# Machine Learning for Particle Astrophysics
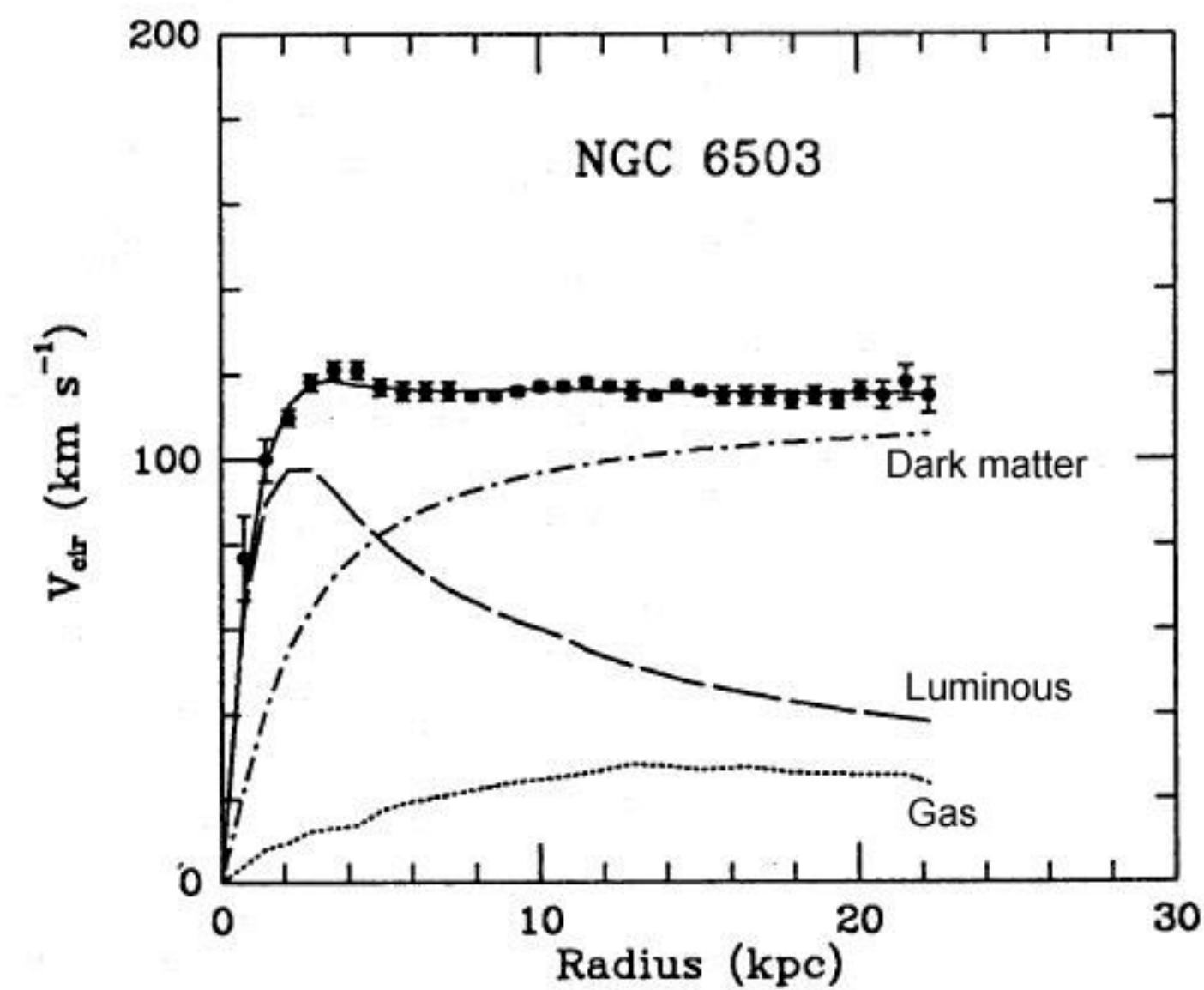
Matthew R Buckley
Rutgers University

- David Shih (Rutgers), Lina Necib (MIT)

- Sung Hak Lim (Rutgers), Claudius Krause (Rutgers/Heidelberg)

- Eric Putney (Rutgers), Anna Hallin (Rutgers/Hamburg), John Tamanas (UCSC)
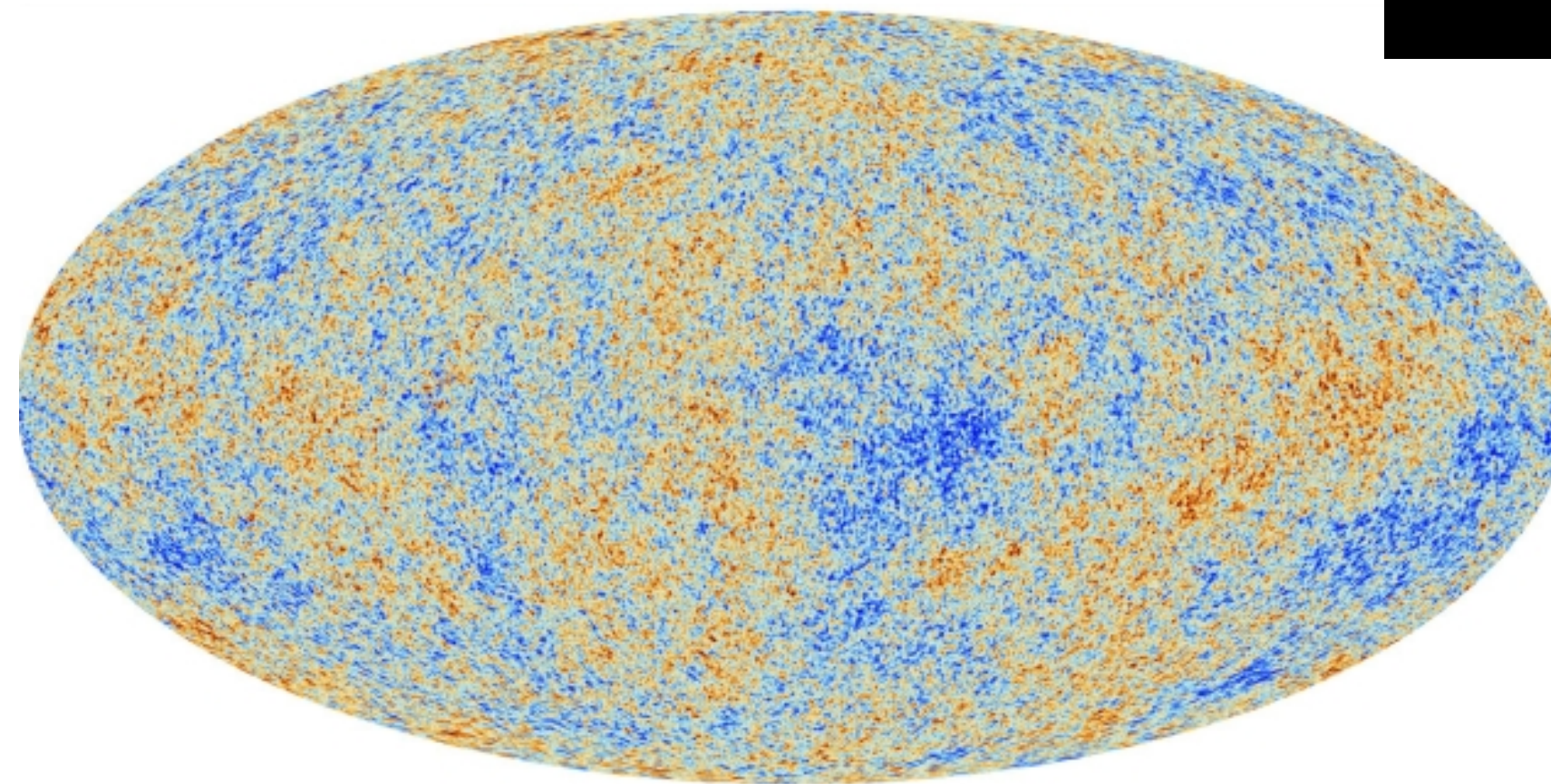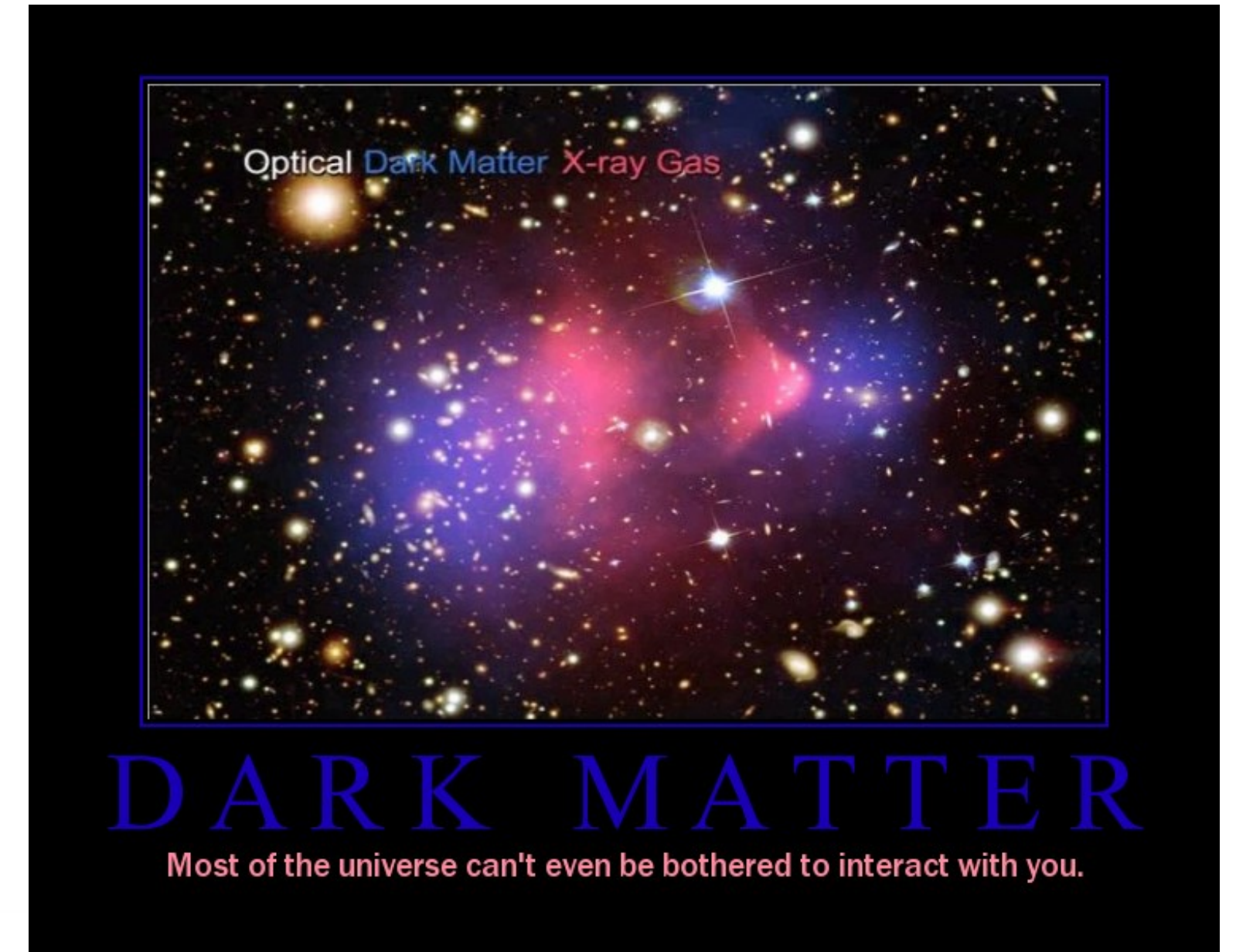
- Kailash Raman (Rutgers)

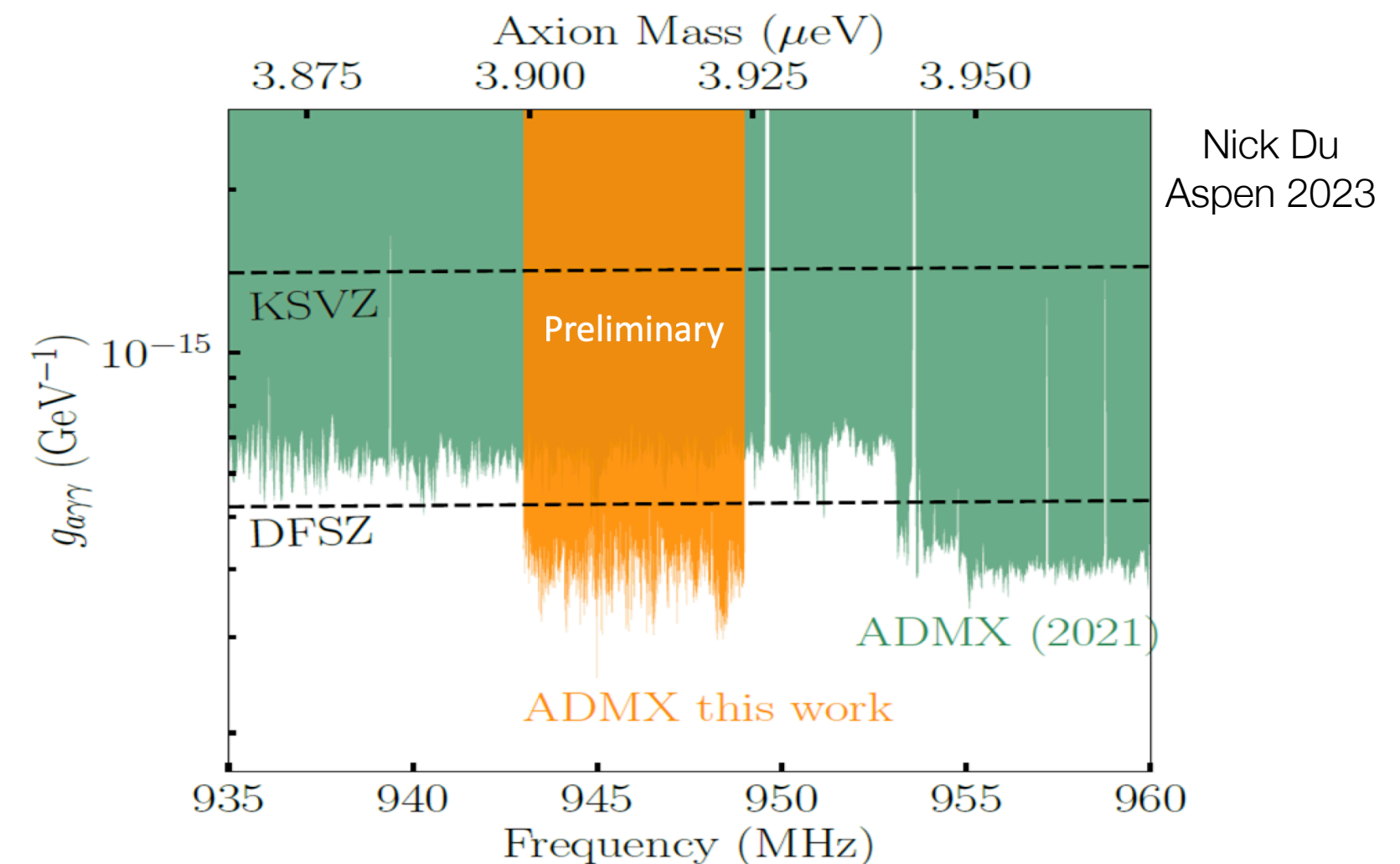- We know dark matter exists, but our evidence is purely astrophysical:



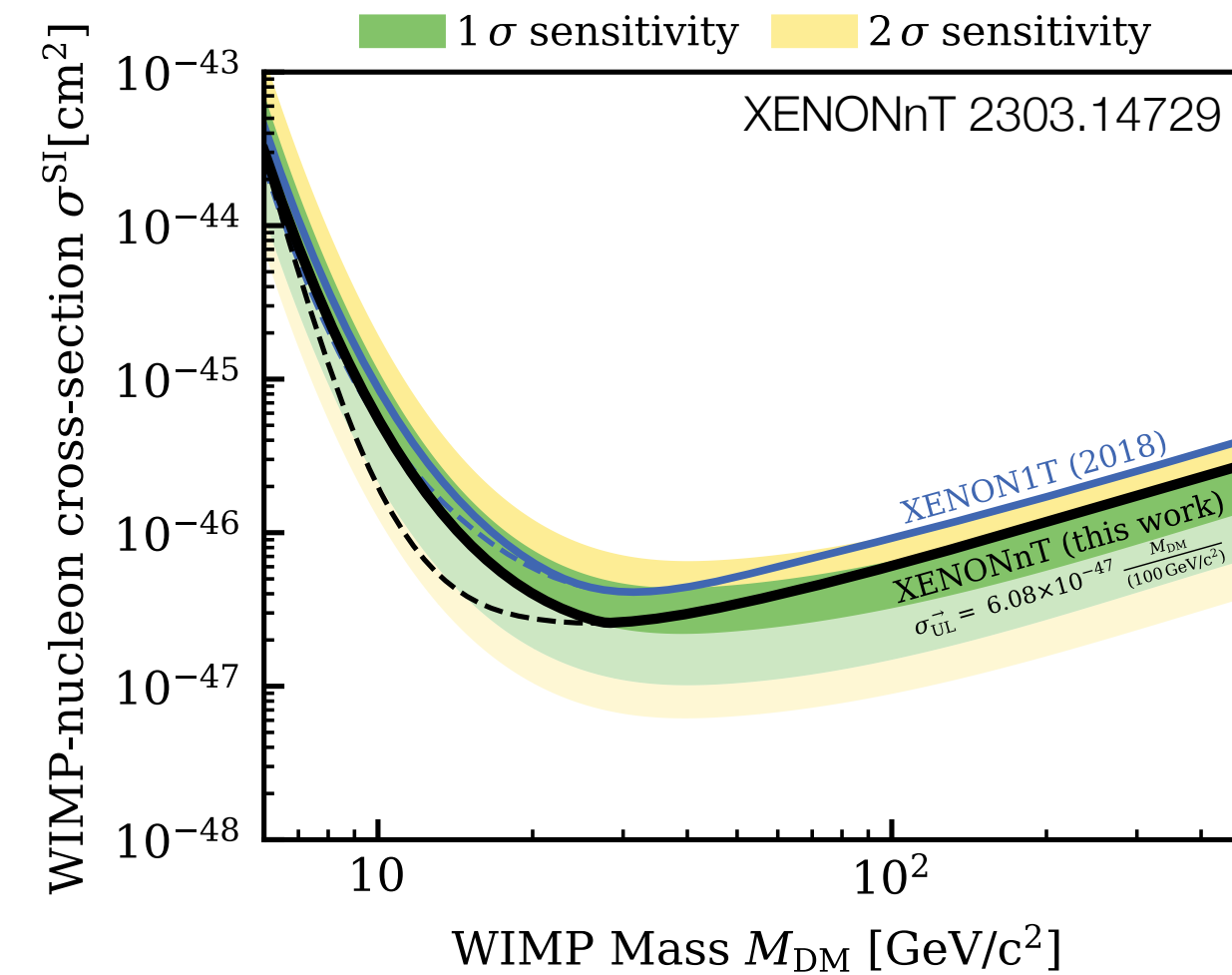K.G. Begeman, A.H. Broels, R.H. Sanders. 1991. Mon.Not.RAS 249, 523.

- Particle Physics experiments are motivated and important, but so far give only negative results

- Large-scale distribution of baryonic matter in the Universe and structure of galaxies can reveal hints of dark matter particle physics.



Buckley *et al* 1405.2075

$m_D = 1$ TeV
$\alpha_D = 0.009$
$B_D = 1$ keV
$\xi_0 = 0.5$

— Self−Interacting DM
--- Warm DM Analogue
--- Silk Damping Envelope
— Cold DM

Illustris Simulation

Sameie *et al* 1904.07872

Draco − like

IC
CDM
$3$ cm$^2$/g
$5$ cm$^2$/g
$10$ cm$^2$/g

Dwarf 1
$M_{200}^{IC} : 2 \times 10^9 M_\odot$
$c_{200}^{IC} : 29.5$
$r_{peri} : 26$ kpc

# Theoretical Motivations

- Large-scale distribution of baryonic matter in the Universe and structure of galaxies can reveal hints of dark matter particle physics.

$M_{\rm halo}$          Probes

**Final frontier** — $10^2 \, M_\odot$
- Gravitational nanolensing (time domain)
- Gravitational waves from compact-object DM (multi-messenger)
- Microlensing of compact-object DM (time domain)
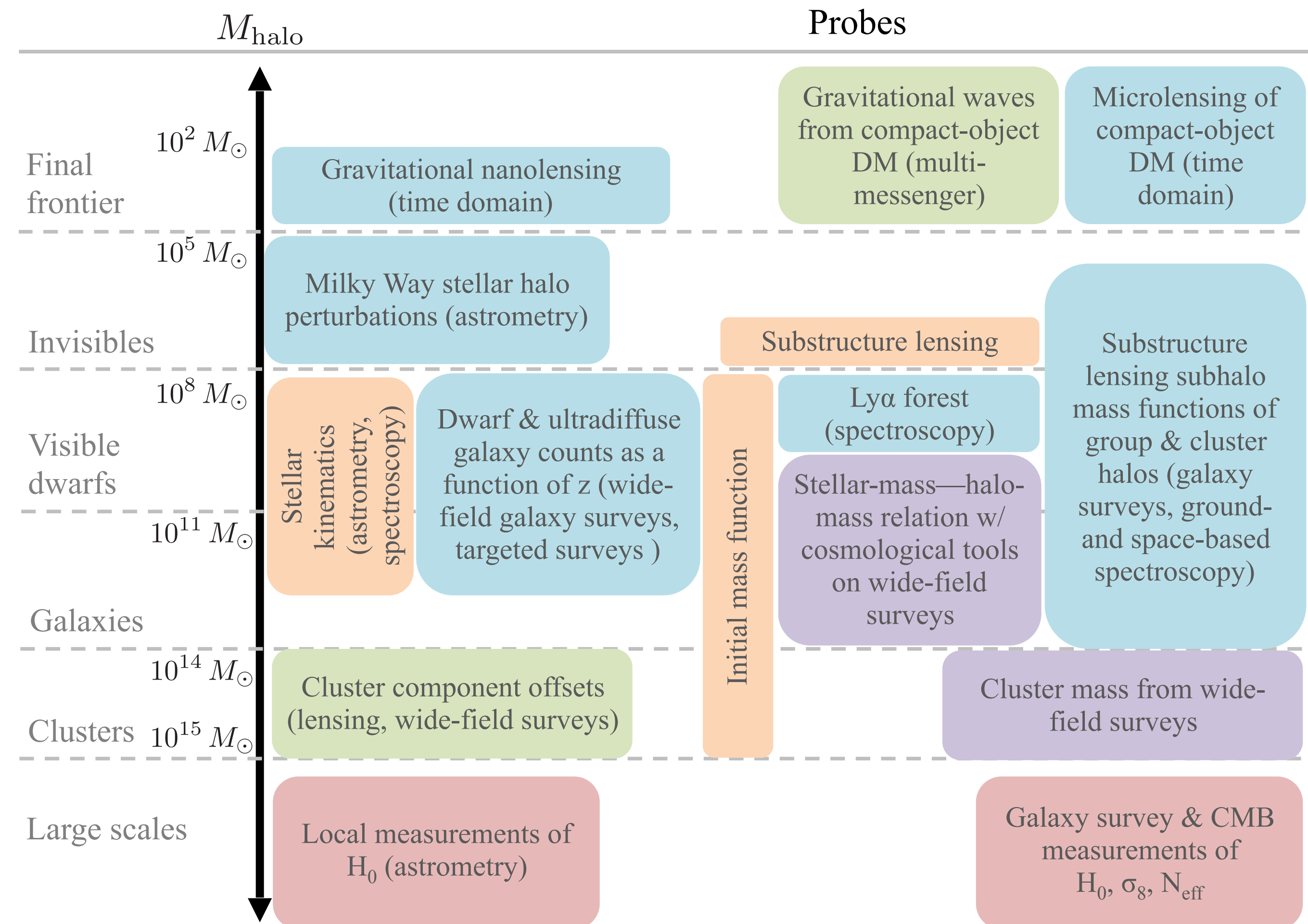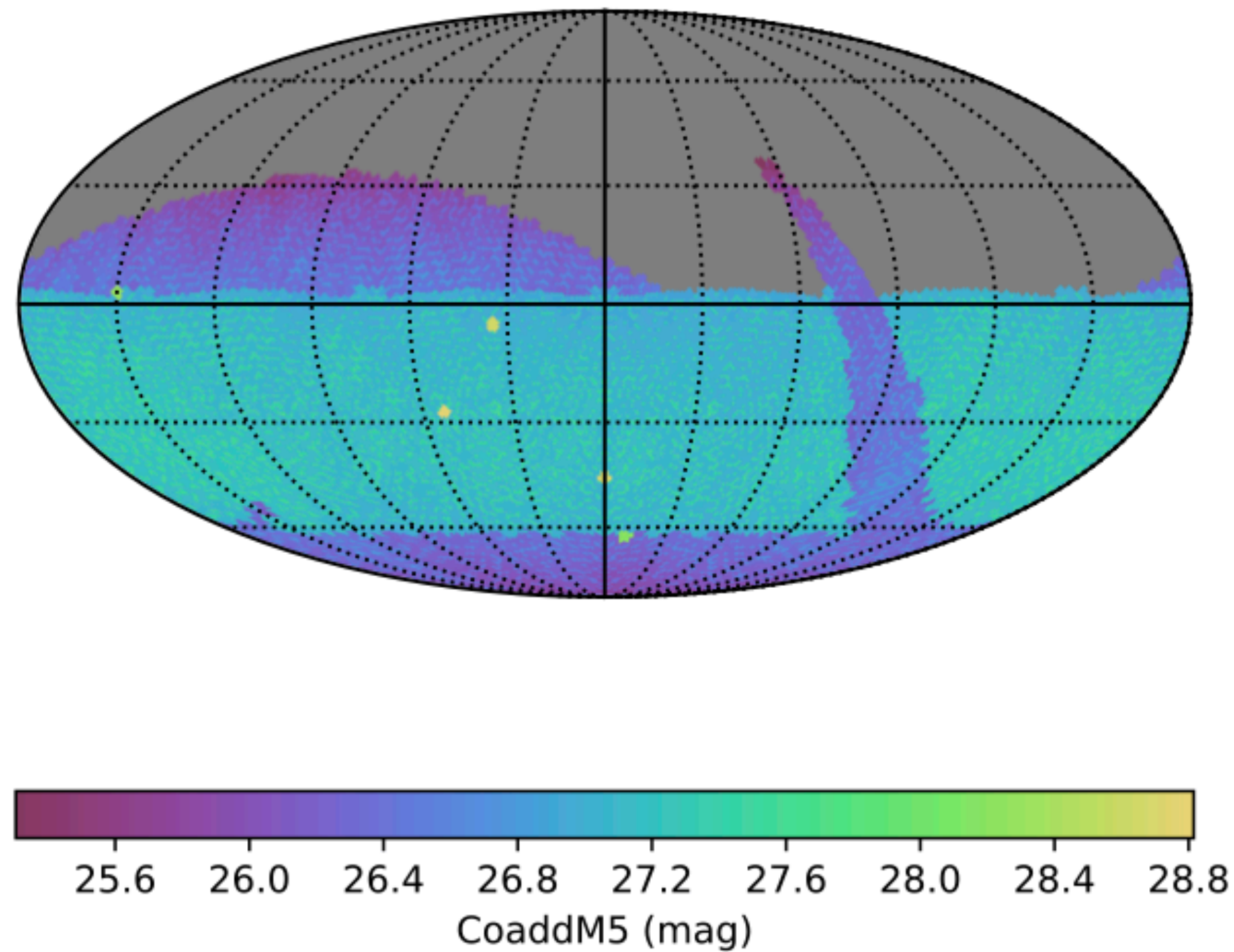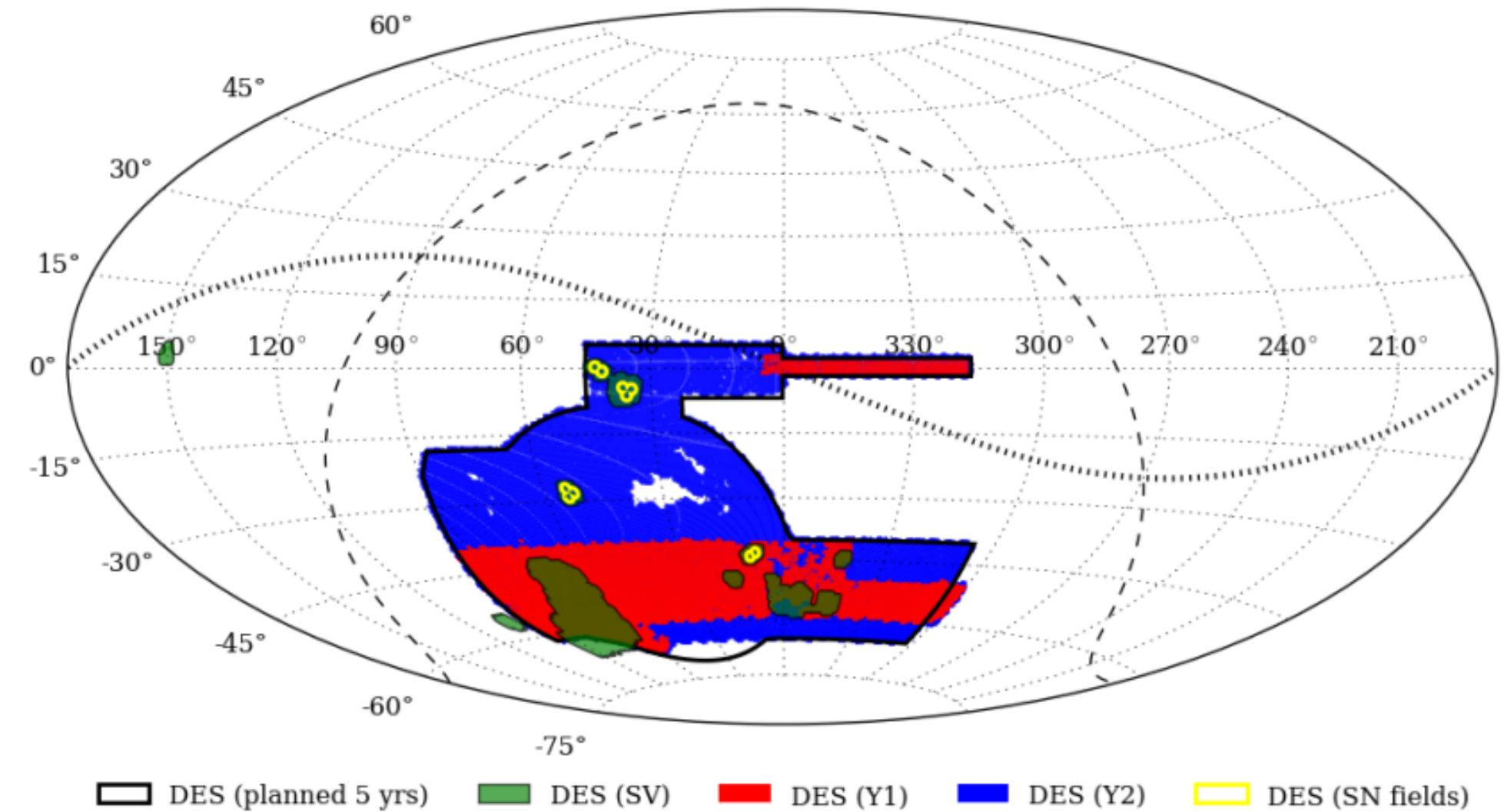
**Invisibles** — $10^5 \, M_\odot$
- Milky Way stellar halo perturbations (astrometry)
- Substructure lensing

**Visible dwarfs** — $10^8 \, M_\odot$
- Stellar kinematics (astrometry, spectroscopy)
- Dwarf & ultradiffuse galaxy counts as a function of z (wide-field galaxy surveys, targeted surveys )
- Lyα forest (spectroscopy)
- Substructure lensing subhalo mass functions of group & cluster halos (galaxy surveys, ground- and space-based spectroscopy)

**Galaxies** — $10^{11} \, M_\odot$
- Initial mass function
- Stellar-mass—halo-mass relation w/ cosmological tools on wide-field surveys

**Clusters** — $10^{14} \, M_\odot$ , $10^{15} \, M_\odot$
- Cluster component offsets (lensing, wide-field surveys)
- Cluster mass from wide-field surveys

**Large scales**
- Local measurements of $H_0$ (astrometry)
- Galaxy survey & CMB measurements of $H_0$, $\sigma_8$, $N_{\rm eff}$

Buckley & Peter 1712.06615

Vera Rubin/LSST

opsim  g: CoaddM5



CoaddM5 (mag)

DES OBSERVING STRATEGY



DES (planned 5 yrs)    DES (SV)    DES (Y1)    DES (Y2)    DES (SN fields)

DESI Legacy

GAIA

- Gaia satellite measures the 3D positions and proper motions of ~1.5 billion stars in the Galaxy.

  - N.B: Gaia measures *parallax*, not *distance.*

  - Provides *photometry* (color and magnitude) and limited *spectroscopy*

  - Line-of-sight motion for ~34 million stars (DR3)

    - This will be ~150 million by end-of-mission

- A huge mine of data for the study of Galactic substructure.

- In this talk, I'm interested in Gaia data as processed locations of stars within 4/5/6D kinematic space — not as individual images/spectra (lots of analysis here!)
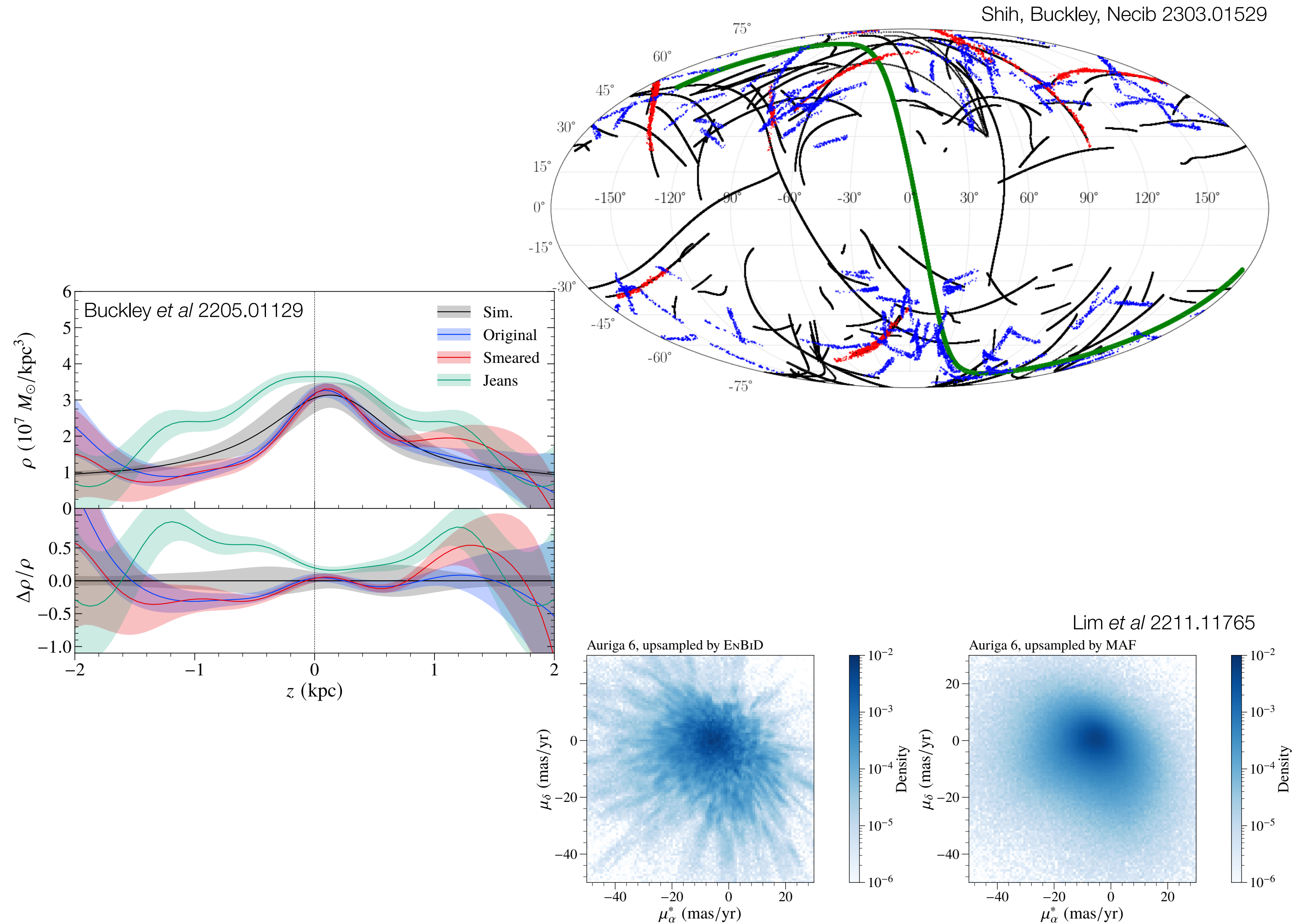
| | # sources in Gaia DR3 | # sources in Gaia DR2 | # sources in Gaia DR1 |
|---|---|---|---|
| **Total number of sources** | **1,811,709,771** | **1,692,919,135** | **1,142,679,769** |
| | Gaia Early Data Release 3 | | |
| Number of sources with full astrometry | 1,467,744,818 | 1,331,909,727 | 2,057,050 |
| Number of 5-parameter sources | 585,416,709 | | |
| Number of 6-parameter sources | 882,328,109 | | |
| Number of 2-parameter sources | 343,964,953 | 361,009,408 | 1,140,622,719 |
| Gaia-CRF sources | 1,614,173 | 556,869 | 2191 |
| Sources with mean G magnitude | 1,806,254,432 | 1,692,919,135 | 1,142,679,769 |
| Sources with mean $G_{BP}$-band photometry | 1,542,033,472 | 1,381,964,755 | - |
| Sources with mean $G_{RP}$-band photometry | 1,554,997,939 | 1,383,551,713 | - |
| | New in Gaia Data Release 3 | Gaia DR2 | Gaia DR1 |
| Sources with radial velocities | 33,812,183 | 7,224,631 | - |
| Sources with mean $G_{RVS}$-band magnitudes | 32,232,187 | - | - |
| Sources with rotational velocities | 3,524,677 | - | - |
| Mean BP/RP spectra | 219,197,643 | - | - |
| Mean RVS spectra | 999,645 | - | - |

- The Milky Way's Mass Density

- Stellar Streams
  - Via Machinae (ANODE)
  - CATHODE

- Synthetic *Gaia* Observations

Shih, Buckley, Necib 2303.01529

Buckley *et al* 2205.01129

Lim *et al* 2211.11765
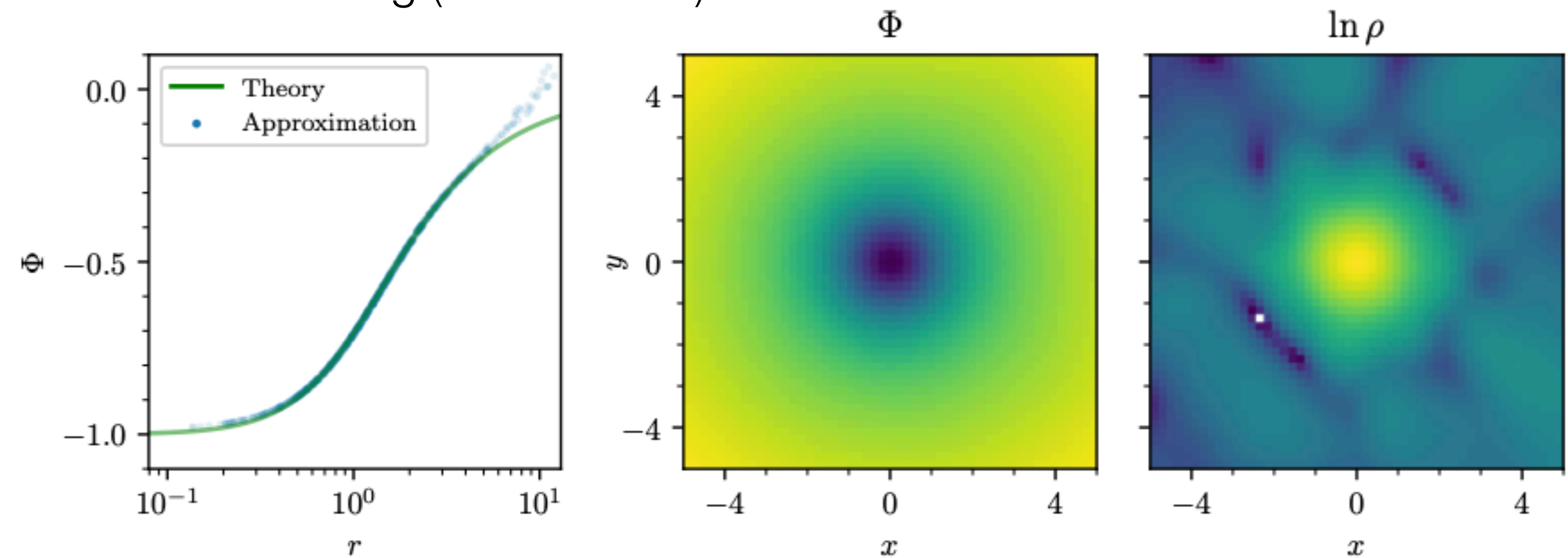
Auriga 6, upsampled by ENBID

Auriga 6, upsampled by MAF

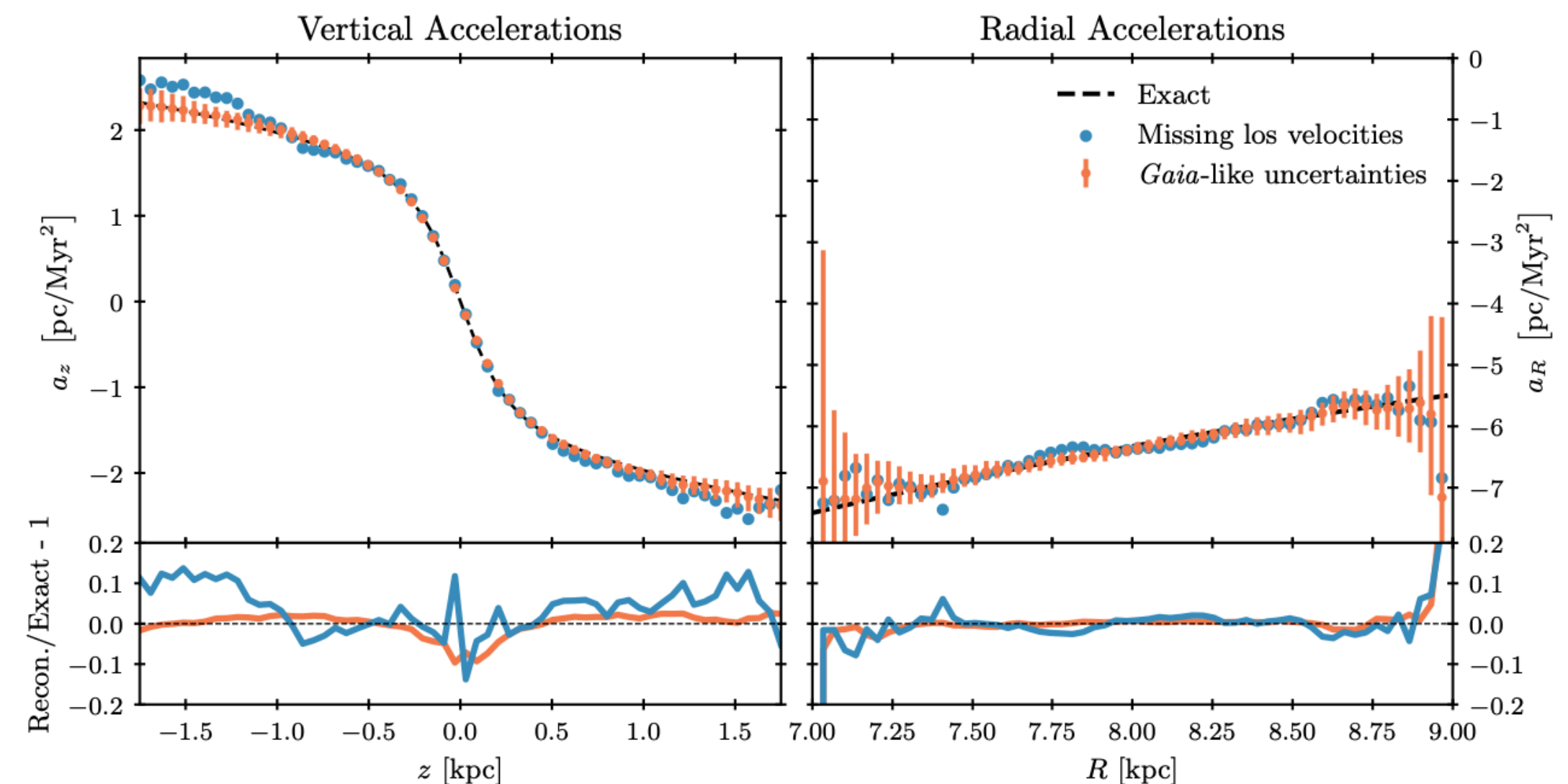- The phase space density of stars in equilibrium is related to the underlying Galactic potential

$$\frac{\partial f}{\partial t} + v_i \frac{\partial f}{\partial x_i} = \frac{\partial \Phi}{\partial x_i} \frac{\partial f}{\partial v_i}$$

- Curse of dimensionality makes it very hard to measure $f$ and derivatives from stellar motions. Traditionally, take moments of the Boltzmann Equation and assume symmetries

- Normalizing flows can do a much better job in estimating $f$ and its derivatives from the available data.

Green & Ting (2011.04673)



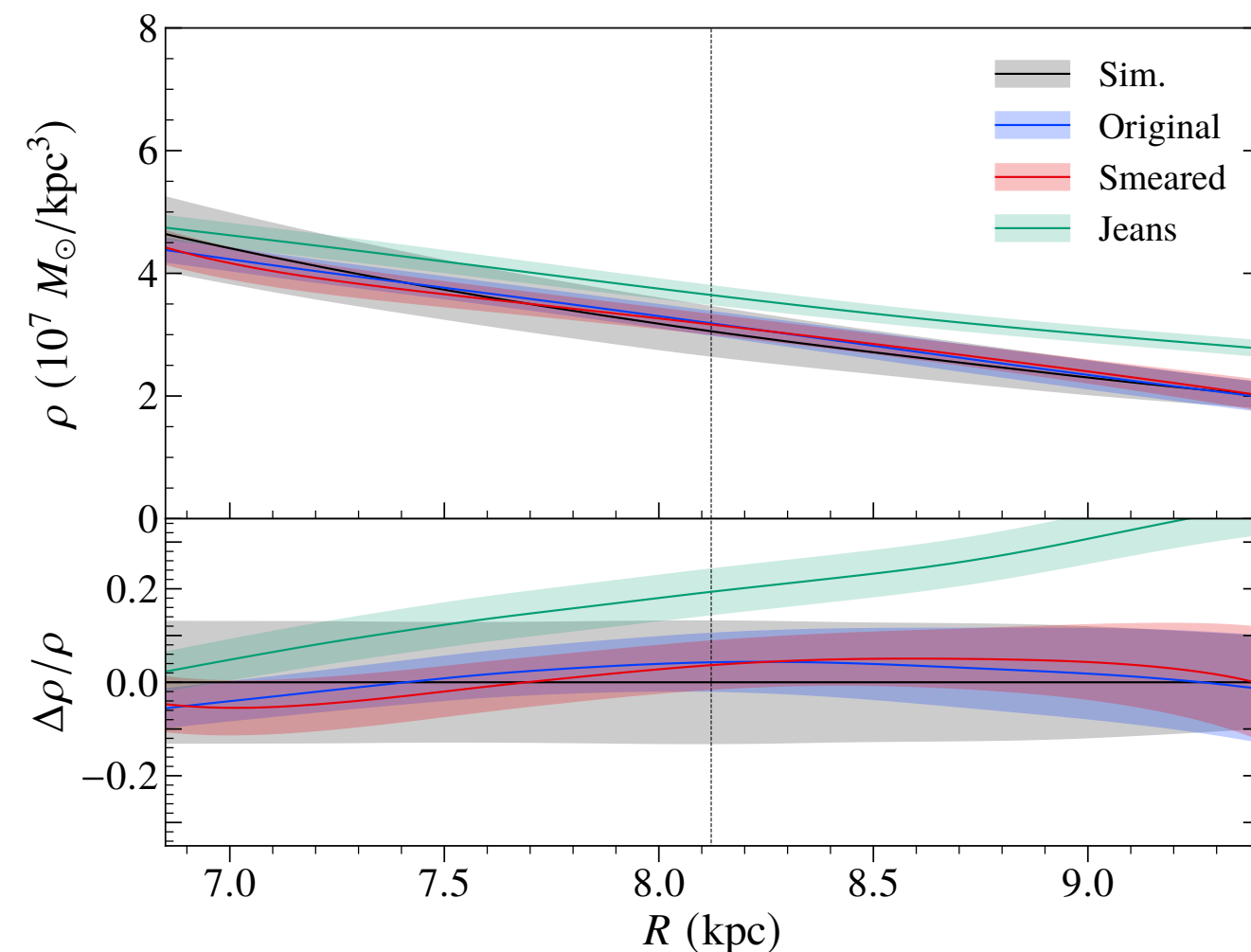An *et al* (2106.05981) and Naik *et al* (2112.07657)

- The real Galaxy is not in equilibrium:

$$\frac{\partial f}{\partial t} \neq 0$$

- Is real data sufficiently precise to get good estimates of $f$?

- First with a simulated Milky Way-like galaxy:
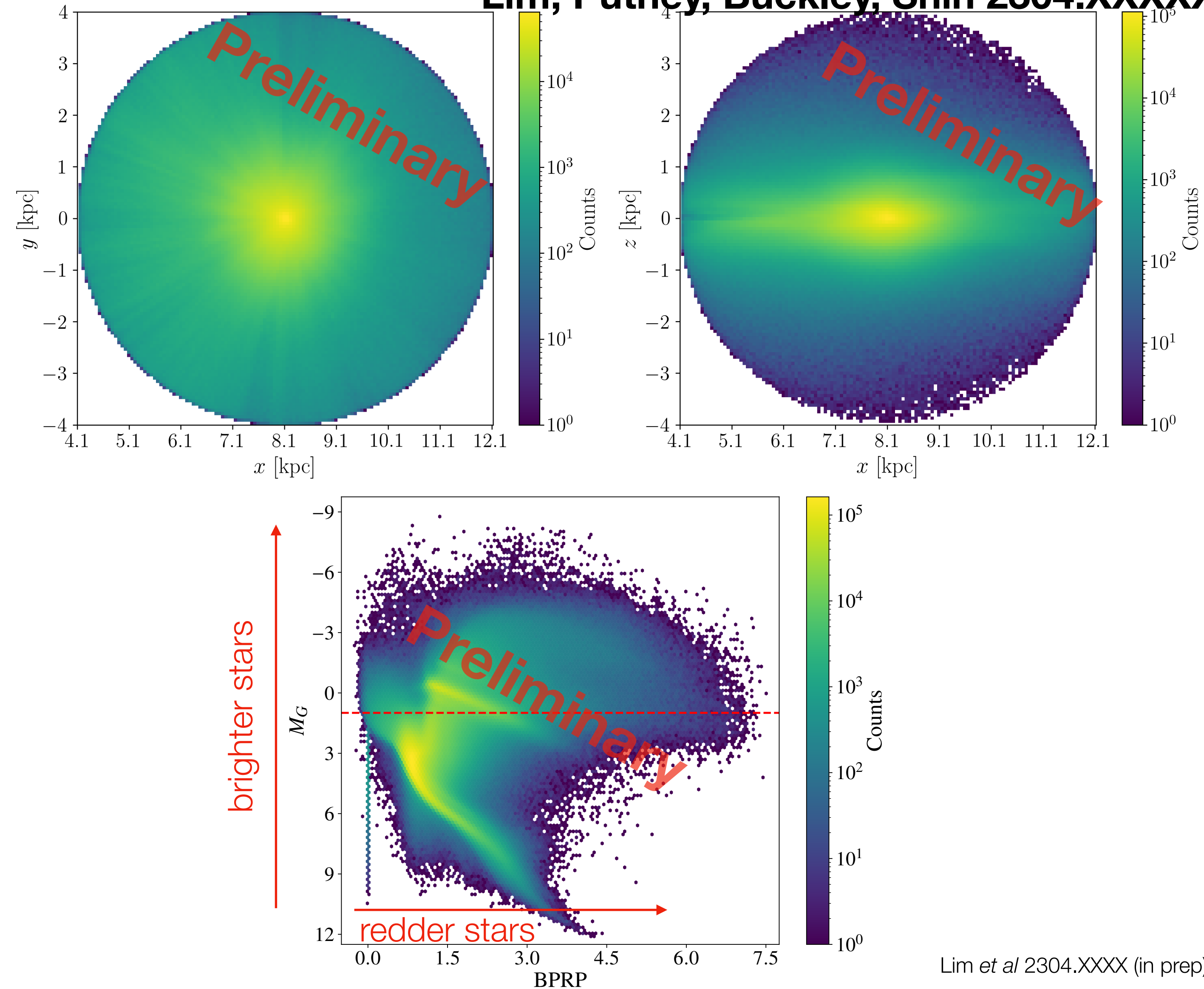
**Buckley, Lim, Putney, Shih 2205.01129**

- Can we do this with real Gaia data?

- Real data is complicated:

  - Observations are not complete, and this completeness varies as a function of distance

  - And with which kinematic parameters are measured, and/or stellar properties

- The goal: get low-error measurements off of the Galactic disk, to regions where dark matter dominates the mass density.
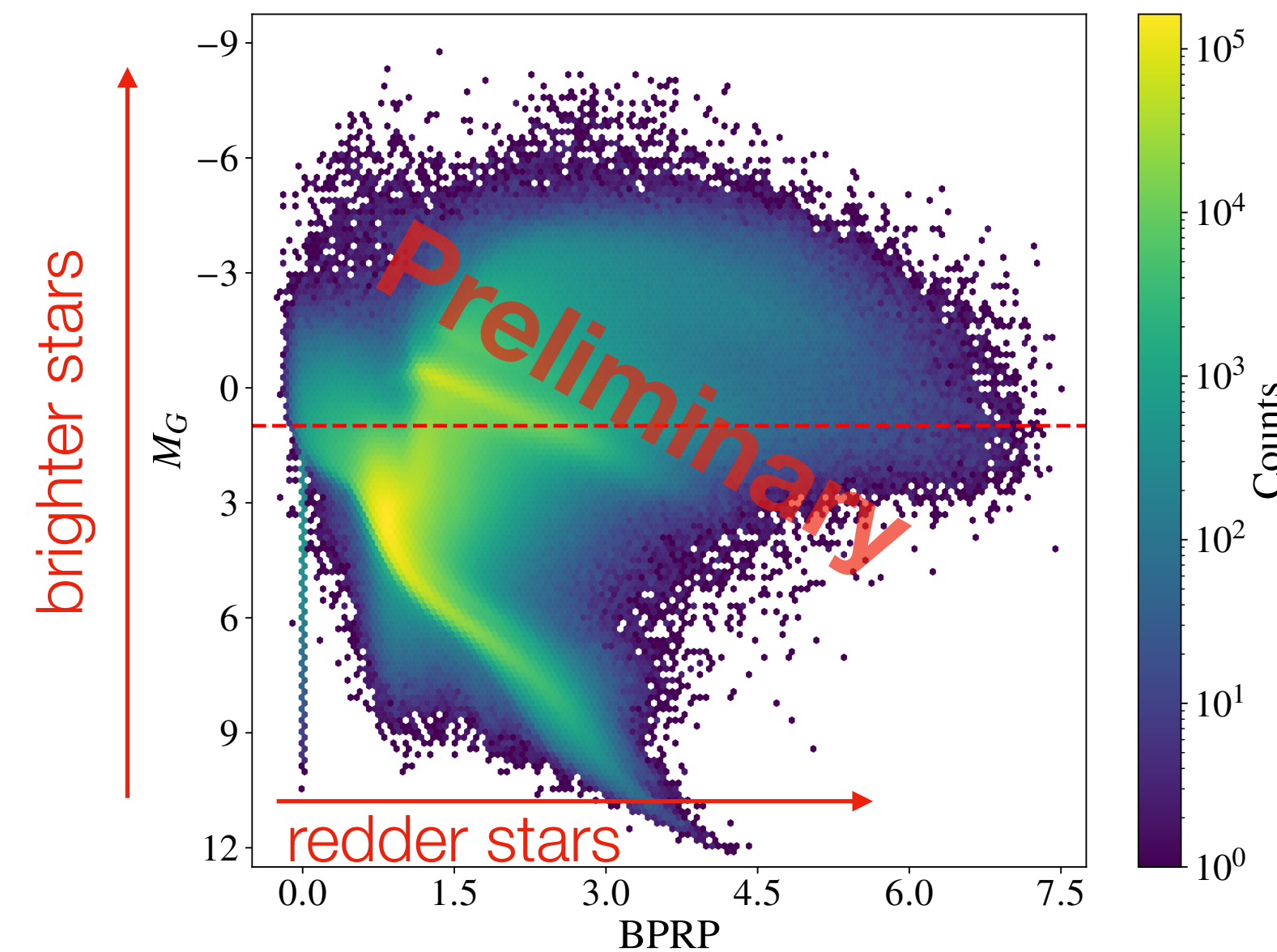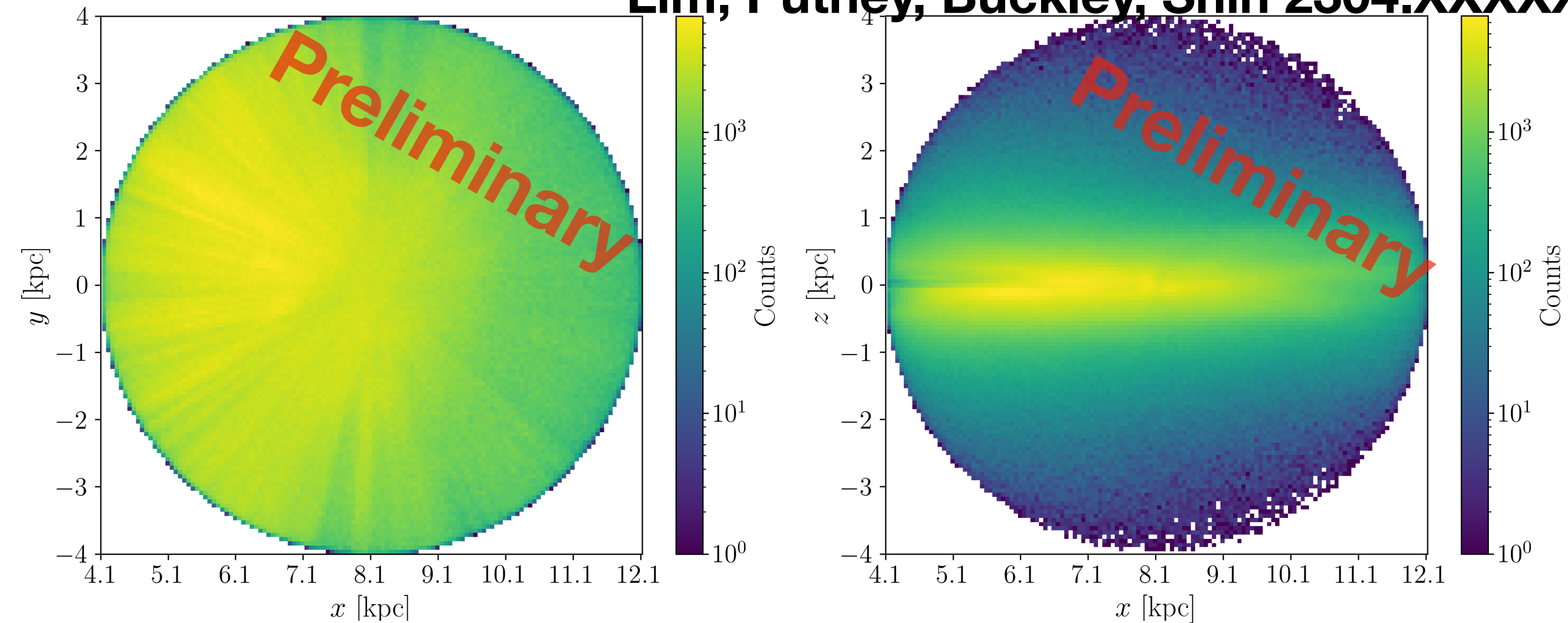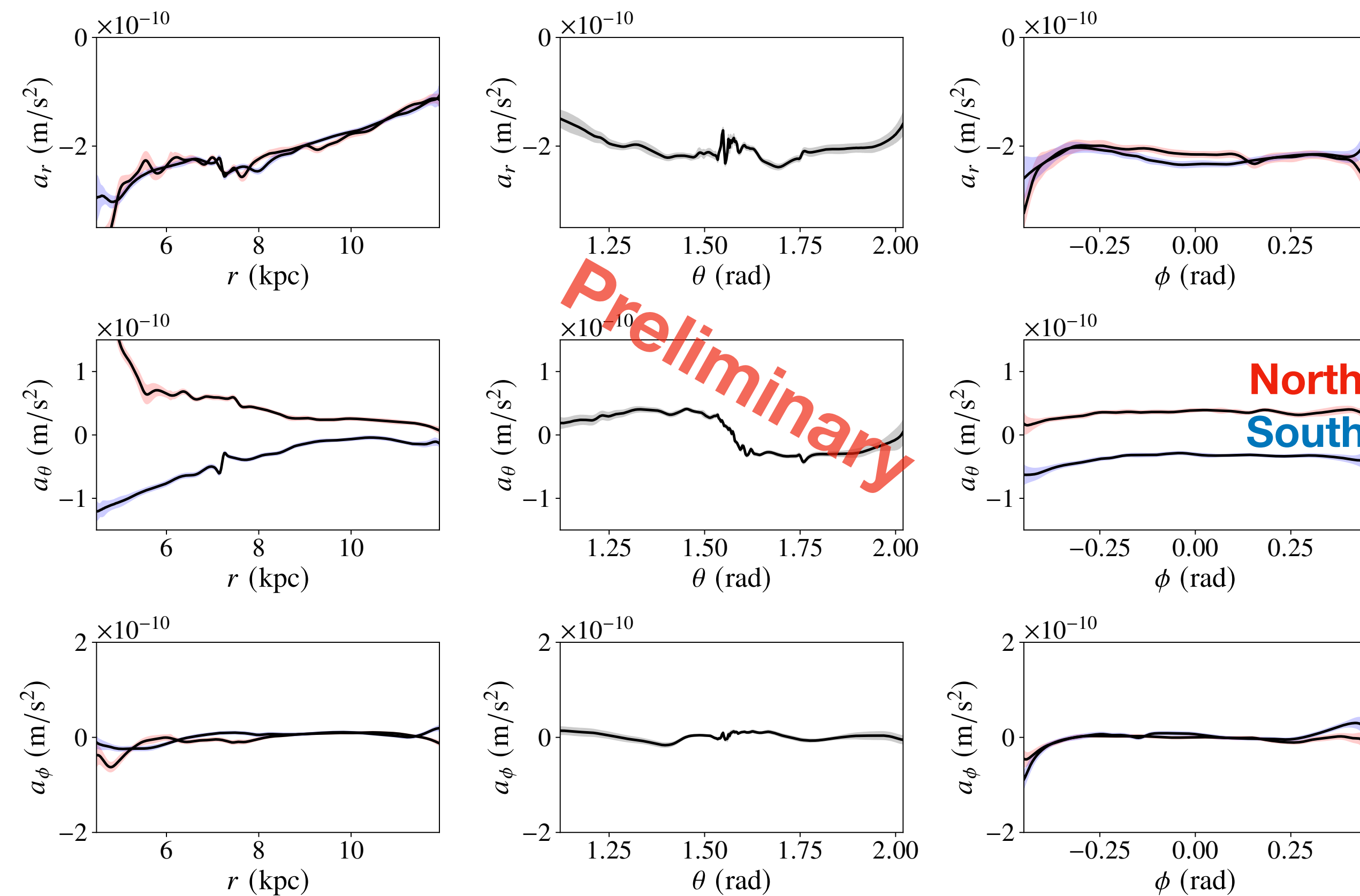
**Lim, Putney, Buckley, Shih 2304.XXXXX**



Lim *et al* 2304.XXXX (in prep)

- Can we do this with real Gaia data?

- Real data is complicated:

  - Observations are not complete, and this completeness varies as a function of distance

  - And with which kinematic parameters are measured, and/or stellar properties

- The goal: get low-error measurements off of the Galactic disk, to regions where dark matter dominates the mass density.

**Lim, Putney, Buckley, Shih 2304.XXXXX**



Lim *et al* 2304.XXXX (in prep)

- 1st: Calculate accelerations:

$$v_i \frac{\partial f}{\partial x_i} = \frac{\partial \Phi}{\partial x_i} \frac{\partial f}{\partial v_i}$$

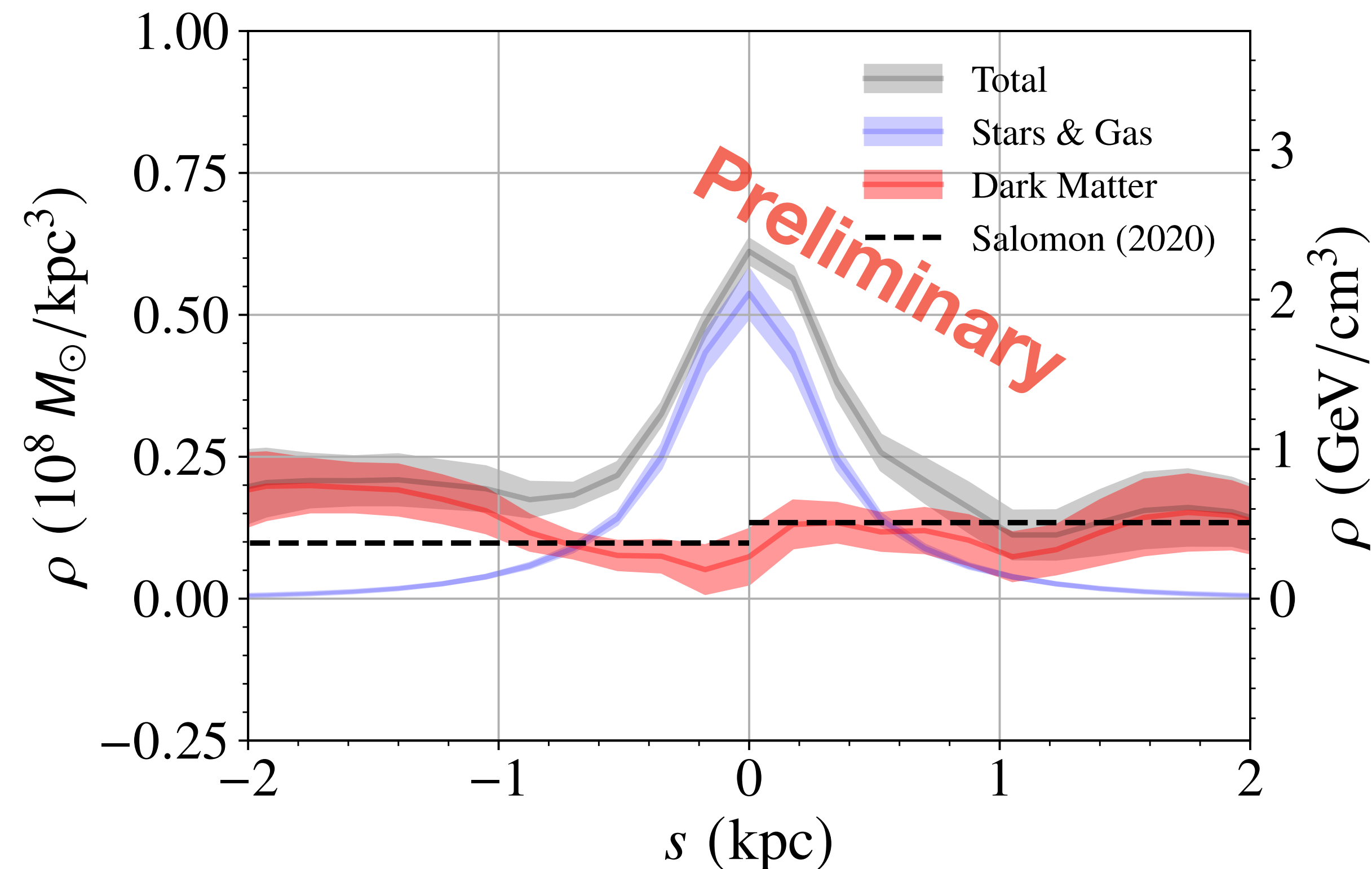- Errors include multiple MAFs, bootstrap, measurement errors



**Lim, Putney, Buckley, Shih 2304.XXXXX**

- Next, calculate mass density by integration by parts over a truncated Gaussian kernel

$$\nabla^2 \Phi = 4\pi G \rho$$

  - Baryonic model is a major source of uncertainty at the Solar location. Much less important away from the disk
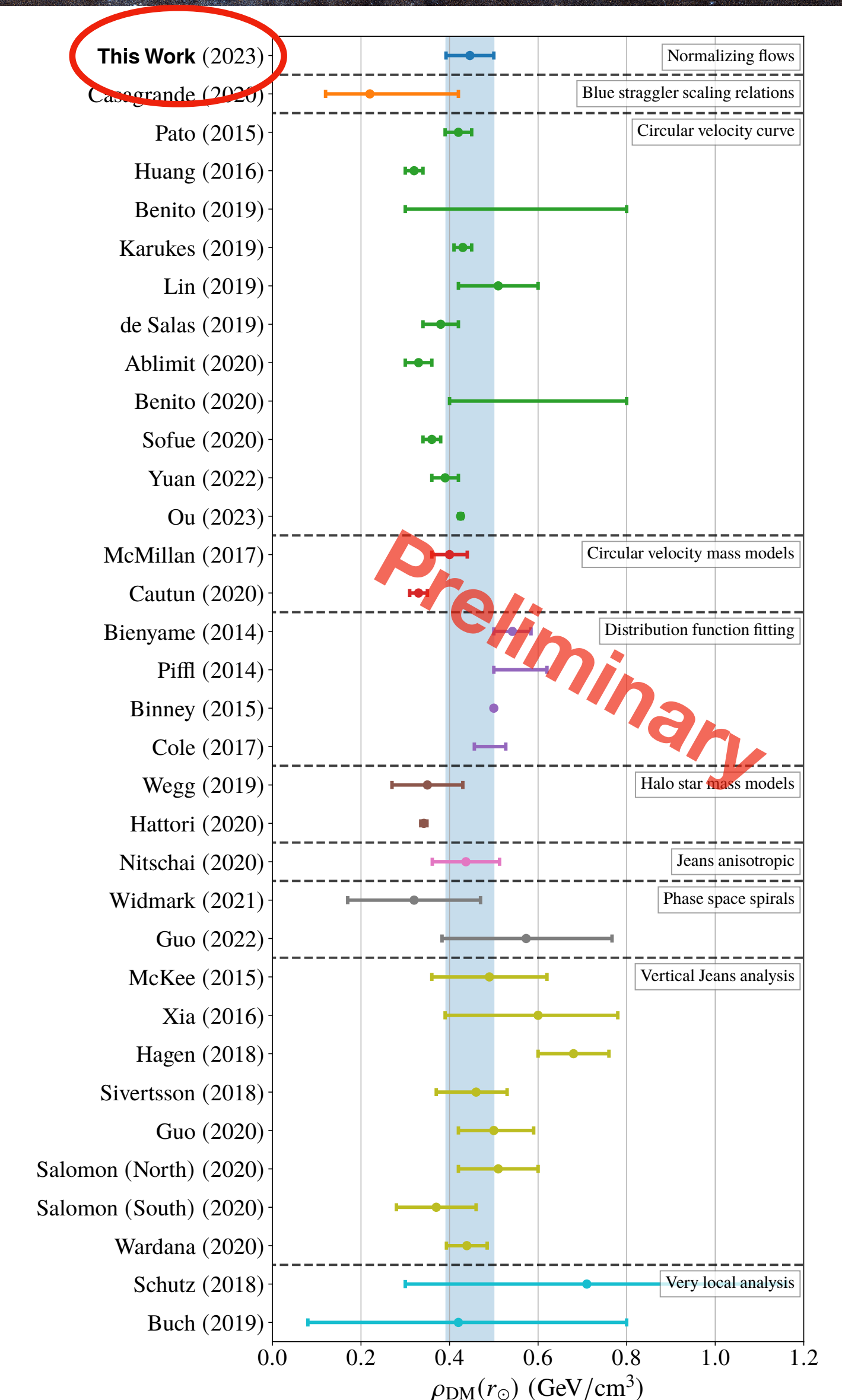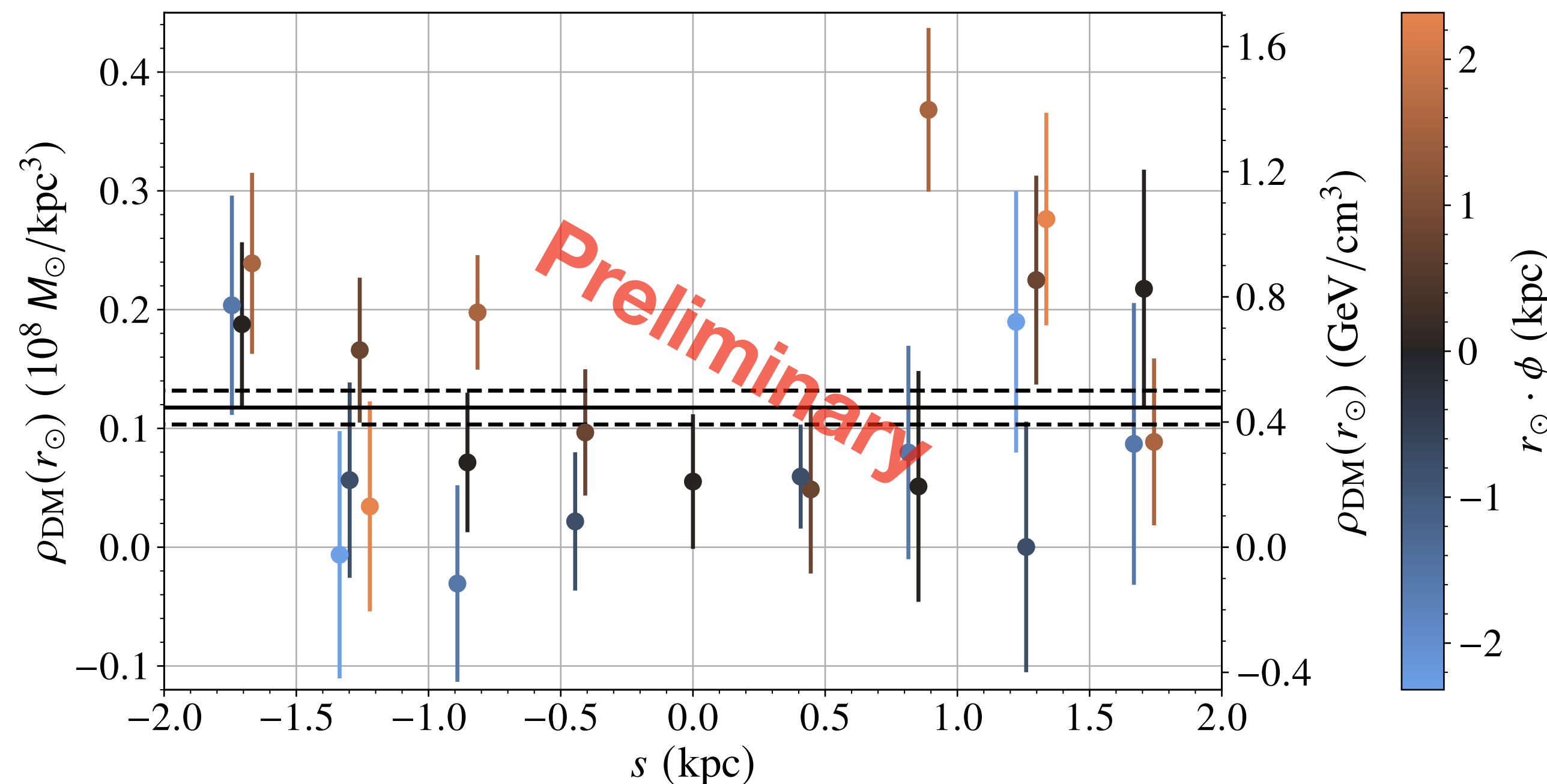


**Lim, Putney, Buckley, Shih 2304.XXXXX**

- Next, calculate mass density using finite differences (averaging over truncated Gaussian kernel)

  - Baryonic model is a major source of uncertainty at the Solar location. Much less important away from the disk
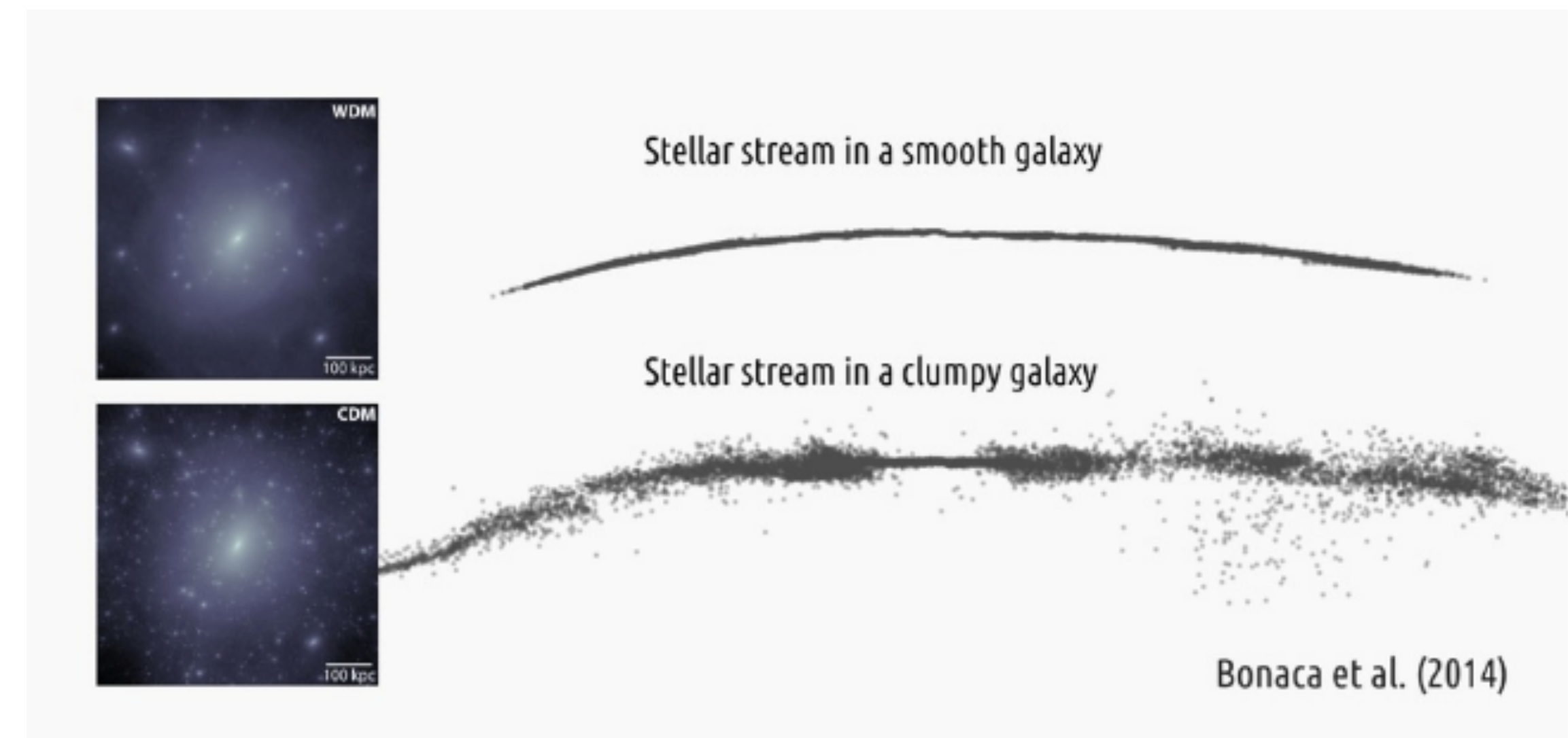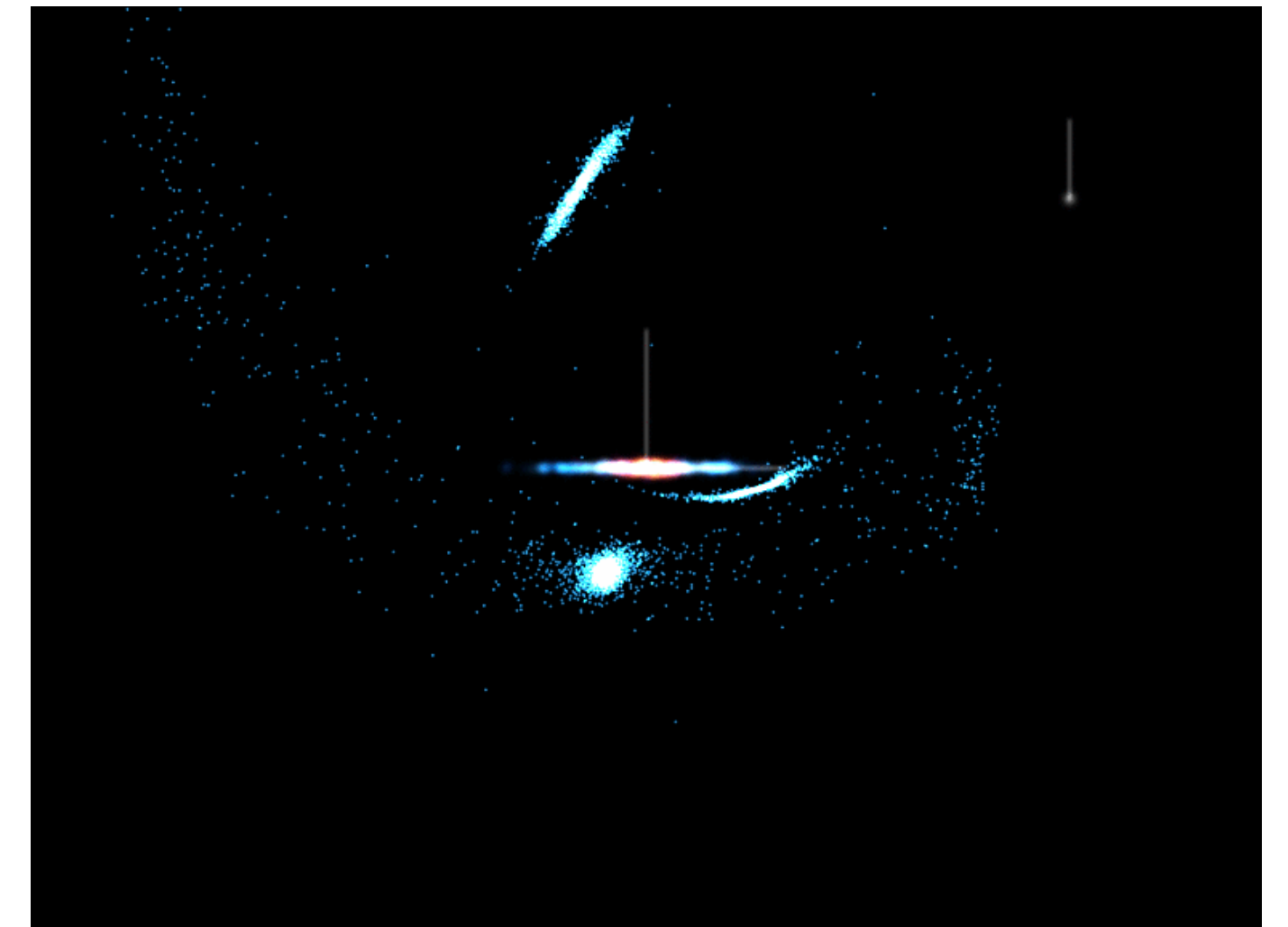
$$\rho_{\rm DM}(r = r_\odot) = 0.446 \pm 0.054 \ {\rm GeV/cm^3}$$



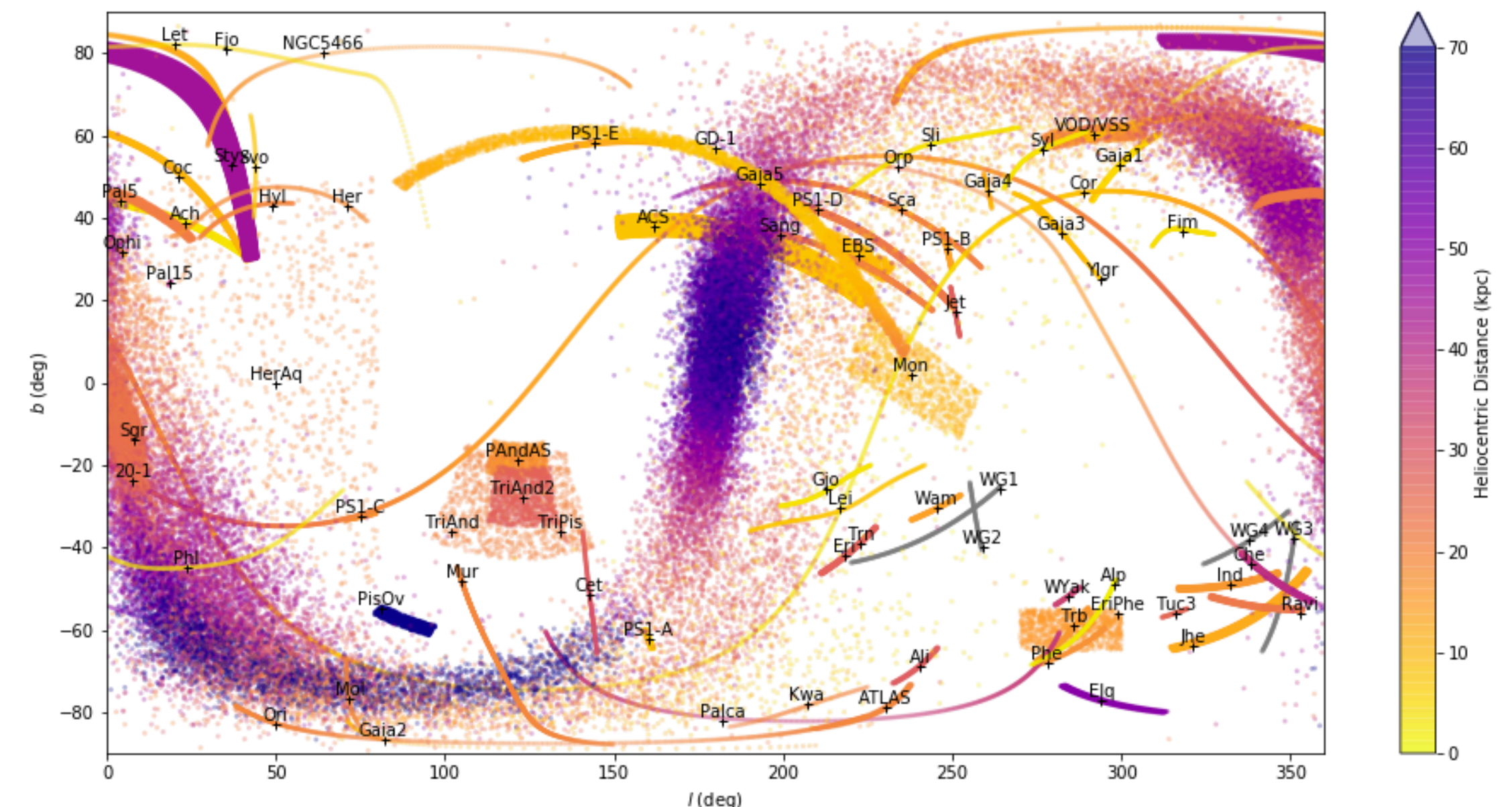**Lim, Putney, Buckley, Shih 2304.XXXXX**

- The Milky Way is built from the merger of smaller objects.

- Compact collections of stars (dwarf galaxies & globular clusters) get tidally stripped during infall and form **stellar streams**, then become **tidal debris**, before becoming completely mixed.



- Streams provide a probe into the Galactic potential through the stream's orbit.

  - Can reveal dark matter substructure through gravitational interactions with the stream itself.

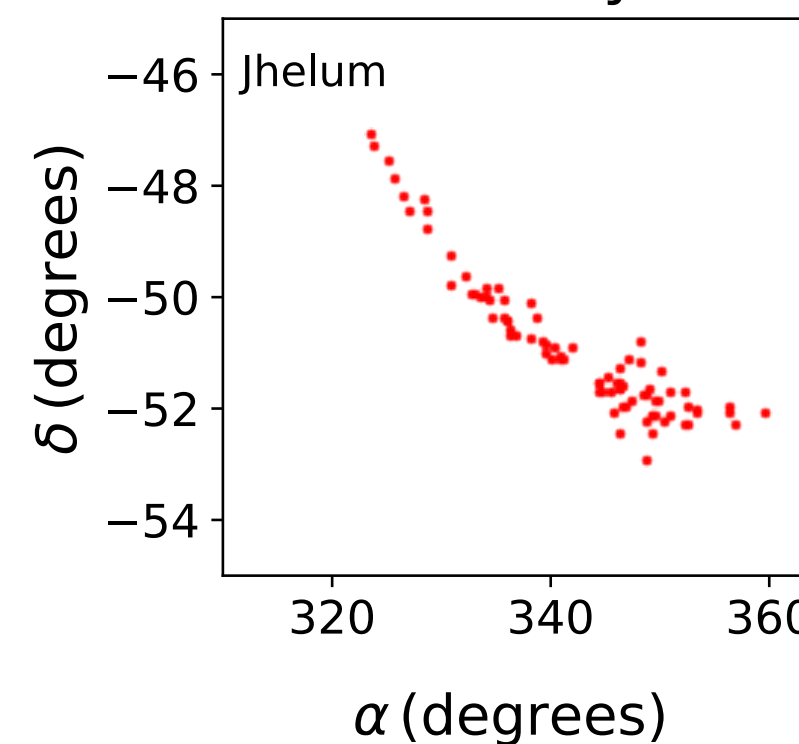- Both streams and debris give a glimpse into the Galaxy's merger history.



WDM

100 kpc

CDM

100 kpc

Stellar stream in a smooth galaxy
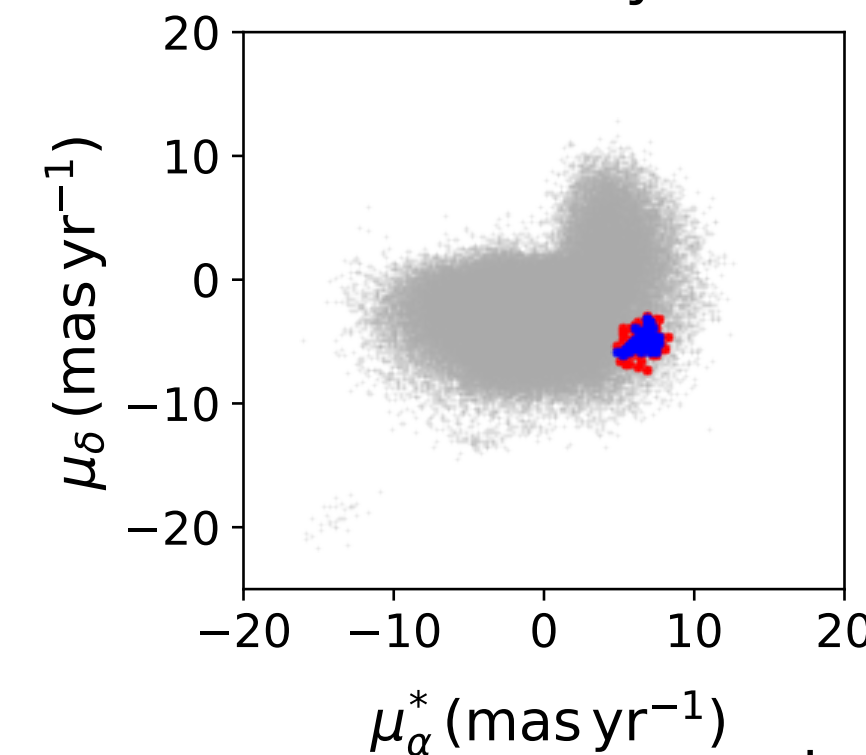
Stellar stream in a clumpy galaxy

Bonaca et al. (2014)

- Narrow & kinematically cold stellar streams are tracers of the Milky Way potential, merger history, imprint of dark matter substructure…

- A stellar stream is a narrow line of stars, compact in proper motion, and with all stars typically of similar age and composition.

- Use ML to build a stream-finding algorithm that:

  - Uses only Gaia data

  - Does not assume a Galactic potential or orbit

  - Does not assume stream stars lie on a particular isochrone.

  - Uses the fact that streams are compact in proper motion space.
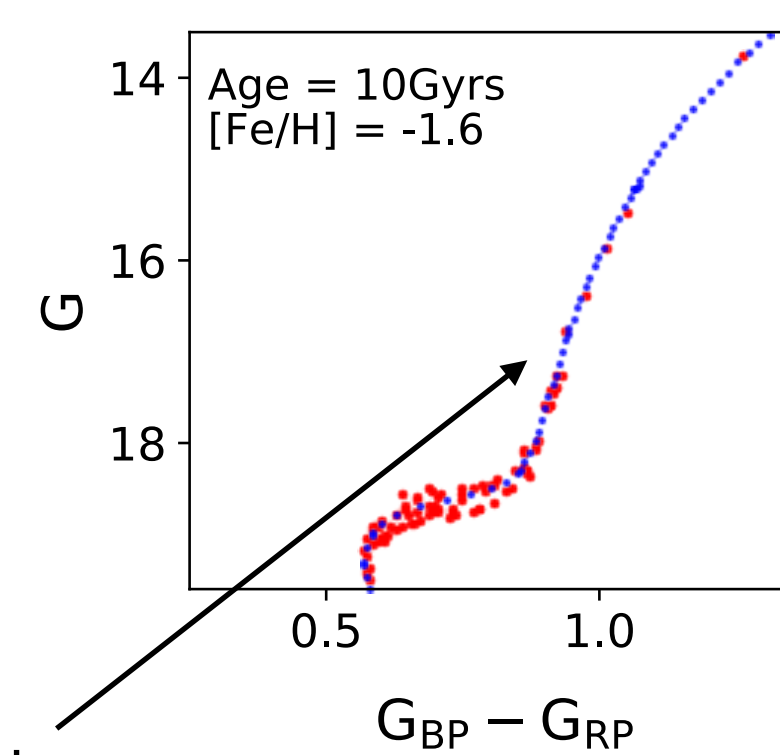


Angular position on sky
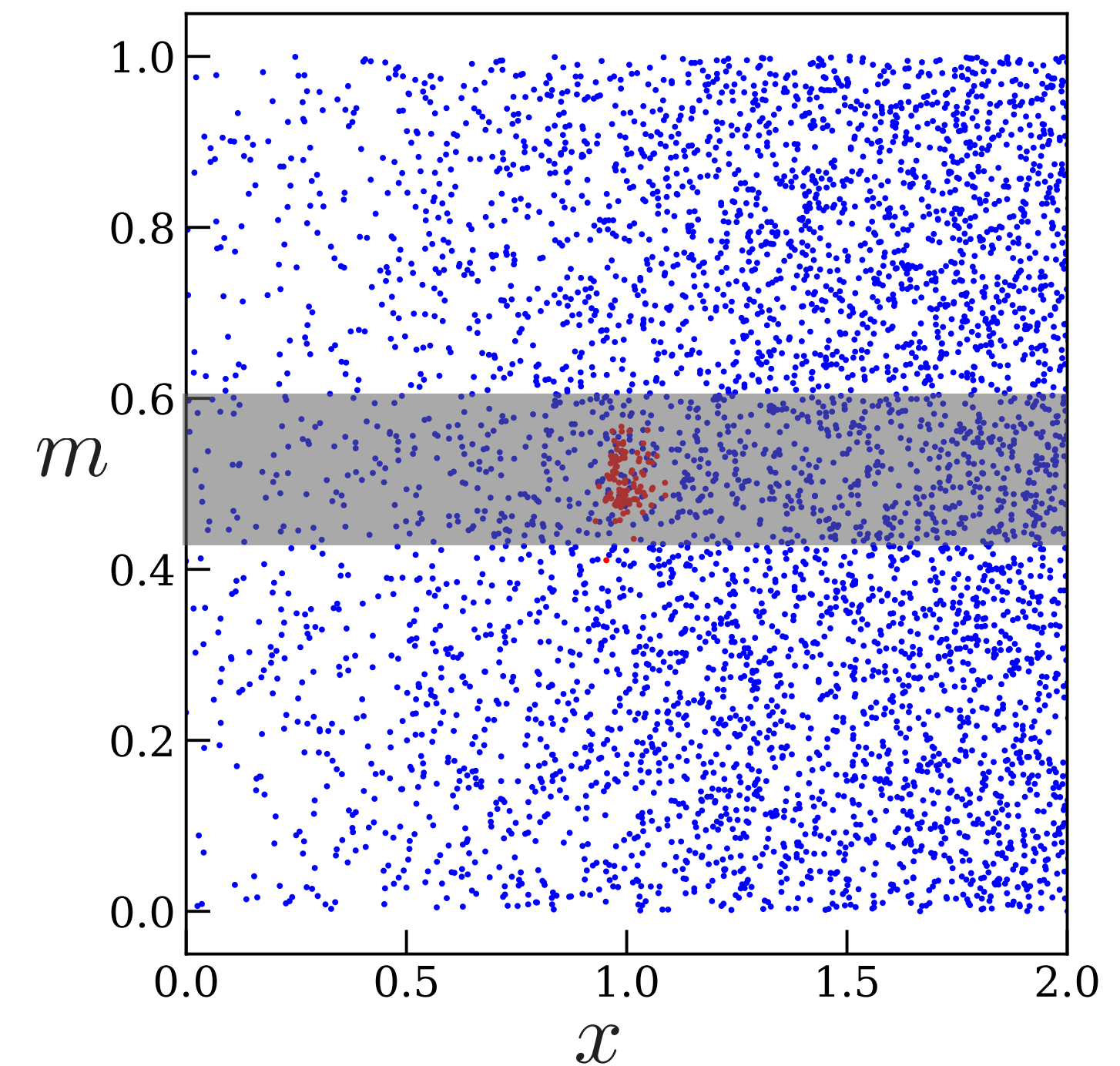
Angular motion on sky

Stellar brightness and color

isochrone

Malhan et al 2018

- Want to find stars that are anomalous based on their position in position, proper motion, and photometry. Use ANODE anomaly detection (Nachman & Shih 2001.04990) to calculate anomaly score $R$ for stars in proper motion Search Regions (SRs)

- Learn the probability distribution with $m \in [m_0 \pm \frac{\Delta m}{2}]$ in two ways:

  - 1st by training directly on the data in the region: $\approx P(\vec{x}|m)$

  - 2nd by training outside in a control region, then interpolating in: $\approx P_{\text{bkg}}(\vec{x}|m)$

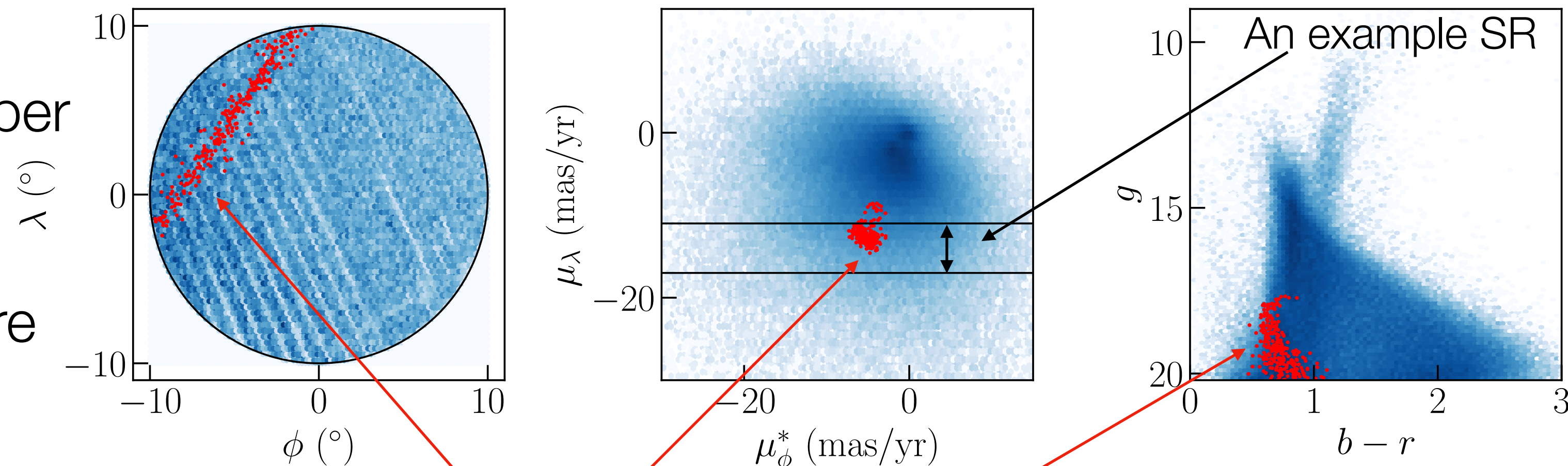- Allows direct estimation of the ratio $R$ inside the SR.

$$R(\vec{x}|m \in \text{SR}) = \frac{P(\vec{x}|m \in \text{SR})}{P_{\text{CR}}(\vec{x}|m \in \text{SR})}$$
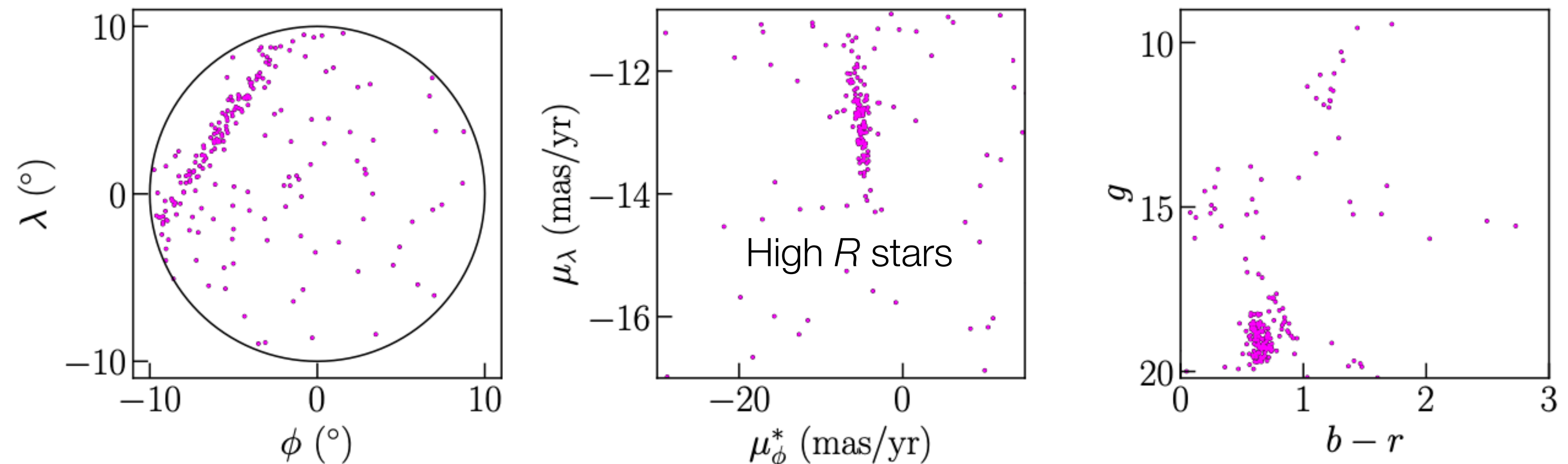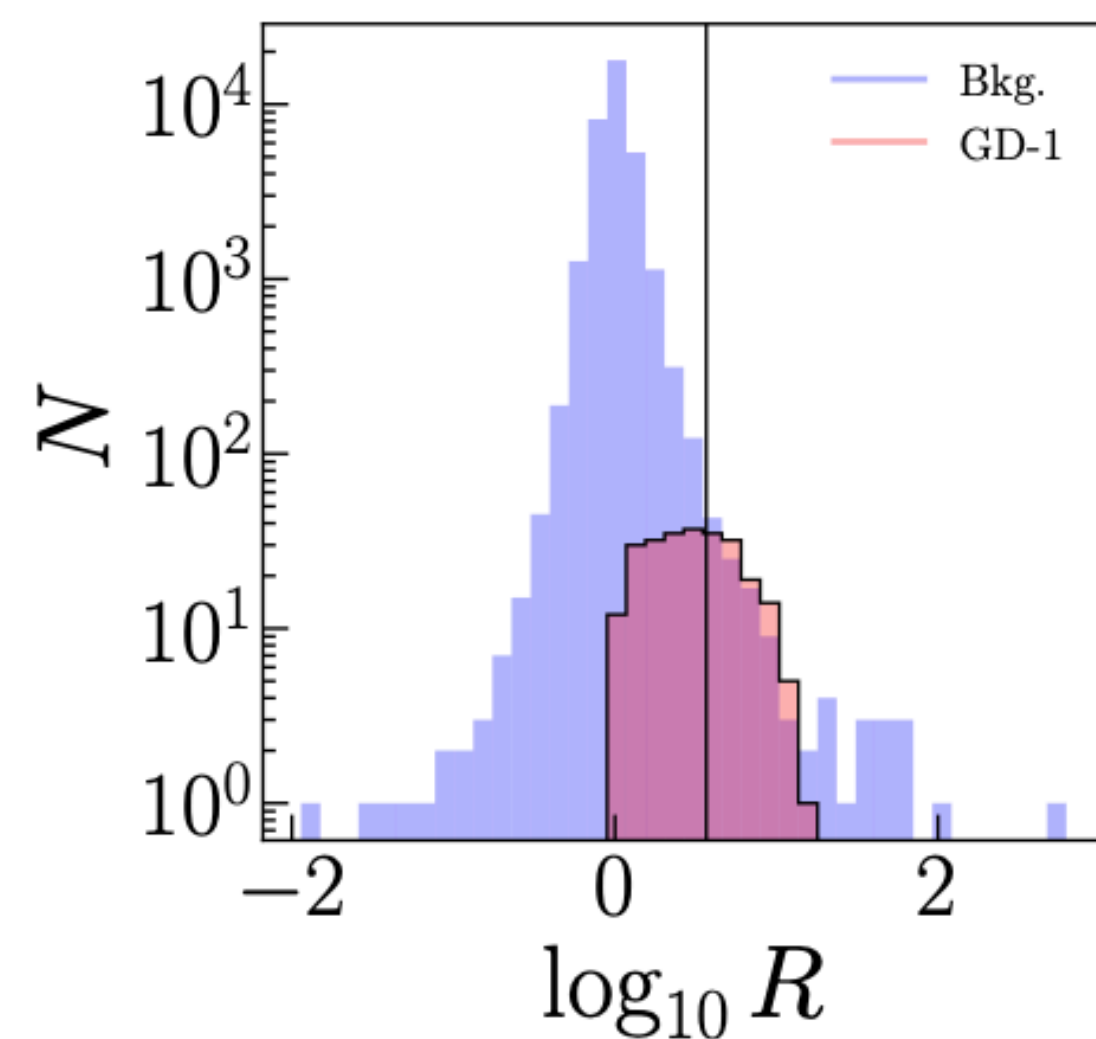
Shih *et al* (2104.12789)

- Want to find stars that are anomalous based on their position in position, proper motion, and photometry. Use ANODE anomaly detection (Nachman & Shih 2001.04990) to calculate anomaly score $R$ for stars in proper motion Search Regions (SRs)



An example SR

**Shih, Buckley, Necib, Tamanas (2104.12789)**

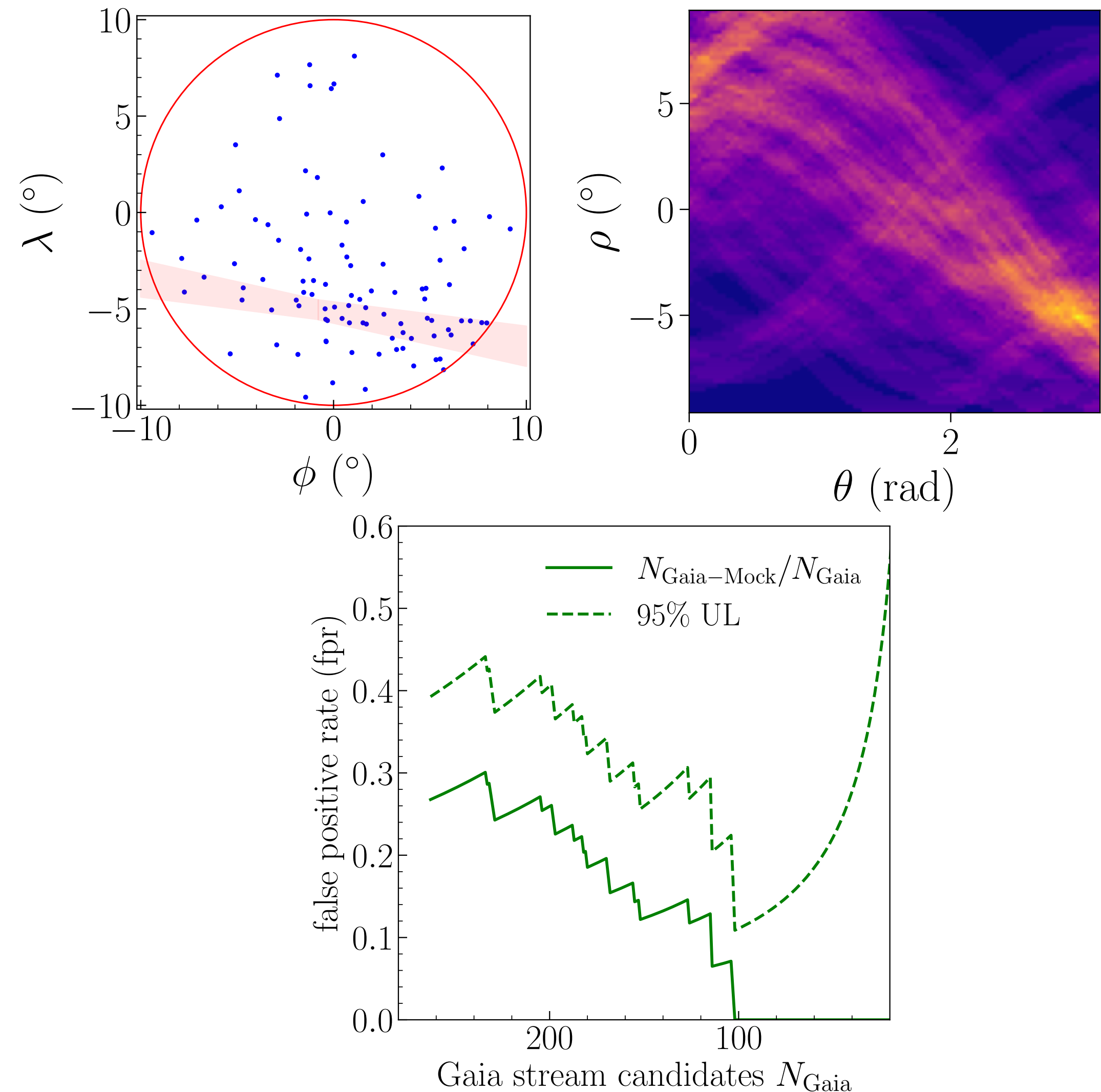Stars identified as likely GD-1 members by Price-Whelan & Bonaca



High $R$ stars

- There are a *lot* of stars in Gaia. Lots of reasons for them to be anomalous

  - Dust lanes, clusters, ...

- The ML architecture only needs to automatically identify linear features in the distributions of position and proper motion.

  - Many hyperparameters needed identify stellar streams at high confidence

- Use a smooth analytic simulation of the Milky Way (totally devoid of streams) to build an estimate of a false positive rate
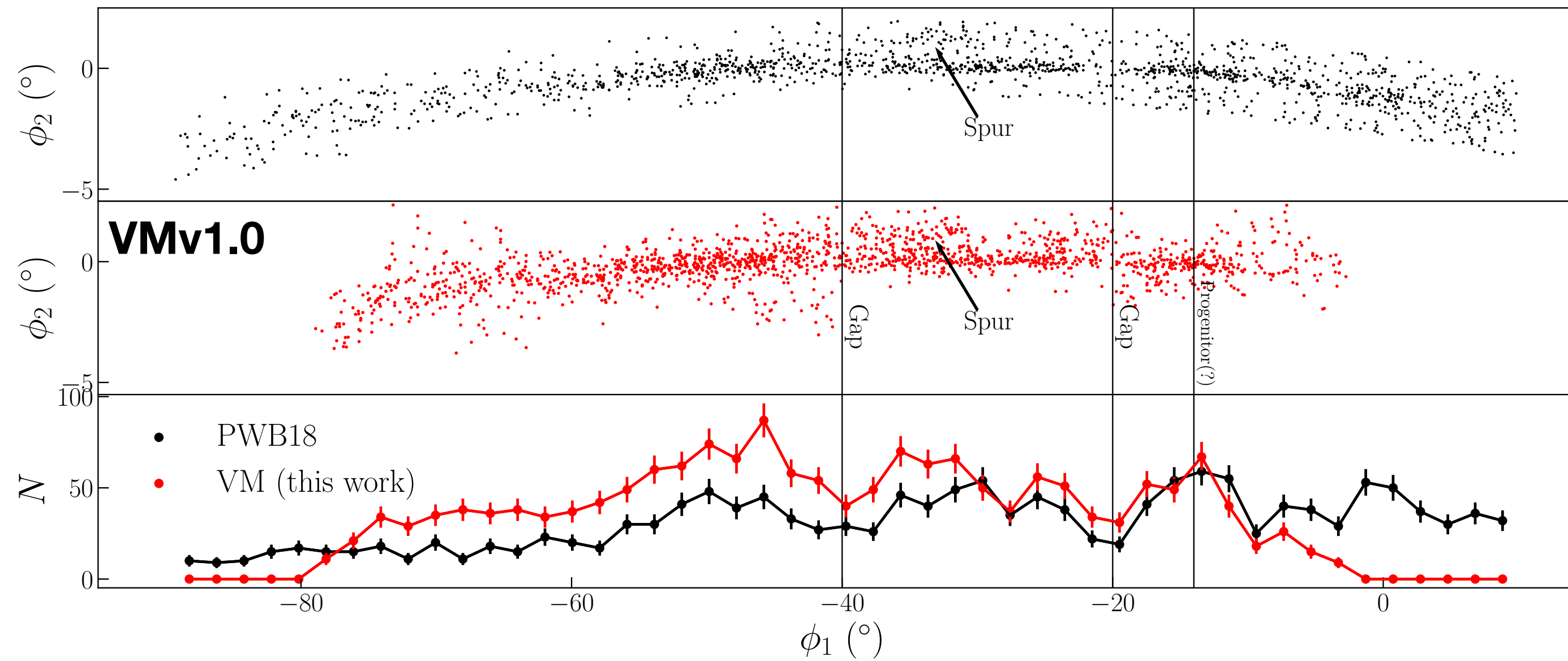
**Shih, Buckley, and Necib 2303.01529**
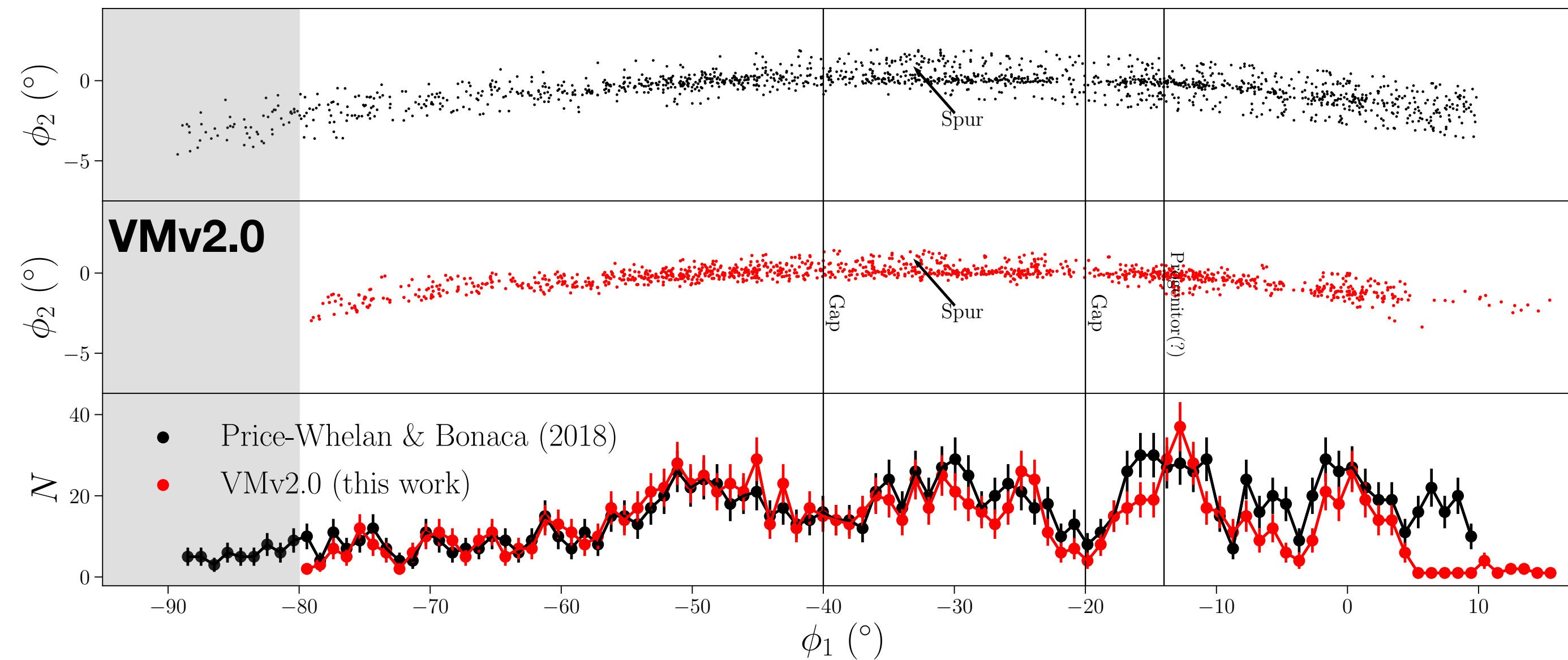
**Hough transform for line-finding**

**Shih, Buckley, Necib, and Tamanas 2104.12789**

**Shih, Buckley, and Necib 2303.01529**

- We identify 82 stream candidates, expect a false-positive rate of ~10%.
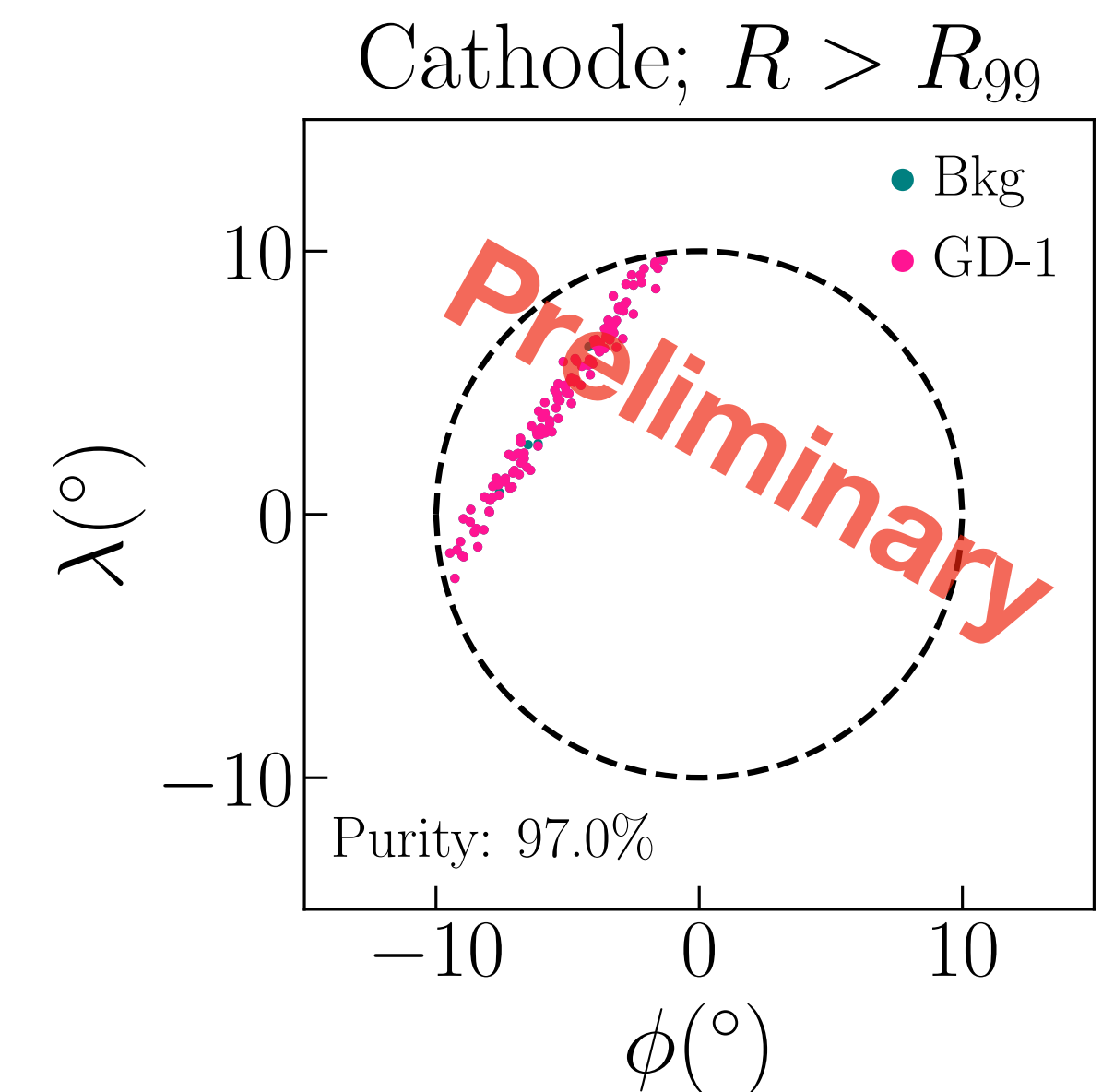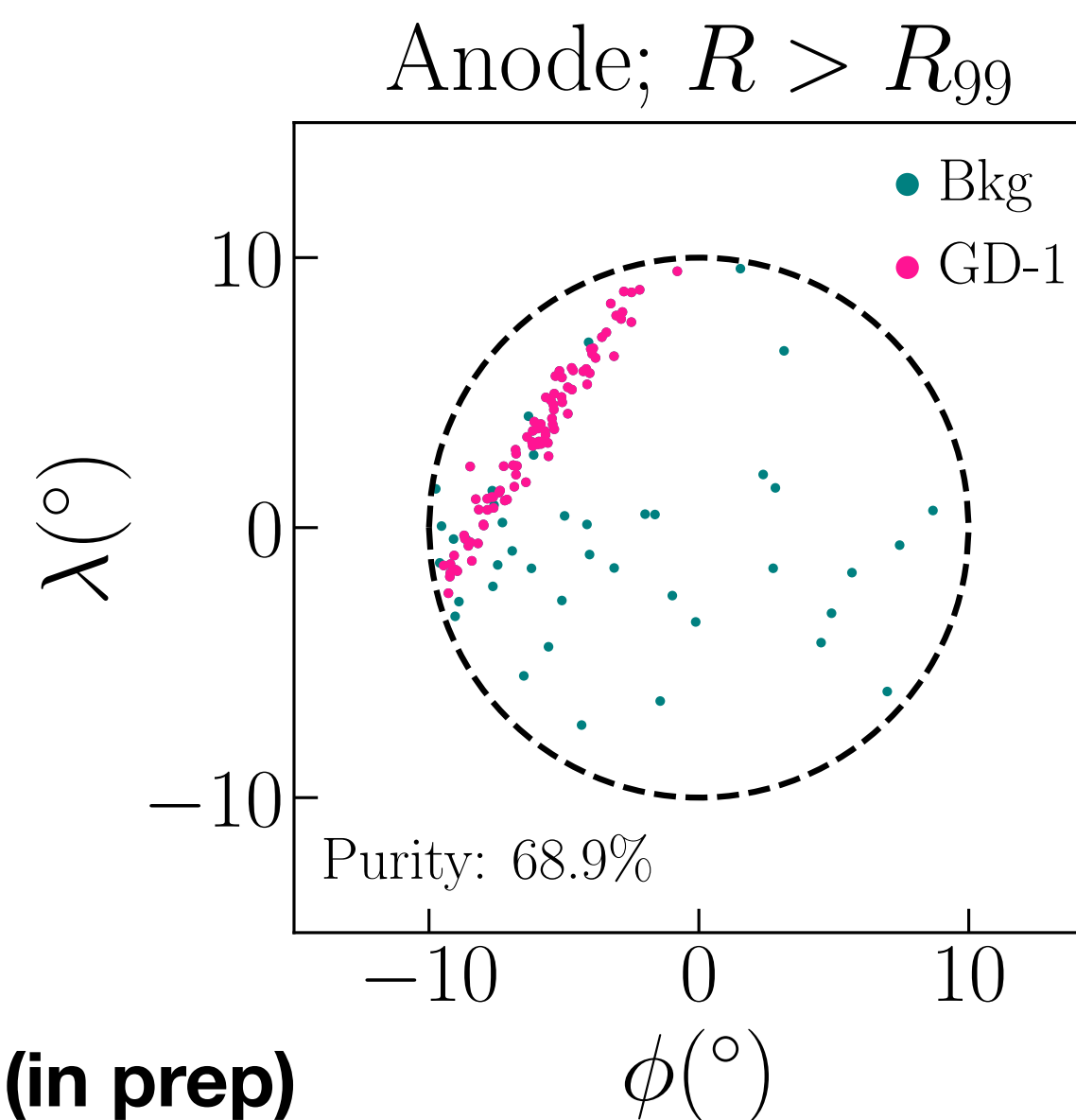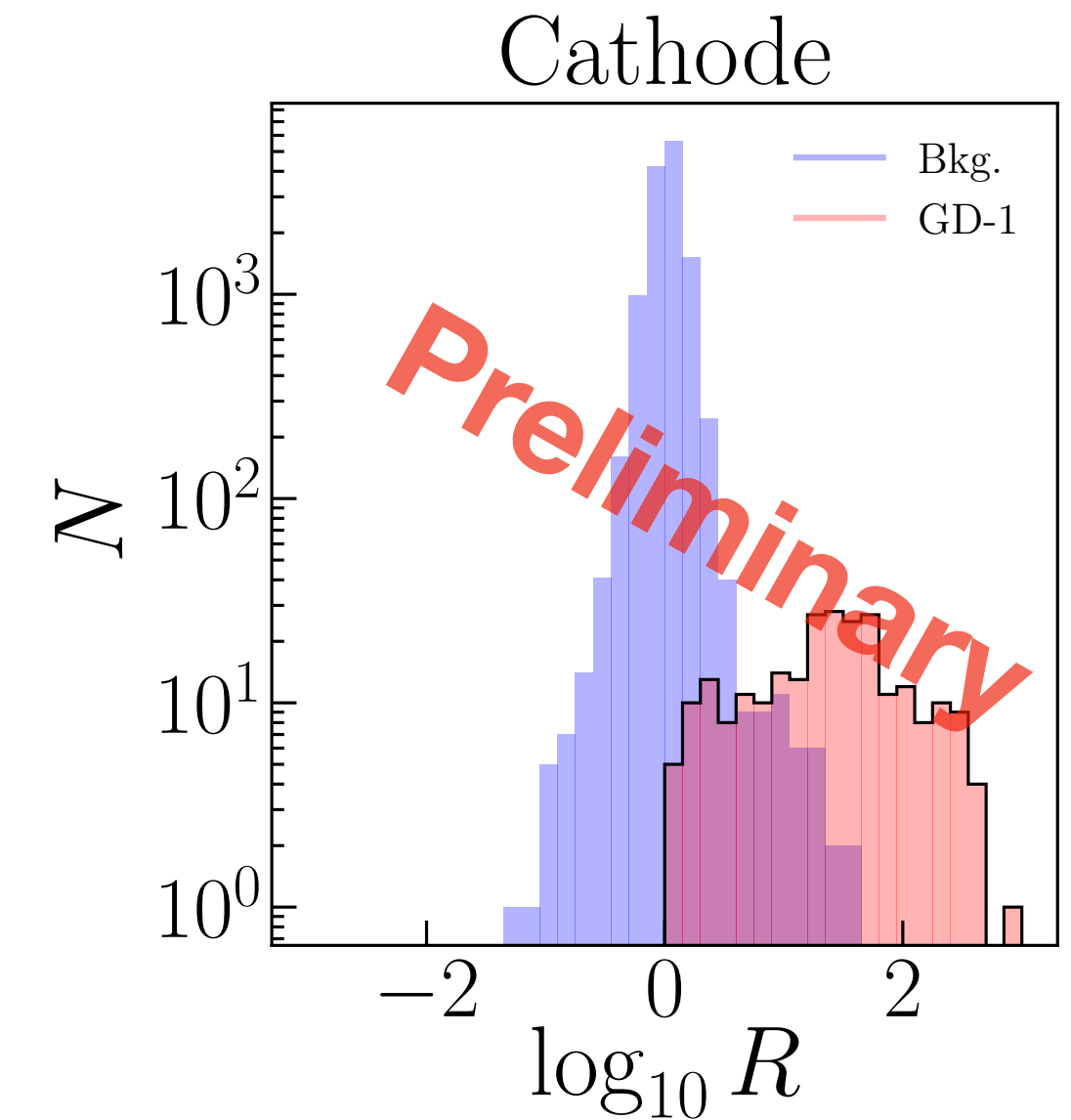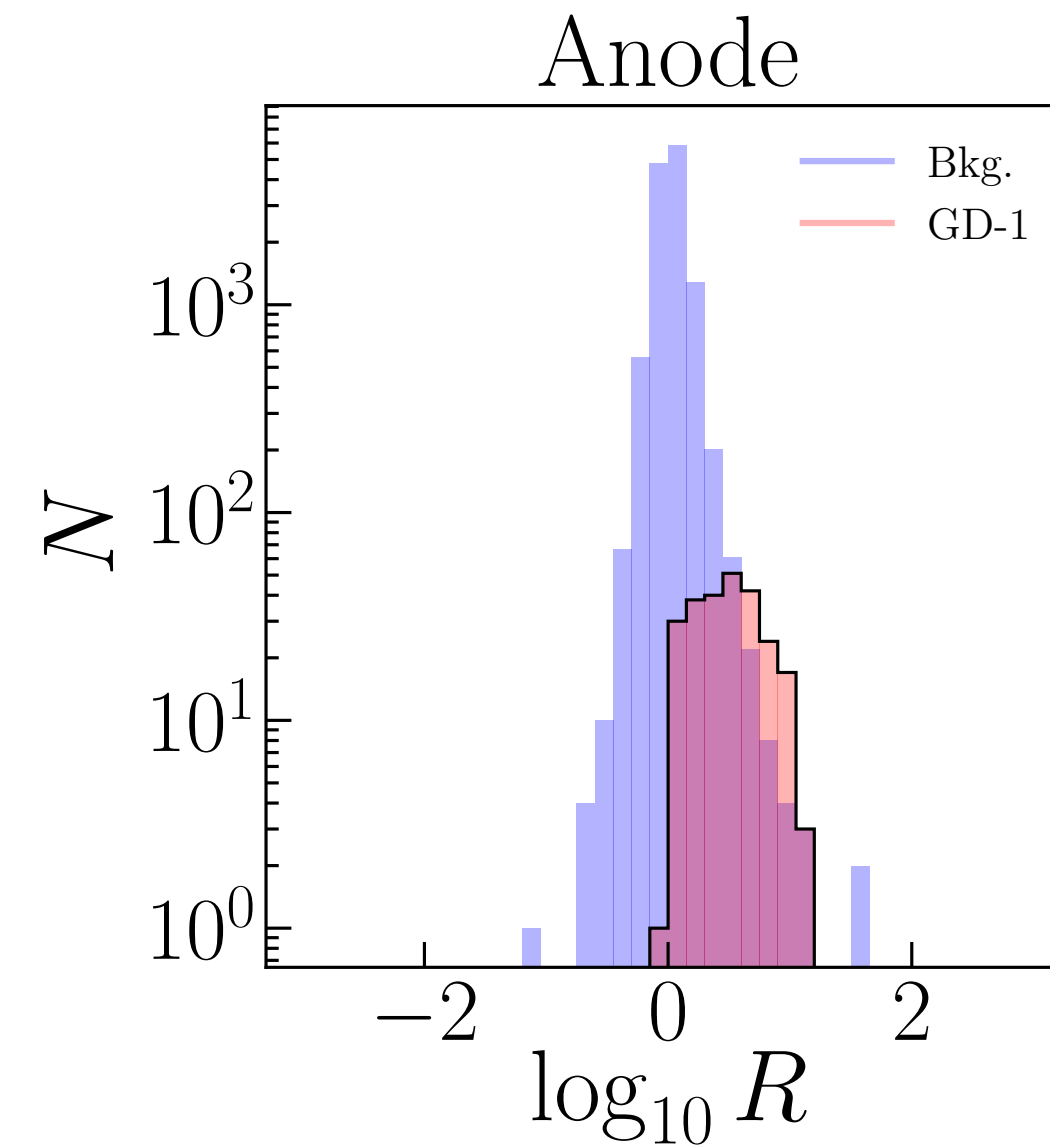
  - Here are the top 15.



- How to confirm stellar streams?

  - Spectroscopic follow-ups with other telescopes.

  - Do the stars have consistent metallicity, age, distance, radial velocity…?

**Shih, Buckley, and Necib 2303.01529**

- The input for the stream-finding is the ML-derived anomaly score $R$

  - Existing version from ANODE, using normalizing flows to learn conditional probabilities in proper motion SR and backgrounds from control regions.

- What if we could do this better?

  - CATHODE (Hallin *et al* 2109.00546)

  - Train a classifier to distinguish events generated in signal region from density estimator trained on control-region.

  - Use this as input for rest of Via Machinae

**Hallin *et al* (in prep)**



Anode

Cathode

Anode; $R > R_{99}$

Cathode; $R > R_{99}$

- Tools exist that can create "theorist-level" simulation for LHC machine learning.

- Much trickier for astrophysics. Can either:

  - Create by-hand analytic smooth models of the Galaxy or,

  - Use *N*-body hydrodynamical simulations

- But in the latter case, there complications:

  - Every galaxy is unique.

  - Simulations work on the level of tens of millions of "star particles," not hundreds of billions of *stars*.

- Upsampling required!

Galaxy h277 (N-Body Shop)



Lim *et al* 2211.11765
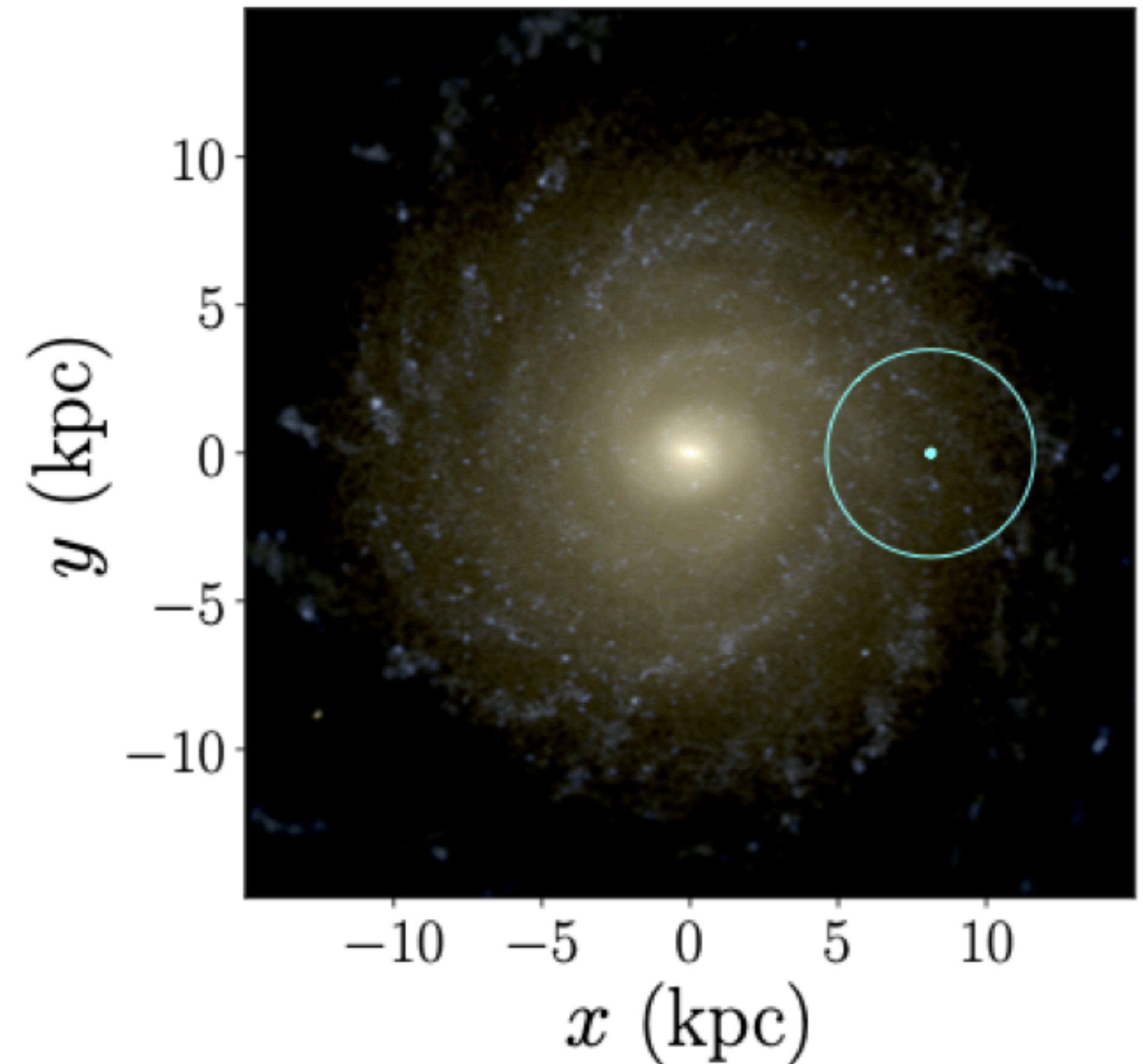
- Tools exist that can create "theorist-level" simulation for LHC machine learning.

- Much trickier for astrophysics. Can either:

  - Create by-hand analytic smooth models of the Galaxy or,

  - Use *N*-body hydrodynamical simulations

- But in the latter case, there complications:

  - Every galaxy is unique.

  - Simulations work on the level of tens of millions of "star particles," not hundreds of billions of *stars*.
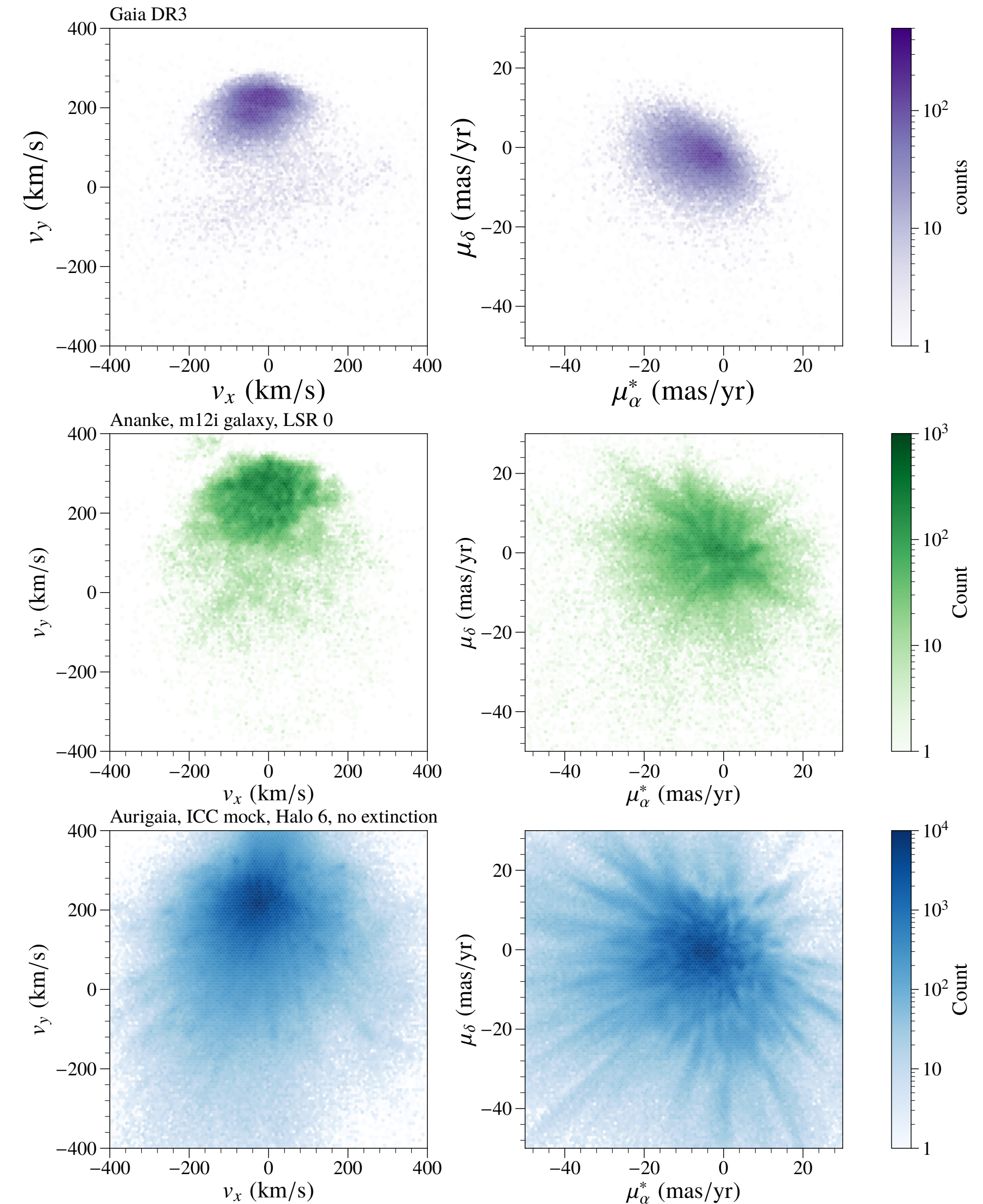
- Upsampling required!

  - But existing upsamplers are "clumpy"

Lim *et al* 2211.11765

# Upsampling Simulations

- Use normalizing flows (CNFs) to learn the density distribution of simulation star particles, then generate synthetic stars from the flow.

  - Demonstrating with stars near the "Sun"

  - Much smoother than stars drawn from existing upsamplers (EnBid)

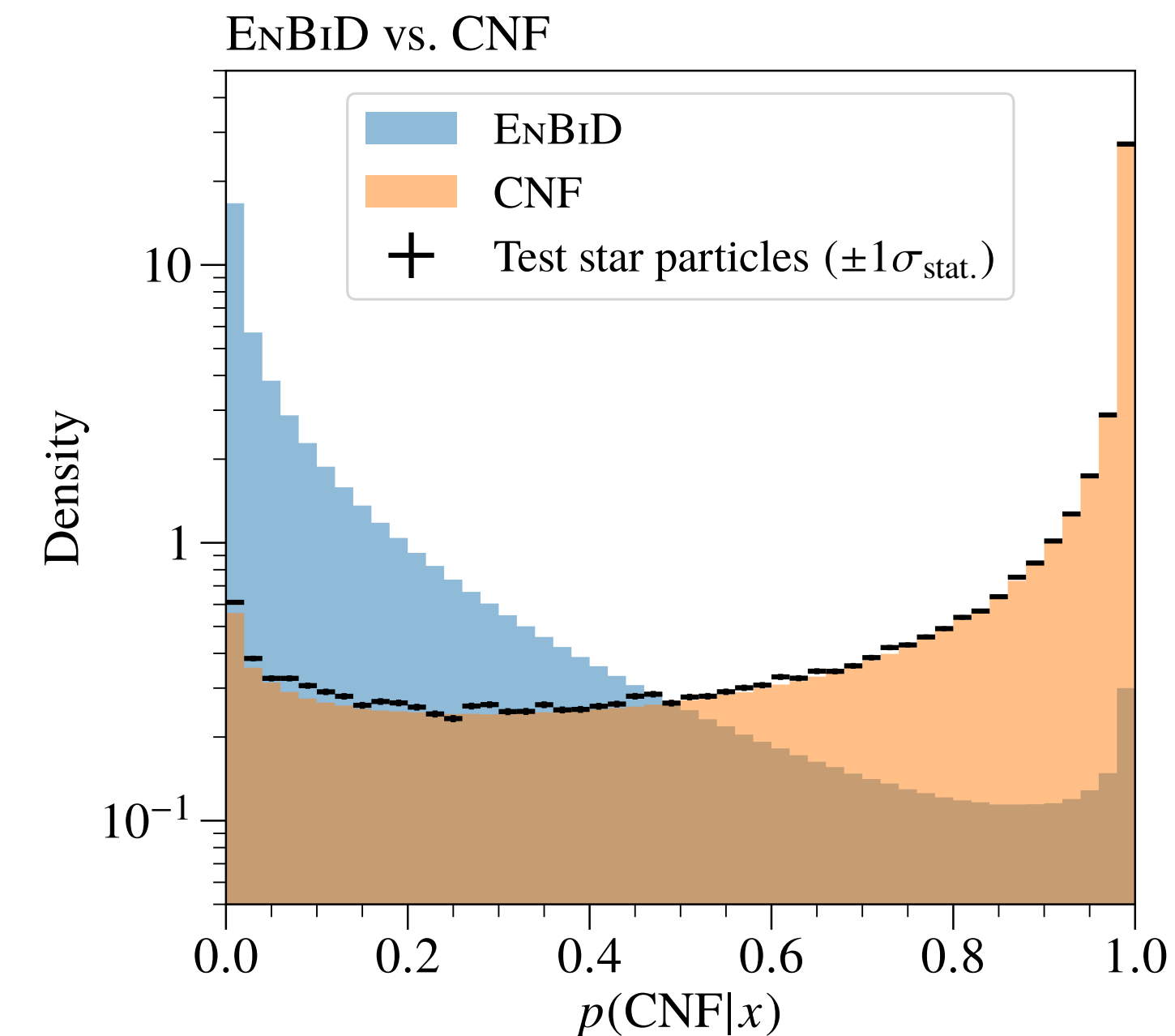  - Confirmed with classifier tests comparing CNF and EnBid

**Real Gaia Data**





Auriga 6, upsampled by ENBID



Auriga 6, upsampled by CNF

Lim *et al* 2211.11765

- 3-sample classifier: we are statistics-limited on the star particles

  - Construct CNF and EnBid datasets from a training subset of the star particles, reserving some star particles for validation

  - Train classifier between a subset of the CNF and EnBid datasets

  - Compare validation star particles with CNF and with EnBid separately

| network | classification target | AUC |
|---|---|---|
| trained on | EnBid vs. CNF | 0.952 |
| applied to | EnBid vs. Star particles | 0.950 |
| | Star particles vs. CNF | 0.508 |



EnBid vs. CNF

EnBid
CNF
+ Test star particles ($\pm 1\sigma_{stat.}$)

Density

$p(\text{CNF}|x)$

Lim *et al* 2211.11765

Shih, Buckley, Necib 2303.01529

- Astrophysical datasets contain information relevant to particle physics questions

  - …and intrinsically interesting on their own merits!

- The datasets are massive and complicated, with lots of systematic effects to deal with.

  - Often harder to simulate exactly what you'd need to test your technique. Interesting ML problems here in transfer learning, generation, quantifying errors.

  - Unsupervised techniques very useful.

- *Gaia* data in particular has lots to say about dark matter and Galaxy structure/history.

  - Lots of need for new techniques, opportunities for ML to help!



Buckley *et al* 2205.01129

Lim *et al* 2211.11765