



# **Towards reconstructing the halo clustering and halo mass function of N-body simulations using neural ratio estimation**

**Androniki Dimitriou, Christoph Weniger, Camila Correa**

arXiv:2206.11312

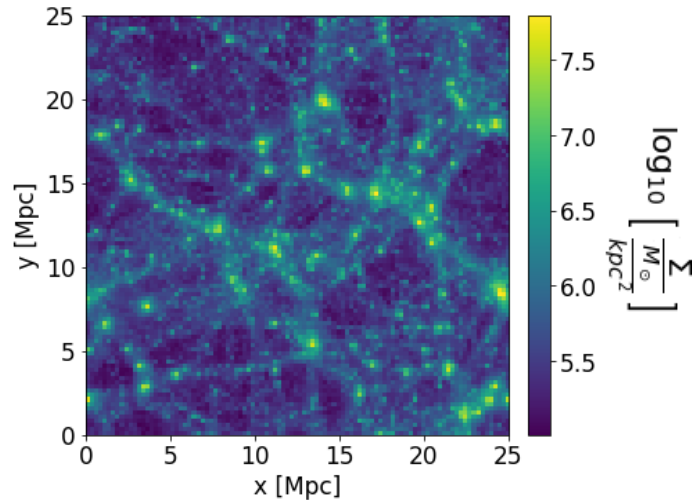
Dark Matters  
1 December 2022, Brussels

- **Goal:** reconstruction of **halo clustering** and **halo mass function** of DM-only cosmological simulations generated by the **EAGLE project**

- **Goal:** reconstruction of **halo clustering** and **halo mass function** of DM-only cosmological simulations generated by the **EAGLE project**
- **Approach:** direct calibration of **a analytical halo** model/simulator **based on a toy implementation of two body correlation functions** on DM-only **simulations**

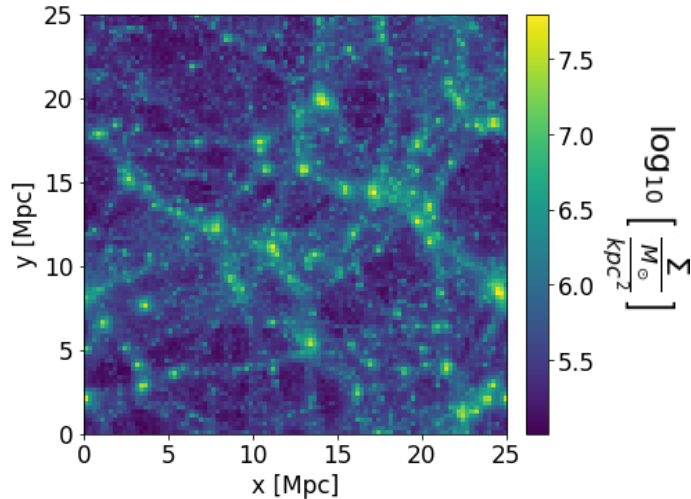
- **Goal:** reconstruction of **halo clustering** and **halo mass function** of DM-only cosmological simulations generated by the **EAGLE project**
- **Approach:** direct calibration of **a analytical halo** model/simulator **based on a toy implementation of two body correlation functions** on DM-only **simulations**
- **The EAGLE project**

(25 Mpc)<sup>3</sup> box with 376<sup>3</sup> particles



- **Goal:** reconstruction of **halo clustering** and **halo mass function** of DM-only cosmological simulations generated by the **EAGLE project**
- **Approach:** direct calibration of **a analytical halo model/simulator based on a toy implementation of two body correlation functions** on DM-only **simulations**
- **The EAGLE project**

(25 Mpc)<sup>3</sup> box with 376<sup>3</sup> particles



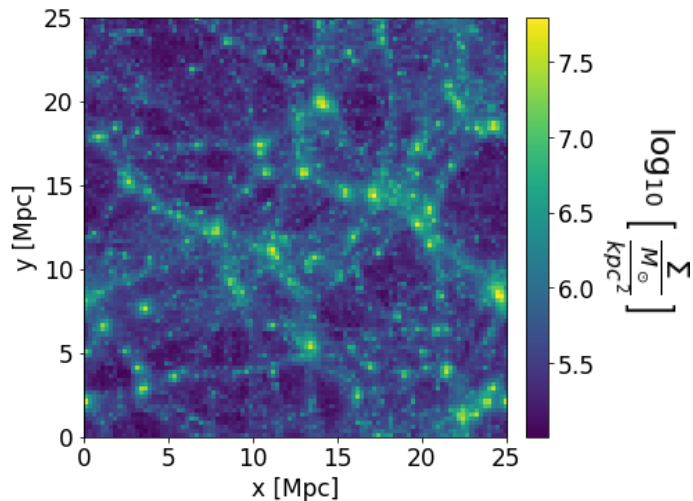
FOF halo finder



Particle data

- **Goal:** reconstruction of **halo clustering** and **halo mass function** of DM-only cosmological simulations generated by the **EAGLE project**
- **Approach:** direct calibration of **a analytical halo model/simulator based on a toy implementation of two body correlation functions** on DM-only **simulations**
- **The EAGLE project**

(25 Mpc)<sup>3</sup> box with 376<sup>3</sup> particles

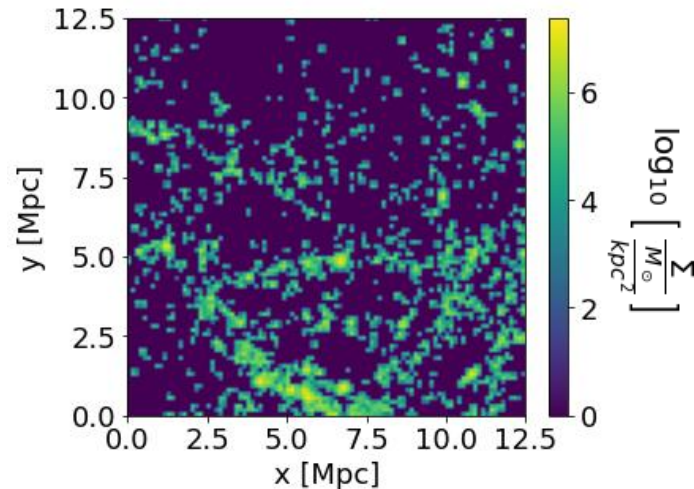


FOF halo finder

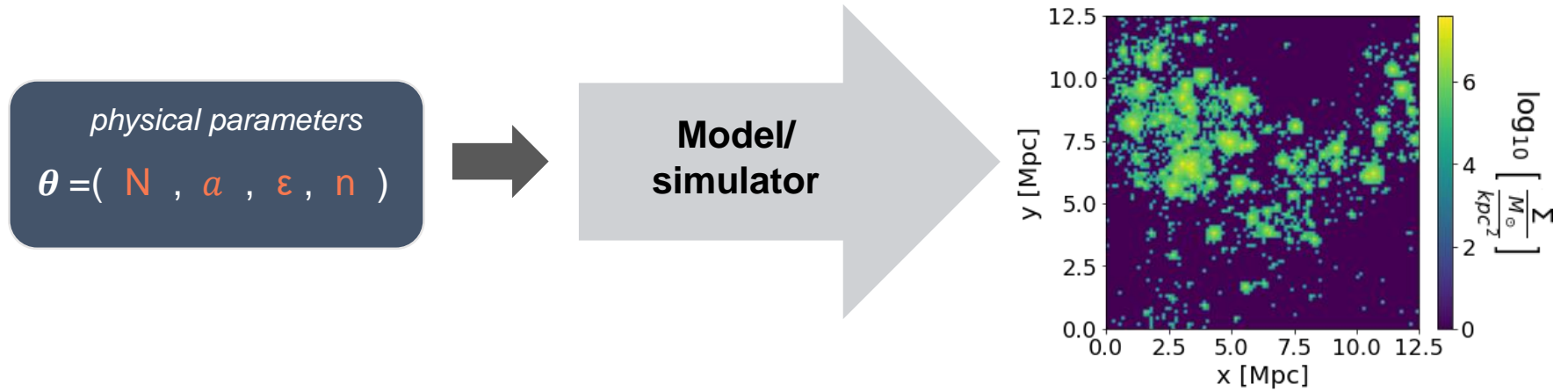


Particle data

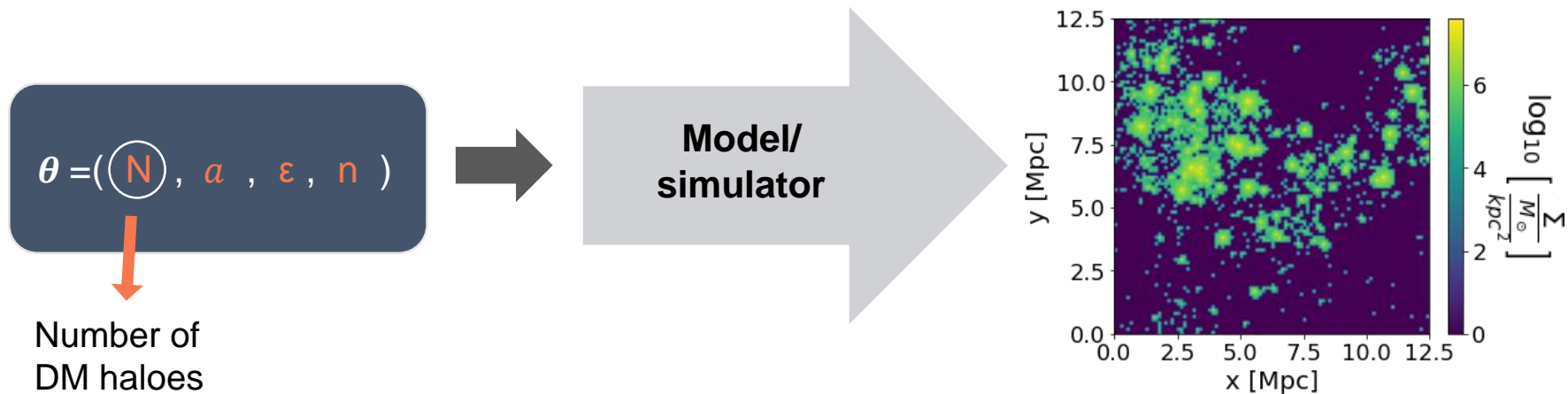
$M_h \in (10^9, 10^{12}) M_\odot$ , (12.5 Mpc)<sup>3</sup> box



- The toy halo model



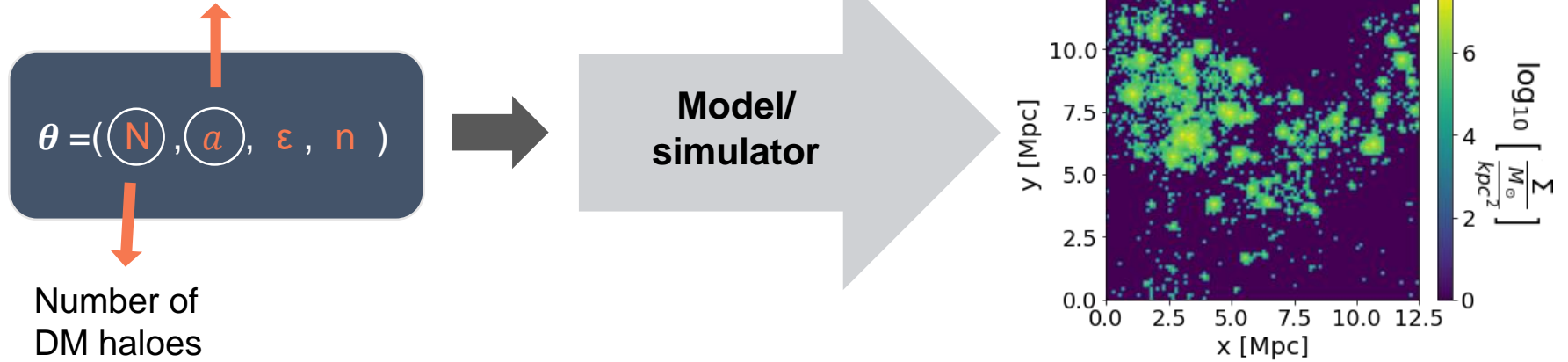
- The toy halo model





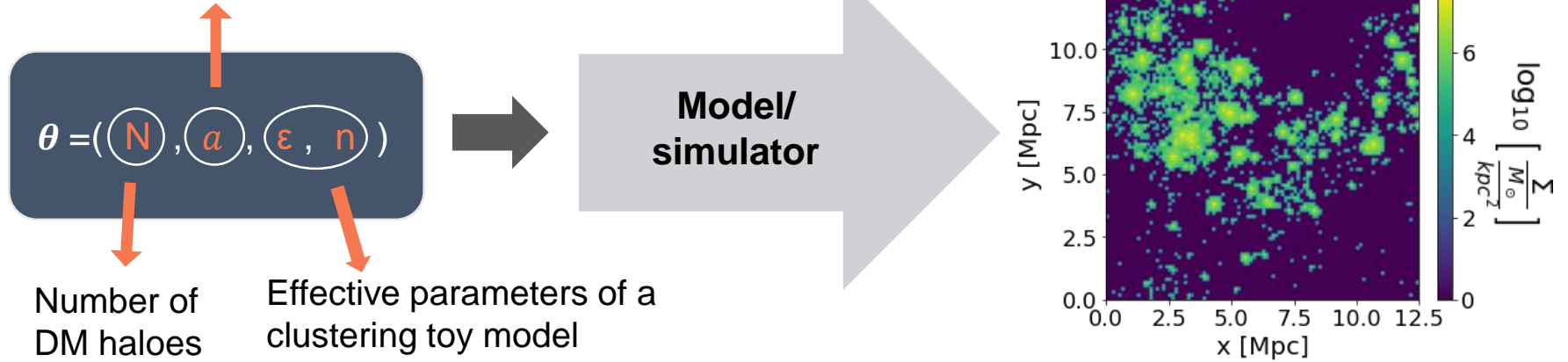
- **The toy halo model**

Slope of HMF:  $\frac{dn}{dM} \propto M^{-a}$ , where  $a \cong 1.9$  for low masses

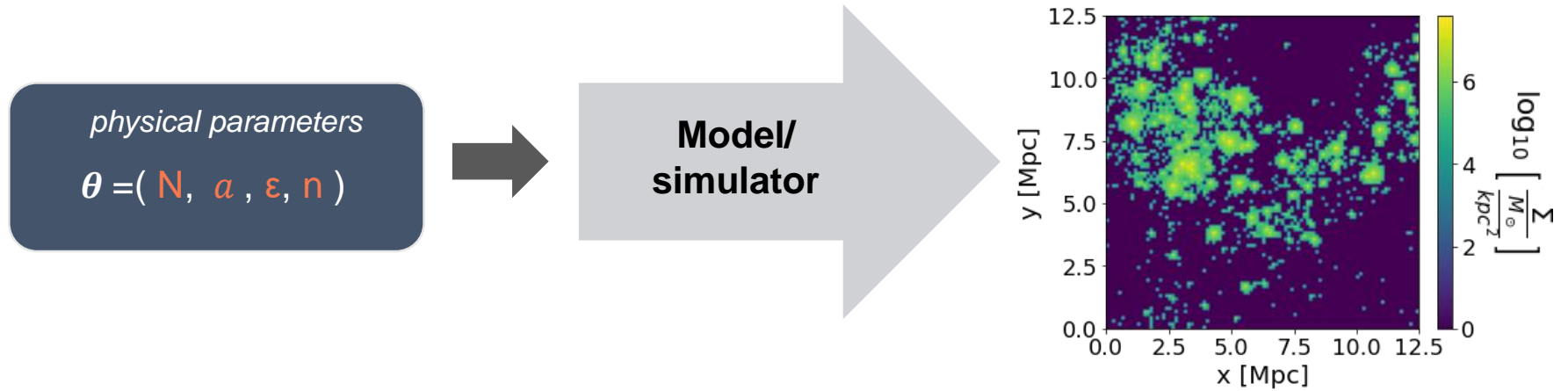


- **The toy halo model**

Slope of HMF:  $\frac{dn}{dM} \propto M^{-a}$ , where  $a \cong 1.9$  for low masses

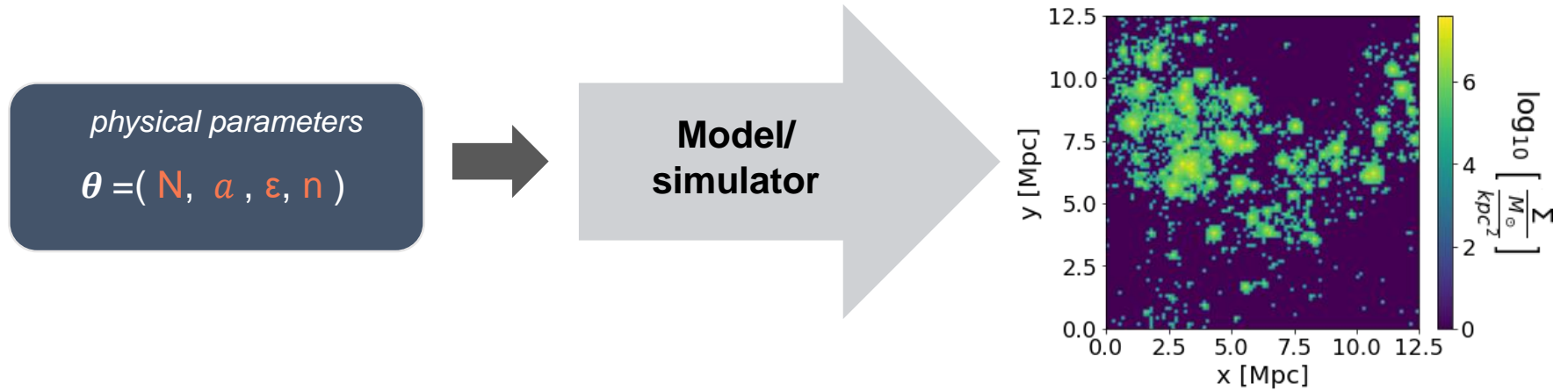


- The toy halo model



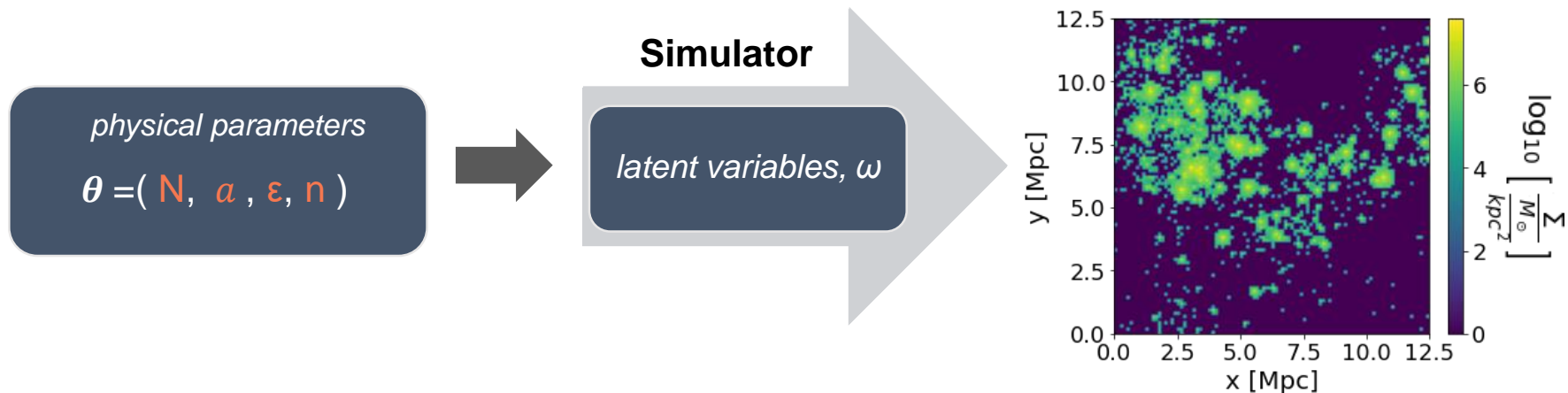
- Given an image, we want to know which parameters,  $\theta$ , generated it

- **The toy halo model**



- Given an image, we want to know which parameters,  $\theta$ , generated it
- **What is inside the simulator?**

- **The toy halo model**

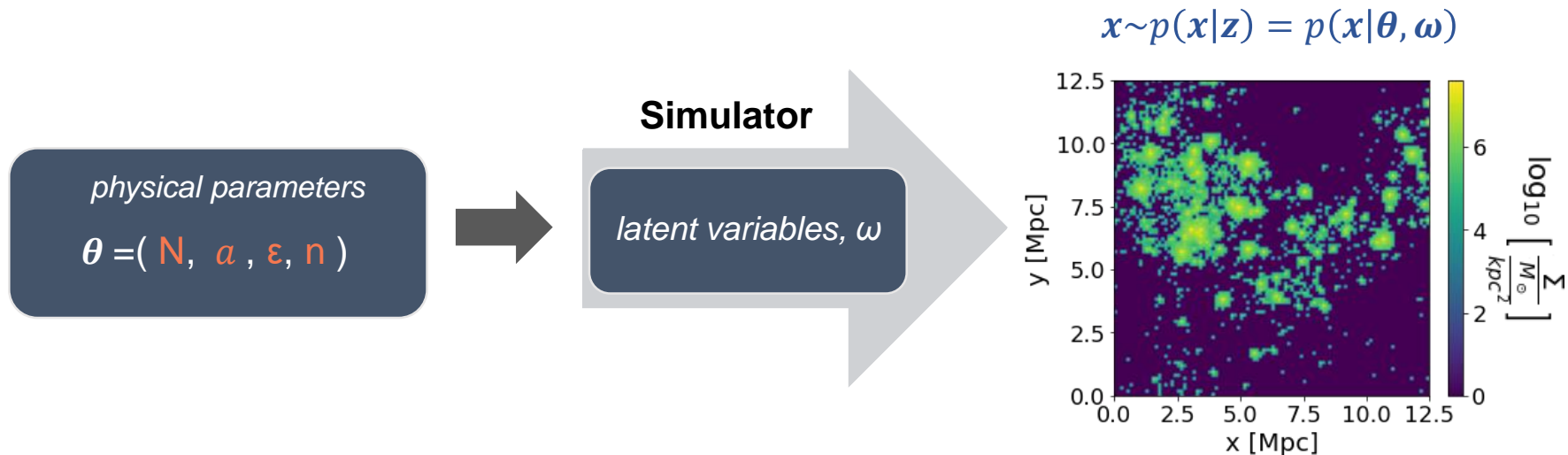


- Given an image, we want to know which parameters,  $\theta$ , generated it

- **What is inside the simulator?**

latent variables

- **The toy halo model**



- Given an image, we want to know which parameters,  $\theta$ , generated it
- **What is inside the simulator?**

latent variables

- Given some simulation  $x$ , and a model with parameters  $z$ , what can we learn about  $z$ ?

- Given some simulation  $\mathbf{x}$ , and a model with parameters  $\mathbf{z}$ , what can we learn about  $\mathbf{z}$ ?

## Bayes Theorem

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$$

- Posterior:  $p(\mathbf{z}|\mathbf{x})$
- Likelihood:  $p(\mathbf{x}|\mathbf{z})$
- Prior:  $p(\mathbf{z})$
- Evidence:  $p(\mathbf{x}) = \int d\mathbf{z} p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$



- Given some simulation  $\mathbf{x}$ , and a model with parameters  $\mathbf{z}$ , what can we learn about  $\mathbf{z}$ ?

## Bayes Theorem

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$$

- Posterior:  $p(\mathbf{z}|\mathbf{x})$
  - Likelihood:  $p(\mathbf{x}|\mathbf{z})$
  - Prior:  $p(\mathbf{z})$
  - Evidence:  $p(\mathbf{x}) = \int d\mathbf{z} p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$
- Bayes theorem, gives the **joint posterior** for the full parameter space


- Given some simulation  $\mathbf{x}$ , and a model with parameters  $\mathbf{z}$ , what can we learn about  $\mathbf{z}$ ?

## Bayes Theorem

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$$

- Posterior:  $p(\mathbf{z}|\mathbf{x})$
  - Likelihood:  $p(\mathbf{x}|\mathbf{z})$
  - Prior:  $p(\mathbf{z})$
  - Evidence:  $p(\mathbf{x}) = \int d\mathbf{z} p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$
- Bayes theorem, gives the **joint posterior** for the full parameter space
  - Almost always, we want to calculate **marginal posteriors** of the parameters of interest,  $\boldsymbol{\vartheta}$

$$p(\boldsymbol{\vartheta}|\mathbf{x}) = \int d\boldsymbol{\eta} p(\boldsymbol{\eta}, \boldsymbol{\vartheta}|\mathbf{x}) = \frac{p(\mathbf{x}|\boldsymbol{\vartheta})p(\boldsymbol{\vartheta})}{p(\mathbf{x})}$$

nuisance parameters 

- Given some simulation  $\mathbf{x}$ , and a model with parameters  $\mathbf{z}$ , what can we learn about  $\mathbf{z}$ ?


## Bayes Theorem

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$$

- Posterior:  $p(\mathbf{z}|\mathbf{x})$
- Likelihood:  $p(\mathbf{x}|\mathbf{z})$
- Prior:  $p(\mathbf{z})$
- Evidence:  $p(\mathbf{x}) = \int d\mathbf{z} p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$

- Bayes theorem, gives the **joint posterior** for the full parameter space
- Almost always, we want to calculate **marginal posteriors** of the parameters of interest,  $\boldsymbol{\vartheta}$

$$p(\boldsymbol{\vartheta}|\mathbf{x}) = \int d\boldsymbol{\eta} p(\boldsymbol{\eta}, \boldsymbol{\vartheta}|\mathbf{x}) = \frac{p(\mathbf{x}|\boldsymbol{\vartheta})p(\boldsymbol{\vartheta})}{p(\mathbf{x})}$$

nuisance parameters 

### Intractable due to high dimensionality

This problem appears in the context of likelihood-based inference methods, e.g., MCMC

# Marginal Neural Ratio Estimation (MNRE)

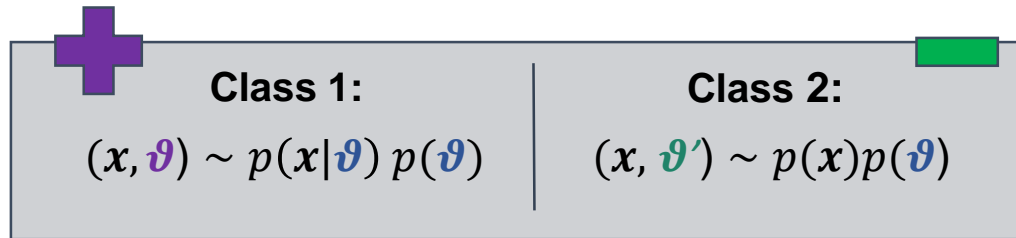
arXiv: 2107.01214

# Marginal Neural Ratio Estimation (MNRE)

arXiv: 2107.01214

**Estimates posteriors through a binary classification problem:**

“Given a (parameter:  $\vartheta$ , image:  $x$ ) pair, is the image,  $x$ , actually generated by the parameter  $\vartheta$ ?”

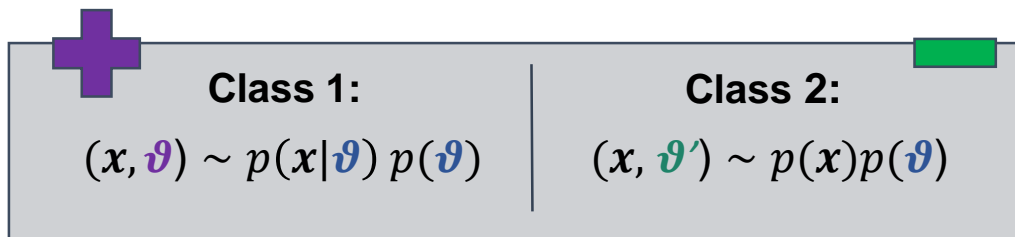


# Marginal Neural Ratio Estimation (MNRE)

arXiv: 2107.01214

**Estimates posteriors through a binary classification problem:**

“Given a (parameter:  $\vartheta$ , image:  $x$ ) pair, is the image,  $x$ , actually generated by the parameter  $\vartheta$ ?”



Train a classifier with mock data to directly estimate:

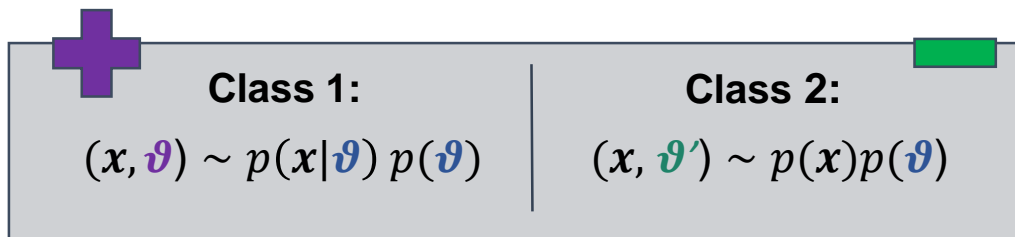
$$r(x, \vartheta) \cong \frac{p(x|\vartheta)}{p(x)}$$

# Marginal Neural Ratio Estimation (MNRE)

arXiv: 2107.01214

**Estimates posteriors through a binary classification problem:**

“Given a (parameter:  $\vartheta$ , image:  $x$ ) pair, is the image,  $x$ , actually generated by the parameter  $\vartheta$ ?”



Train a classifier with mock data to directly estimate:

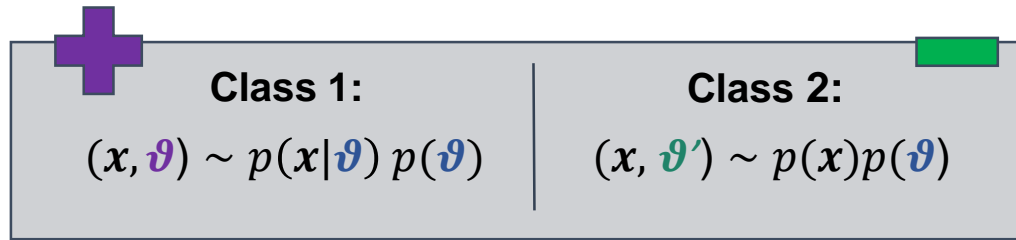
$$r(x, \vartheta) \cong \frac{p(x|\vartheta)}{p(x)} \stackrel{\text{bayes theorem}}{=} \frac{p(\vartheta|x)}{p(\vartheta)}$$

# Marginal Neural Ratio Estimation (MNRE)

arXiv: 2107.01214

**Estimates posteriors through a binary classification problem:**

“Given a (parameter:  $\vartheta$ , image:  $x$ ) pair, is the image,  $x$ , actually generated by the parameter  $\vartheta$ ?”



Train a classifier with mock data to directly estimate:

$$r(x, \vartheta) \cong \frac{p(x|\vartheta)}{p(x)} = \frac{p(\vartheta|x)}{p(\vartheta)}$$

Once we have trained the network, we can **estimate the posterior**:

$$p(\vartheta|\vec{x}) = r(x, \vartheta)p(\vartheta)$$



# Training with swyft

arXiv:2107.01214

# Training with swyft

arXiv:2107.01214

- **Physical parameters:**

- **N**: Number of halos, where  $N \in (100, 2100)$
- **$\epsilon$** : Exponent of the density field, where  $\epsilon \in (0, 2)$
- **a**: Inner slope of the halo mass function, where  $a \in (1, 3)$
- **n**: Slope of the power spectrum, where  $n \in (0, 10)$

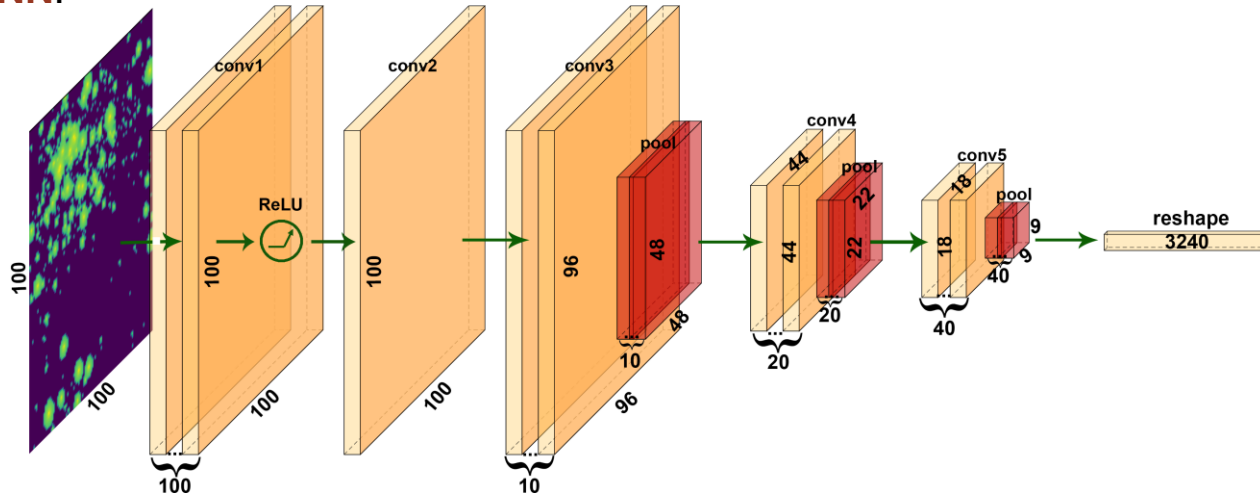
# Training with swyft

arXiv:2107.01214

- **Physical parameters:**

- **N**: Number of halos, where  $N \in (100, 2100)$
- **$\epsilon$** : Exponent of the density field, where  $\epsilon \in (0,2)$
- **a**: Inner slope of the halo mass function, where  $a \in (1, 3)$
- **n**: Slope of the power spectrum, where  $n \in (0,10)$

- We define a **CNN**:



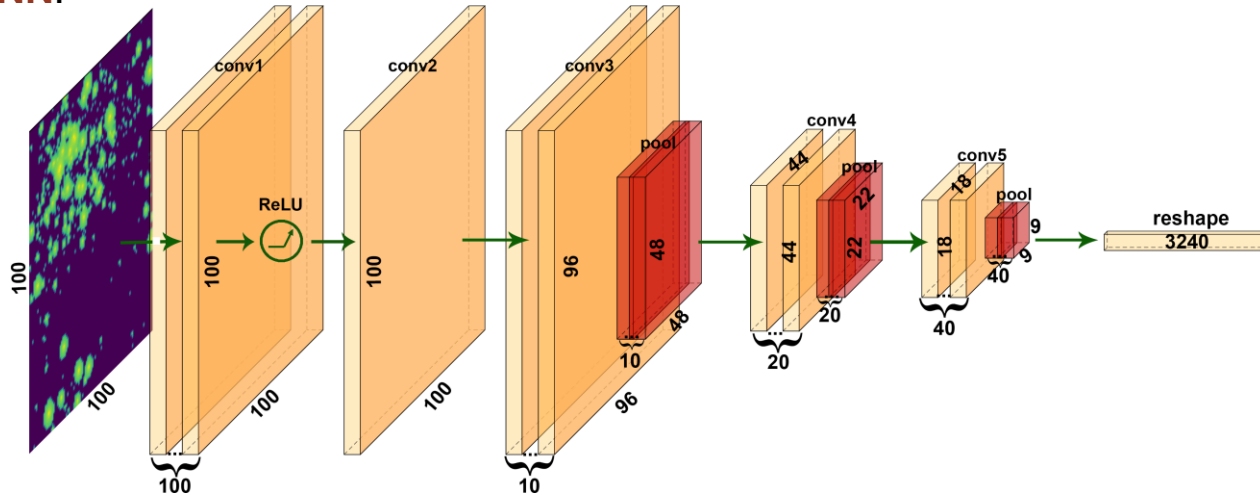
# Training with swyft

arXiv:2107.01214

- **Physical parameters:**

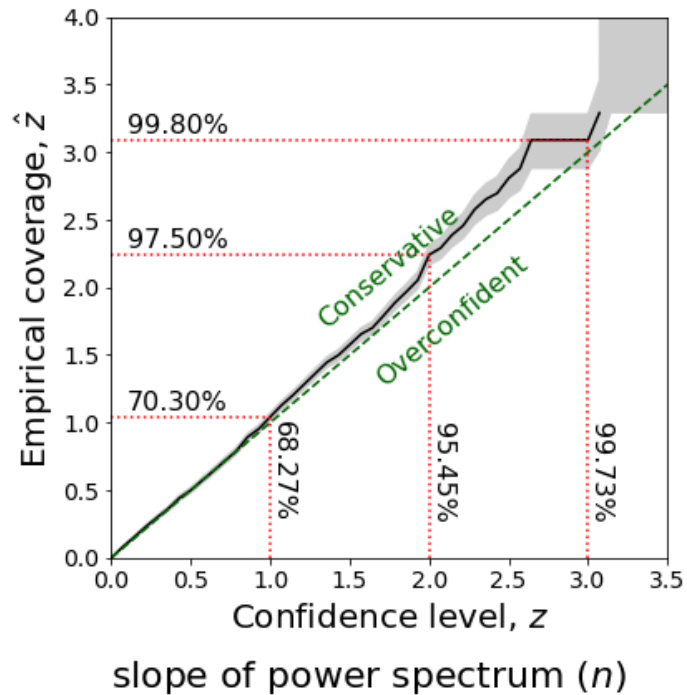
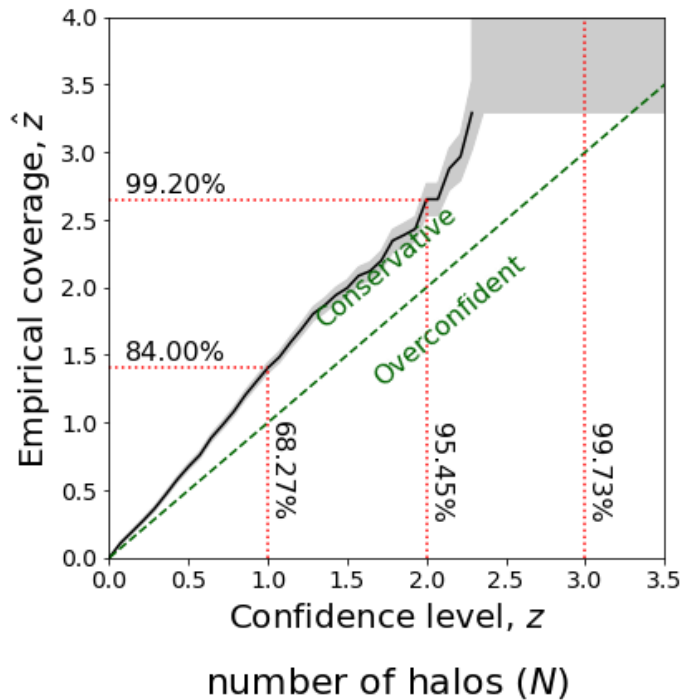
- **N**: Number of halos, where  $N \in (100, 2100)$
- $\epsilon$ : Exponent of the density field, where  $\epsilon \in (0,2)$
- **a**: Inner slope of the halo mass function, where  $a \in (1, 3)$
- **n**: Slope of the power spectrum, where  $n \in (0,10)$

- We define a **CNN**:



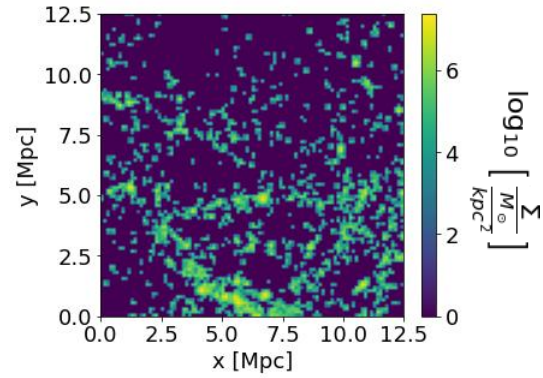
- We **train** using 200.000 mock images

# Results on mock data



# Results on actual N body simulations

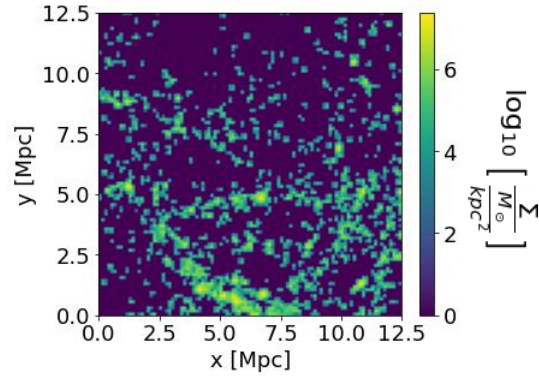
- one simulation box



trained NN

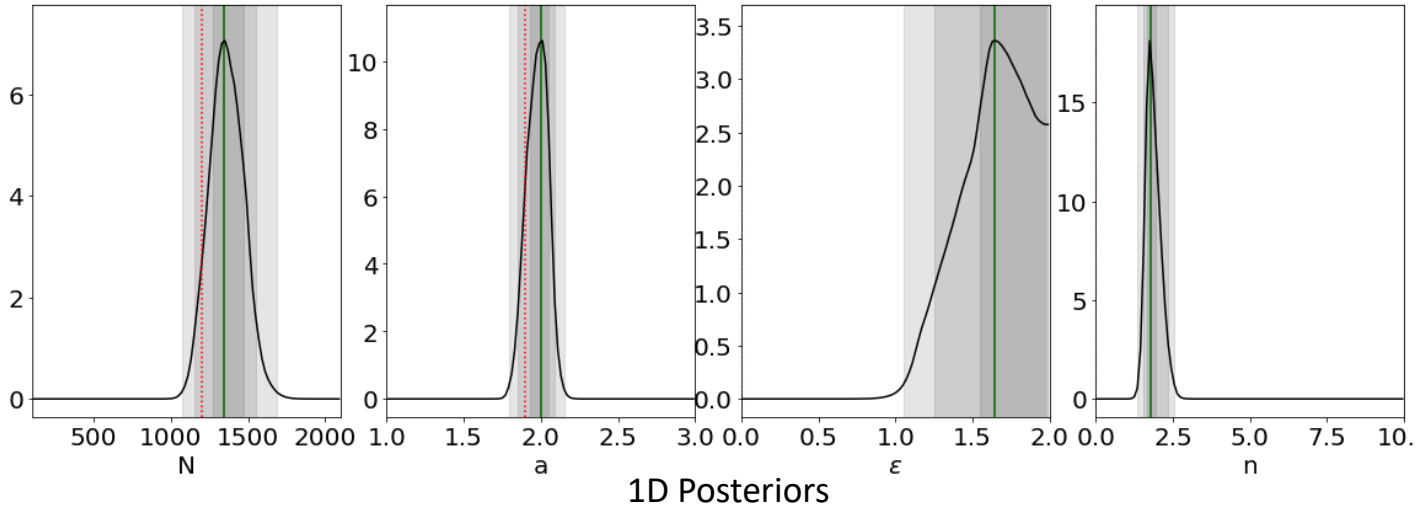
# Results on actual N body simulations

- one simulation box



trained NN

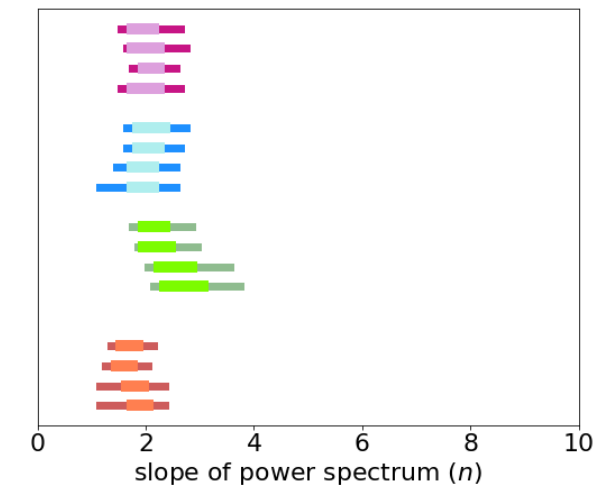
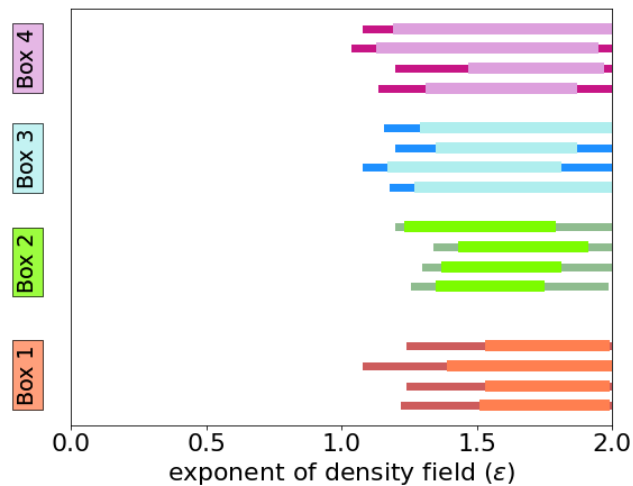
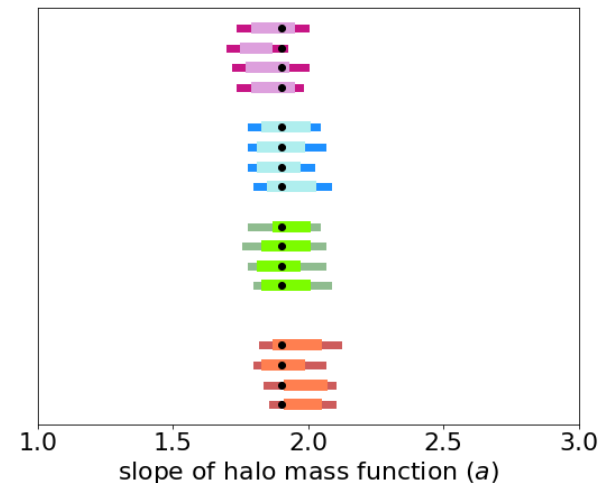
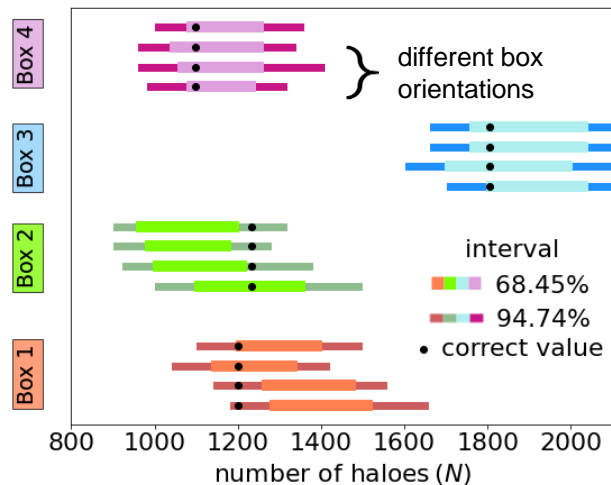
— mode      - - - correct value



# Results on actual N body simulations

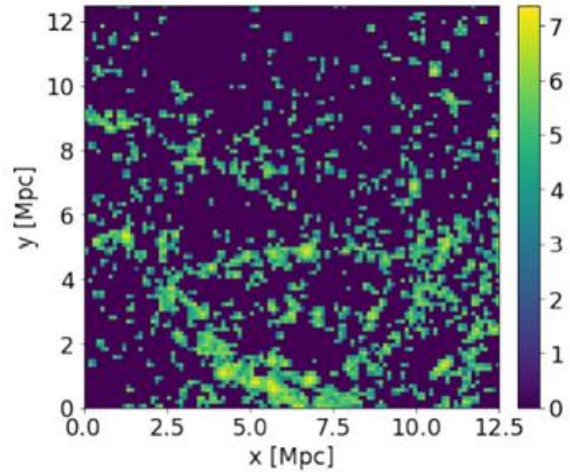
- 4 simulation boxes and 3 rotations of them

Accurate marginal posteriors!

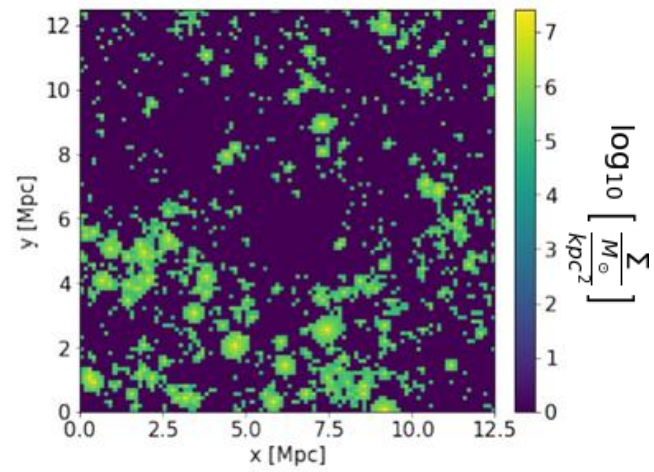




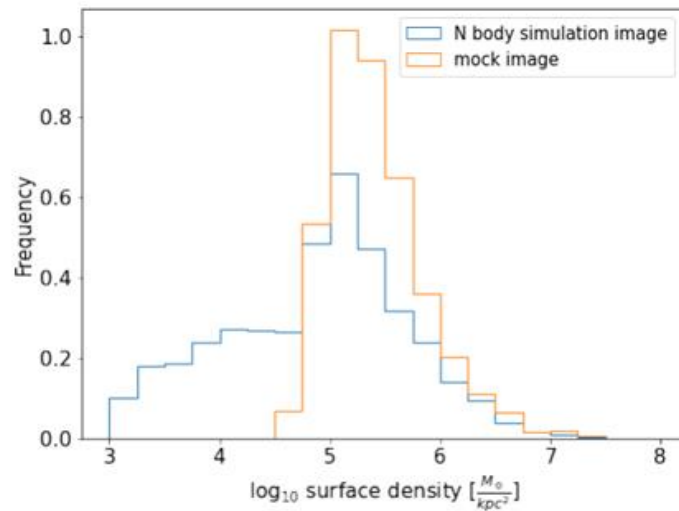
# Comparison of a N-body simulation and a mock image



(a) N body simulation image



(b) mock image



- **Using a toy halo model and MNRE**
  - we reconstructed the halo mass function
  - we generated images similar to DM-only N-body simulations

- **Using a toy halo model and MNRE**

- we reconstructed the halo mass function
- we generated images similar to DM-only N-body simulations

- **Long-term goal**

a analytical model for haloes, subhaloes, clustering and baryonic matter that generates actual N-body simulations

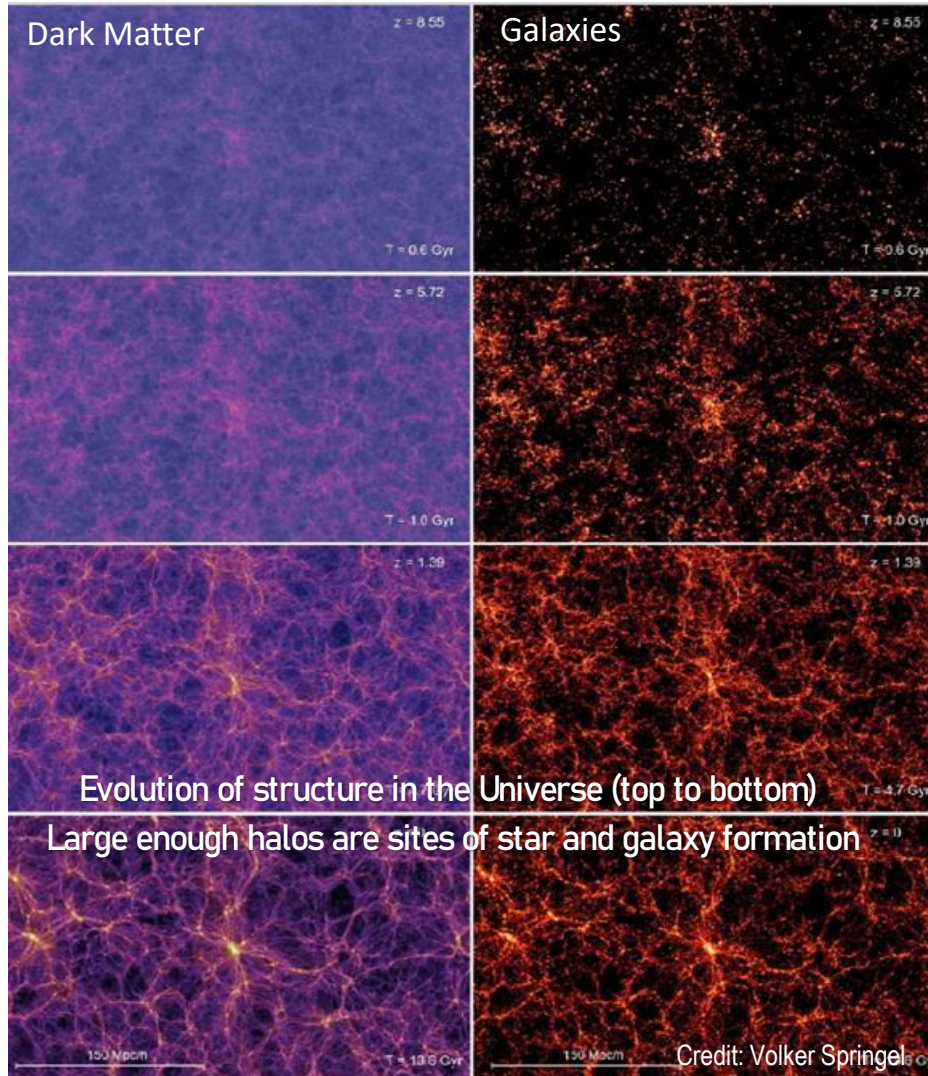
- **Using a toy halo model and MNRE**
  - we reconstructed the halo mass function
  - we generated images similar to DM-only N-body simulations

- **Long-term goal**

a analytical model for haloes, subhaloes, clustering and baryonic matter that generates actual N-body simulations

**Thank you!**

Backup slides



# Identifying Dark Matter Haloes

## FOF Halo Finders

Connect particles that are close to each other



Haloes with arbitrary shapes, i.e., it's difficult to assign a mass to them



Define a boundary: the virial radius  $r_h$  within which the mean internal matter density is:  $\rho_h = 200\rho_c^0$



The mass of the halo,  $M_h$ , is defined as the total mass contained within the radius  $r_h$

# The Eagle Project

- **Using the FOF halo finder:**

- Read halo masses.
- Pick haloes with masses that belong to a specific mass range, e.g.  $(10^9, 10^{12}) M_{\odot}$ .
- Define sub volumes of the  $25 \text{ Mpc}^3$  box, e.g.,  $12.5 \text{ Mpc}^3$ ,  $12.5 \times 12.5 \times 5 \text{ Mpc}^3$  boxes.
- Read the centers of potential (x, y, z coordinates) of those haloes.
- Pick the haloes that belong to those boxes and count their number.



# The Eagle Project

- **Using the Particle data:**

- Read particle group numbers
- Select particles that belong to the haloes that we picked before from the FOF halo finder
- Read the coordinates of those particles (x,y,z coordinates)
- Construct a heatmap with 100 bins by projecting the z direction
- Multiply the counts of the histogram with the mass of the particles:  $1.15 \cdot 10^7 M_{\odot}$  and divide with  $x_{\text{edge}} \cdot y_{\text{edge}}$  of the histogram i.e.,  $125 \cdot 10^2 \text{ kpc}^2$

# The Model

- The **first physical parameter** is the number,  $N$ , of the haloes.
- The masses of the haloes,  $M_h \in (10^9, 10^{12}) M_\odot$ , can be sampled from a halo mass function:  $\frac{dn}{dM} = bM^a$ .
- The **second physical parameter** is the slope,  $a$ , of the halo mass function, while:

$$b = (1 - a) \cdot \frac{N}{((10^{12})^{1-a} - (10^9)^{1-a})V}$$

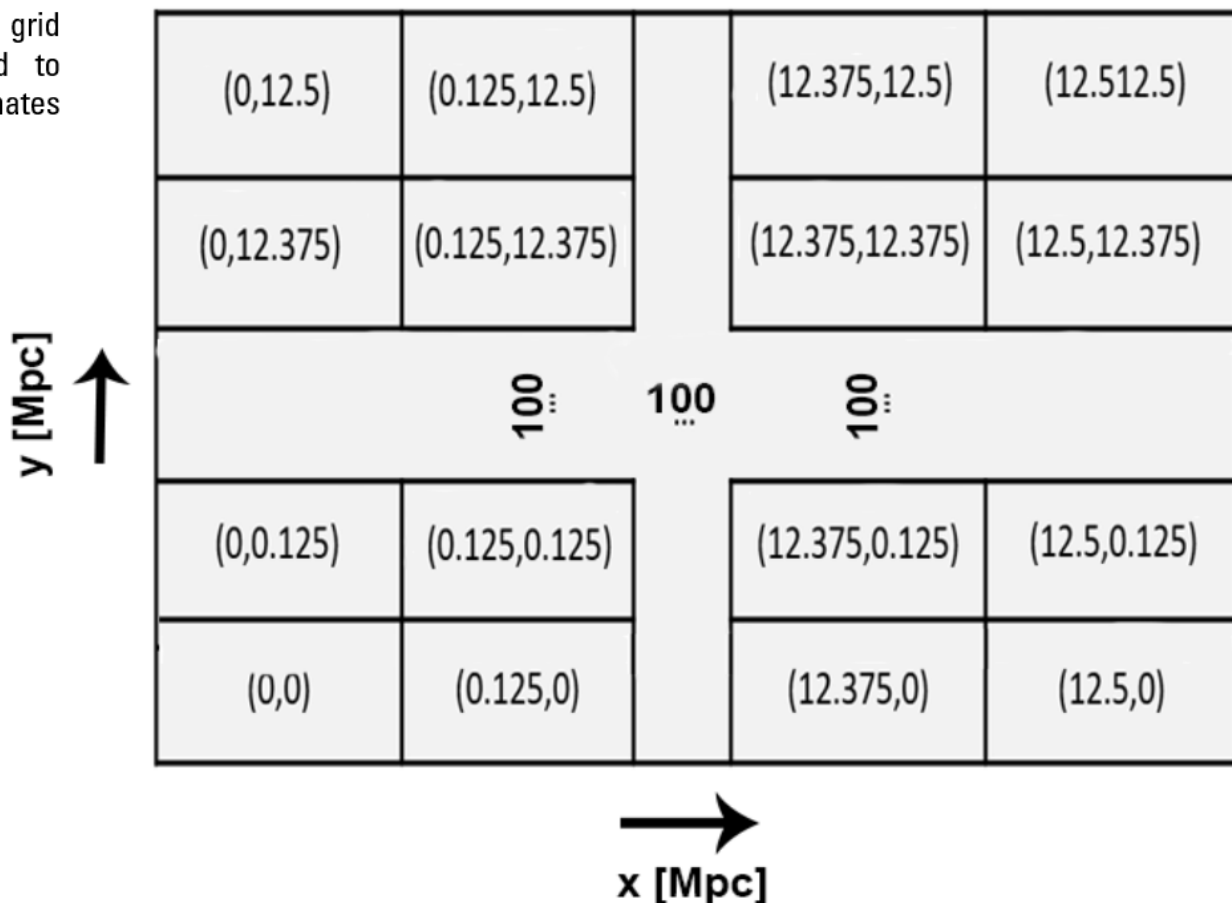
- From the masses of the haloes, we can calculate their concentrations  $c$ :

$$\log_{10}c = 1.4986 - 0.02499\log_{10}(M/M_\odot)[1 + 0.00565(\log_{10}(M/M_\odot)^2)] \quad (\text{Correa et al., 2015b})$$


- Now, we want to place the haloes in the 2D sky:



We construct a 100x100 grid whose values correspond to pairs of x and y coordinates where  $(x,y) \in (0,12.5)$  Mpc.



- Now, we want to place the haloes in the 2D sky:



➤ We construct a 100x100 grid whose values correspond to pairs of x and y coordinates where  $(x, y) \in (0, 12.5)$  Mpc.


# Adding Clustering to the Model

- Will sample the positions according to distributions generated from 2D realizations of gaussian random fields on an 100x100 grid.
- The gaussian fields will be specified by a power-law power spectrum:

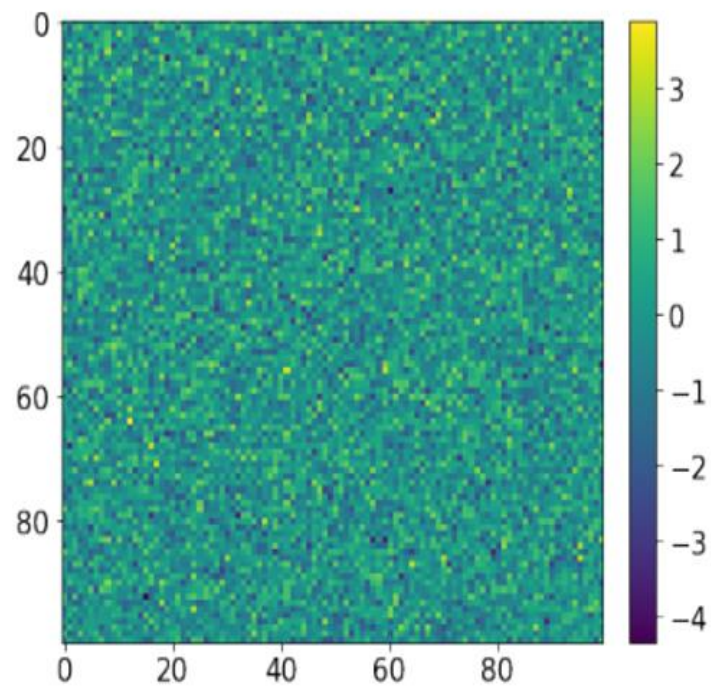
$$P(k) = \frac{1}{k^n}$$

- The slope of the power spectrum,  $n$ , is the **third physical parameter** of our model.

# Constructing the Realizations of the Gaussian Fields

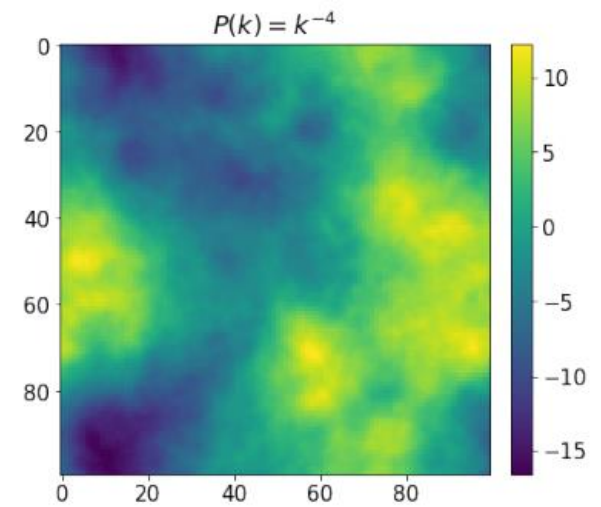
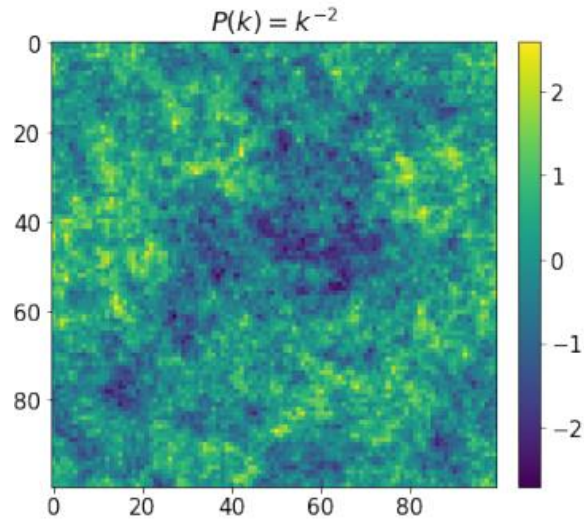
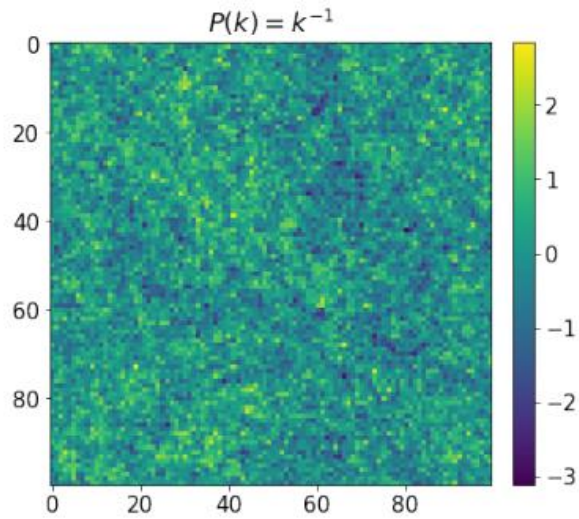
- We generate position space realization of a white noise field,  $\varphi_{ab}$ , with unit amplitude, on a 100x100 grid, i.e.,  $a, b \in \{0, \dots, 99\}$ .
- We Fourier transform the white noise realization:  $\varphi_{ab} \rightarrow \varphi_{k_a k_b}$ , where  $k_a, k_b \in \frac{2\pi}{N} \{0, \dots, 99\}$ .
- We want to multiply  $\varphi_{k_a k_b}$  with  $\sqrt{P(k)}$  to get  $\delta_{k_a k_b}$ .
- **Naive way:** Calculate  $P(k)$  at points  $k = \sqrt{k_a^2 + k_b^2}$   leads to imaginary fields
- **Alternatively:** Calculate  $P(k)$  at points  $k = \sqrt{k_a'^2 + k_b'^2}$ , where  $k'_a, k'_b \in \frac{2\pi}{N} \{0, \dots, 50, -49, \dots, -1\}$ .
- $\delta_{k_a k_b} = \sqrt{P(k)} \varphi_{k_a k_b}$
- $\delta_{k_a k_b} \rightarrow \delta_{ab}$

## Realization of White Noise Field



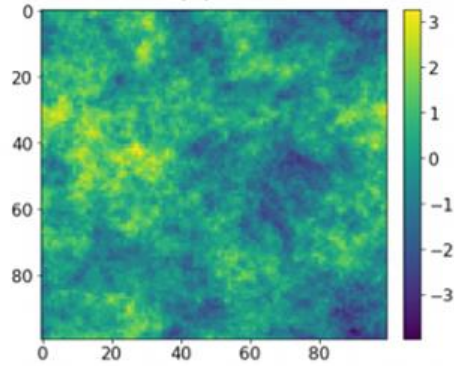


# Realizations of the Gaussian Fields



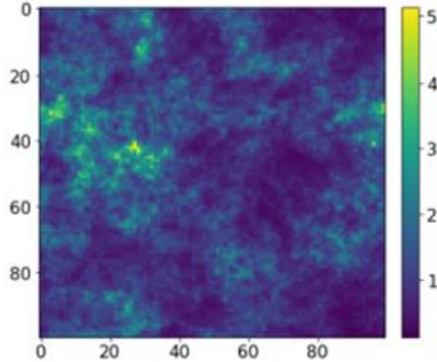
# The effect of parameter $\varepsilon$

$$P(k) = k^{-2.4}$$

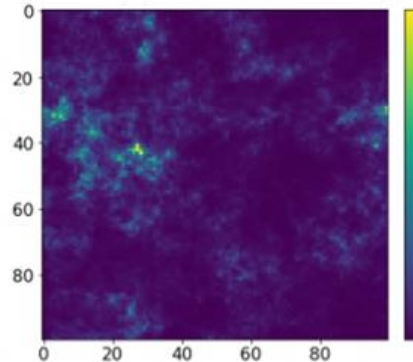


add parameter  $\varepsilon$

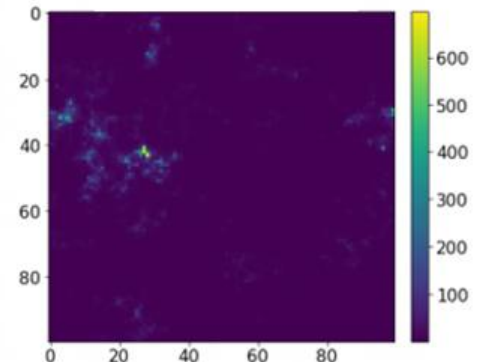
$\varepsilon=0.5$

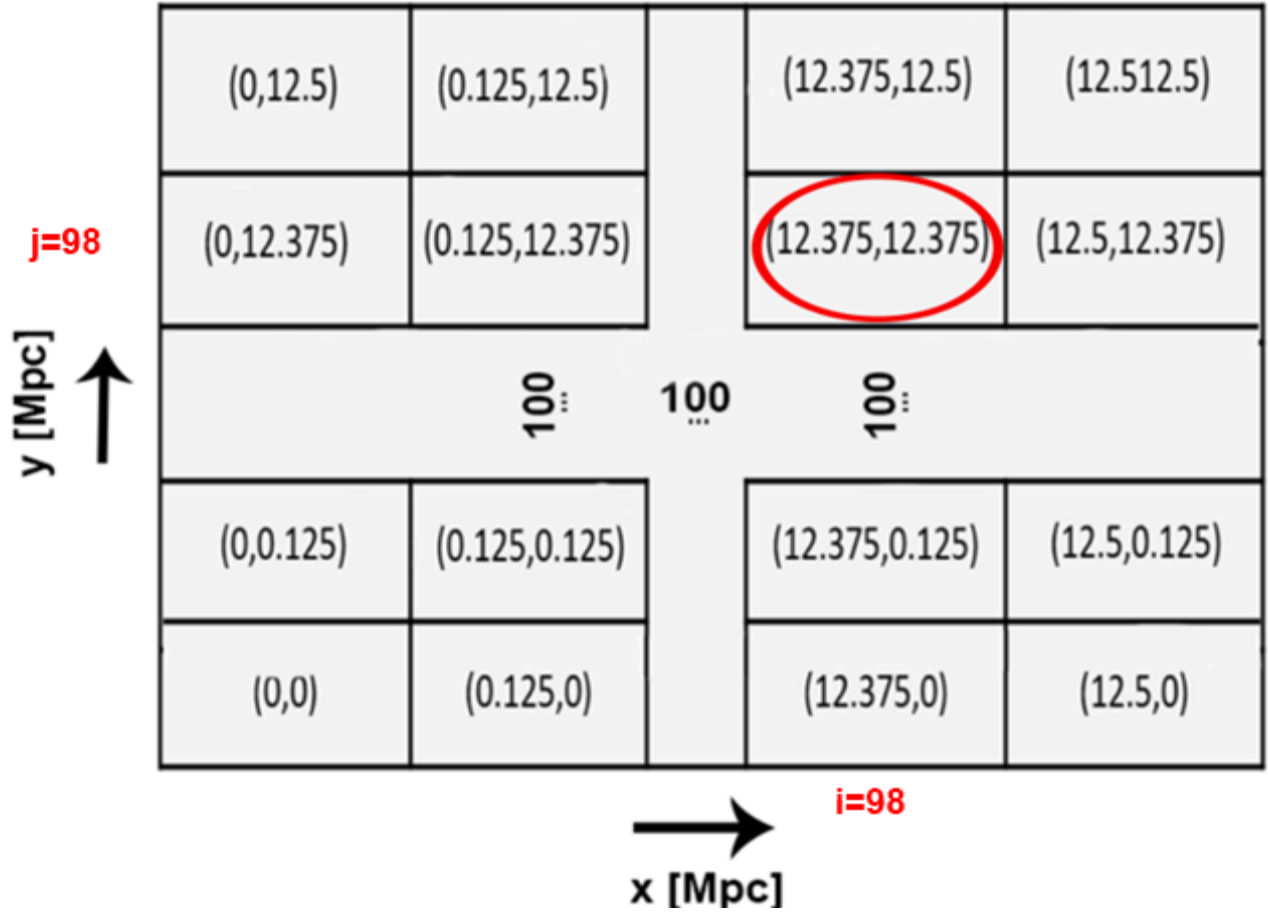


$\varepsilon=1$




$\varepsilon=2$

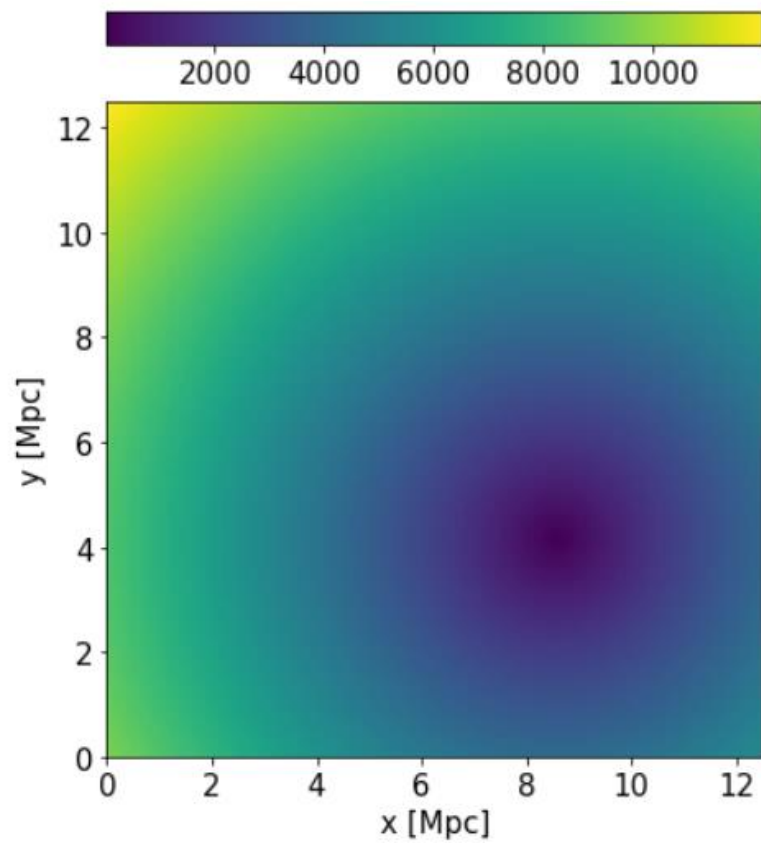




- Now, we want to place the haloes in the 2D sky:

- 
- We construct a 100x100 grid whose values correspond to pairs of  $x$  and  $y$  coordinates where  $(x, y) \in (0, 12.5)$  Mpc.
  - For each one of the  $N$  haloes, we sample its coordinates  $X, Y$  coordinates from distributions generated from 2D realizations of gaussian random fields.
  - We subtract the coordinates of each halo (as pairs of  $X, Y$  values) from the values  $x, y$  of the grid and we end up with  $N$  grids.
  - For each grid, we calculate the root sum square of the two values in each one of its cells, i.e., the projected radius of the halo.

Projected radius  $r'$  [kpc]



- Now, we want to place the haloes uniformly in the 2D sky:

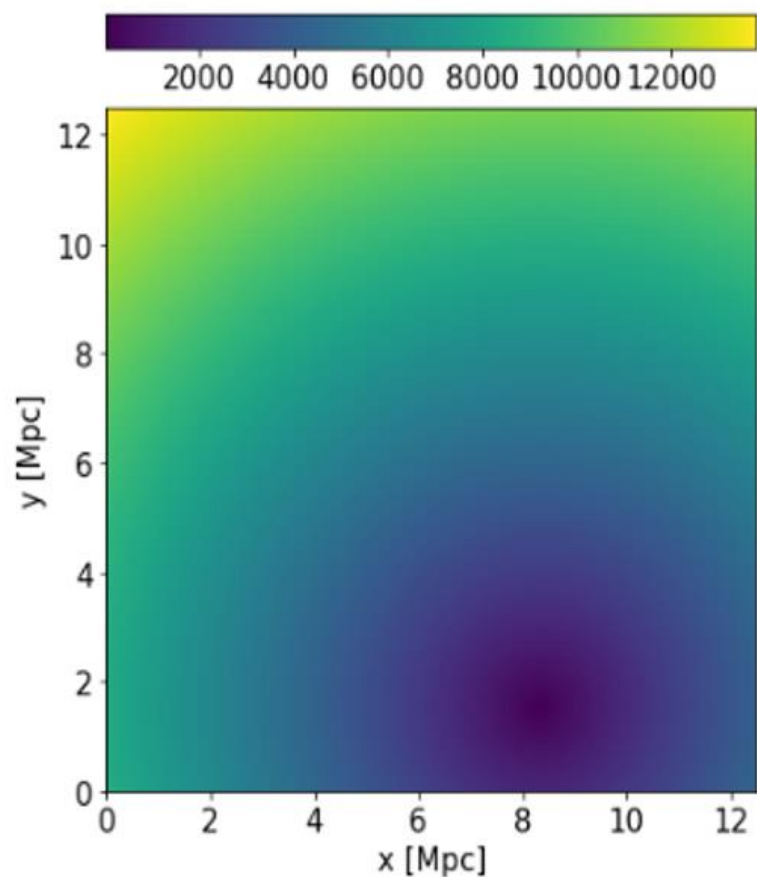
- For each halo we calculate the surface density:

$$f(r') = \frac{2r_s \left[ r' \cdot (r_s - r'^2) + 2r_s \cdot r' \cdot \arctan \left( \frac{r' \cdot \sqrt{r' - r_s}}{r' \cdot \sqrt{r' + r_s}} \right) \cdot \sqrt{r'^2 - r_s^2} \right]}{r' \cdot (r_s^4 - 2r_s^2 r'^2 + r'^4)}$$

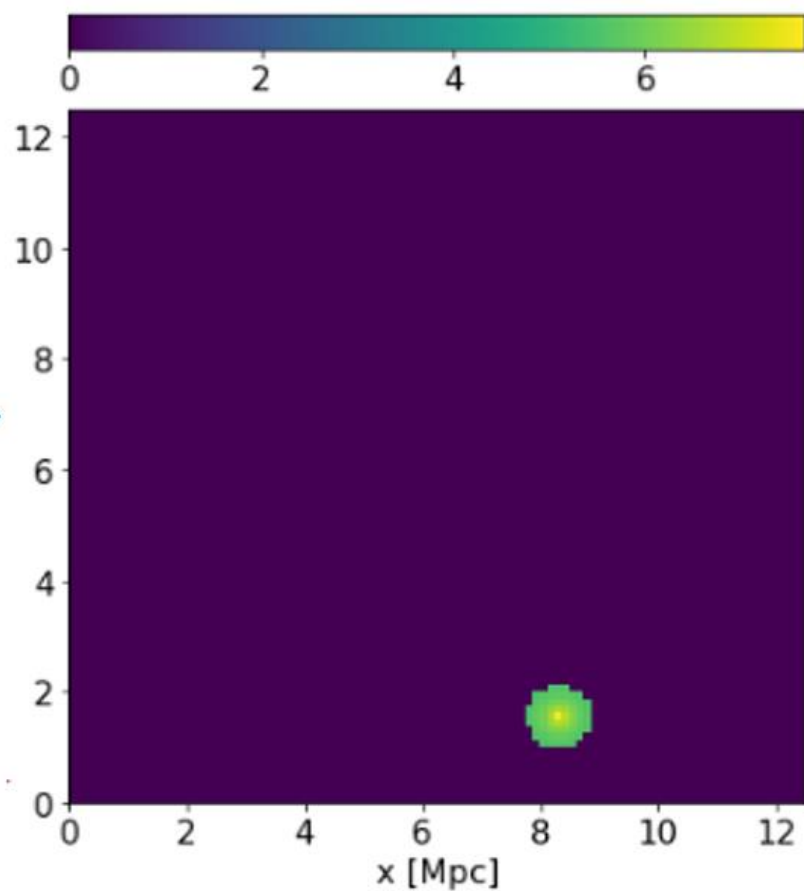
where:  $r_s = \frac{(3M_h)^{1/3}}{(800\pi\rho_c^0)^{1/3}c}$

- We set the values of the pixels that correspond to  $r' > 2.7r_h$  equal to 1.
- We add all the images of the individual haloes together to obtain the total surface density field
- We add poisson noise to the final image.

Projected radius  $r'$  [kpc]



$\log_{10}$  Surface Density  $[\frac{M_{\odot}}{\text{kpc}^2}]$  of a halo with mass  
 $M_h = 10^{12} M_{\odot}$



# Joint vs marginal samples

1) Examples for  $H_0$ , jointly sampled from  $(x, \vartheta) \sim p(x|\vartheta) p(\vartheta)$



2) Examples for  $H_1$ , marginally sampled from  $(x, \vartheta) \sim p(x) p(\vartheta)$



Data:  $x$ =Image, Label:  $\vartheta \in \{\text{Cat}, \text{Donkey}\}$



# Loss function

- **Strategy:** We train a neural network  $d_\varphi(x, \vartheta) \in [0, 1]$  as binary classifier to estimate the probability of hypothesis  $H_0$  or  $H_1$ . The Network output can be interpreted, for a given input pair  $x$  and  $\vartheta$ , as probability that  $H_0$  is true.
  - $H_0$  is true:  $d_\varphi(x, \vartheta) \simeq 1$
  - $H_1$  is true:  $d_\varphi(x, \vartheta) \simeq 0$
- The corresponding loss function is the so-called “binary cross-entropy”:

$$L[d(x, \vartheta)] = - \int dx d\vartheta [p(x, \vartheta) \ln(d(x, \vartheta)) + p(x)p(\vartheta) \ln(1 - d(x, \vartheta))]$$

- Minimizing that function w.r.t the network parameters  $\varphi$  yields:

$$d_\varphi(x, \vartheta) \approx \frac{p(x, \vartheta)}{p(x, \vartheta) + p(x)p(\vartheta)}$$

# Likelihood-to-evidence ratio

Training binary classification networks yield true Bayesian posterior estimates!

- With a bit of math one can show that:

$$r(\mathbf{x}, \boldsymbol{\vartheta}) \equiv \frac{d_{\varphi}(\mathbf{x}, \boldsymbol{\vartheta})}{d_{\varphi}(\mathbf{x}, \boldsymbol{\vartheta}) - 1} \approx \frac{p(\mathbf{x}|\boldsymbol{\vartheta})}{p(\mathbf{x})} = \frac{p(\boldsymbol{\vartheta}|\mathbf{x})}{p(\boldsymbol{\vartheta})}$$

- Once we have trained the network  $d_{\varphi}(\mathbf{x}, \boldsymbol{\vartheta})$ , we can **estimate the posterior**:

$$p(\boldsymbol{\vartheta}|\mathbf{x}) \simeq r(\mathbf{x}, \boldsymbol{\vartheta})p(\boldsymbol{\vartheta})$$

- Swyft: a flexible and powerful tool for efficient marginal posterior estimation using NN, designed by B. Miller et al. (2020, 2021).

- **Until now, using our model and Swyft we were able:**
  - to reconstruct the halo mass function
  - and to produce images like those of N body simulations!
- As a **next step** we can test if we can identify the lowest mass of the haloes in these images.

To do that:

- Instead of sampling haloes with masses:  $M_h \in (10^9, 10^{12}) M_\odot$ , we will sample masses:  $M_h \in (10^c, 10^{12}) M_\odot$ ,
- where  $c$  is a **new parameter** of our model.
- We set the values of the parameters  $a$ ,  $\varepsilon$  and  $n$  equal to the modes of their combined posteriors.

# Results on actual N body simulations

