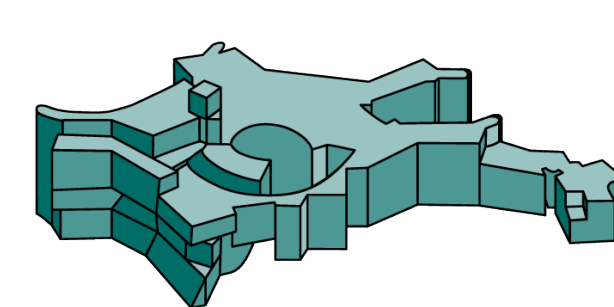# *Explainable deep learning models for cosmological structure formation*

**Luisa Lucie-Smith**

*Postdoctoral Research Fellow @ Max-Planck-Institute for Astrophysics*

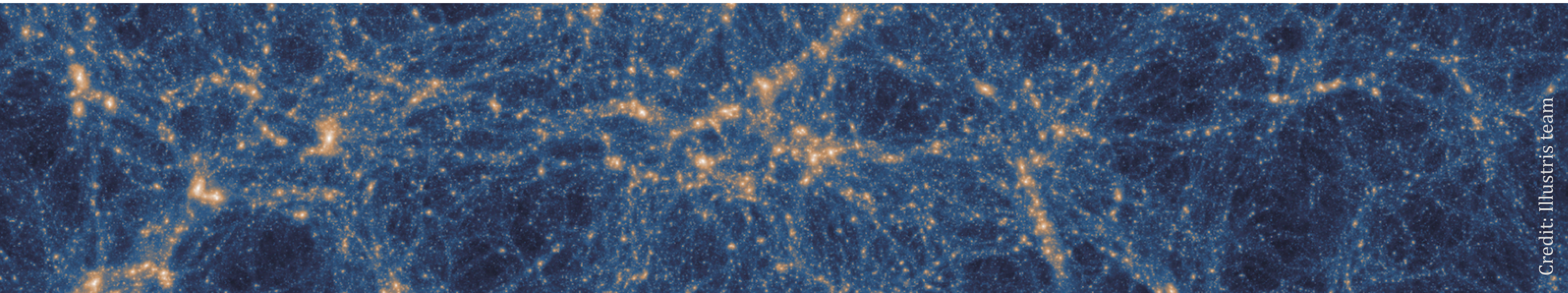European AI for Fundamental Physics Conference
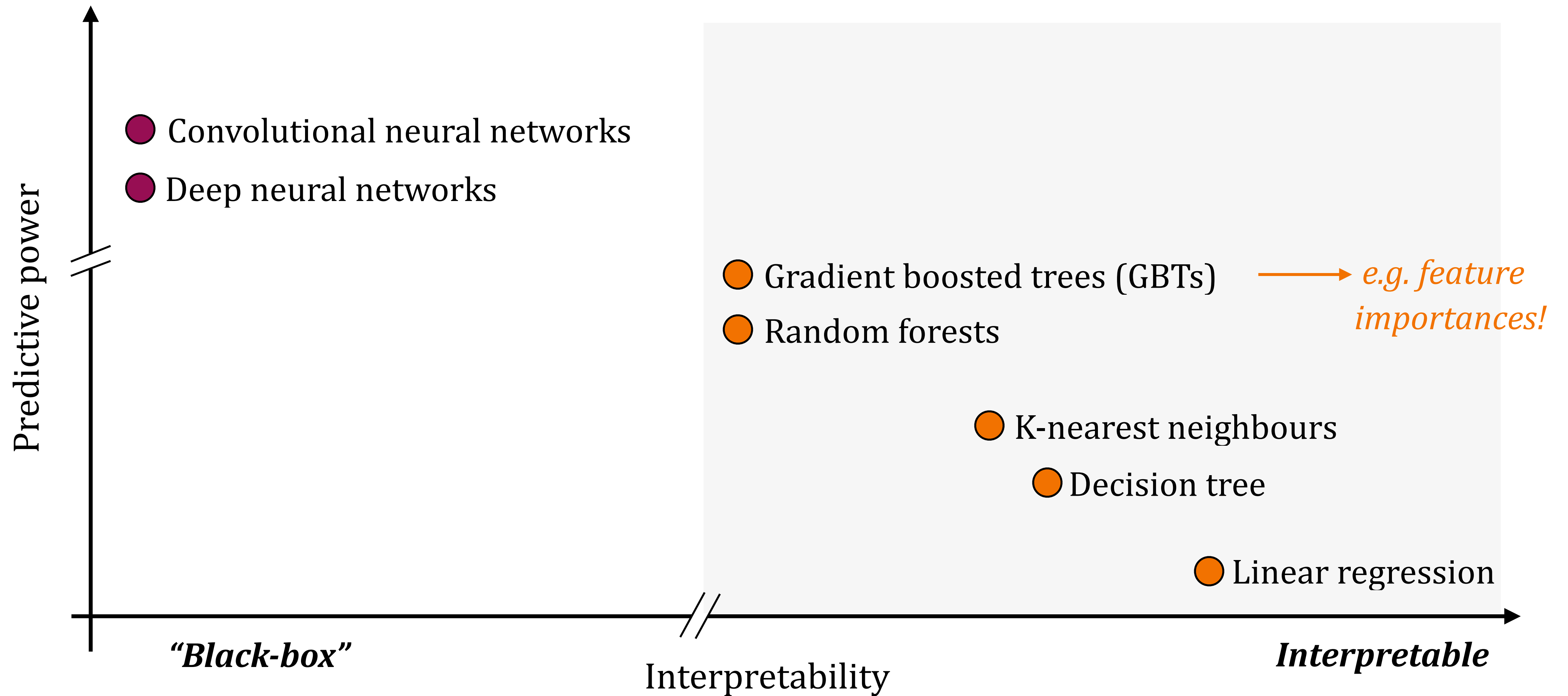Amsterdam, 1st May 2024

**MAX PLANCK INSTITUTE**
FOR ASTROPHYSICS

# Machine learning in (Astro)physics

- *ML successful at automating/accelerating known physical models (e.g. emulators)*

- *Can we **extract new knowledge** about the underlying physics from deep learning models by interpreting their outputs? Requires **explainable AI***
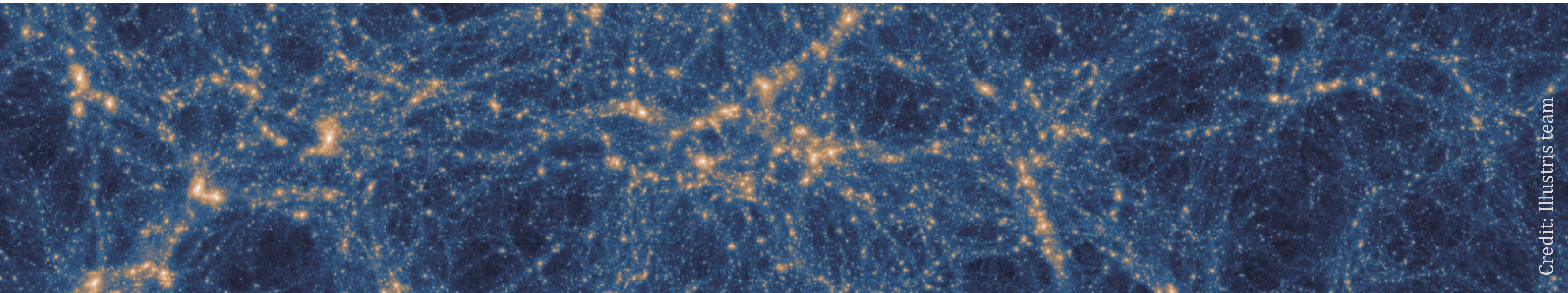
# Current landscape for explainable AI

# Requirements for explainable AI
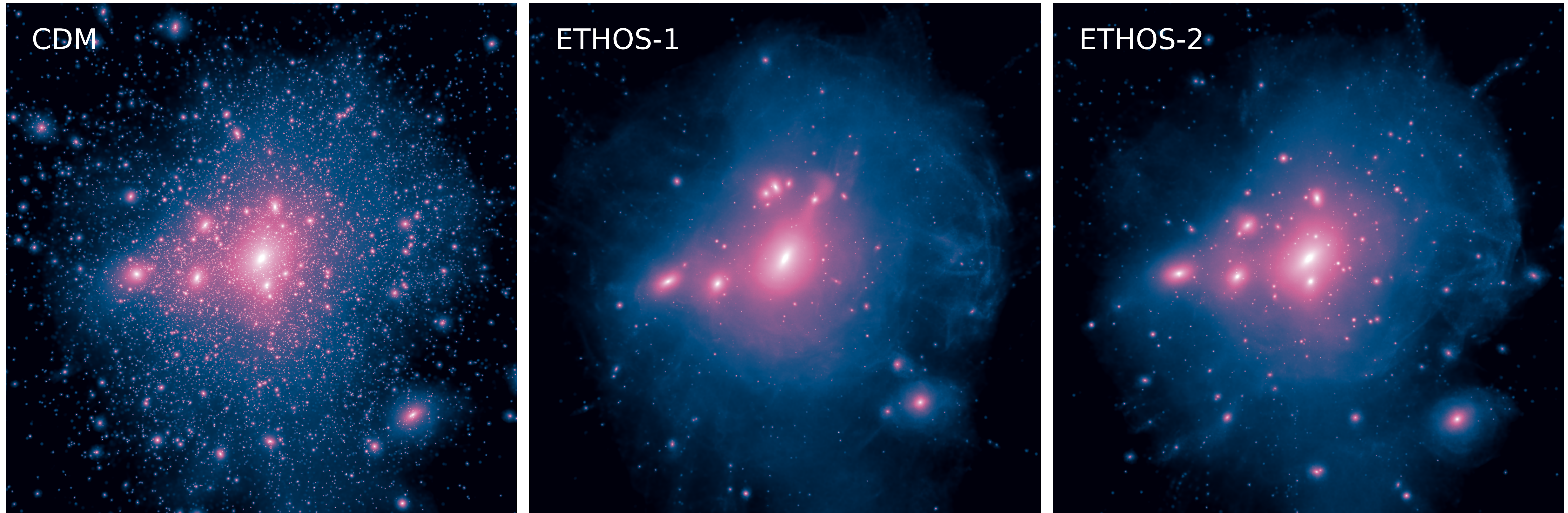
1. ***Interpretability:*** account for why the ML model reaches its predictions

2. ***Explainability:*** map this account onto existing knowledge in the relevant science domain

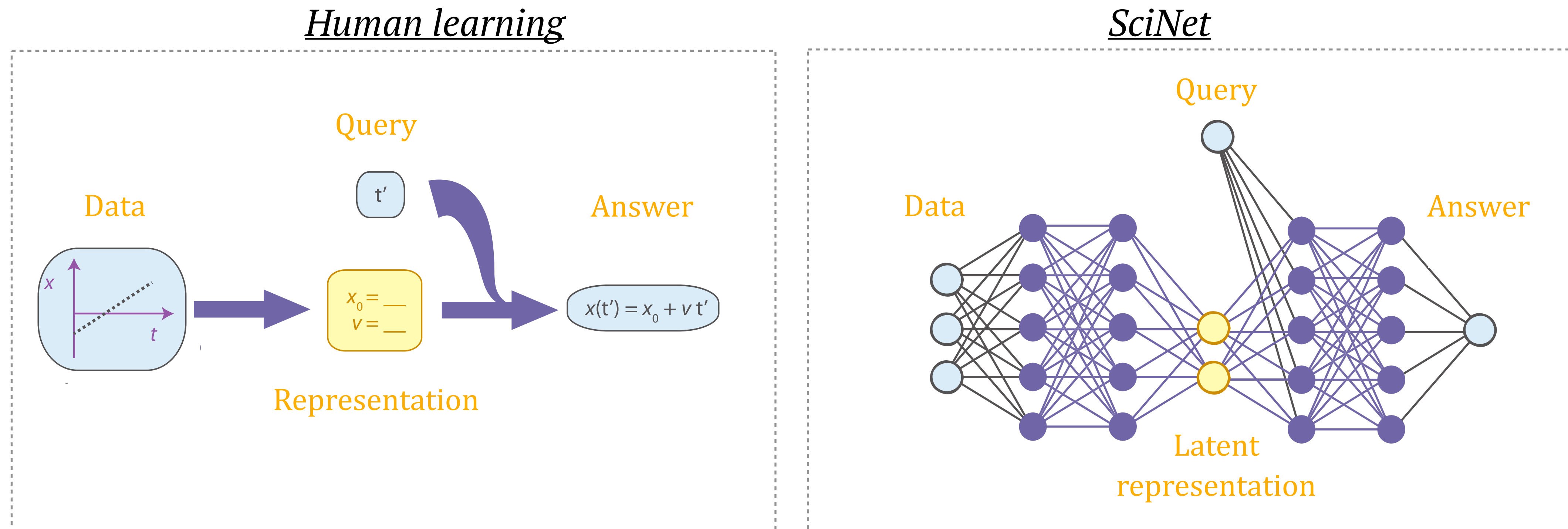N.B.: *many physical models in cosmology are also not explainable!*

# Dark matter halo structure contains signatures of nature of dark matter

CDM

ETHOS-1

ETHOS-2

*Current models based on 'universal' empirical fitting functions (e.g. NFW) lack **explainability***

# SciNet model

### Human learning

Query

Data

$t'$

Answer

$x_0 = \underline{\quad}$
$v = \underline{\quad}$

$x(t') = x_0 + v\,t'$

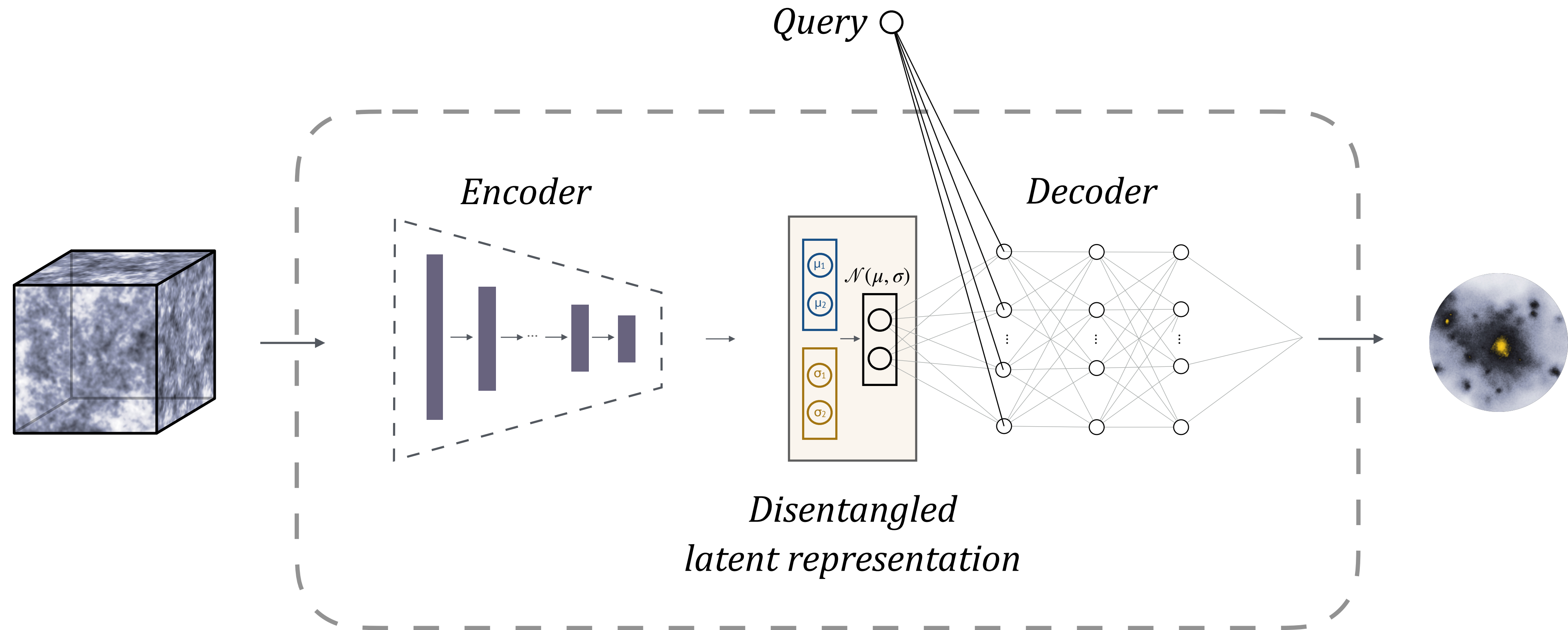Representation

### SciNet



Query

Data

Answer

Latent
representation

- SciNet learns relevant physical parameters in toy 1D problems

- Relies on comparing latents with already-known physical parameters
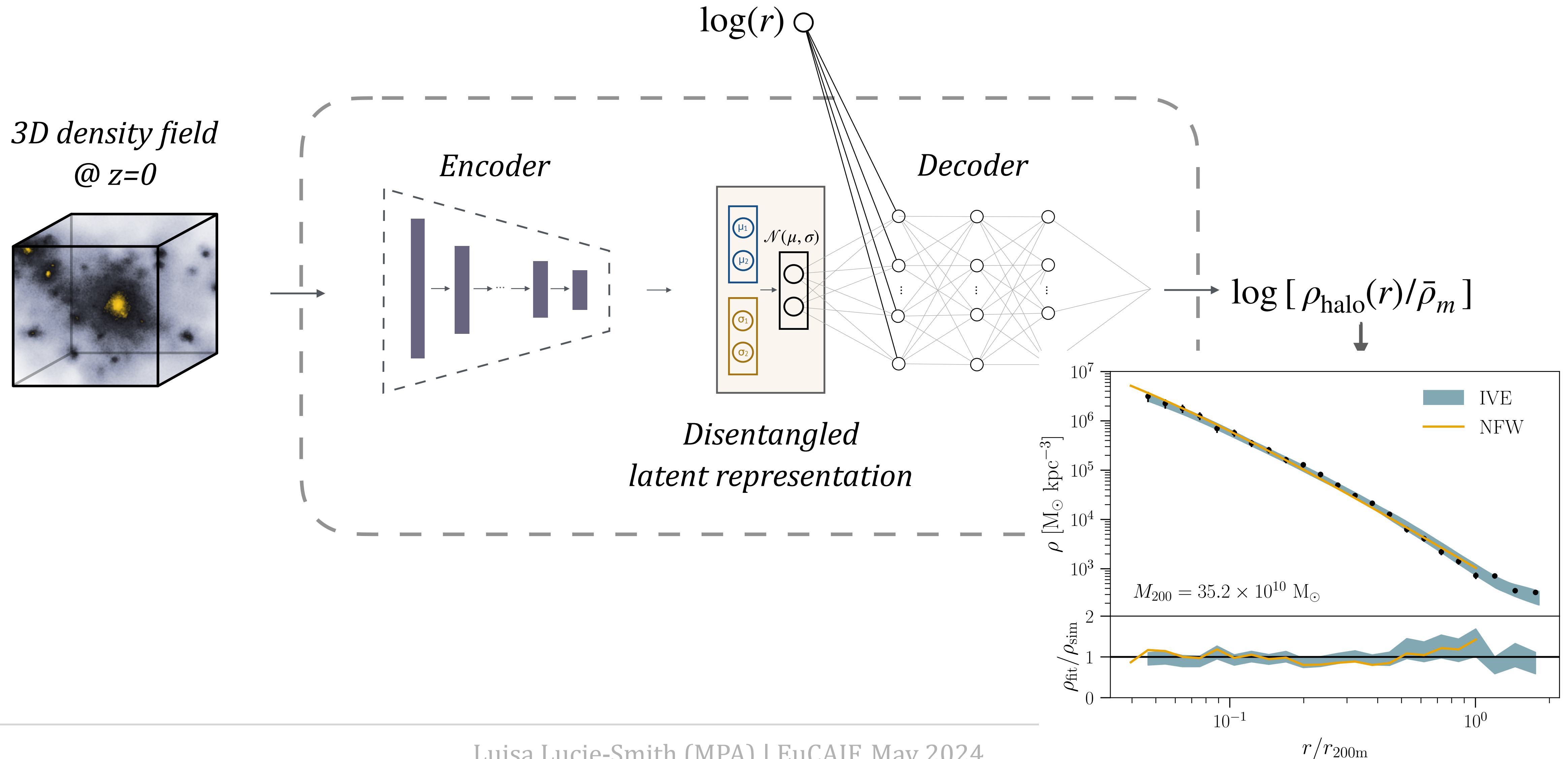
Iten et al. (PRL, 2020)

# Interpretable Variational Encoder (IVE) for explainable AI



Query

Encoder

Decoder

$\mathcal{N}(\mu, \sigma)$

$\mu_1$
$\mu_2$

$\sigma_1$
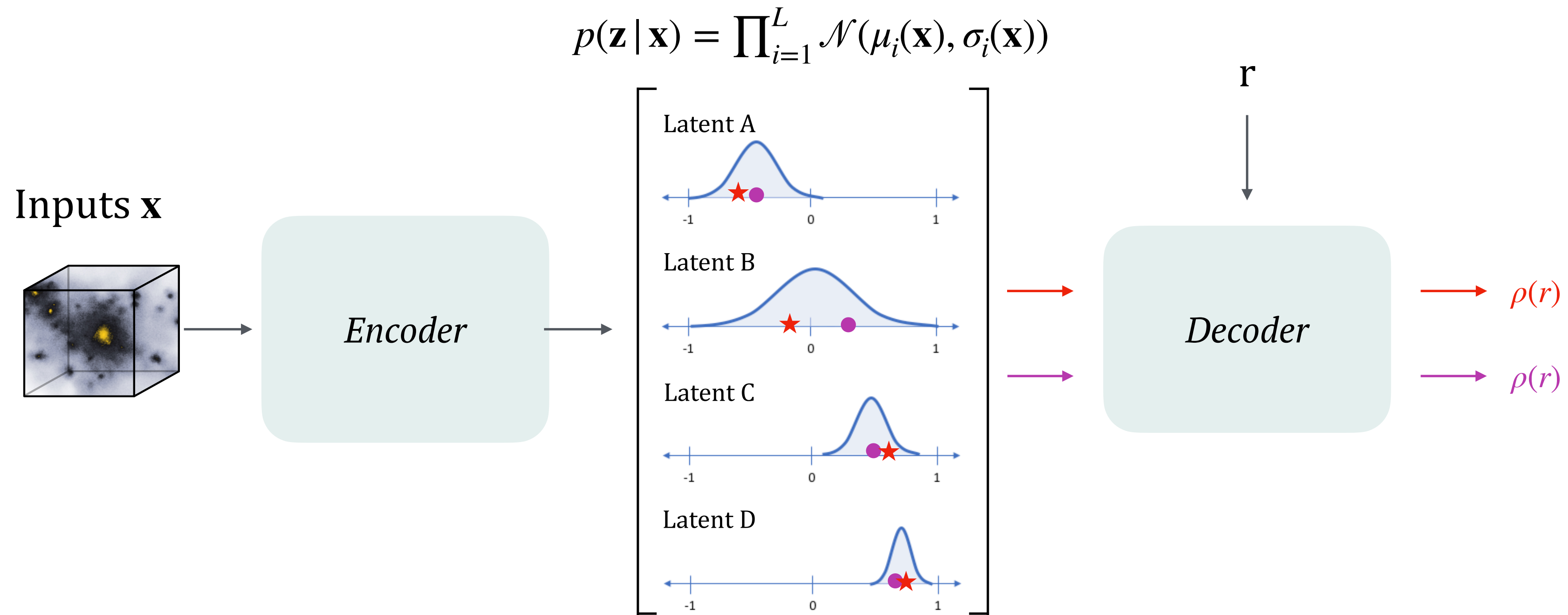$\sigma_2$

Disentangled
latent representation

## Model compression enables "explainability"

Iten et al. (PRL, 2020); Lucie-Smith et al. (PRD, 2022); Lucie-Smith et al. (PRL, 2024)

# Discovering the building blocks of halo density profiles out to the halo outskirts
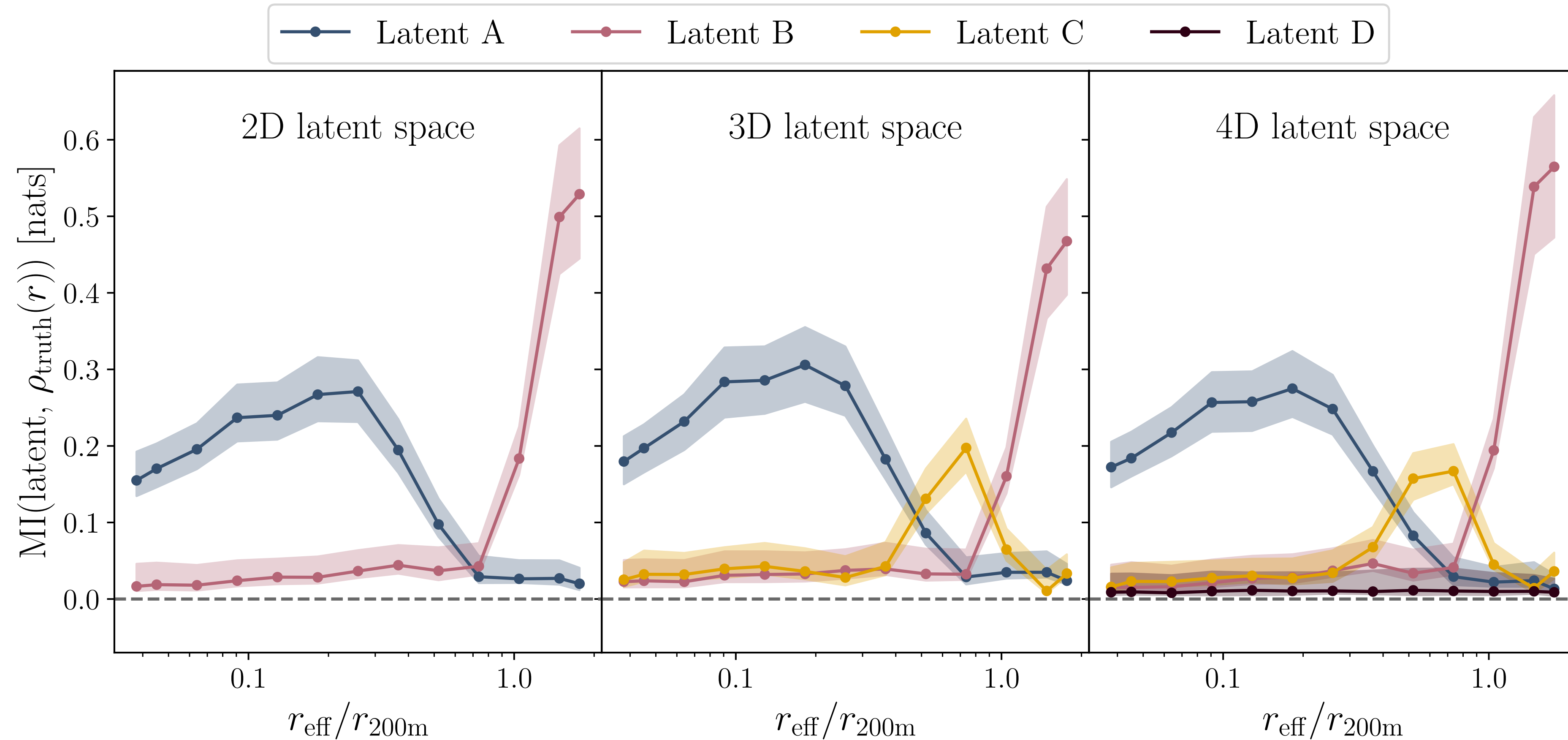


$\log(r)$

*3D density field @ z=0*

*Encoder*

*Decoder*

$\mu_1$ $\mu_2$ $\mathcal{N}(\mu, \sigma)$ $\sigma_1$ $\sigma_2$

*Disentangled latent representation*

$\log \left[ \rho_{\mathrm{halo}}(r)/\bar{\rho}_m \right]$

$\rho \; [\mathrm{M}_\odot \; \mathrm{kpc}^{-3}]$

IVE

NFW

$M_{200} = 35.2 \times 10^{10} \; \mathrm{M}_\odot$

$\rho_{\mathrm{fit}}/\rho_{\mathrm{sim}}$

$r/r_{200\mathrm{m}}$

# Desired latent representation properties for interpretability



$$p(\mathbf{z}\,|\,\mathbf{x}) = \prod_{i=1}^{L} \mathcal{N}(\mu_i(\mathbf{x}), \sigma_i(\mathbf{x}))$$

Inputs $\mathbf{x}$

Encoder

Latent A

Latent B

Latent C

Latent D

r

Decoder

$\rho(r)$

$\rho(r)$

- **Interpretability** can be achieved if latent space is **disentangled:** independent factors of variation in profiles captured by different, independent latents

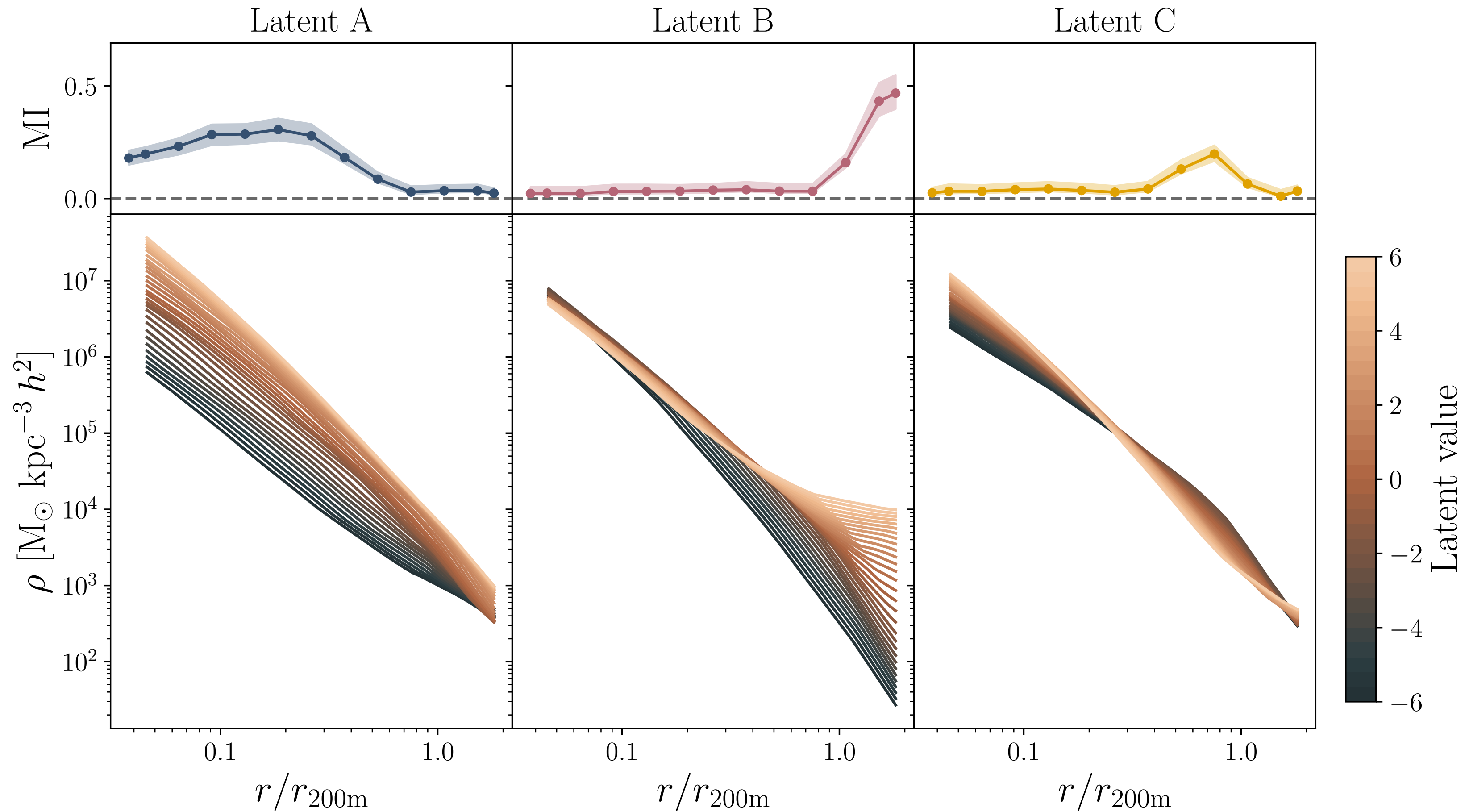- Disentanglement encouraged via **loss function** optimised during training

# Interpreting the latent representation using *mutual information*



*Explainability* achieved by evaluating MI between latents and density profile

Lucie-Smith, Peiris, Pontzen, Nord et al. (Phys. Rev. D, 2022)
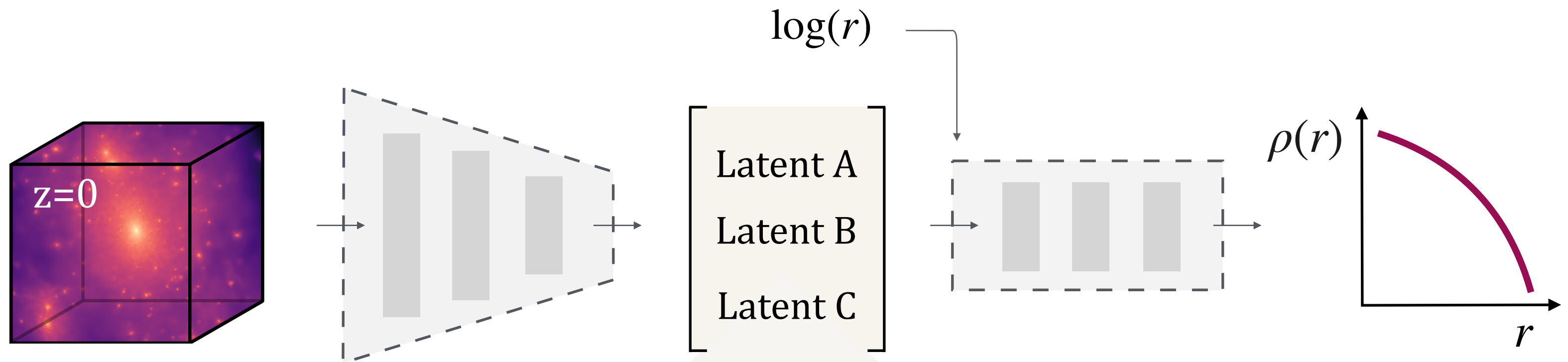
# Systematically varying one latent at a time
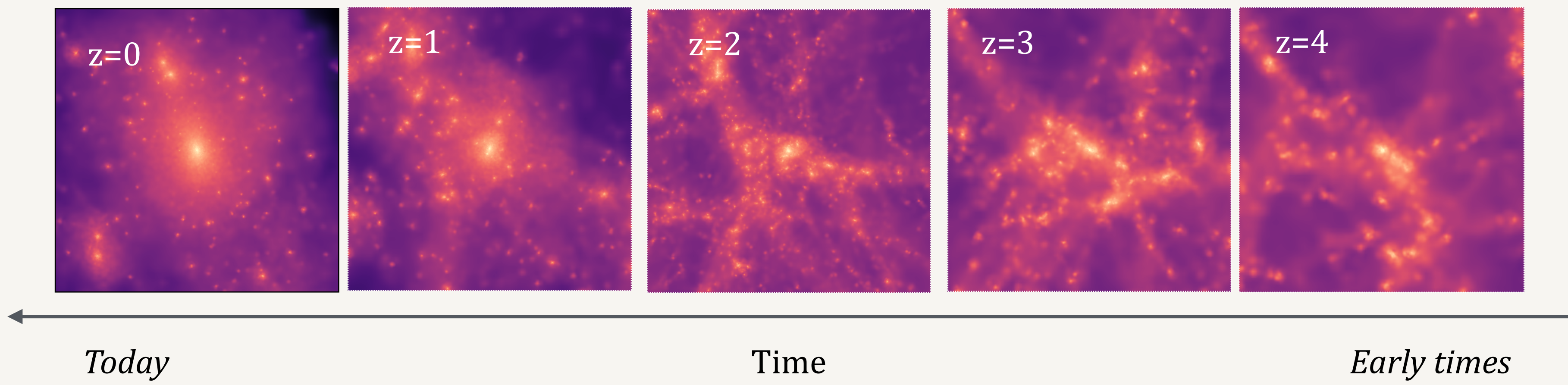


Latent A = **normalisation**;   Latent B = **outer slope**;   Latent C = **inner slope**

Lucie-Smith, Peiris, Pontzen, Nord et al. (Phys. Rev. D, 2022)

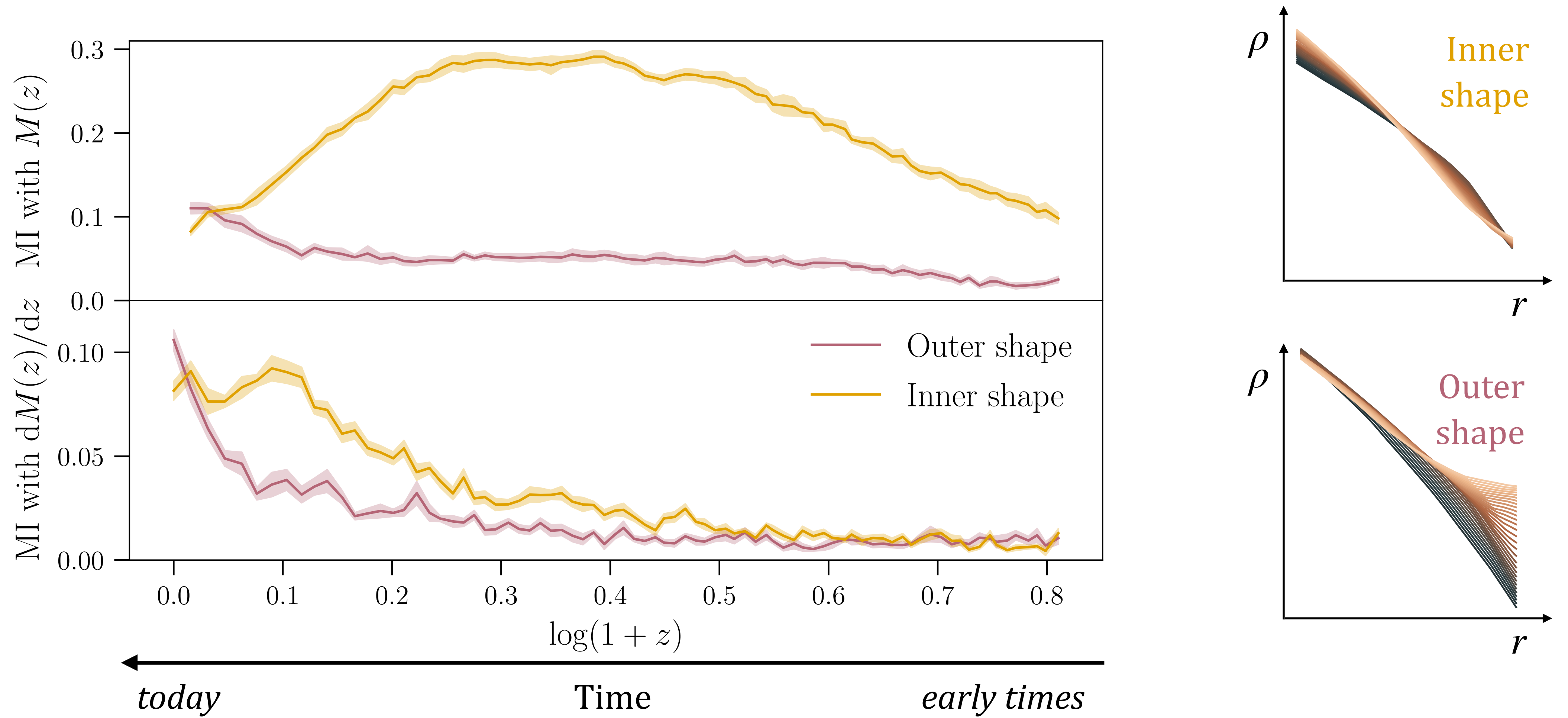# Exploiting the latent representation beyond its original training task



Does the latent space contain information about the origin of the halo density structure?

Today

Time

Early times

Lucie-Smith, Peiris, Pontzen (Phys. Rev. Lett., 2024)

# Connection between the latents and the *halo evolution history*



$\rho$

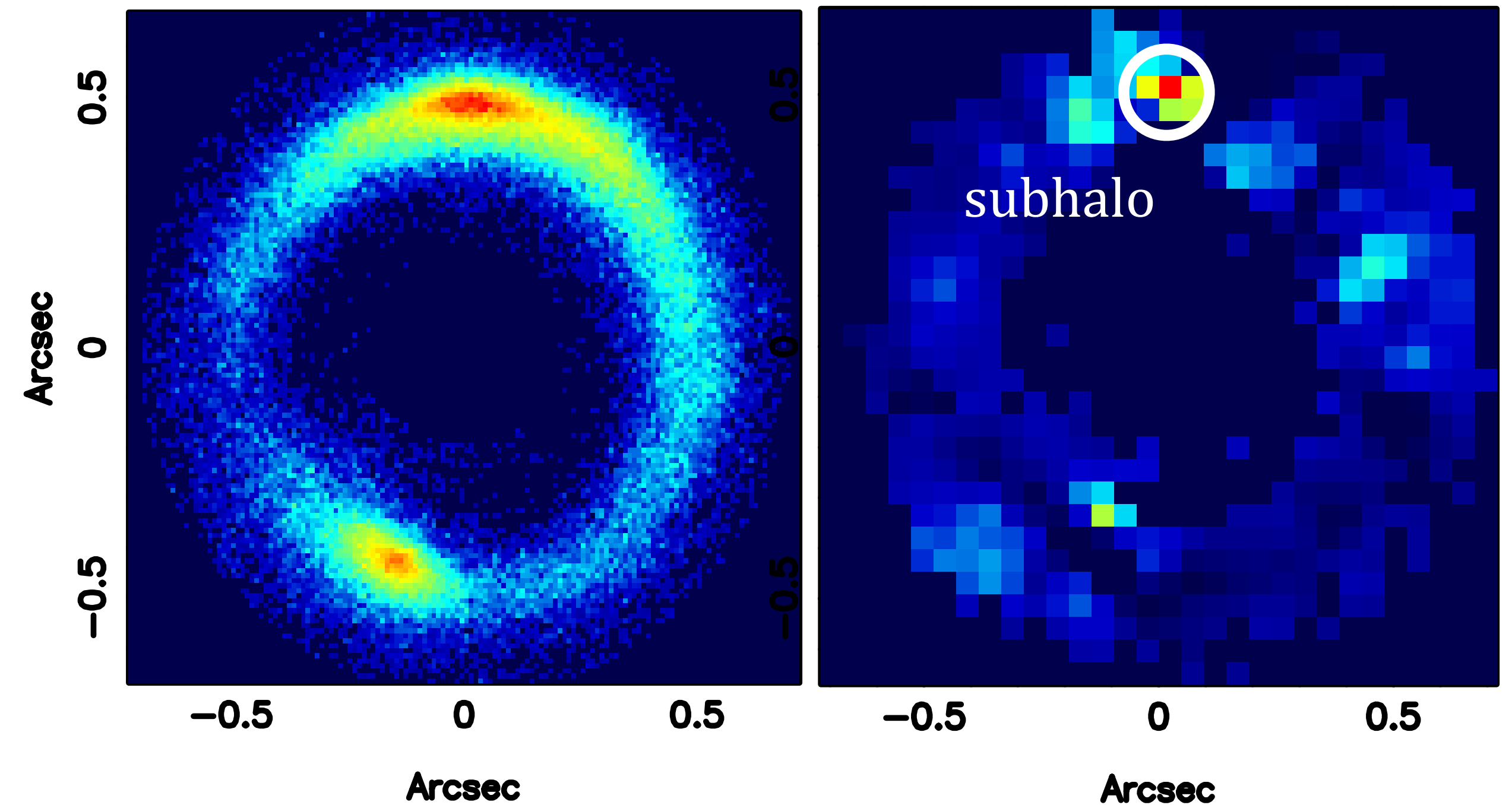Inner shape

$\rho$

Outer shape

$r$

MI with $M(z)$

MI with $\mathrm{d}M(z)/\mathrm{d}z$

0.3

0.2

0.1

0.0

0.10

0.05

0.00

— Outer shape
— Inner shape

0.0   0.1   0.2   0.3   0.4   0.5   0.6   0.7   0.8

$\log(1+z)$

*today*   Time   *early times*

# *What about* *dark matter substructures?*

**Simulations**

**Strong lensing observations**

CDM

subhalo
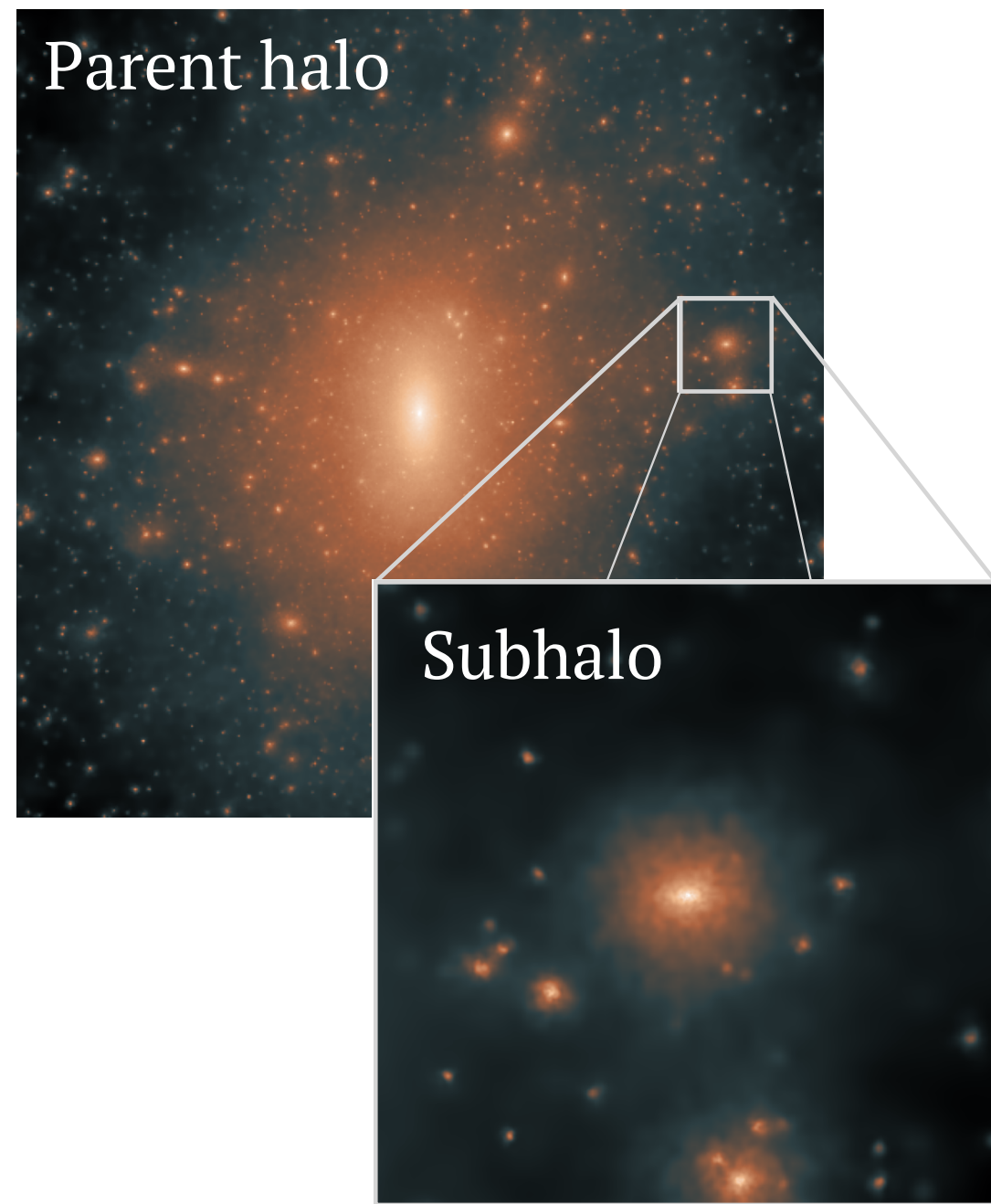
subhalo
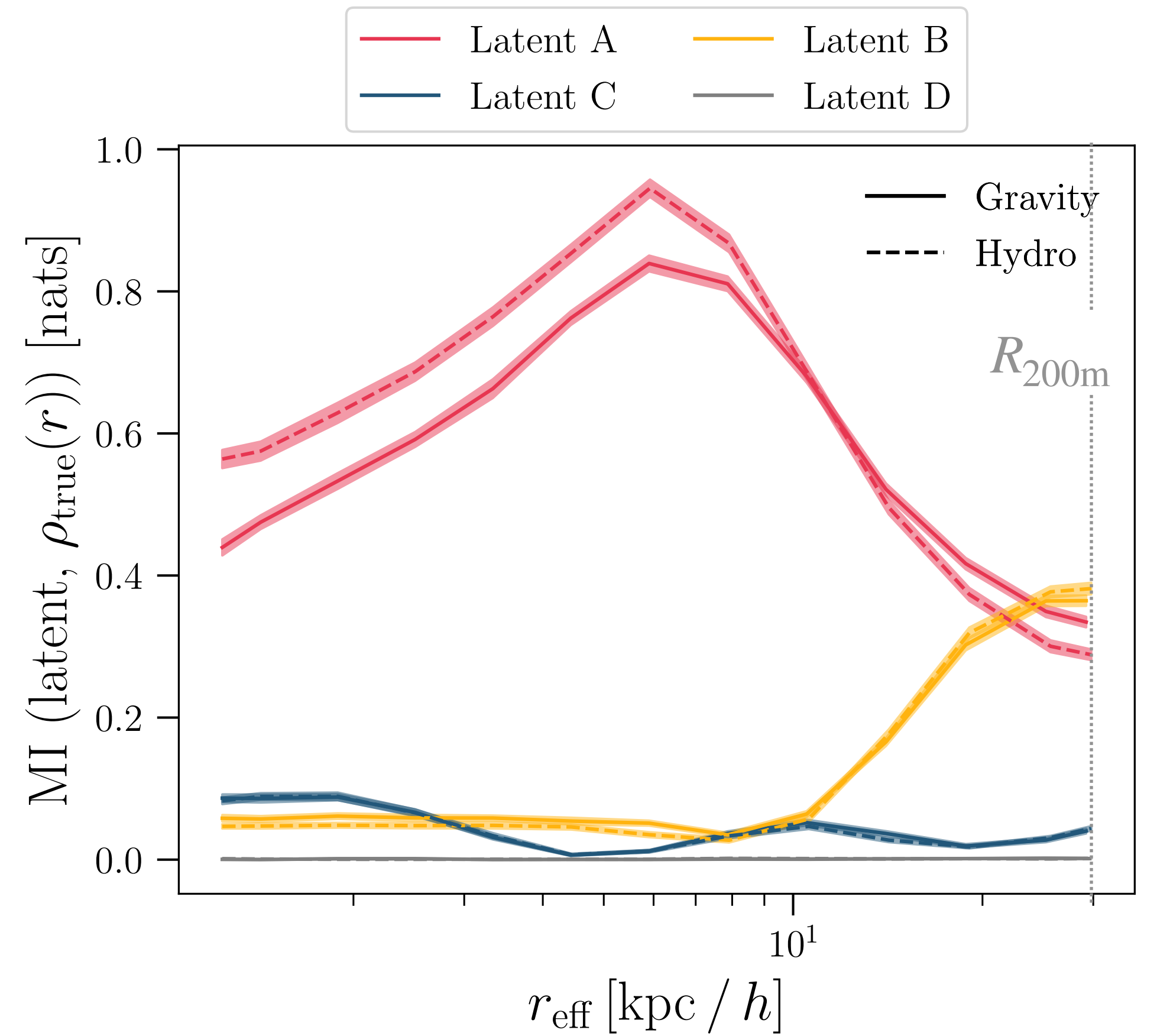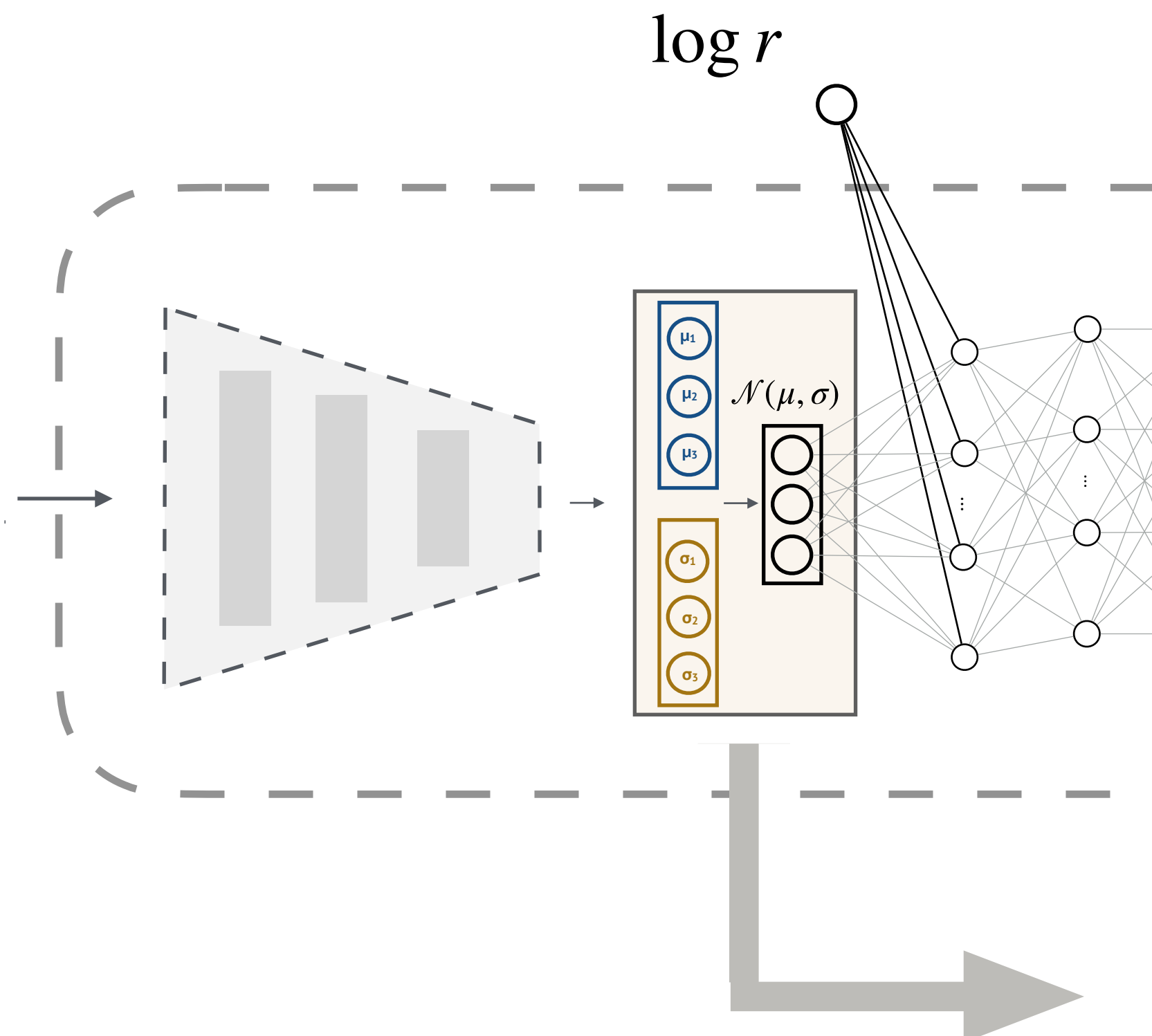
0.5

0

−0.5

−0.5   0   0.5

**Arcsec**

−0.5   0   0.5

**Arcsec**

JVAS B1938+666; Vegetti et al. (2012)

*Inferred subhalo properties* depends strongly on *assumed density profiles*

Luisa Lucie-Smith (MPA)

# *What about* **subhalos** *($r < R_{200\mathrm{m}}$)?*



Parent halo

Subhalo

*Illustris-TNG100*
*(Gravity & Hydro)*

$\log r$
$\log(r)$

$\mu_1$
$\mu_2$
$\mu_3$
$\mathcal{N}(\mu, \sigma)$
$\sigma_1$
$\sigma_2$
$\sigma_3$

Latent A — Latent B
Latent C — Latent D

— Gravity
-- Hydro

$R_{200\mathrm{m}}$

MI (latent, $\rho_{\mathrm{true}}(r)$) [nats]

$r_{\mathrm{eff}}$ [kpc / $h$]

Lucie-Smith, Despali, Springel (MNRAS, 2024)

# Subhalos require additional latent capturing *tidal truncation*



'Normalization'     'Truncation'/'Flattening'     'Inner shape'

Gravity

Hydro

$\log[\rho/\bar{\rho}_{\mathrm{M}}]$

$r/r_{\mathrm{vir}}$

Latent value — Max / Min

Lucie-Smith, Despali, Springel (MNRAS, 2024)
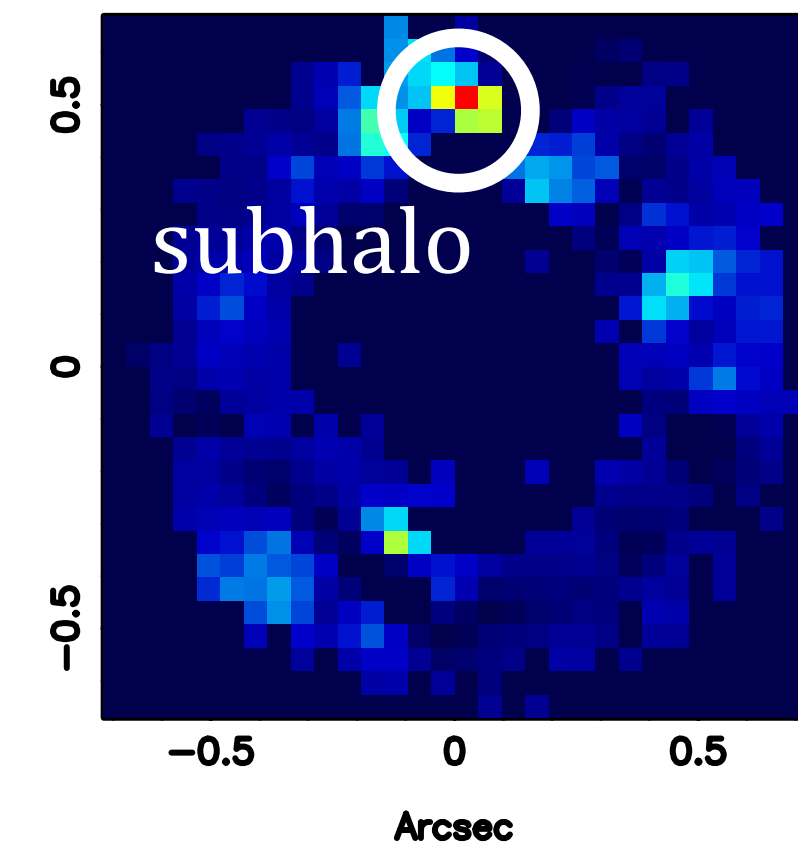
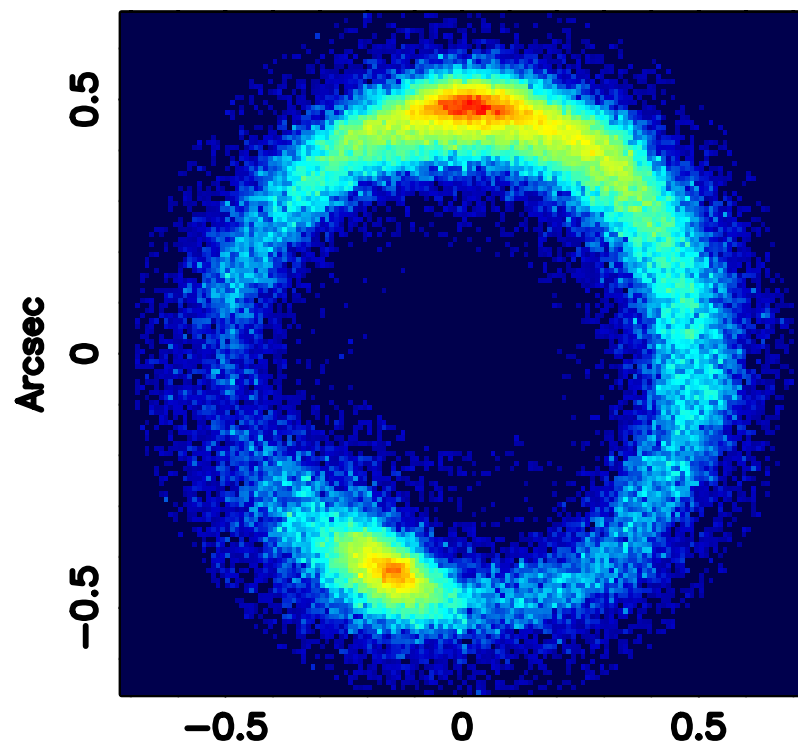# *Mutual information between latents and M(z)*



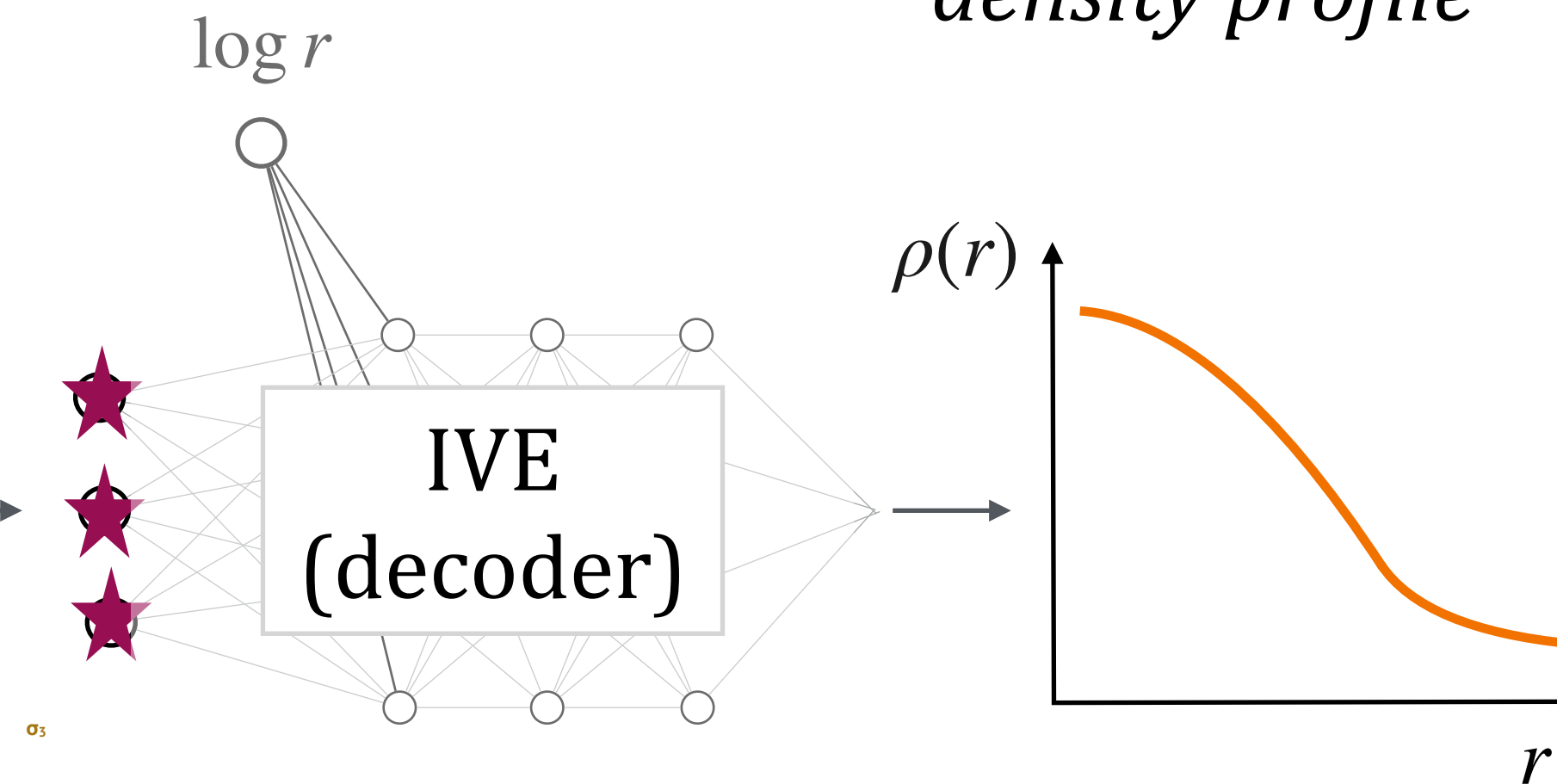'**Normalization**' *latent sensitive to formation history* **before infalling** *into the main host halo*
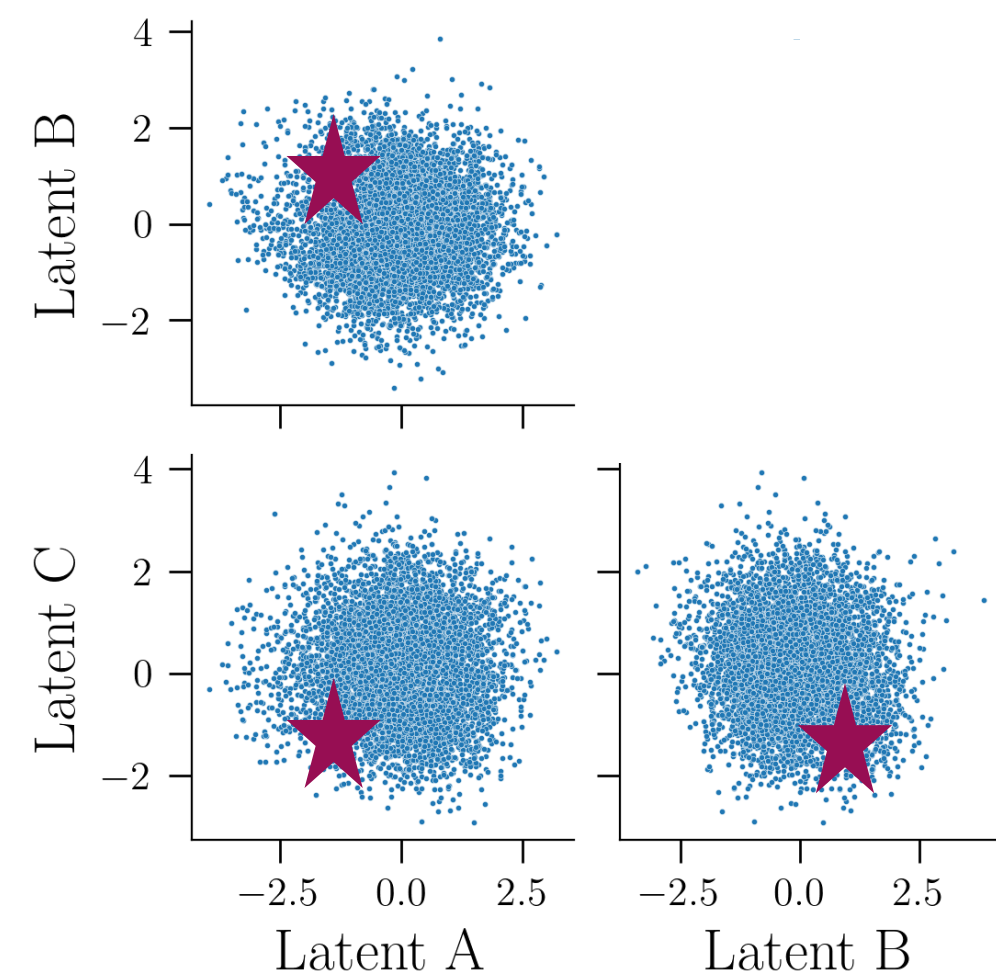
'**Truncation**' *latent sensitive to formation history* **after infalling** *into the main host halo*

'**Inner shape**' *latent sensitive to half-mass formation time*

Lucie-Smith, Despali, Springel (MNRAS, 2024)

A ... er ... *halo density profile* ... or ... *al lensing*

subhalo

*Sample from the*

*Infer subhalo density profile*

IVE (decoder)

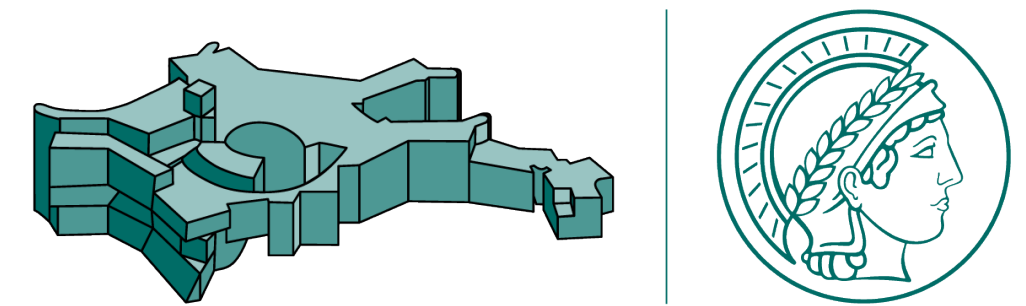$\rho(r)$

$r$

*Next step*: Integrate IVE model within strong gravitational lensing pipeline
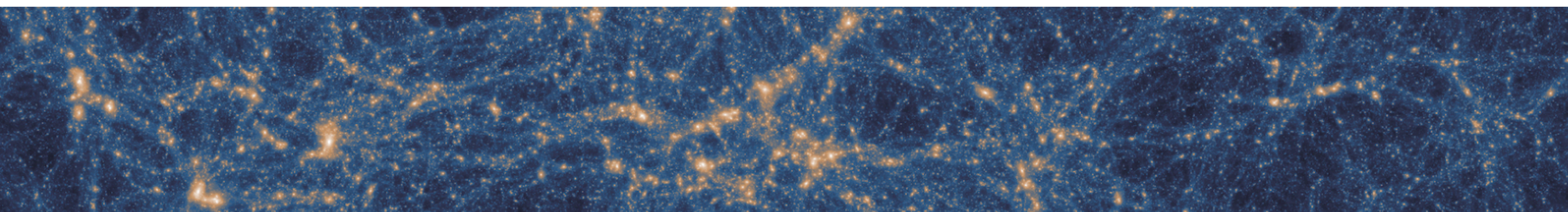
# *Conclusions*

- Interpretable variational encoders (IVE) provide new avenue to provide robust, physically interpretable models

- IVE disentangles different physical effects in minimal set of ingredients

- Explainable AI shows promise in enabling new, data-driven scientific discoveries

*Luisa Lucie-Smith, luisals@mpa-garching.mpg.de*

**MAX PLANCK INSTITUTE**
FOR ASTROPHYSICS

# IVE loss function

$\beta$ must be carefully fine-tuned
to balance accuracy with disentanglement

*Predictive term*      *KL-divergence term*

$$\mathcal{L} = \mathcal{L}_{\text{pred}}(\rho_{\text{true}}, \rho_{\text{pred}}) + \beta \, \mathcal{D}_{\text{KL}}(p(\mathbf{z}\,|\,\mathbf{x}); q(\mathbf{z}))$$

(Higgins+, 2017)

MSE/Gaussian likelihood:

$$\mathcal{L}_{\text{pred}} = \frac{1}{N}\sum_{i=1}^{N}\left[\log_{10}\rho_{i,\text{true}} - \log_{10}\rho_{i,\text{pred}}\right]^2$$

*How close are the predictions
to the ground truths*

Learnt latent distribution:      Prior:

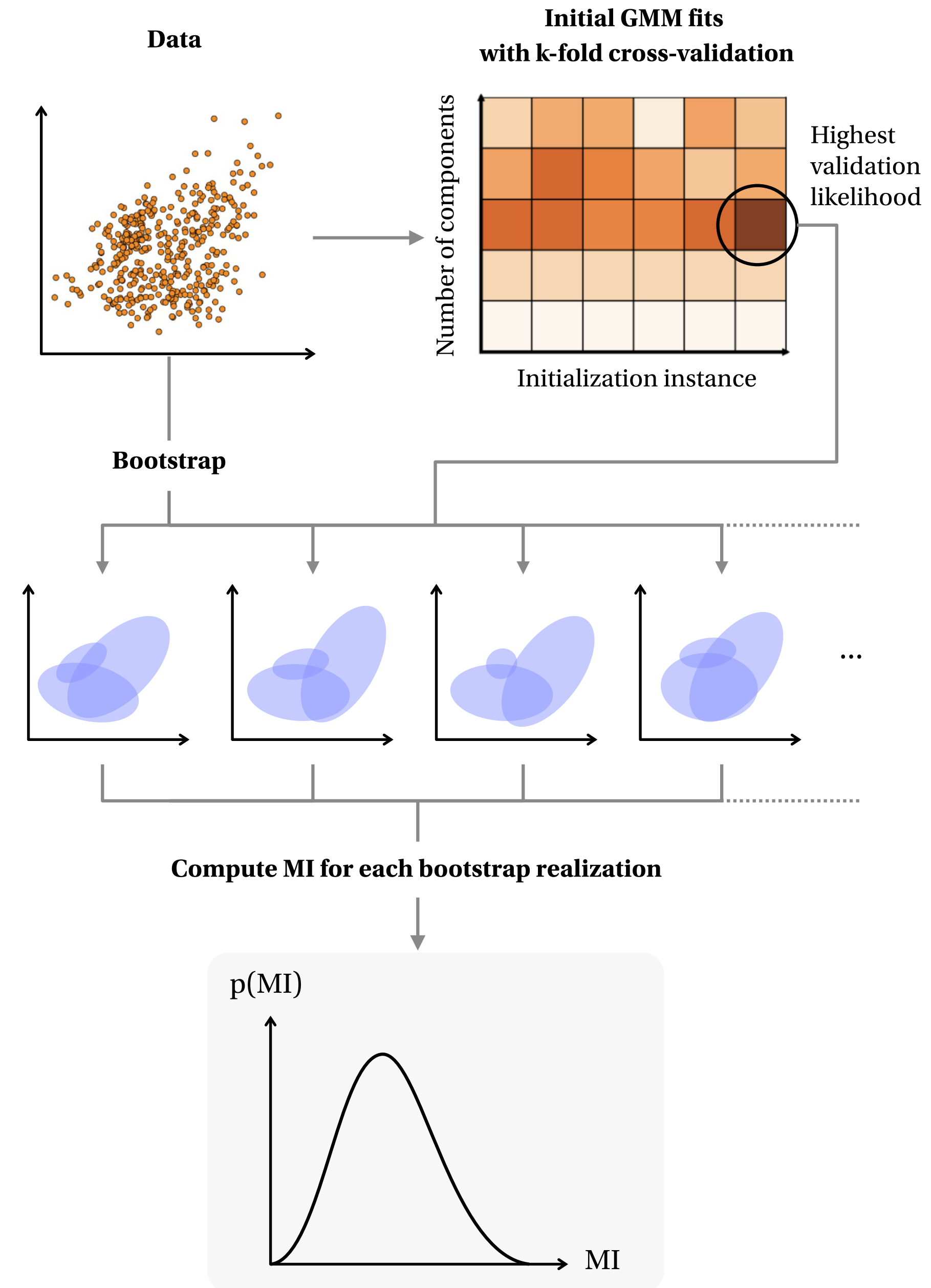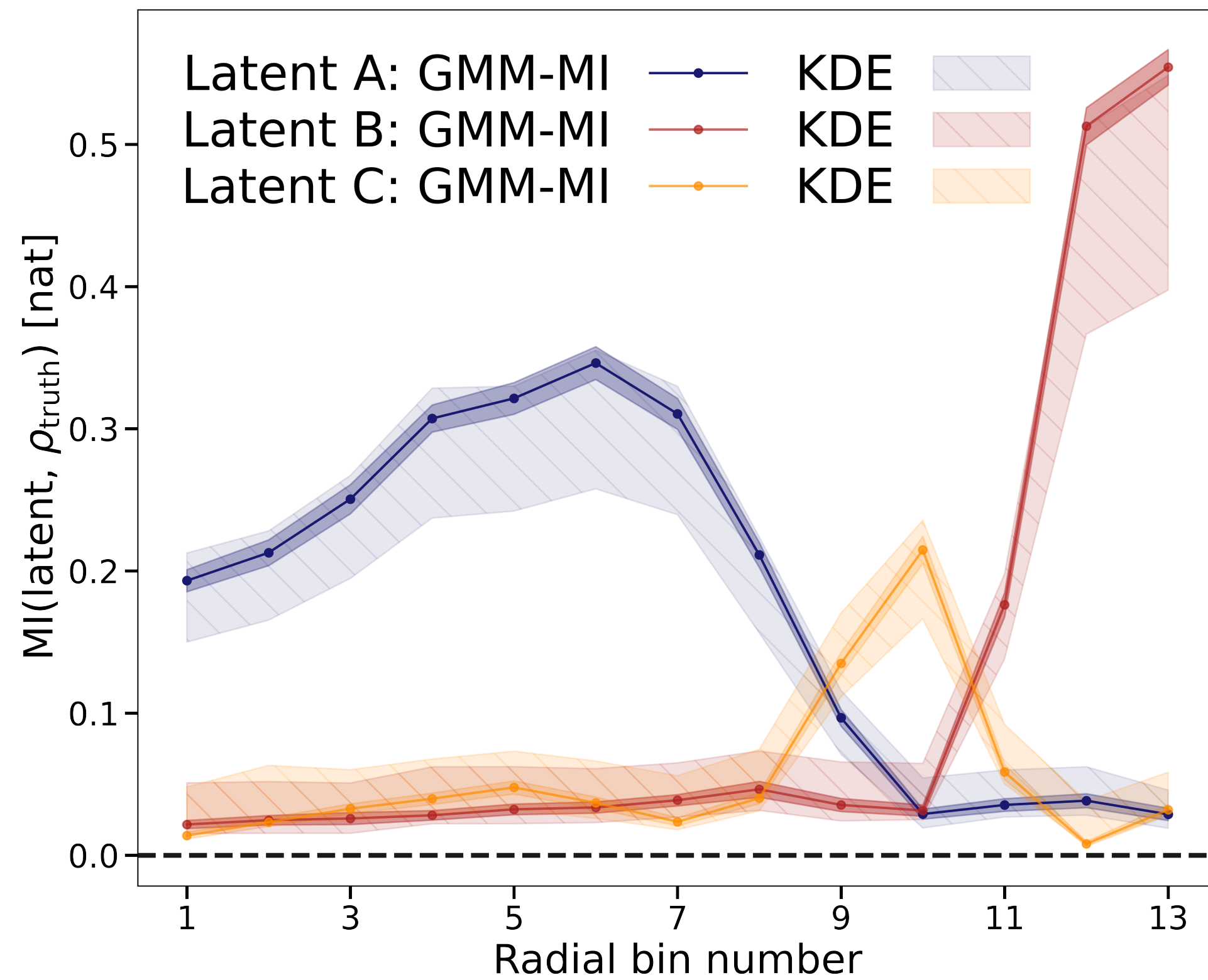$$p(\mathbf{z}\,|\,\mathbf{x}) = \prod_{i=1}^{L}\mathcal{N}(\mu_i(\mathbf{x}), \sigma_i(\mathbf{x}))$$      $$q(\mathbf{z}) = \prod_{i=1}^{L}\mathcal{N}(0,1)$$

*How close is the latent distribution to
set of independent unit Gaussians*

# *Mutual information*

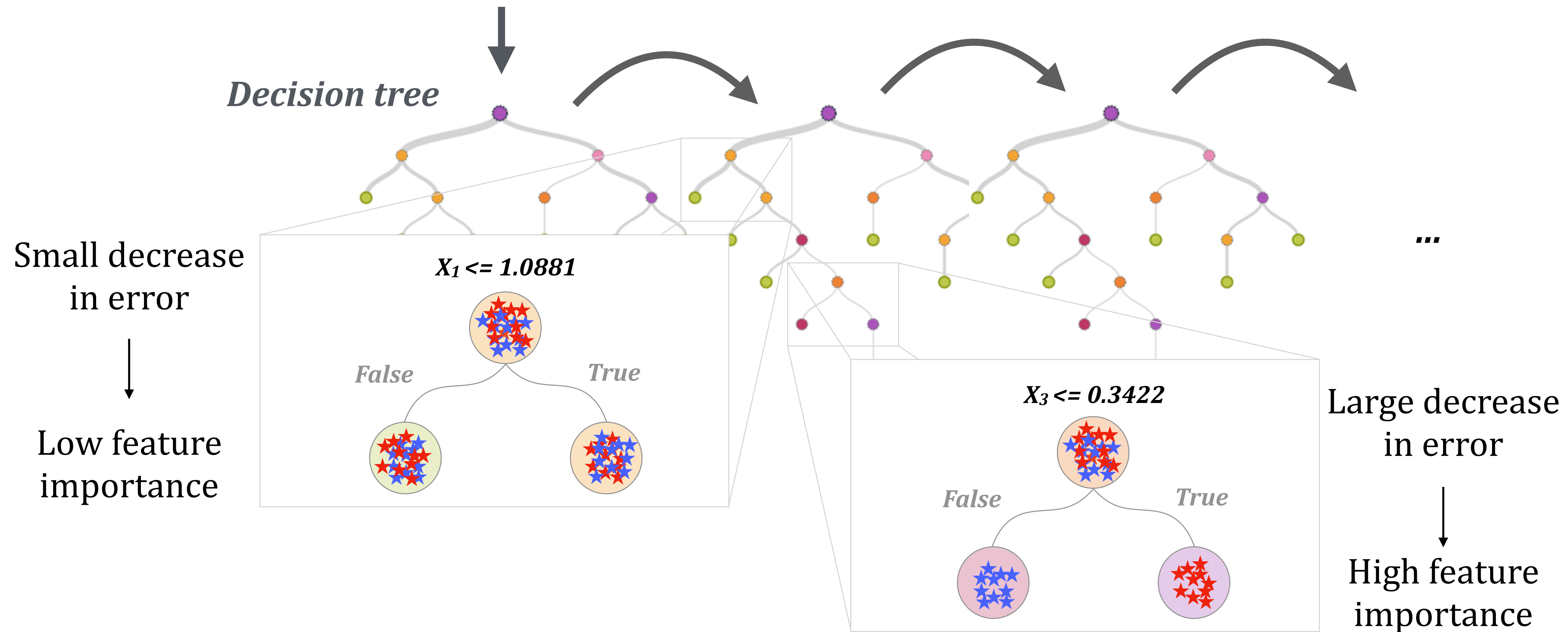$$\mathrm{MI}(X, Y) = \iint p(x, y) \, \log \left[ \frac{p(x, y)}{p(x)p(y)} \right] \mathrm{d}x \mathrm{d}y$$

# ML algorithm: gradient boosted trees (GBTs)

GBTs add new decisions trees to correct mistakes of previous trees



*Decision tree*

Small decrease
in error

$\downarrow$

Low feature
importance

$X_1 <= 1.0881$

*False*          *True*

$X_3 <= 0.3422$

*False*          *True*

Large decrease
in error

$\downarrow$

High feature
importance

...

**Feature importance** $\propto$ *decrease in error due to splits made by feature*

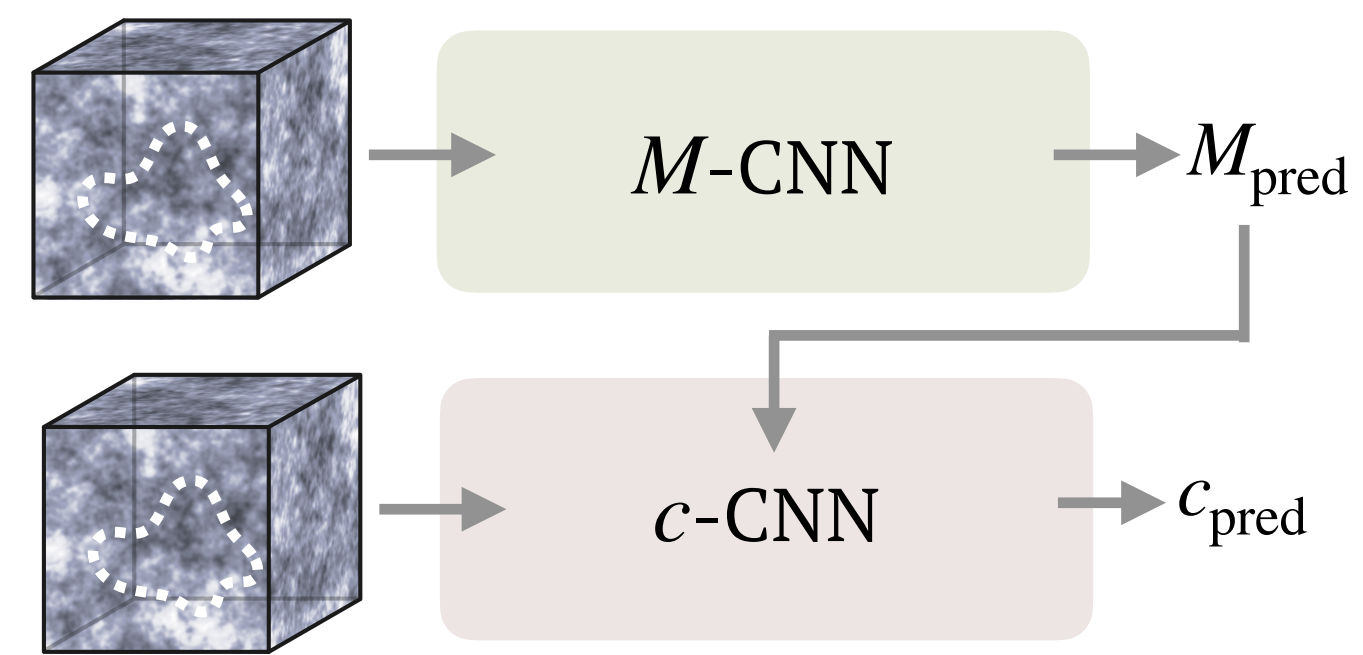Friedman, 2001; Friedman, 2002

# Importance of mass accretion history in predicting cluster mass profiles



Lucie-Smith, Adhikari, Wechsler (MNRAS, 2022)

# *Interpretability methods for deep learning*

- Saliency methods



*My view:* visualization methods give some insight but lack quantitative conclusions