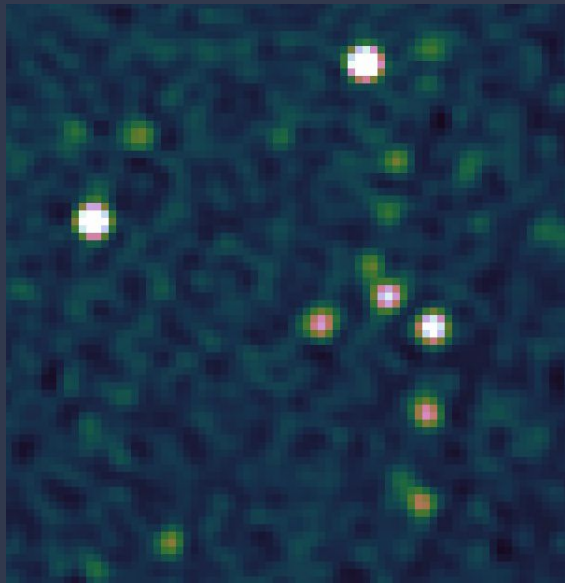
- a.k.a  $\chi^2$  minimisation

For linear models and data that is independent and identically distributed, least squares minimisation is unbiased, and has minimum variance.

# Radio Images

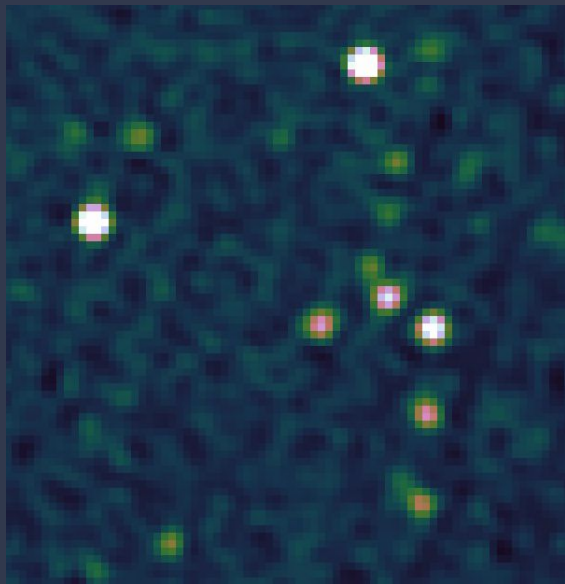
We fit with a source model that is Gaussian

$$f(x, y) = Ae^{-\left(\frac{(x-x_0)^2}{2\sigma_x^2} + \frac{(y-y_0)^2}{2\sigma_y^2}\right)}$$



# Radio Images

We fit with a source model that is Gaussian



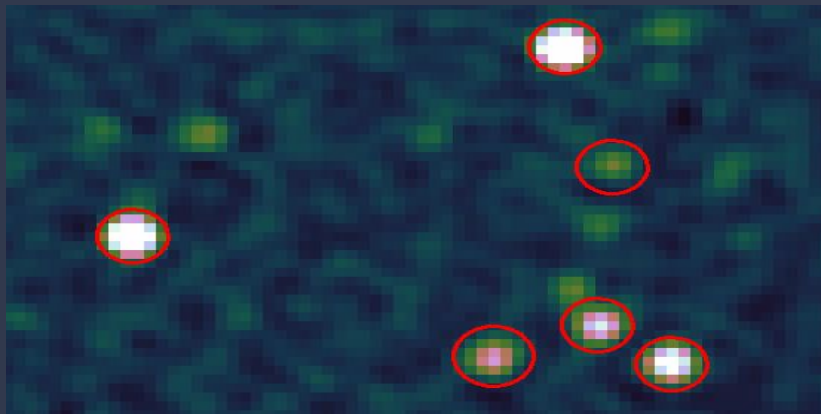
$$f(x, y) = A e^{-\left( \frac{(x-x_0)^2}{2\sigma_x^2} + \frac{(y-y_0)^2}{2\sigma_y^2} \right)}$$

Not linear, not even close

# (non linear) Least squares fitting

Given:

- $x$  - data
- $f(\theta; x)$  - model data with parameters  $\theta$



Goal:

- Minimise the sum of the square of the residuals

$$\arg \text{Min} \sum (f(\theta; x) - x)^2$$

- a.k.a  $\chi^2$  minimisation

For non linear models least squares minimisation gives a biased result.

All parameters are biased, even the 'linear ones' like amplitude

# Quantifying Bias

Refreiger & Brown 1998 (arXiv:9803279)

describe the expected **bias** as:

$$\langle a_i \rangle = \hat{a}_i - \frac{1}{2} \sigma_N^2 B_{lkj} D_{li} D_{kj} + O(\text{SNR}_s^{-3})$$

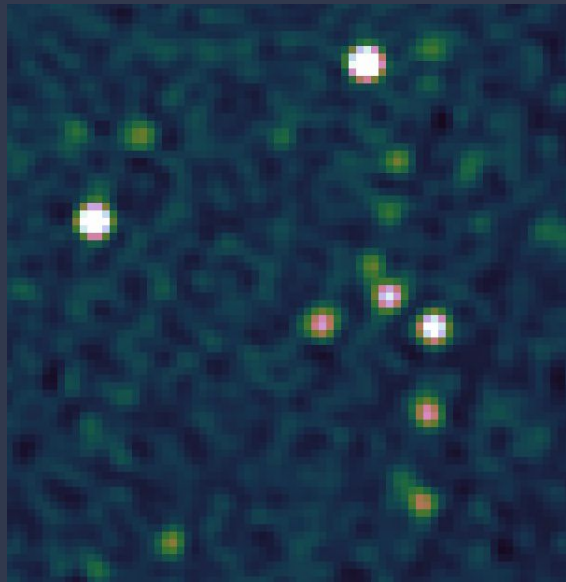
Where

$$\begin{aligned} D_{ij} &= (H^{-1})_{ij}, \\ H_{ij} &= \sum_p \frac{\partial F}{\partial a_i}(\mathbf{x}^p; \hat{\mathbf{a}}) \frac{\partial F}{\partial a_j}(\mathbf{x}^p; \hat{\mathbf{a}}), \\ B_{ijk} &= \sum_p \frac{\partial F}{\partial a_i}(\mathbf{x}^p; \hat{\mathbf{a}}) \frac{\partial^2 F}{\partial a_j \partial a_k}(\mathbf{x}^p; \hat{\mathbf{a}}), \end{aligned}$$

\* Math is for demonstration purposes only - Do not try this at home



# Radio Images Again



Data are correlated:

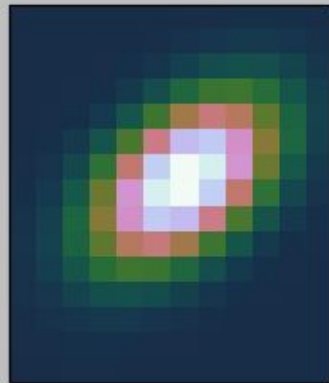
$\text{corr}(x,y) = \text{Dirty Beam} / \text{Point Spread Function}$

Even worse:

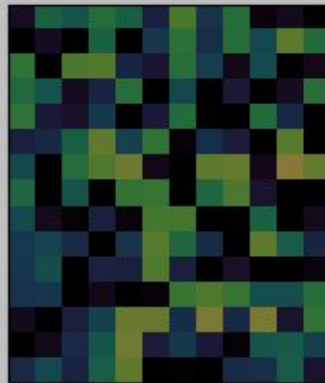
CLEAN-ing modifies the correlation function

# Our data

What our  
fitting  
algorithms  
**assume**  
we have



signal



noise



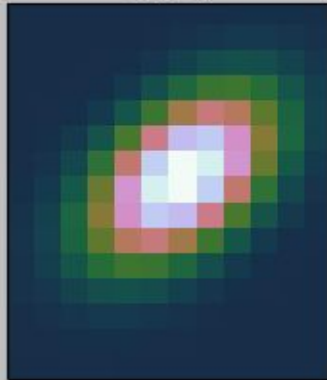
data

# Our data

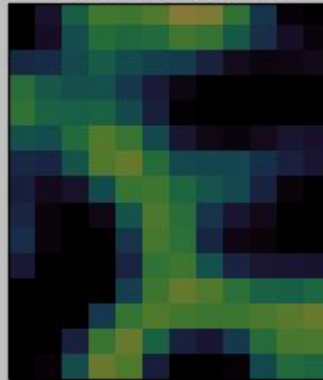
What we  
**actually**  
have

What our  
fitting  
algorithms  
**assume**  
we have

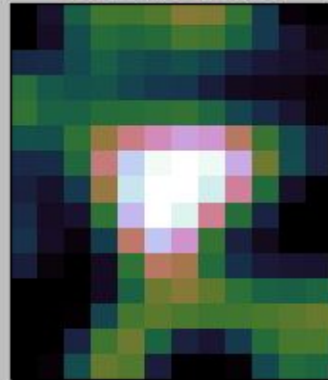
signal



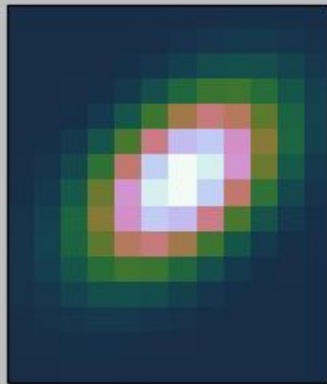
correlated noise



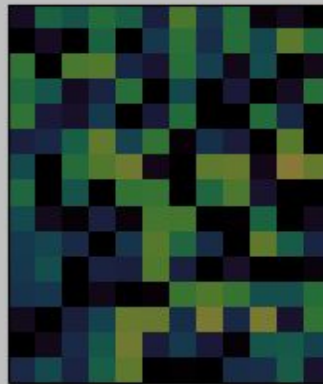
correlated data



signal



noise



data



# Correlated Data

Increases bias in all parameters

Additional bias towards local noise peaks at low SNR

Nearby sources now have correlated parameters

# Approaches



Ignore correlations completely



Fit as usual and then fiddle errors to account for correlations

- How many DoF do we “really” have?



Fit as usual and then replace errors with empirically derived quantities

- Condon 1997



Fit as usual and report errors based on analytical solution

- Refreiger & Brown 1998 (arXiv)



Fit with a cost function that incorporates a correction for correlated data

- Aegean 2.0, Hancock et al. 2018



# How do we do better?

Given:

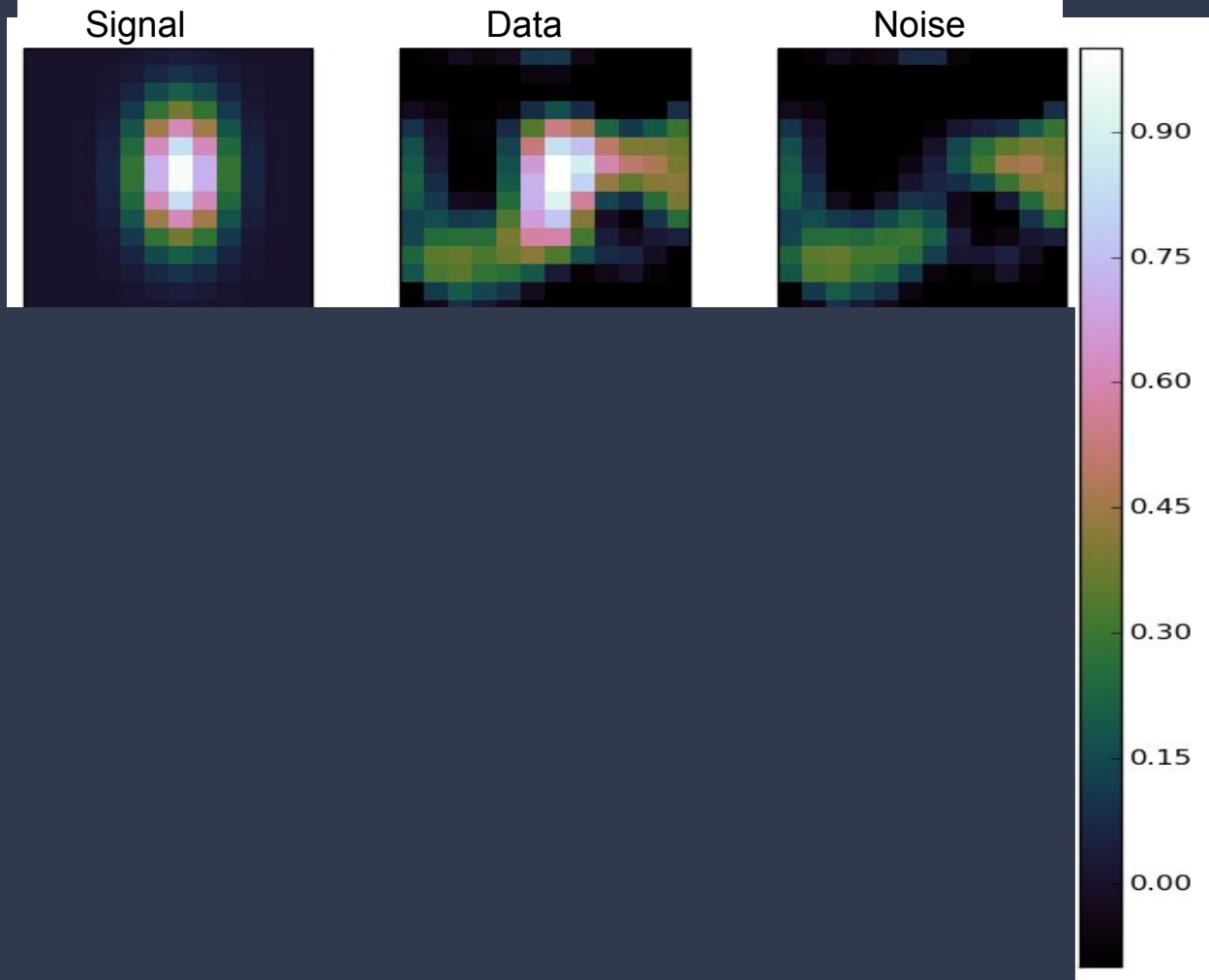
- $x$  - data
- $f(\theta; x)$  - model data with parameters  $\theta$
- Covariance matrix  $C$

Goal:

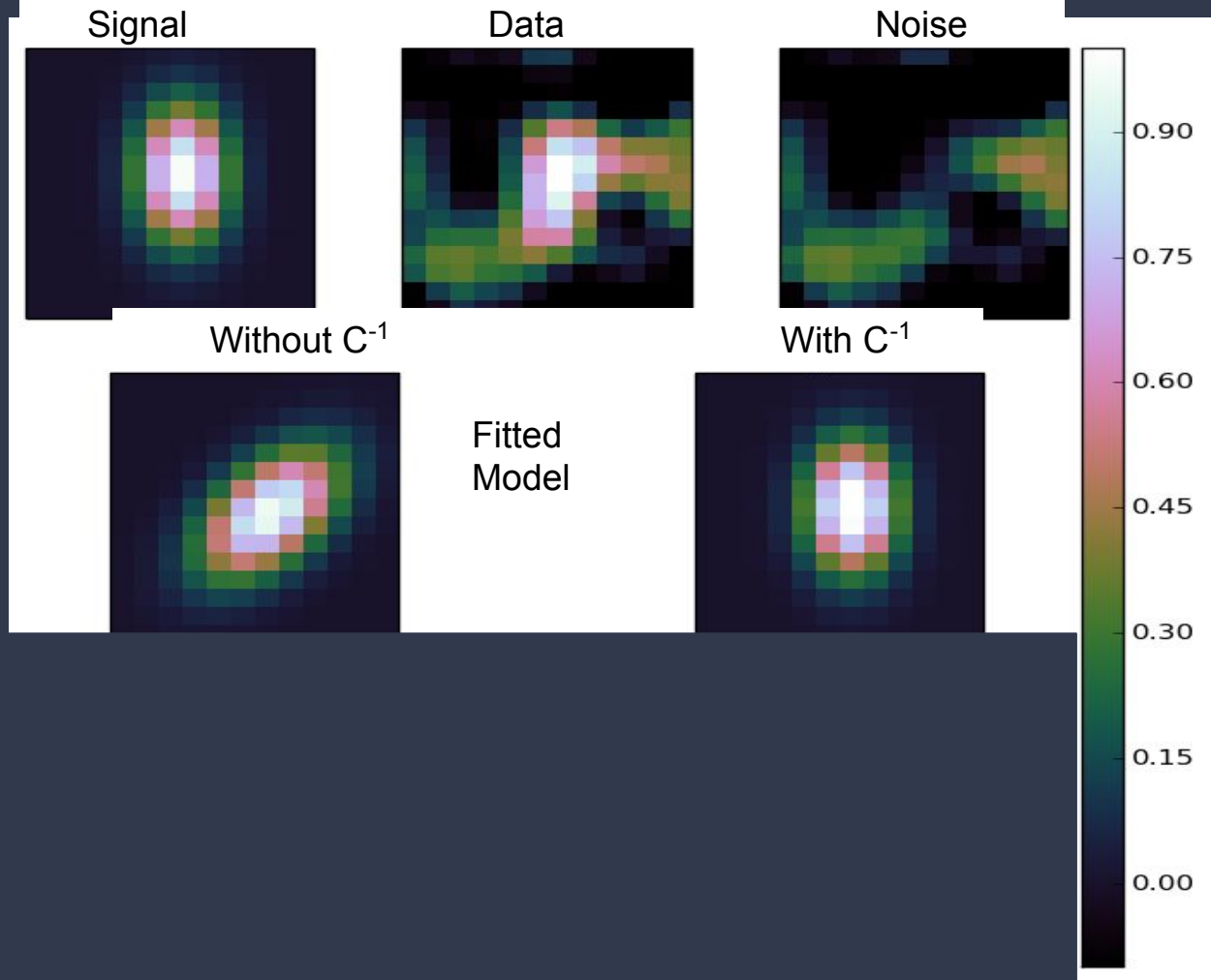
- Minimise the sum of the square of the residuals modified by the inverse covariance matrix

$$\text{Min } \{ (f(\theta; x) - x)^T C^{-1} (f(\theta; x) - x) \}$$

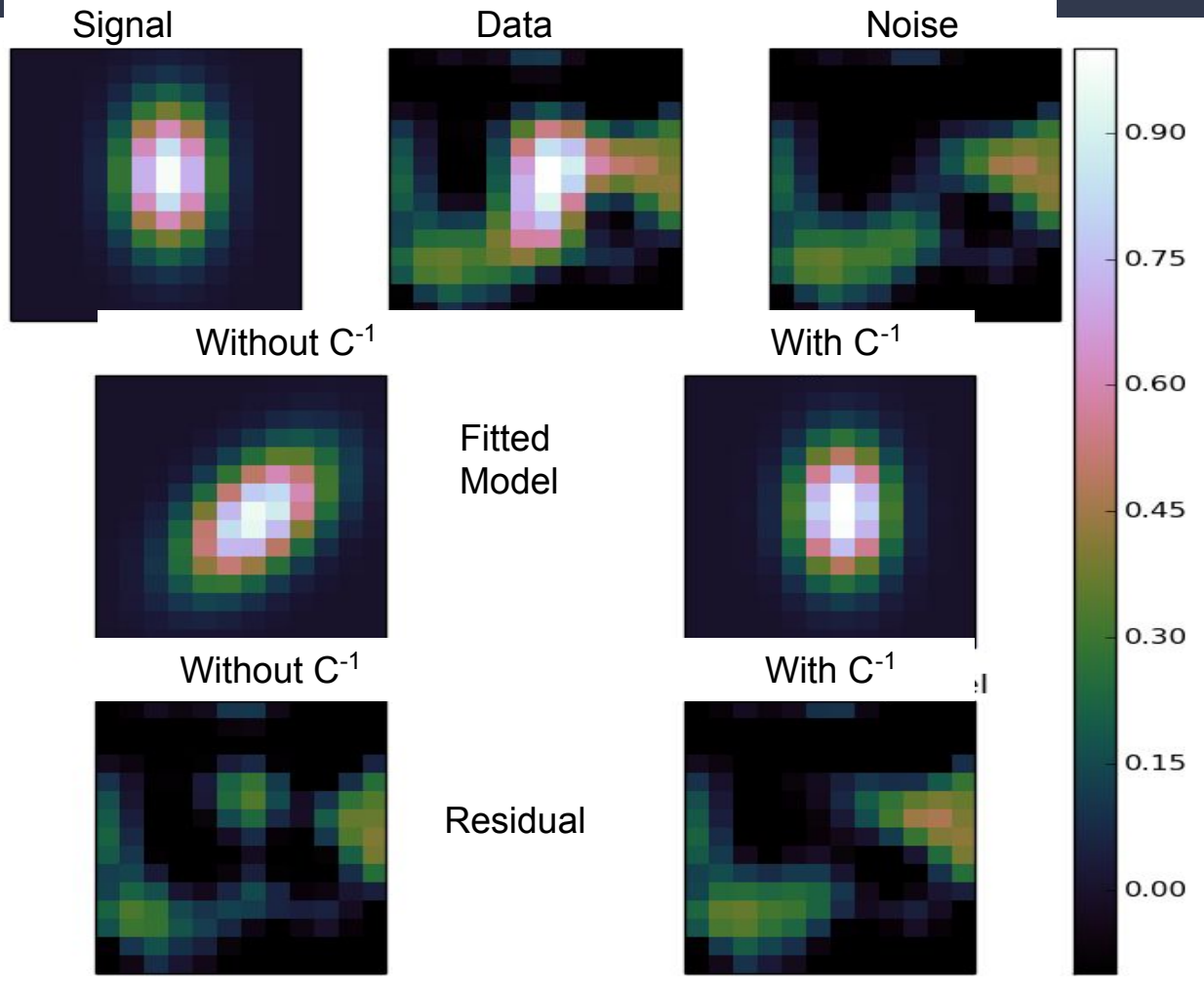
# Fitting with $C^{-1}$



# Fitting with $C^{-1}$

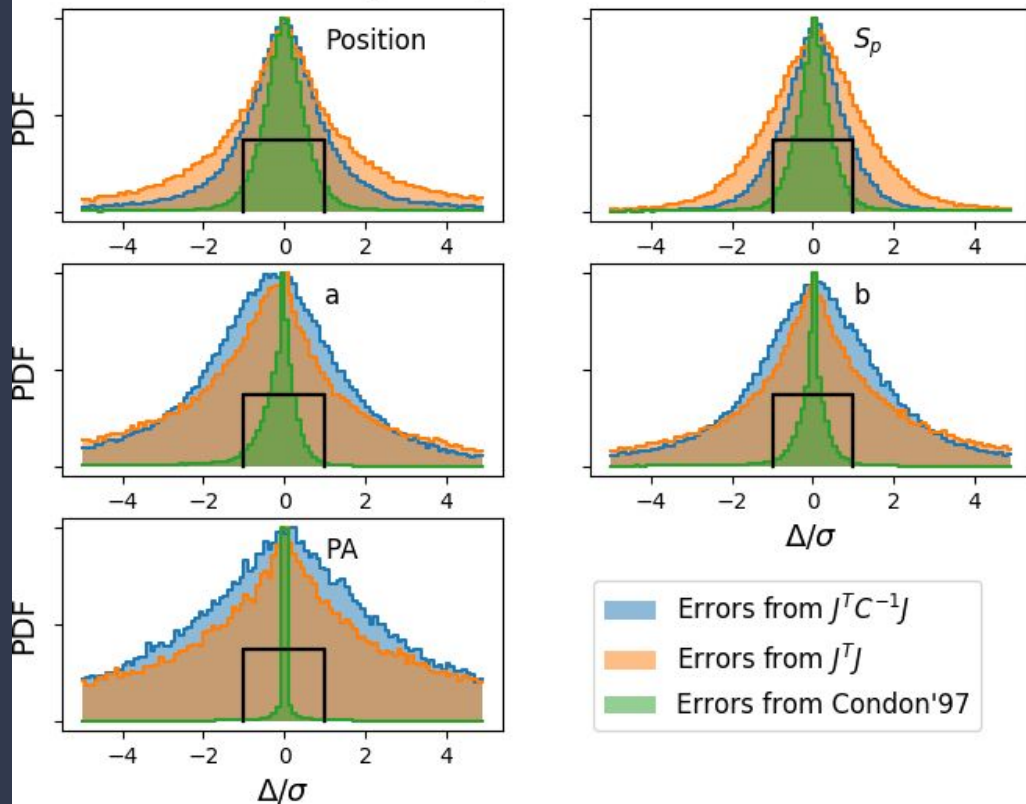


# Fitting with $C^{-1}$

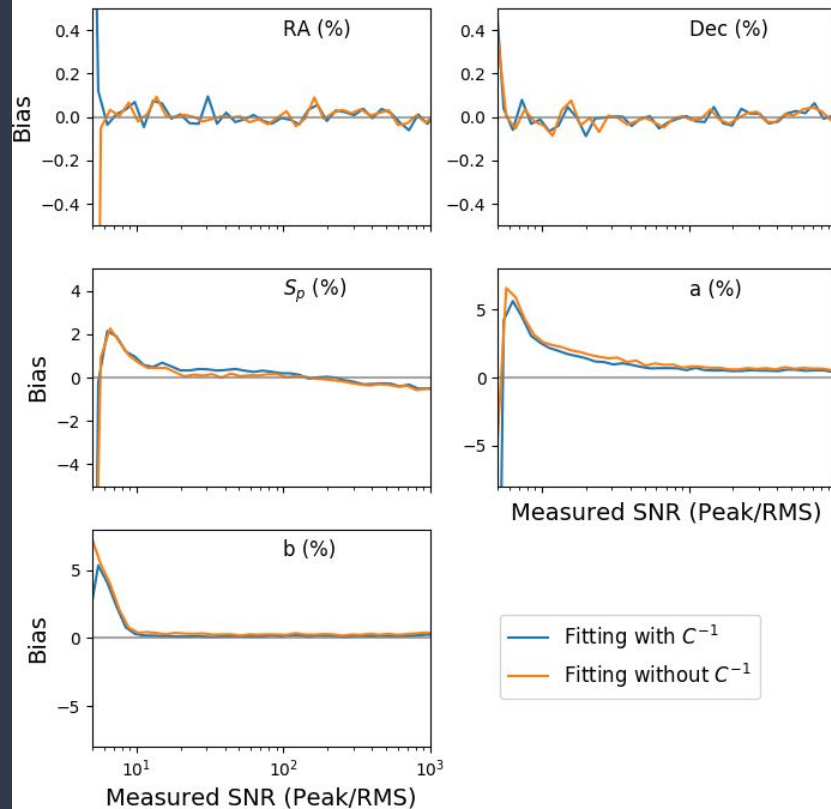


# Accurate errors, biases

Accuracy of reported uncertainties



Fractional bias in fit parameters

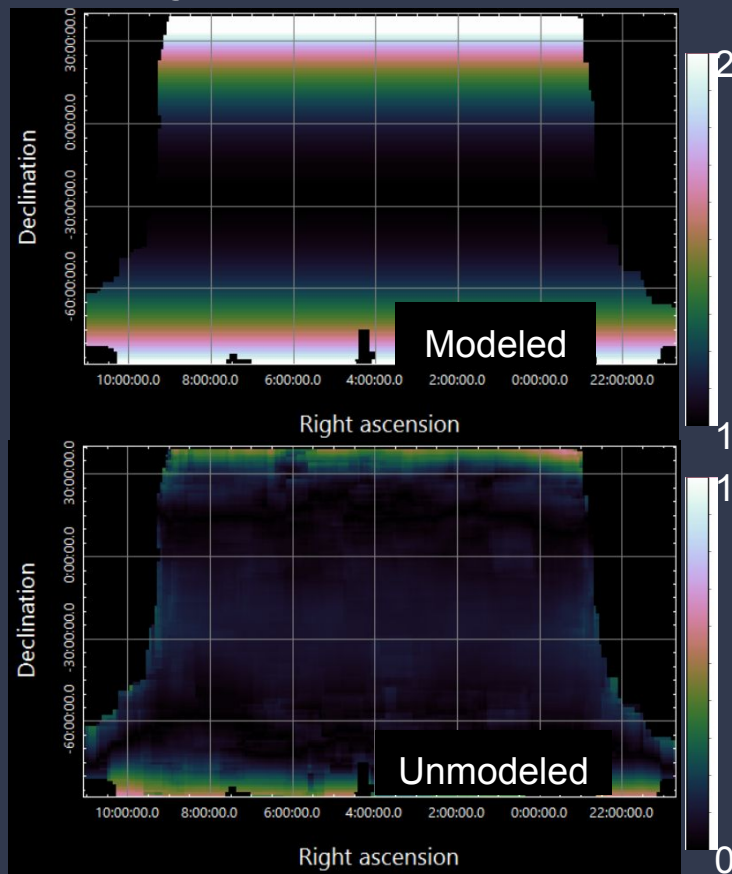
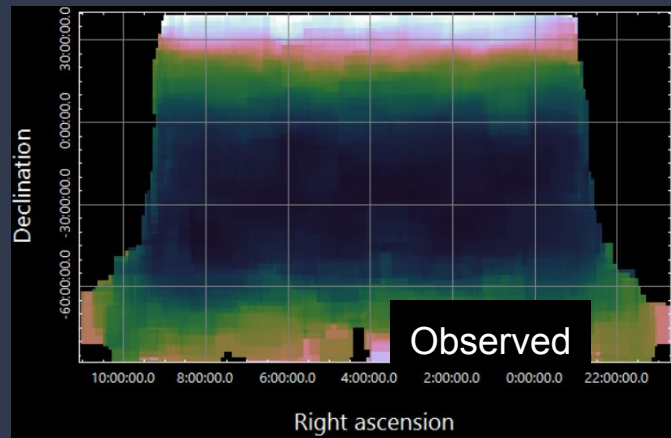




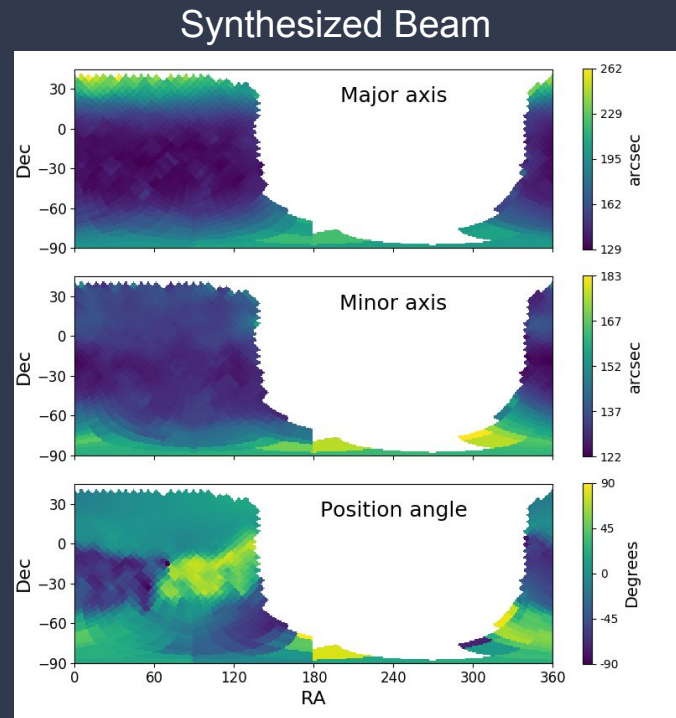
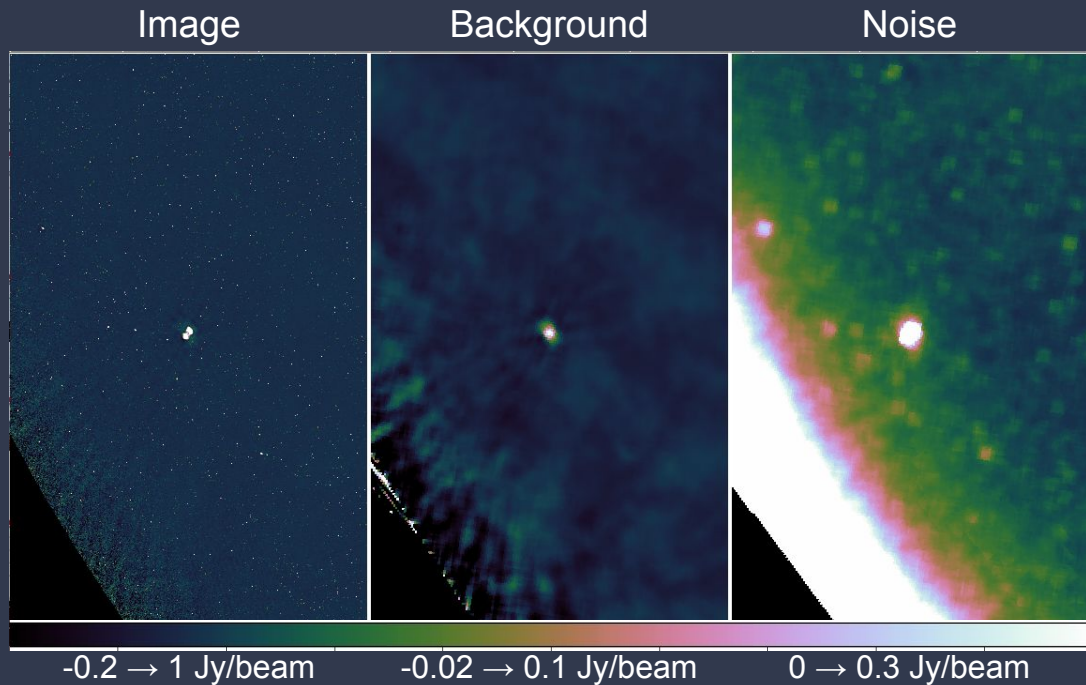
# Making

- Catalogues
- Light-curves
- SEDs

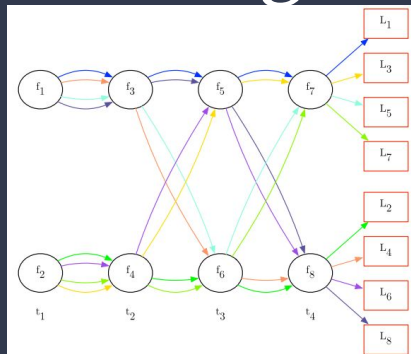
# Direction Dependant Synthesized Beam



# Catalogues at large FoV

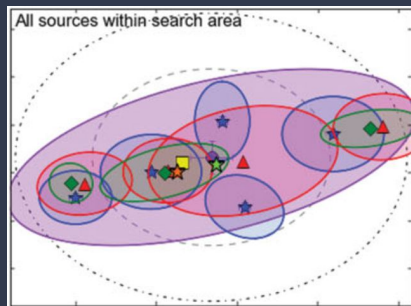


# Catalogues at large $\Delta\nu$



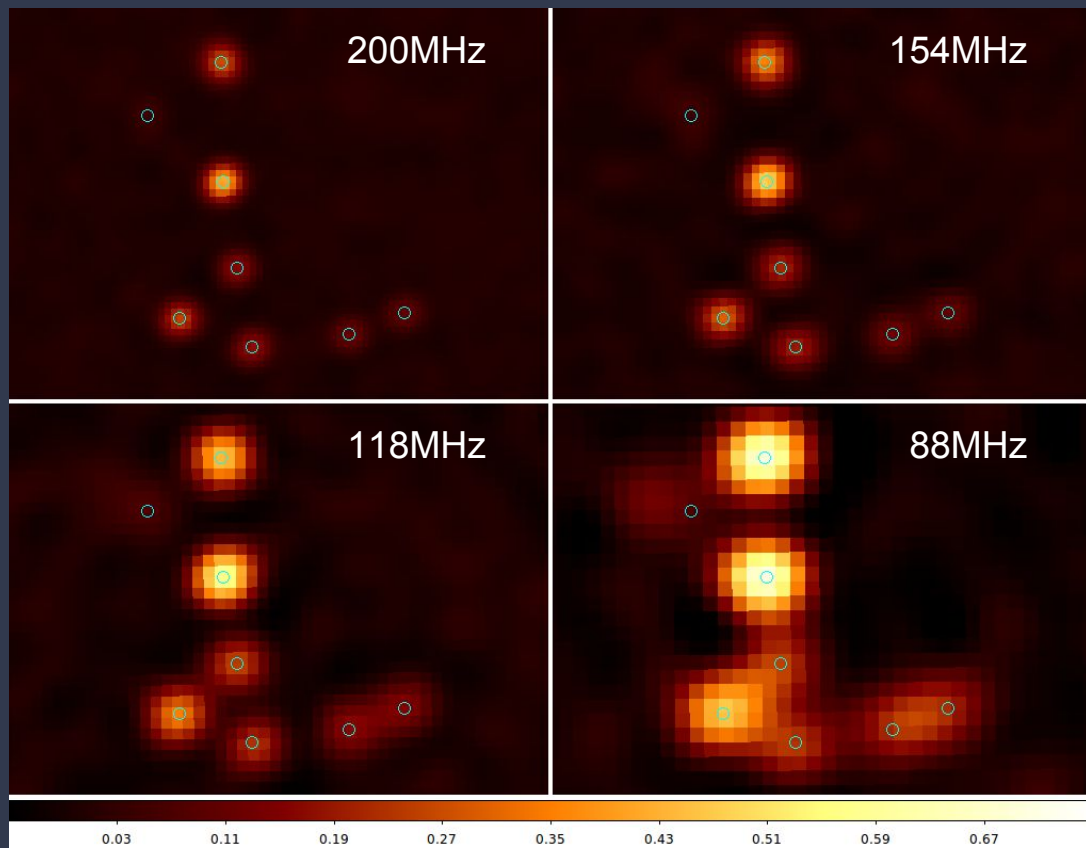
Catalogue  
and  
X-match?

Swinbank et al. 2015



Hierarchical  
association?

Line et al. 2017

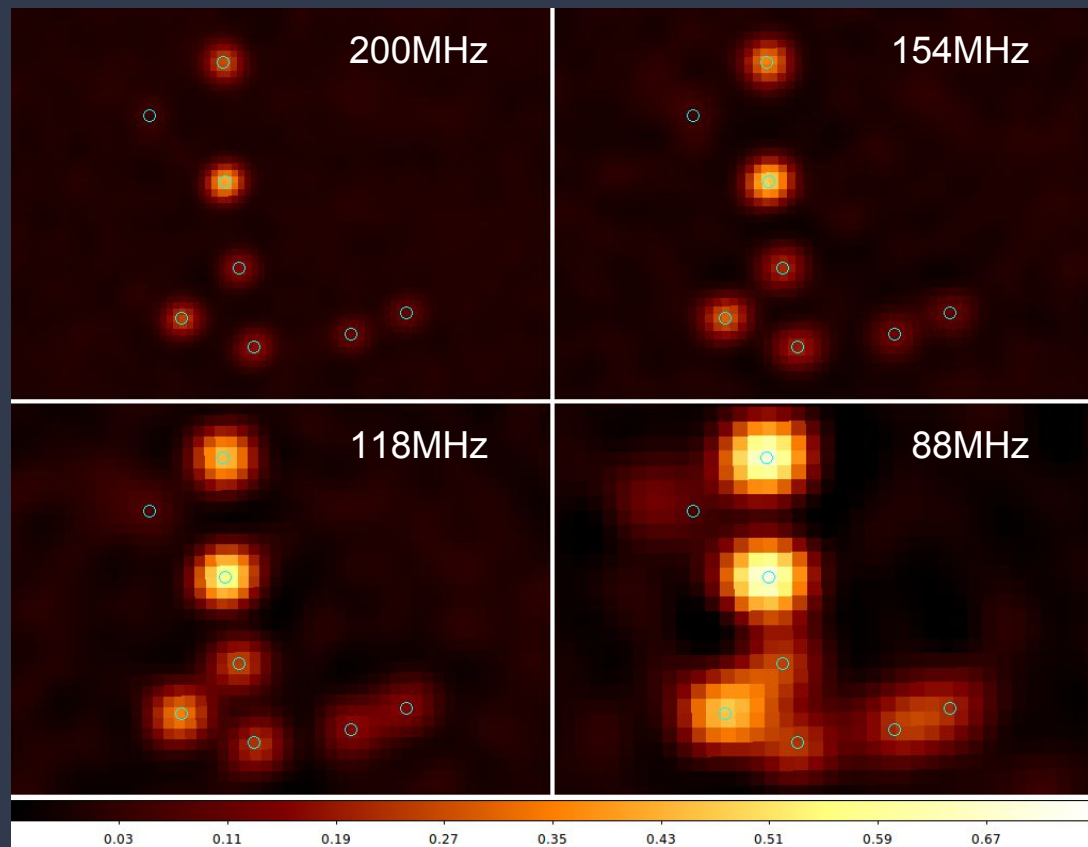


The diagram shows a directed graph with nodes  $i_1, i_3, i_5, i_7, i_2$  and labels  $l_1, l_3, l_5, l_7$ . A large green V-shape is overlaid on the graph, with its vertex pointing towards the bottom right. The nodes are arranged in a sequence from left to right, with  $i_1$  and  $i_2$  at the bottom left,  $i_3$  and  $i_5$  in the middle, and  $i_7$  at the top right. The labels  $l_1, l_3, l_5, l_7$  are arranged in a column on the far right. Directed edges connect the nodes and labels, with colors (blue, green, yellow, purple) indicating different types of edges. The green V-shape is a large, solid green area that covers the lower half of the diagram, with its vertex pointing towards the bottom right.

**Swinbar Too Hard!**

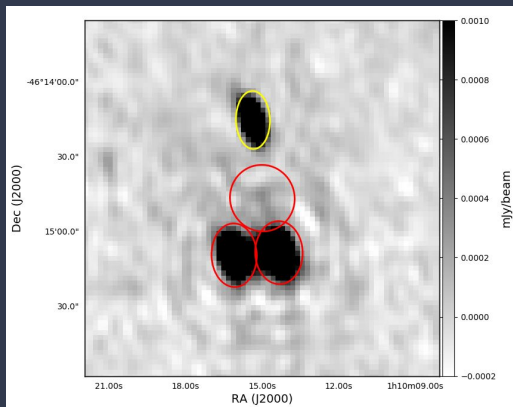
# Too Hard!

Line et al. 2017

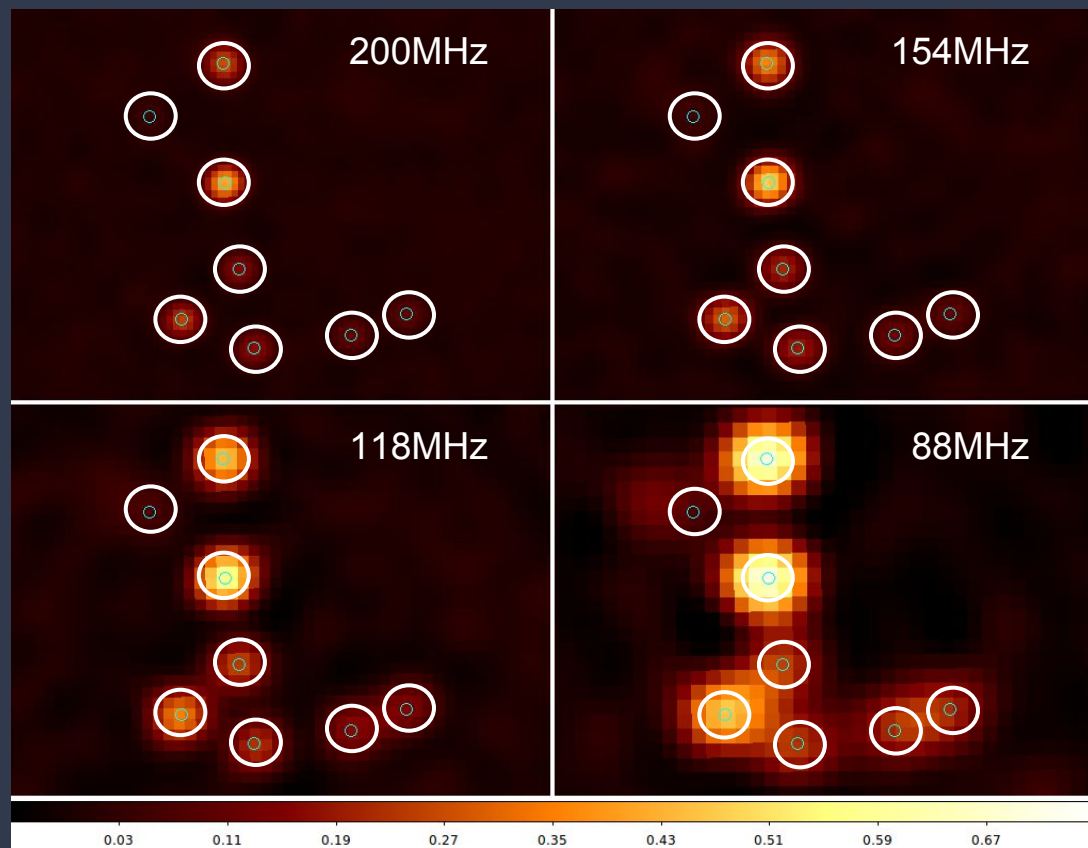




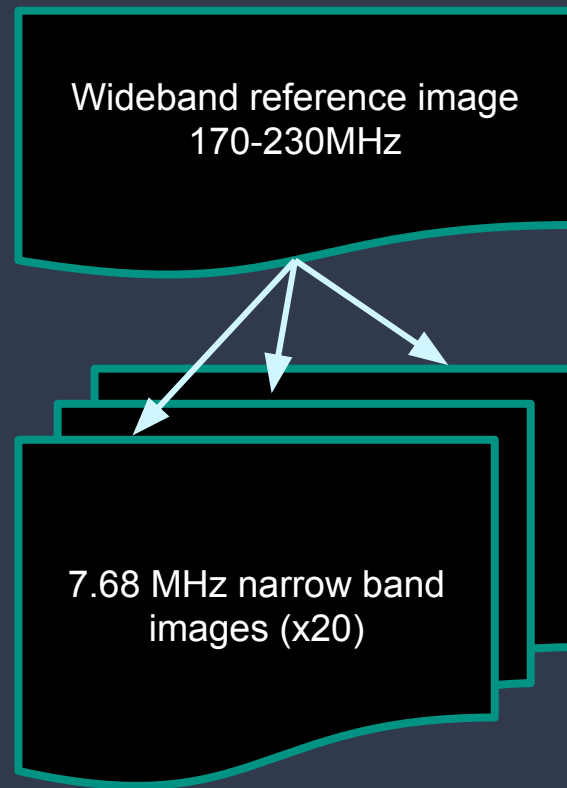
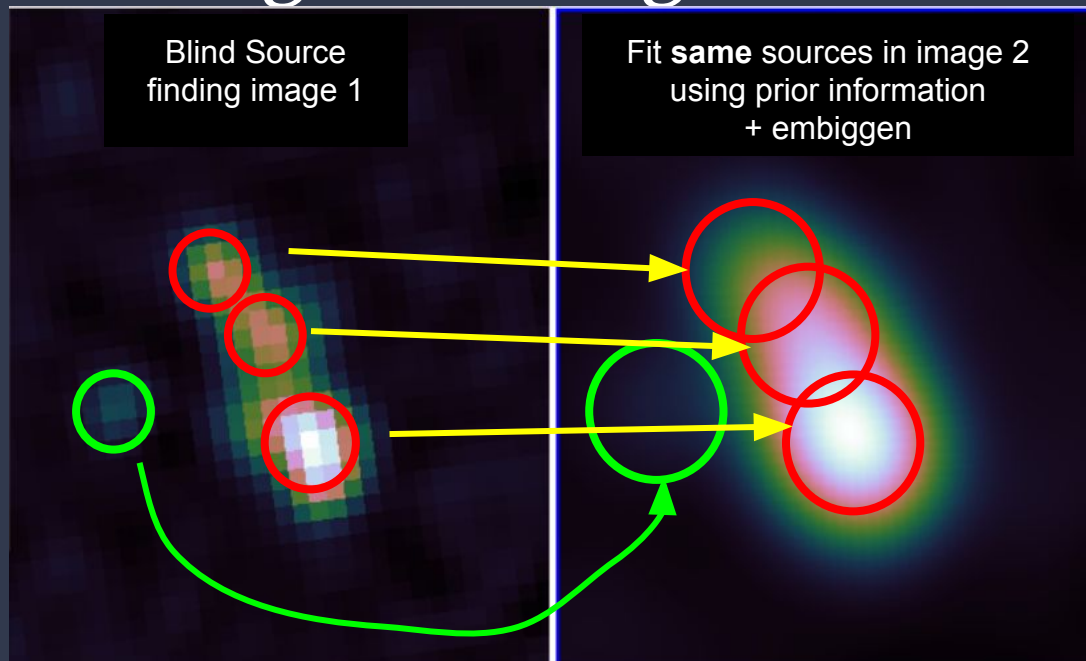
# Catalogues at large $\Delta\nu$



Priorized fitting with Aegean  
(Hancock et al. 2012/18)  
(now also pyBDSF)



# Catalogues at large $\Delta\nu$



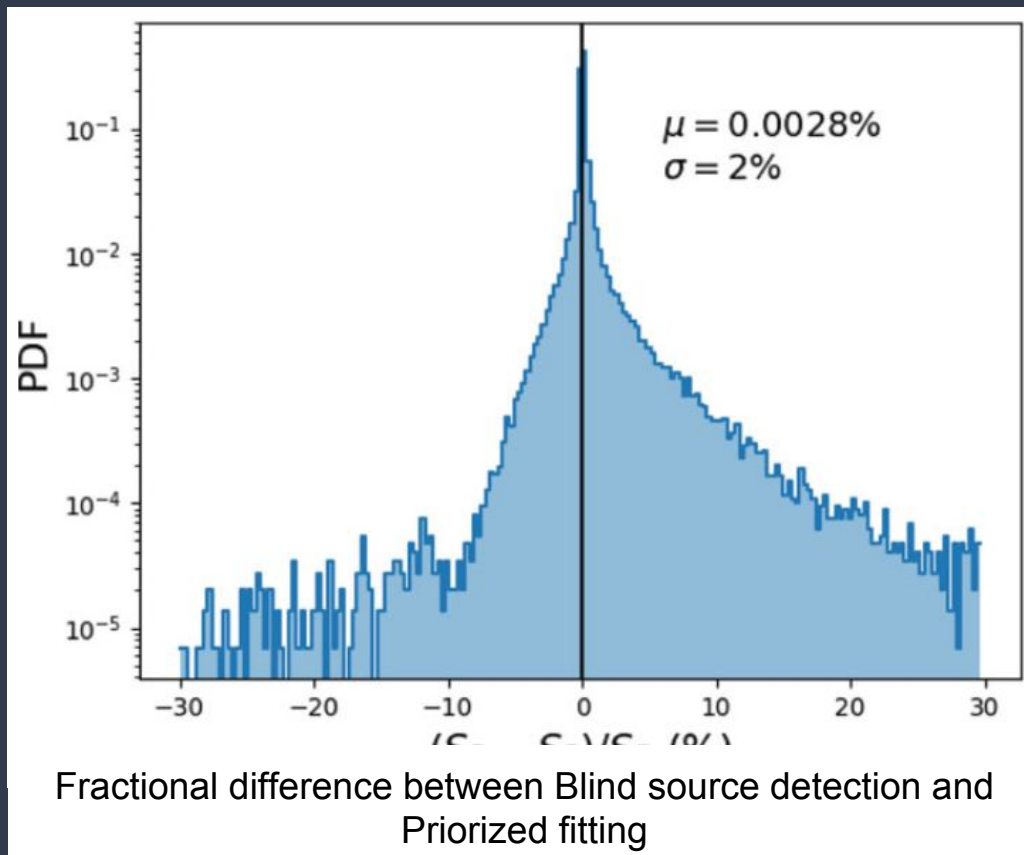
# Priorized fitting

Swapping

a **detection** experiment for  
a **measurement** experiment  
reduces uncertainties

Good astrometry is essential so use  
`fits_warp`:

Hurley-Walker & Hancock [2018arXiv180808017H](#)

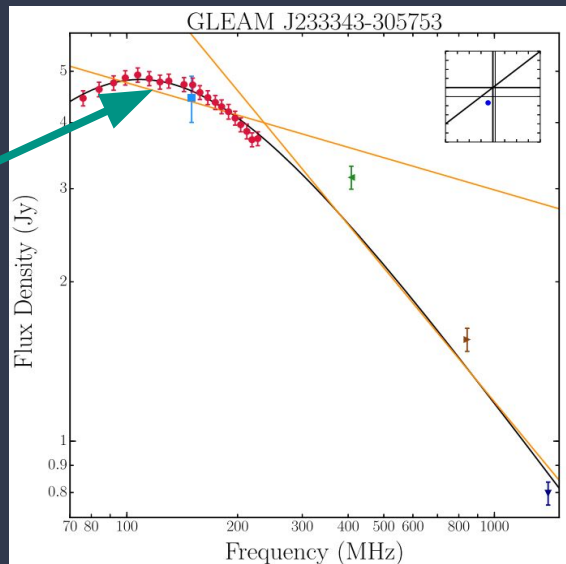


# Catalogues at large $\Delta\nu$

Catalog contains

- **all** sources from deep image
- fluxes from each narrow band for each source
- sub-threshold fluxes
- ZERO false cross ids

GLEAM prioritized fits  
at 20 frequencies



Wideband reference image  
170-230MHz

7.68 MHz narrow band  
images (x20)

# Source Finding Solution: Aegean

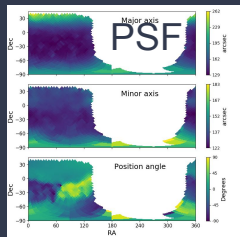
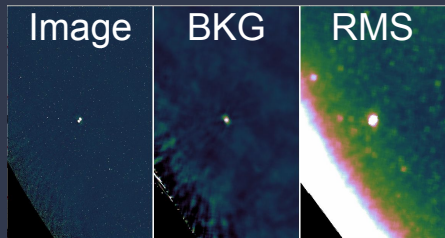
<https://github.com/PaulHancock/Aegean>



This repository

Search

$$\text{Min} \{ (f(\theta;x) - x)^T C^{-1} (f(\theta;x) - x) \}$$



Wideband  
reference image  
170-230MHz

7.68 MHz  
narrow band  
images (x20)

## Aegean

- Find sources
- Characterise sources
- Sources can be gaussians or blobs



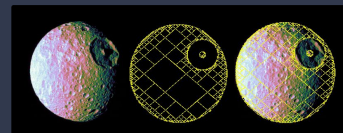
## BANE

- Characterise background
- Characterise noise
- Do it right
- Do it fast



## MIMAS

- Describe regions
- Combine regions
- Mask images
- Constrain Aegean
- Write MOC files



# Other solutions:

## Good ones:

- Selavy - Whiting & Humphries [2012PASA...29..371W](#)
- PyBDSF - Mohan & Rafferty [2015ascl.soft02007M](#)
- PySE - Carbone et al. [2018A&C....23...92C](#)

## Not good ones:

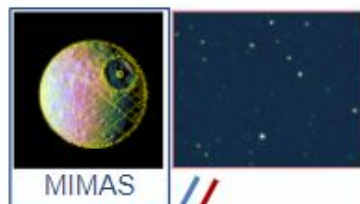
- imsad (miriad)
- SAD/VSAD (aips)
- SExtractor
- Blobcat

All-in-one  
solutions



# Survey image processing with the VAST pipeline

1. Create images. Exclude areas that are not of interest by: trimming, masking, or using MIMAS\* region files. (not part of pipeline)



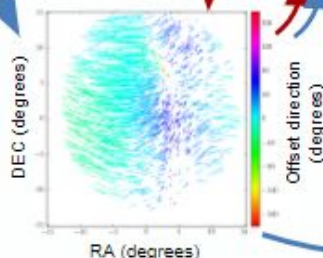
Two approaches can be used: a stream processing approach (blue flow) designed for real-time ingestion of images, and a batch processing approach (red flow) designed for already completed surveys.

2. Create background and noise images using BANE\*. Source finding and characterisation using Aegean\*. If existing catalogues are found, then source finding is not duplicated.



\*see [github.com/PaulHancock/Aegean](https://github.com/PaulHancock/Aegean)  
Hancock et al 2012, MNRAS, 422, 1812

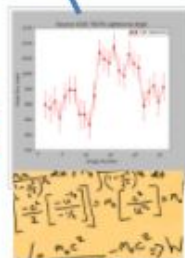
3. (optional, but recommended!) Crossmatch new measurements with a reference catalogue, and perform **astrometry and gain corrections**. This can reduce the ionospheric induced positional offsets from as much as 1arcmin, down to 5arcsec.



7. **Publish** awesome results! (not part of pipeline, awesomeness subject to input data)



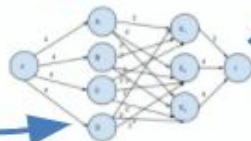
6. Generate **light curves** and calculate **variability statistics** for each source. Statistics include a measure of the magnitude and confidence of variability.



5. **Flux monitoring**. For any sources which have a measurement missing from an image, replace the missing measurement with a prioritized fit.

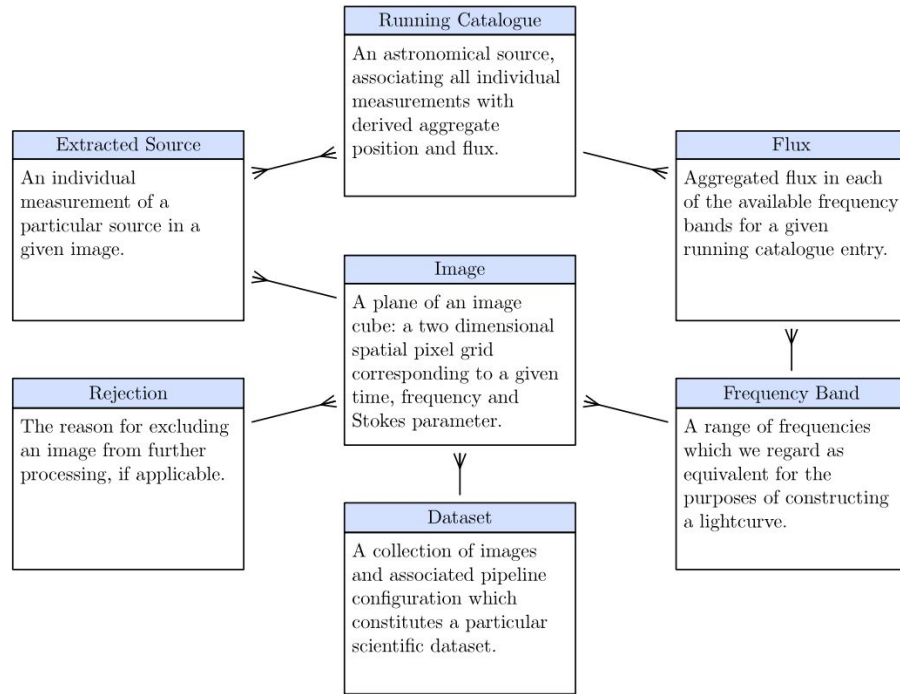


4. **Source association**: Regroup all the individual flux measurements into sources. A source will have at most one flux measurement per image.

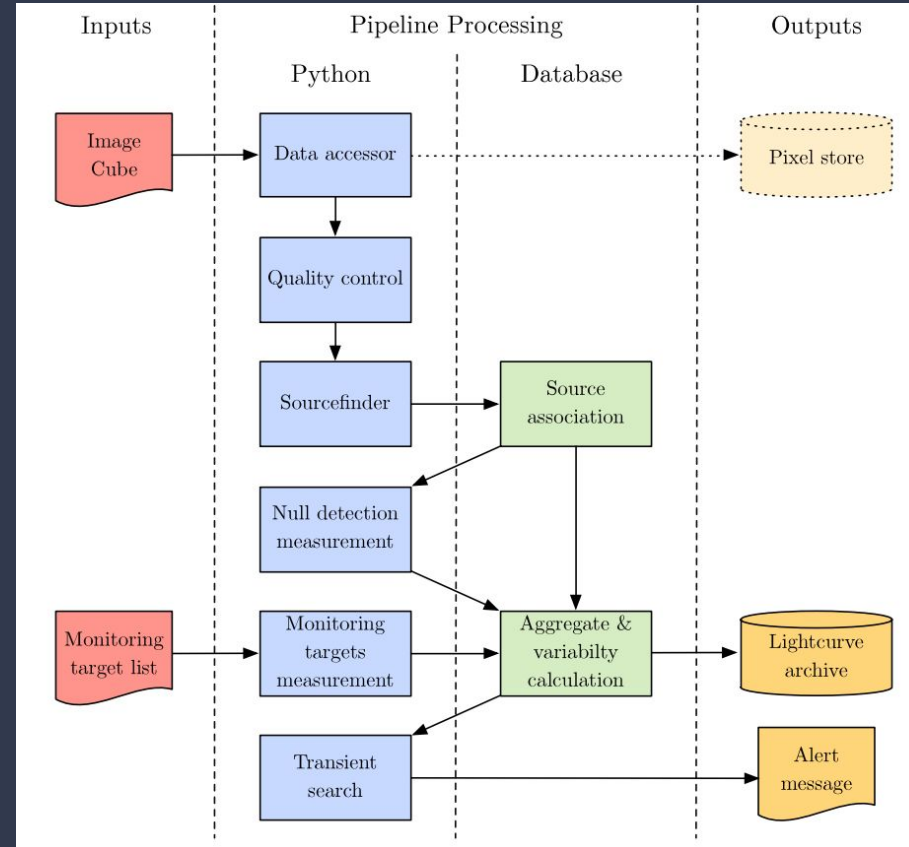




# LOFAR – TraP



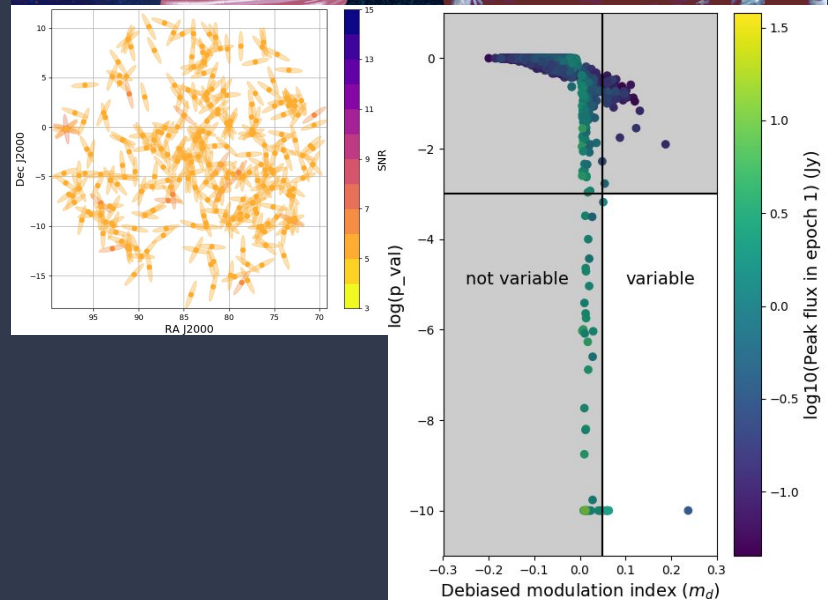
Swinbank et al. [2015A&C...11...25S](#)



# Robbie

1. Astrometry correct each epoch
2. Stack to form mean image
3. Find persistent source in mean image
4. Mask persistent sources in single epoch
5. Create light curves for persistent sources
6. Blind search for transient candidates in single epochs
7. Identify transients and characterise variability

<https://github.com/PaulHancock/Robbie>  
(Astronomy & Computing, Submitted)



# Further reading

Condon [1997PASP..109..166C](#) Empirical measure of errors  
Refreiger & Brown 1998 [arXiv:9803279](#) analytical treatment of uncertainty and bias  
Hancock et al. 2012 [2012MNRAS.422.1812H](#) Source finding with Aegean  
Hancock et al. 2018 [2018PASA...35...11H](#) Source finding on correlated data  
Whiting & Humphries [2012PASA...29..371W](#) ASKAP soft source finder  
Mohan & Rafferty [2015ascl.soft02007M](#) LOFAR source finder  
PySE - Carbone et al. [2018A&C....23...92C](#) LOFAR source finder (for TraP)  
Hurley-Walker & Hancock 2018 [2018arXiv180808017H](#) Correcting ionospheric effects in the image plane  
Banyer et al [2012ASPC..461..725B](#) VAST pipeline  
Hancock et al. [2018ascl.soft08011H](#) Robbie (= vast lite / vast ++)  
Swinbank et al. [2015A&C....11...25S](#) LOFAR TraP