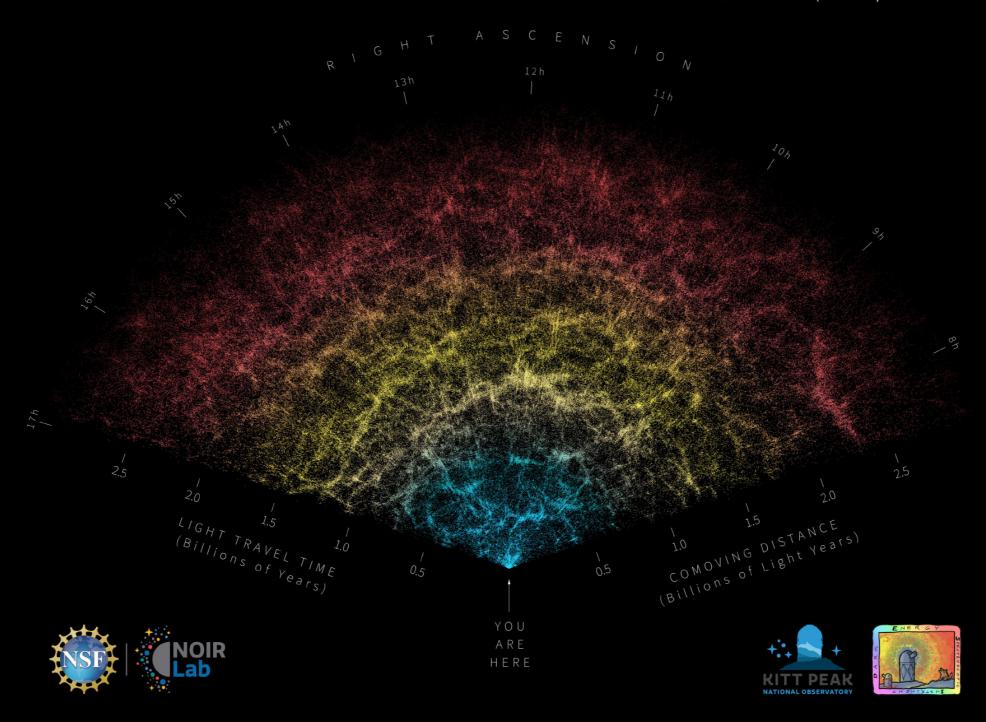


A BIT ABOUT ME

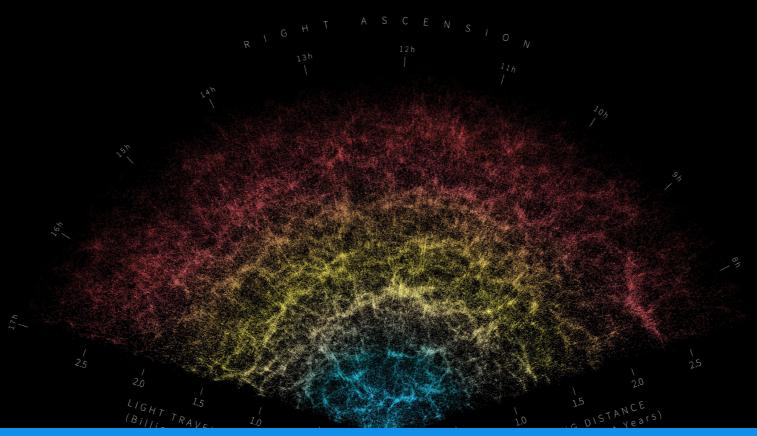
- PhD in Astrophysics from Edinburgh, UK
- Worked at Kaggle developing machine learning competitions, this included
 - Observing dark worlds
 - Galaxy Zoo
 - Titanic
- Post-doc Switzerland
- Fellowship Leiden, NL
- Machine Learning consultant at Terres des Hommes
- New position at EPFL
- Machine Learning consultant at Prophy
- Machine Learning consultant at TruthEngine

Astrophysics differs from almost any science

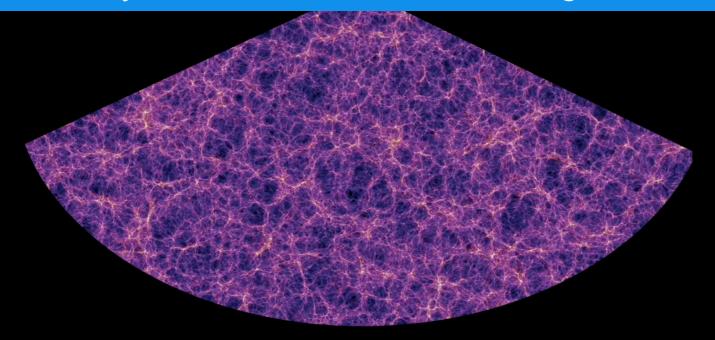
DISTRIBUTION OF NEARBY GALAXIES MAPPED BY THE DARK ENERGY SPECTROSCOPIC INSTRUMENT (DESI)



DISTRIBUTION OF NEARBY GALAXIES MAPPED BY THE DARK ENERGY SPECTROSCOPIC INSTRUMENT (DESI)

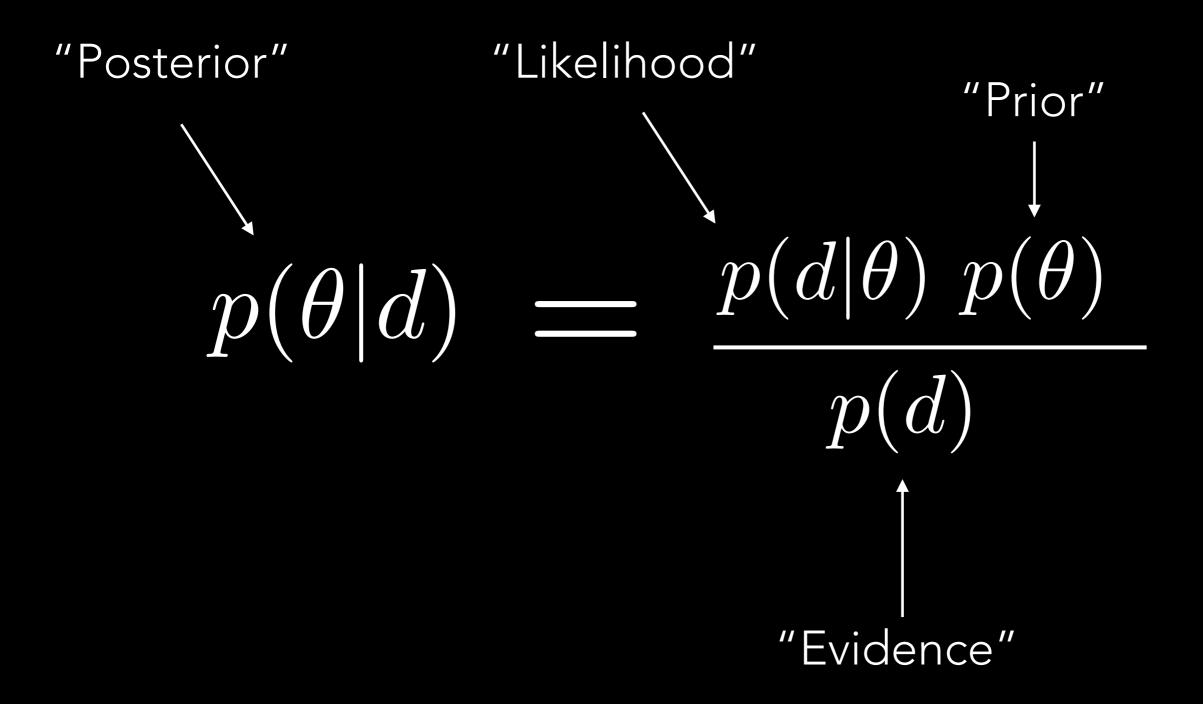


What is the probability of our model of the Universe given what we observe?

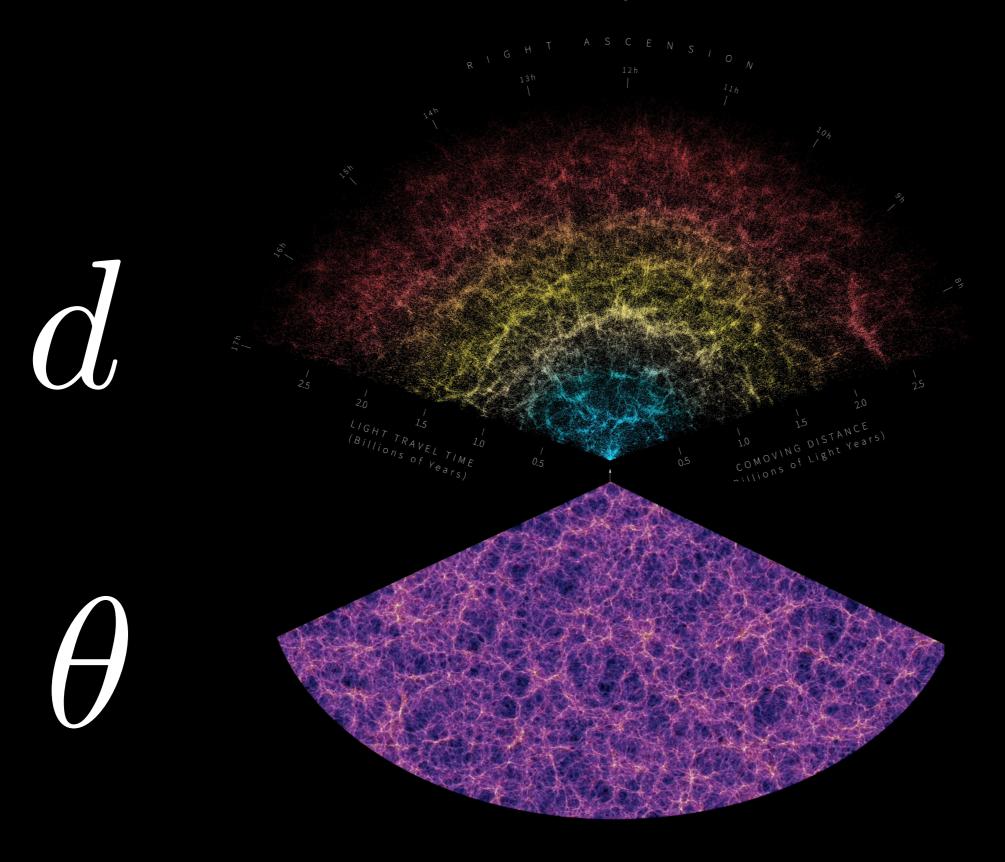


Distribution of matter as seen by the millennium simulation

Bayes theorem as a way to compare what we model to what we see.

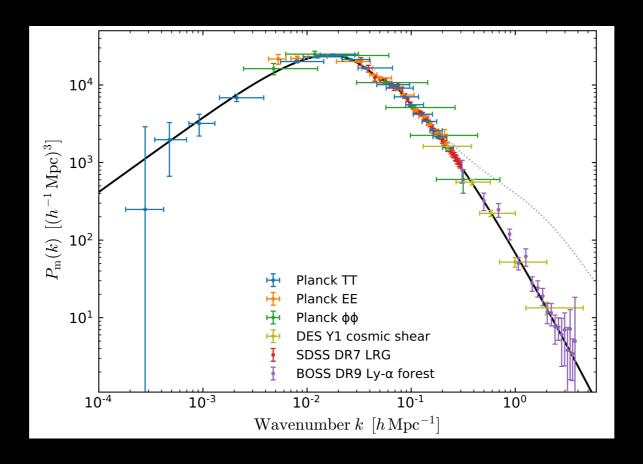


Classical astrophysical inference

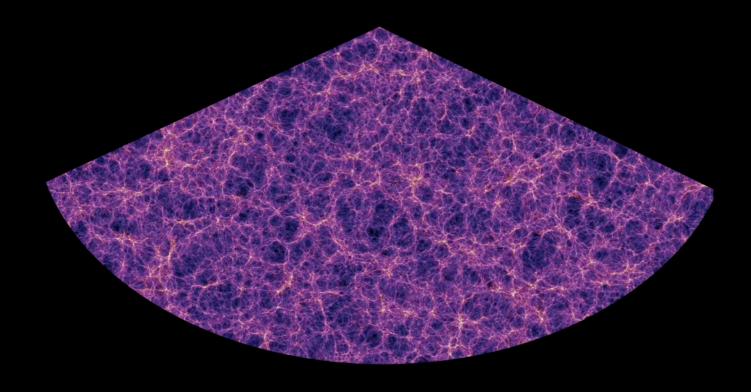


Distribution of matter as seen by the millennium simulation

d

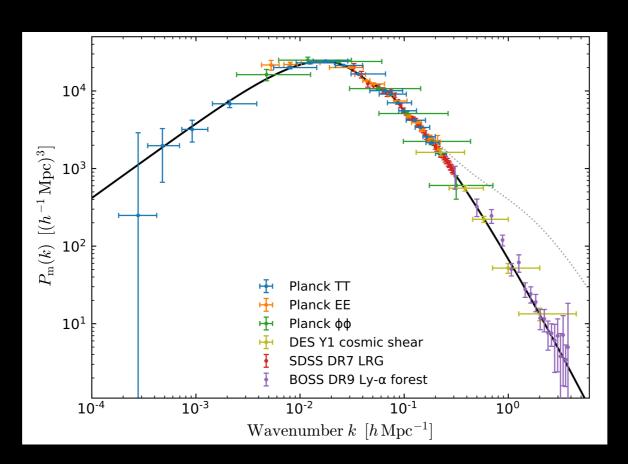


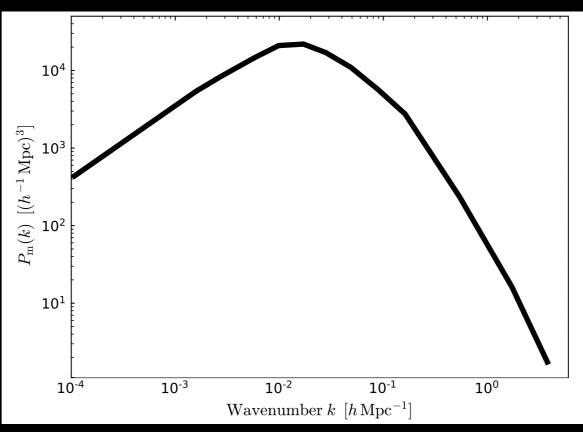




d

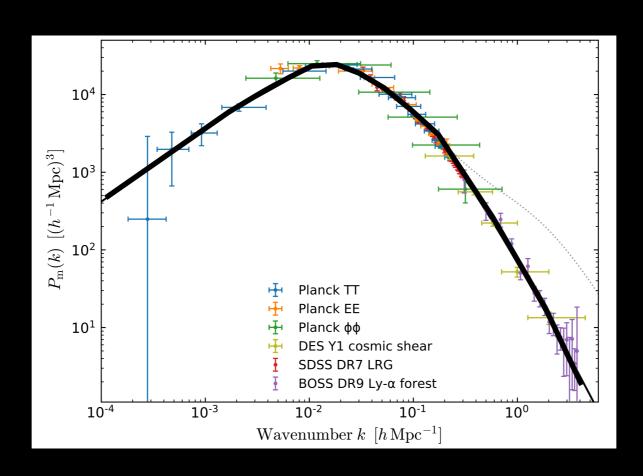






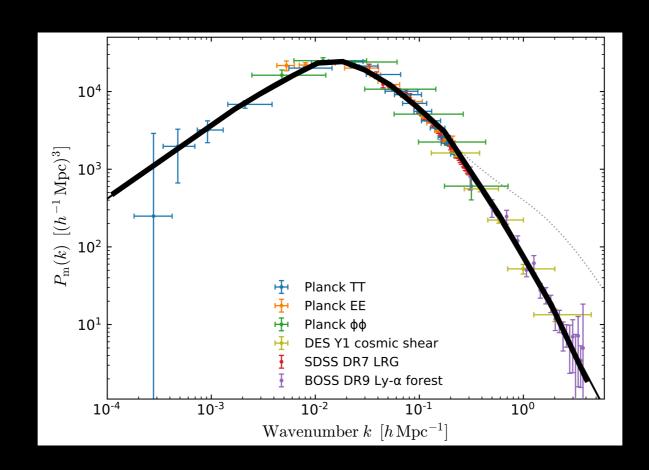
Compare the likelihood of the model given the data to find the best fitting model





Compare the likelihood of the model given the data to find the best fitting model





Assuming an analytical model
Assuming a Gaussian likelihood function.
With covariances from theory or simulations

- Classical inference requires analytical models that can be passed in to a MCMC
- The number of posterior calls is normally very large.
- Likelihoods are not Gaussian.
- Direct comparison of simulations to observations would be impossible in the current situation.
- Astrophysical simulations are often computationally expensive

- Classical inference requires analytical models that can be passed in to a MCMC
- The number of posterior calls is normally very large.
- Likelihoods are not Gaussian.
- Direct comparison of simulations to observations would be impossible in the current situation.
- Astrophysical simulations are often computationally expensive

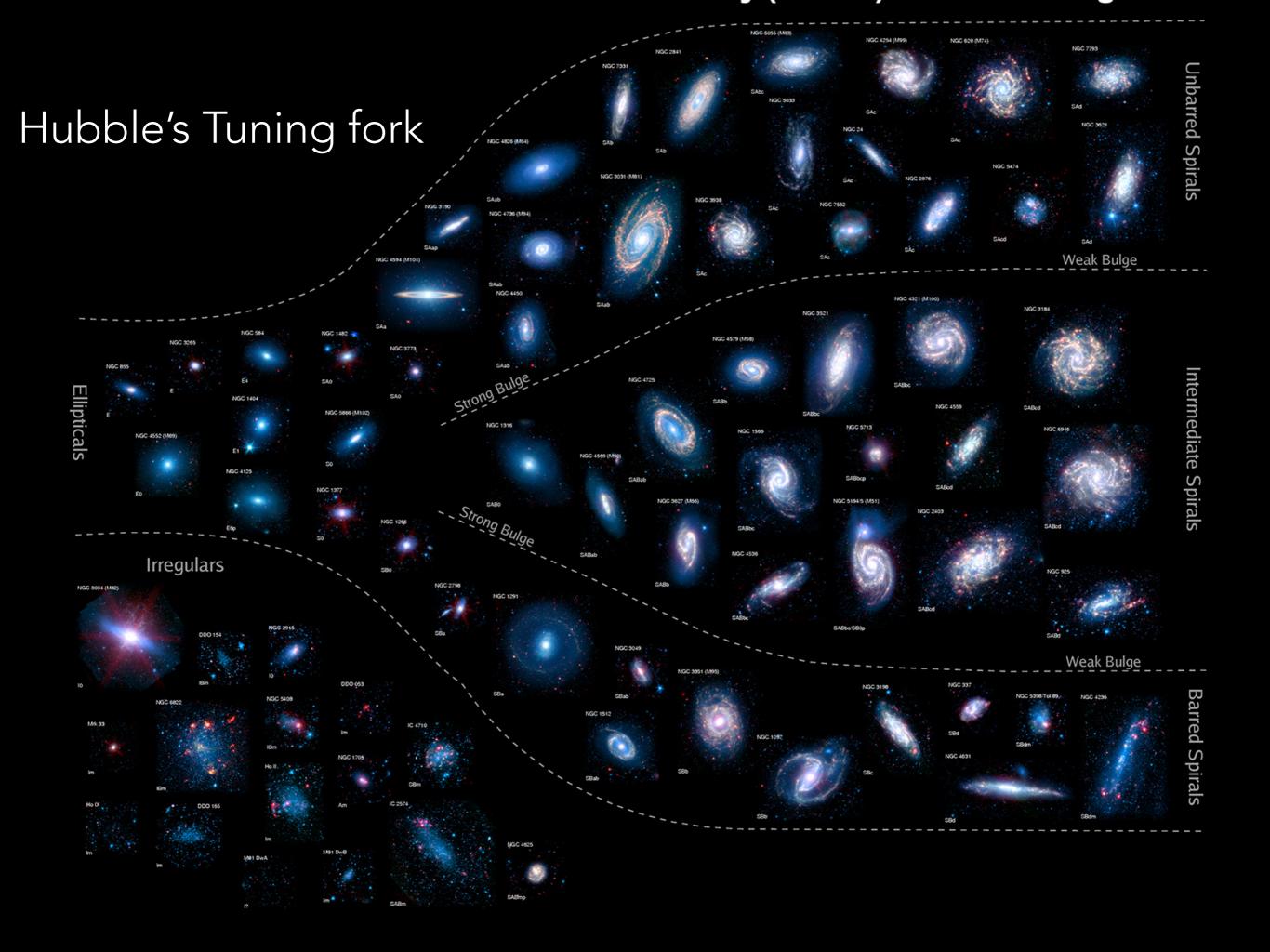
- Classical inference requires analytical models that can be passed in to a MCMC
- The number of posterior calls is normally very large.
- Likelihoods are not Gaussian.
- Direct comparison of simulations to observations would be impossible in the current situation.
- Astrophysical simulations are often computationally expensive

- Classical inference requires analytical models that can be passed in to a MCMC
- The number of posterior calls is normally very large.
- Likelihoods are not Gaussian.
- Direct comparison of simulations to observations would be impossible in the current situation.
- Astrophysical simulations are often computationally expensive

- Classical inference requires analytical models that can be passed in to a MCMC
- The number of posterior calls is normally very large.
- Likelihoods are not Gaussian.
- Direct comparison of simulations to observations would be impossible in the current situation.
- Astrophysical simulations are often computationally expensive

CLASSICAL INFERENCE REQUIRES ANALYTICAL MODELS
THAT CAN BE PASSED IN TO A MCMC WITH MANY CALLS
TO A LIKELIHOOD FUNCTION

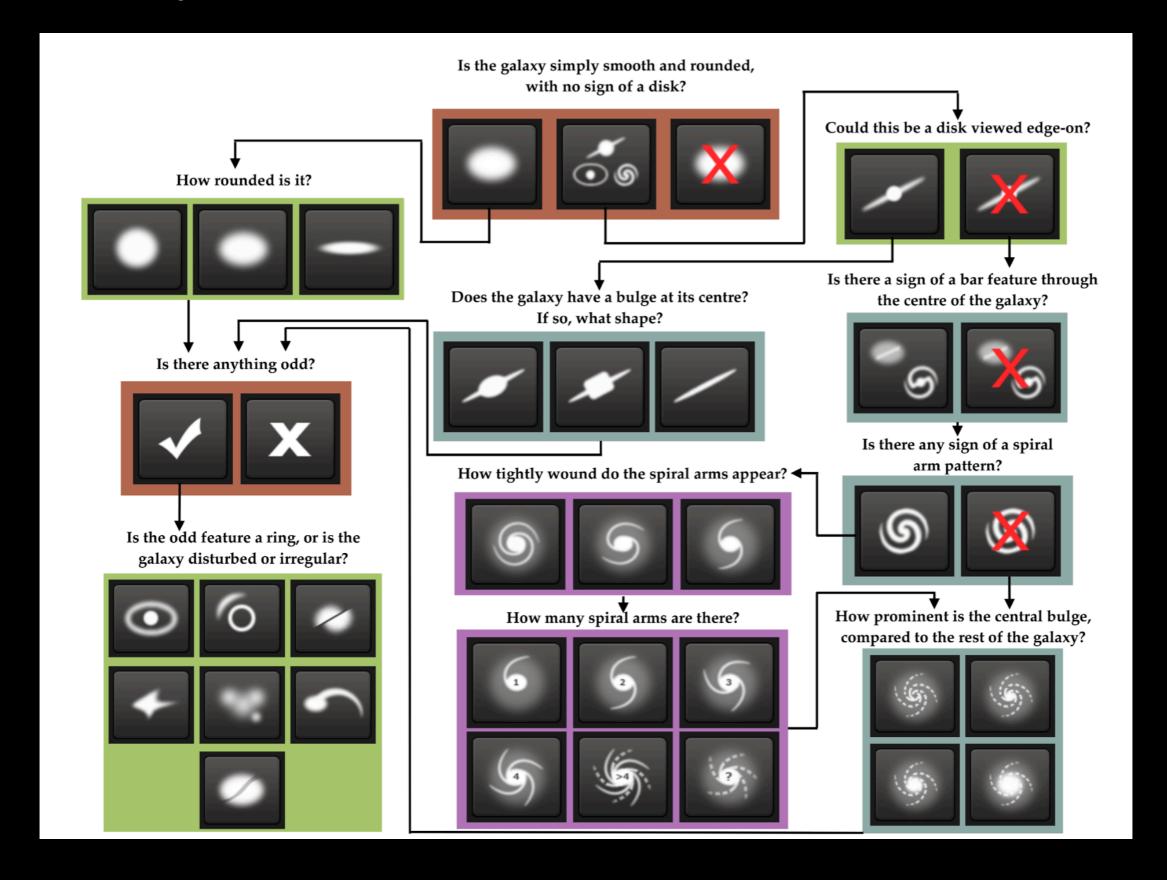
EXAMPLE 1: GALAXY CLASSIFICATION

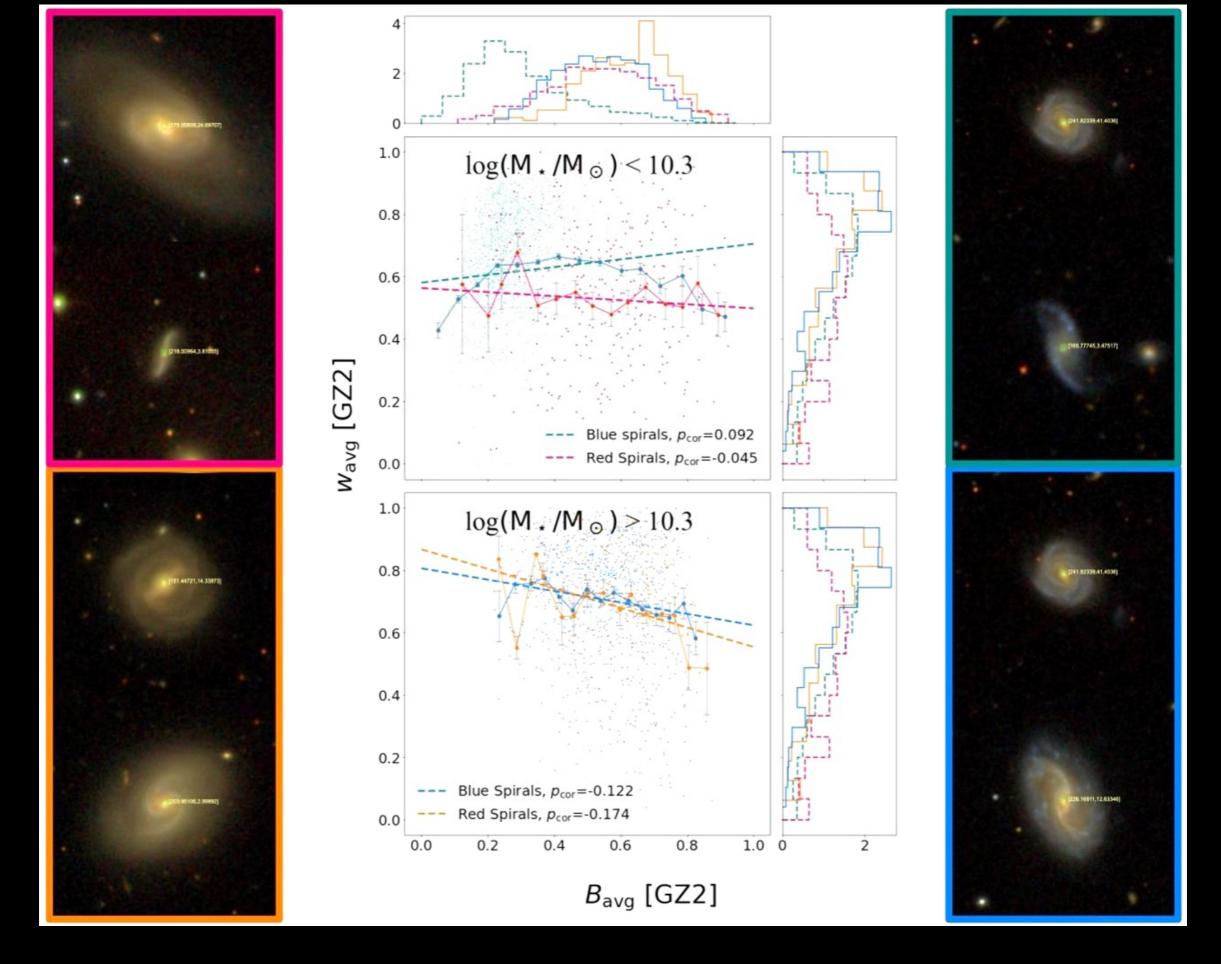




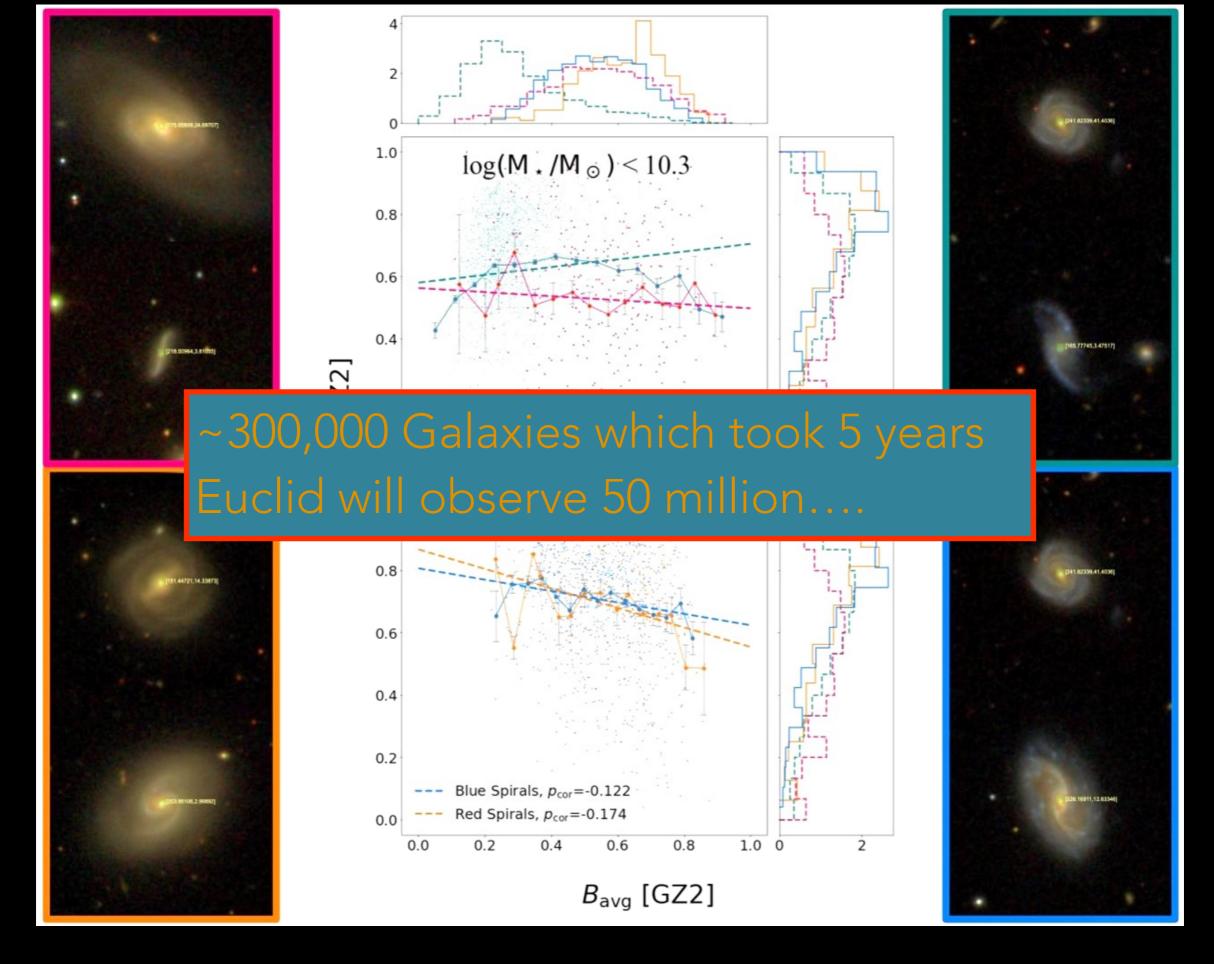
Lintott et al 2008

The Galaxy Zoo decision tree.





Mengistu & Master 2023



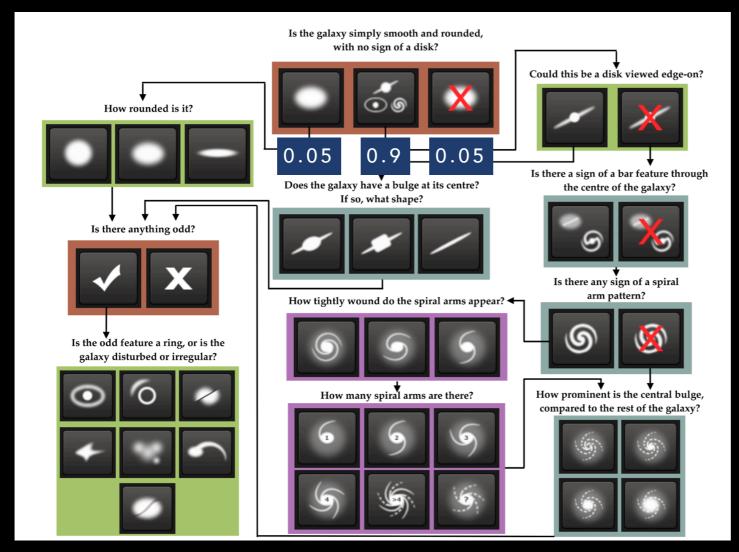
Mengistu & Master 2023

kaggle

https://www.kaggle.com/c/galaxy-zoo-the-galaxy-challenge/data

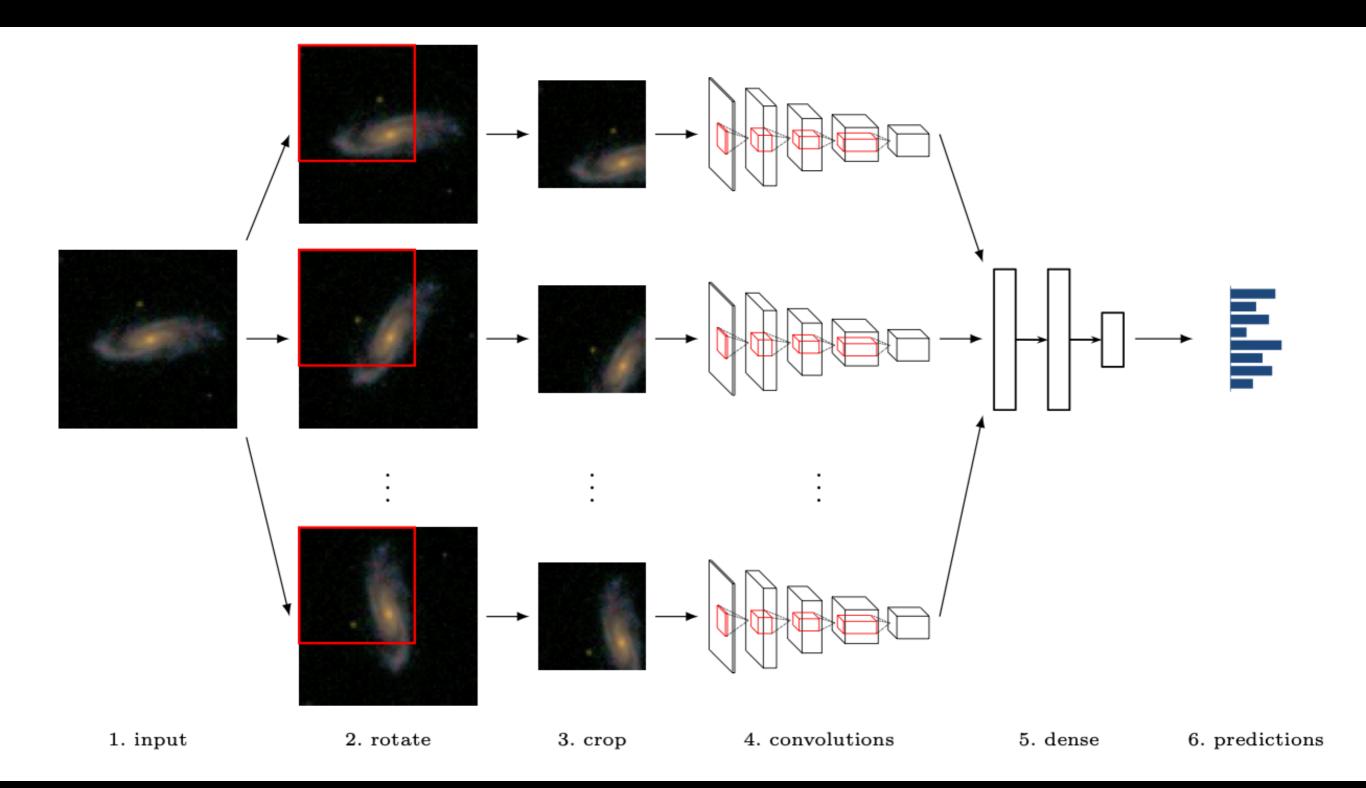
Rotationally invariant convolutional neural networks to predict what a classifier would measure.



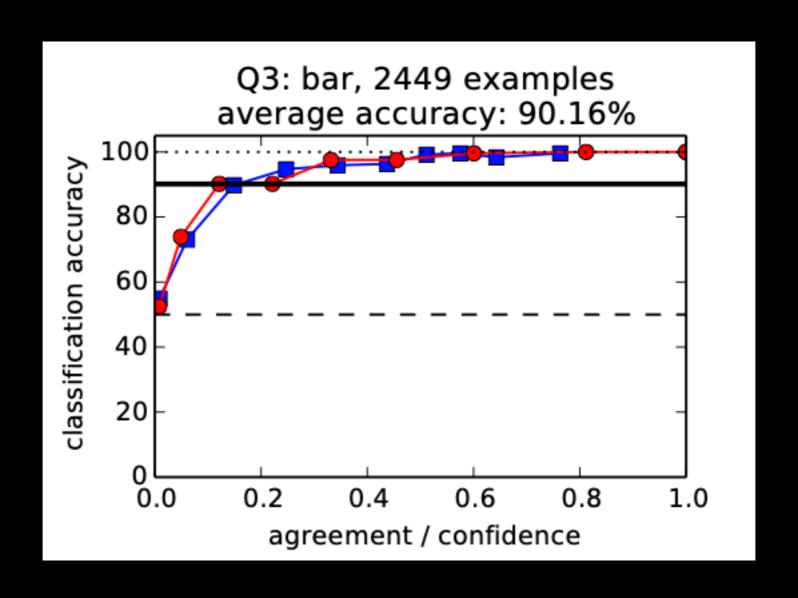


$$e(\hat{p}, p) = \sqrt{\sum_{k=1}^{37} (\hat{p} - p)^2}$$

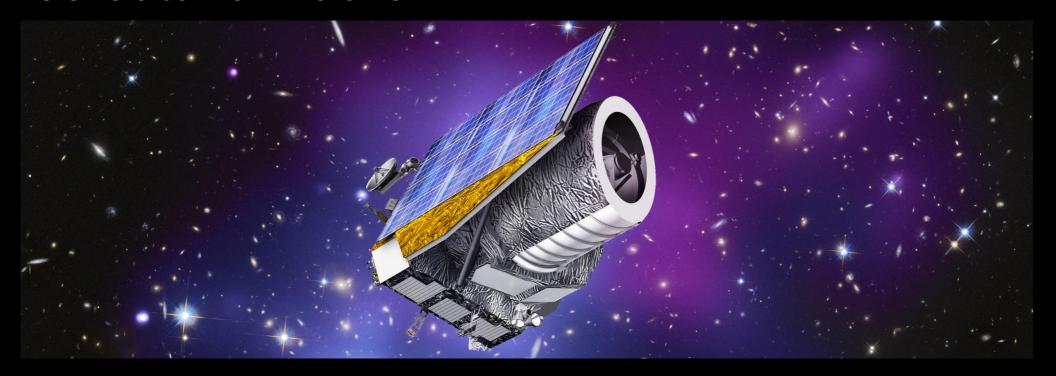
The winning algorithm



Winning the game or doing science?

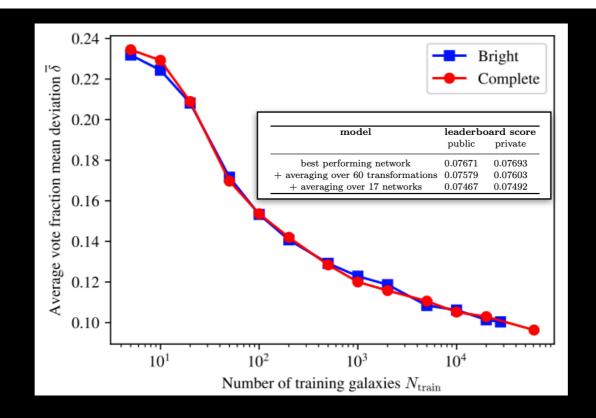


Labelled data for Euclid



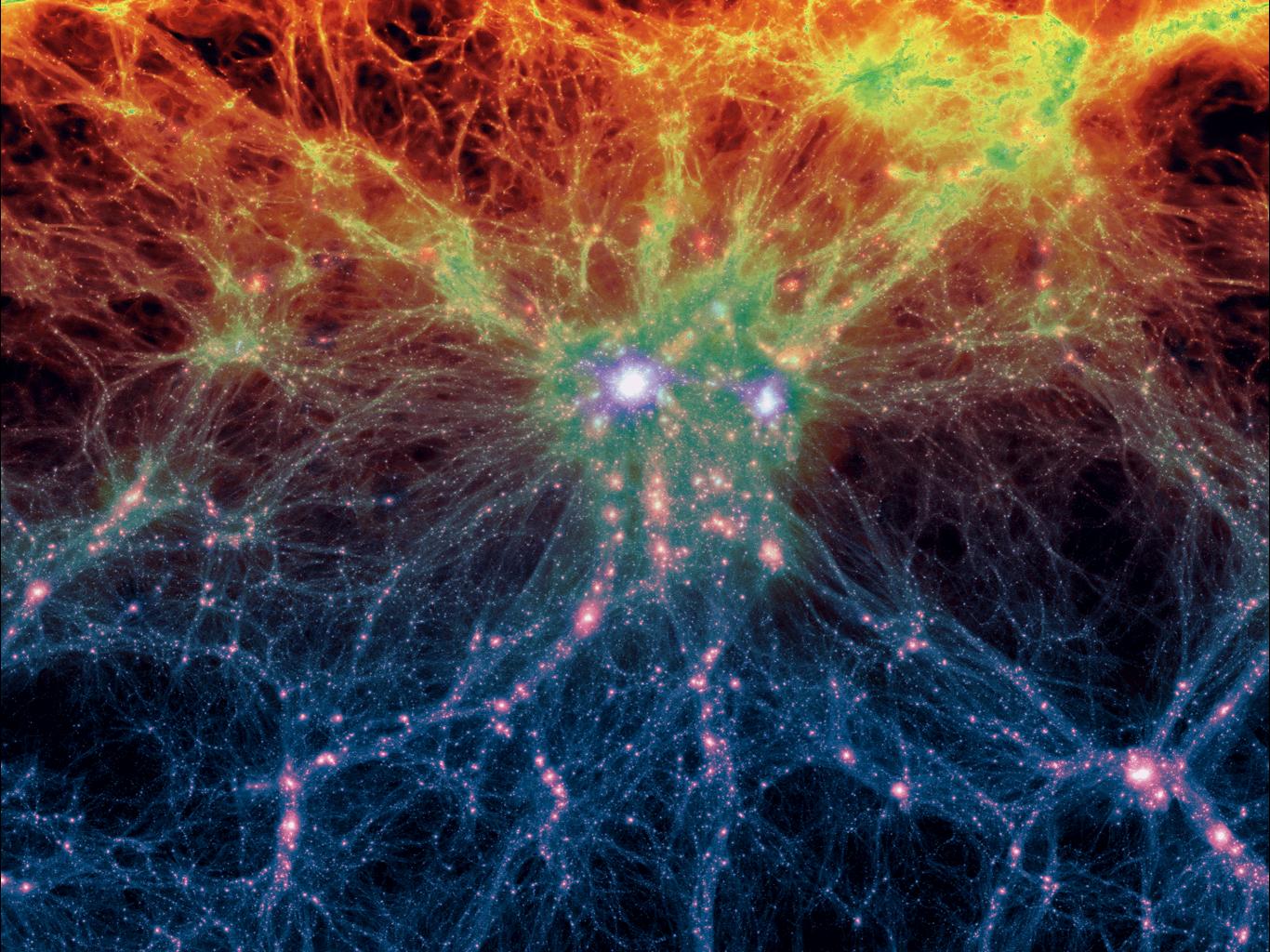
Euclid preparation

XLIII. Measuring detailed galaxy morphologies for Euclid with machine learning

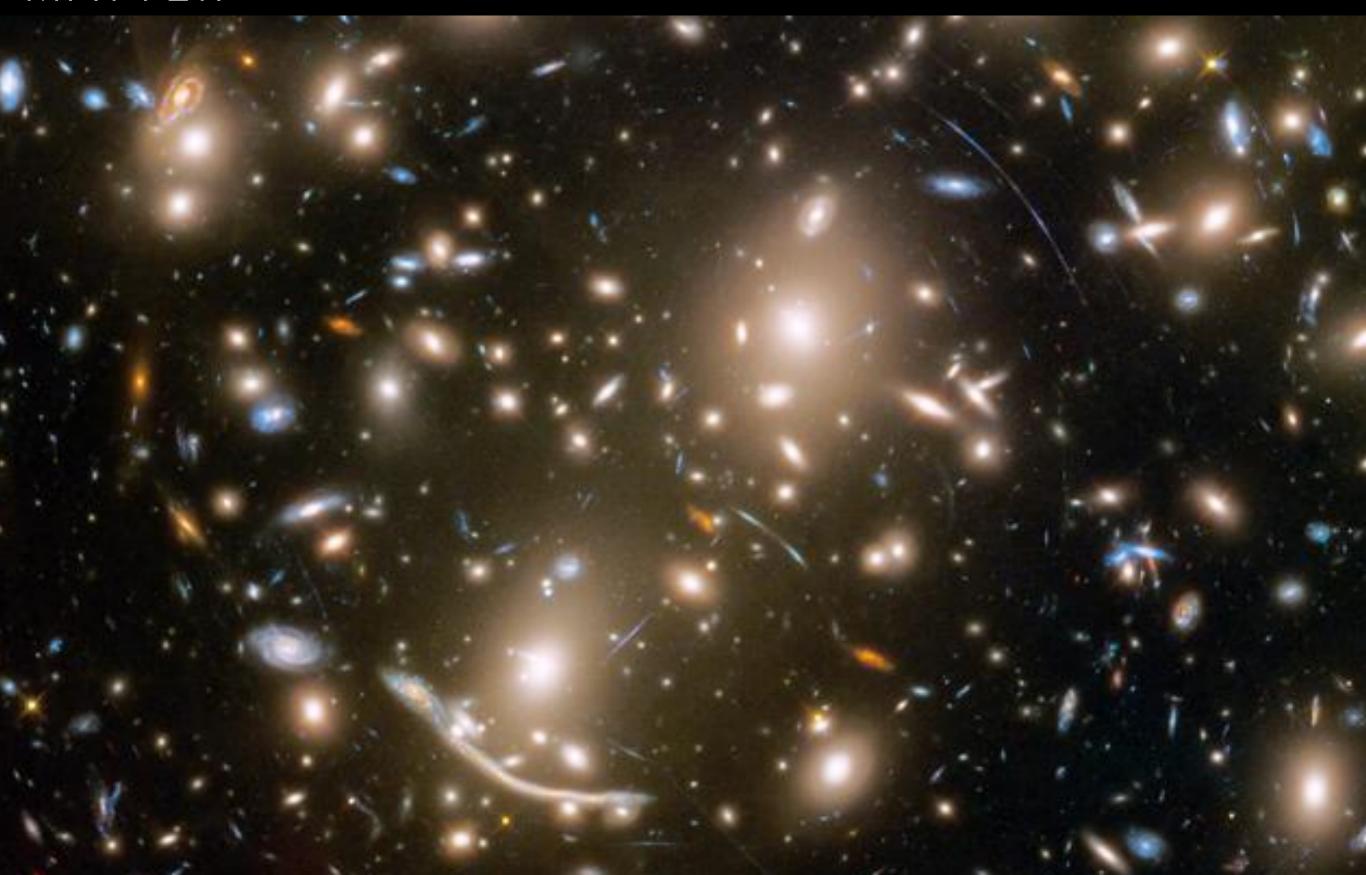


CLASSICAL INFERENCE REQUIRES ANALYTICAL MODELS
THAT CAN BE PASSED IN TO A MCMC WITH MANY CALLS
TO A LIKELIHOOD FUNCTION

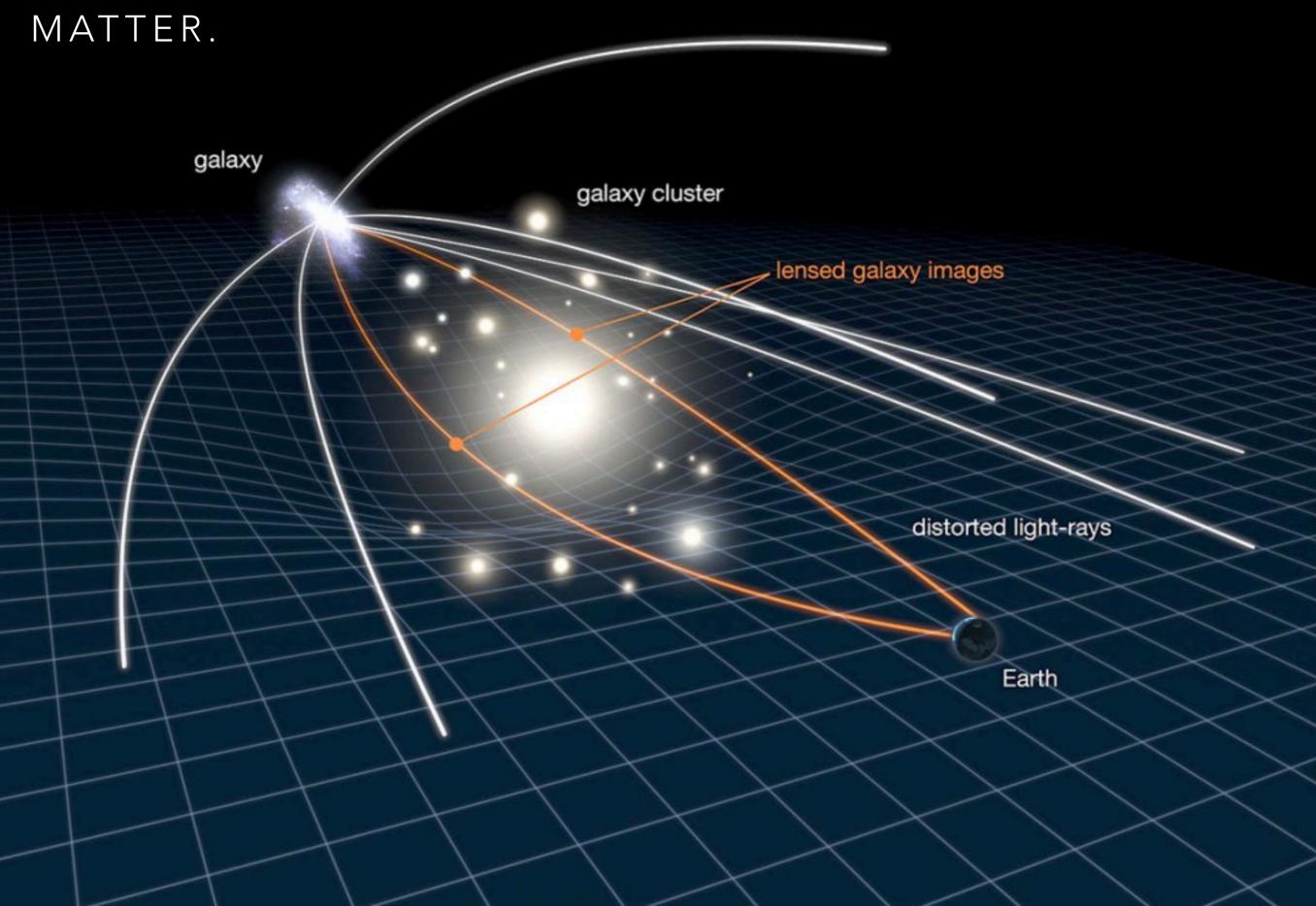
EXAMPLE 2: DARK MATTER



GALAXY CLUSTERS ARE DOMINATED BY DARK MATTER



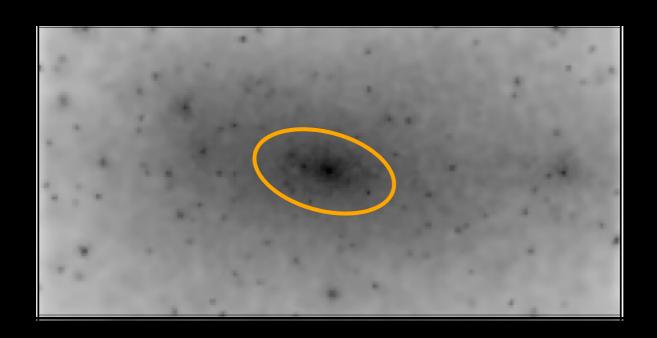
GRAVITATIONAL LENSING HELPS US TRACE THE DARK MATTER



CLASSICAL MODELLING TO INFER PARAMETERS

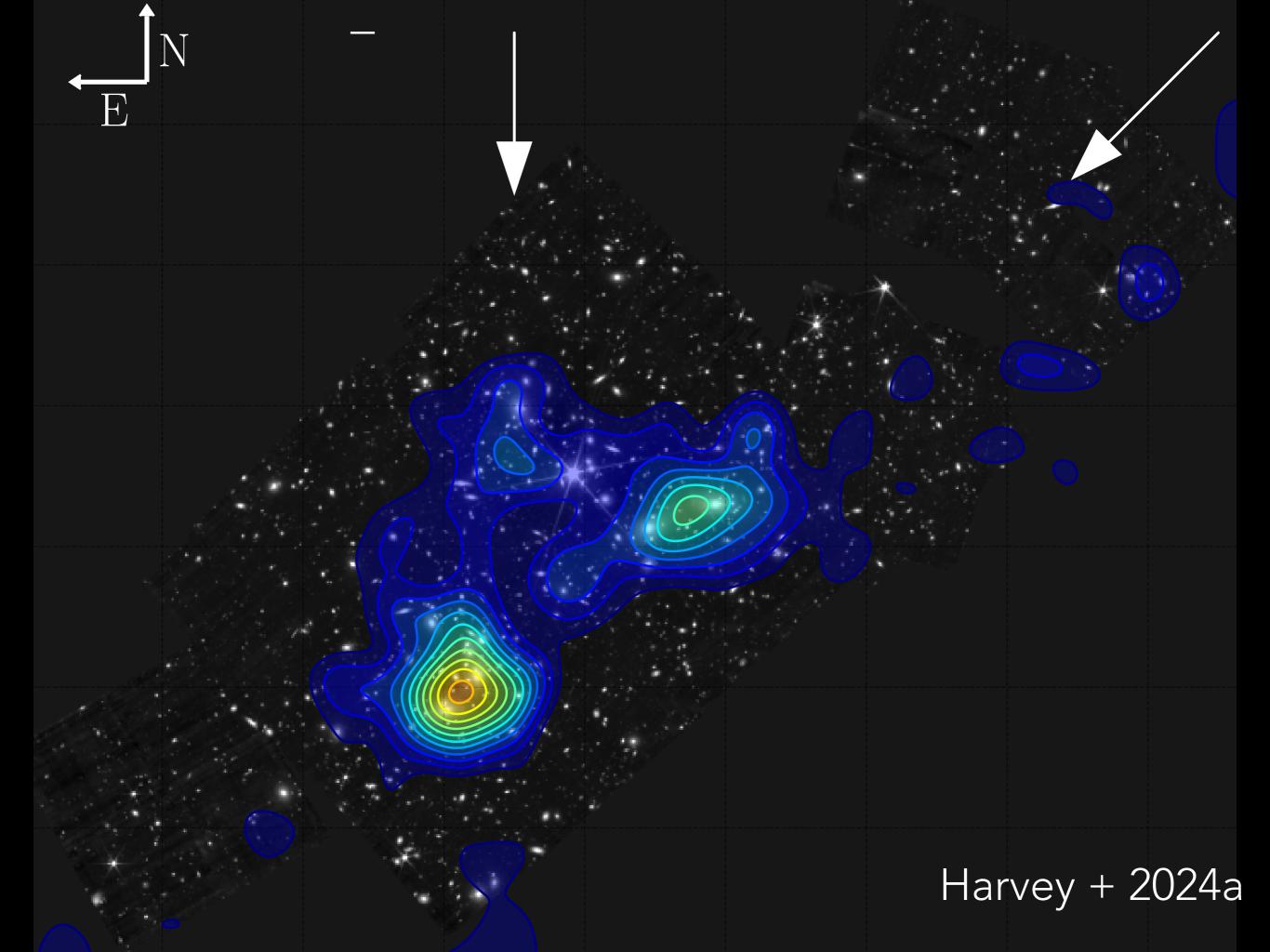


COMPARE THE PARAMETERS TO SIMULATIONS OF GALAXY CLUSTERS WITH DIFFERENT MODELS OF DARK MATTER

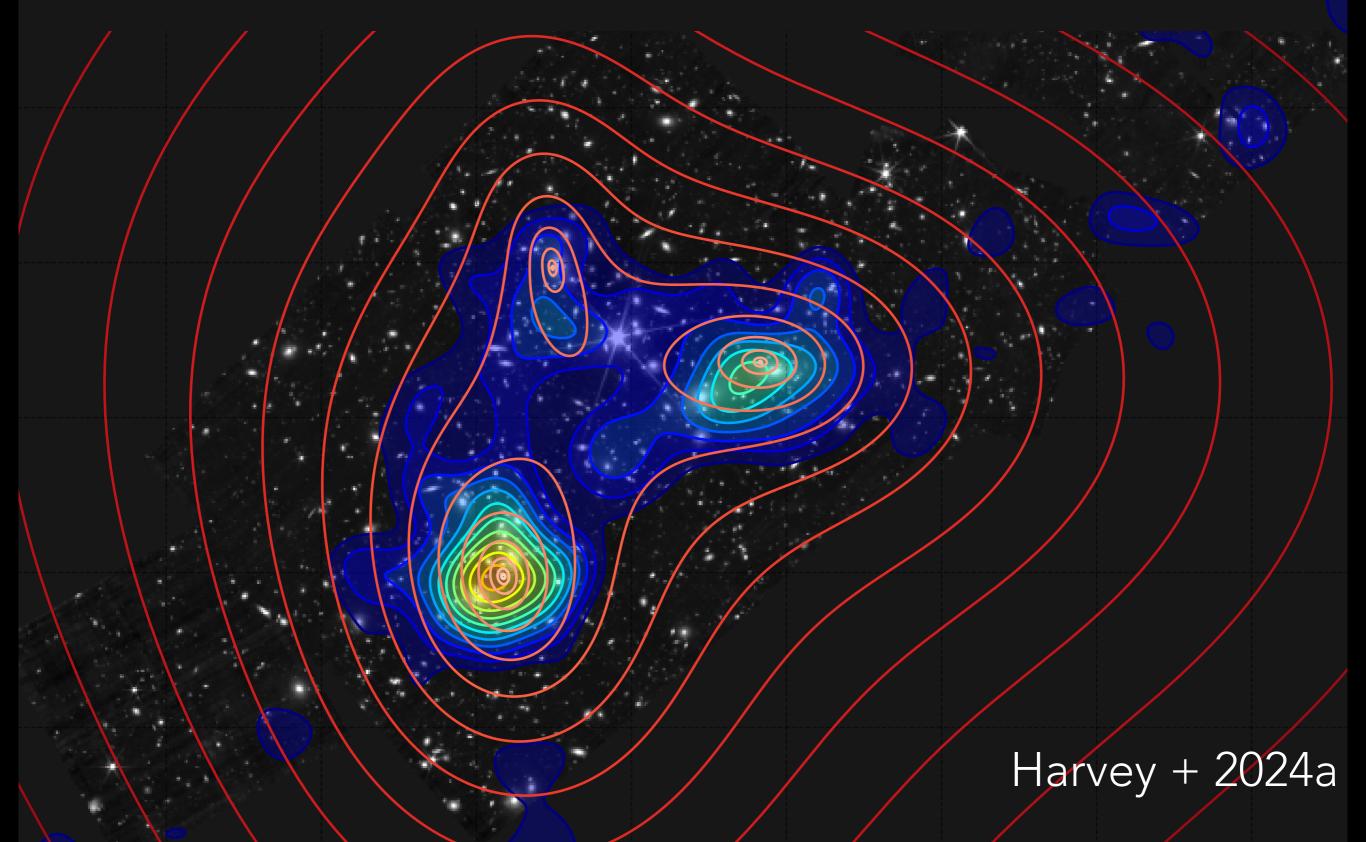


MODELLING IS SLOW

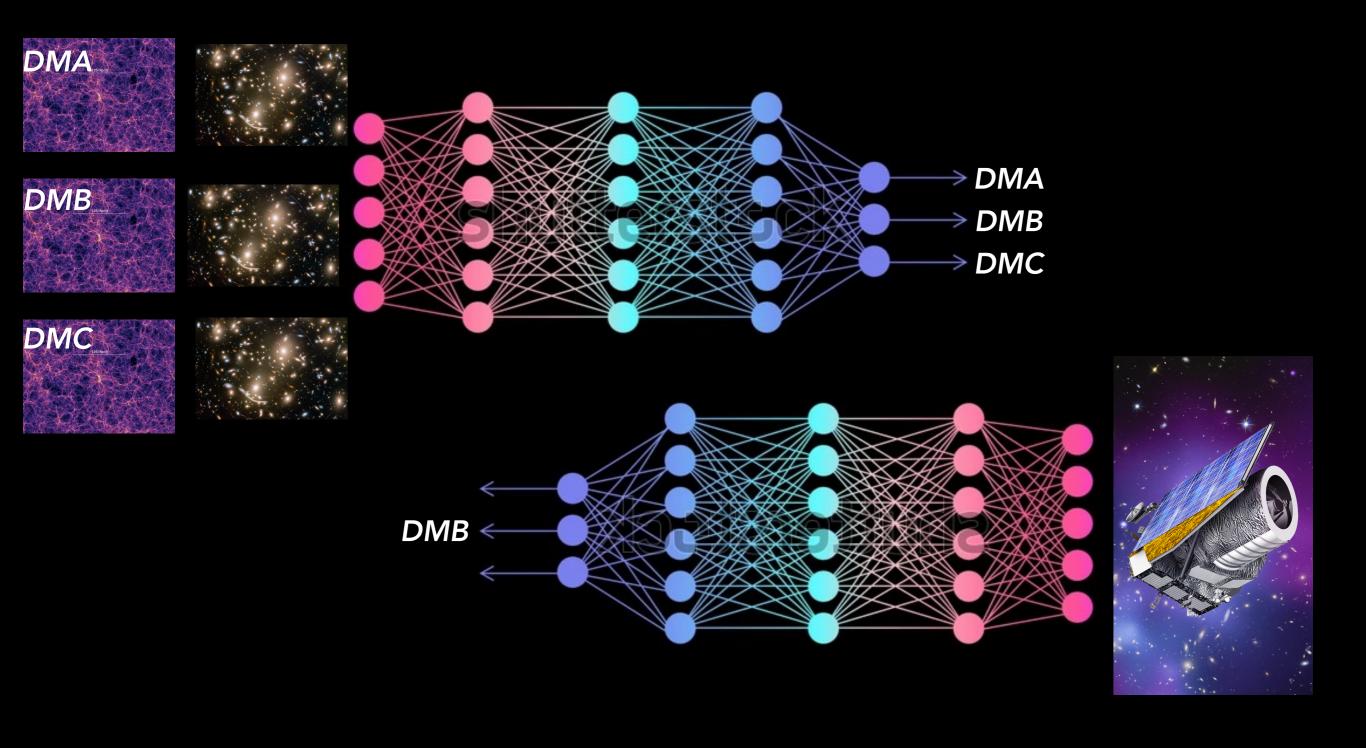




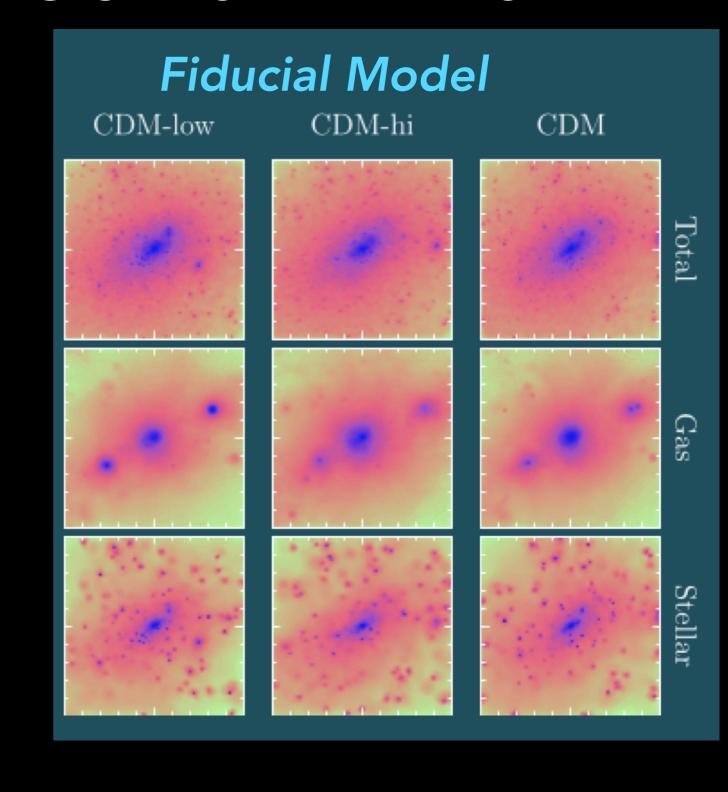
BUT INFORMATION CAN BE LOST IN SIMPLIFICATION.



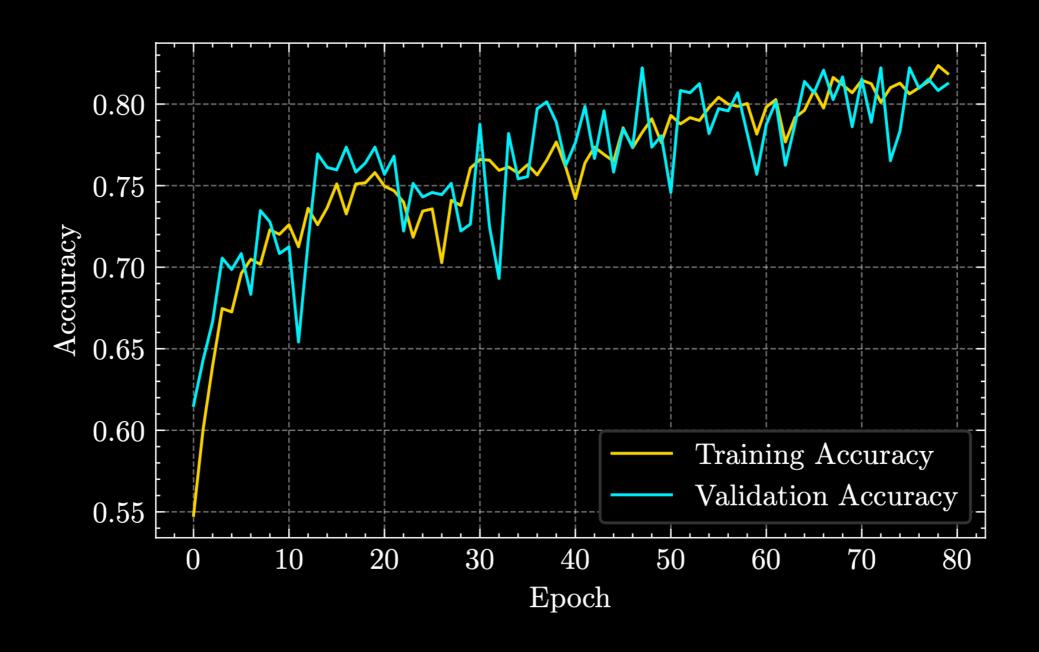
DEEP NETS CAN HELP US AGNOSTICALLY PROBE DARK MATTER WHILST SPEEDING UP THE PROCESS



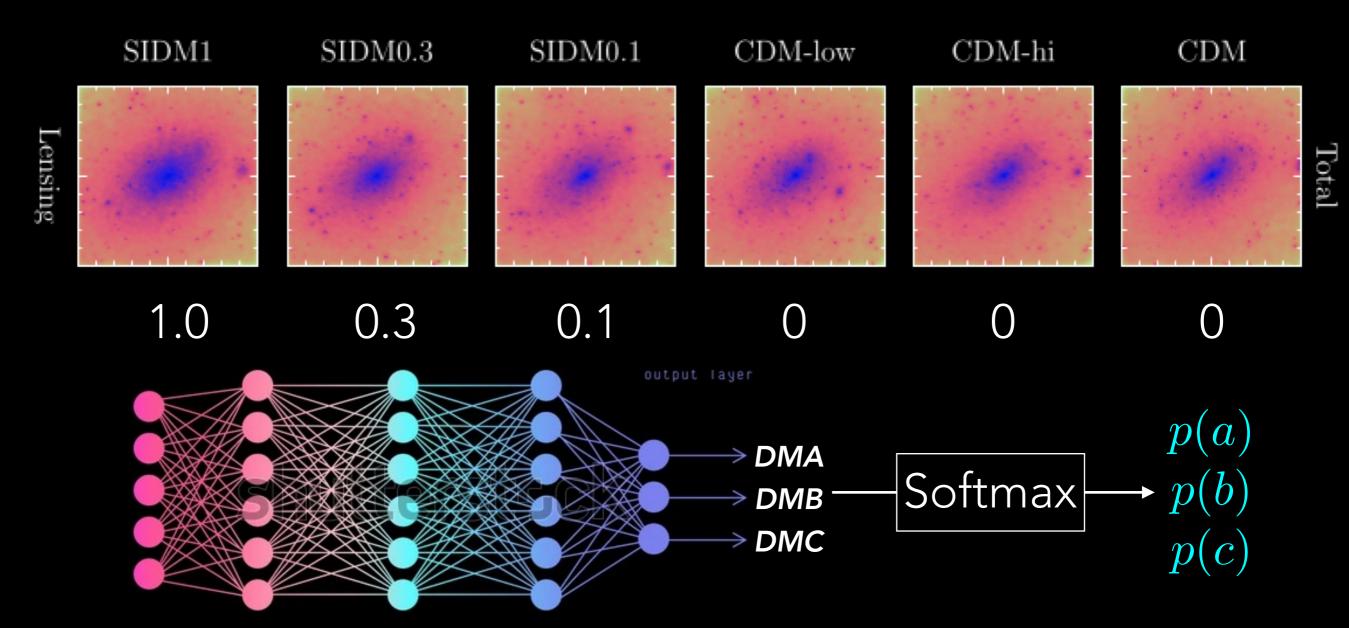
FIRST WE NEED OUR SAMPLES



WE CAN REACH 80% ACCURACY



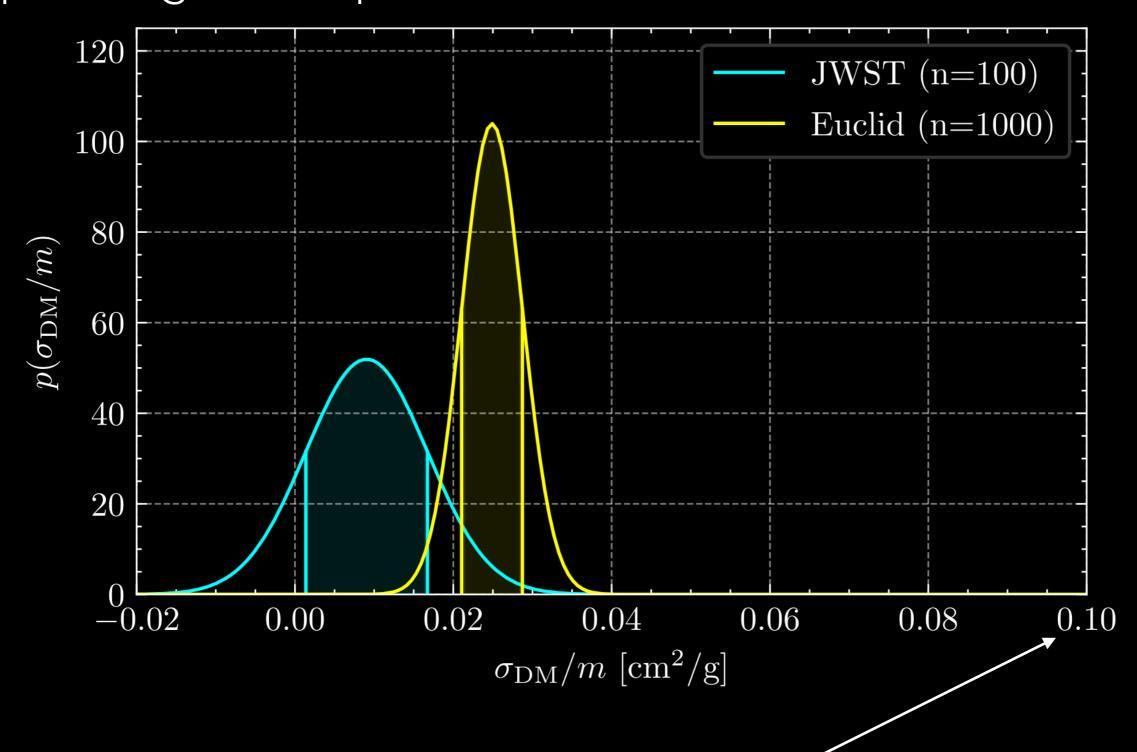
CLASSIFICATION TO REGRESSION



$$\hat{x} = \sum p(x)x$$

Harvey 2024

Predictions after forward modelling on upcoming telescopes



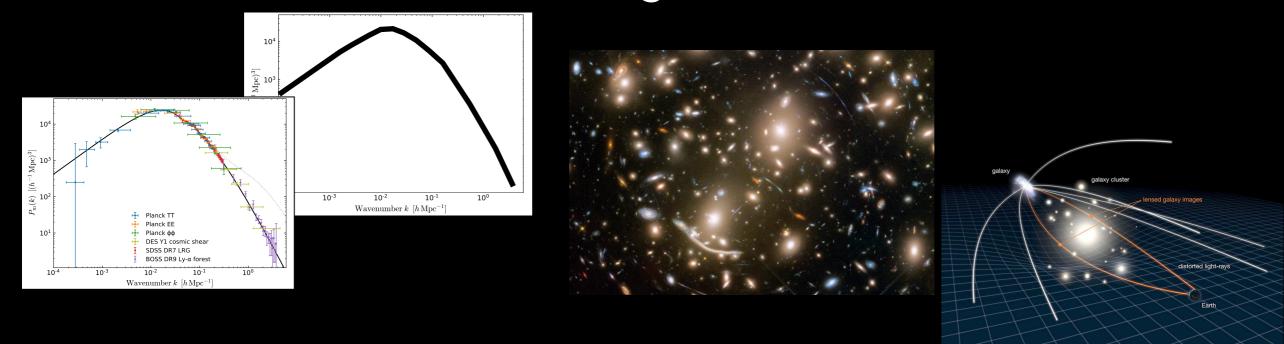
Limit with classic methods

Harvey + 2024b

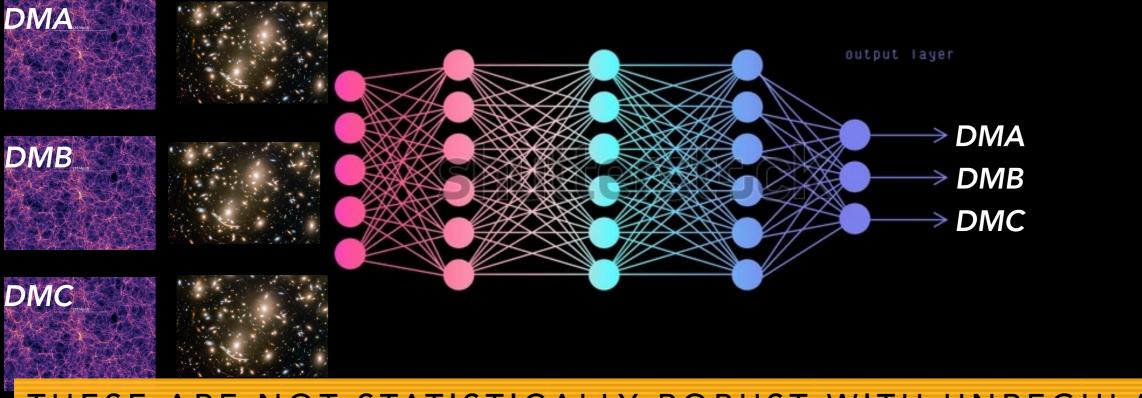
Why astronomy needs machine learning

- Classical inference requires analytical models that can be passed in to a MCMC
- The number of posterior calls is normally very large.
- Likelihoods are not Gaussian.
- Direct comparison of simulations to observations would be impossible in the current situation.
- Astrophysical simulations are often computationally expensive

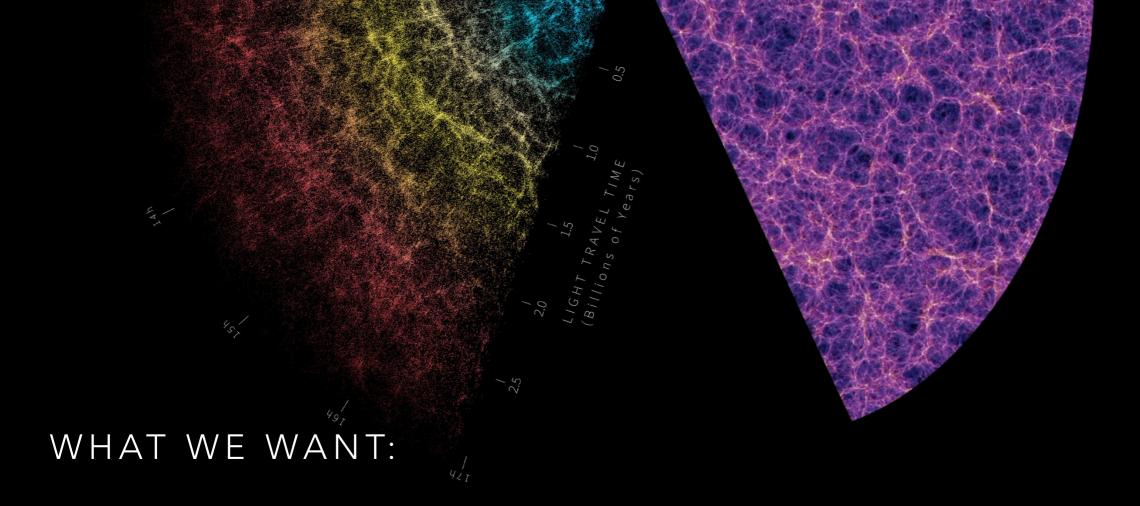
Classical inference is slow and requires analytical models that often lose information resulting in loose and biased inference.



Use machine learning directly on simulations to compare.



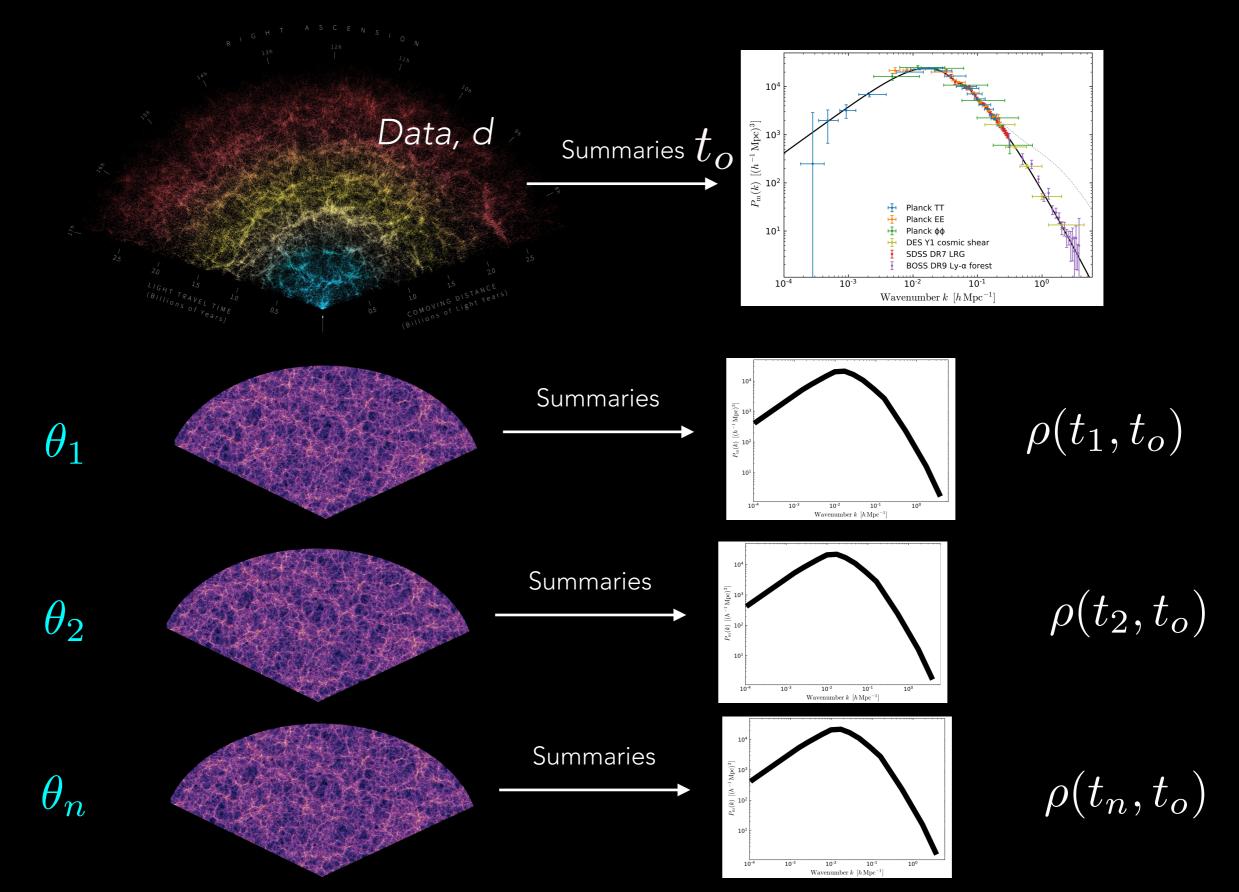
THESE ARE NOT STATISTICALLY ROBUST WITH UNREGULATED
UNCERTAINTY



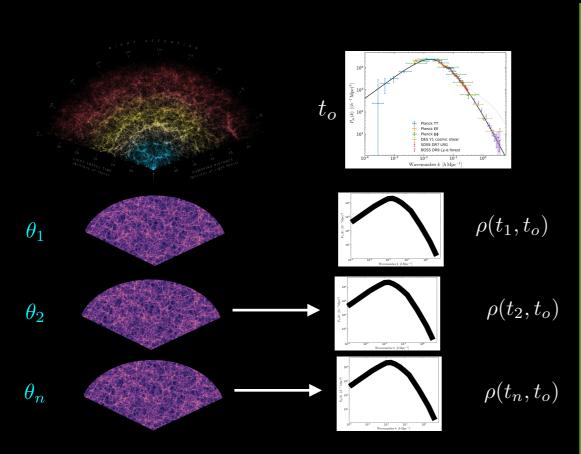
- TO DIRECTLY COMPARE OUR COMPLEX FORWARD MODELS TO THE DATA
- THAT DO NOT ASSUME HOW THE LIKELIHOOD OF THE PARAMETERS OF THE MODEL ARE DISTRIBUTED
- QUICKLY, EFFICIENTLY AND UNBIASED

Simulation Based Inference

BASIC PRINCIPLES OF SBI

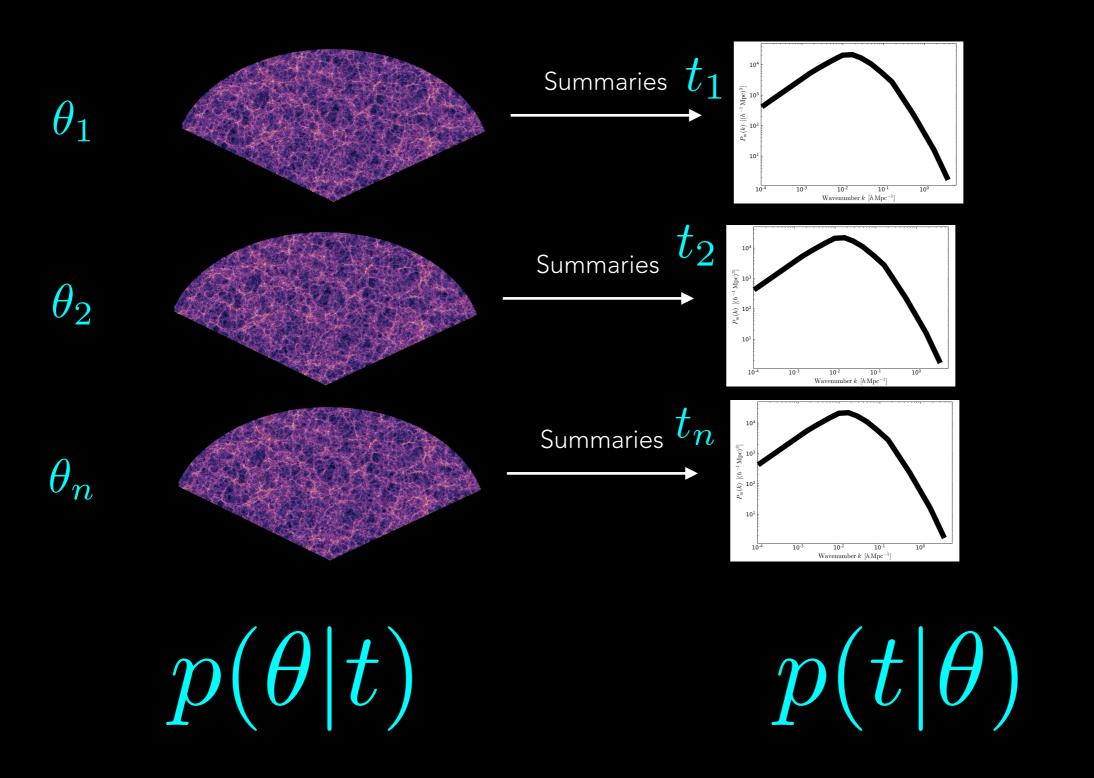


APPROXIMATE BAYESIAN COMPUTATION



- Directly compare complex non-linear models
- Forward modelled to resemble the data
- Unbiased, conservative
- Evidence is calculated
- Take some delta in "agreement" not robust still
- Requires a lot of simulations
- Scales exponentially with number of parameters

What I would rather: Estimate the true posterior or likelihood of the parameters in a model free way.



Then train some density estimator to predict the likelihood / posterior at non-sampled points in the parameter space -

Then evaluate at the location of the observations...

$$p(\theta|t)$$

$$p(\theta|t=t_o)$$

Does not scale with parameters very well

$$p(t|\theta)$$

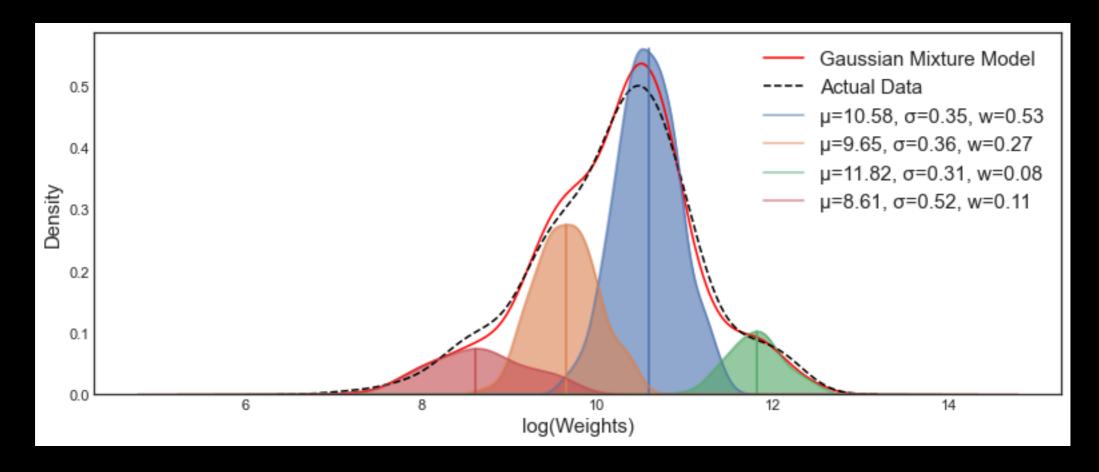
$$p(t_o|\theta)p(\theta)$$

Nice that different priors can be placed

How do we estimate the likelihood for the sampled distribution Neural Density Estimators

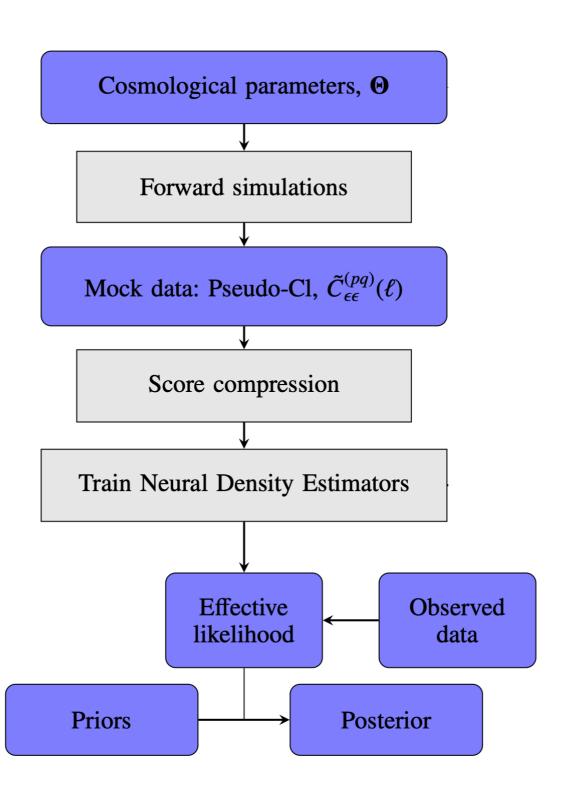
 $p(t|\theta;w)$

Gaussian mixture model

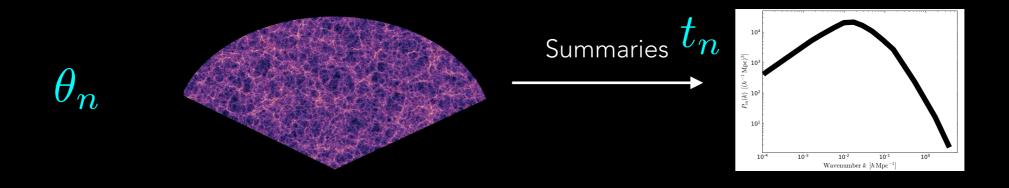


Train using a discrete, monte carlo estimation of the Kullback-Leibler divergence.

HOW THIS WORKS IN PRACTICE



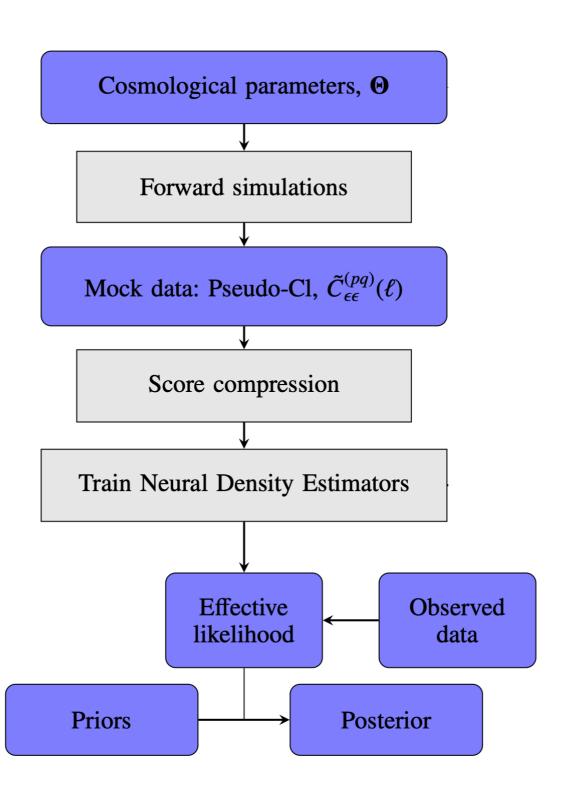
Simulations are very expensive to run.



Why astronomy needs machine learning

- Classical inference requires analytical models that can be passed in to a MCMC
- The number of posterior calls is normally very large.
- Likelihoods are not Gaussian.
- Direct comparison of simulations to observations would be impossible in the current situation.
- Astrophysical simulations are often computationally expensive

HOW THIS WORKS IN PRACTICE



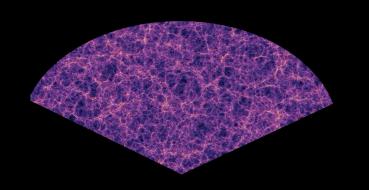
Active learning (with Sequential Neural Likelihood)

First -> more neural density estimators the better

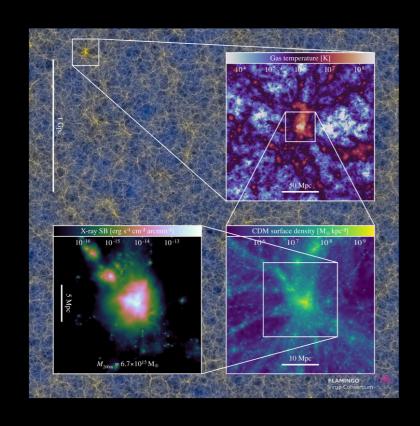
$$p(t|\theta;w) = \sum_{1}^{N_{\mathrm{NDE}}} \beta_{\alpha} p_{\alpha}(t|\theta;w)$$
 Ensembles NDEs

Secondly -> this tells you where the likelihood is important and where there is uncertainty.

SPEED UP SIMULATIONS WITH GENERATIVE METHODS.

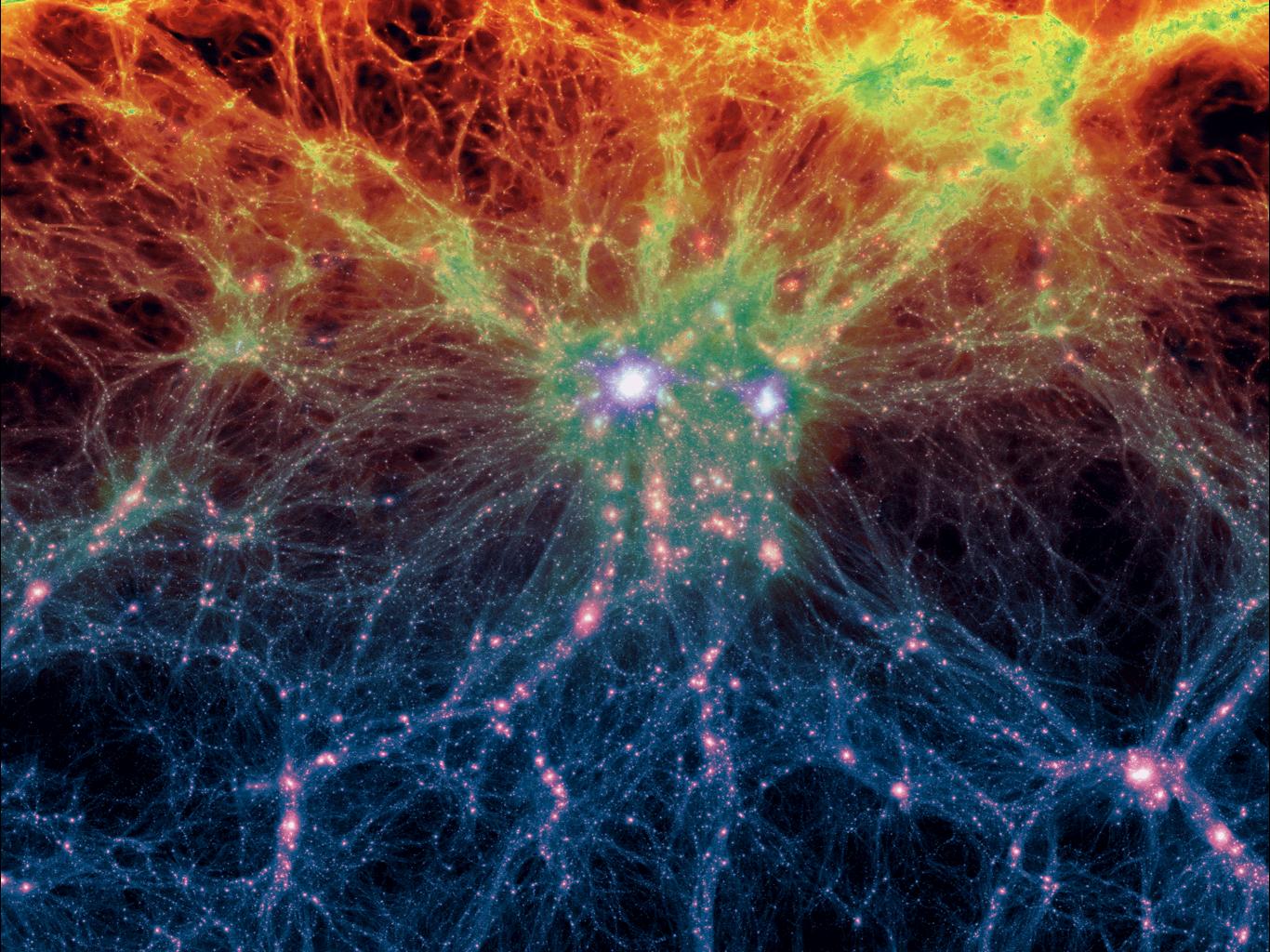


Dark matter only simulations are quick(ish)



Full simulations are slow

TRAIN A METHOD TO LEARN THE CONNECTION BETWEEN DARK MATTER AND ASTROPHYSICS



SPEED UP SIMULATIONS WITH GENERATIVE METHODS: VARIABLE AUTOENCODERS

X-ray emission

Dark matter distribution

Some parameterised model

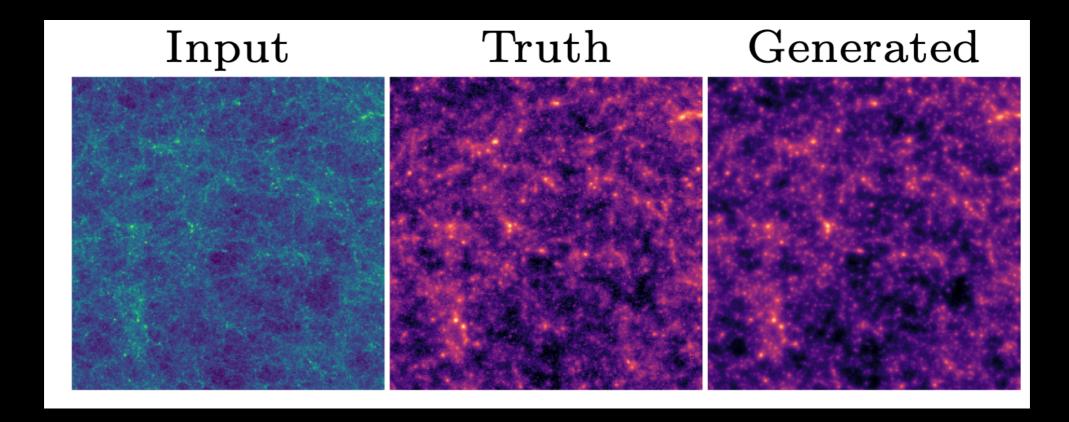
$$p(x|y) = \int dz \ p(x|y,z)p(z|y)$$

Prior information on the weights.

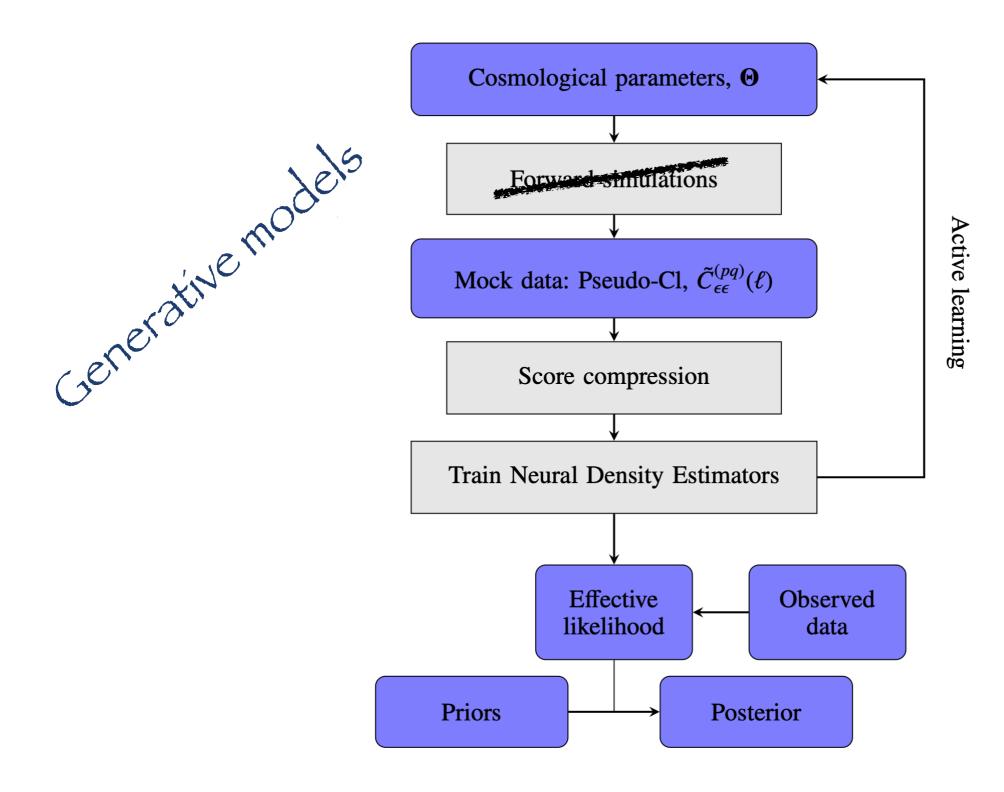
ASSUMING GAUSSIAN MIXTURES

$$p(x|y) = \int dz \ p(x|y,z)p(z|y)$$

Predict the weights of the mixture model using a CNN then sampled from the resulting distributions to create new samples



HOW THIS WORKS IN PRACTICE



OTHER METHODS FOR GENERATIVE MODELS IN ASTRONOMY:

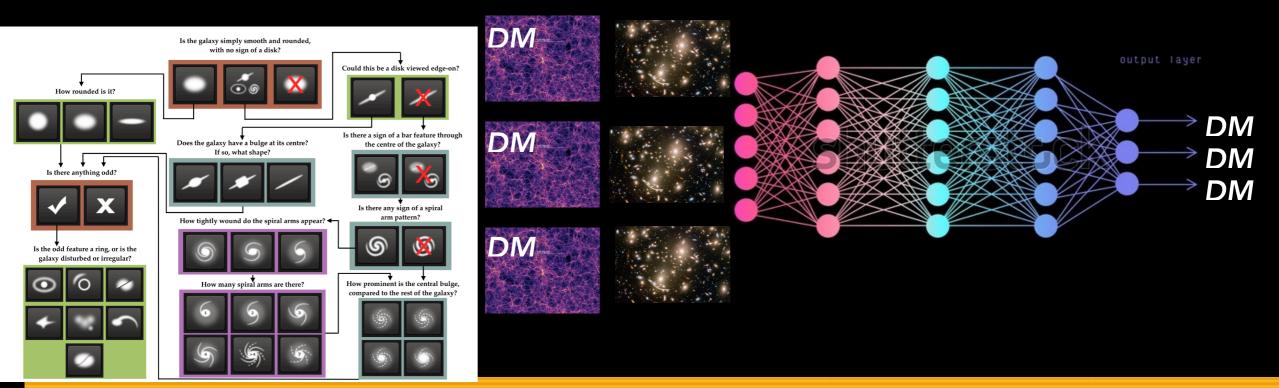
NORMALISING FLOWS (E.G. <u>HTTPS://ARXIV.ORG/PDF/</u> <u>2211.15161</u>, HTTPS://ARXIV.ORG/ABS/2105.12024)

DIFFUSION MODELS (E.G. <u>HTTPS://ARXIV.ORG/PDF/2311.05217</u>)

CLASSICAL INFERENCE IS SLOW AND REQUIRES ANALYTICAL MODELS THAT OFTEN LOSE INFORMATION RESULTING IN LOOSE AND BIASED INFERENCE.

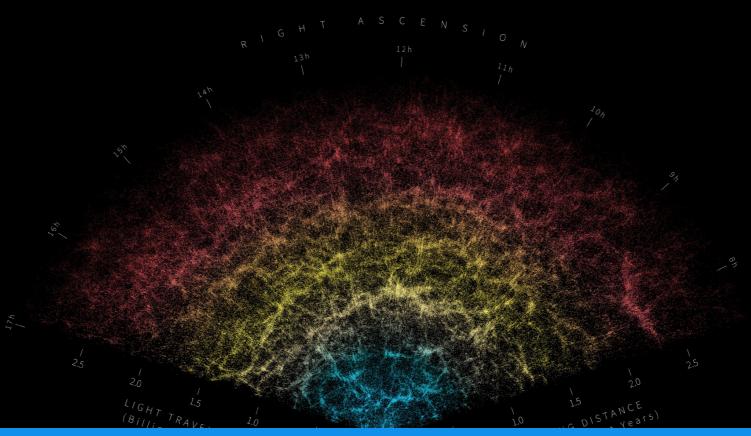


ONE CAN USE MACHINE LEARNING DIRECTLY ON SIMULATIONS TO COMPARE AVOIDING LOSING INFORMATION.

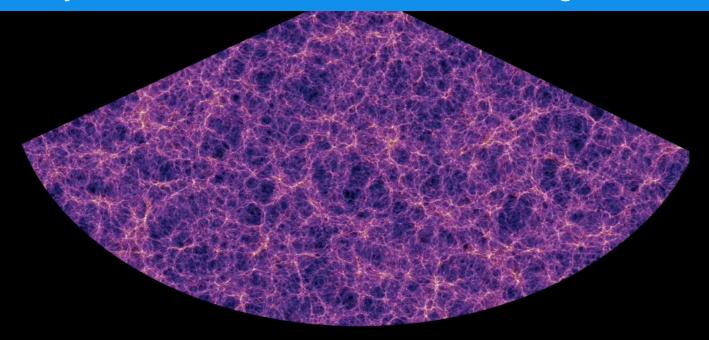


SIMULATION BASED INFERENCE PRESENTS A WAY TO COMPARE DIRECTLY TO SIMULATIONS IN A ROBUST STATISTICAL WAY WITH ADVANCE IN GENERATIVE NETWORKS AIDING INFERENCE

DISTRIBUTION OF NEARBY GALAXIES MAPPED BY THE DARK ENERGY SPECTROSCOPIC INSTRUMENT (DESI)



What is the probability of the our model of the Universe given what we observe?



Distribution of matter as seen by the millennium simulation