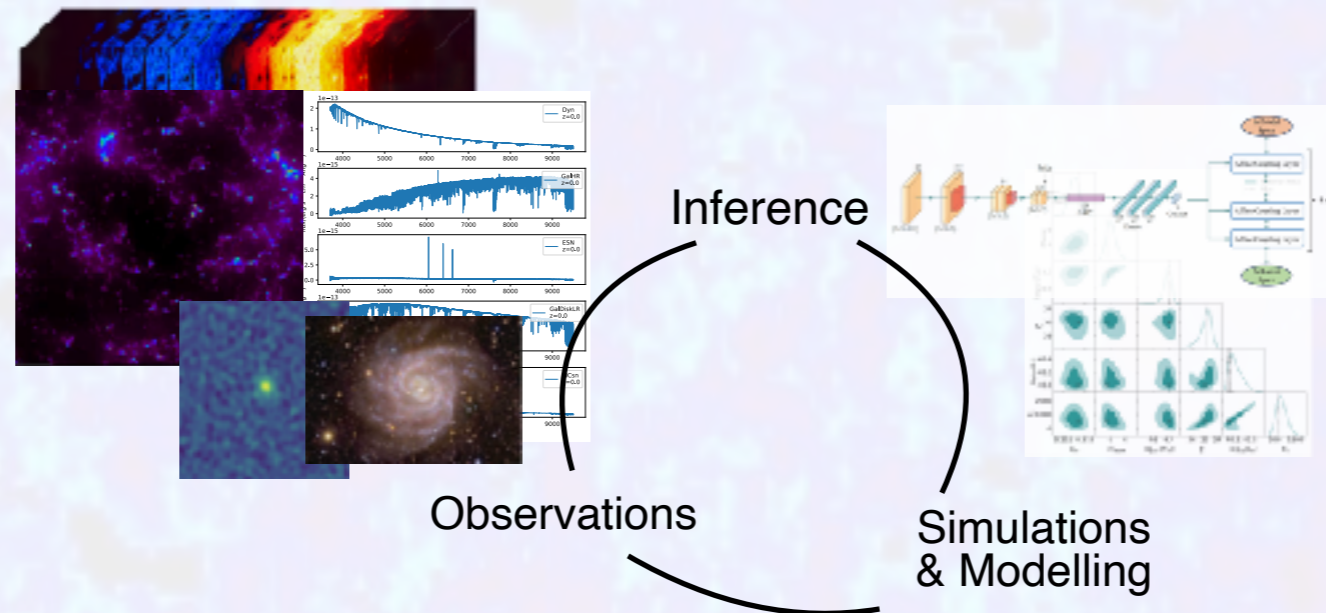


Machine Learning for Astrophysics



Caroline Heneka, ITP Heidelberg

Group 'Computer Vision Astrophysics and Cosmology'

Advanced School on Applied Machine Learning, ICTP Trieste, May 28th 2024

Who am I? Caroline Heneka Institute for Theoretical Physics, Heidelberg University



- B.Sc. and M.Sc. Physics in Heidelberg (+ Erasmus NPAC Paris)
- Ph.D. (2017) at Copenhagen University,
DARK Cosmology Centre, Niels Bohr Institute
- 2017-2019: DFG Transregio 33 Fellowship (Heidelberg), Postdoc SNS Pisa
- ca. 1.5 yrs: DLR (German Aerospace Center) Headquarters Cologne,
Executive Board Area Space, Programme Strategy Space
- 2020-2022: Senior Postdoc Hamburg University
- Since Oct 2022: back in HD
Junior Group Leader & Freigeist (Volkswagen Foundation) Fellow
'Computer Vision Astrophysics and Cosmology'



My Research Interests

- **Line Intensity Mapping**

(also: radio galaxy clustering, cross-correlation studies, galaxy clusters, ..)



- **High-redshift astrophysics & Epoch of Reionization:**

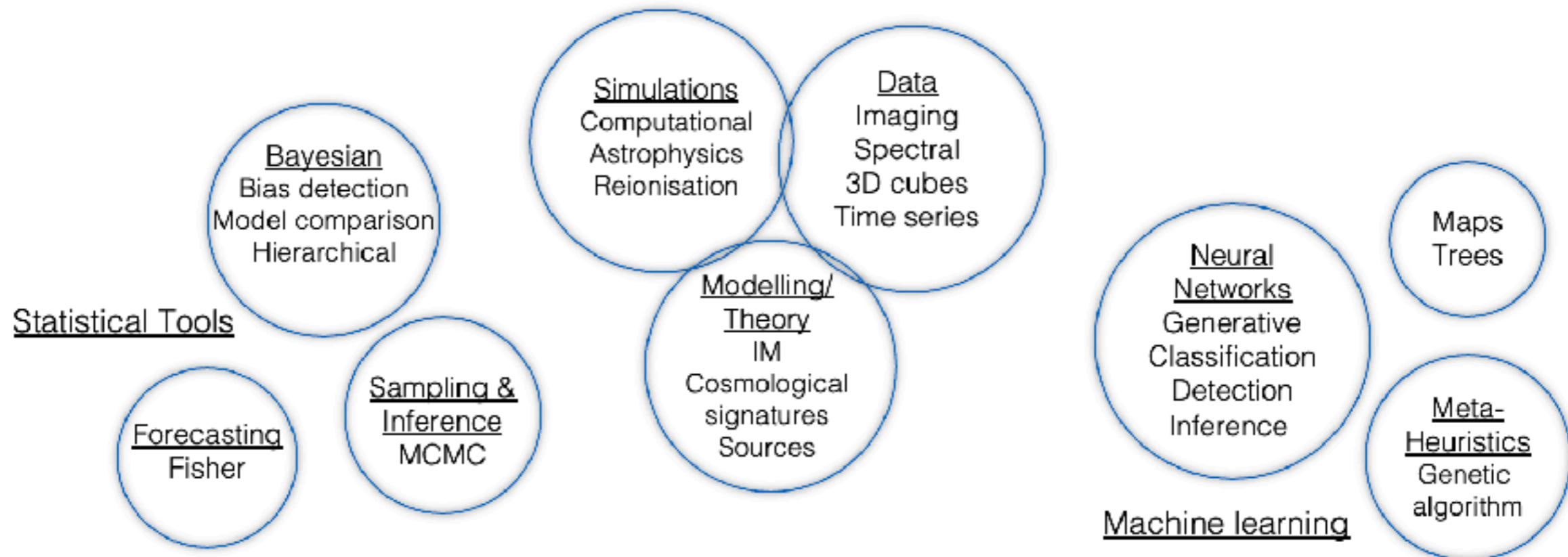
Modelling of 21cm background and further high-redshift lines (Lya, Ha, ..)

- The modern machine learning toolkit with '**Computer Vision Astrophysics + Cosmology**', specifically for intensity mapping of large-scale backgrounds:

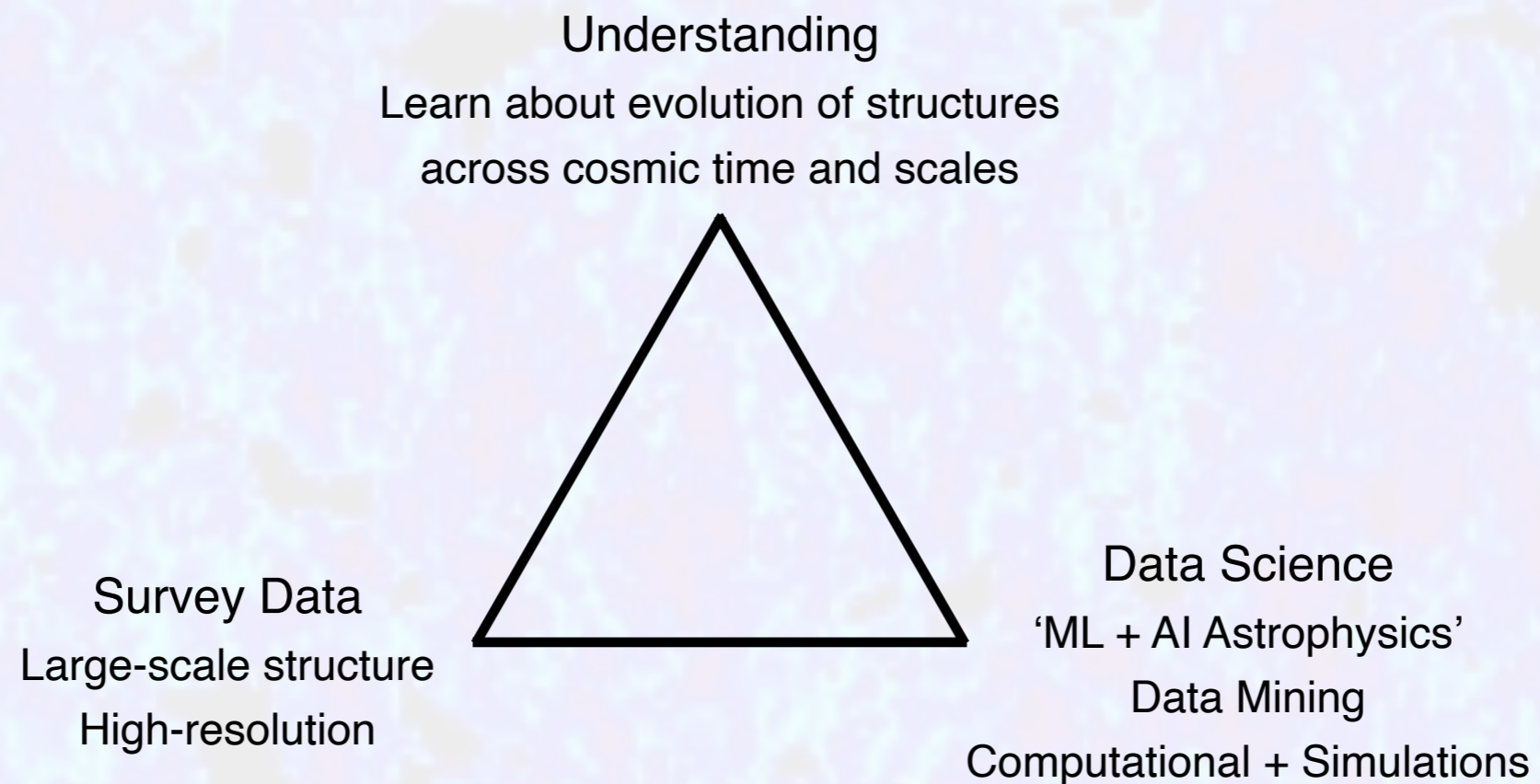
- Emulation, generation
- Inference

Also:

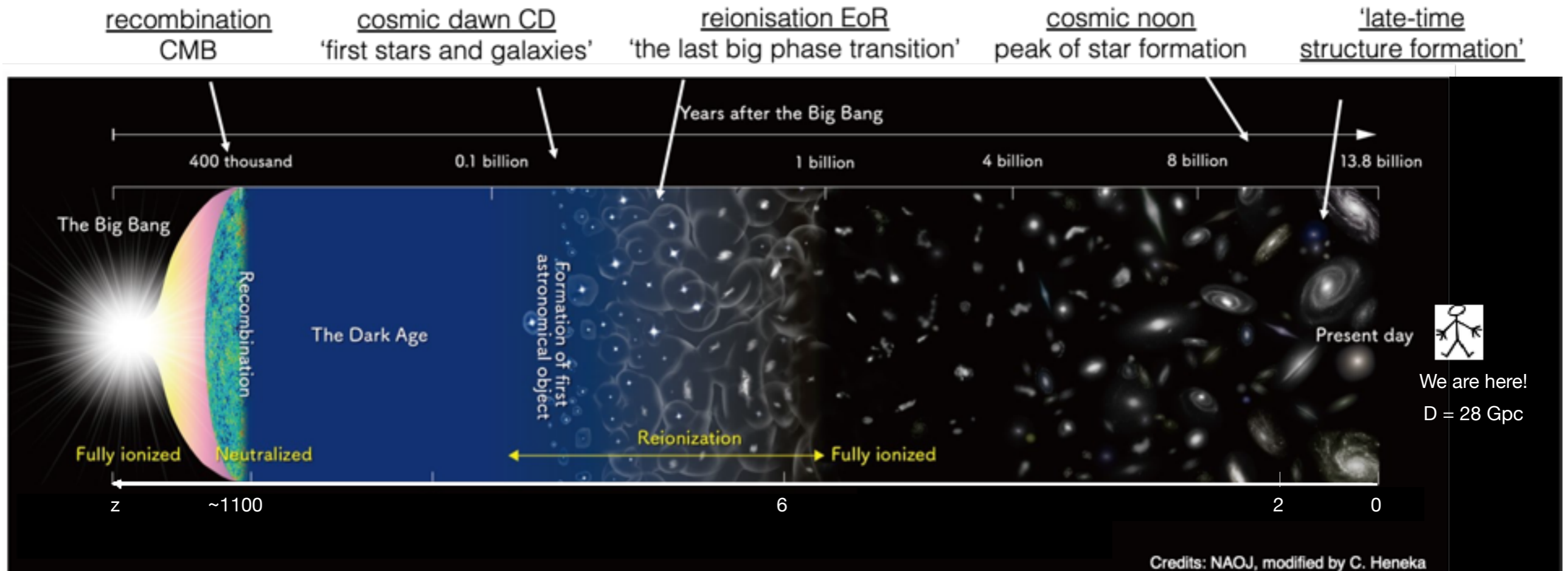
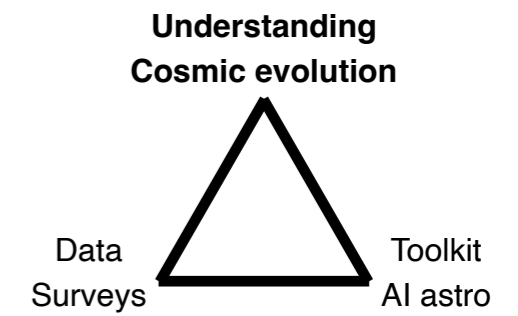
- Classification (e.g. 4MOST spectra)
- Detection (SKA preparation)



How will Astronomy and Astrophysics advance in coming years?



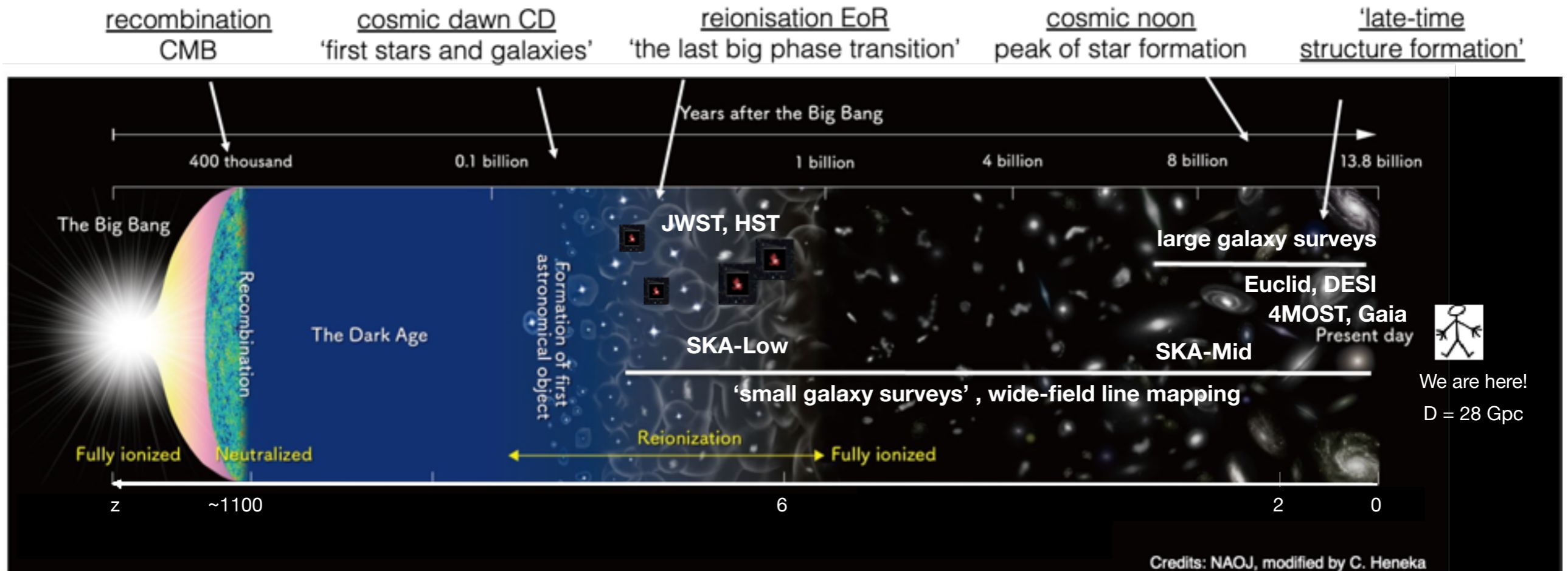
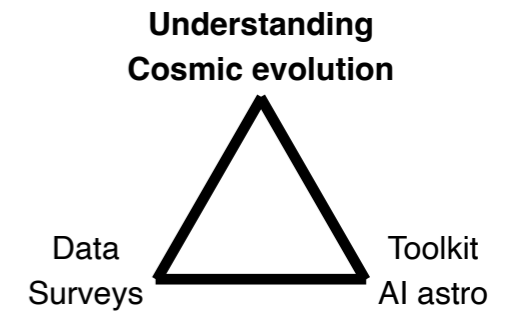
Where we stand: Data revolution and cosmic evolution



Our goal:

Learn about astrophysical & cosmological evolution
across cosmic time and scales

Where we stand: Data revolution and cosmic evolution



Our goal:

Learn about astrophysical & cosmological evolution
across cosmic time and scales

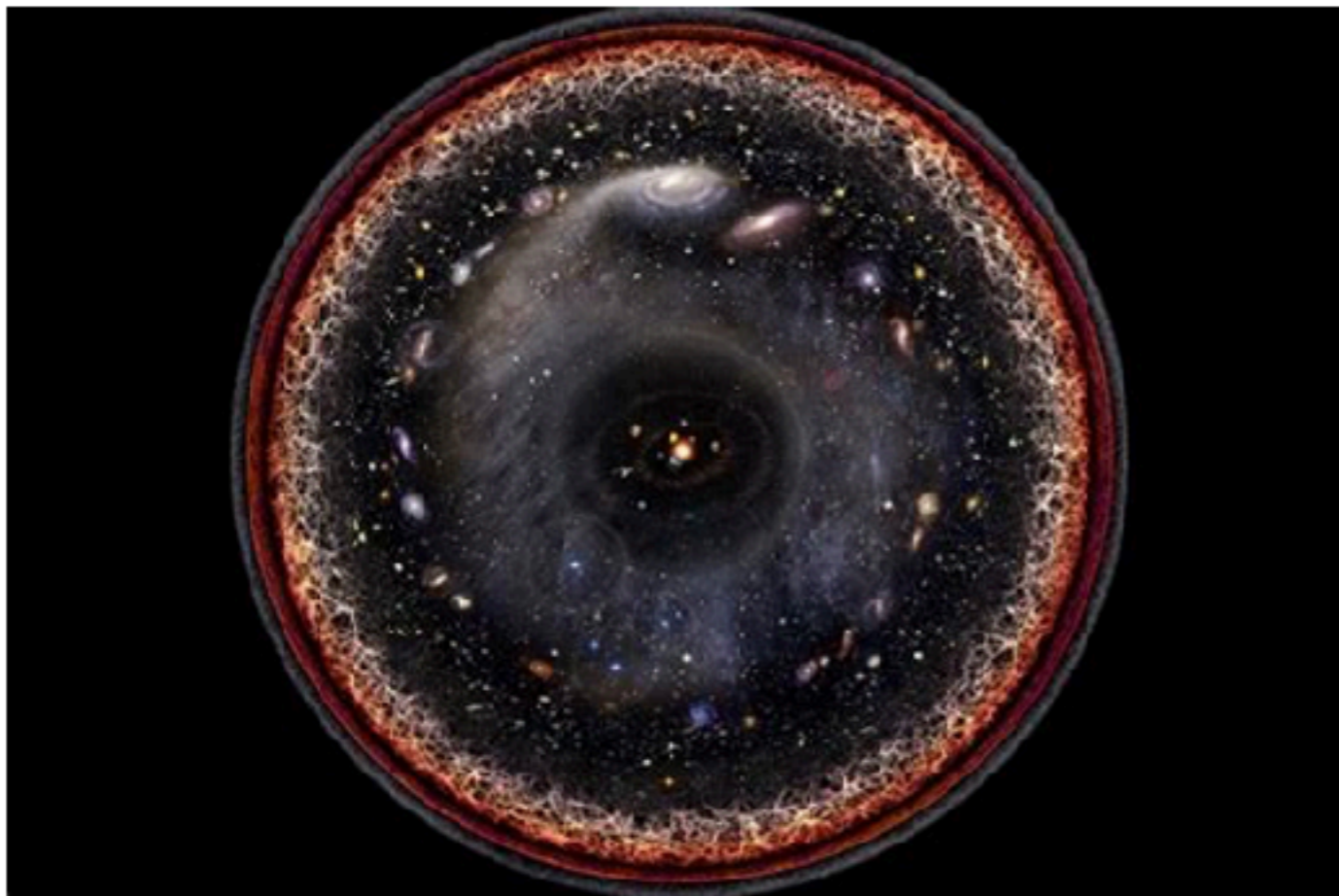
Coming decade: push to map up to **80% of the observable Universe**

... what does 80% of the observable Universe mean?

Modelling challenges

True LSS probes \longrightarrow orders of magnitude of scales up to the ultra-large

...what does 80% of the observable Universe even mean?



APOD, NASA, License & Credit: Wikipedia, Pablo Carlos Budassi

Observable Universe:

$d \sim 28$ Gpc (x3 Glyr)

80% if this:

$d \sim 22$ Gpc

Let's say we resolve (only) \sim Mpc

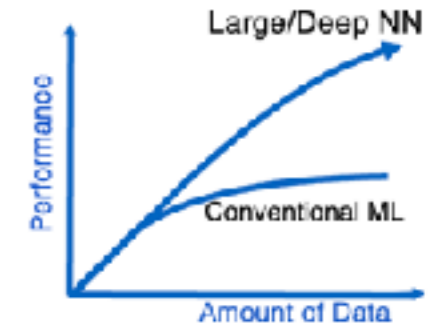
\longrightarrow about 3-4 orders of magnitude

\longrightarrow about 10^9 - 10^{10} modes!

... at some point we sub-grid model and/or change modelling approach

Advances in Data Science, for Science

Deep Learning
Driven by ability to improve
with large datasets



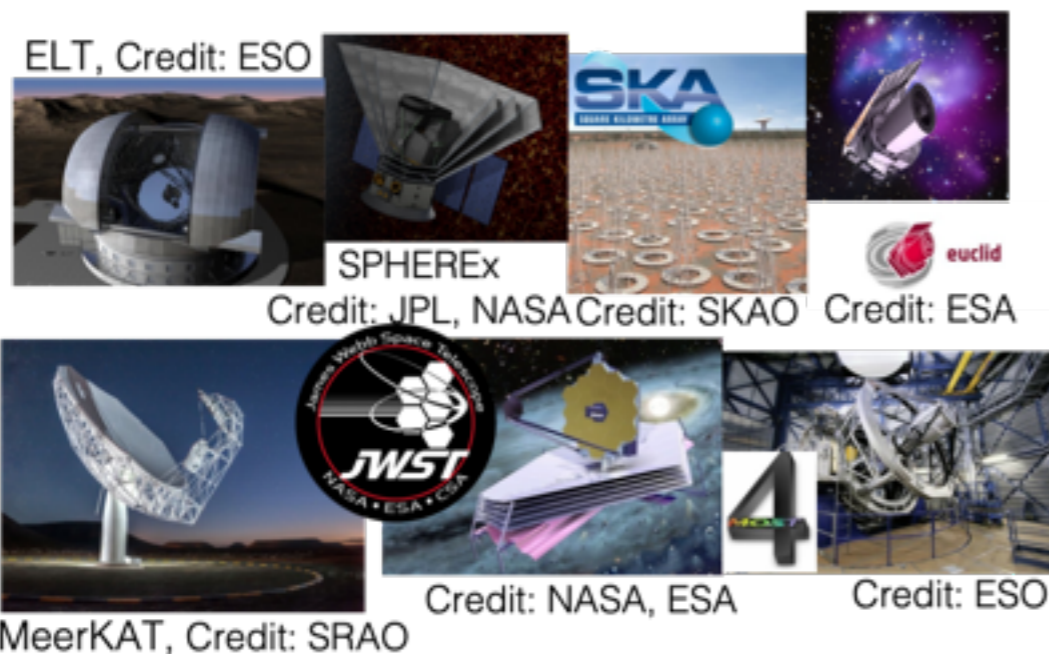
One big goal? Learn about astrophysical & cosmological evolution
across cosmic time and scales

No, really - many goals, ...

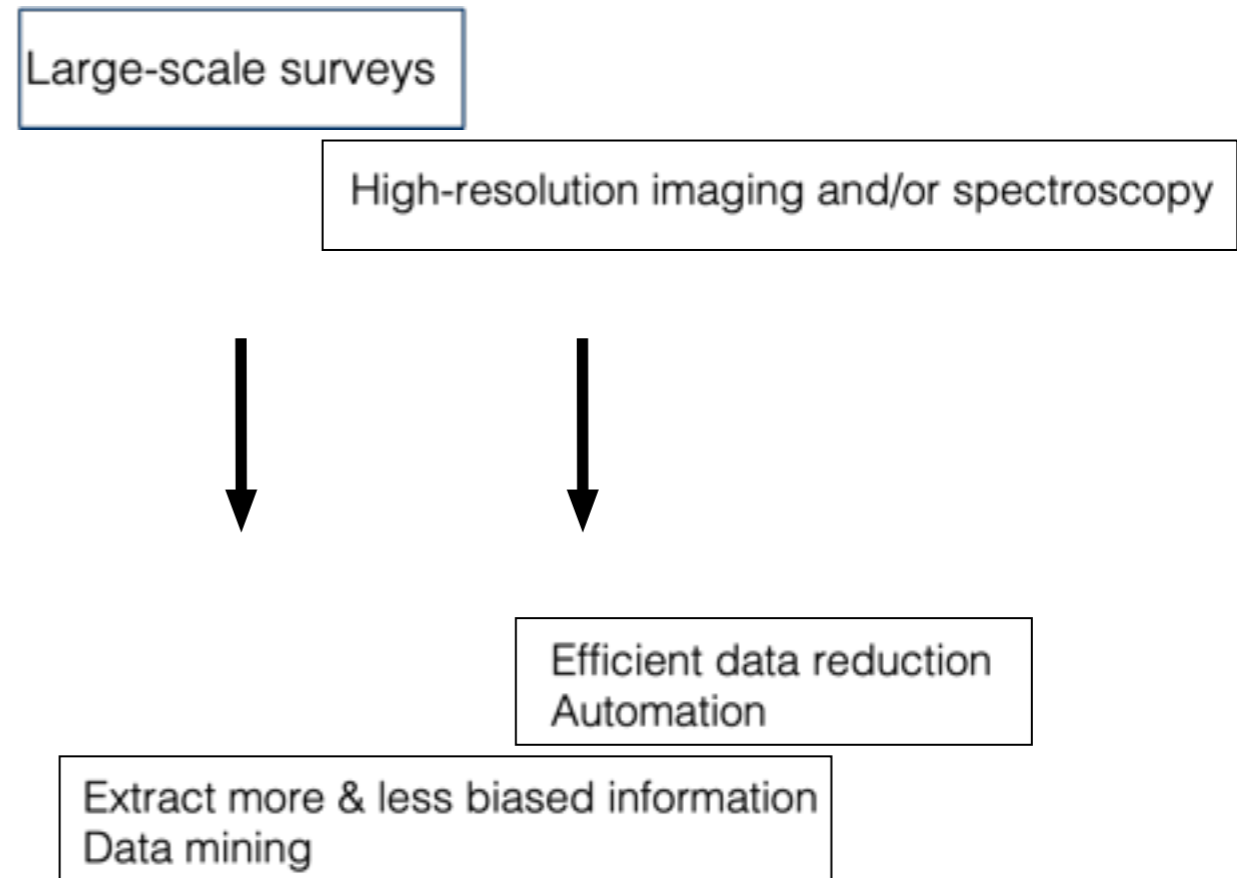
... and instruments

... types of data

... datarates, scales, signal-to-noise



** non-comprehensive list*



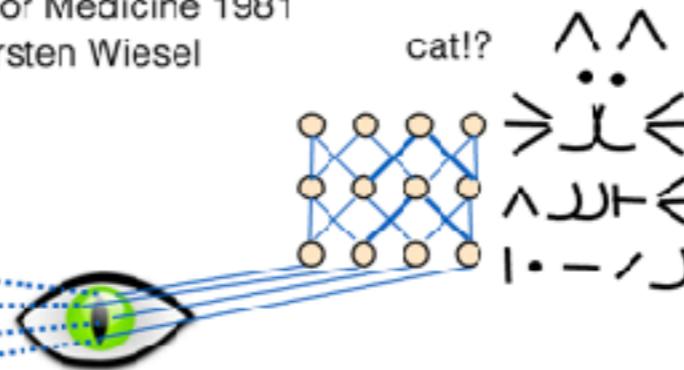
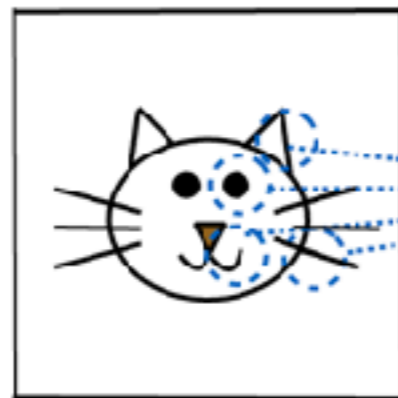
Why ML/DL for Astrophysics?



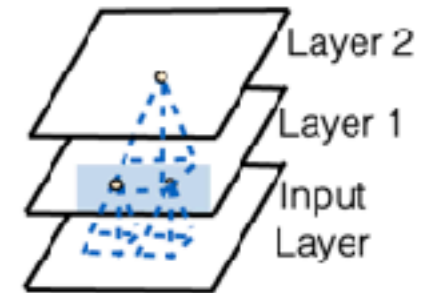
<https://creintour.harvard.edu/archives/portfolio-items/the-first-neurobiology-department>

Neuron response to visual stimuli

Nobel Prize in Physiology or Medicine 1981
David Hubel and Torsten Wiesel



Hyper-complex
Complex cells
Simple cells



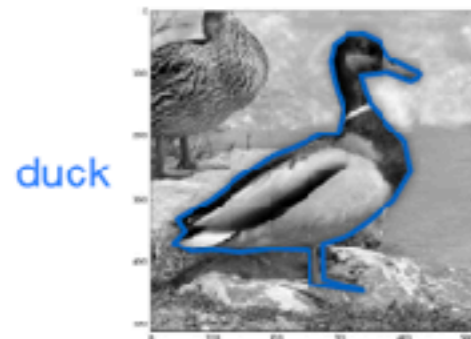
Convolutional neural network
'Neccogniton' Kunihiro Fukushima 1979

Representation learning

Hierarchical learning

Non-linear, non-Gaussian

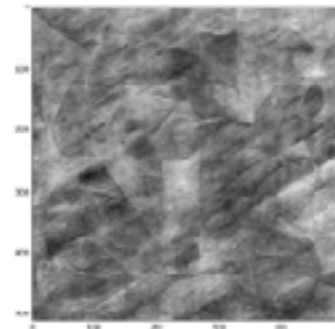
High-dim. correlations



size
color
bkg



randomise
phases



vs. "The famous
Gaussian duck"

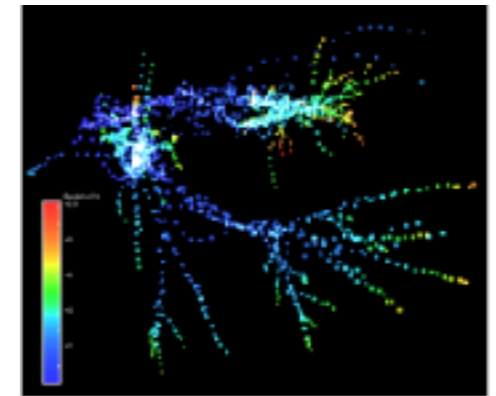
Same 2D
power spectrum

Where we stand: ML and AI for Astronomy and Astrophysics

We need a versatile ML/AI toolkit, for:

- detection, segmentation
- classification
- regression, inference
- anomaly detection
- generation, emulation, ...

Hierarchical

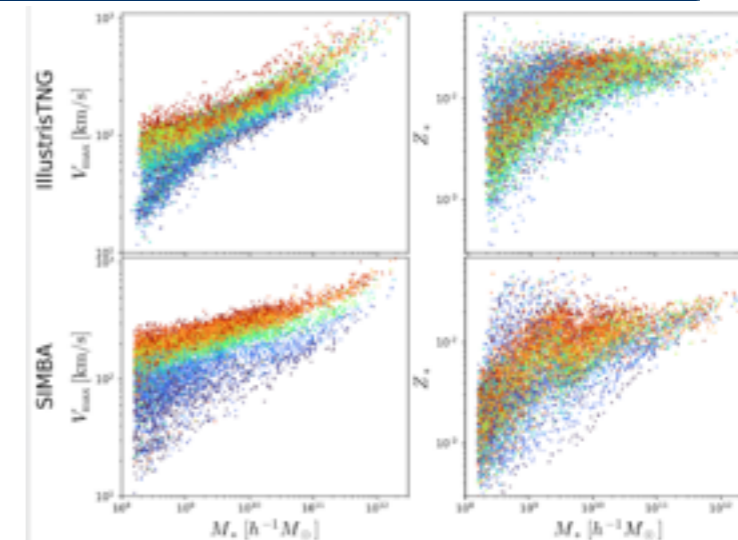


@Chris Fluke, Swinburne University of Technology

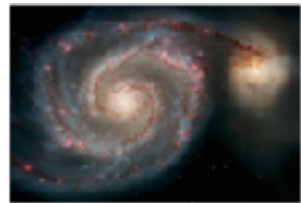
+ plenty of inverse problems

$$I^D(x, y) = R \times I(x, y) + n$$

High-dim. correlations



arXiv:2201.02202

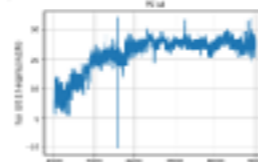
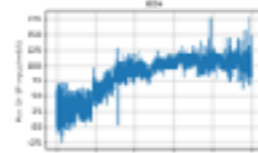


@Hubble, NASA

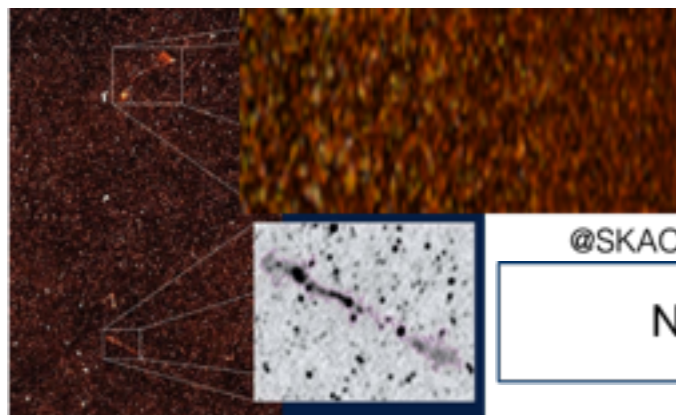


@Hubble, NASA

Representation learning



@SDSS



@SKAO

Non-linear, non-Gaussian

@MeerKAT

Where we stand: ML and AI for Astronomy and Astrophysics

We need a versatile ML/AI toolkit, for:

- detection, segmentation
- classification
- regression, inference
- anomaly detection
- generation, emulation, ...

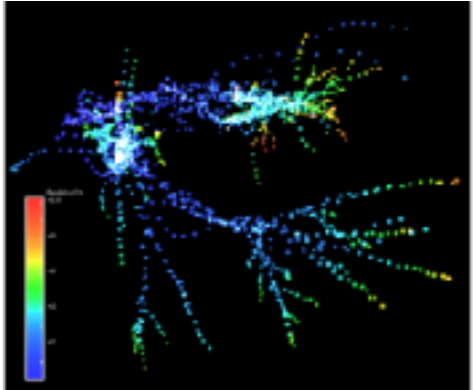
Examples:

Few sec: Classification 40.000 spectra

Few sec: 7-parameter inference ~100MB cube

Few sec: detection, segmentation & flux measurement O(100-1000) sources

Hierarchical

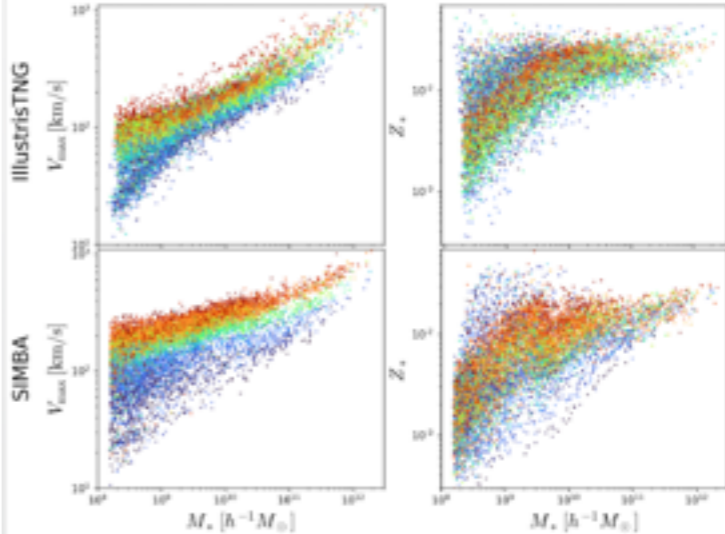


@Chris Fluke, Swinburne University of Technology

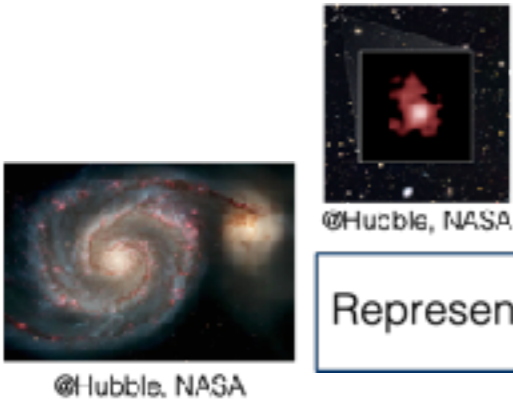
+ plenty of inverse problems

$$I^D(x, y) = R \times I(x, y) + n$$

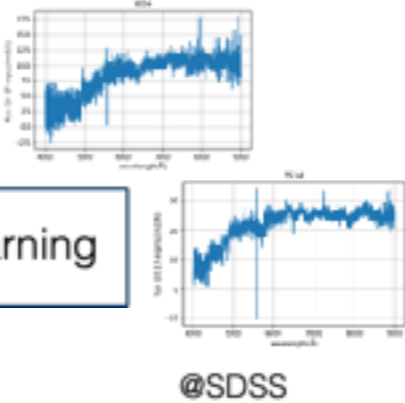
High-dim. correlations



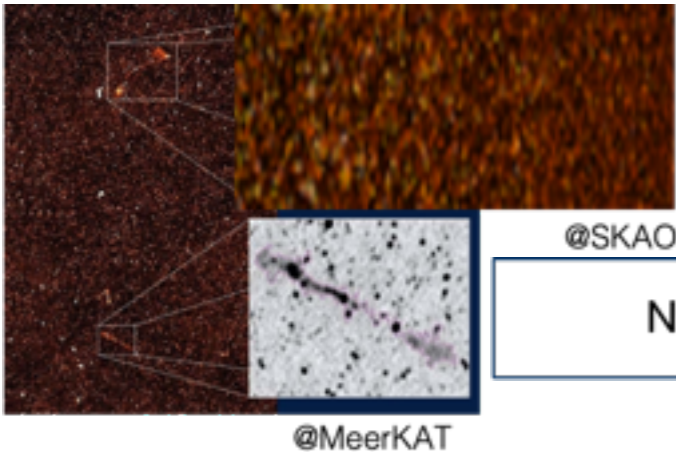
arXiv:2201.02202



Representation learning



@SDSS

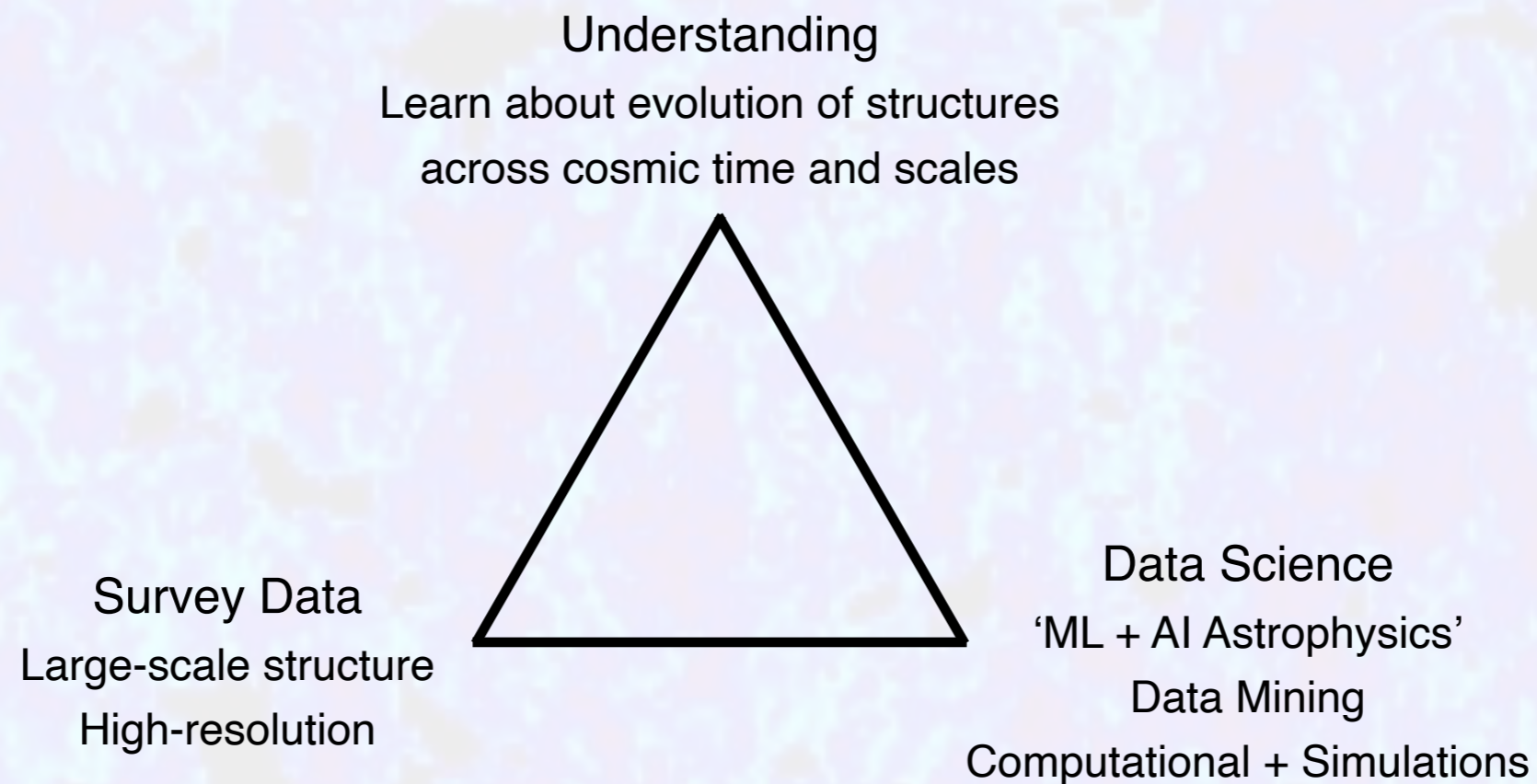


Non-linear, non-Gaussian

@MeerKAT

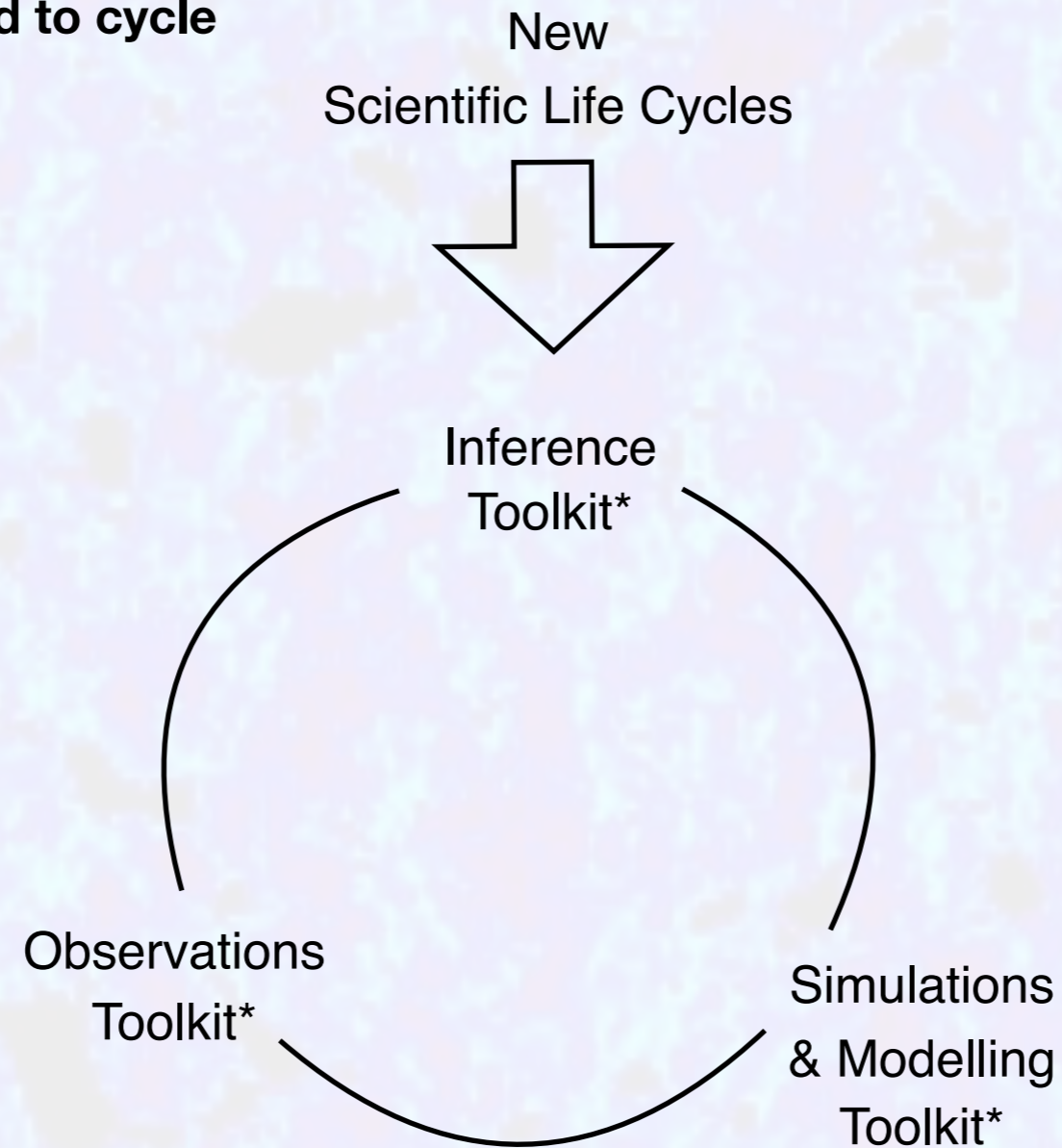
@SKAO

How will Astronomy and Astrophysics advance in coming years?



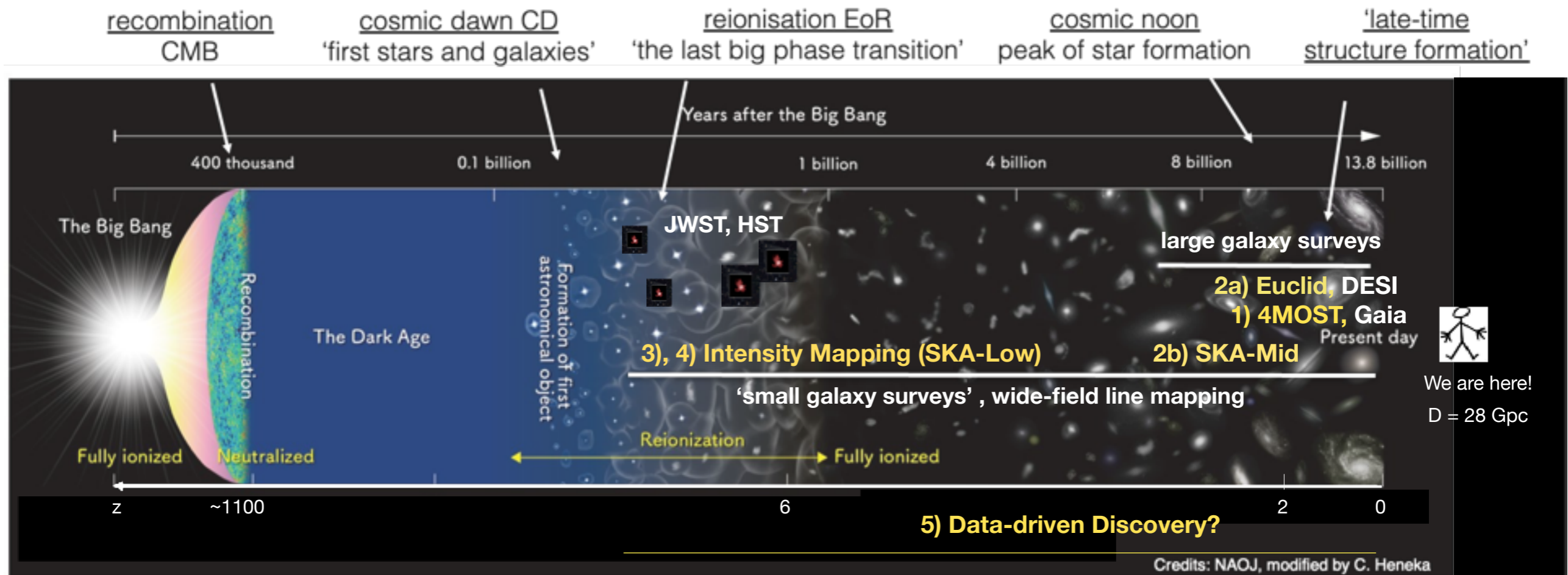
Astronomical and Astrophysical Machine Learning

In practice: From pyramid to cycle



*The Toolkit: Statistics, Machine & Deep Learning, Data Mining

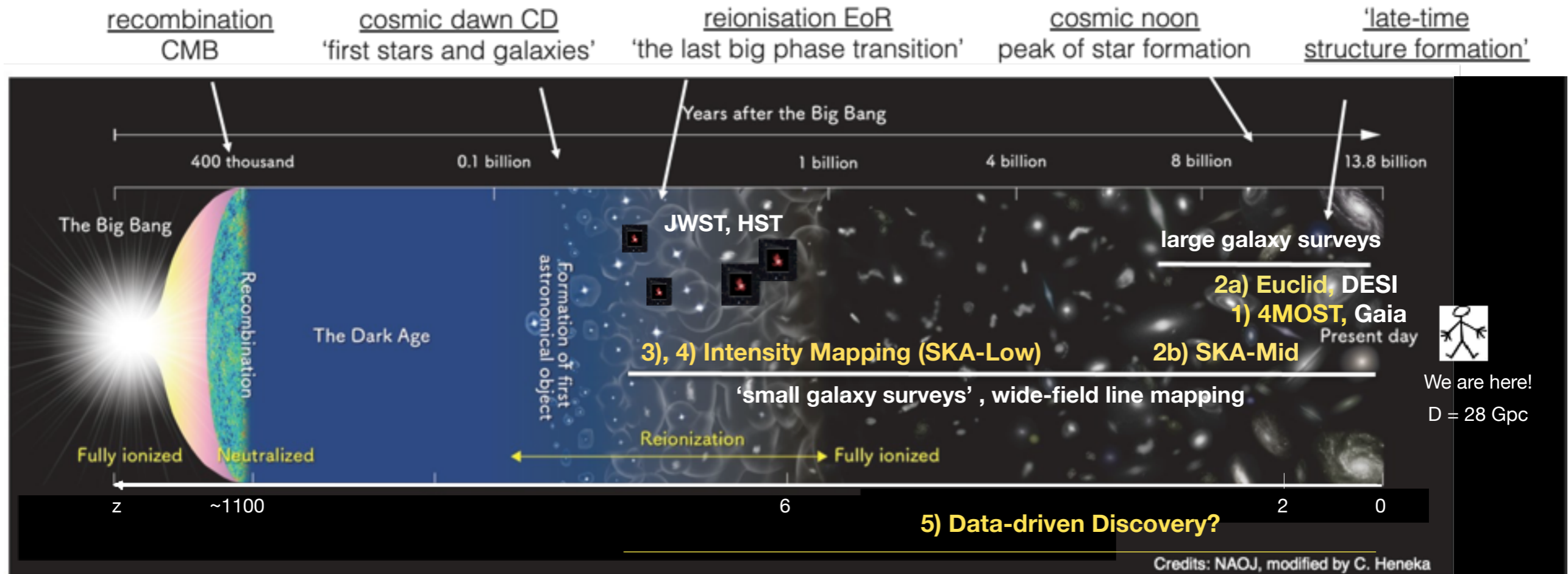
Astronomical and Astrophysical Machine Learning



Highlights in this Lecture

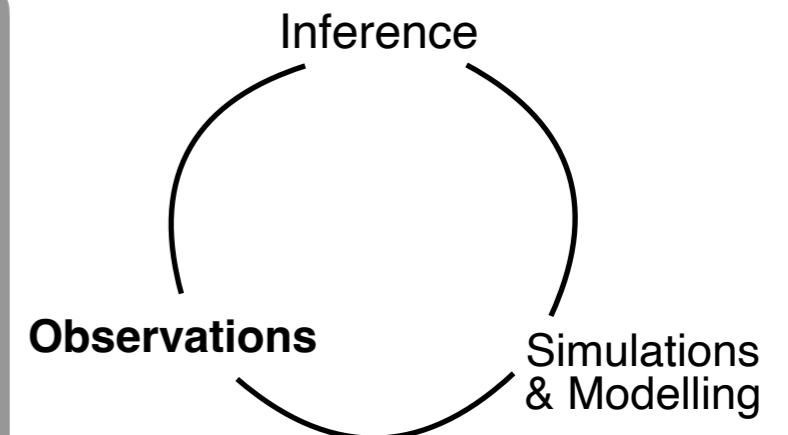
- 1) Classification / Triggering
- 2) Source detection & characterisation
- 3) Simulation-based inference (SBI) in 3D
- 4) Generative methods
- 5) Data-driven Discovery

Astronomical and Astrophysical Machine Learning



Highlights in this Lecture

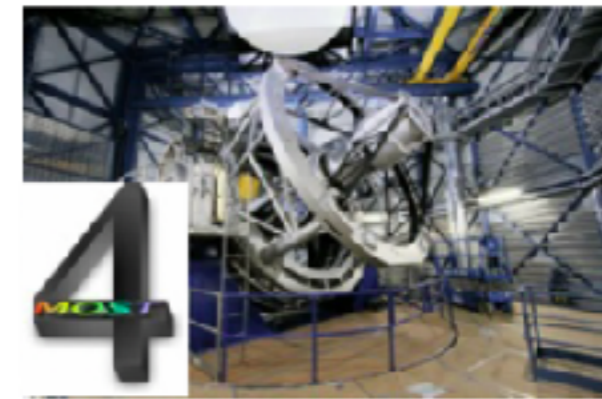
- 1) Classification / Triggering
- 2) Source detection & characterisation
- 3) Simulation-based inference (SBI) in 3D
- 4) Generative methods
- 5) Data-driven Discovery



1) Classification and triggering for large astronomical surveys

4MOST: On-the-fly classification of spectra (1D)

- 5-year survey
- wide-field, fibre-fed, optical spectroscopy
- on ESO's 4-m-class telescope VISTA
- 2.5-degree diameter field-of-view, 2436 fibres
- HRS $R \approx 18000 - 21000$, LRS $R \approx 4000 - 7500$
- 20mio. (LRS), 3mio. (HRS) sources



<https://www.4most.eu> Credit: ESO

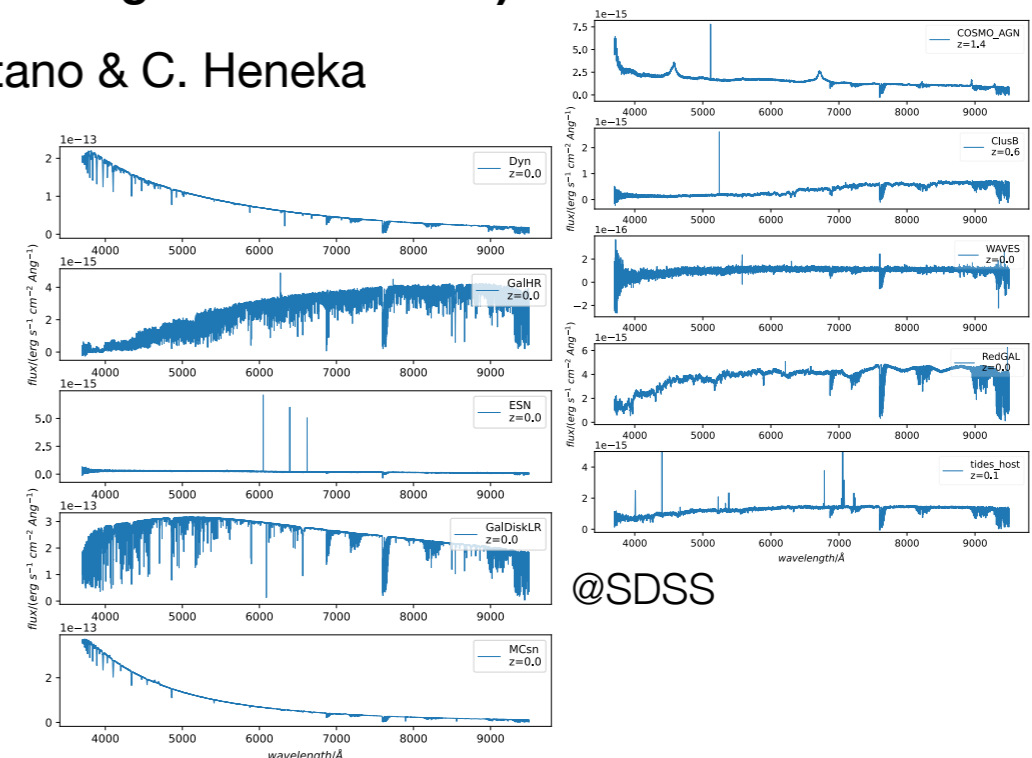
Goal: Data-driven classification pipeline layer (galactic & extragalactic sources)

Classification infrastructure working group, led by: N. Napolitano & C. Heneka



*Benchmark with
SDSS archival spectra:*

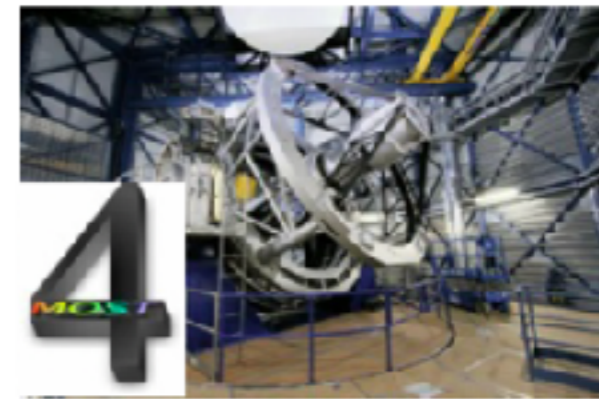
See also our tutorial!



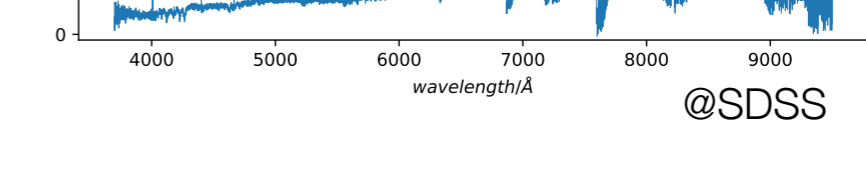
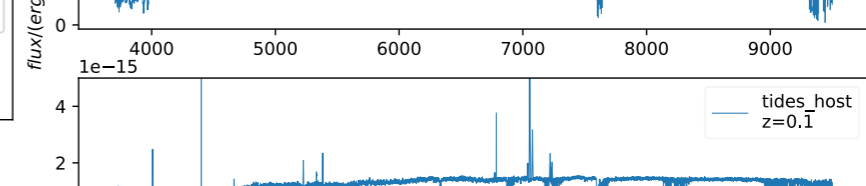
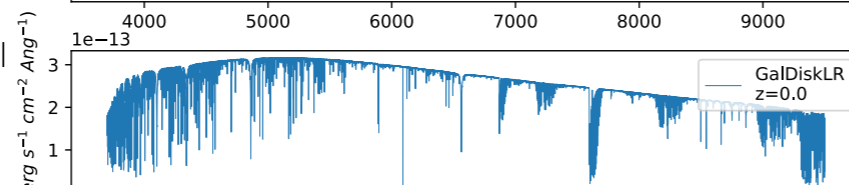
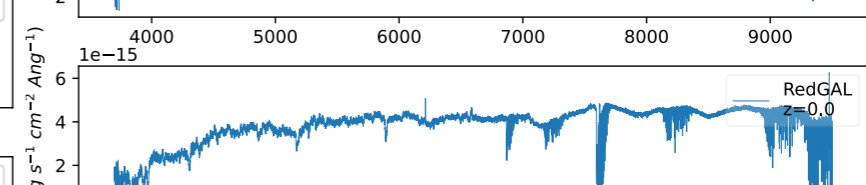
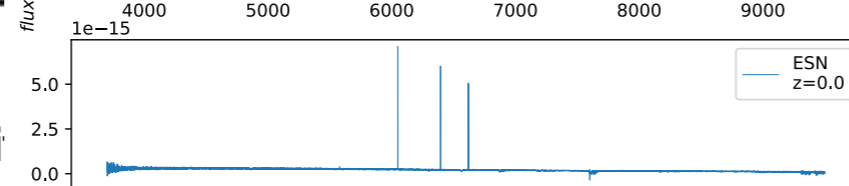
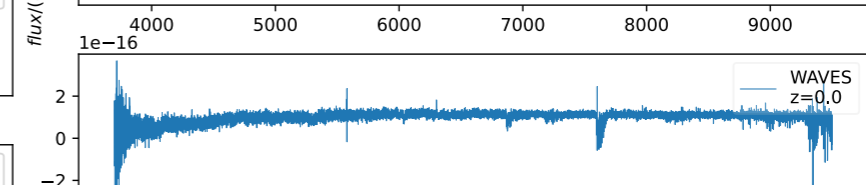
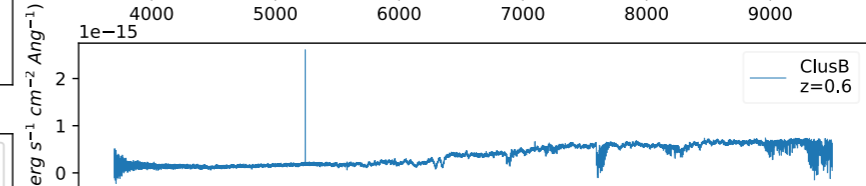
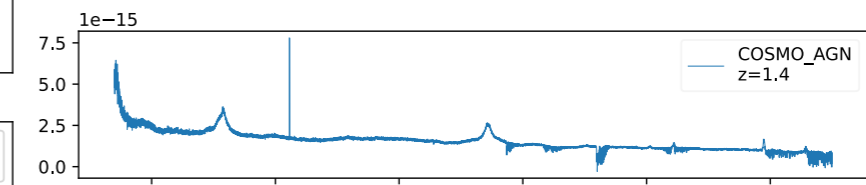
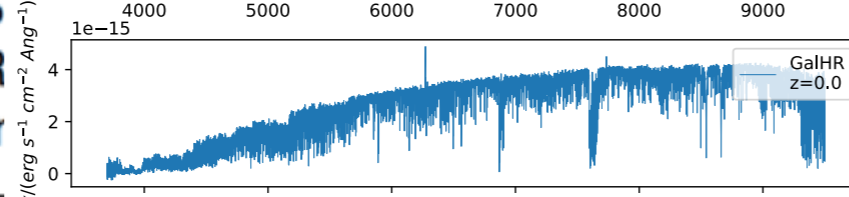
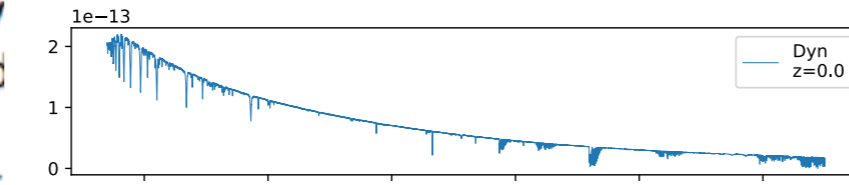
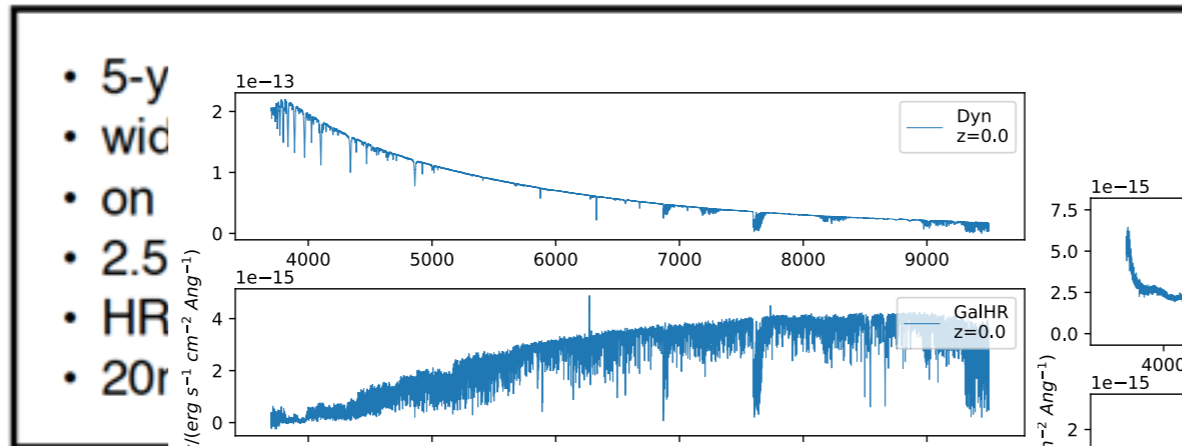
@SDSS

1) Classification and triggering for large astronomical surveys

4MOST: On-the-fly classification of spectra (1D)



<https://www.4most.eu> Credit: ESO



Goal: Data-driven classification

Classification infrastructure

*Galactic vs. Extragalactic
See also our tutorial!*

@SDSS

1) Classification and triggering for large astronomical surveys



<https://www.4most.eu> Credit: ESO

4MOST: On-the-fly classification of spectra (1D)

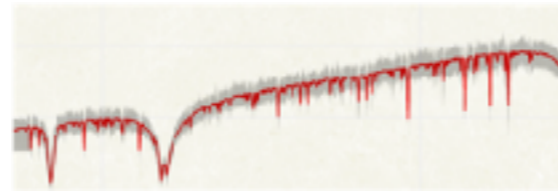
Goal: Data-driven classification pipeline layer (galactic & extragalactic sources)

Classification infrastructure working group, led by: N. Napolitano & C. Heneka

→ Probabilistic multi-classifier

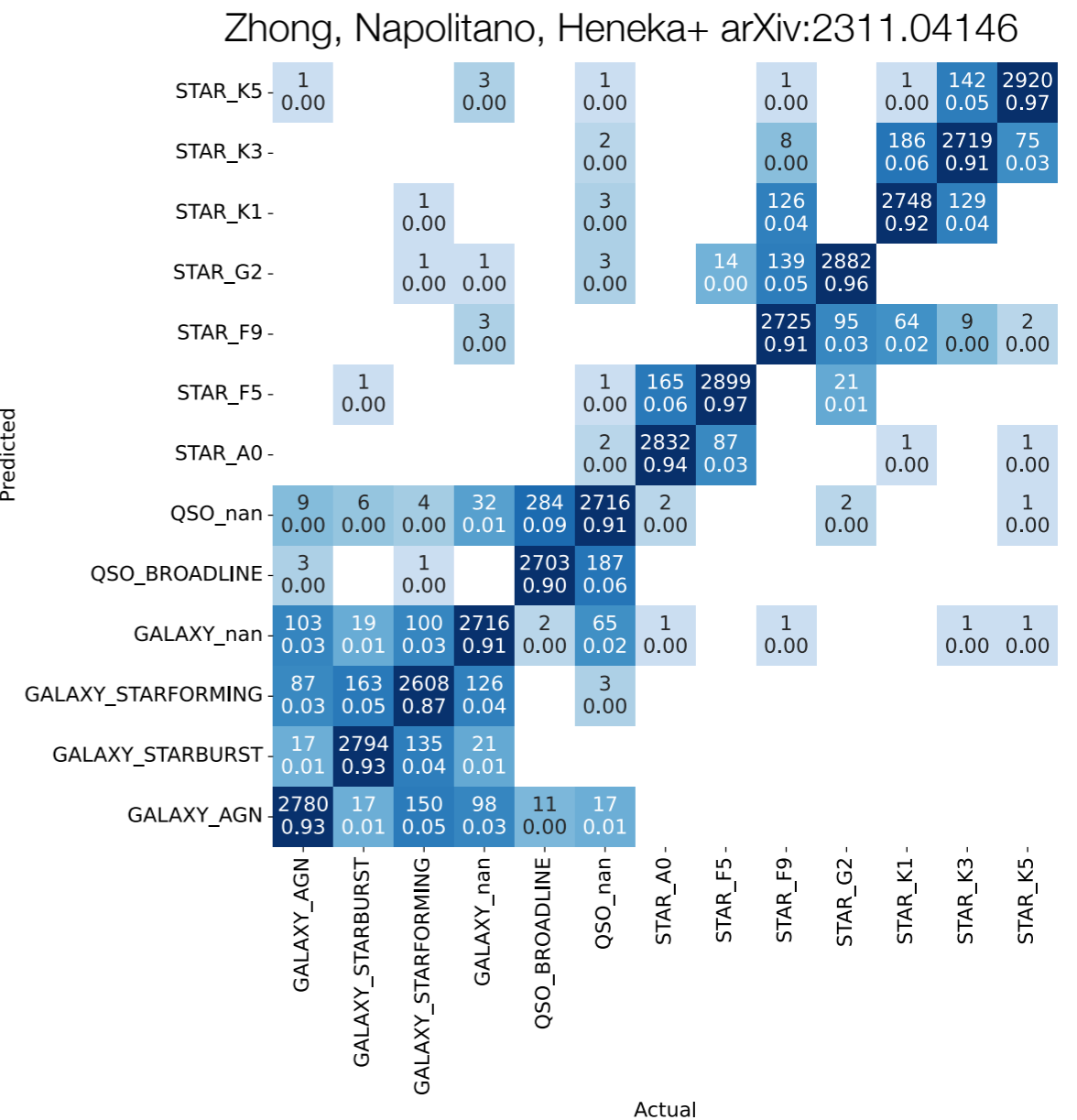
For class:
Convolutional
network variants

For class uncertainties:
Bayesian neural networks
and contrastive learning

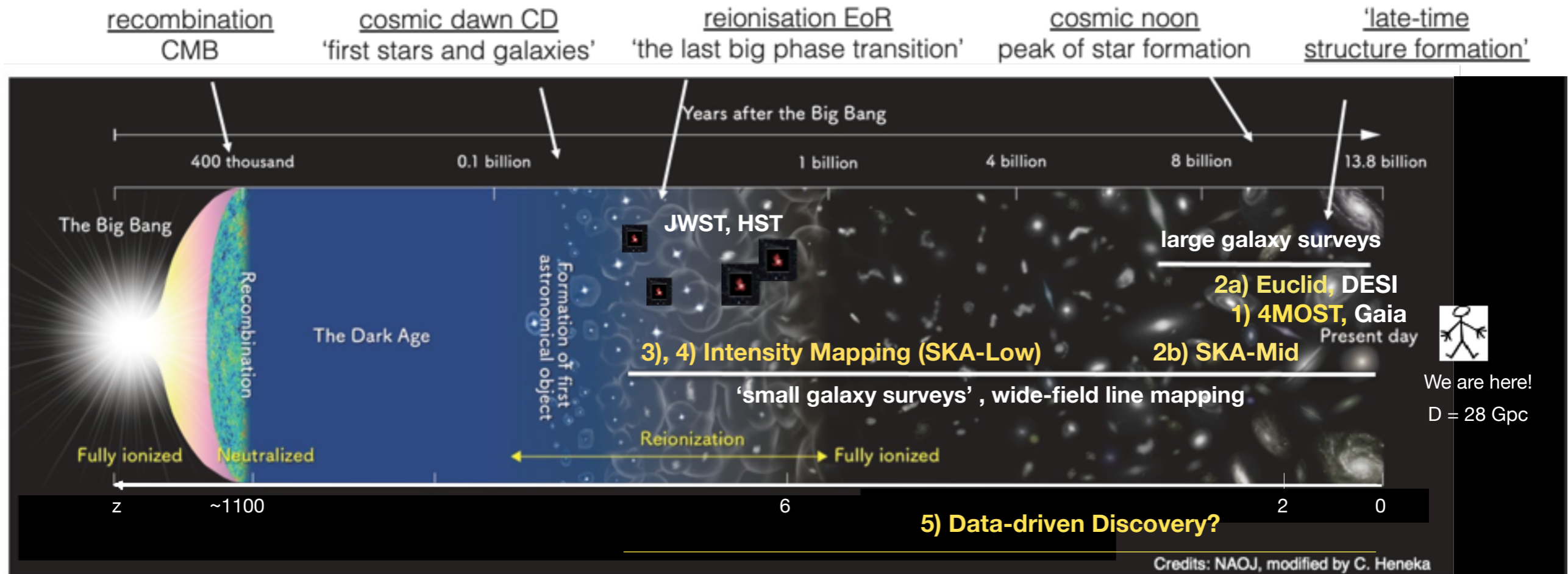


++ competitive with template fitting

Examples:
Few sec: Classification 40.000 spectra
 Few sec: 7-parameter inference ~100MB cube
 Few sec: detection, segmentation & flux measurement O(100-1000) sources

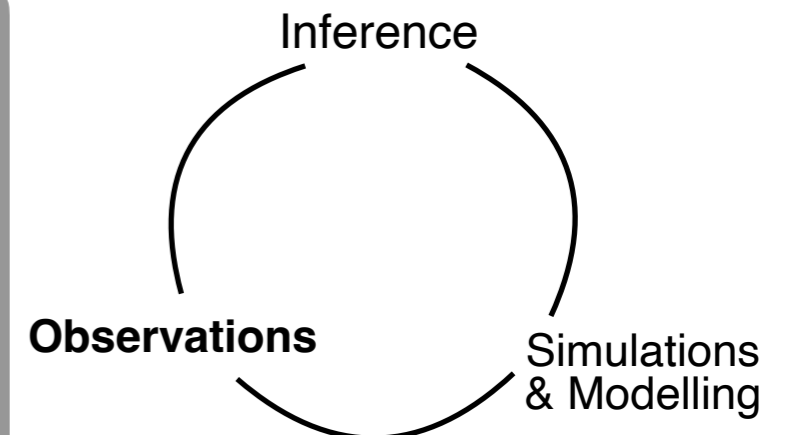


Astronomical and Astrophysical Machine Learning



Highlights in this Lecture

- 1) Classification / Triggering
- 2) **Source detection & characterisation**
- 3) Simulation-based inference (SBI) in 3D
- 4) Generative methods
- 5) Data-driven Discovery



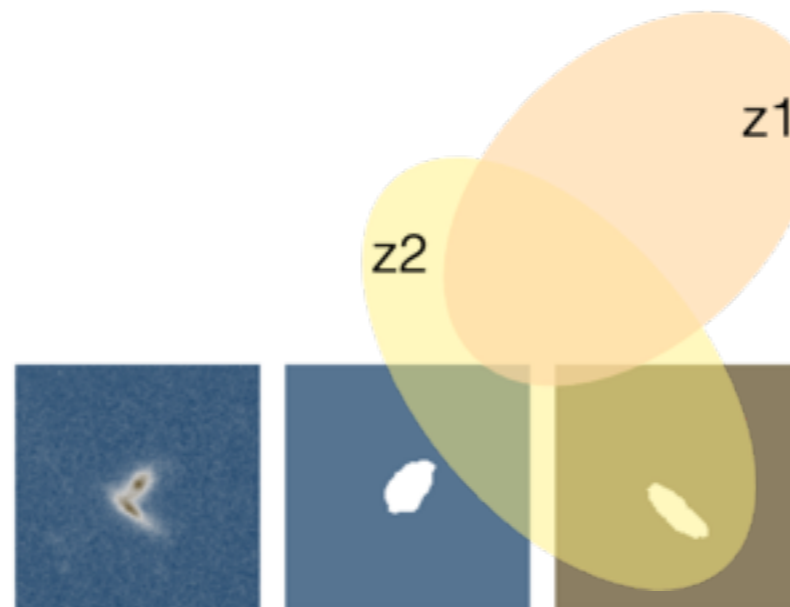
2a) The deblending problem

Example: Optical source detection & characterisation

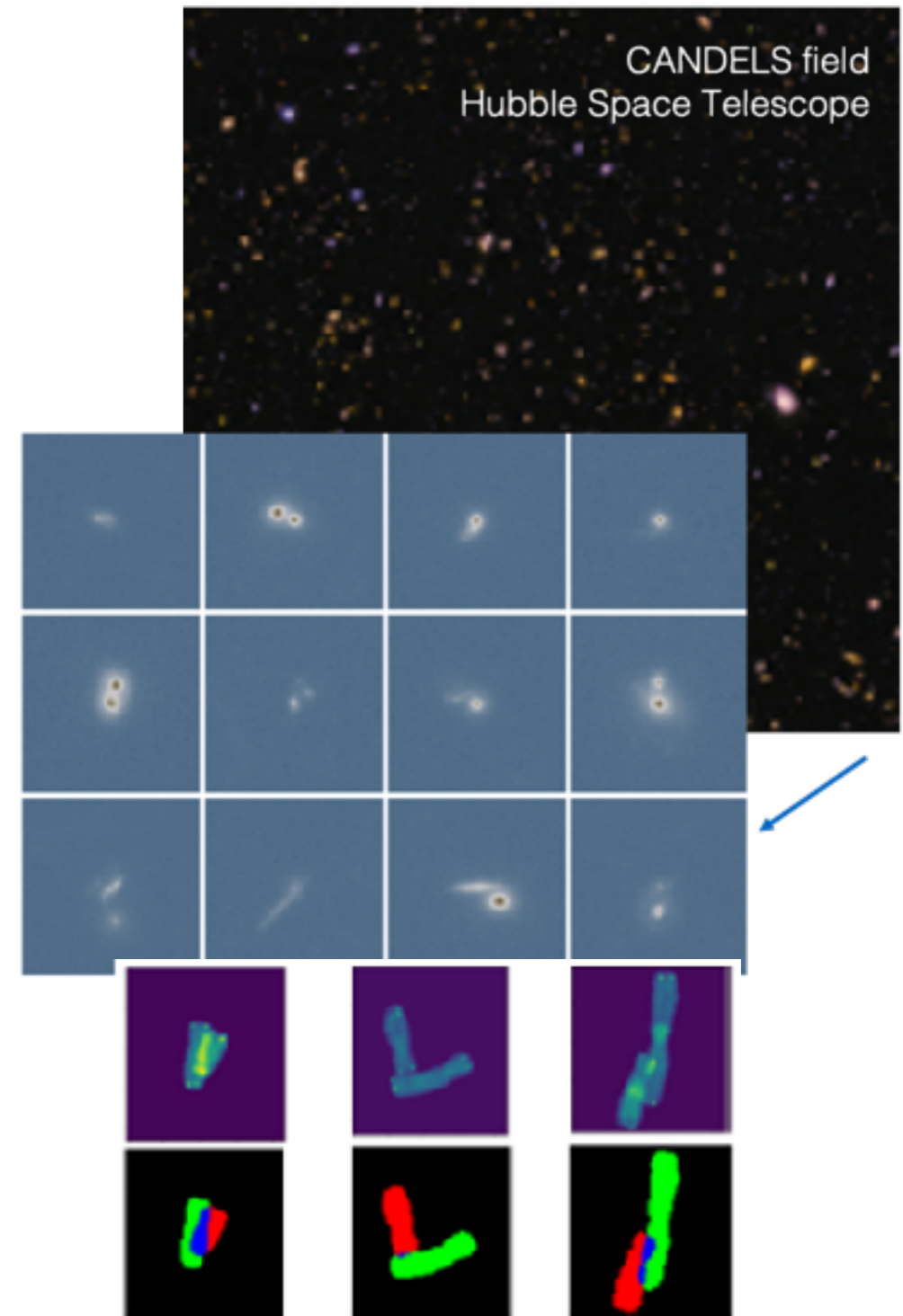
Goal: 'Good' photometry for surveys with high blended fraction
- avoid bias!

Galaxy morphology

Challenge: Galaxies are 'transparent'



Boucaud, Huertas-Company, Heneka+ 20,
arXiv:1905.01324



Lily Hu+ 2017

Similar challenge:
Overlapping chromosomes

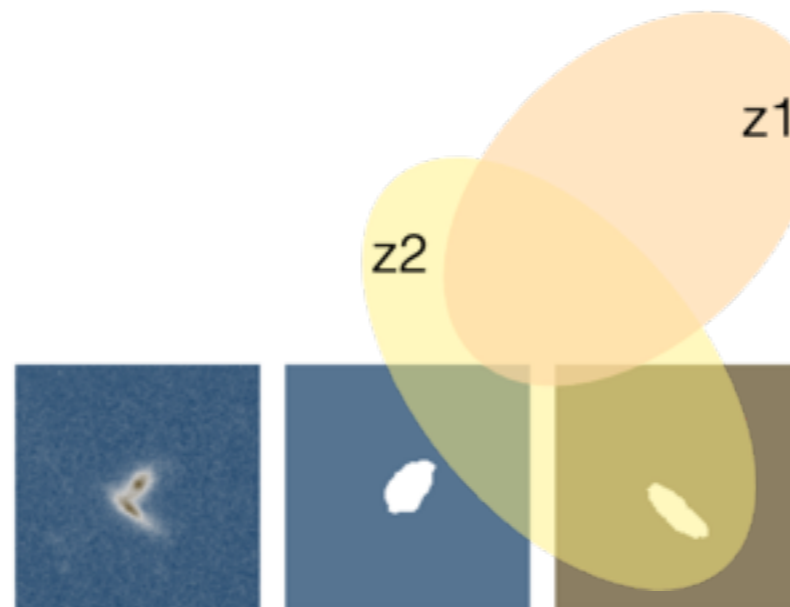
2a) The deblending problem

Example: Optical source detection & characterisation

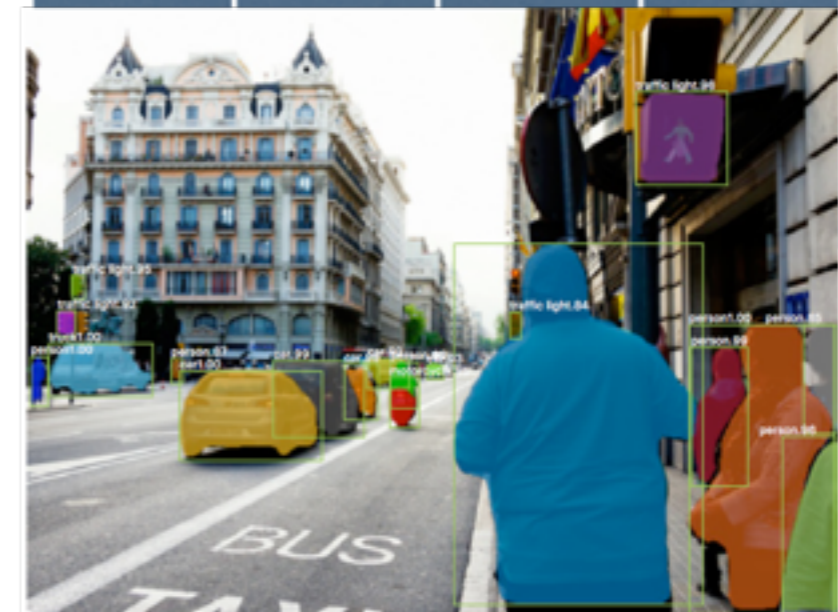
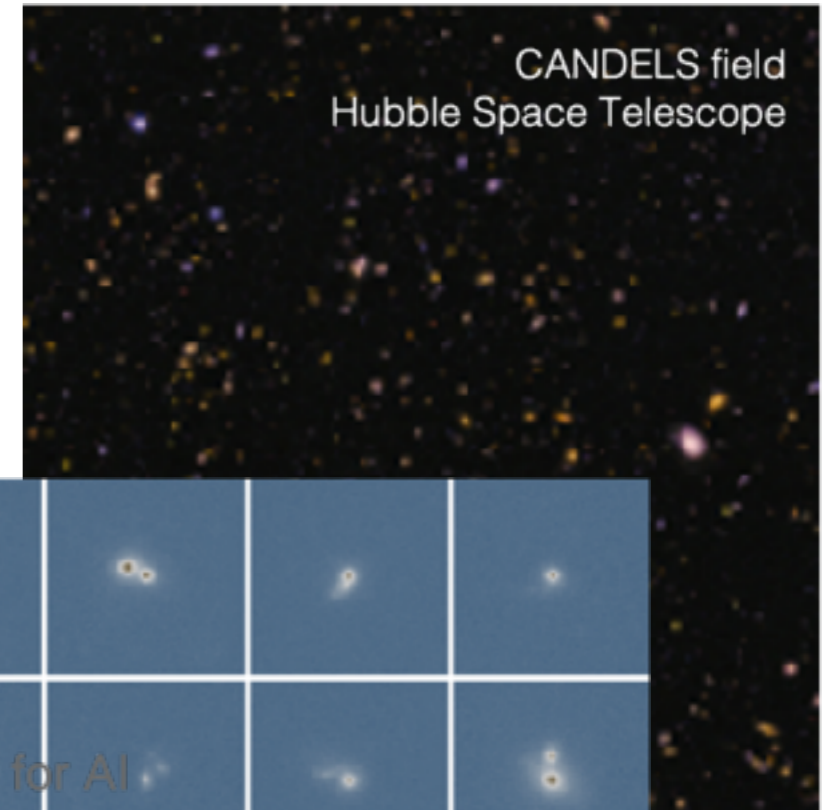
Goal: 'Good' photometry for surveys with high blended fraction
- avoid bias!

Galaxy morphology

Challenge: Galaxies are 'transparent'



Boucaud, Huertas-Company, Heneka+ 20,
arXiv:1905.01324



https://medium.com/@umerfarooq_26378/from-r-cnn-to-mask-r-cnn-d6367b196cfd

Another challenge: Object detection

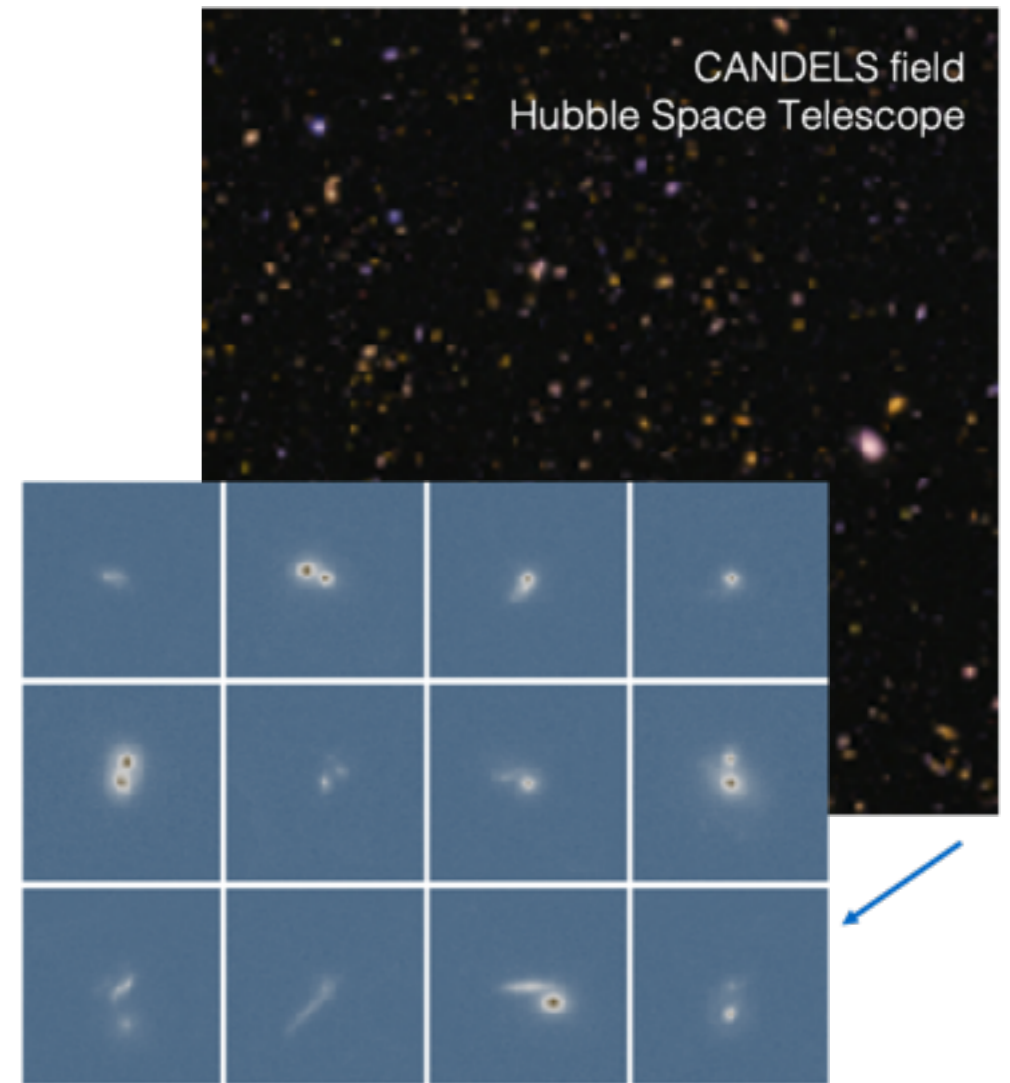
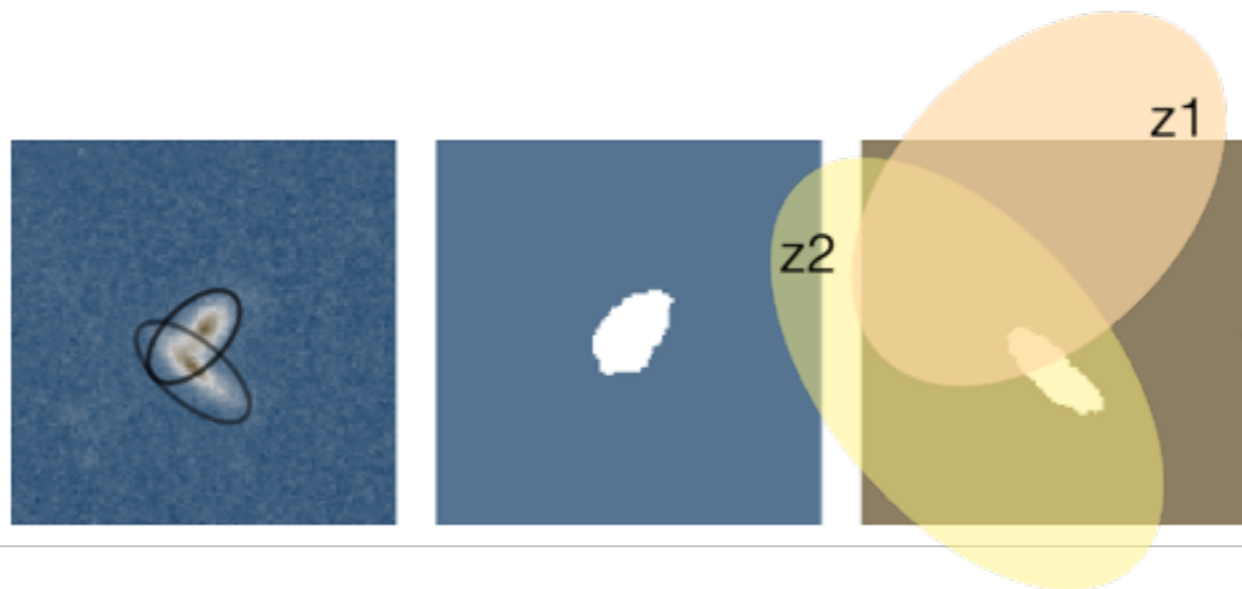
2a) The deblending problem

Example: Optical source detection & characterisation

Goal: 'Good' photometry for surveys with high blended fraction
- avoid bias!

Galaxy morphology

Challenge: Galaxies are 'transparent'



'Classic':
Fit ellipse(s)
and profile(s)

e.g. Einasto ('65):

$$\frac{d \log(\rho)}{d \log(r)} = -2 \left(\frac{r}{r_s} \right)^\alpha$$



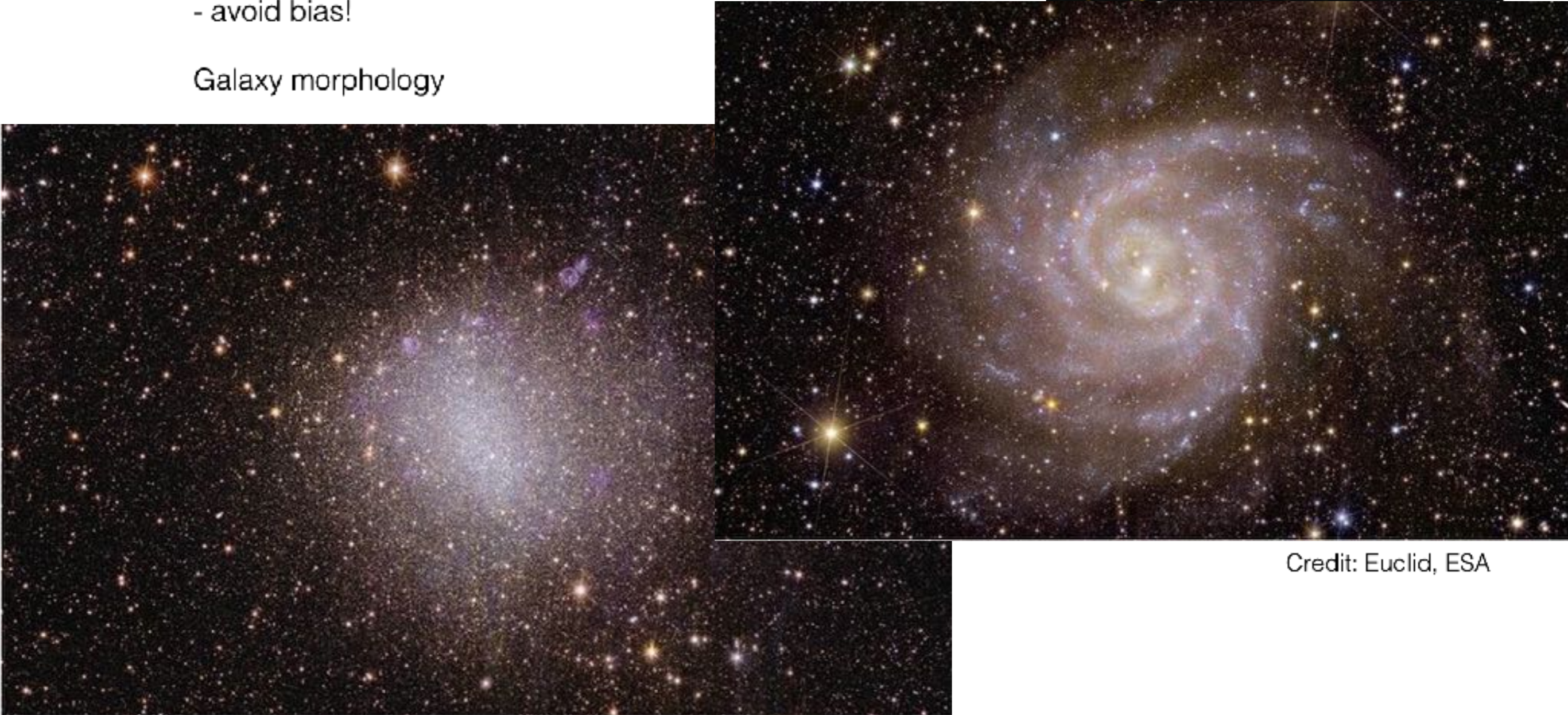
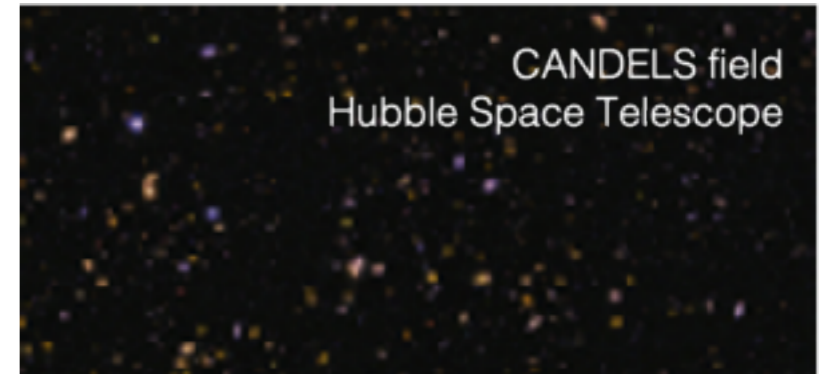
Boucaud, Huertas-Company, Heneka+ 20,
arXiv:1905.01324

2a) The deblending problem

Example: Optical source detection & characterisation

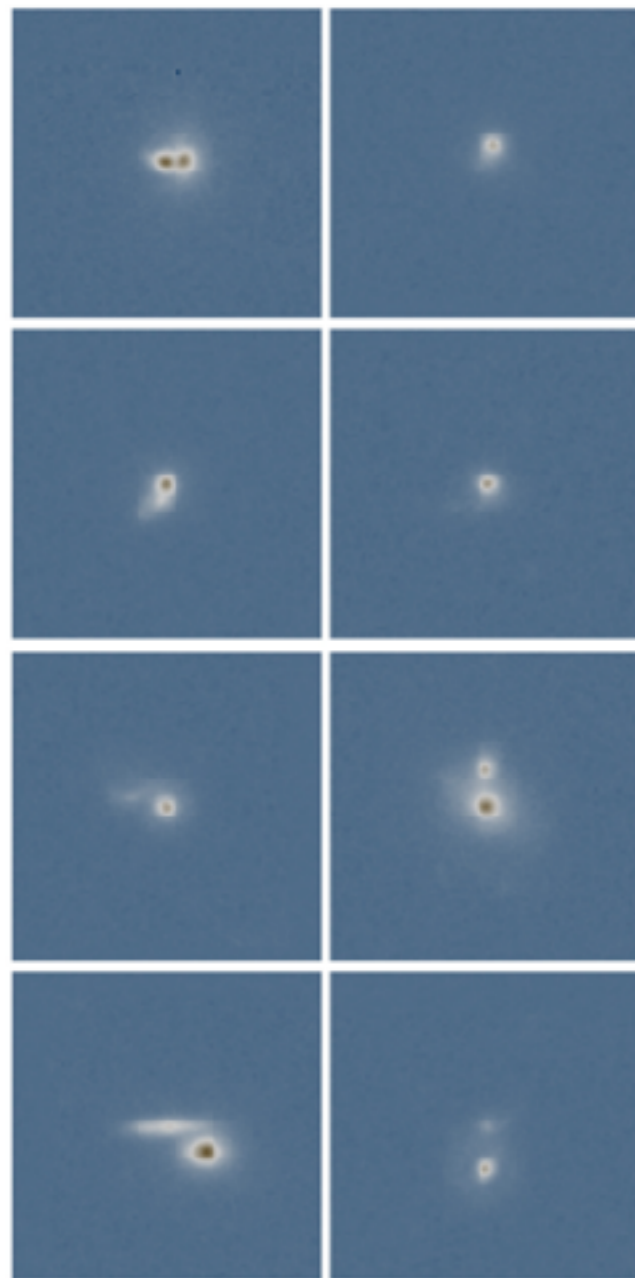
Goal: 'Good' photometry for surveys with high blended fraction
- avoid bias!

Galaxy morphology



2a) The deblending problem

Example: Optical source detection & characterisation

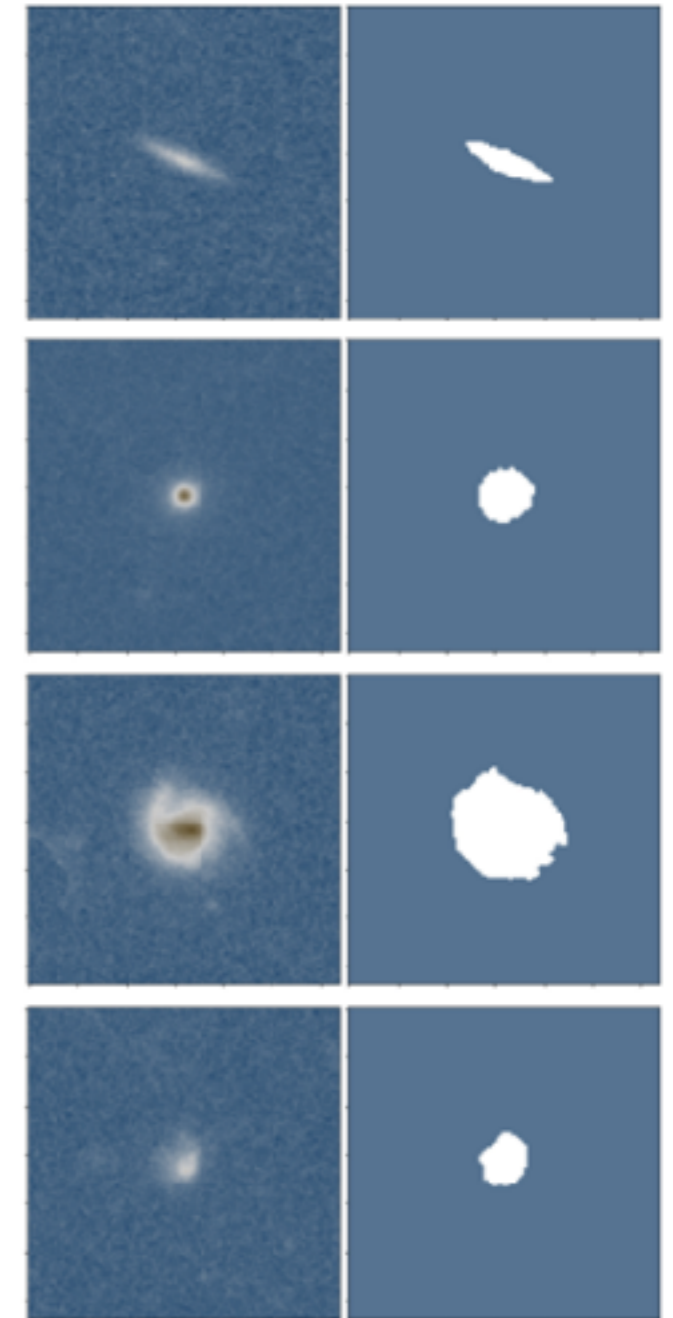


Get precise
Photometry

Derive masks

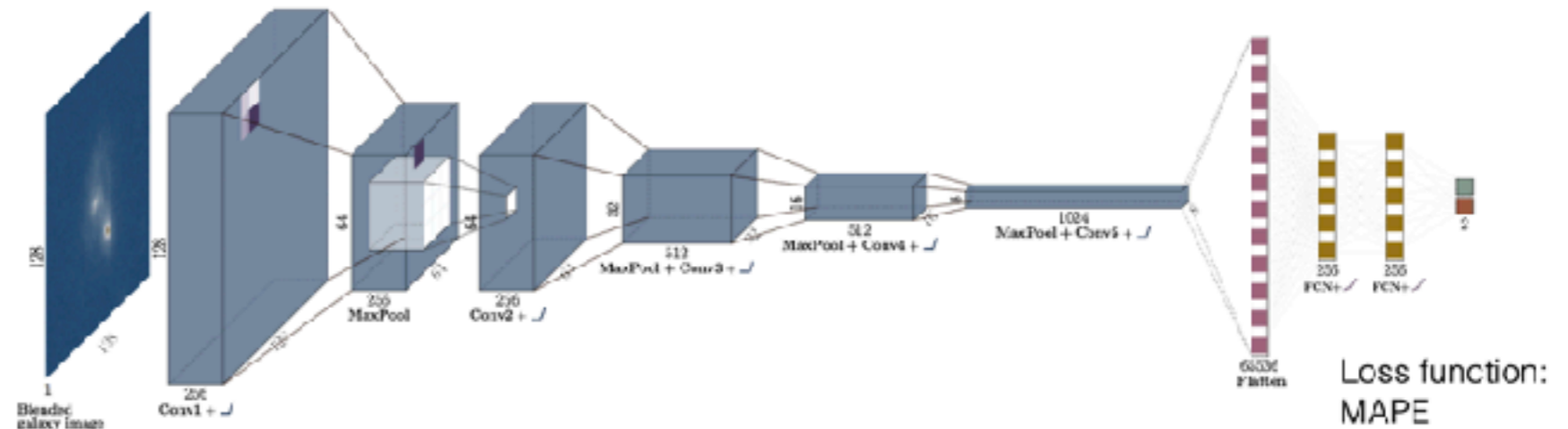
..and do so bias-free

CANDELS Field
Hubble Space Telescope

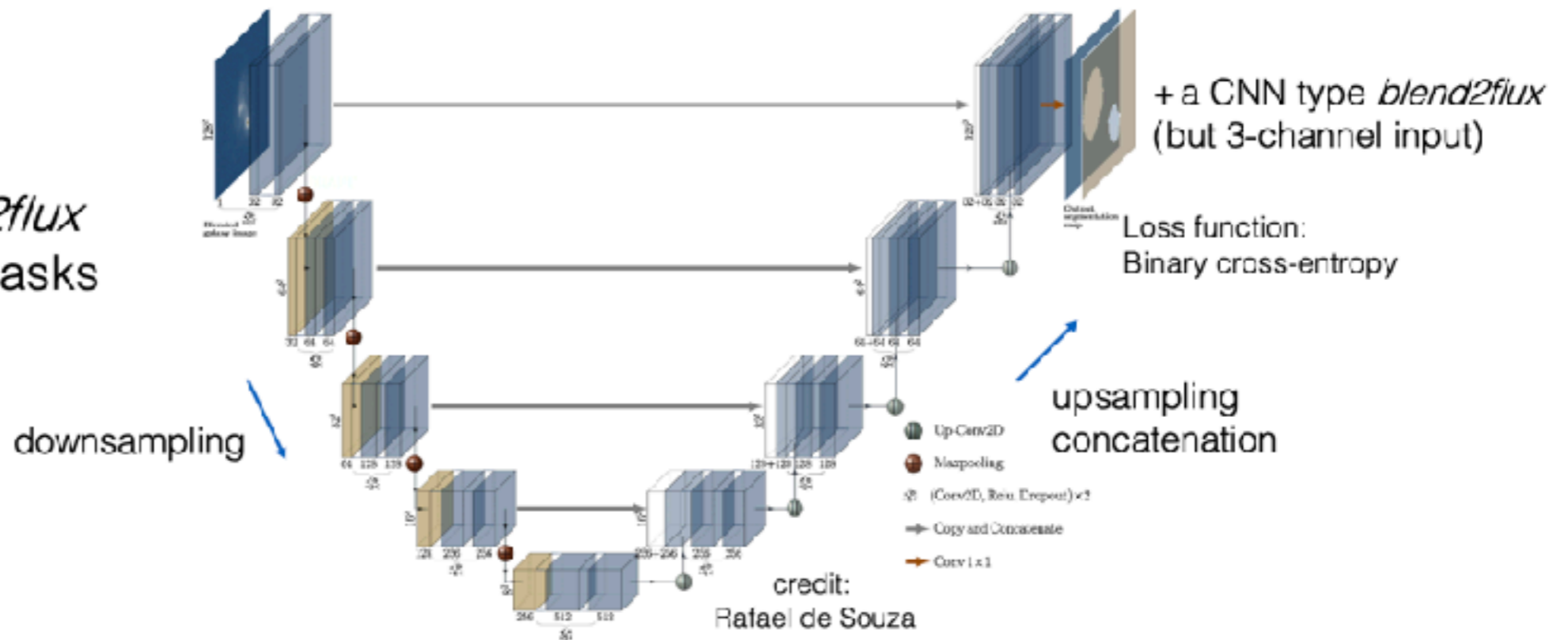


2a) The deblending problem

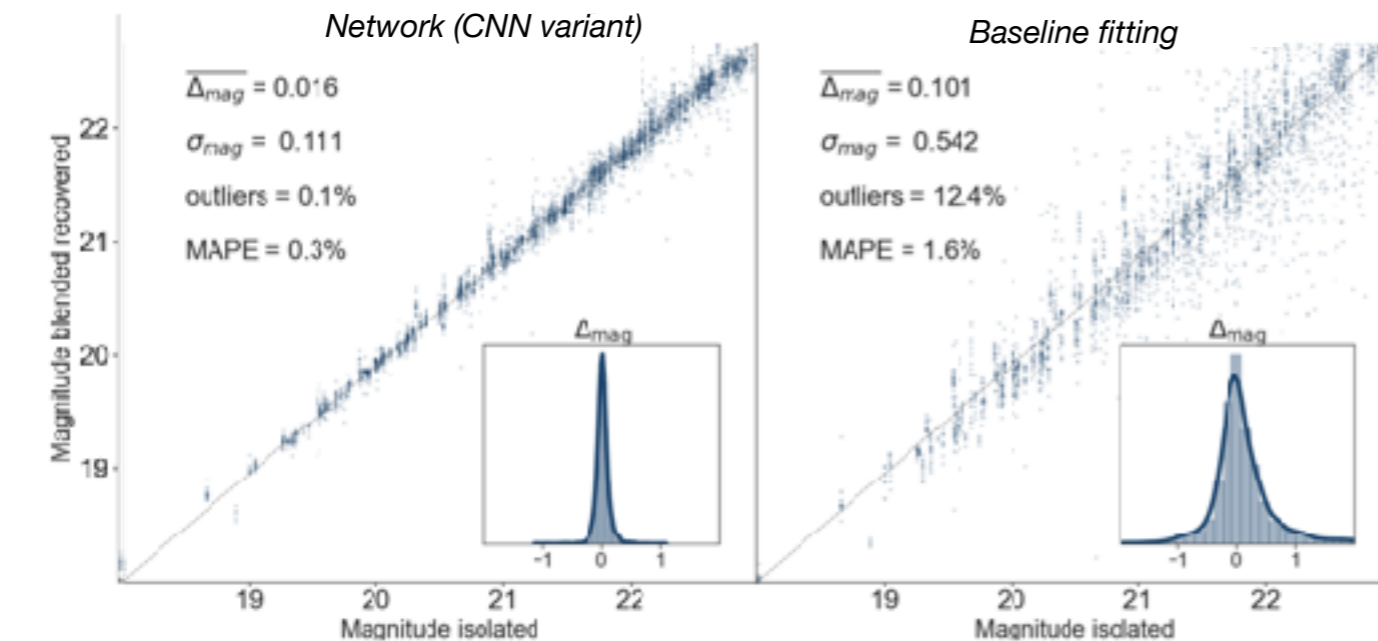
1) *blend2flux* a CNN for photometry



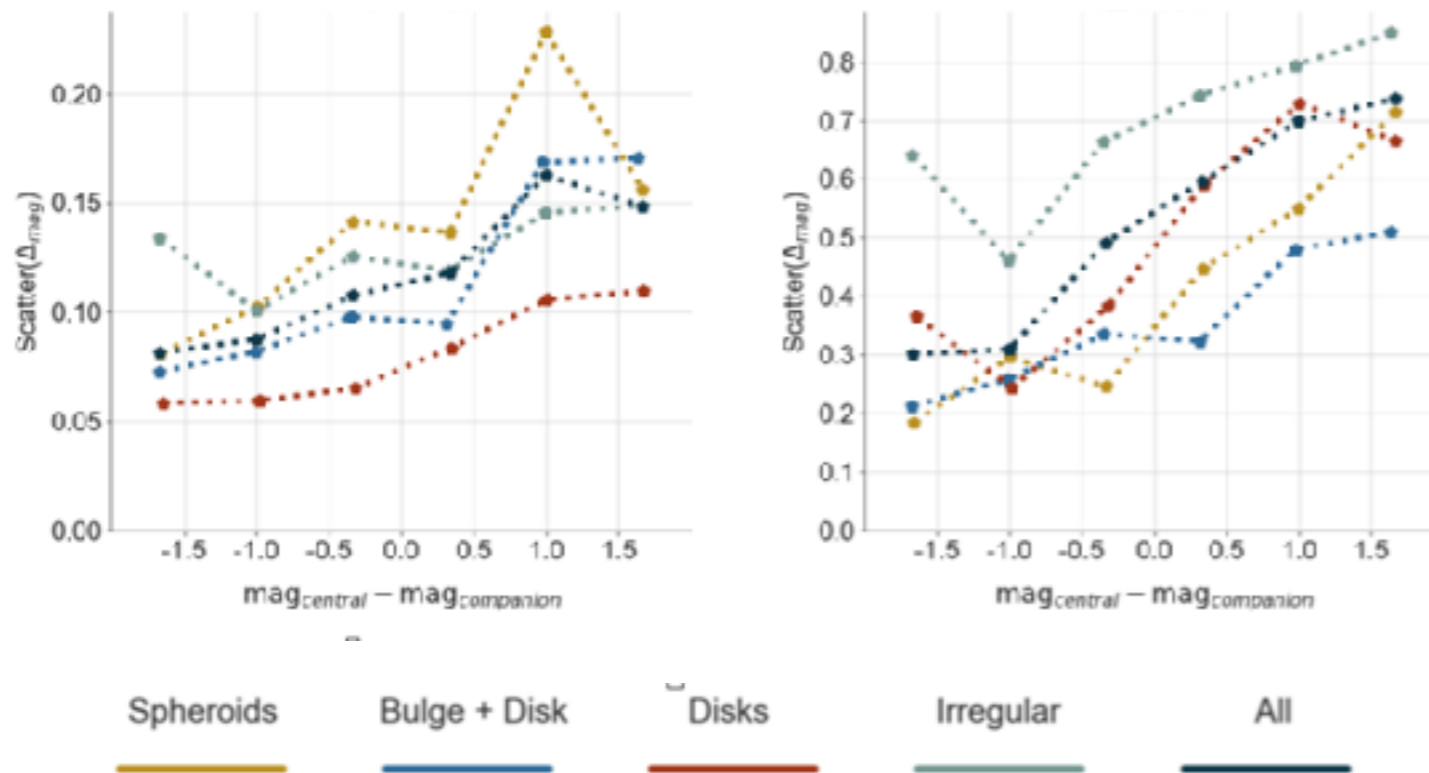
2) *blend2mask2flux* photometry + masks



2a) Optical source detection and characterisation



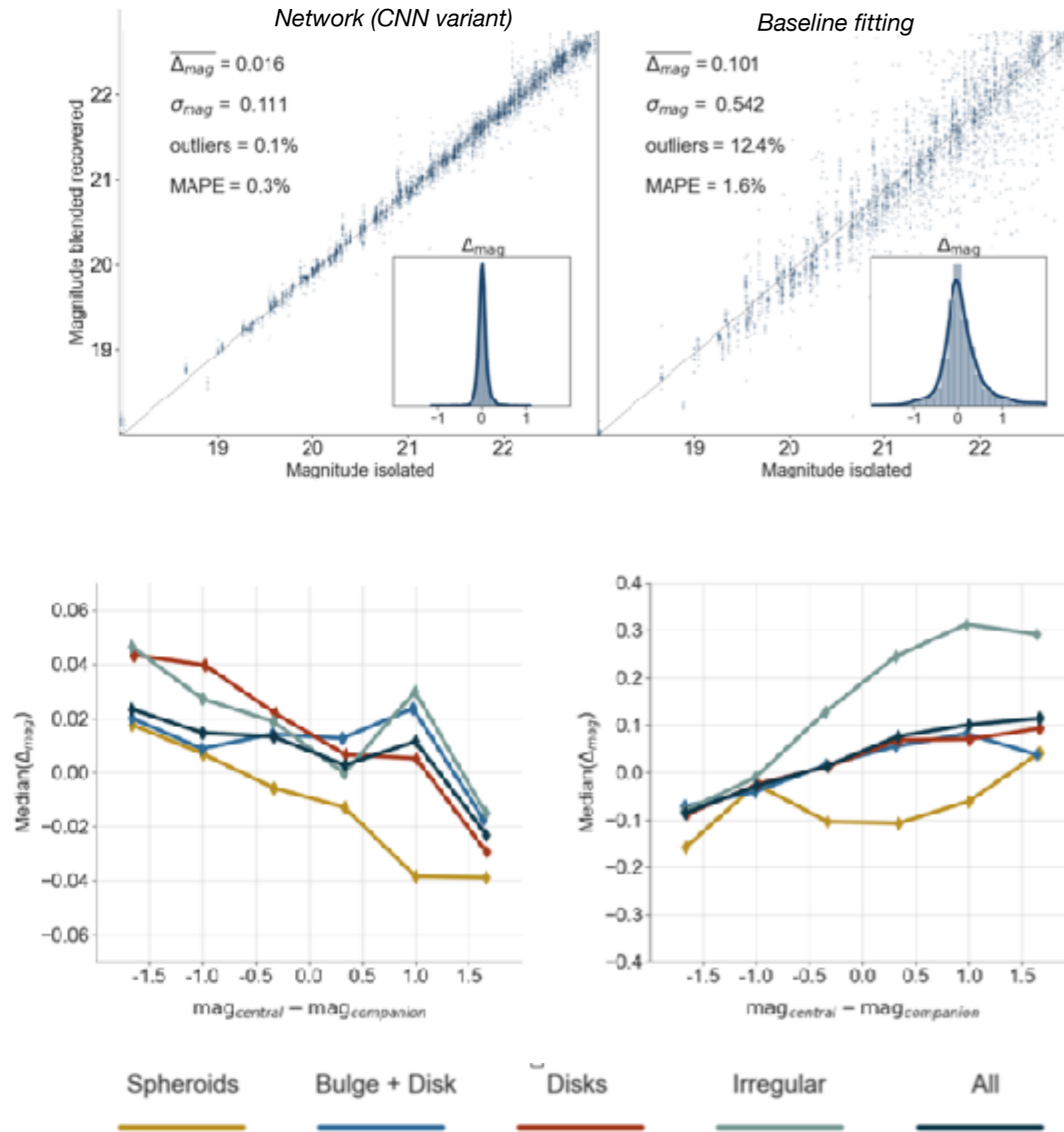
Get precise
Photometry



For irregular
shapes



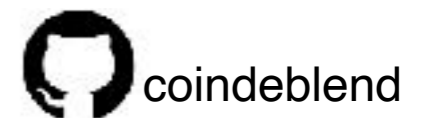
2a) Optical source detection and characterisation



Get precise
Photometry



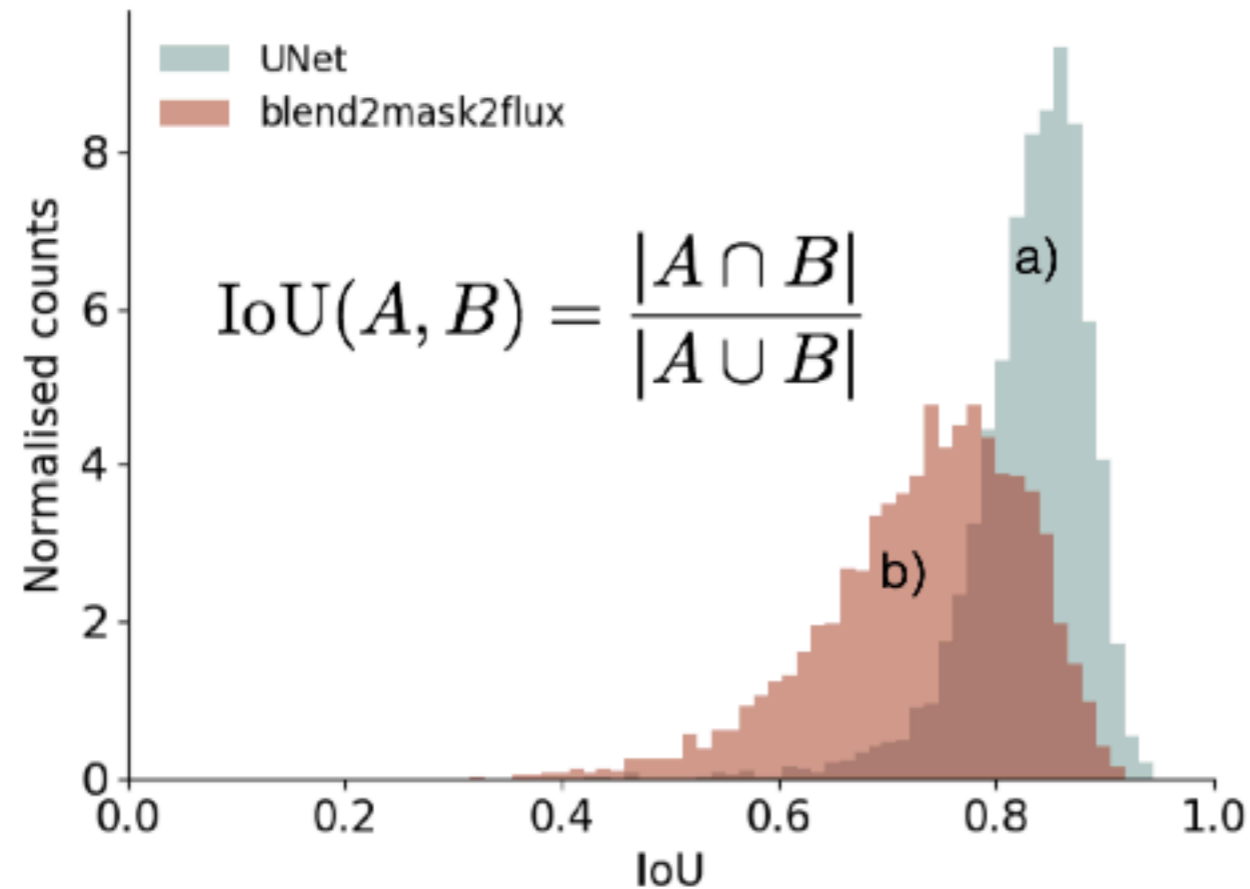
..and do so
bias-free



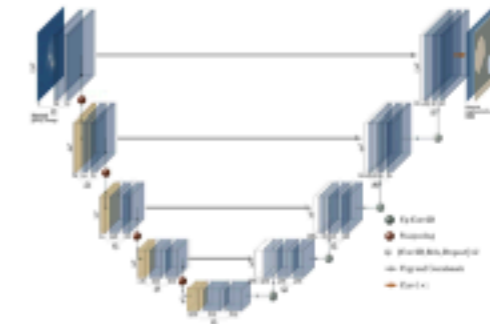
Boucaud, Huertas-Company, Heneka+ 20

2a) Optical source detection and characterisation

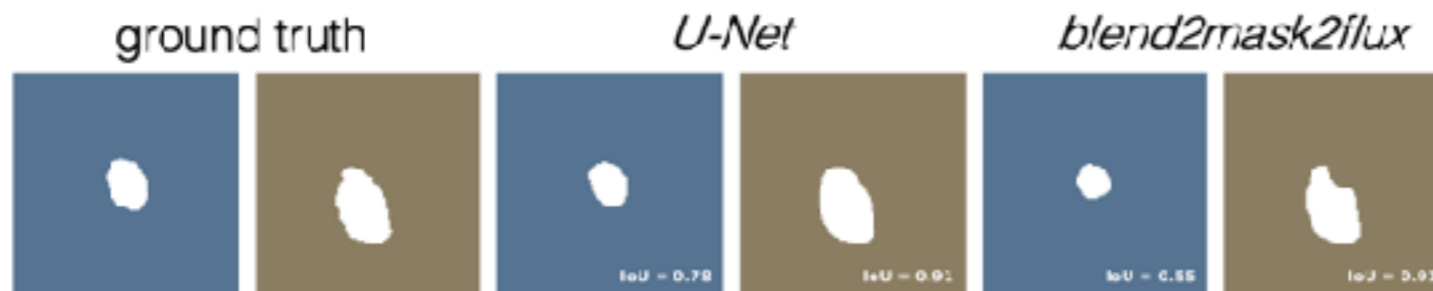
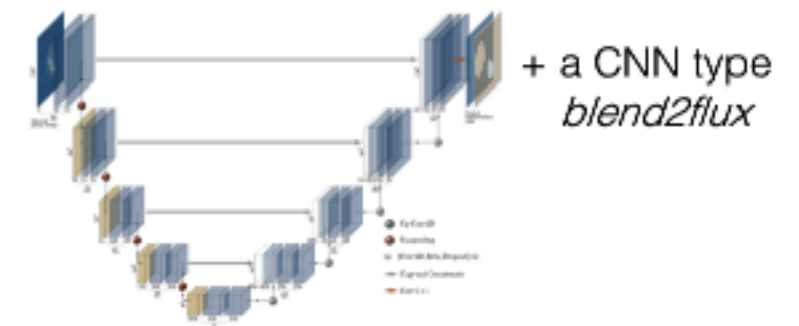
Optimal design is goal-dependent



a) *U-net*
masks, no photometry



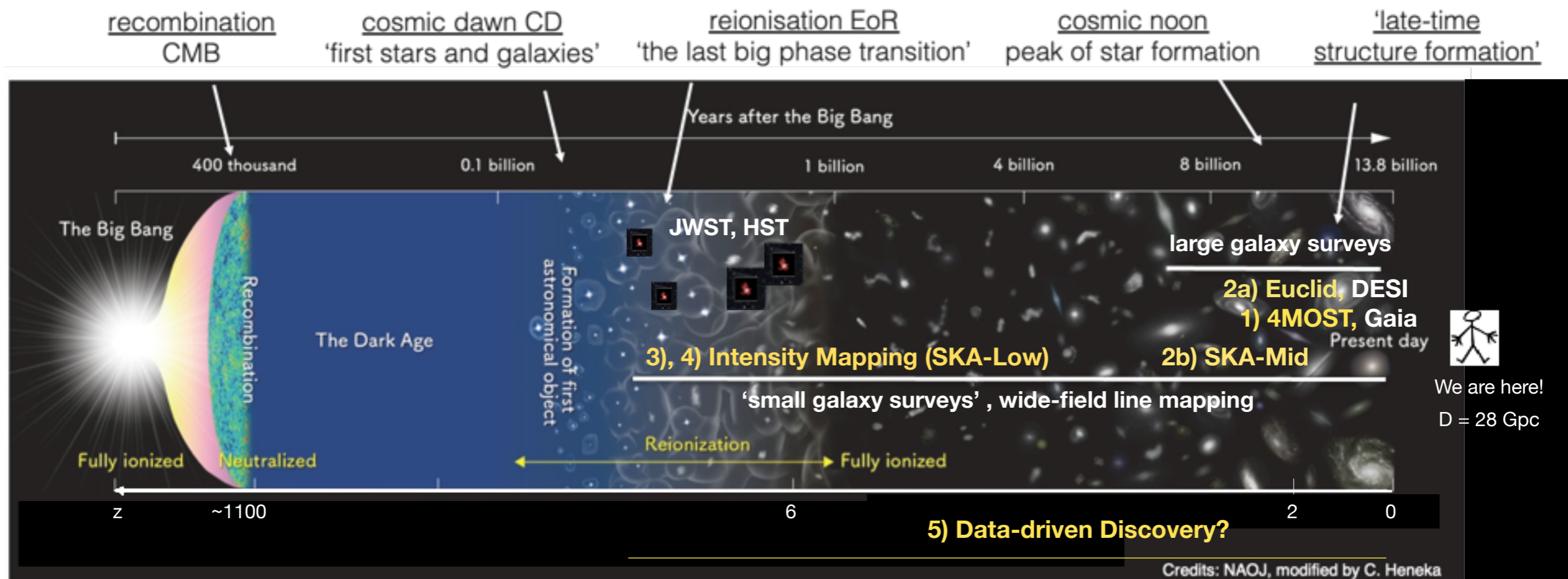
b) *blend2mask2flux*
photometry + masks



→ Dispersion broadens when optimised for photometry

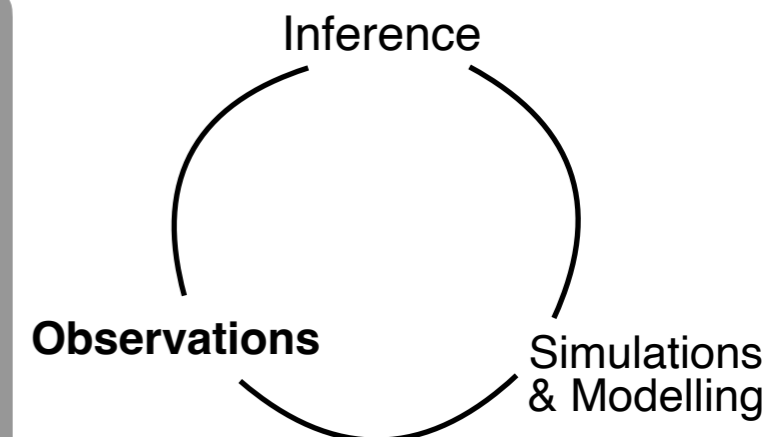
→ Tailor to your research question

Astronomical and Astrophysical Machine Learning



Highlights in this Lecture

- 1) Classification / Triggering
- 2) **Source detection & characterisation**
- 3) Simulation-based inference (SBI) in 3D
- 4) Generative methods
- 5) Data-driven Discovery



The Square Kilometre Array (SKA) in one slide

SKA in numbers

- Currently 16 member countries, >100 member organisations
- Routine science observations are expected to start in the late 2020s
- Consists of thousands of dishes and up to 1 million antennas, >1km² collecting area
- Expected data rate in full operation: 1 TB/s



SKA-LOW

SKA-MID

Credits: SKAO

SKA1-mid

the SKA's mid-frequency instrument



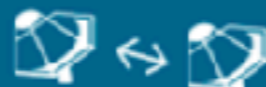
Location:
South Africa



Frequency range:
350 MHz
to
15.3 GHz
with a goal of 24 GHz



197 dishes
(including 84 MeerKAT dishes)



Maximum baseline:
150km

SKA1-low

the SKA's low-frequency instrument



Location: Australia



Frequency range:
50 MHz
to
350 MHz



~131,000
antennas spread between
512 stations



Maximum baseline:
~65km

The Square Kilometre Array (SKA) in one slide

<https://doi.org/10.1093/mnras/staa3837>

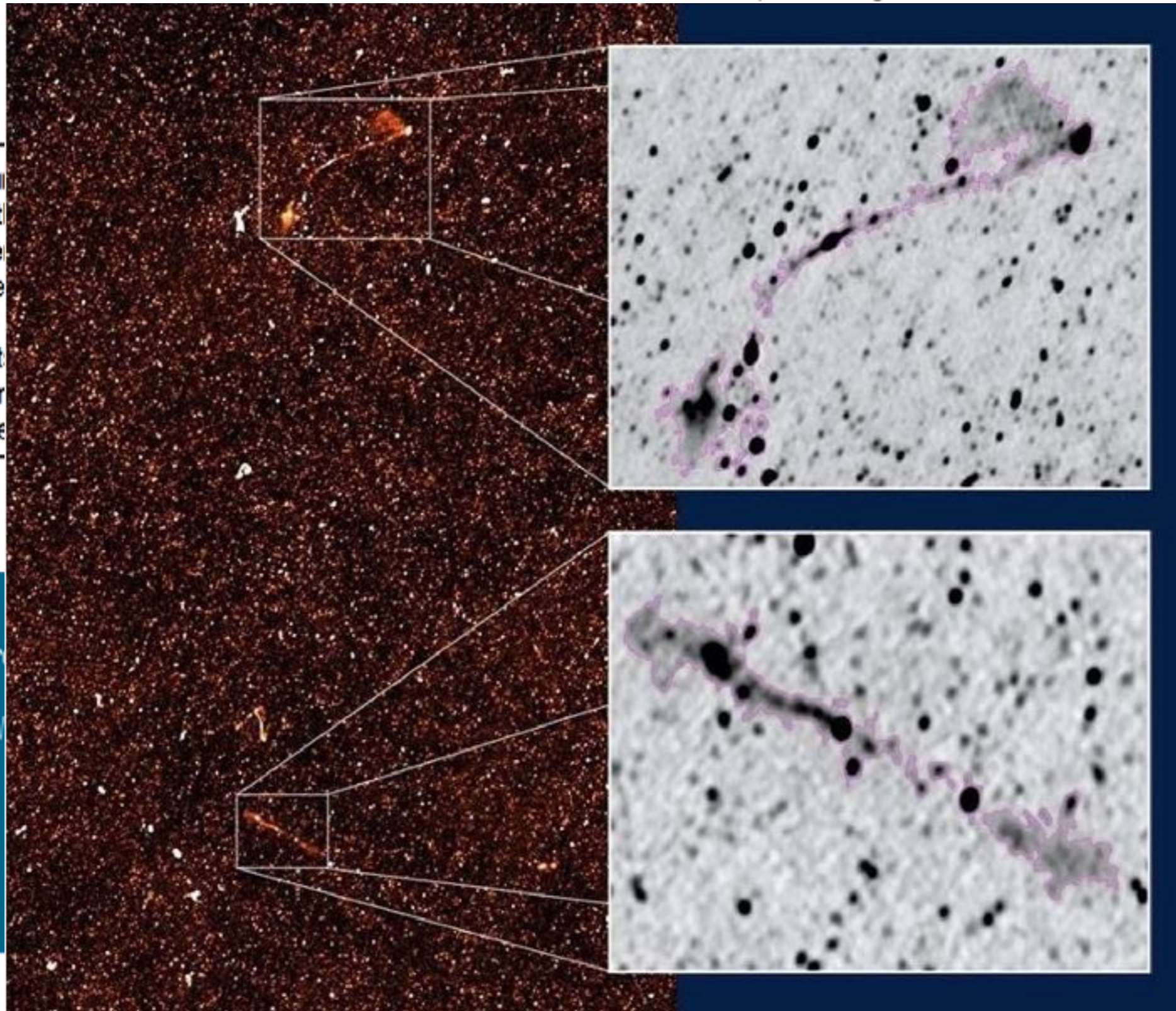
SKA in numbers

- Current membership
- Routine to start
- Consist
- 1 million
- Expecte

SKA1-mid
the SKA's mid-frequency



Location:
South Africa



2b) Radio source detection and characterisation

Example: Source detection in 3D tomographic data

Source finding

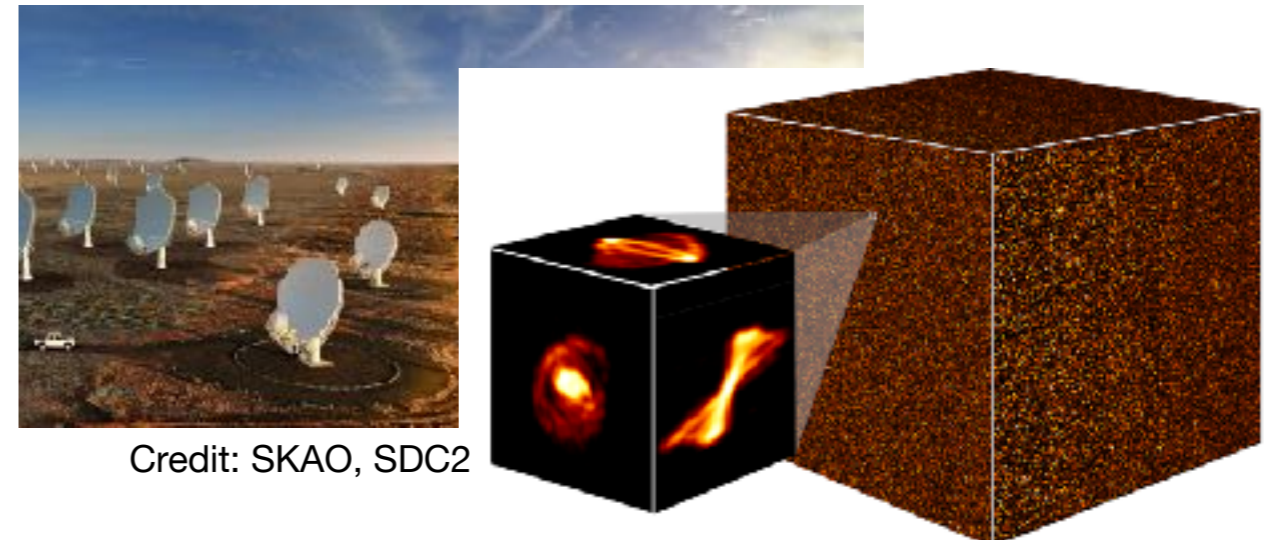
Location in RA, Dec,
central frequency (Hz)

Characterisation

- Integrated line flux (Jy Hz)
- Line width (km/s)
- HI major axis diameter (arcsec)
- Position angle (degrees)
- Inclination angle (degrees)

The challenging HI sources:

- low S/N
- small spatial size
- systematics

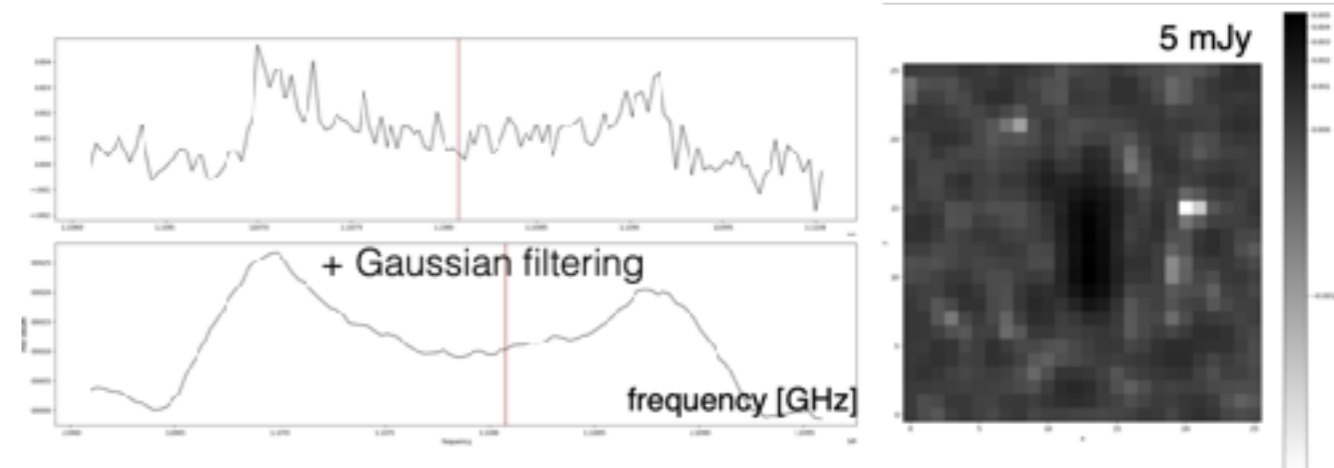


Credit: SKAO, SDC2

Total dimensions: (25,714 x 25,714 x 6,667) vox



The brightest HI source

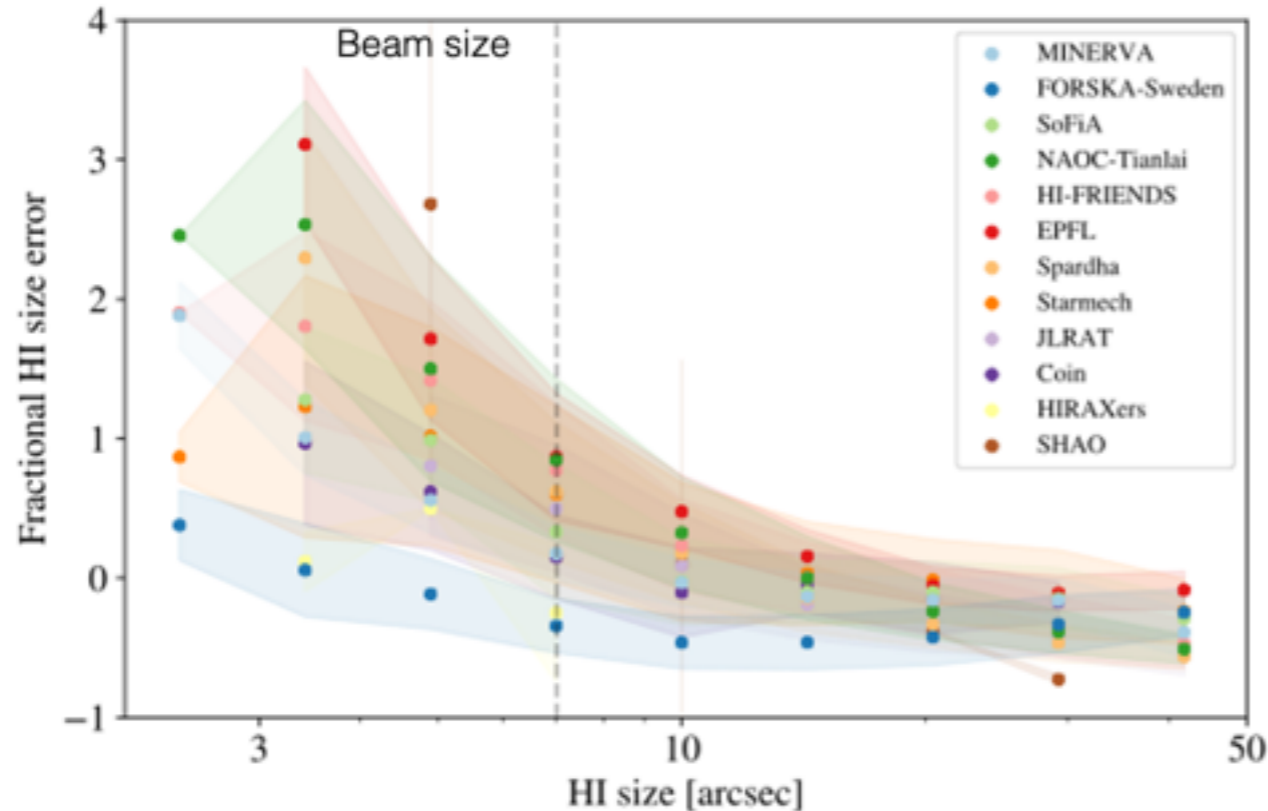
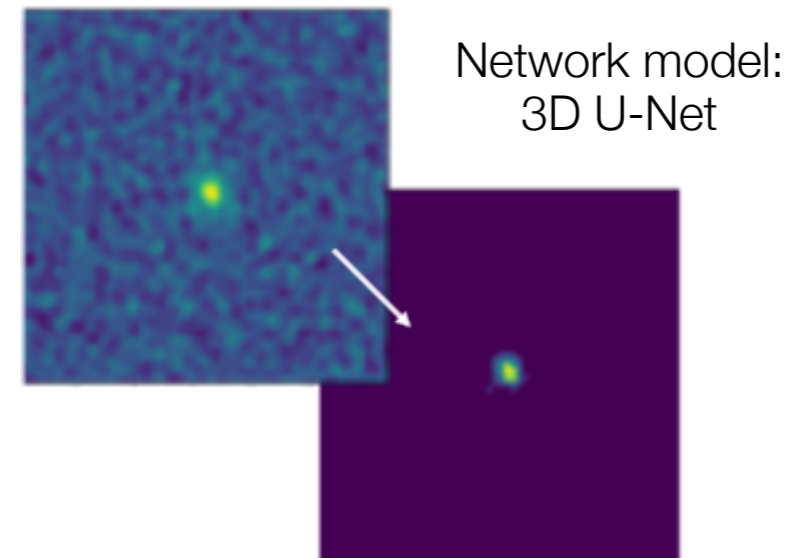


Hartley+ 23 (incl. Heneka), arXiv:2303.07943
Heneka 23, arXiv:2311.17553

2b) Radio source detection and characterisation

Example: Source detection in 3D tomographic data

- 3D better than stitching of 2D + 1D
- High-fidelity 3D reconstructions
- **Good prior for characterisation tasks via nets:**

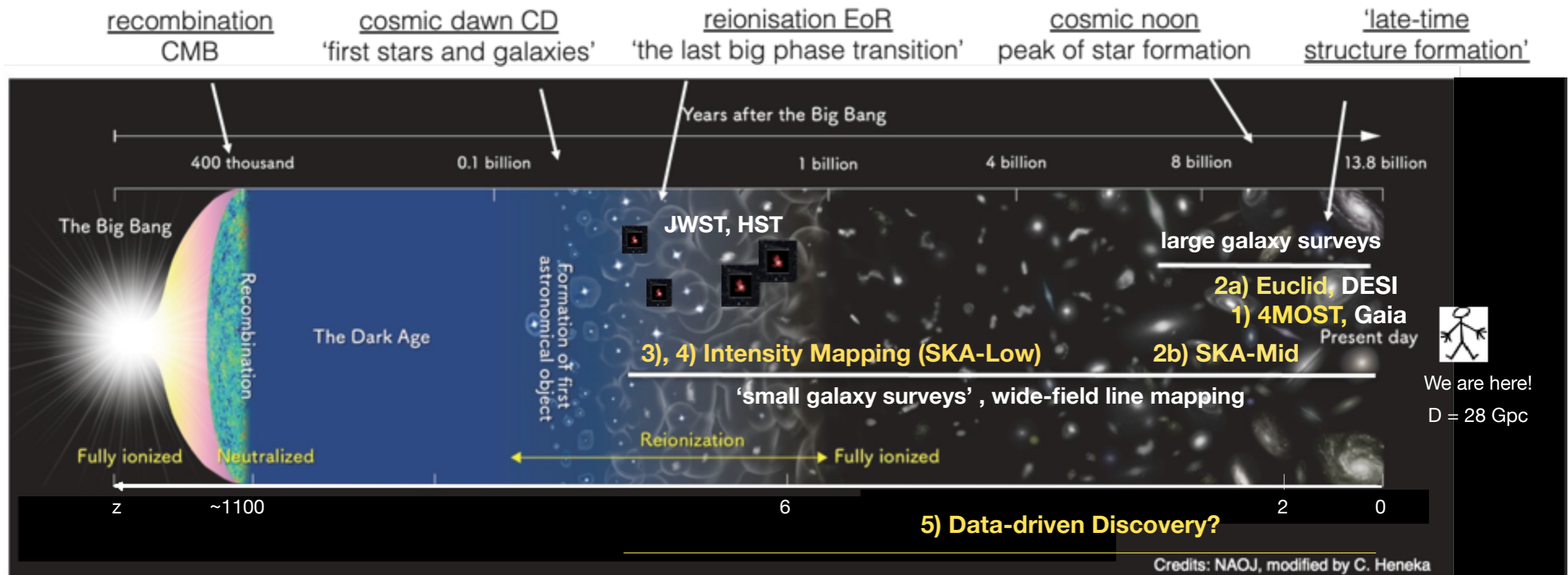


Recovery across wide range in HI flux and size
Pushing to low S/N recovery came at a cost (FPs)

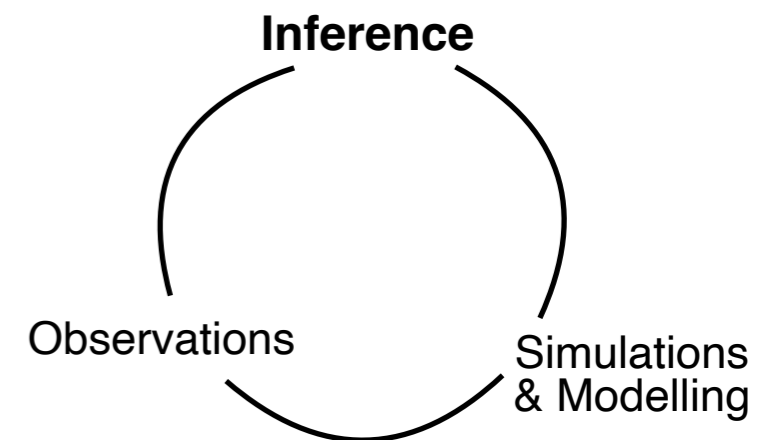
Uncertainty?

Hartley+ 23 (incl. Heneka), arXiv:2303.07943
Heneka 23, arXiv:2311.17553

Astronomical and Astrophysical Machine Learning



- ### Highlights in this Lecture
- 1) Classification / Triggering
 - 2) Source detection & characterisation
 - 3) **Simulation-based inference (SBI) in 3D**
 - 4) Generative methods
 - 5) Data-driven Discovery



3) Simulation-based inference (SBI) for intensity mapping (3D)

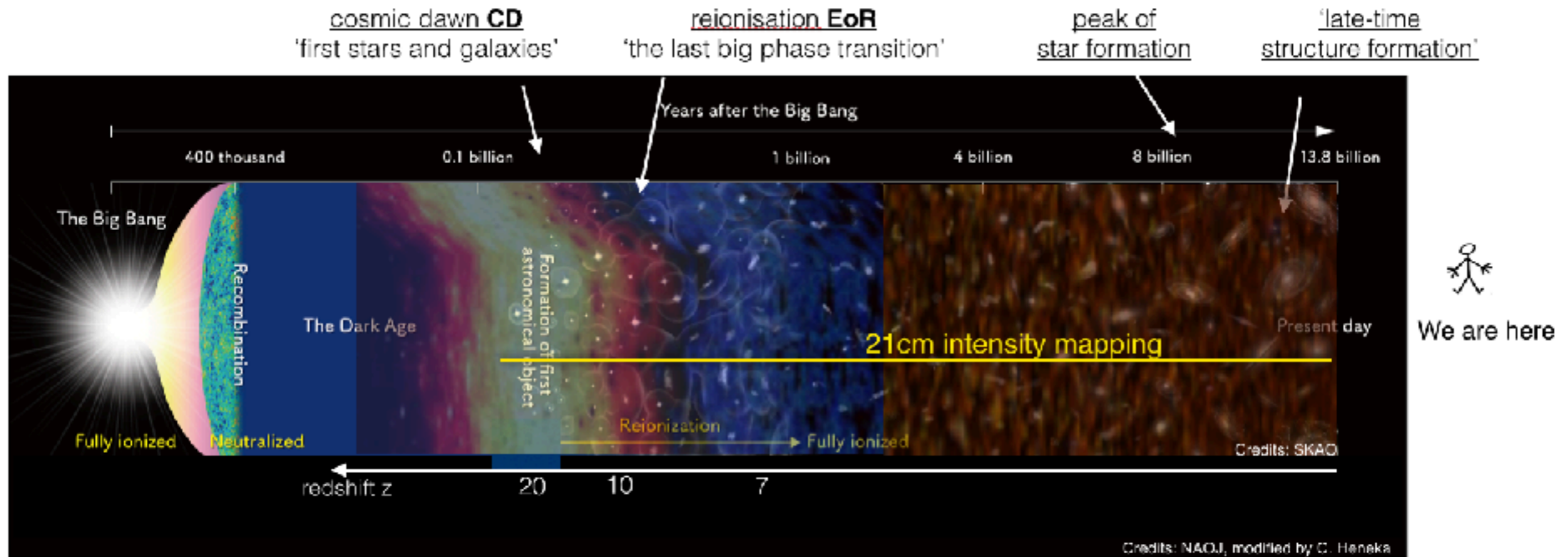
21cm signal

a tracer of neutral hydrogen:

$$\delta T_b(\nu) = \frac{T_S - T_\gamma}{1+z} (1 - e^{-\tau_{\nu 0}})$$

$$\propto x_{\text{HI}} (1 + \delta_{\text{nl}}) \left(\frac{H}{dv_r/dr + H} \right)$$

Why care?
Tomography of >80% of the Universe

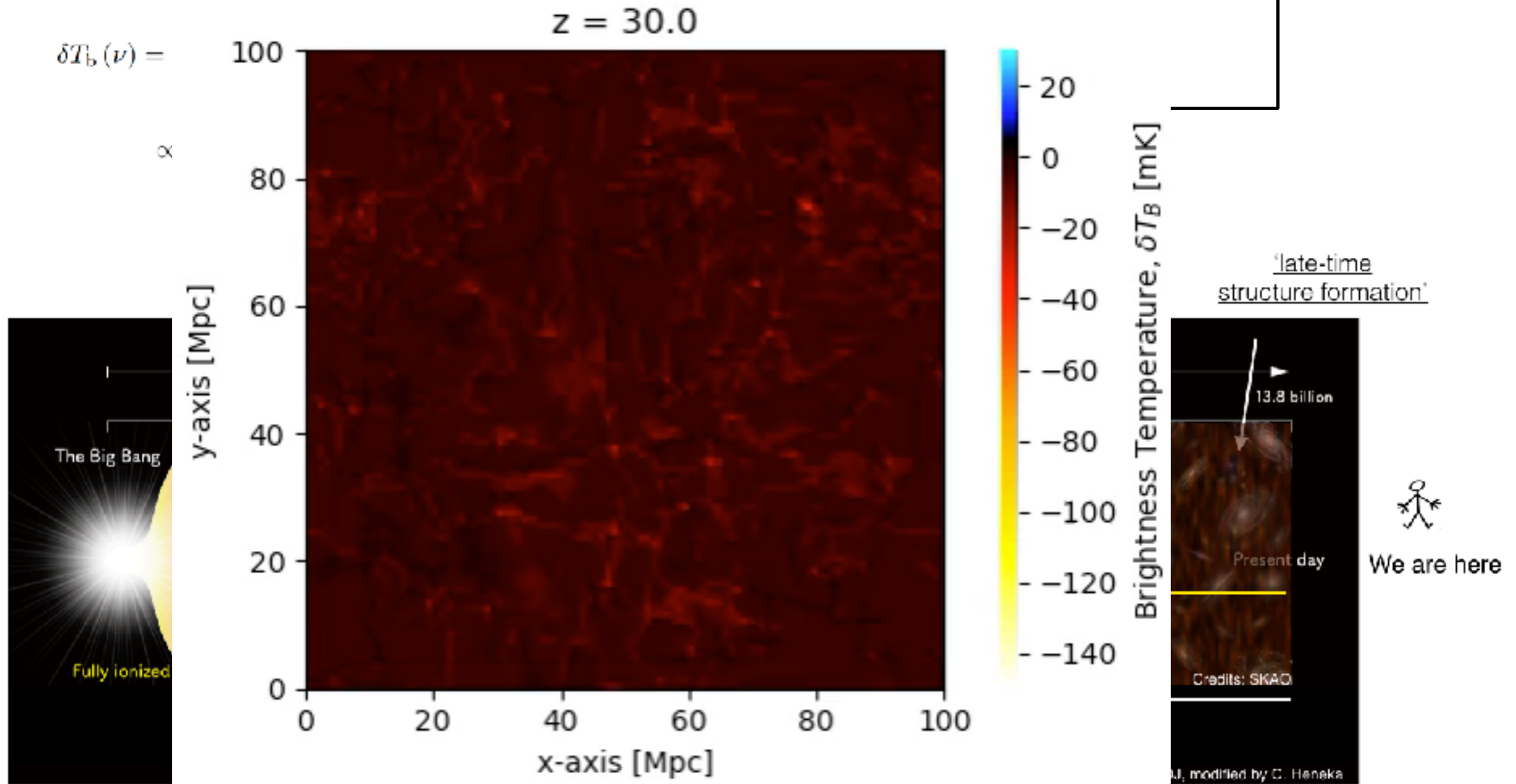


3) Simulation-based inference (SBI) for intensity mapping (3D)

21cm signal

a tracer of neutral hydrogen:

Why care?
Tomography of >80% of the Universe



3) Simulation-based inference (SBI) for intensity mapping (3D)

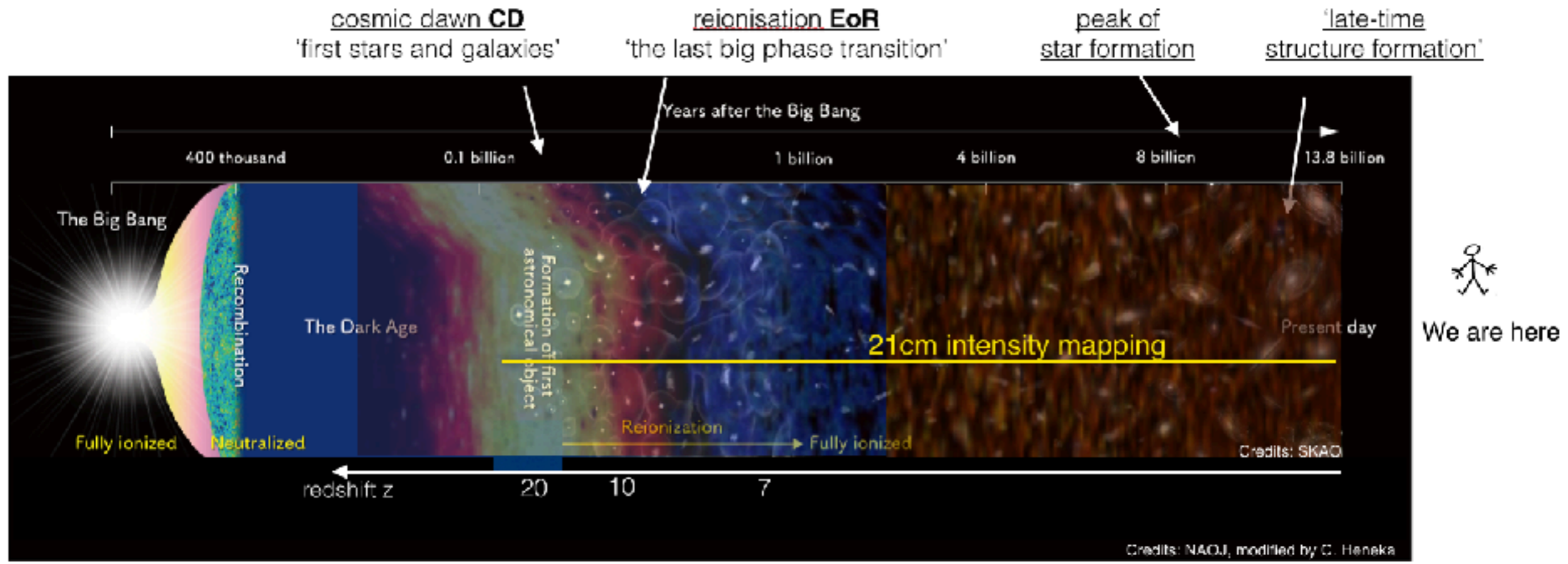


Why care?
 Tomography of >80% of the Universe
 Square Kilometre Array - true 'Big Data'
 non-linear, non-Gaussian signal

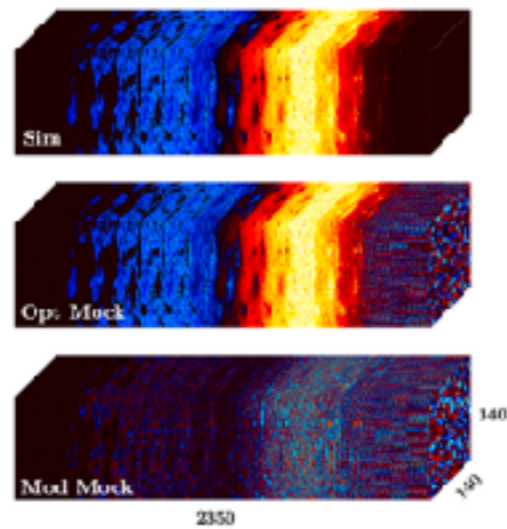
Expected data SKA rate:
 TB/s, few EB/day
 Archive: ~700 PB/yr



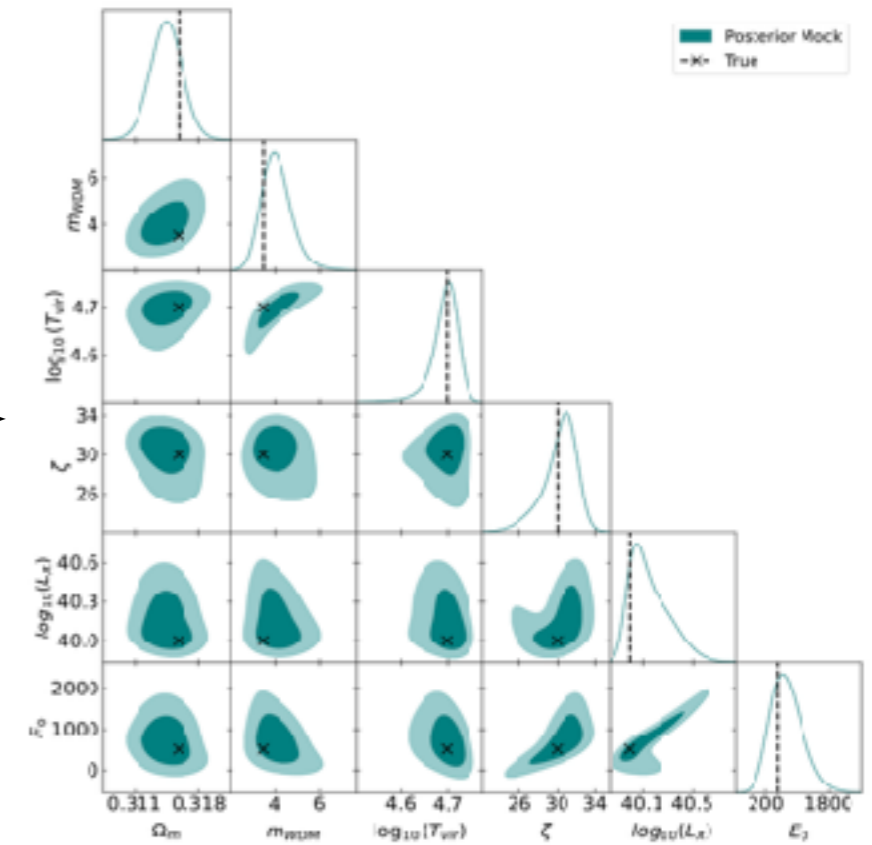
Move to full likelihood inference with networks



3) Simulation-based inference (SBI) for intensity mapping (3D)



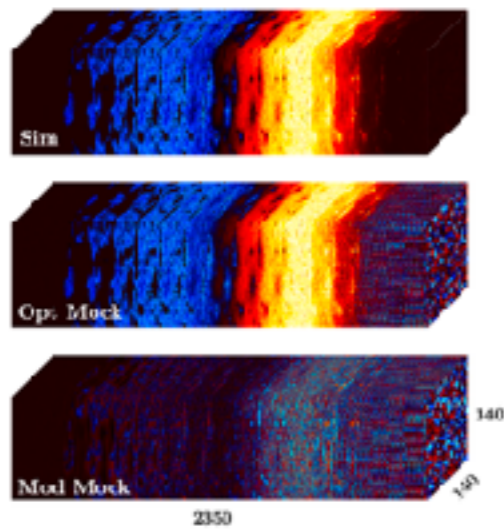
?



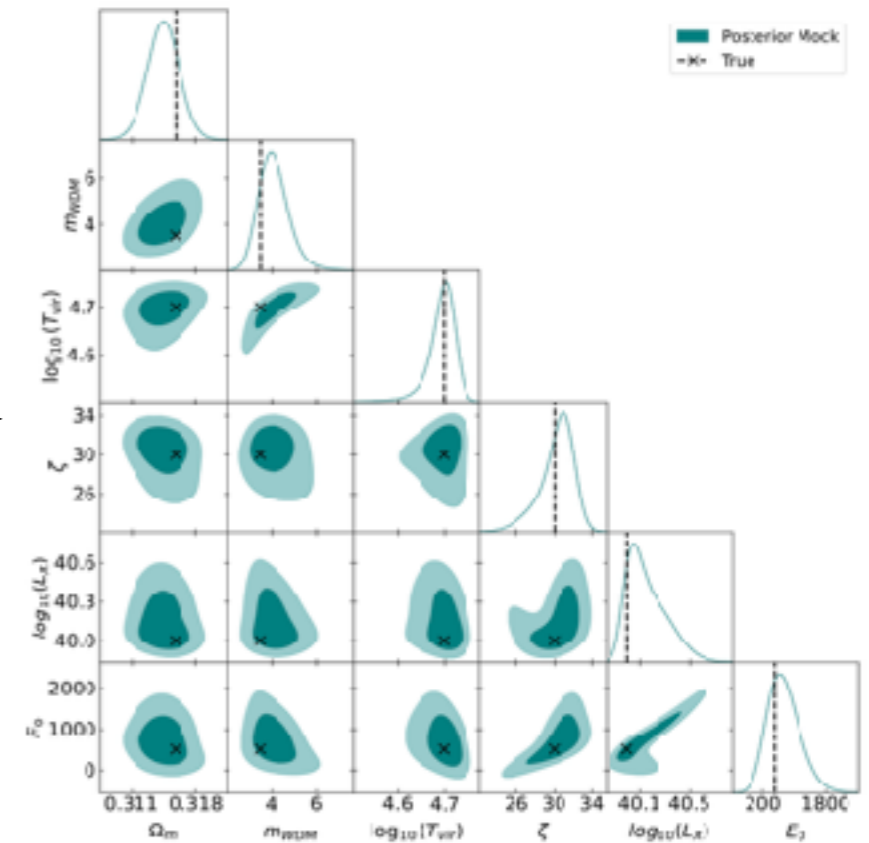
Neutsch, Heneka, Brüggem (2022), arXiv:2201.07587

Schossler, Heneka, Plehn, arXiv:2401.04174

3) Simulation-based inference (SBI) for intensity mapping (3D)



SBI with flows/cINN
based on BayesFlow
(arXiv:2003.06281)



3D-21cmPIE-Net (public)

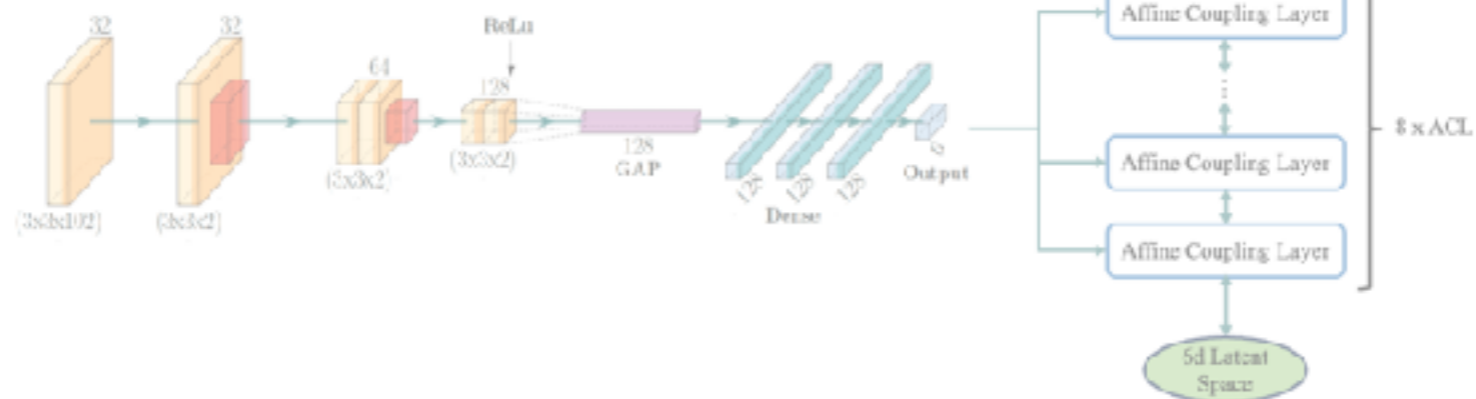
Neutsch, Heneka, Brüggén (2022)
arXiv:2201.07587

+

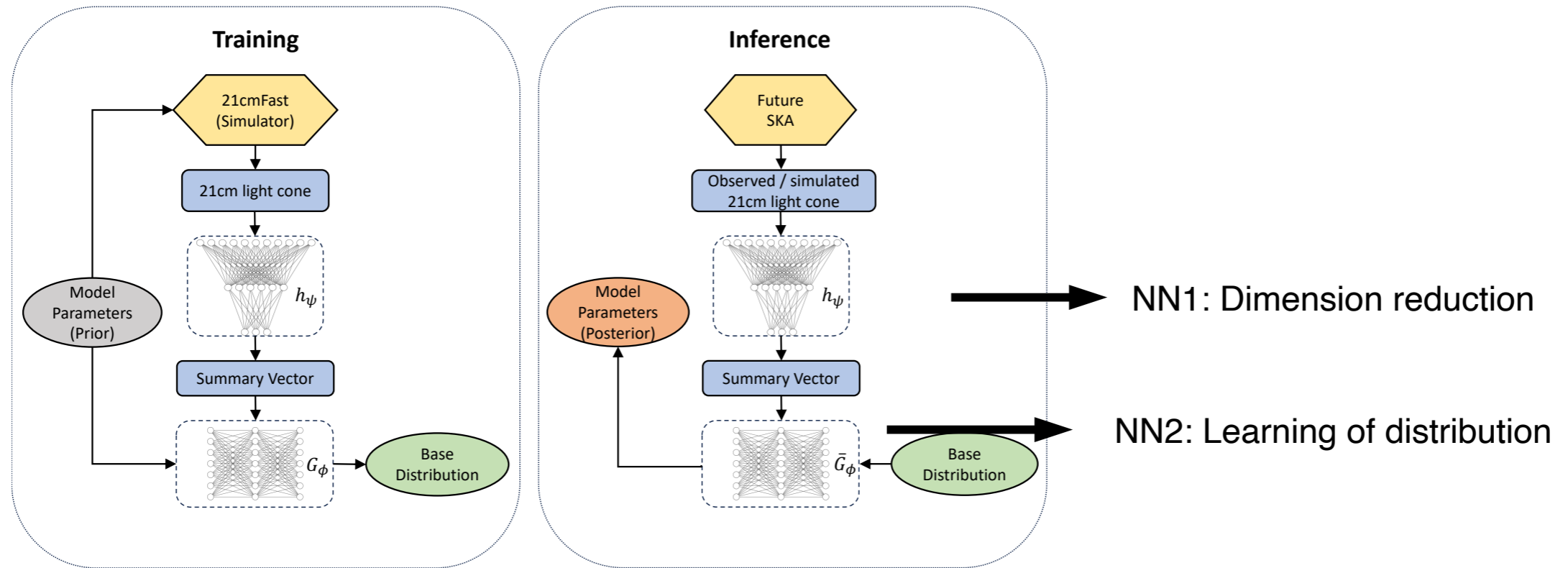


21cm-cINN (public when published)

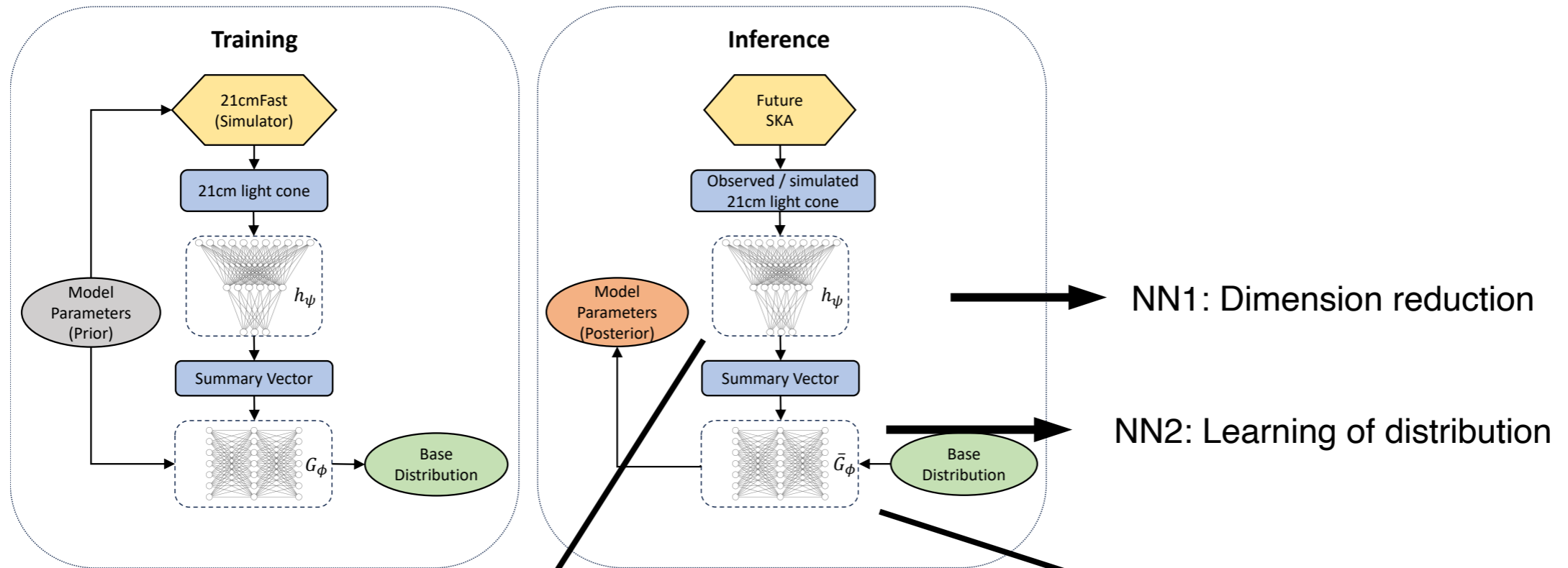
Schosser, Heneka, Plehn, arXiv:2401.04174



3) Simulation-based inference (SBI) for intensity mapping (3D)

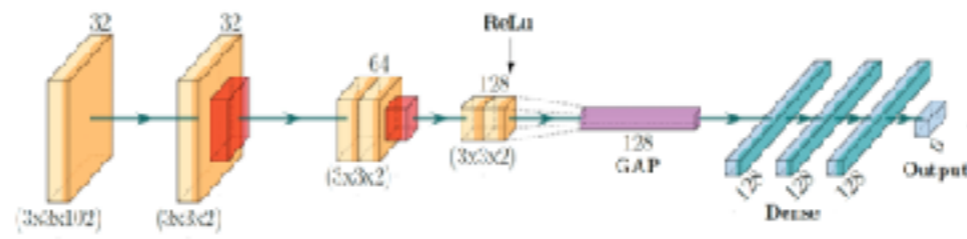


3) Simulation-based inference (SBI) for intensity mapping (3D)



Network Model:

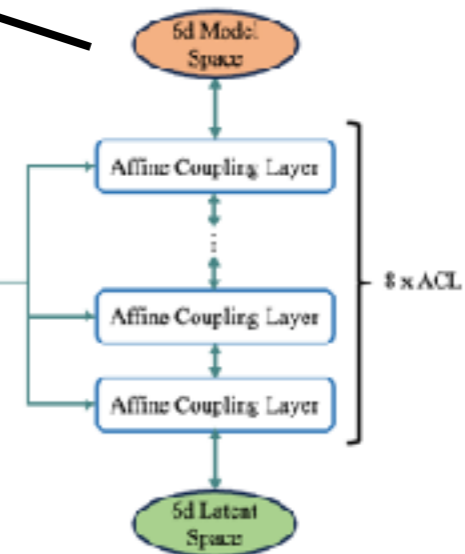
Small networks,
comparably few parameters
= fast training and convergence



NN1: 3D-CNN

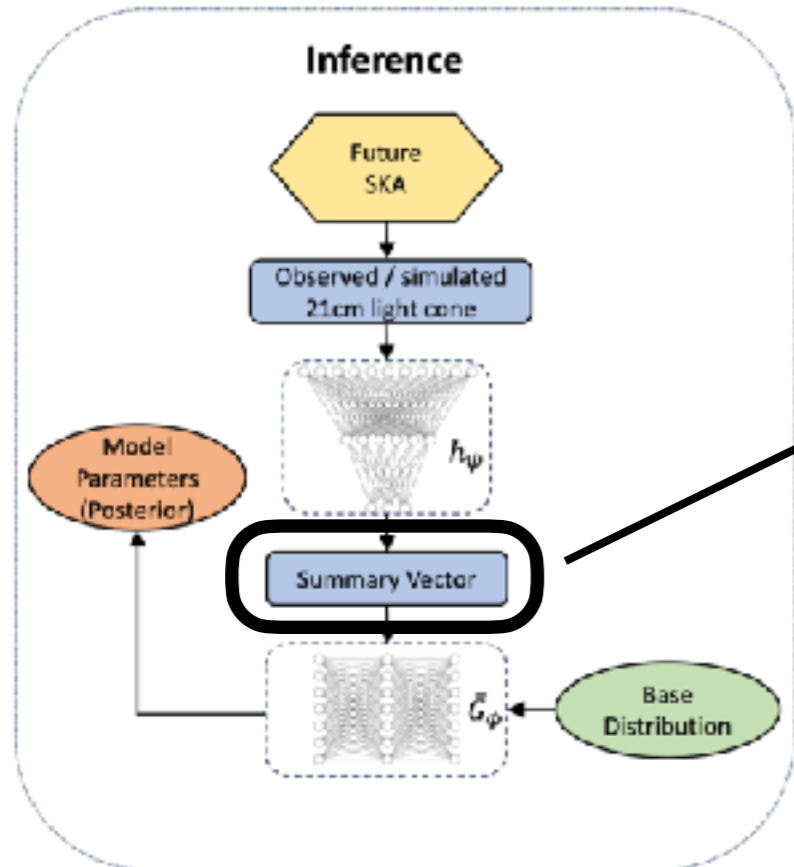


3D-21cmPIE-Net (public)
Neutsch, Heneka, Brüggem (2022)
arXiv:2201.07587



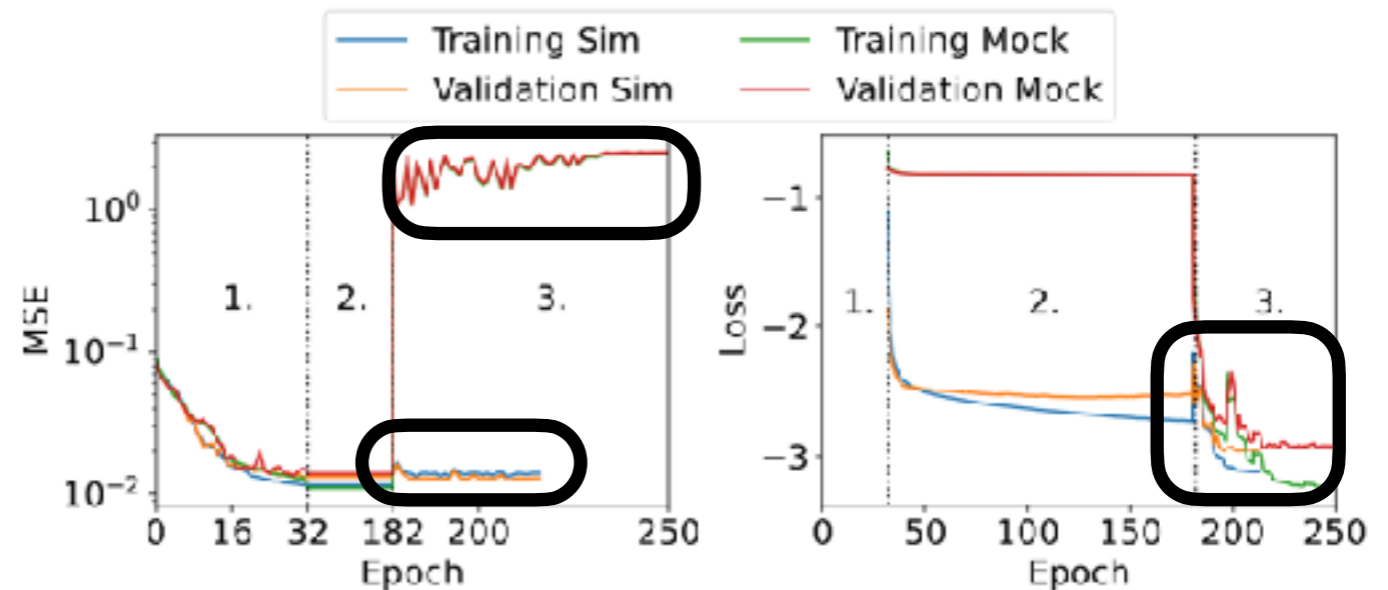
NN2: Invertible Network
21cm-cINN (public when published)
Schosser, Heneka, Plehn, arXiv:2401.04174

3) Simulation-based inference (SBI) for intensity mapping (3D)



Summary vector versus summary statistics:
 Free adjustment with scheduled training
 = move away from 'parameter interpretation' of vector

Key difference between Sim and Mock:



We profit from learned summary in presence of noise (more).

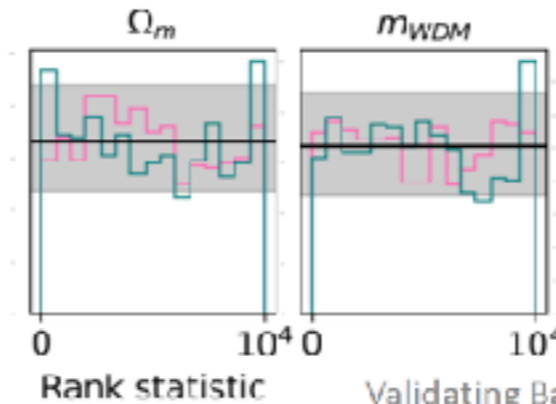
Sim: Summary stays close to original
 Mock: Heavy adjustment of summary vector

Schossler, Heneka, Plehn, arXiv:2401.04174

3) Simulation-based inference (SBI) for intensity mapping (3D)

Performance validation via:

- Distribution of latent variables
- Simulation-based calibration
- Parameter recovery

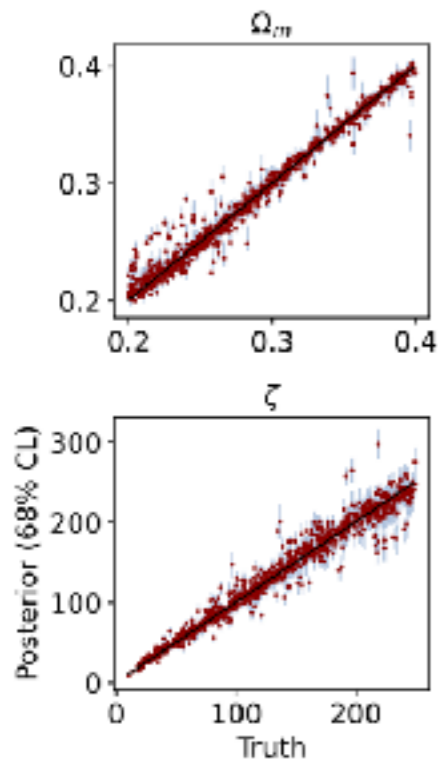


Uniform distribution



Self-consistent sampling

Validating Bayesian Inference Algorithms with Simulation-Based Calibration, arXiv:1804.06788



Check Posterior vs. True label

3) Simulation-based inference (SBI) for intensity mapping (3D)

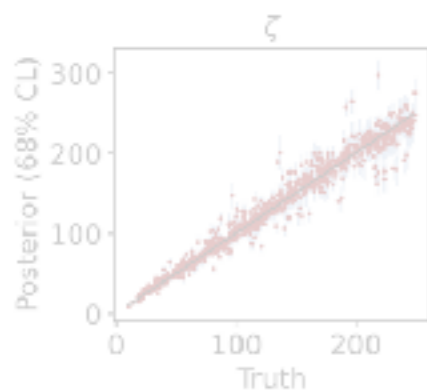
Performance validation via:

- Distribution of latent variables
- Simulation-based calibration
- Parameter recovery

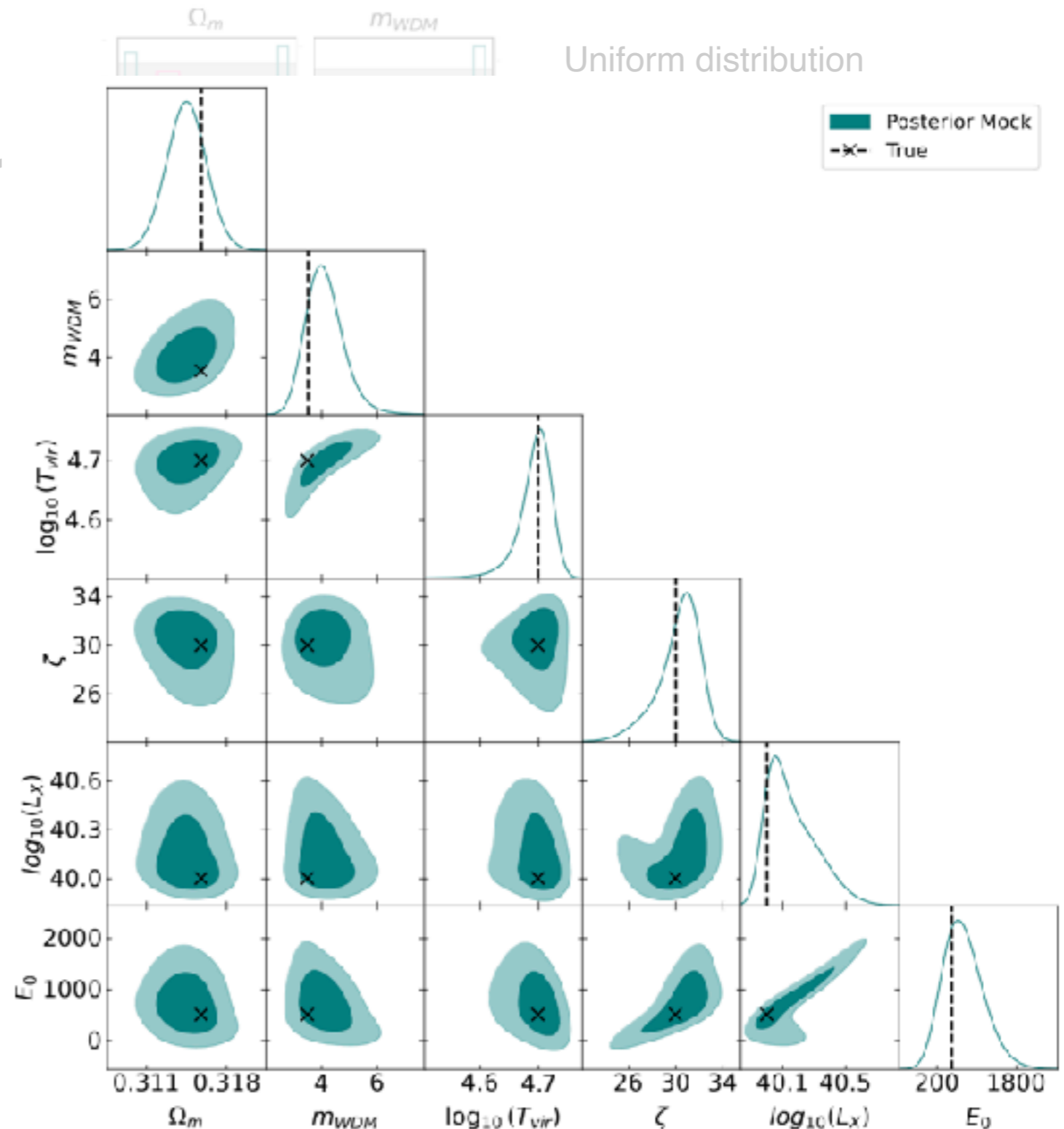
Likelihood contours:

Explore any model in prior range!

+ similar performance
Mock vs. Sim
(except for E_0)



Check Posterior vs. True label



Schosser, Heneka, Plehn, arXiv:2401.04174

3) Simulation-based inference (SBI) for intensity mapping (3D)

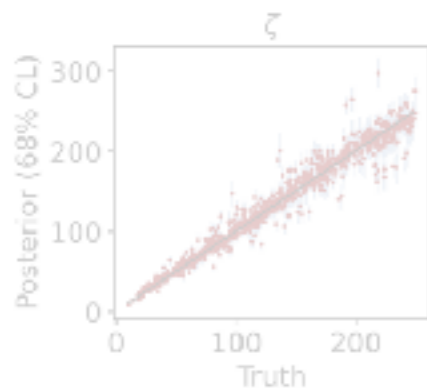
Performance validation via:

- Distribution of latent variables
- Simulation-based calibration
- Parameter recovery

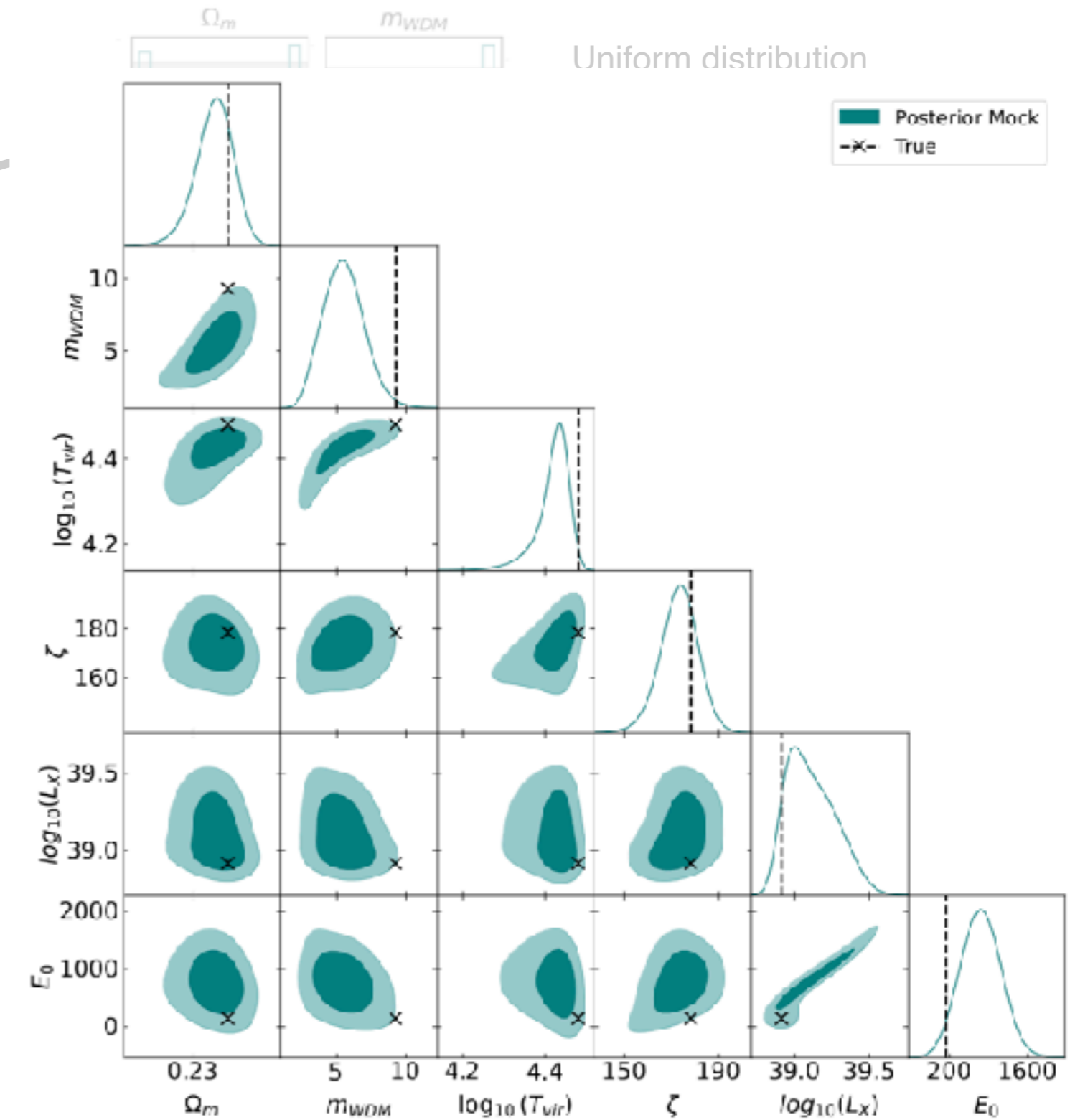
Likelihood contours:

Explore any model in prior range!

+ similar performance
Mock vs. Sim
(except for E_0)



Check Posterior vs. True label



Schosser, Heneka, Plehn, arXiv:2401.04174

3) Simulation-based inference (SBI) for intensity mapping (3D)

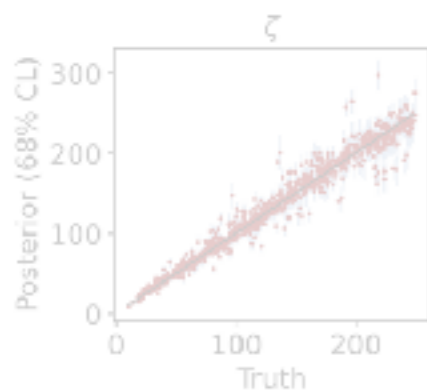
Performance validation via:

- Distribution of latent variables
- Simulation-based calibration
- Parameter recovery

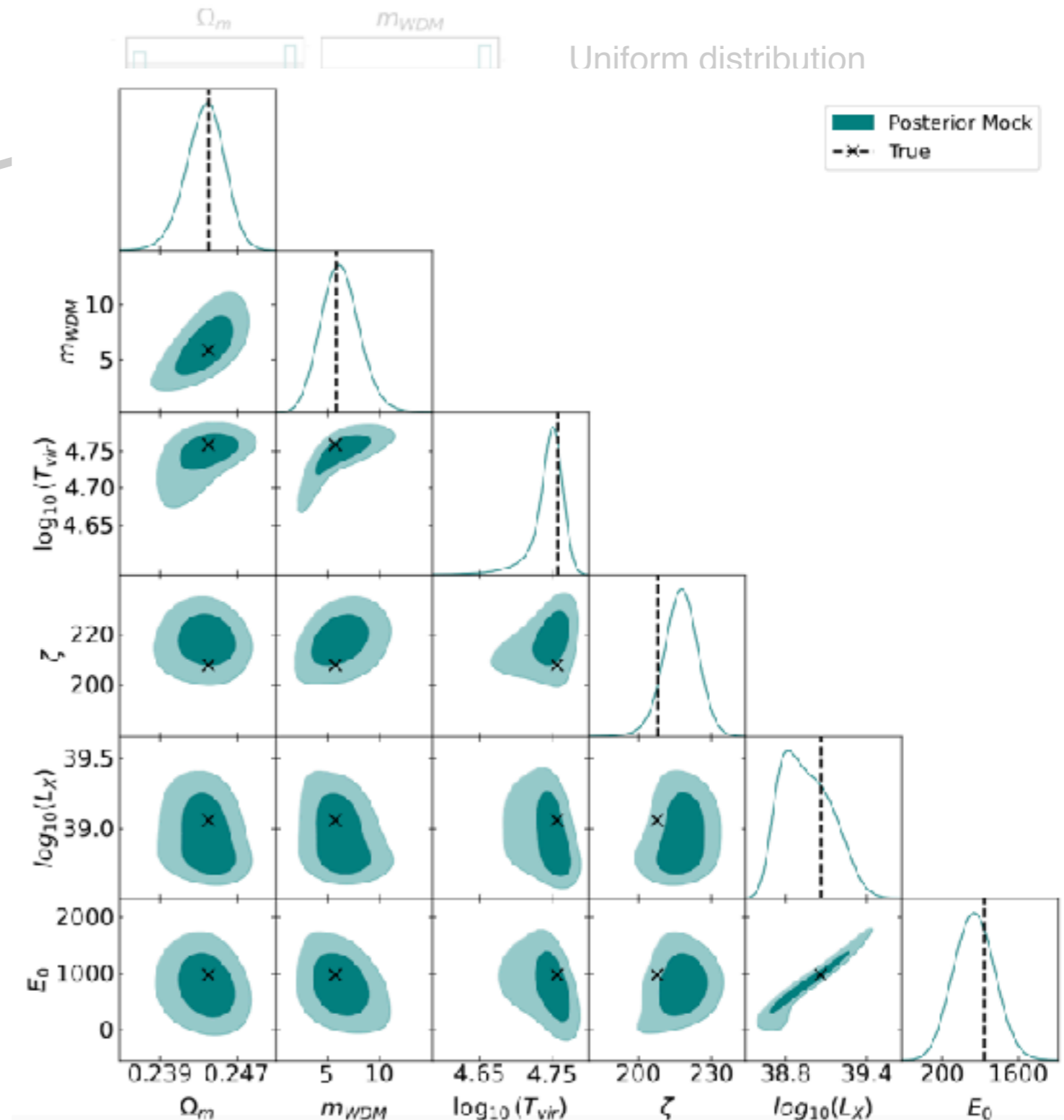
Likelihood contours:

Explore any model in prior range!

+ similar performance
Mock vs. Sim
(except for E_0)

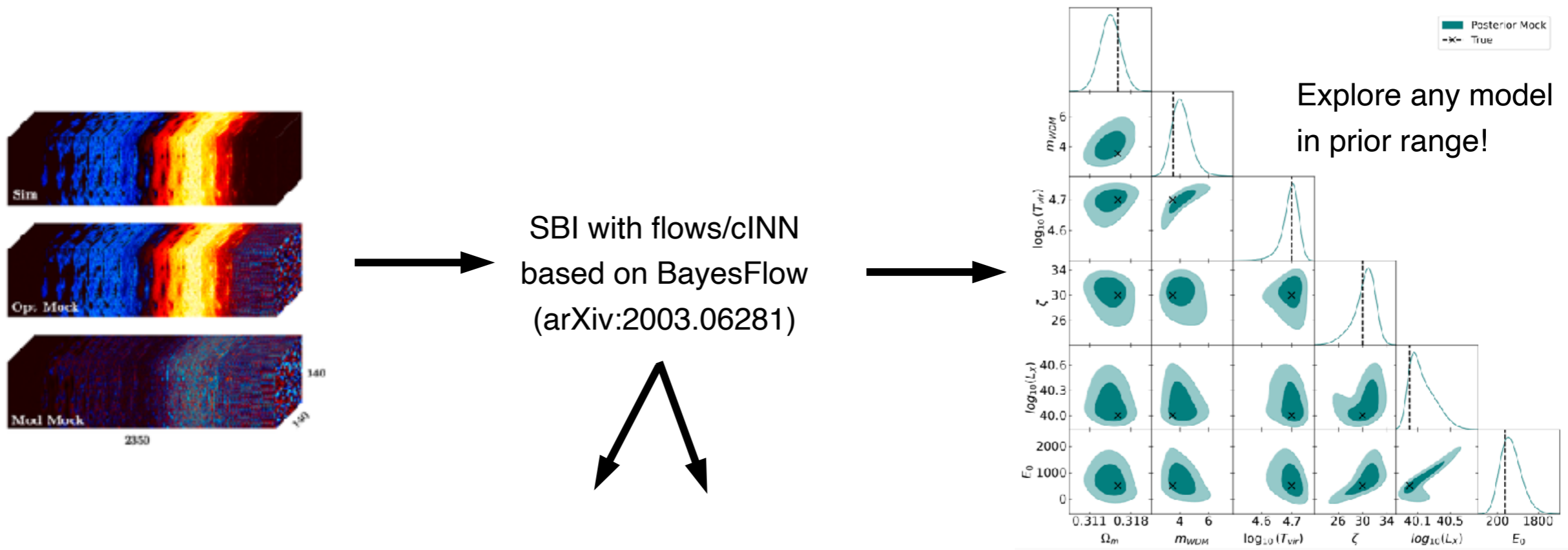



Check Posterior vs. True label



Schossler, Heneka, Plehn, arXiv:2401.04174

3) Simulation-based inference (SBI) for intensity mapping (3D)

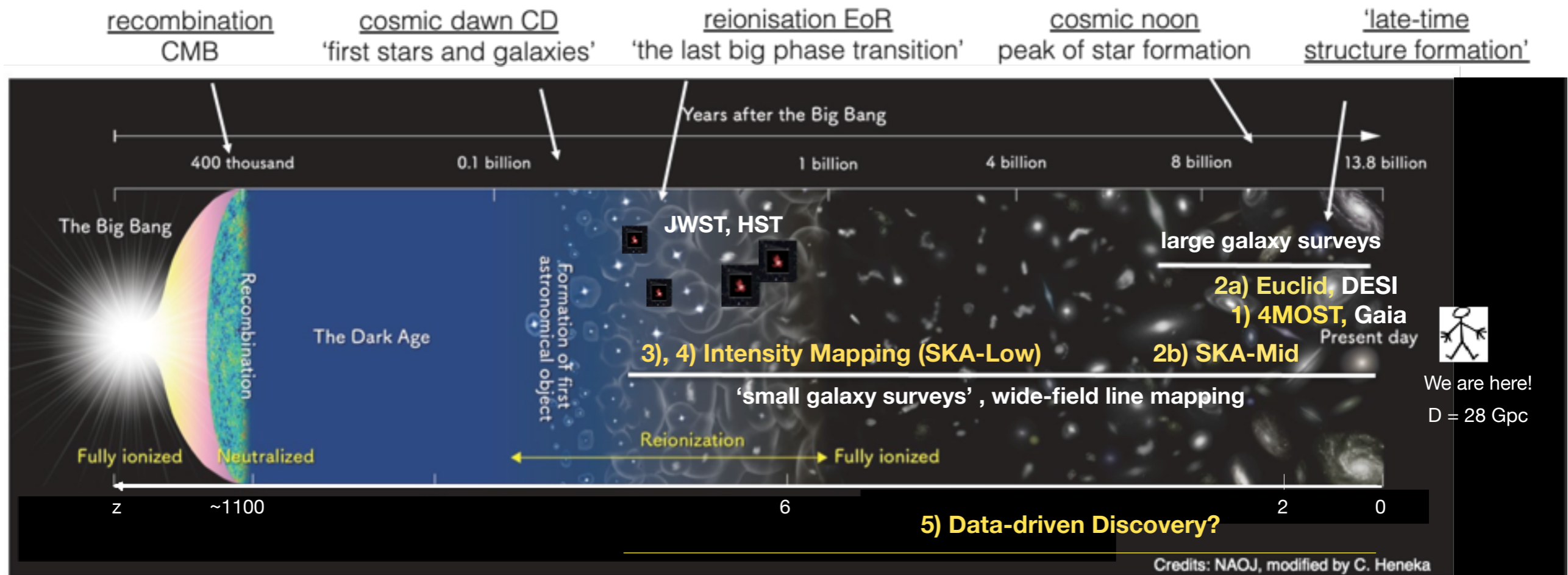


 3D-21cmPIE-Net (public)
 Neutsch, Heneka, Brüggem (2022)
 arXiv:2201.07587

+  21cm-cINN (public when published)
 Schosser, Heneka, Plehn (2024), arXiv:2401.04174

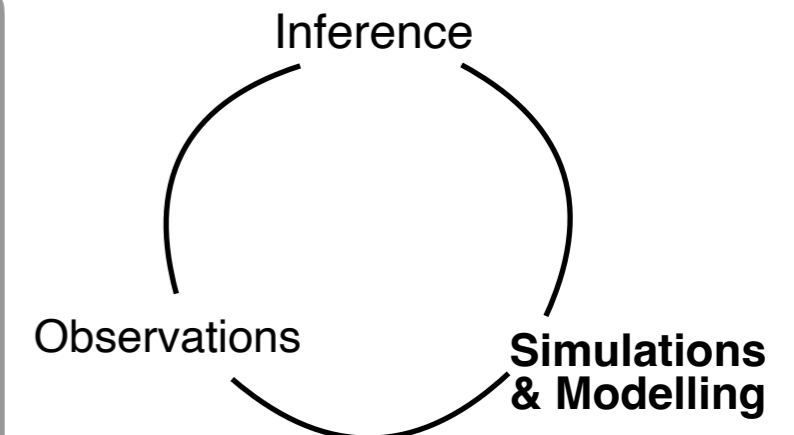
'Optimal, fast, and robust inference of reionization-era cosmology with the 21cmPIE-INN'

Astronomical and Astrophysical Machine Learning



Highlights in this Lecture

- 1) Classification / Triggering
- 2) Source detection & characterisation
- 3) Simulation-based inference (SBI) in 3D
- 4) **Generative methods**
- 5) Data-driven Discovery



3) Generative methods for simulation

Reionization simulations are costly

Is there a fast way to emulate whole simulations?

+ Radiative transfer

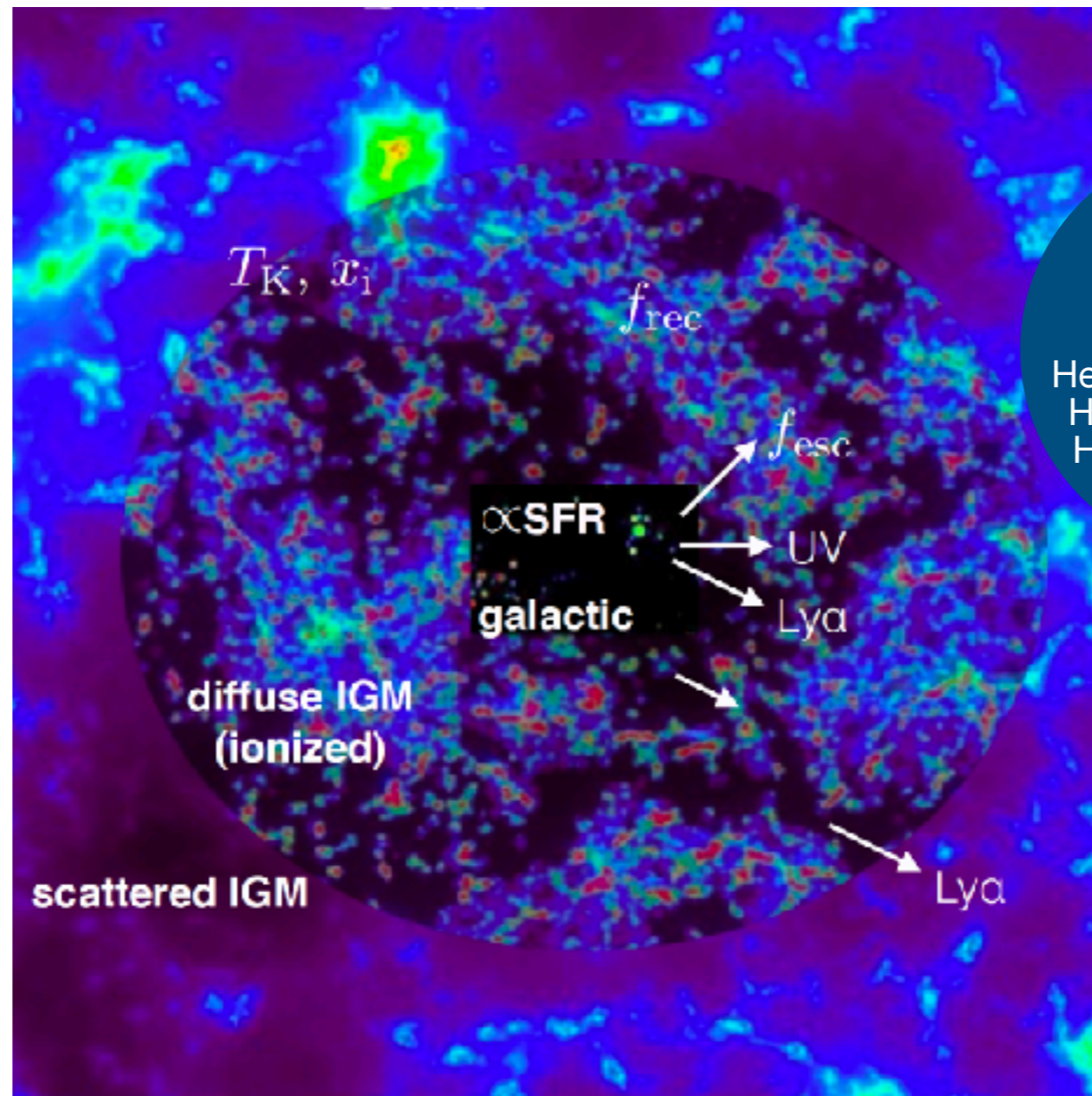
Hydrodynamical

Semi-numerical

Semi-analytical

Empirical scaling

Simulation
Cost



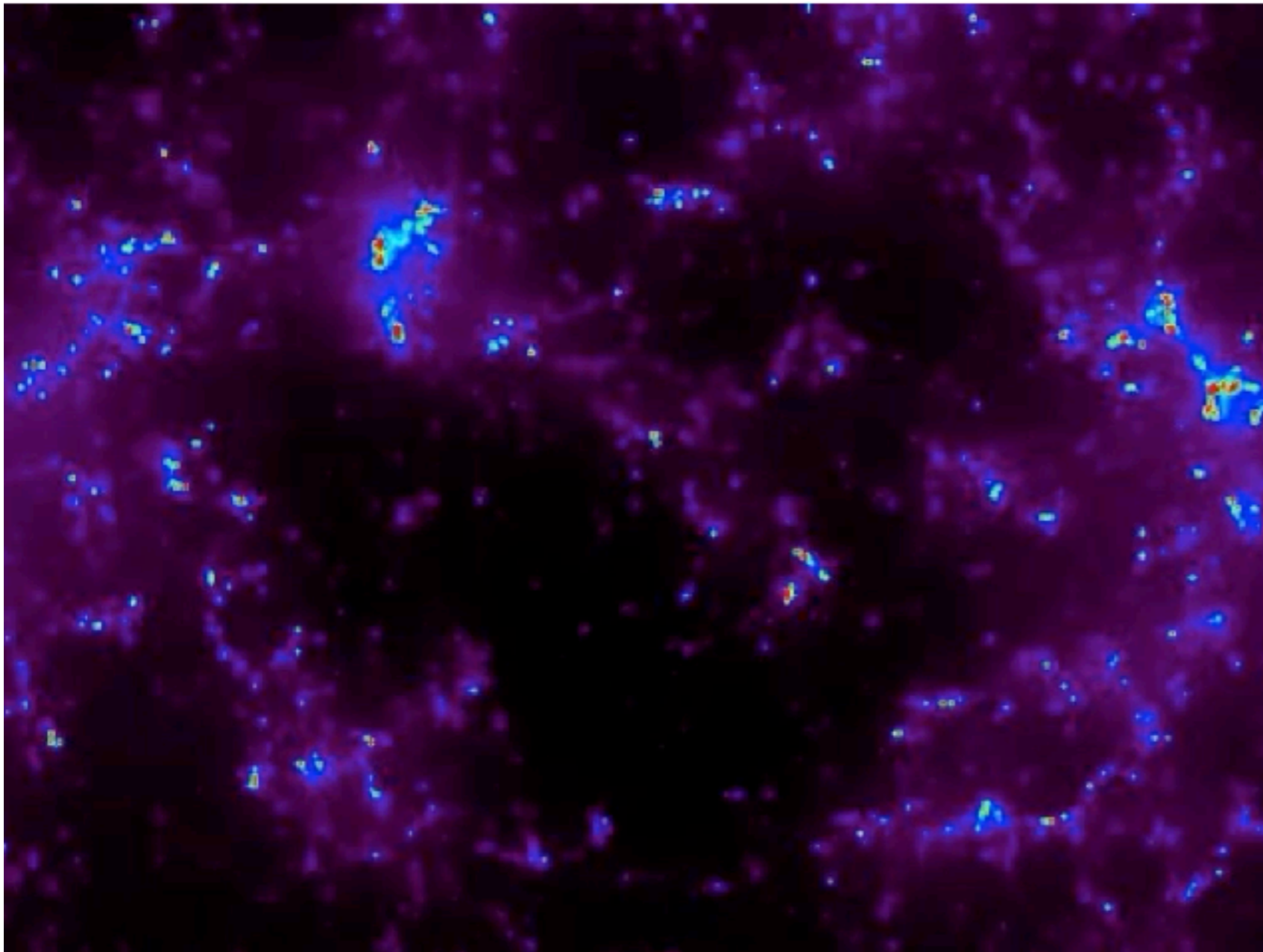
Multi-line modelling

Heneka+17
Heneka, Mesinger 20
Heneka, Cooray 21
Hutter, Heneka+23

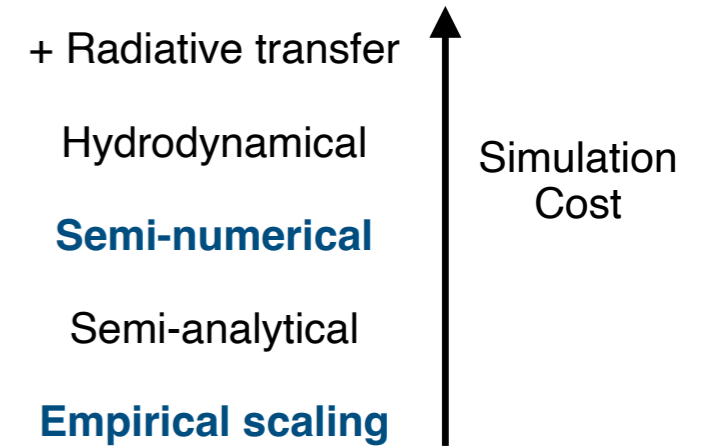
Example: Lyman-alpha

3) Generative methods for simulation

Reionization simulations are costly ... and model a rich signal
Is there a fast way to emulate whole simulations?



Example: Lyman-alpha to 21cm



Multi-line modelling

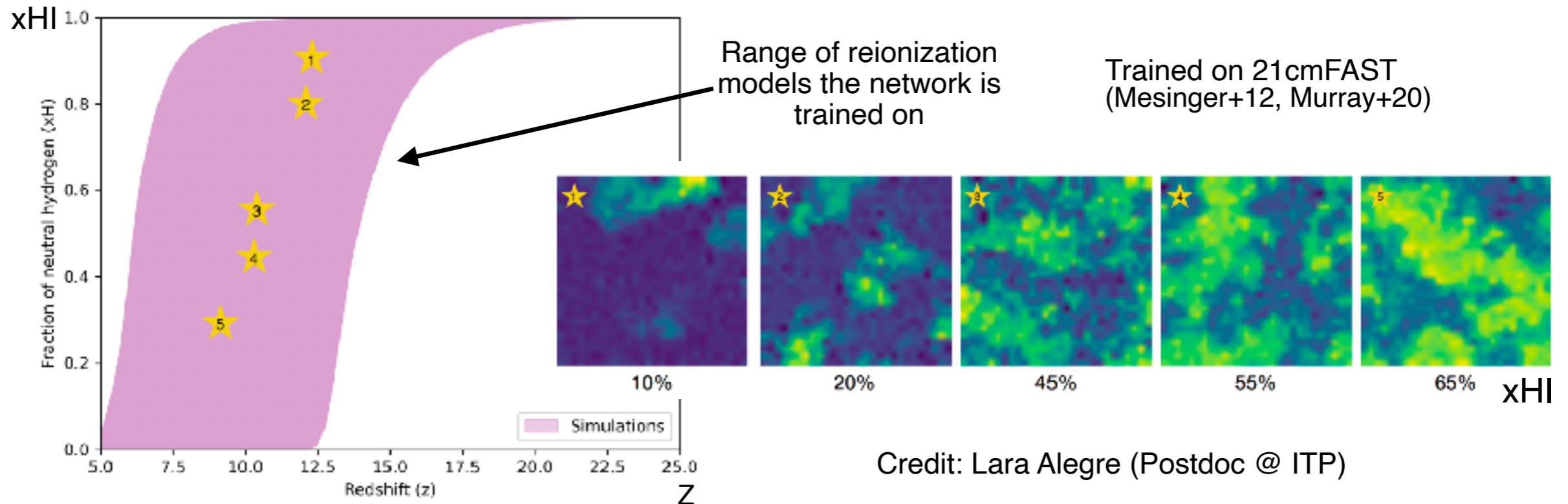
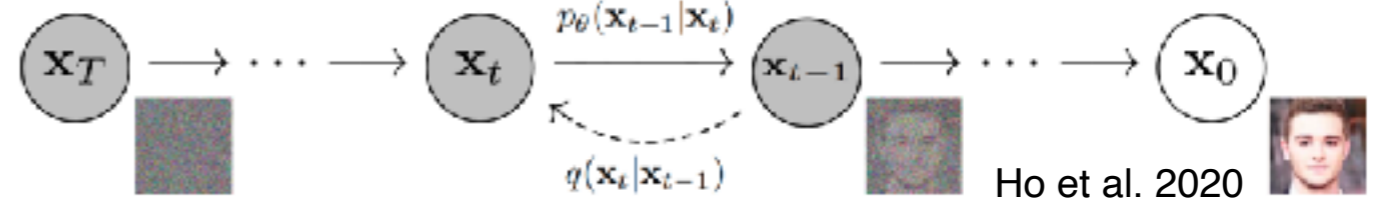
- Heneka+17
- Heneka, Mesinger 20
- Heneka, Cooray 21
- Hutter, Heneka+23

3) Generative methods for simulation (of the 21cm signal)

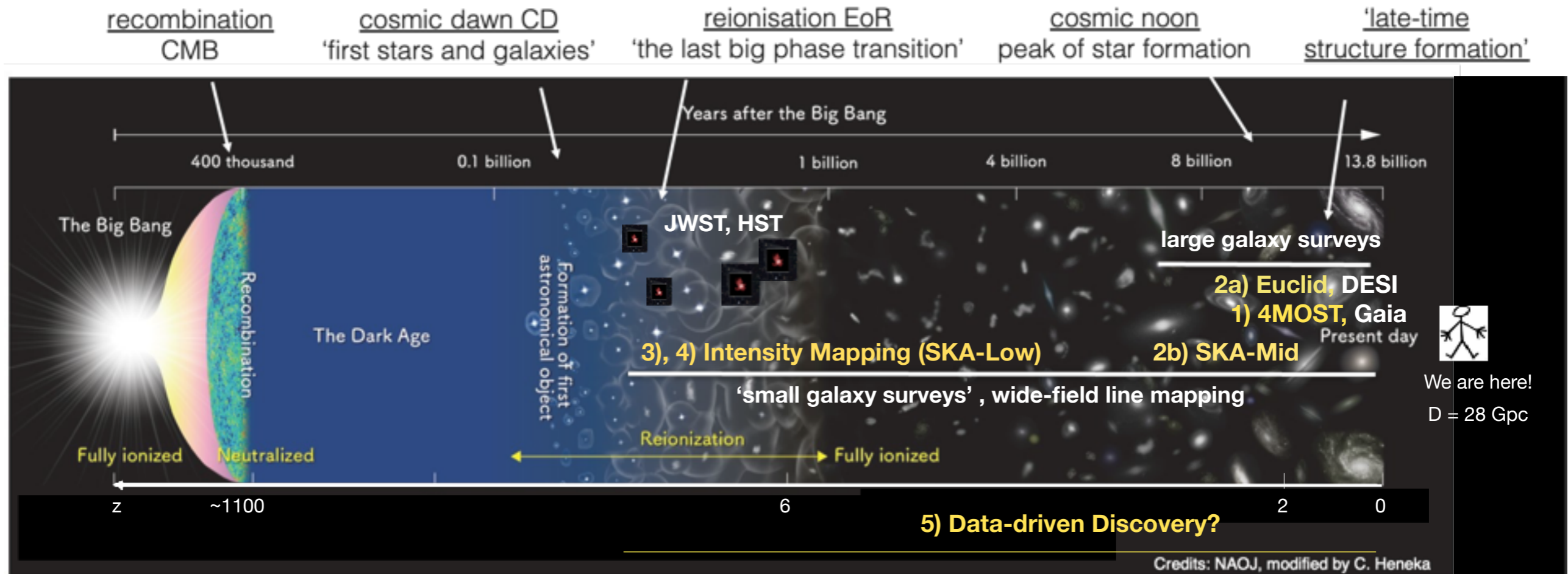
Reionization simulations are costly

Is there a fast way to emulate whole simulations?

Our DL solution: Generation of slices of the 21cm brightness temperature using diffusion models

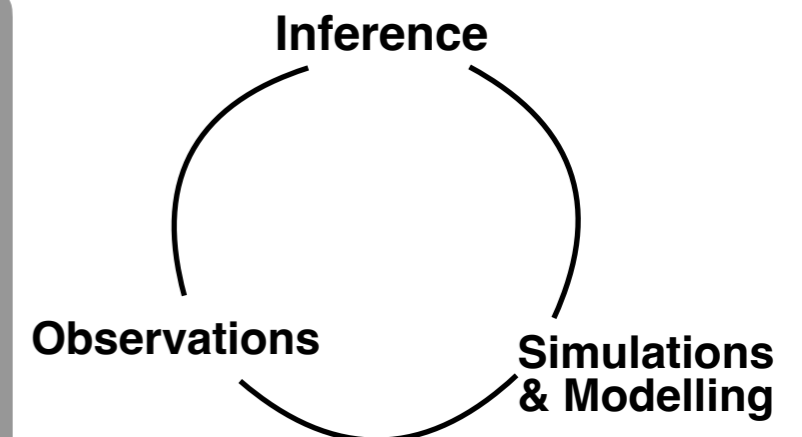


Astronomical and Astrophysical Machine Learning



Highlights in this Lecture

- 1) Classification / Triggering
- 2) Source detection & characterisation
- 3) Simulation-based inference (SBI) in 3D
- 4) Generative methods
- 5) **Data-driven Discovery**



5) Already now: Data-driven discovery

CAMELS

= Cosmology and Astrophysics with Machine Learning Simulations

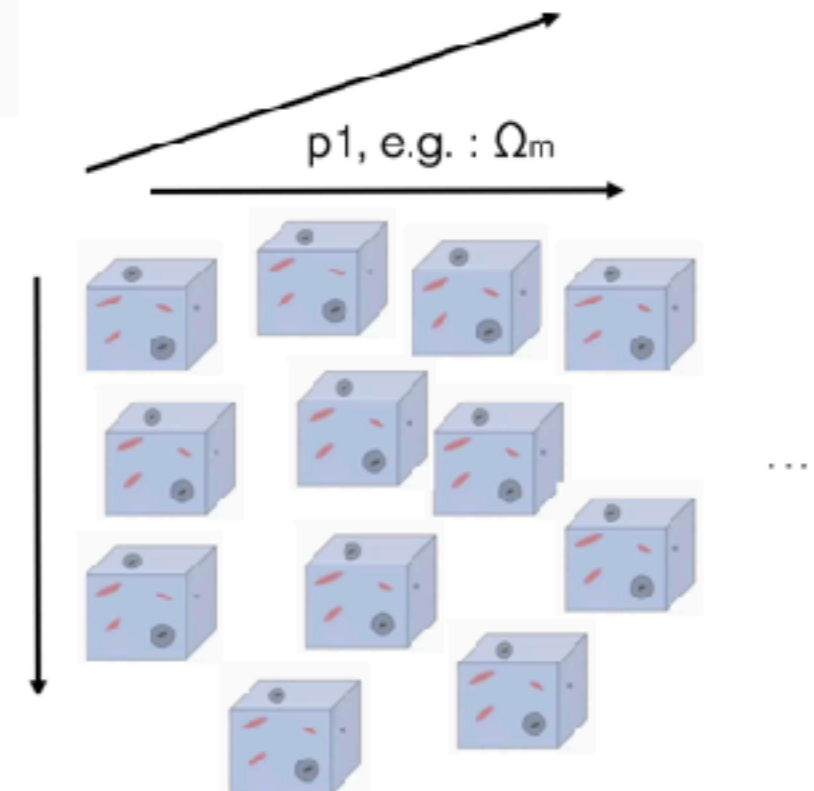
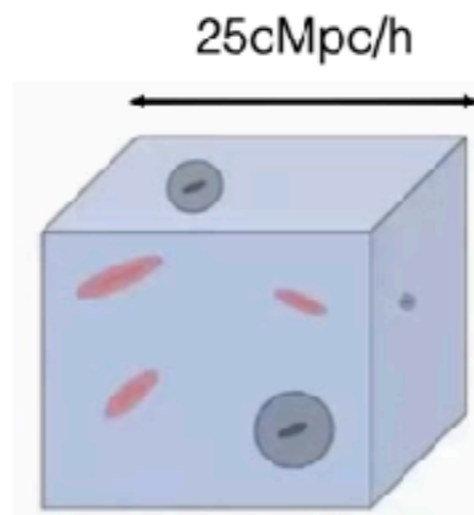
Type	Code	Subgrid model	Simulations
Hydrodynamic	Arepo	IllustrisTNG	1,092
Hydrodynamic	Gizmo	SIMBA	1,092
Hydrodynamic	MP-Gadget	Astrid	1,092
N-body	Gadget-III	—	3,049



<https://camels.readthedocs.io>

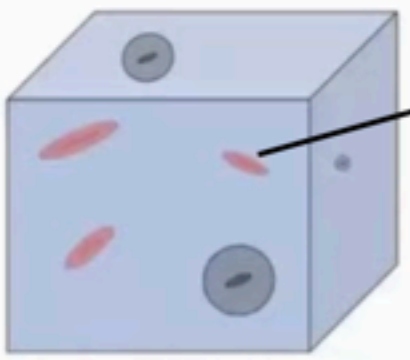
Parameter set:

$$\begin{aligned}
 0.1 &\leq \Omega_m &\leq 0.5 \\
 0.6 &\leq \sigma_8 &\leq 1.0 \\
 0.25 &\leq A_{\text{SN}1} &\leq 4.0 \\
 0.50 &\leq A_{\text{SN}2} &\leq 2.0 \\
 0.25 &\leq A_{\text{AGN}1} &\leq 4.0 \\
 0.50 &\leq A_{\text{AGN}2} &\leq 2.0
 \end{aligned}$$



5) Already now: Data-driven discovery

for each galaxy:



$O(10^4)$ galaxies

$\Omega_m = 0.32$
 $\sigma_8 = 0.79$
 $A_{SN1} = 1.2$
 $A_{SN2} = 0.8$
 $A_{AGN1} = 1.5$
 $A_{AGN2} = 0.7$

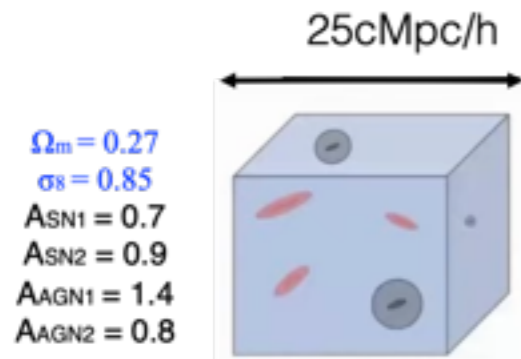
M_*
K
 M_g
 Z_g
 V_{max}
 Z_*
g
 σ_v
 R_*
 M_t
U
 R_t
 R_{max}
SFR
J
V
 M_{bh}

(SIMBA)

How many galaxies do we need to constrain e.g. Ω_m ?
Let's start with one!

5) Already now: Data-driven discovery

How many galaxies do we need to constrain e.g. Ω_m ?

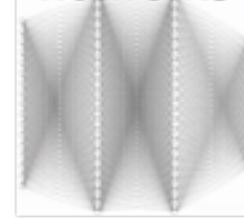


$O(10^4)$ galaxies per cube

random galaxy:
 \longrightarrow
 galaxy properties

- M_*
- K
- M_g
- Z_g
- V_{max}
- Z_*
- g
- σ_v
- R_*
- M_t
- U
- R_t
- R_{max}
- SFR
- J
- V
- M_{bh}

Moment density networks



arXiv:2011.05991

\longrightarrow
 ?

Ω_m
 $\delta\Omega_m$

~10% uncertainty

recovery of matter density for very different masses (and environments)
also: holds for redshifts other than $z=0$

connections on a high dimensional manifold?

= Cosmology and Astrophysics with Machine Learning Simulations



<https://camels.readthedocs.io>

5) Already now: Data-driven discovery

Loss based on moment density networks (MDN) Jeffrey & Wandelt 2011
arXiv:2011.05991

MDN idea: hierarchy of neural regression models (mean \rightarrow variance \rightarrow skewness \rightarrow ...)

We begin by noting that if we find some function of our data $\mathcal{F}(x)$ that minimizes an L_2 loss over the distribution of possible training examples $\{x_i, \theta_i\}$,

$$J_0 = \int \|\theta - \mathcal{F}(x)\|^2 p(x, \theta) dx d\theta, \quad (4)$$

then \mathcal{F} , which we represent as a neural network, evaluated for the observed data is the mean of the posterior distribution $\mathcal{F}(x_{obs}) = \langle \theta \rangle_{\theta|x_{obs}}$. It is therefore possible to create a hierarchy of networks

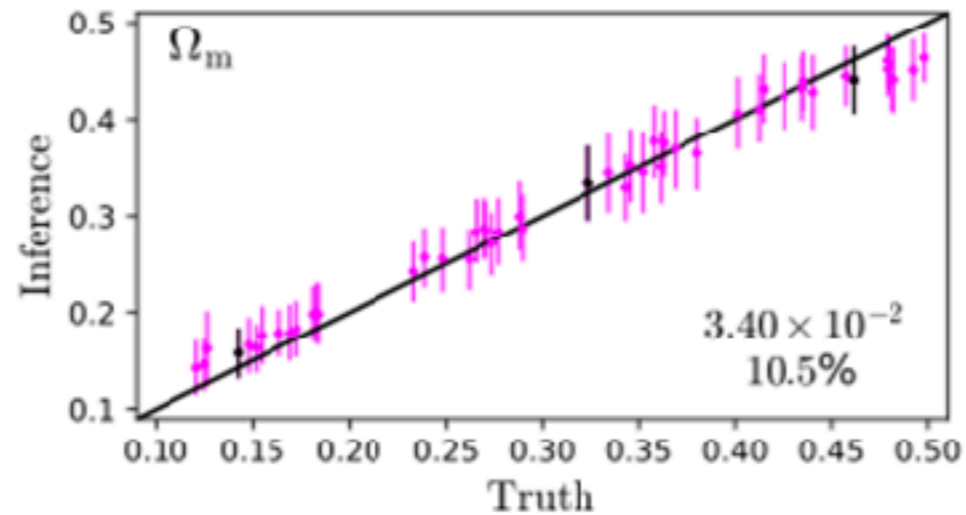
In practice we minimise the following loss function:

$$\mathcal{L} = \sum_{i=1}^6 \log \left(\sum_{j \in \text{batch}} (\theta_{i,j} - \mu_{i,j})^2 \right) + \sum_{i=1}^6 \log \left(\sum_{j \in \text{batch}} \left((\theta_{i,j} - \mu_{i,j})^2 - \sigma_{i,j}^2 \right)^2 \right)$$

Our model $F(x)$:
CNN layers (19)
+ dense (2)

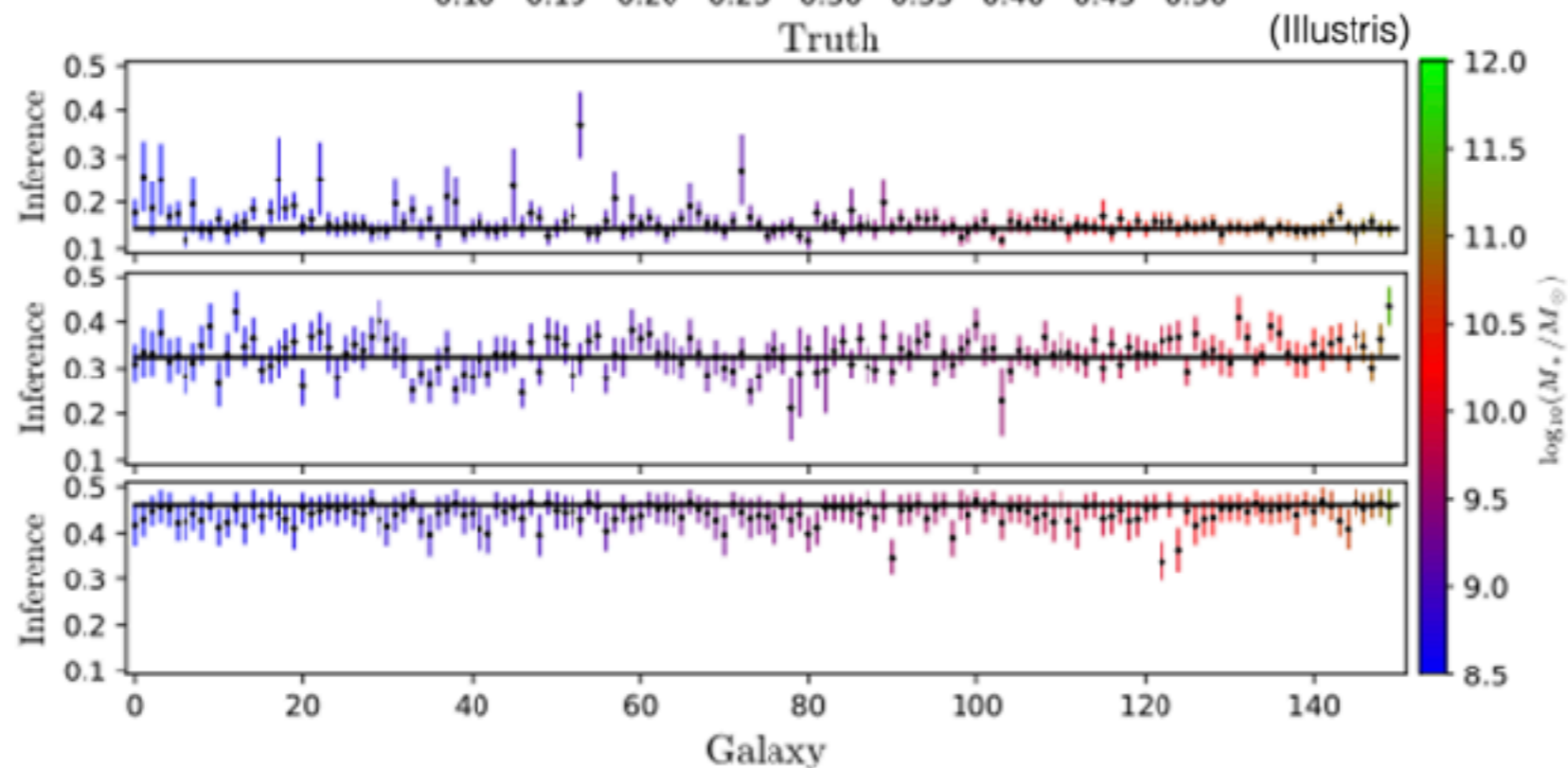
Villaescusa-Navarro+
arXiv:2109.10915

5) Already now: Data-driven discovery



Introspection via:

- Attribution, one-by-one retraining
- SHAP values (Shapley Additive exPlanations)
- Principal Component Analysis (PCA)
- Correlation strengths



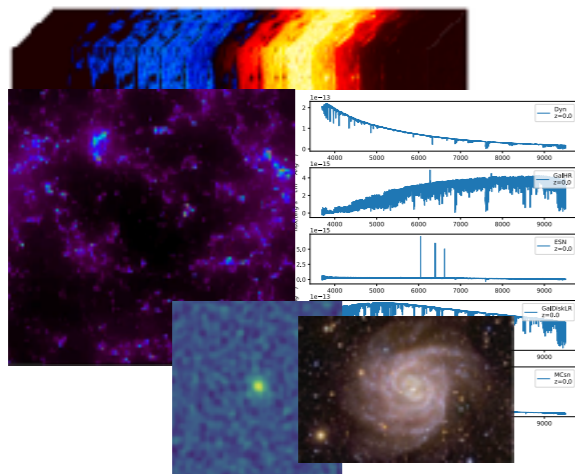
→ recovery of matter density for very different masses (and environments)
also: holds for redshifts other than $z=0$

Summary: Where we stand

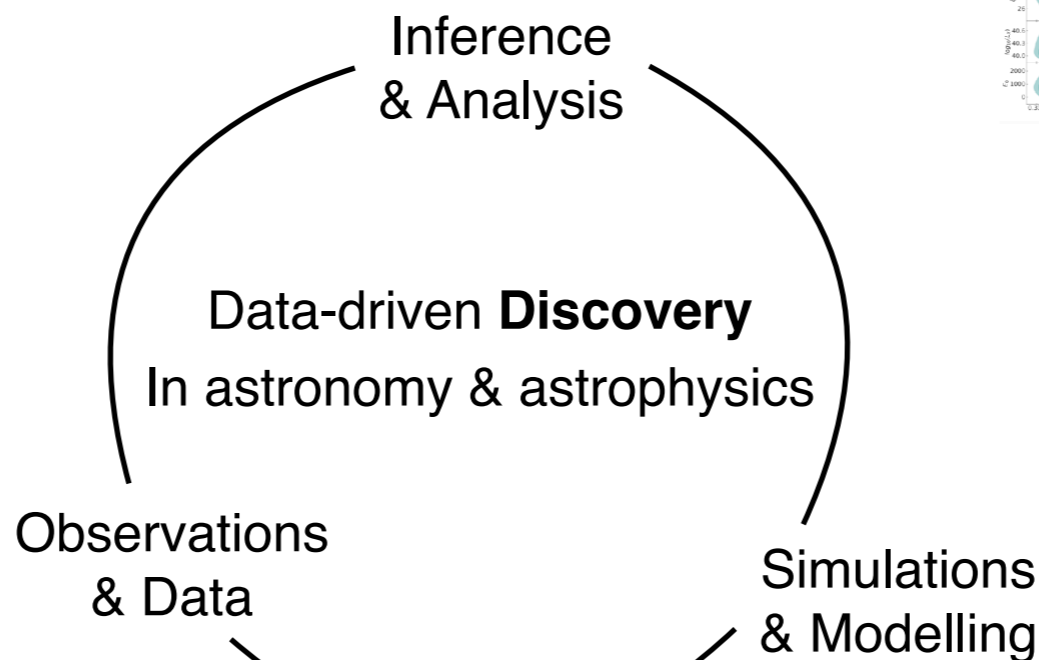
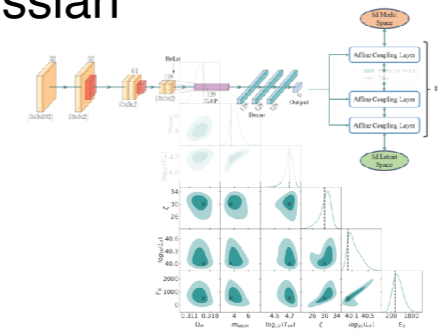
Goal: Understanding & discovery

Detection & Characterisation
Unbiased measurements from diverse sources (galaxies)

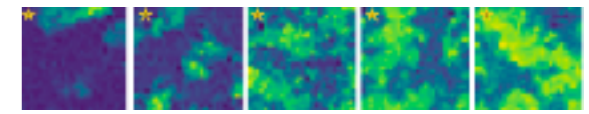
Classification
Online classifier and triggering



Inference
SBI in 3D from non-Gaussian tomographic data



Simulation
Produce large range of reionization topologies

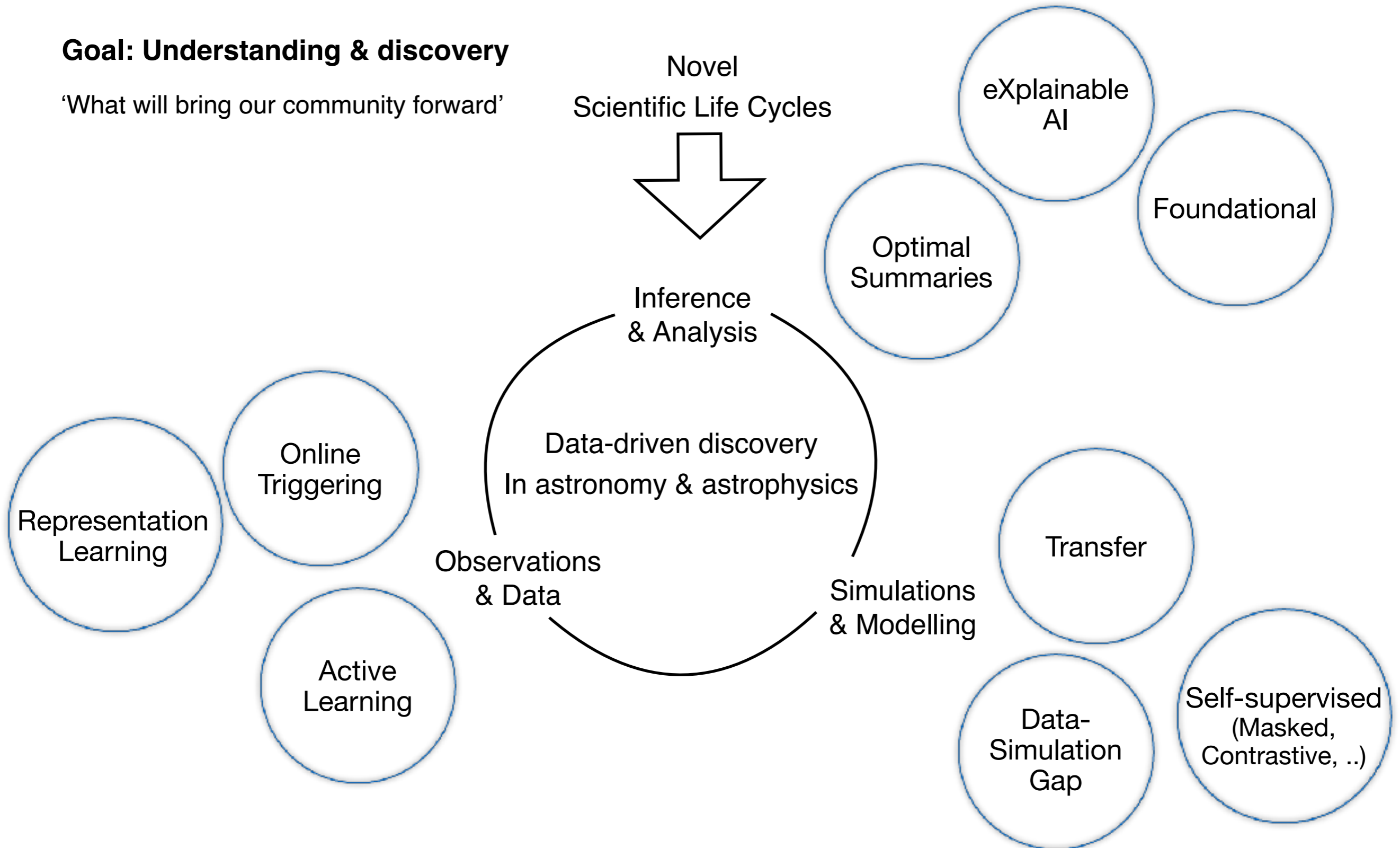


- Select publications:
- Neutsch, Heneka, Brüggem 22, arXiv:2201.07587
 - Schosser, Heneka, Plehn, arXiv:2401.04174
 - Hartley+ 23 (incl. Heneka), arXiv:2303.07943
 - Heneka 23, arXiv:2311.17553
 - Boucaud, Huertas-Company, Heneka+ 20, arXiv:1905.01324
 - Zhong, Napolitano, Heneka+24, arXiv:2311.04146

Research plans: Where we are going

Goal: Understanding & discovery

'What will bring our community forward'



Example: Robust Foundational models for inference

Data-Simulation Gap

Generation & modelling

1) simulations, mock catalogues + empirical relations

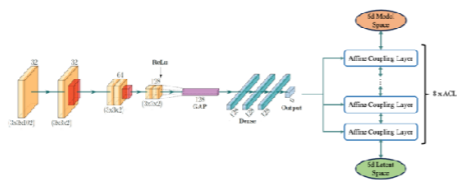
Mock observations

Augmentation, semi-supervised, importance weighing

Optimal Summaries

Foundational ?

Simulation-based inference



2) Transfer

3) Inference

LOFAR observations



LoTSS maps

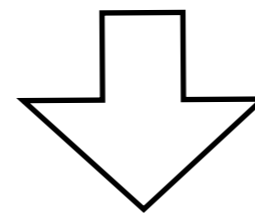
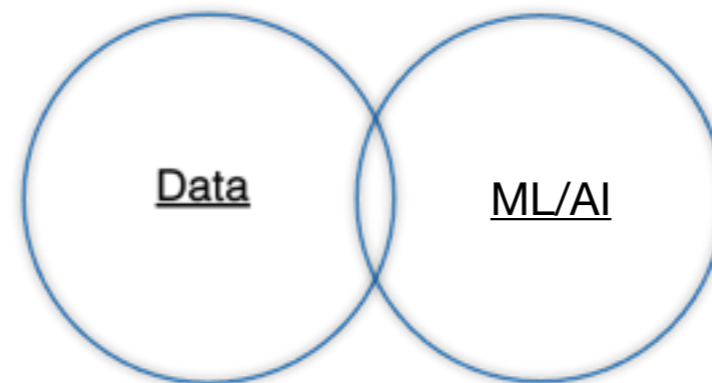
Currently:

Comparison to 'random mocks'

$$C_l^{gg} + \text{jackknife} + \text{MCMC}$$

Research plans: where we are moving towards

Use data from astronomical surveys for **data-driven discovery**.



Goals:

Map-based, multi-channel, '3D' approach to astrophysics
Representation learning, beyond summaries
Automated data mining, anomaly search
New signatures & discovery

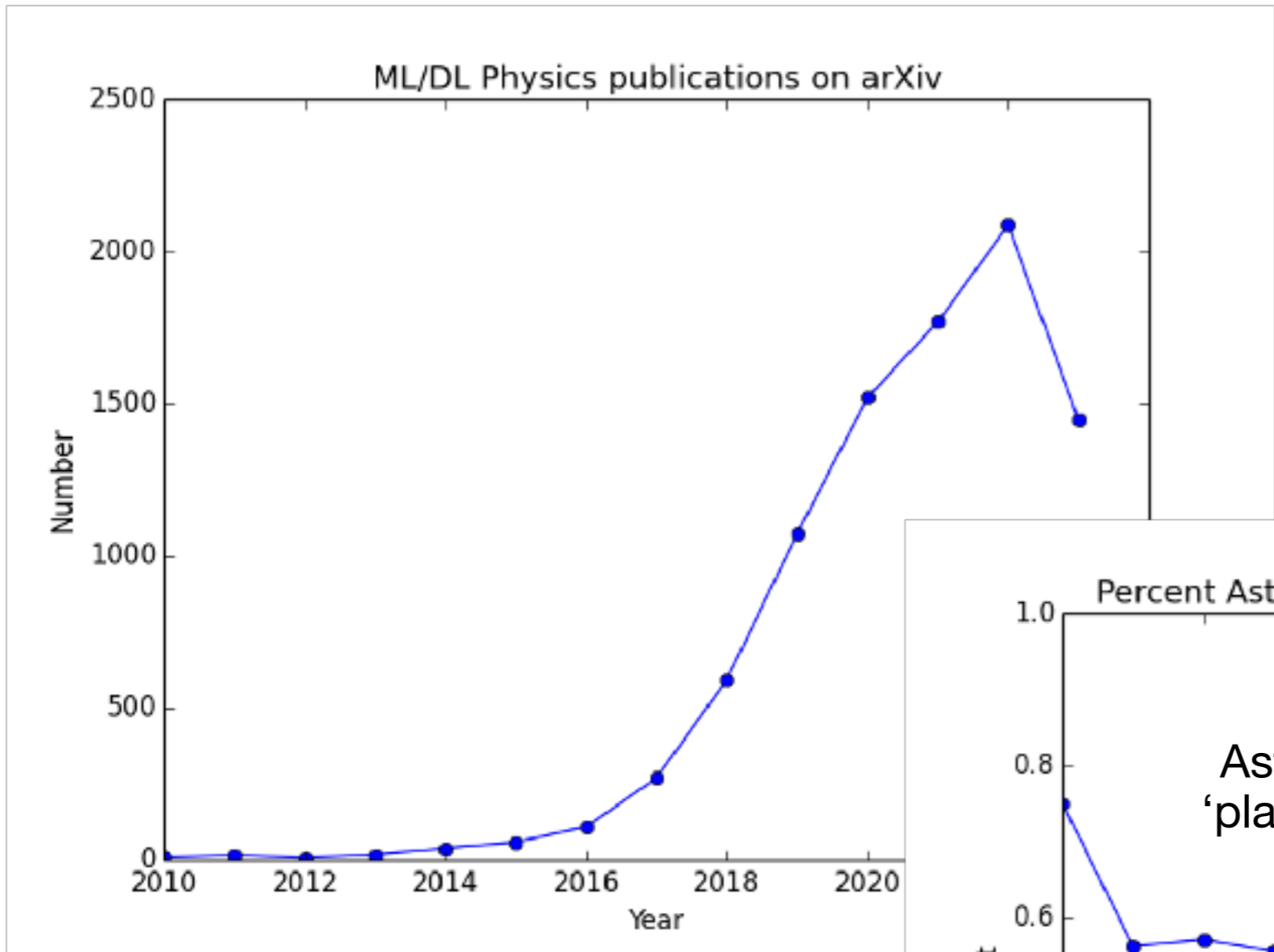
Outcomes:

Fast and expressive simulators
Reliable error estimates & inference
Interpretability
Active interaction and learning
'human in the loop'
New signatures & discovery

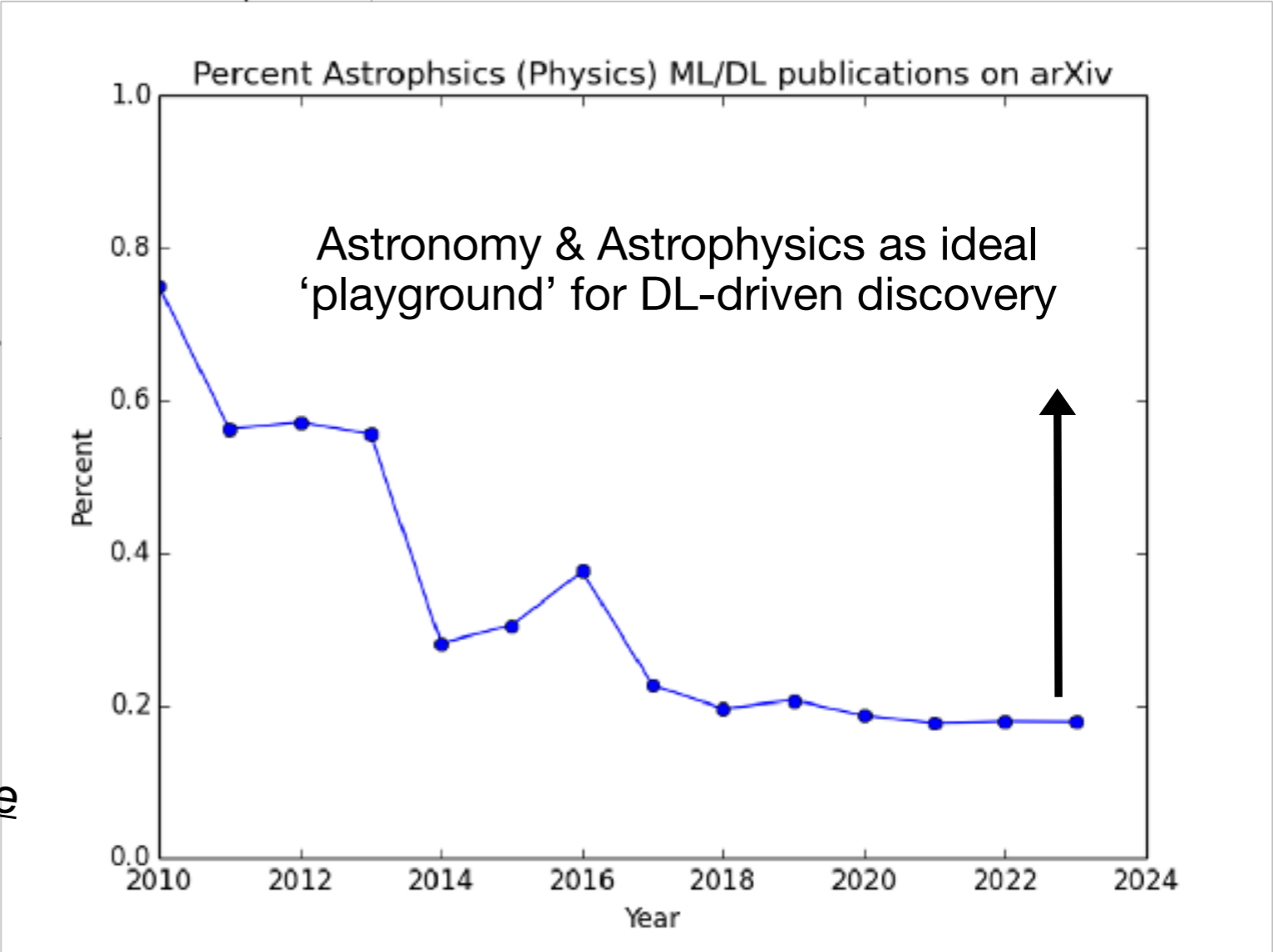
Novel
Scientific
Life Cycles

Based on
foundational models

ML/DL/AI has come to stay when dealing with astronomical data.



Some advertisement:
We have 'endless' open-source public data



Astronomy & Astrophysics as ideal 'playground' for DL-driven discovery

Thank you for your attention!
heneka@thphys.uni-heidelberg.de