

International Centre for Radio Astronomy Research

EXPLORING THE LIMITS OF THE BAYESIAN UNIVERSE: HOW TO TACKLE BREADTH AND DEPTH







Government of Western Australia
Department of the Premier and Cabinet
Office of Science

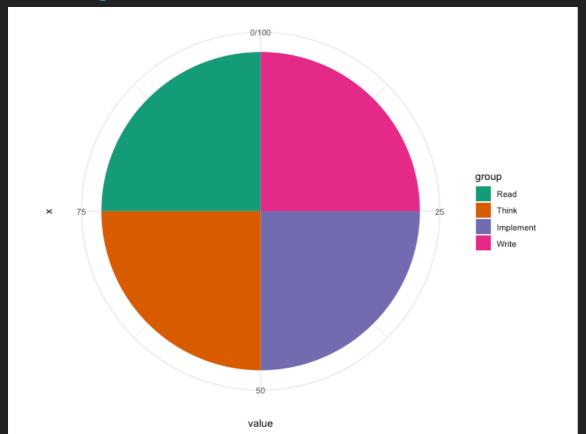
AARON ROBOTHAM

Sabine Bellstedt, Jessica Thorne

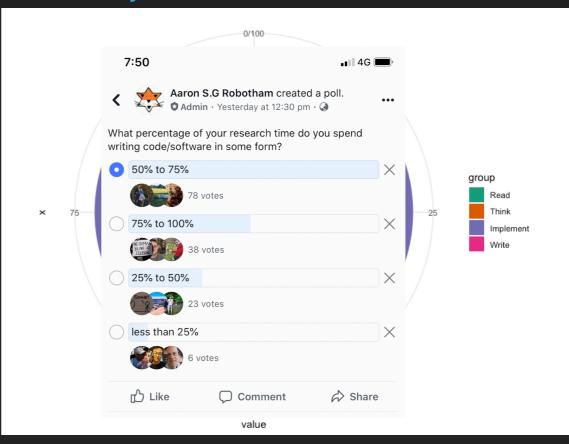
OPEN-SOURCE STATISTICS AND COMPUTING

WHAT WE DO WHEN WE RESEARCH:

The Legend



The Reality



In a poll on the Astrostatistics Facebook group, 80% of people reported spending more than half (average around 70%) of their research time writing code/software on a computer.

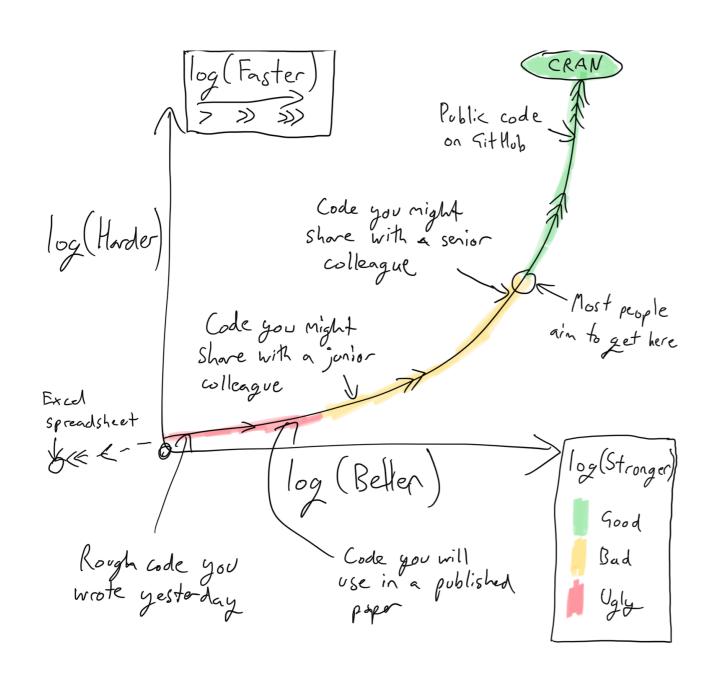
Better: commented, maintainable

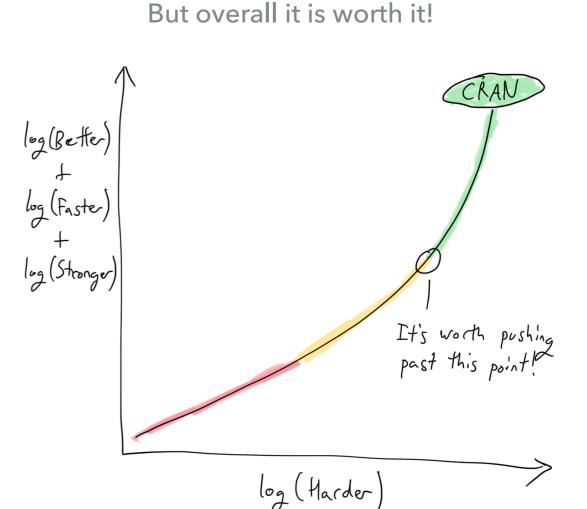
Faster: does the job faster

Stronger: Reliable, fault tolerant, robust

OPEN-SOURCE STATISTICS AND COMPUTING

SO WE MUST BE GREAT AT THIS, RIGHT? WELL...



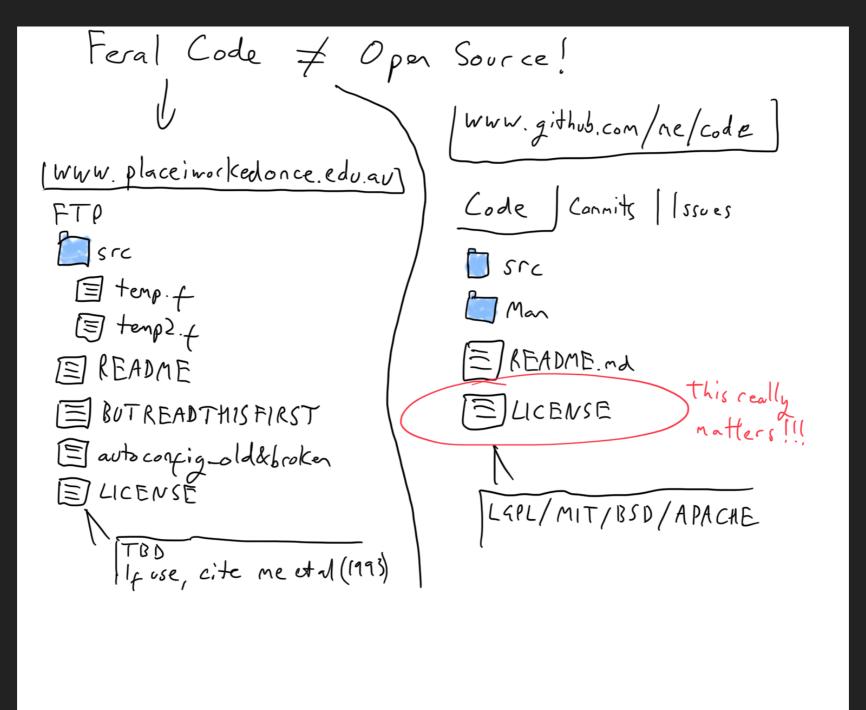


OPEN-SOURCE STATISTICS AND COMPUTING

THE SOLUTION

- ▶ We (as a community) bootstrap the "implement" part wherever feasible.
- ▶ Even when we re-invent the wheel, knowing what wheels are available makes this:
 - A) worthwhile (need a roller-blade wheel in a world of cartwheels)
 - ▶ B) efficient (we can still borrow ideas make it round!).
- ▶ To help our community this means the following:
 - Aggressively public code +
 - Default collaborative efforts =
 - Open-source analysis (statistical methods and software)

BUT REMEMBER!



OPEN-SOURCE STATISTICS AND COMPUTING

WHY WE SHOULD DO OPEN-SOURCE STATISTICS IN ACADEMIA?

- Because it is the new fun thing!
- ▶ I keep seeing talks about GitHub, maybe it will help me get a job...?
- My grant funding insists.
- Better than nothing for backing things up.
- Paying for things sucks.
- ▶ To make the world a better place...
- These are not the best reasons actually... we can do better by appealing to our more ruthless and selfish instincts!

OPEN-SOURCE STATISTICS AND COMPUTING

WHY WE REALLY SHOULD DO OPEN-SOURCE STATISTICS IN ACADEMIA?

- Only way to efficiently collaborate on complex projects, and drops the bar to entry to almost nothing. If
 people cannot see it then they will not ask to help.
- Get a lot more credit and exposure to your work (you will start appearing all over the internet).
- It instills a much higher level quality of work: if you do not feel comfortable sharing it then you should not feel comfortable publishing it! As Russians say: "Trust, but verify".
- It provides free (!) wide-coverage quality checking you cannot replicate alone.
- It allows you to work seamlessly on projects across multiple devices without penalty (I often show people snippets from GitHub on my phone handy in a pinch).
- ▶ It naturally enforces modern version control by any practical route (Git / GitHub).
- ▶ It helps less senior members of our community. What does this mean...?

mSFR 10000 **CHARM Iteration Number**

STATISTICS AND COMPUTING

BARRIERS

TOOLS TO BE A MODERN ASTRONOMER

- A new PhD student in a highly mature field like astronomy (a few thousand years of research and counting) has a lot of catching up to do. They need:
 - Solid coding skills (no matter what they do).
 - Relatively sophisticated knowledge of data analytics and statistics.
 - Broad understanding of the current state of the field, and a deep understanding of their own sub-field.
 - Excellent project management.
 - Good scientific writing (an un-taught skill usually)
 - Functional people skills.

OPEN-SOURCE STATISTICS AND COMPUTING

A CASE STUDY WITH HYPER-FIT

- Astronomy is unusual within physics in how expensive it is to get observations, so we tend to have un-ignorable errors, and available techniques like PCA and SVM rarely support the treatment of errors *properly*.
- Because of some mass-size data, I started wondering about hyperplane fitting with heteroscedastic covariant errors, and published a paper tackling it:

Publications of the Astronomical Society of Australia (PASA), Vol. 32, e033, 14 pages (2015). © Astronomical Society of Australia 2015; published by Cambridge University Press. doi:10.1017/pasa.2015.33

Hyper-Fit: Fitting Linear Models to Multidimensional Data with Multivariate Gaussian Uncertainties

A. S. G. Robotham and D. Obreschkow

ICRAR, M468, University of Western Australia, Crawley, WA 6009, Australia Email: Aaron.robotham@uwa.edu.au

(RECEIVED July 01, 2015; ACCEPTED August 10, 2015)

$$\ln \mathcal{L} = -\frac{1}{2} \sum_{i=1}^{N} \left[\ln \left(\sigma_{\perp}^{2} + \frac{\mathbf{n}^{\mathsf{T}} \mathbf{C}_{i} \mathbf{n}}{\mathbf{n}^{\mathsf{T}} \mathbf{n}} \right) + \frac{(\mathbf{n}^{\mathsf{T}} [\mathbf{x}_{i} - \mathbf{n}])^{2}}{\sigma_{\perp}^{2} \mathbf{n}^{\mathsf{T}} \mathbf{n} + \mathbf{n}^{\mathsf{T}} \mathbf{C}_{i} \mathbf{n}} \right]$$

This involves N-dimensional projections and careful treatment of the data covariance terms. This is all fiddly stuff that took me a lot of sanity checking to implement robustly (e.g. adapting for geometric corner cases). Now the question is, would you confidently trust a new PhD student to implement the above in code you need to use?

A CASE STUDY WITH HYPER-FIT

> So a paper with an equation is often not enough. How about an R CRAN package?

hyper.fit: Generic N-Dimensional Hyperplane Fitting with Heteroscedastic Covariant Errors and Intrinsic Scatter

Includes two main high level codes for hyperplane fitting (hyper.fit) and visualising (hyper.plot2d / hyper.plot3d). In simple terms this allows the user to produce robust 1D linear fits for 2D x vs y type data, and robust 2D plane fits to 3D x vs y vs z type data. This hyperplane fitting works generically for any N-1 hyperplane model being fit to a N dimension dataset. All fits include intrinsic scatter in the generative model orthogonal to the hyperplane.

Version: 1.1.0

Depends: $R (\ge 3.00), \underline{\text{magicaxis}}, \underline{\text{MASS}}, \underline{\text{rgl}}, \underline{\text{LaplacesDemon}}$

Published: 2019-01-31

Author: Aaron Robotham and Danail Obreschkow

Maintainer: Aaron Robotham <aaron.robotham at uwa.edu.au>

License: <u>GPL-3</u> NeedsCompilation: no

Citation: <u>hyper.fit citation info</u>

Materials: <u>NEWS</u>
CRAN checks: hyper.fit results

Downloads:

Windows binaries: r-devel: https://doi.org/10.2121/n.nc/2012/n.n

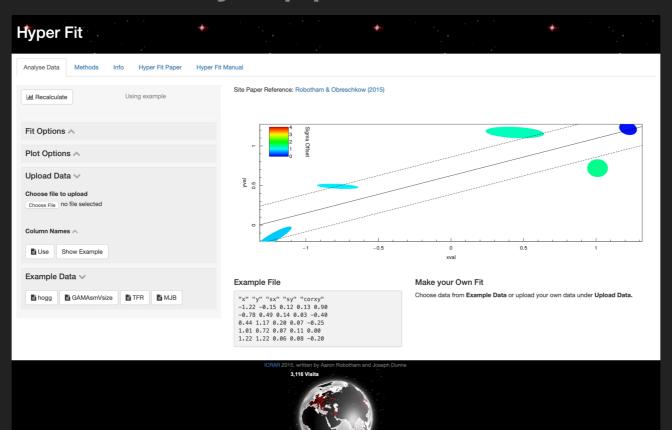
OS X binaries: r-release: hyper.fit 1.1.0.tgz, r-oldrel: hyper.fit 1.1.0.tgz

Old sources: <u>hyper.fit archive</u>

The bar for this is high (*much* higher than a simple code repo like SourceForge or PyPI)- it requires full code documentation (40 pages!) and a full suite of useable examples.

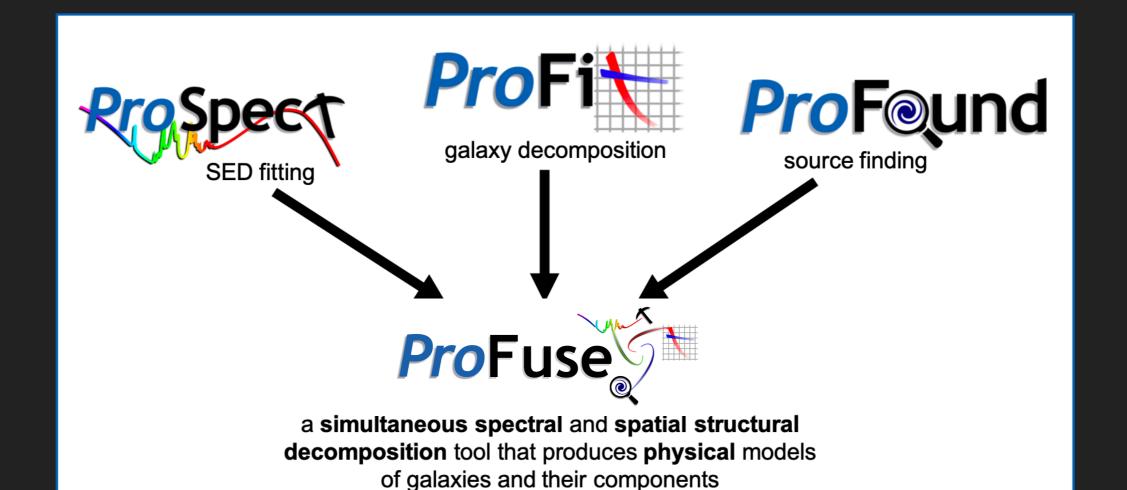
A CASE STUDY WITH HYPER-FIT

Is that enough? Well lots of people do not use R, fear it, and will not read a manual (arguably they do not deserve my help!). Behold Shiny Apps:



hyperfit.icar.org

https://www.icrar.org/our-research/tools/ The HI Fidelity Calculator hifi.icrar.org Cosmology Calculator Cosmocal c.icrar.org Info z, Width, Flux Angles & Beam Flux/Mass Flux/nHI Fluxd/Tb Noise SNR What's Up? ✓ Plot: Survey Design Sky Separation Info B Code whatsup.icrar.org The redshift z is 1 Change the par dependence ple The observed solid angle of the source is 1 arcsec2 Results: The observed flux density is 1e-06 Jy Calculate The redshift z is 1 The area of the source (small angle approx) is 64.14 kpc2 Set Variables The expansion factor a is 0.5 Y2019M11D6, RA 125, Dec 19, Lon 116, Lat -32, Mphase 0.67, Msep 150, Ssep 100 General Options Brightness Temperature vs z Distance: User +19:28:33.71 The Comoving Radial Distance Tra Longitude Log x axis Log y axis 1-OmegaM[0] -31.97333 OmegaR[0] Local Date Save Data 2019-11-06 Hubble's constant H at z is 123.248 Calc fSigma8 The Deceleration parameter (g) at: Custom Calc Altitude 15 16 17 18 19 20 21 22 23 00 01 02 03 04 05 06 07 08 09 10 11 12 13 Look-back Time prospect.icrar.org * **ProSpect SED** astromap.icrar.org Astro Map Description Tool Get Code Contact Redshift **Aaron Robotham** GitHub: asgr/ProSpec Map Projection Light Cones Deep Fields SFH/AGN: **Dust inputs:** hmfcalc.icrar.org 10⁻¹⁰ Projection MW Plane 10-1 Young Mass (log10): MW Cente Tau Screen: Long-Cent (0,180) 0.3 10 ~ 10⁻¹ Coordinate Type Mid Mass (log10): 10⁻¹ Tau AGN: Celestial 10⁻¹ Alpha SF Birth: 10-1 0 10-1 10-19 Alpha SF Screen: 10-20 10⁻² Alpha SF AGN: Wavelength (Angstrom) 10¹⁰ 10^{11} 10¹² 10¹ WAVES-Wide DEVILS Mass $(M_{\odot}h^{-1})$



FROM PIXELS TO SCIENCE

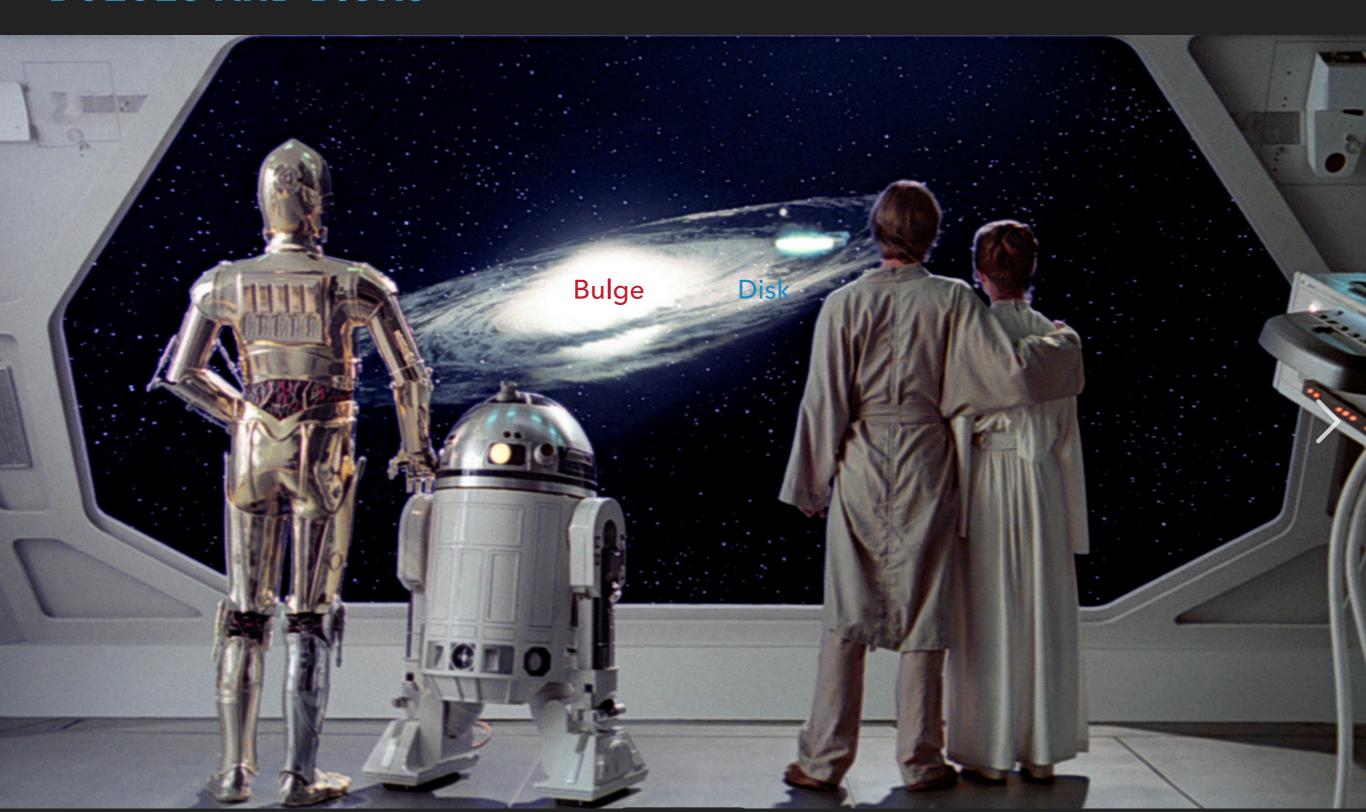
PROFUSE

MNRAS **513**, 2985–3012 (2022) Advance Access publication 2022 April 15 https://doi.org/10.1093/mnras/stac1032

ProFuse: physical multiband structural decomposition of galaxies and the mass-size-age plane

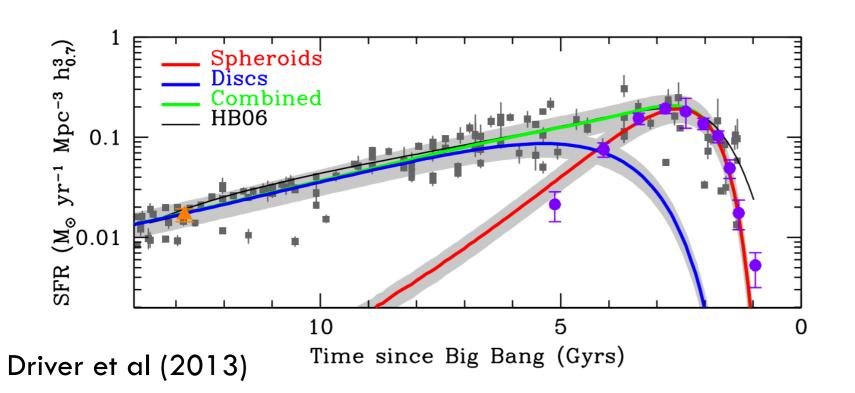
A. S. G. Robotham ⁶, ^{1,2} S. Bellstedt ⁶ and S. P. Driver ⁶

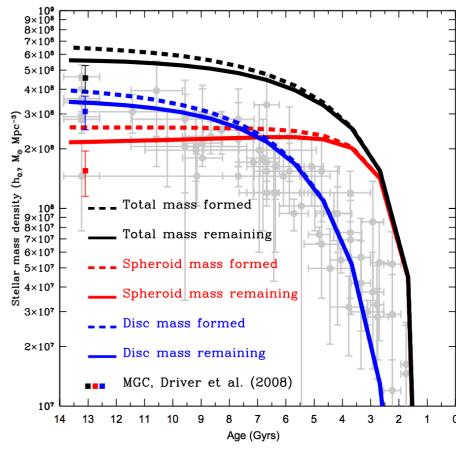
BULGES AND DISKS



Why Decompose Galaxies?

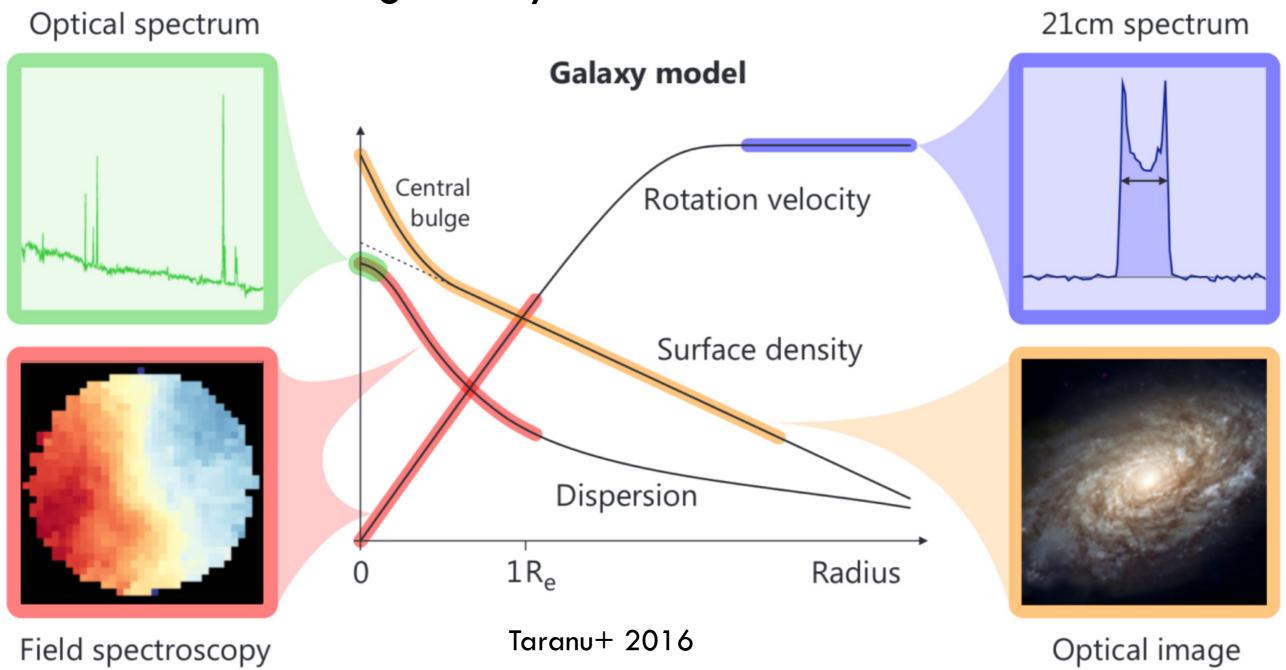
- 2 reasons:
 - But the more important reason is that the components of galaxies contain the evolutionary origin of galaxies.
 - Two dominant paradigms:
 - Rapid assembly of bulges (star bursts and mergers)
 - Slower accretion growth of disks





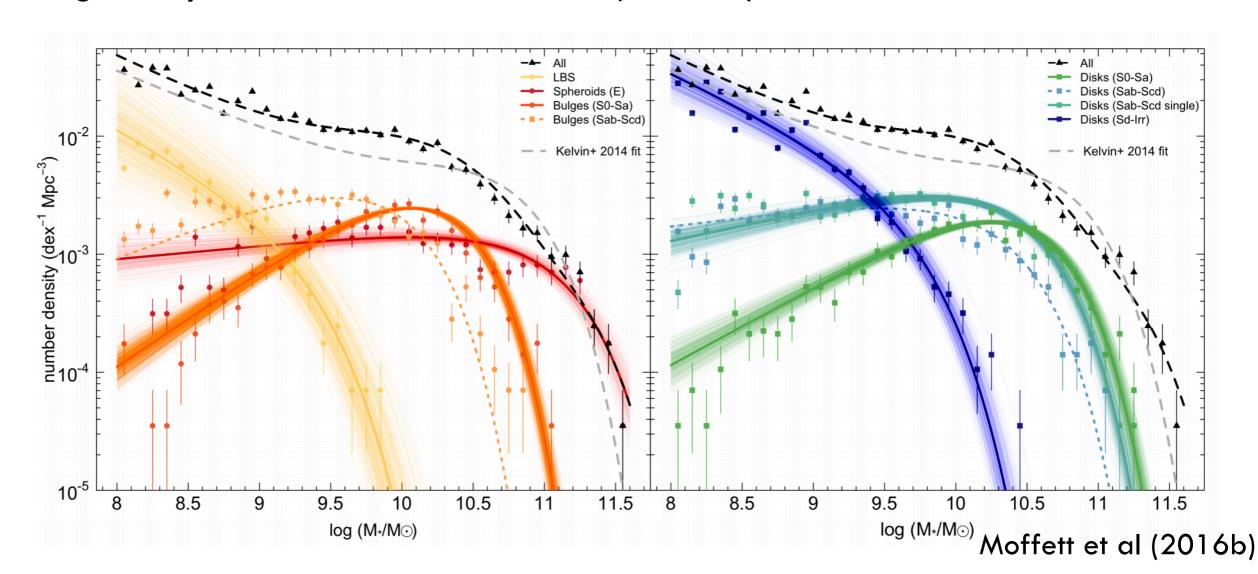
Why Decompose Galaxies?

The optical properties of galaxies are key to understanding the dynamics:



Why Decompose Galaxies?

The two major components also show remarkably different levels of significance for different mass galaxies (high mass = bulge dominated, low mass = disk dominated). Even the galaxy stellar mass function (GSMF) is non-trivial in detail:



THE PROFUSE PIPELINE

- ProFound Detect sources automatically, estimate rough model parameters, create segmentation and error maps. See Robotham+ (18), Bellstedt+ (20), GAMA/DEVILS/WAVES.
- ProFit Parameter driven single image galaxy decomposition. See Robotham+ (17), Cook+ (18/19), Casura+ (22).
- ProSpect Multi-band spectral energy distribution generation and fitting. See Robotham+ (20), Bellstedt+ (20/21), Thorne+ (21/22ab).
- ProFuse = ProFound + ProFit + ProSpect
 See Robotham+ (22)



EXTRACTING SOURCES AND GETTING READY FOR PROFIT

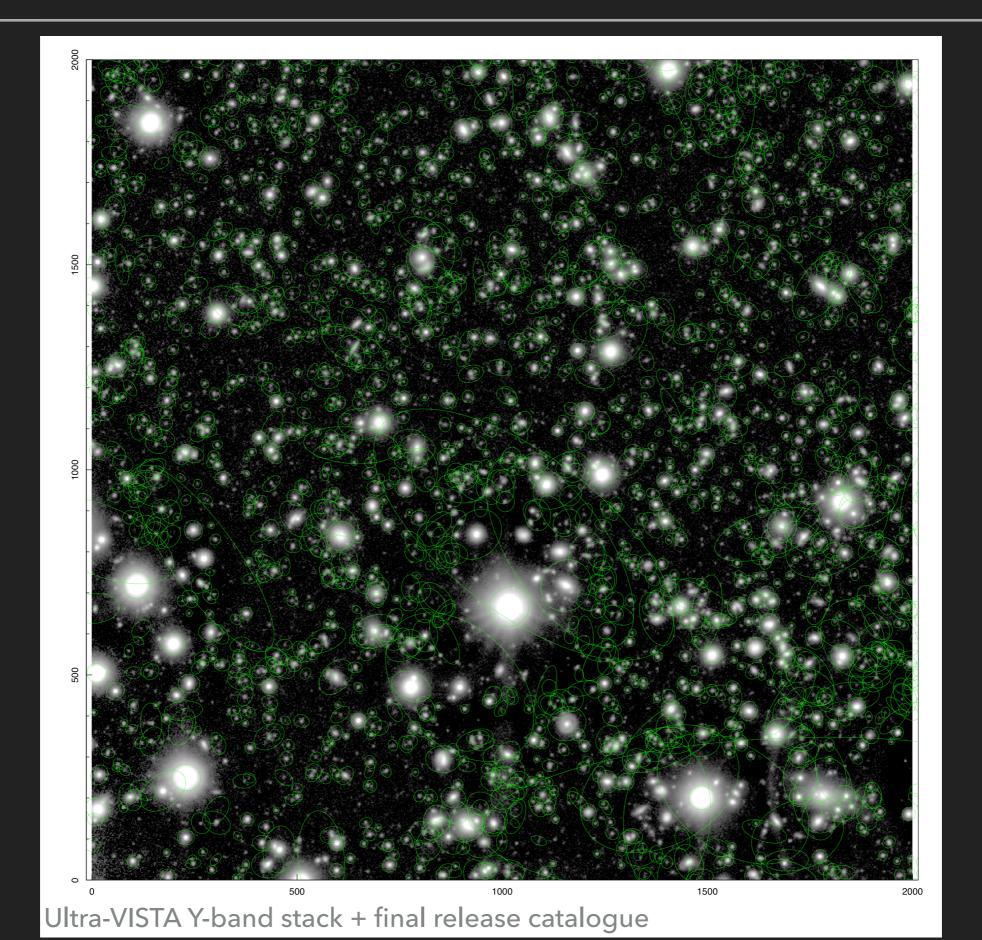
PROFOUND

MNRAS **476**, 3137–3159 (2018) Advance Access publication 2018 February 23 doi:10.1093/mnras/sty440

PROFOUND: SOURCE EXTRACTION AND APPLICATION TO MODERN SURVEY DATA

A. S. G. Robotham, ¹* L. J. M. Davies, ¹ S. P. Driver, ¹ S. Koushan, ¹ D. S. Taranu, ^{1,2} S. Casura³ and I. Liske³

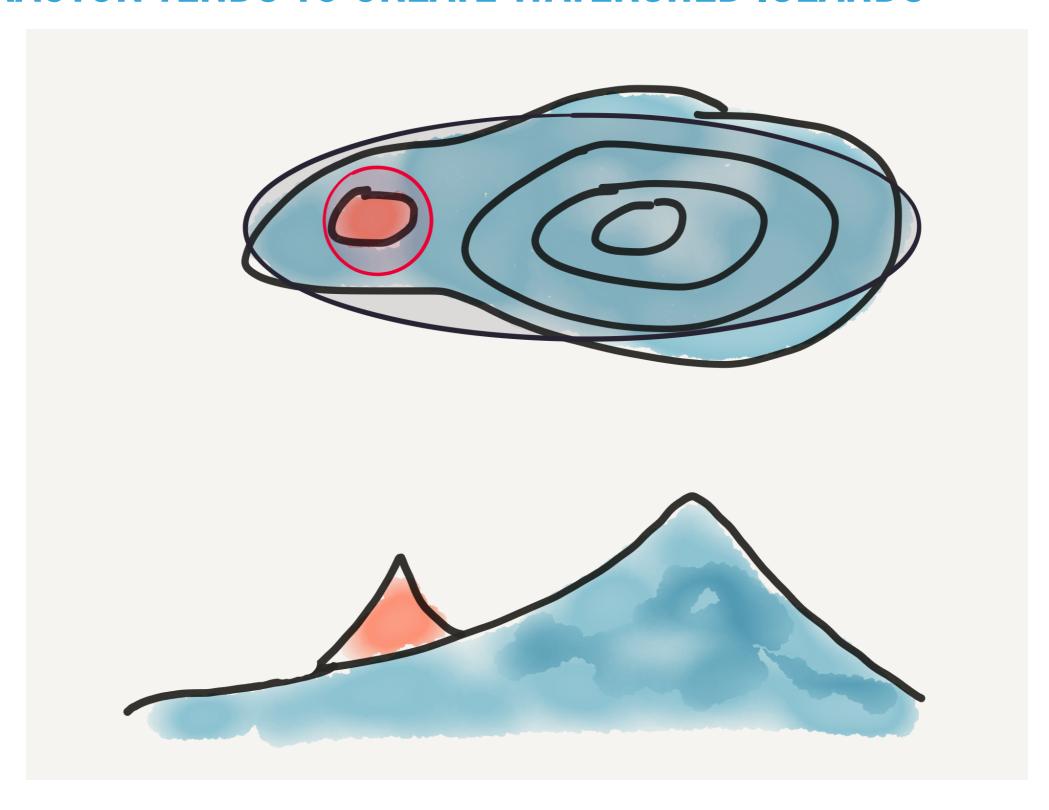
THE PROBLEM WITH SEXTRACTOR APERTURES



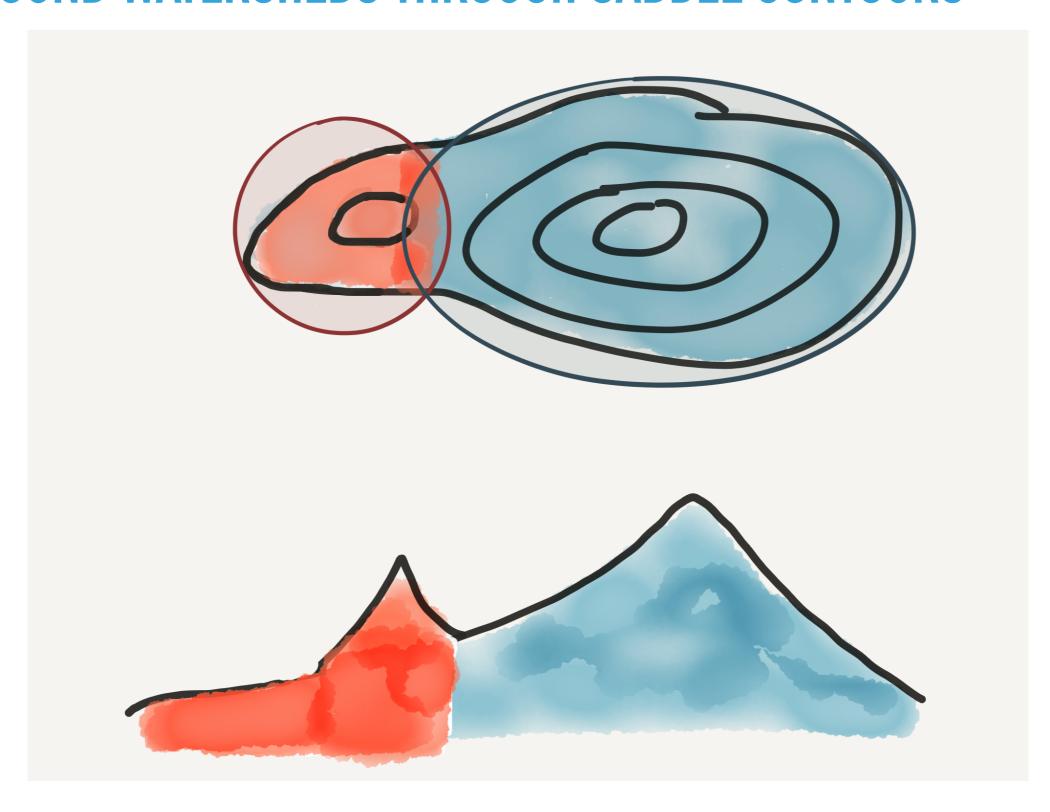
THE NEED TO START AGAIN

- For GAMA we investigated a lot of effort in manual aperture fixing. This was not scaleable or transferable to other data sets, we sought a better solution for WAVES.
- In short, I started again with the source extraction.
- It was not obvious what improvements might be possible over SExtractor (given how well tested and established it is) but two areas quickly came to light:
 - It does not watershed de-blend optimally (the most common failure we see is due to this). It does coarse island-based deblending.
 - It uses strictly elliptical apertures and then tries to distribute overlapping flux using a number of opaque internal schemes.

SEXTRACTOR TENDS TO CREATE WATERSHED ISLANDS



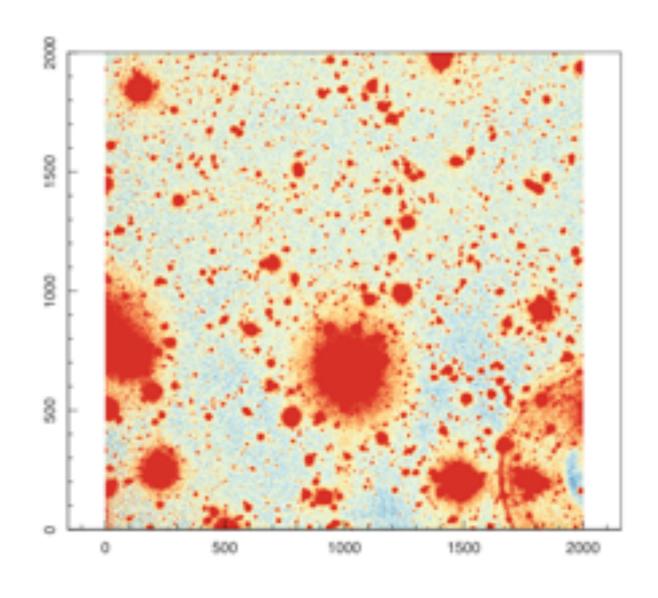
PROFOUND WATERSHEDS THROUGH SADDLE CONTOURS

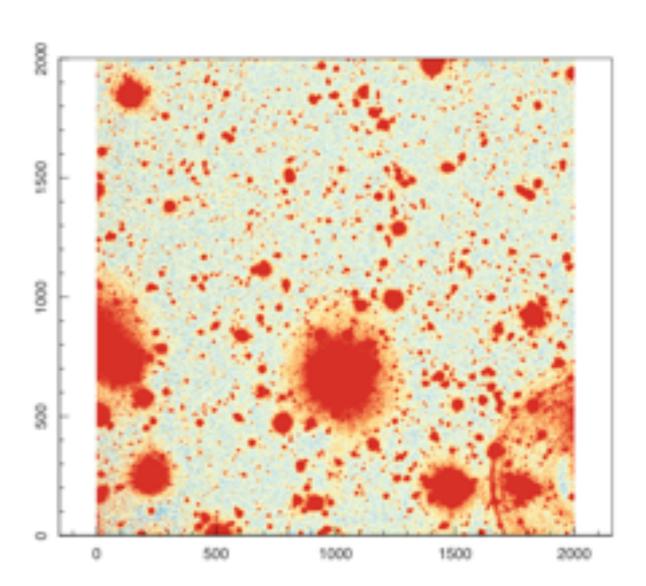


OUR PROFOUND SOLUTION GETS ROUND BOTH THESE ISSUES

- We use a similar approach to find the initial high S/N images segments:
 - Careful sky subtraction (iterative masking and clipping)
 - S/N=1.5 threshold as standard (can change)
 - Segments are de-blended to some tolerance (using a different algorithm to SExtractor- non-discretised surface brightness / sky-RMS thresholds and locally adaptive).
 - > Segments are grown organically- apertures never used.

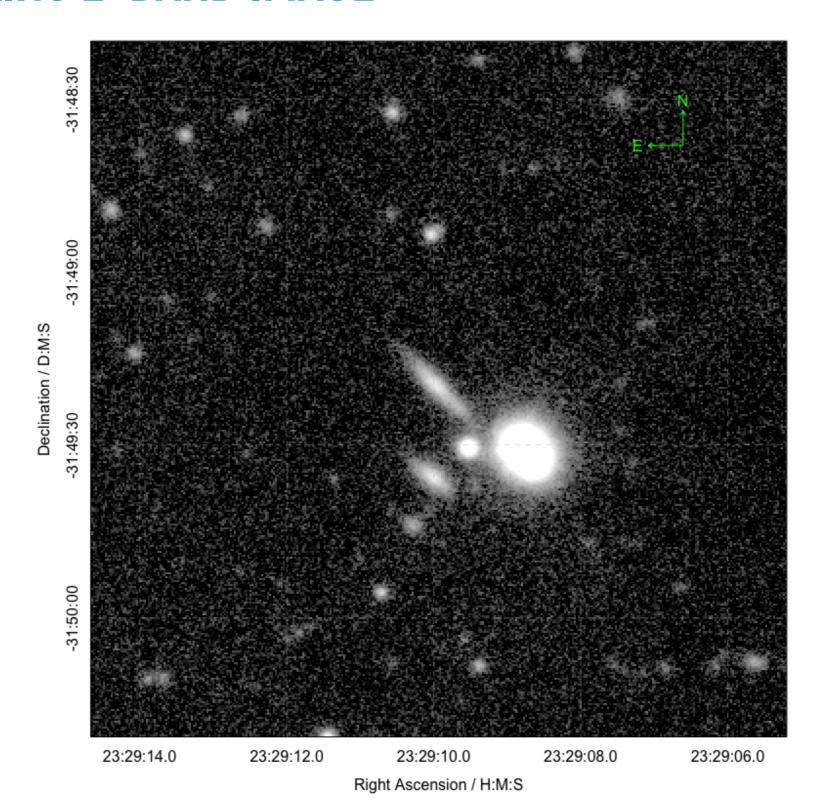
PROFOUND USES AN AGGRESSIVE MESH BASED SCHEME



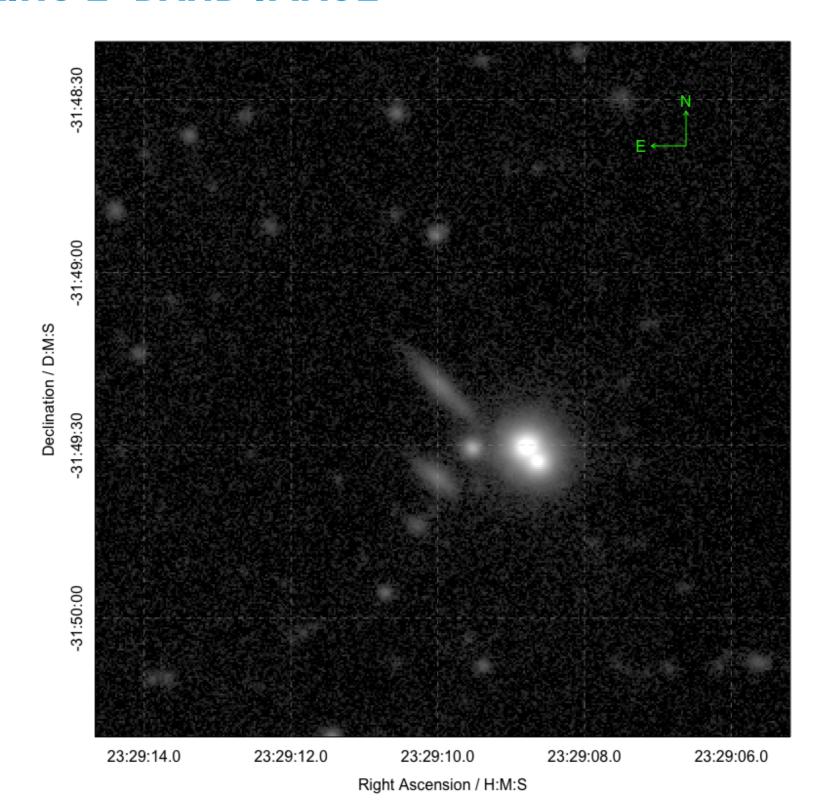


Ultra-VISTA Sextractor sky versus ProFound sky

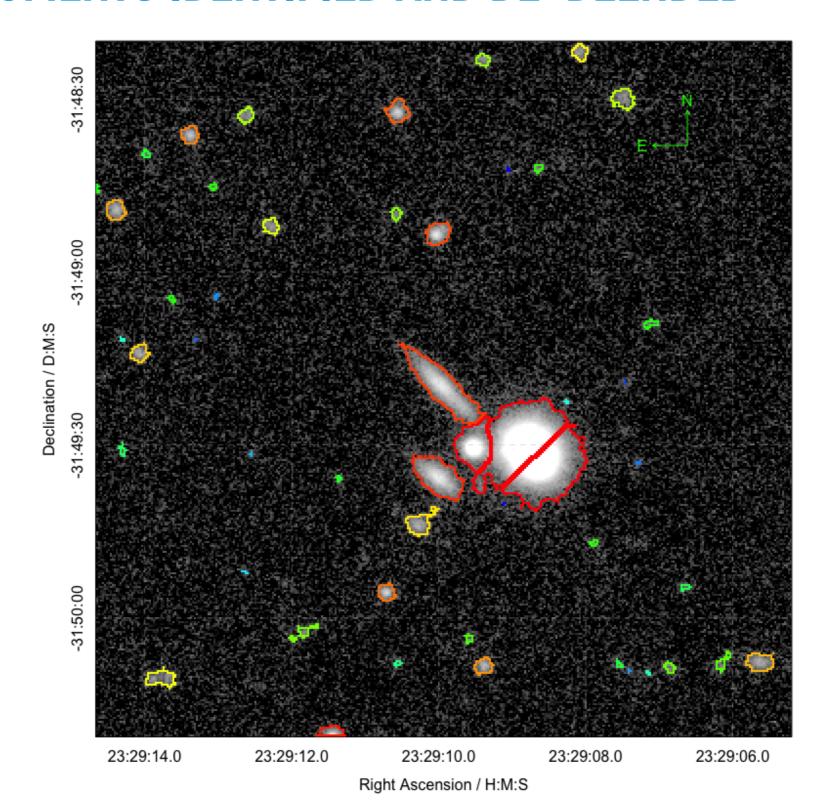
INITIAL VIKING Z-BAND IMAGE



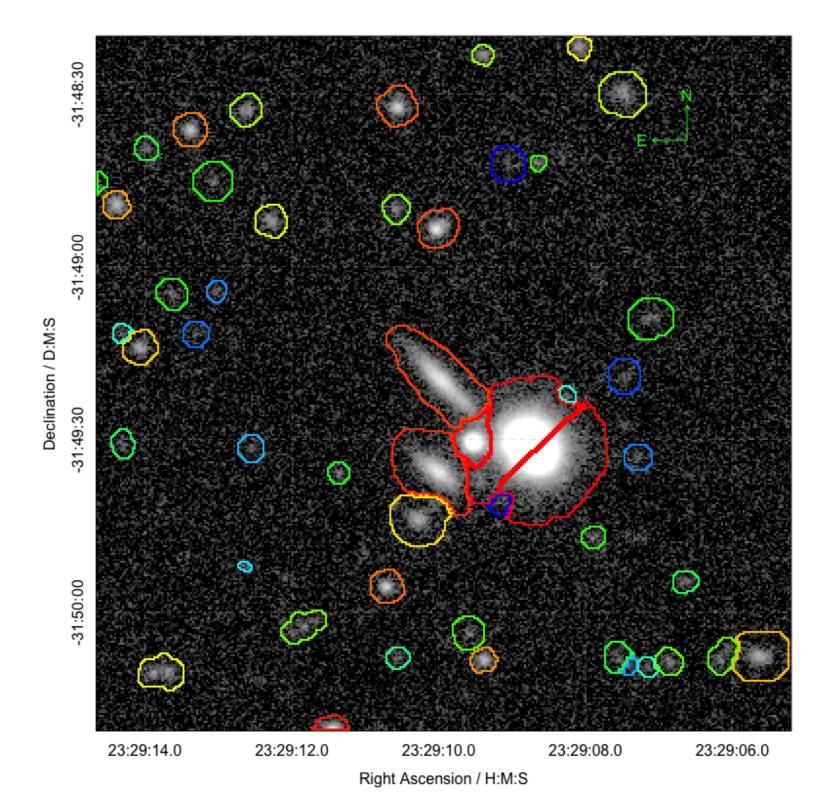
INITIAL VIKING Z-BAND IMAGE



BRIGHT SEGMENTS IDENTIFIED AND DE-BLENDED



SEGMENTS DILATED UNTIL THE FLUX CONTAINED CONVERGES



$PROFOUND \rightarrow PROFIT$

- Here we have focussed on the general photometry aspects of ProFound, but it is at least equally focussed on creating good inputs for ProFit.
- It make everything we need:
 - Source identification (star/gal separation)
 - Good initial conditions (based on photometric analysis)
 - Segmentation maps (both tight and dilated to capture all flux)
 - Sky (can be calculated a number of ways as appropriate)
 - Sigma maps (using sky variance and image gain)



TAKING PHOTOMETRY FURTHER

PROFIT

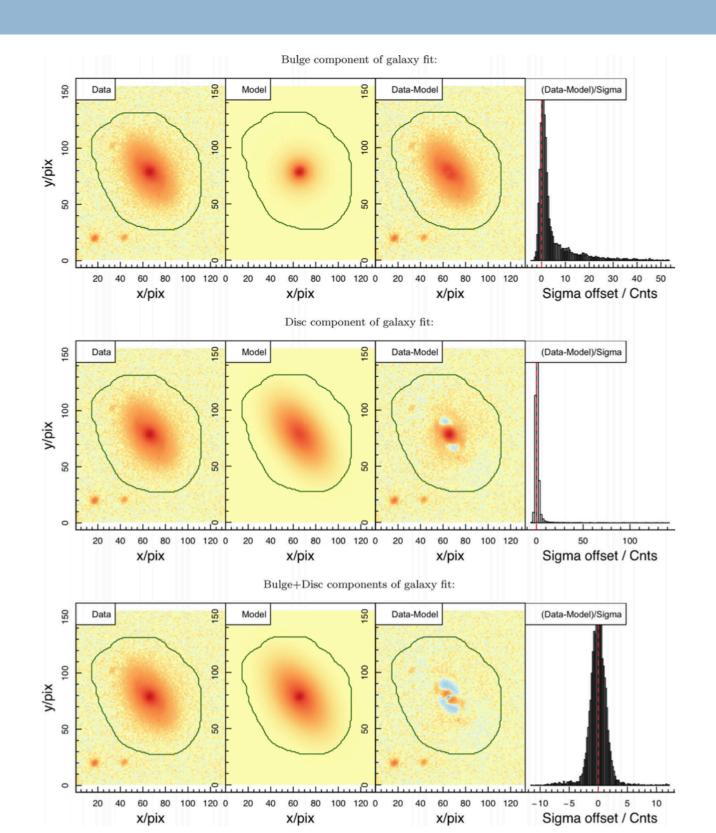
MNRAS **466**, 1513–1541 (2017) Advance Access publication 2016 November 23 doi:10.1093/mnras/stw3039

PROFIT: Bayesian profile fitting of galaxy images

GitHub: ICRAR/ProFit A. S. G. Robotham, ^{1★} D. S. Taranu, ^{1,2} R. Tobar, ¹ A. Moffett ¹ and S. P. Driver ¹

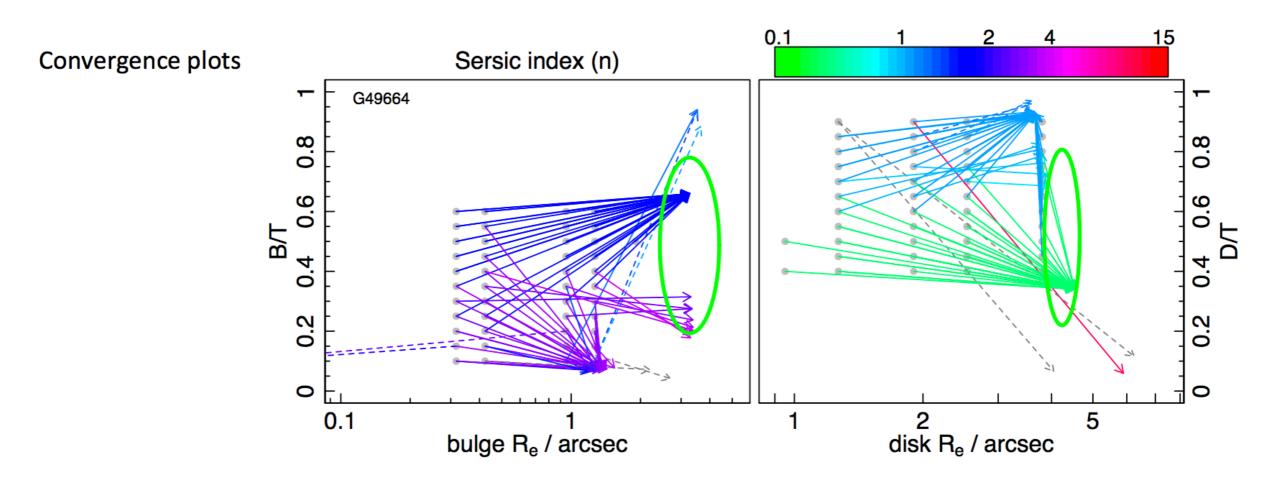
What's the Aim?

- In simple "single band mode" we want to produce a 2D mixture model of a galaxy image.
- This usually separates a galaxy into a compact bulge and and extended disk.



What's the Problem?

Work by ICRAR PhD R. Lange showed that the code of popular community choice (GALFIT) struggles with convergence for many galaxies:



What to do?

- Options were to modify existing code to use more sophisticated samplers, but GALFIT is closed license (binaries only) and other codes available at the time were not battle-tested.
- Vitally for ProFuse (already in the planning in 2017), we need image generation capability, not just 'fitting'.
- Best option all round was to start again and develop a rapid image generation library written in C++ (libprofit) and higher level interfaces to achieve the fitting (ProFit, initially Python and R).
- Initial version was written by me, subsequently R. Tobar (libprofit / pyprofit) and D. Taranu (ProFit) have added a lot of additional features and code re-factors.

What Does libprofit Do For You?

Profiles supported in libprofit:

Sersic: Popular for galaxy components.

Core-Sersic: Popular for early-type galaxies.

Moffat: Popular for stars / PSFs.

Ferrer/s: Popular for galaxy bars.

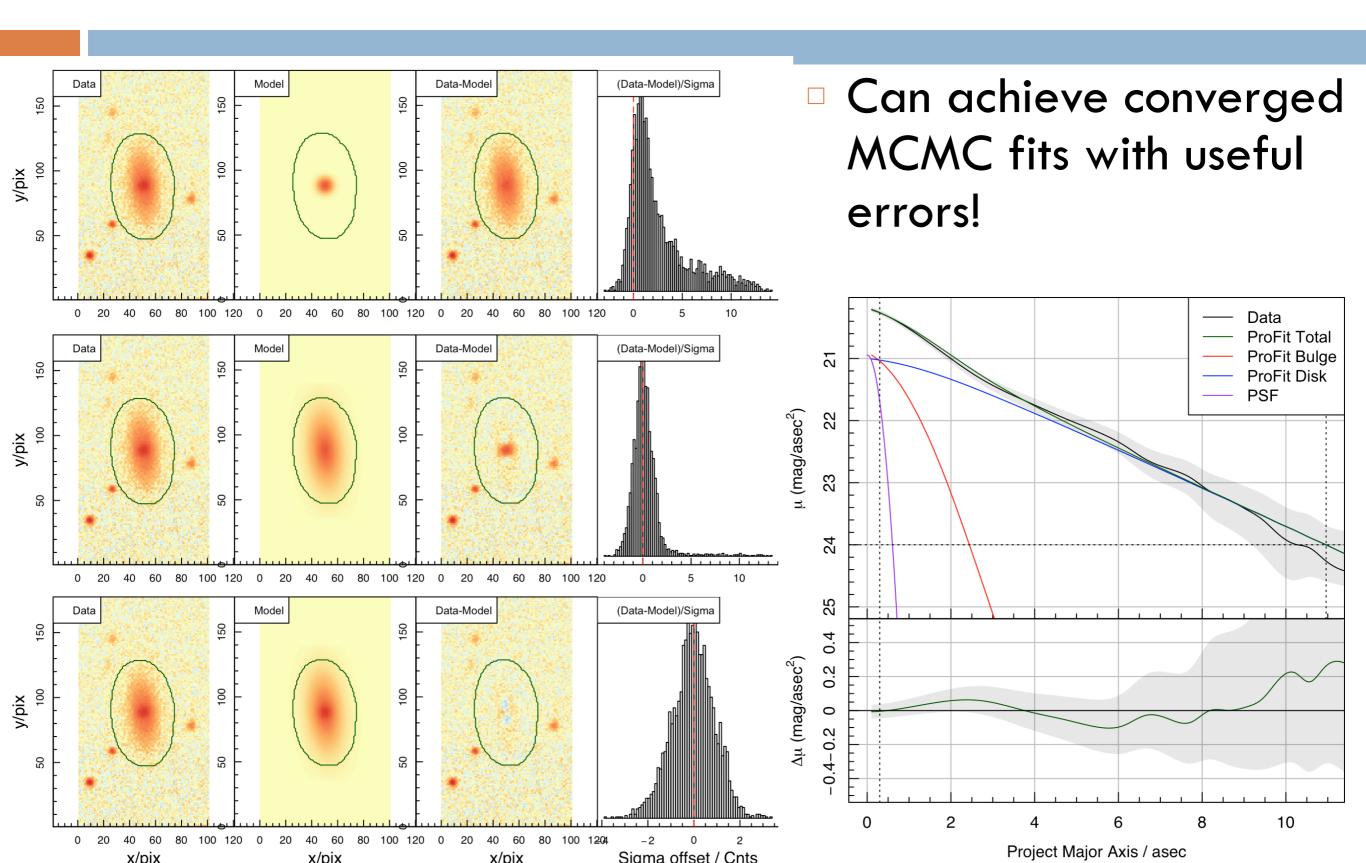
King: Popular for globular clusters.

Empirical PSF: Generic point spread function.

Sky: The sky background.

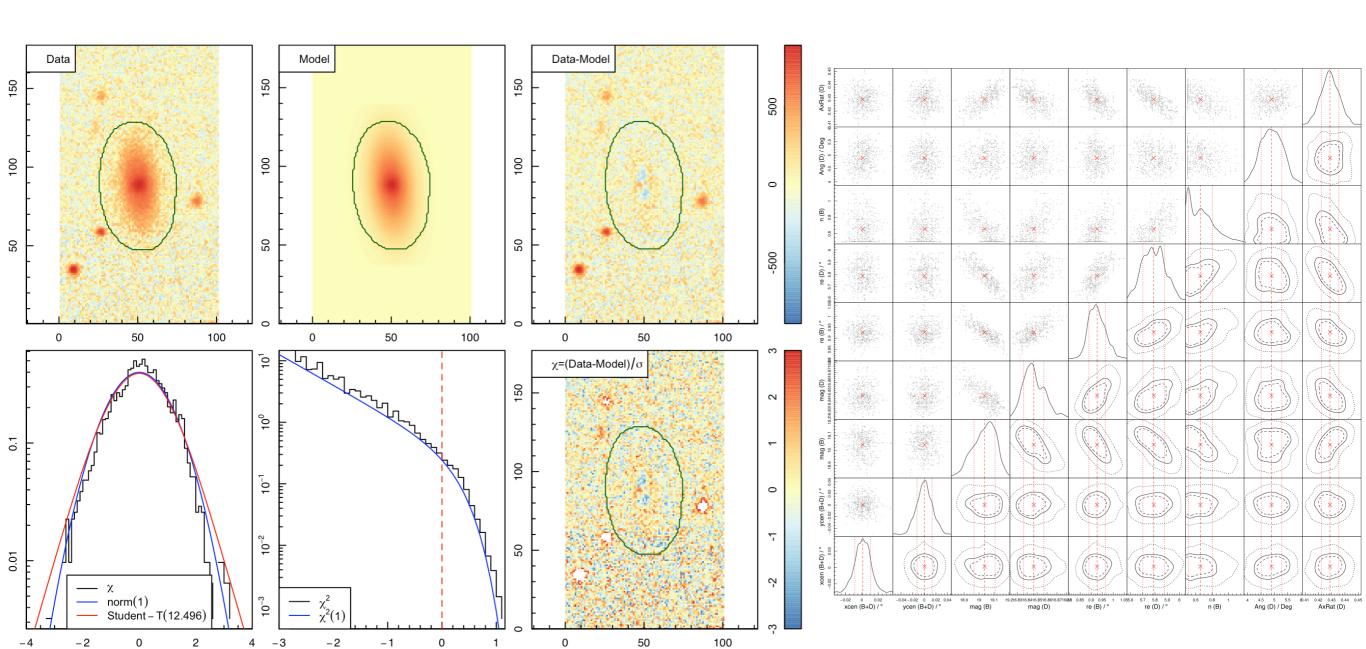
Easy (relatively) to add more profiles: ~30 lines of C++ code to describe the 1D radial profile shape and total flux calculation.

Decomposing Bulges and Disks with ProFit



Decomposing Bulges and Disks with ProFit

Errors are provided via full covariance information.





TAKING SPECTRA TO THE EDGE

PROSPECT

MNRAS **495**, 905–931 (2020) Advance Access publication 2020 April 26 doi:10.1093/mnras/staa1116

PROSPECT: generating spectral energy distributions with complex star formation and metallicity histories

A. S. G. Robotham , , , 2 S. Bellstedt , C. del P. Lagos , J. E. Thorne , L. J. Davies , S. P. Driver and M. Bravo

PROSPECT CODE GOAL

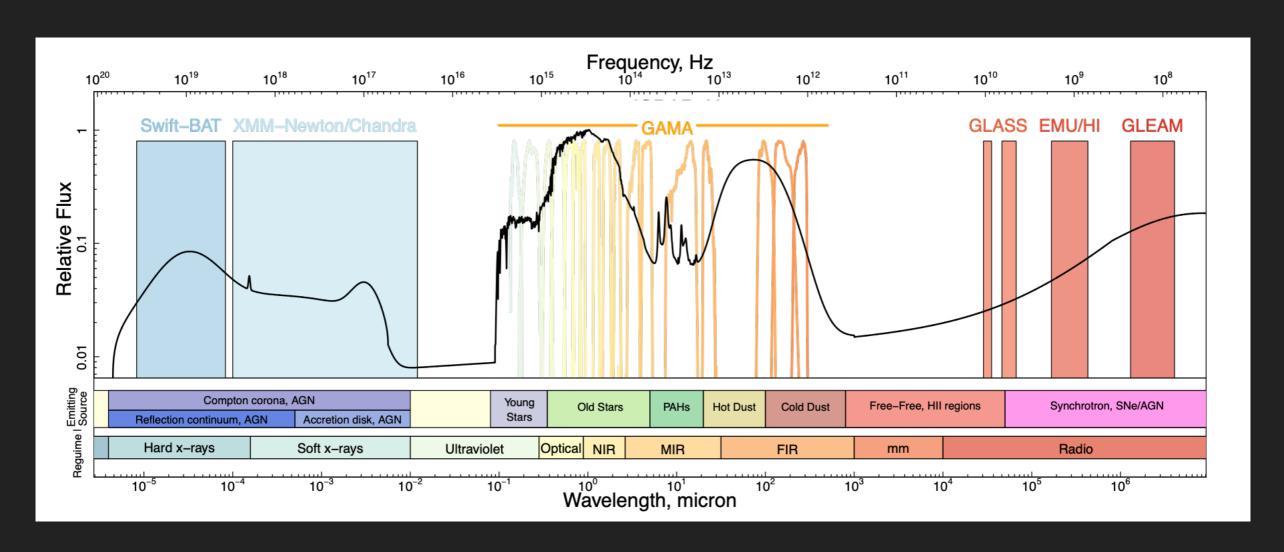
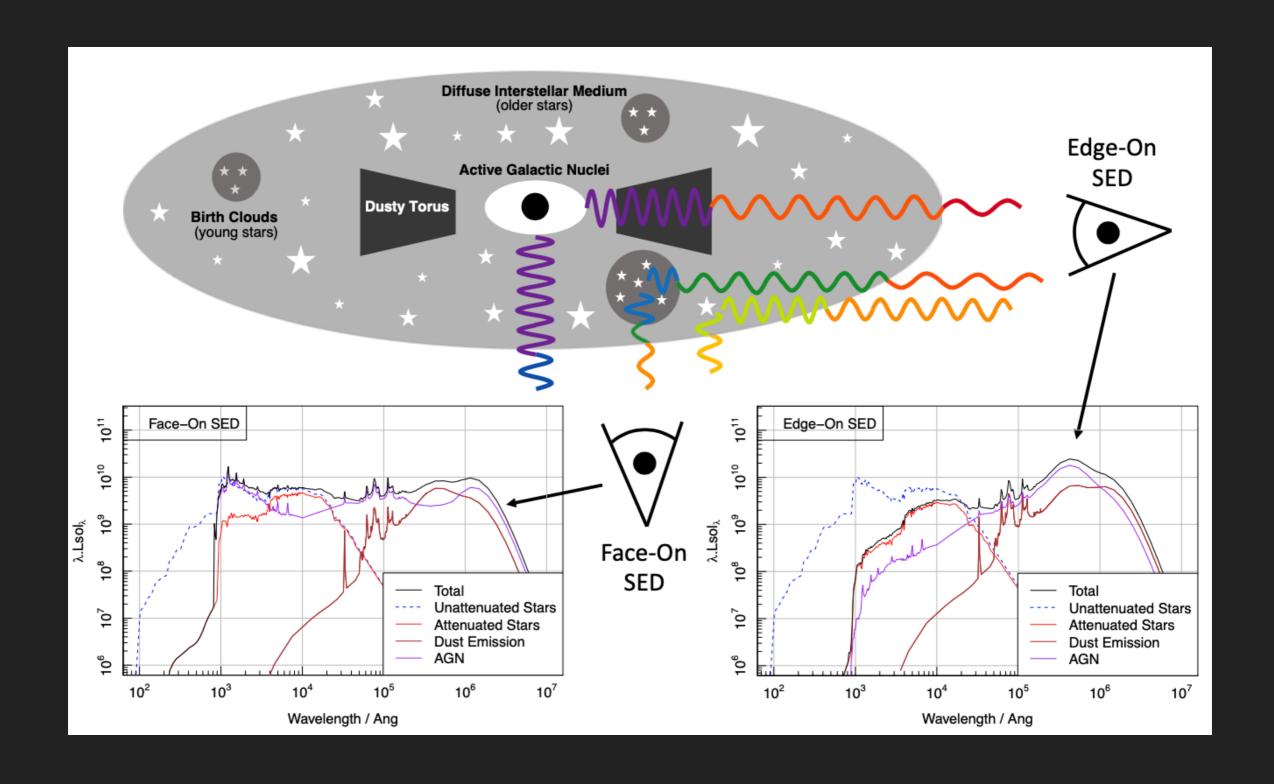
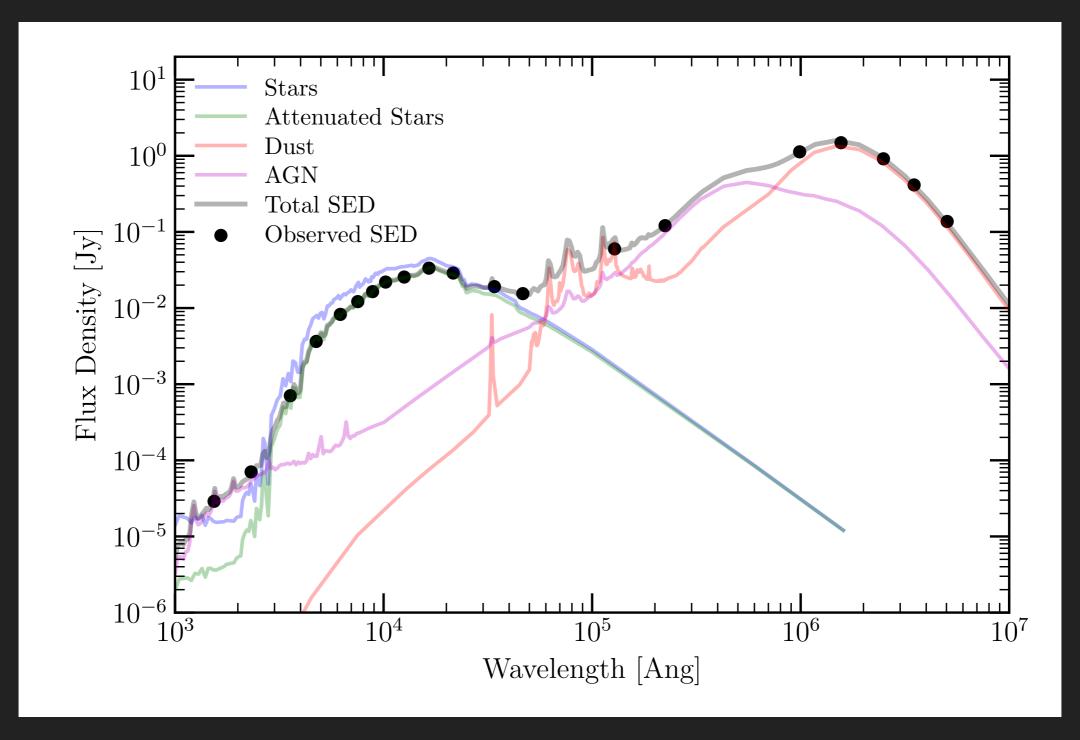


Figure credit: Luke Davies



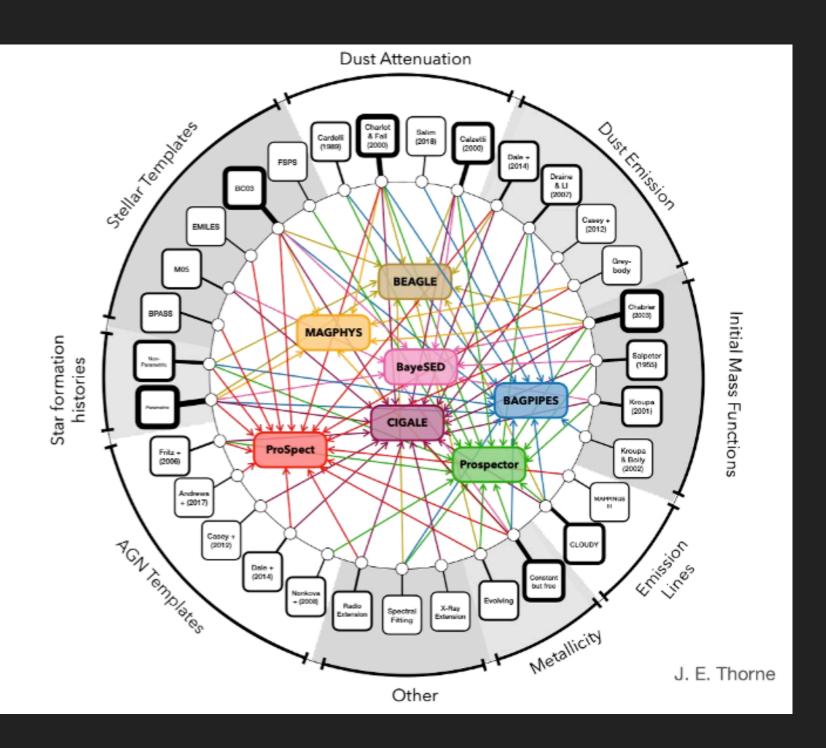
UV - FIR



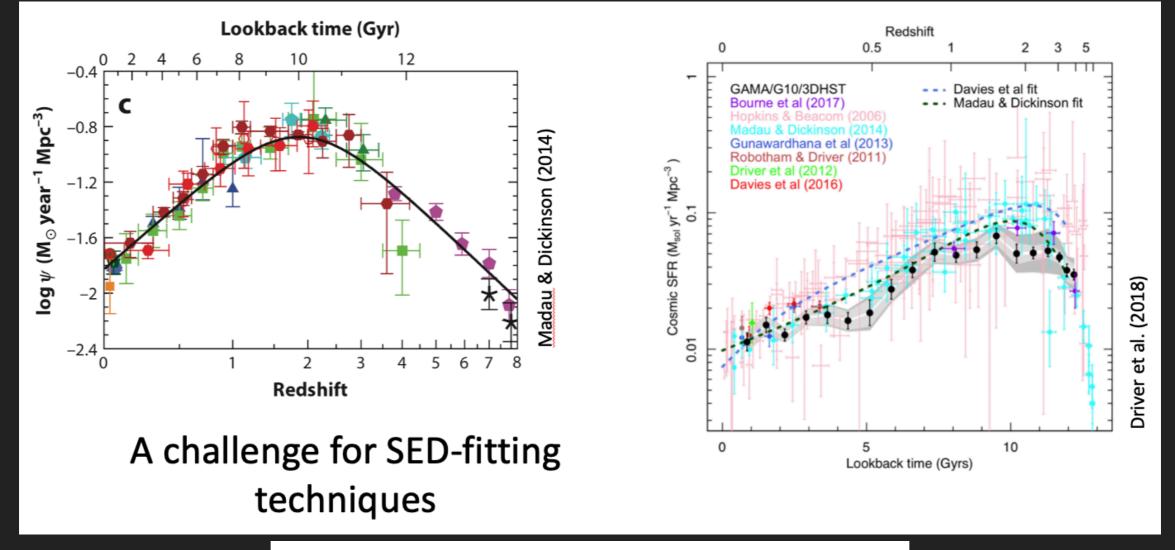
PROSPECT VERSUS

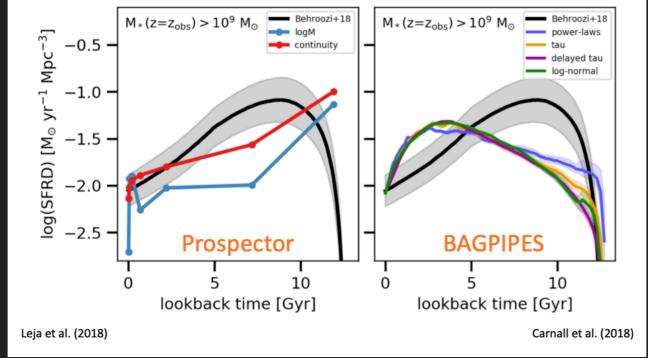
Many SED-fitting tools exist.

They all do roughly the same thing (model an SED via its stellar and dust components), but all in slightly different ways...

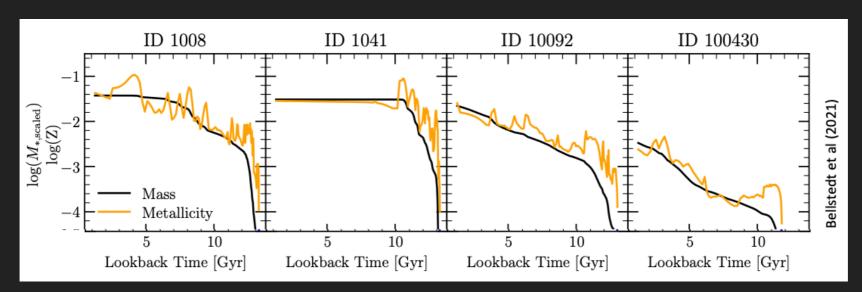


Again, with ProFuse in mind we do not just want an SED 'fitting' code, we need a generative model.

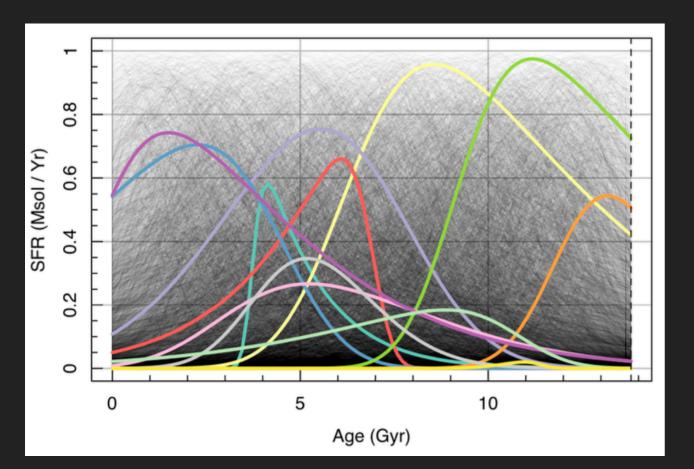


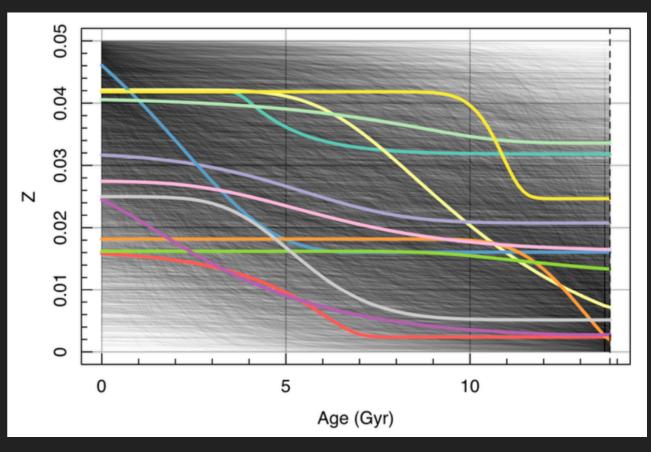


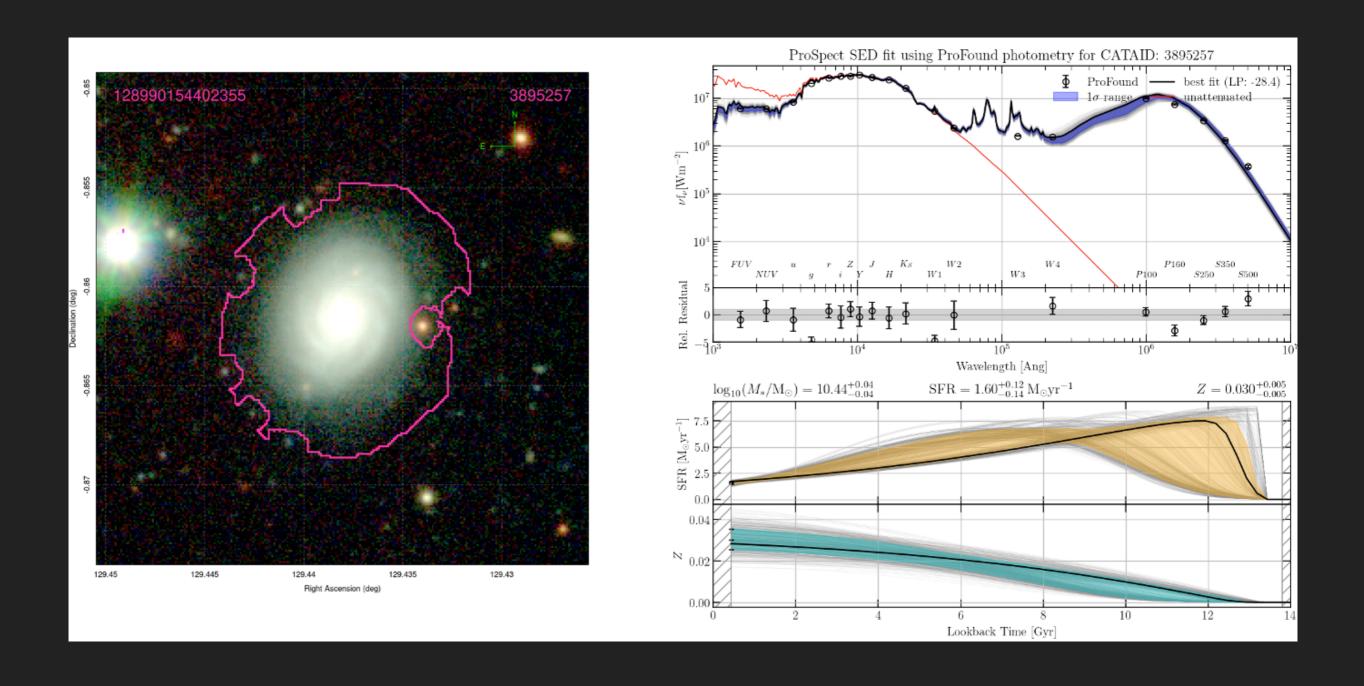
PROSPECT SFH AND ZH



Allows almost any, but we focus on skewed Normal SFH, and linearly mapped ZH:

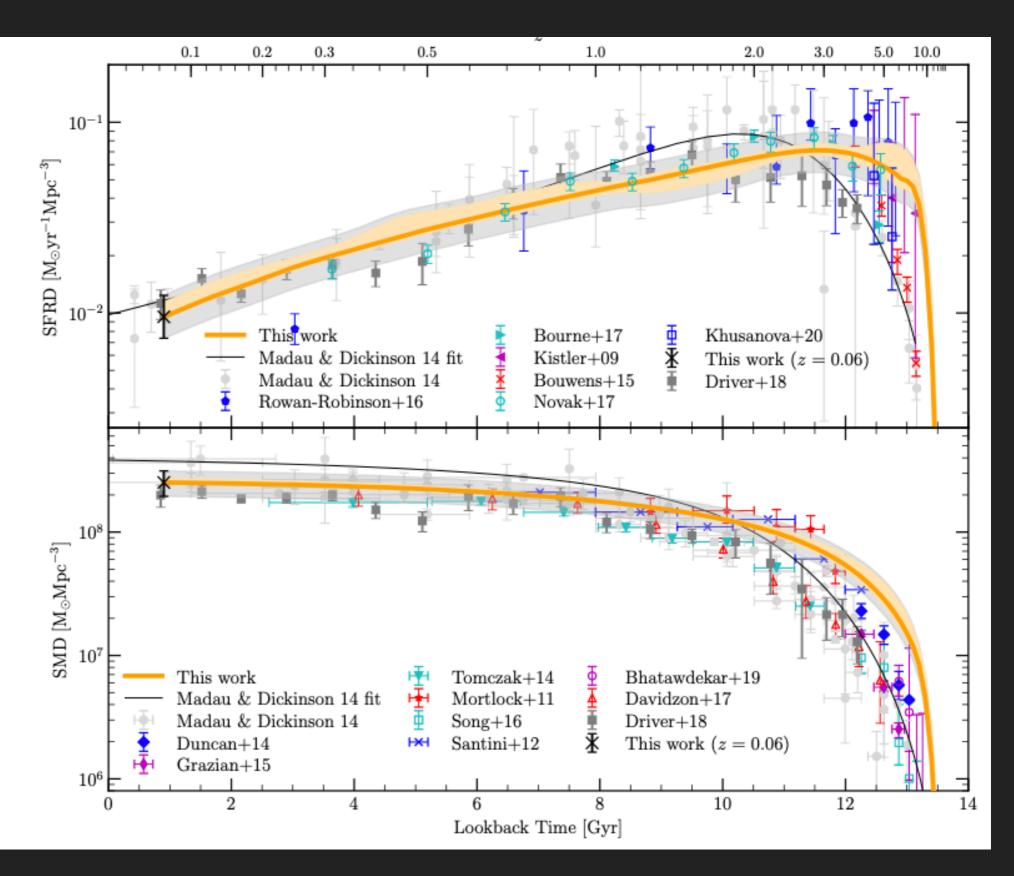


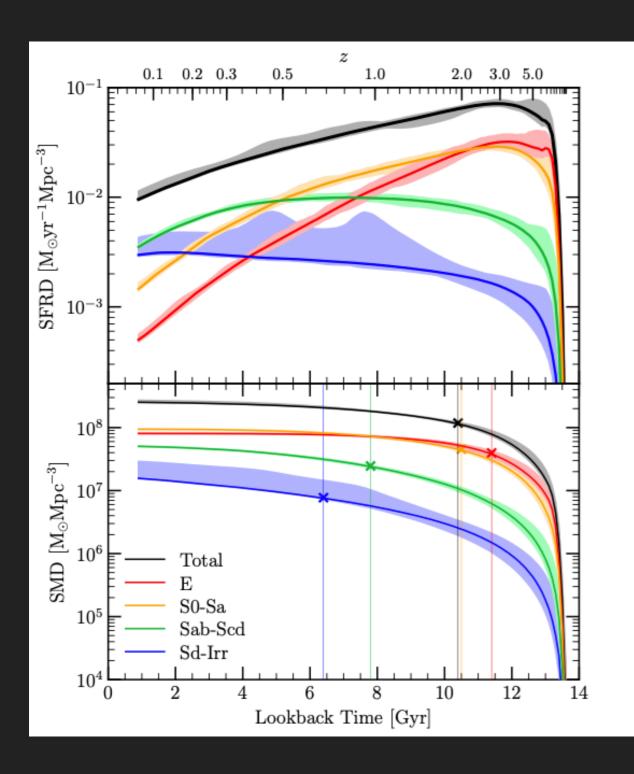




cosmic star formation history

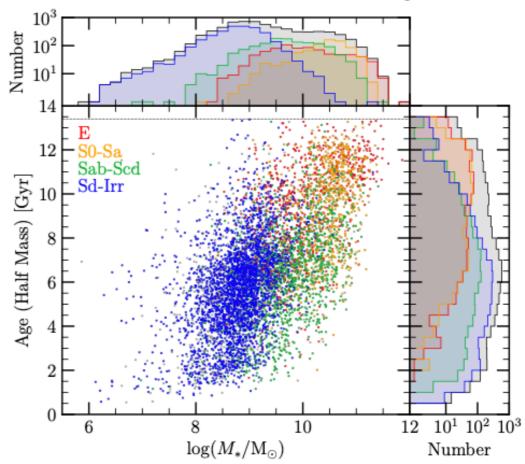
cosmic stellar mass build-up

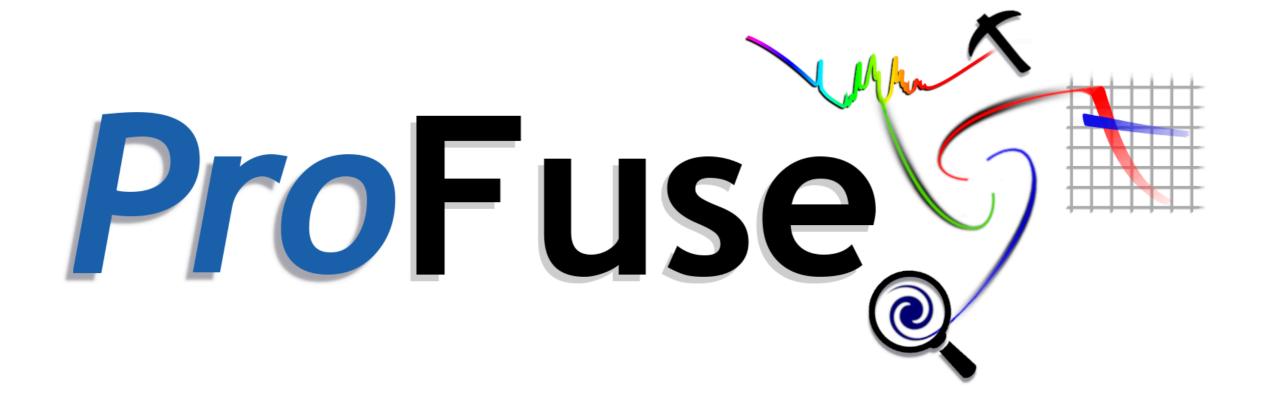




By morphology:

Early-type galaxies are older, whereas late-type galaxies were formed more recently.





PUTTING THIS ALTOGETHER

PROFUSE

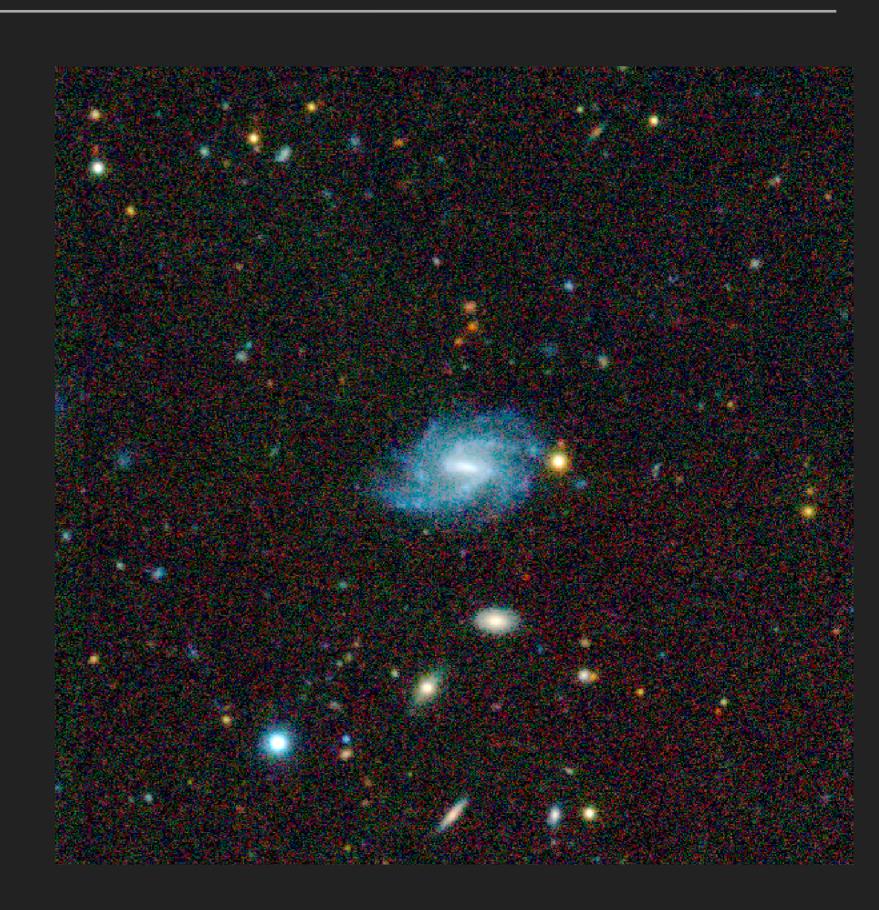
MNRAS **513**, 2985–3012 (2022) Advance Access publication 2022 April 15 https://doi.org/10.1093/mnras/stac1032

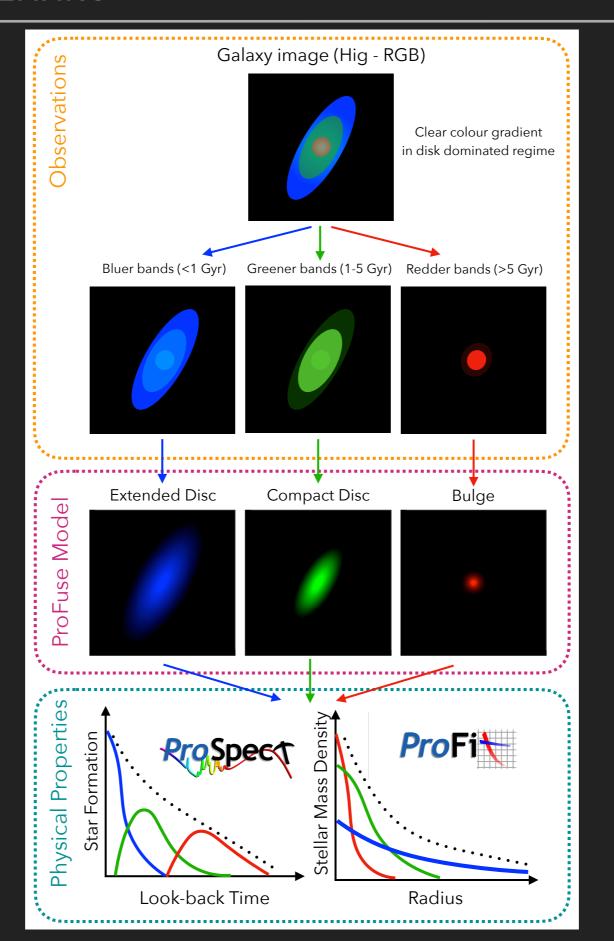
ProFuse: physical multiband structural decomposition of galaxies and the mass-size-age plane

GitHub: asgr/ProFuse

A. S. G. Robotham [®], ^{1,2★} S. Bellstedt ^{®1} and S. P. Driver ^{®1}

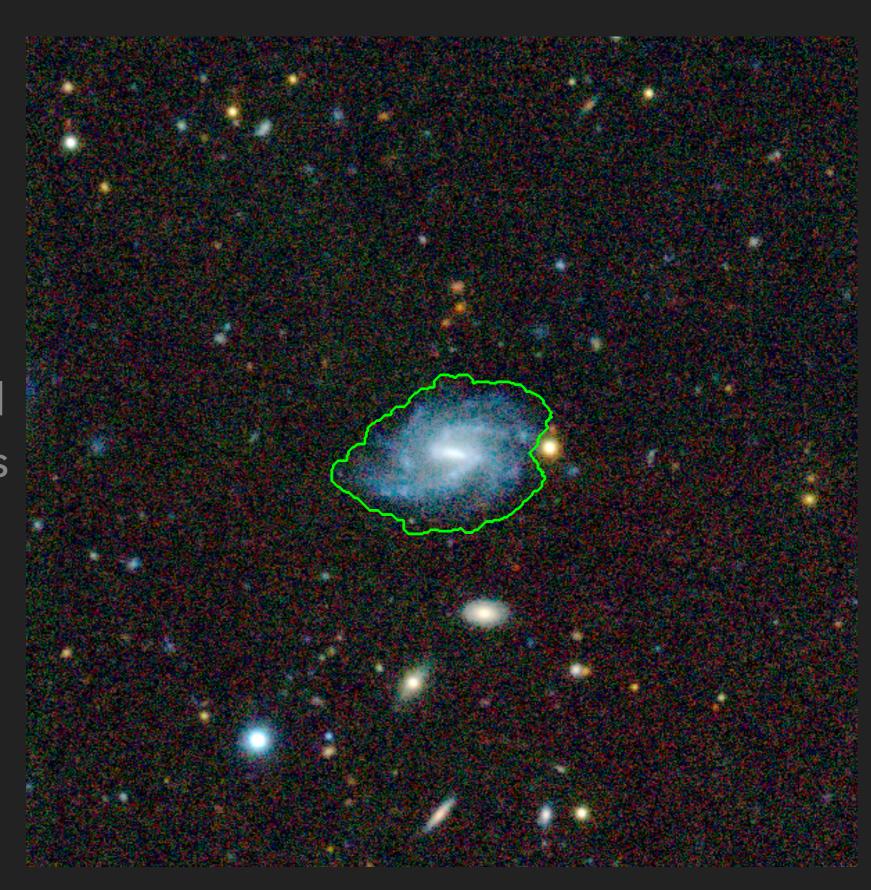
- Starting from just multi-band image date we want to:
 - Detect object
 - Estimate sky
 - Find PSFs
 - Fit bulge + disk
 model with
 distinct SED for
 each component

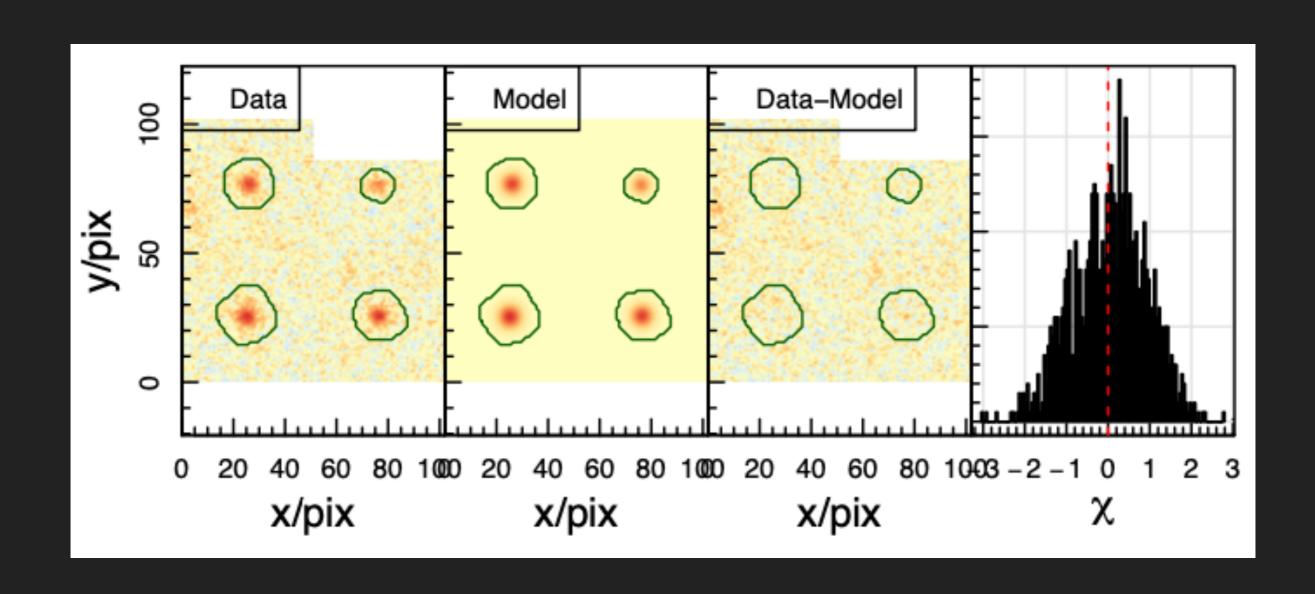




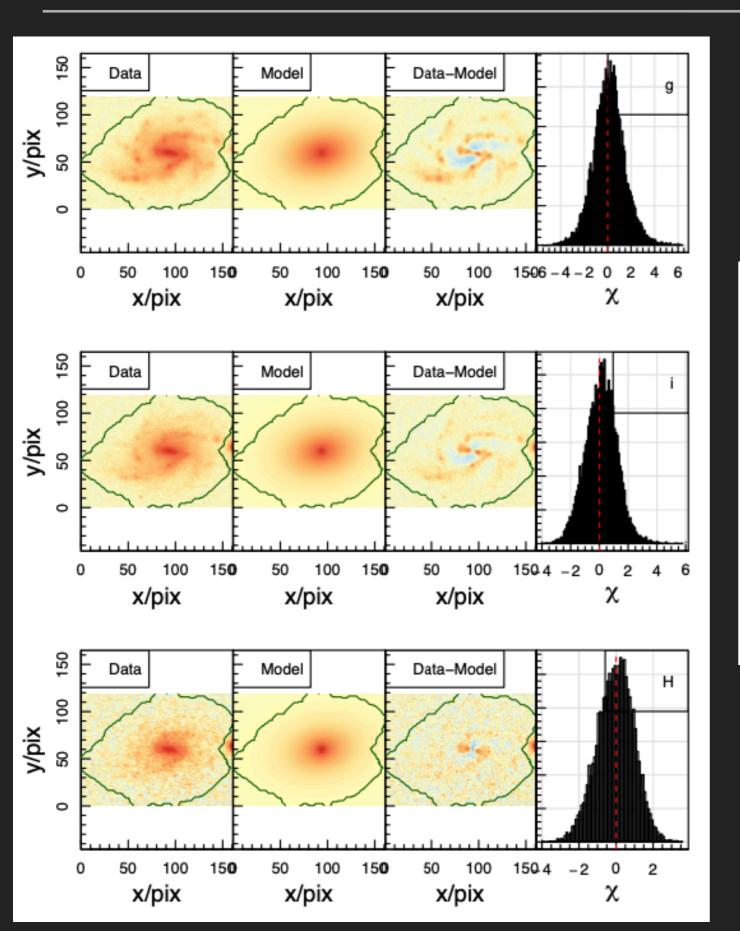
PROFUSE IN ANGER - IMAGE SEGMENTATION AND STAR DETECTION (PROFOUND)

- Objected are roughly detected in all bands, and sky / sky-RMS estimated.
- An inverse variance stack is created, and a global detection is carried out (deep).
- Good candidate
 stars are detected
 per band, and
 Moffat PSFs made.

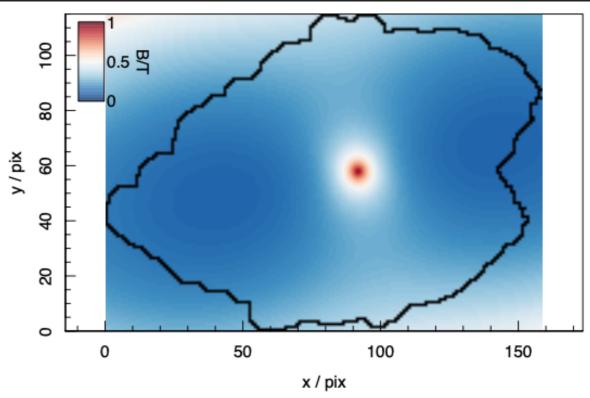




PROFUSE IN ANGER - SPECTRAL + SPATIAL DECOMPOSITION (PROSPECT + PROFIT)



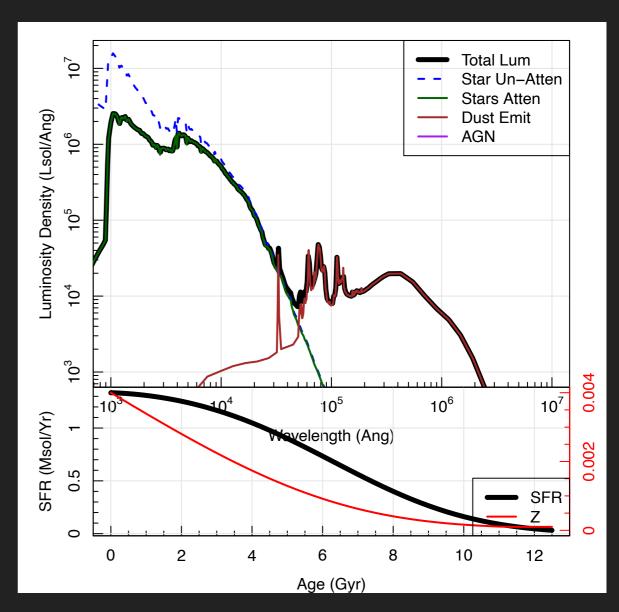
B/T Map





Total Lum Star Un-Atten Stars Atten Luminosity Density (Lsol/Ang) Dust Emit **AGN** 104 4 10² 0.03 10⁵ 10^{4} SFR (Msol/Yr) Wavelength (Ang) 0 8 2 10 12 Age (Gyr)

Disk

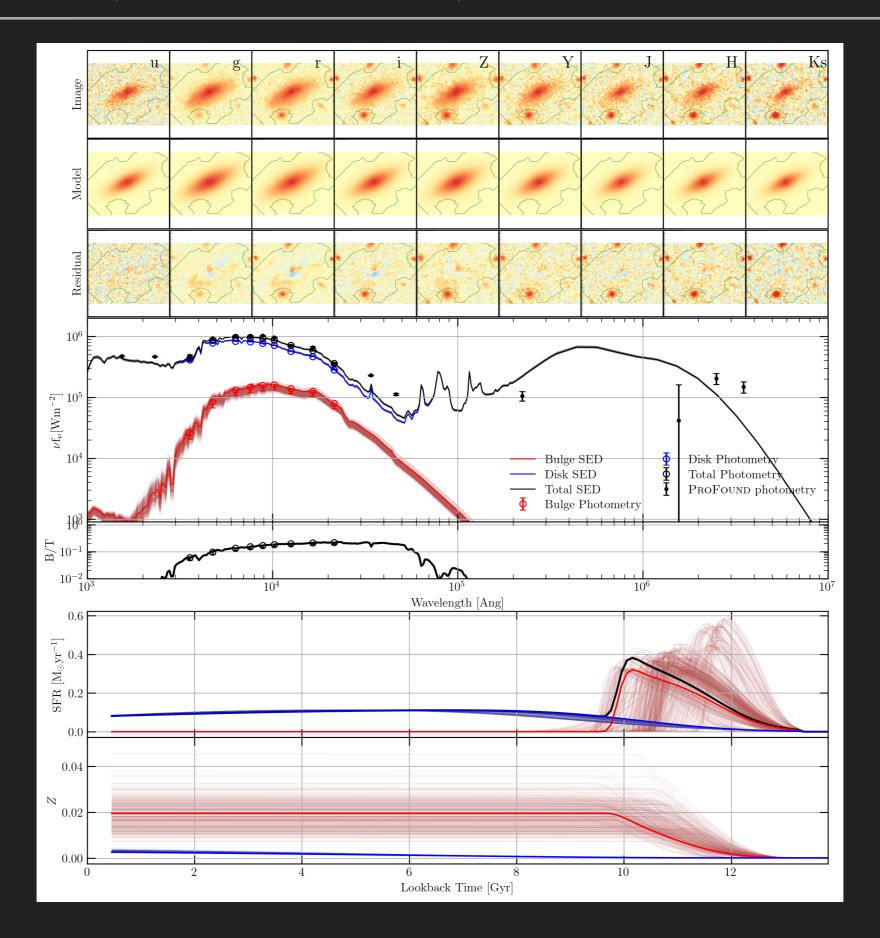


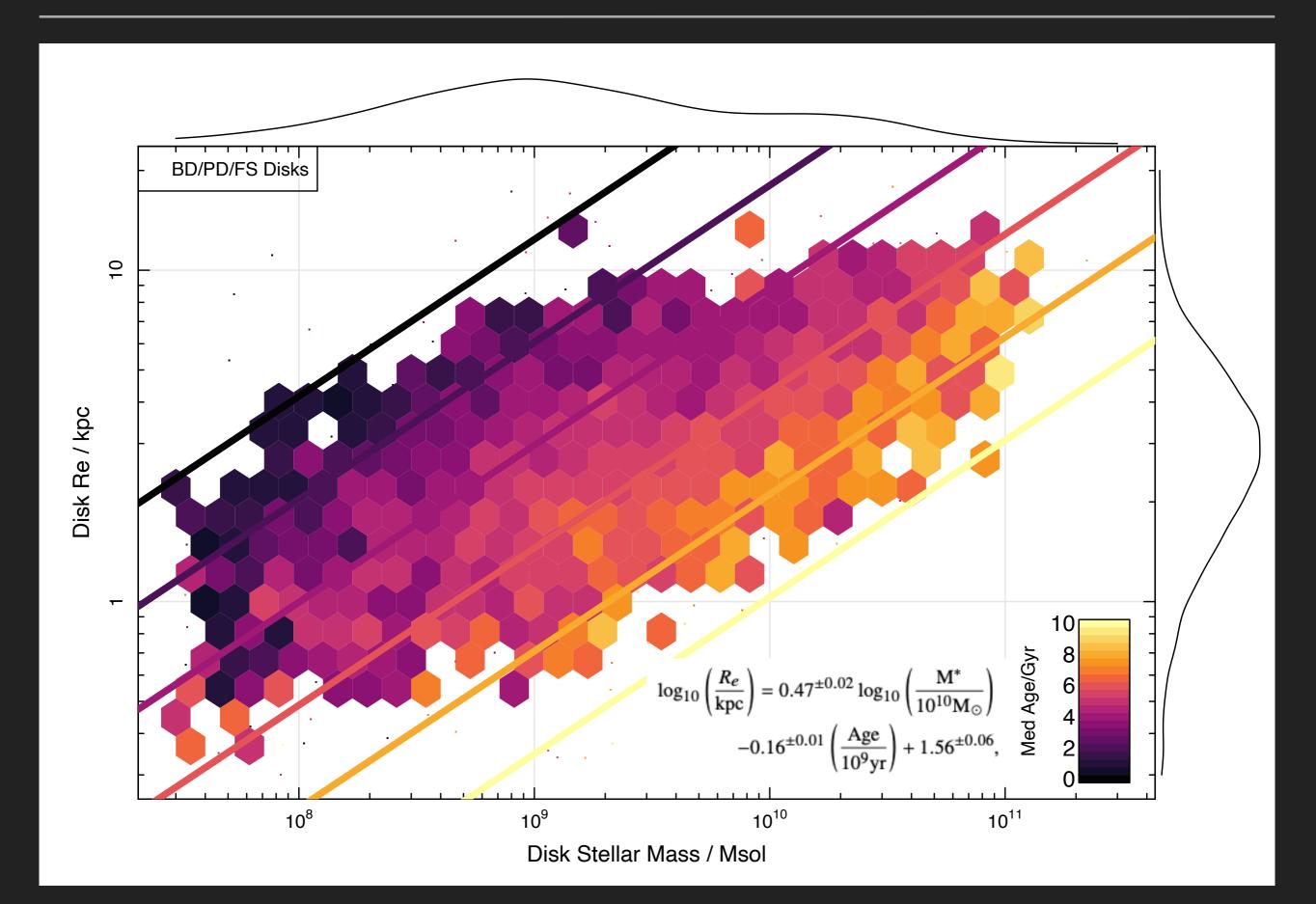
- We get all the usual outputs of ProSpect (SM, SFR, SFH, ZH, and lots more), but per component rather than for the ensemble galaxy.
- ▶ This allows us to generate multi-epoch stellar populations that are not possible in ProSpect with a single skewed Normal SFH.

WITH GREAT POWER...

- In this first application we restrict ourselves to relatively simple models:
 - Free Sersic [FS] a single profile that better suits single component systems with a dominant SFH/ZH (elliptical).
 - Bulge (n=4) + Disk (n=1) [BD] a bulge-disk fit that allows for an extended bulge.
 - Bulge (PSF) + Disk (n=1) [PD] a bulge-disk fit that allows for a PSF (unresolved) bulge.
 - Bulge (n=4) + Disk (n=1) + Disk (n=1) [BDD] a bulge-double-disk fit that allows for an extended bulge and two disk components. These interact to create SFH/ZH and colour gradients.

PROFUSE ON GAMA (~7K GALAXIES, Z<0.06) - ROUGHLY 2-4 HOURS PER FIT





PROFUSE LIVES!

- ProFuse combines all of the ProTools developed to date to provide a fully automated Bayesian Inference engine for the spectral-spatial decomposition of galaxies:
 - ProFound detects sources, segments the galaxy of interest and identifies stars.
 - Stars are then modelled with ProFuse in single-band mode to provide per-band PSFs.
 - ProSpect and ProFit are then run in tandem to produce SFHs and ZHs of bulges and disks simultaneously.
 - In principle arbitrarily complex models are allowed, but with the first application with GAMA we only push as far as a Bulge + Disk +Disk (BDD) model. Mostly FS and BD models are preferred.

PROFUSE EXPERIENCE

- Working to build a code as complicated as ProFuse has taken a significant fraction of my time over the last 7 years.
- Whilst it has supported many students (~10) and post docs (~2) along the way, the core task itself was too major and longterm to be supported via traditional academic means, i.e.:
 - PhDs (still only 3 years in Australia)
 - Post docs (usually only 3 years)
 - Grants (max typically ~4) and my current Future Fellowship position is not to develop any tools- the proposal to develop ProFuse was actually rejected in 2018 (too 'risky' and not directly 'scientific' enough).