# Interpretable machine learning?
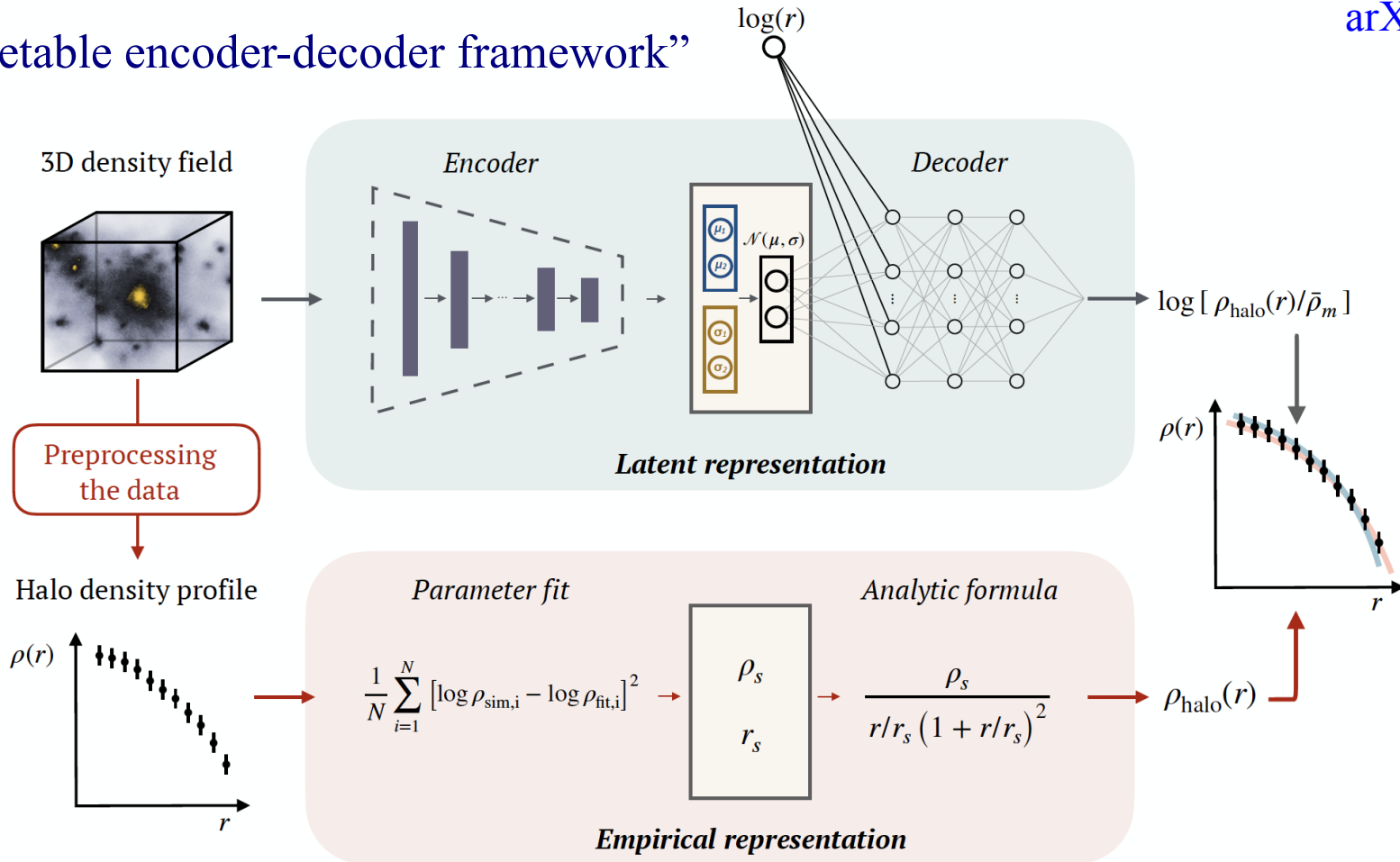
## Discovering the building blocks of dark matter halo density profiles with neural networks

Luisa Lucie-Smith,[1, *] Hiranya V. Peiris,[2, 3] Andrew Pontzen,[2] Brian Nord,[4, 5, 6] Jeyan Thiyagalingam,[7] and Davide Piras[2]

"An interpretable encoder-decoder framework"



$$\mathcal{L} = \mathcal{L}_{\text{pred}}(\boldsymbol{\rho}_{\text{true}}, \boldsymbol{\rho}_{\text{pred}}) + \beta \, \mathcal{D}_{\text{KL}}[p_\phi(\boldsymbol{z}|\boldsymbol{x}); q(\boldsymbol{z})],$$
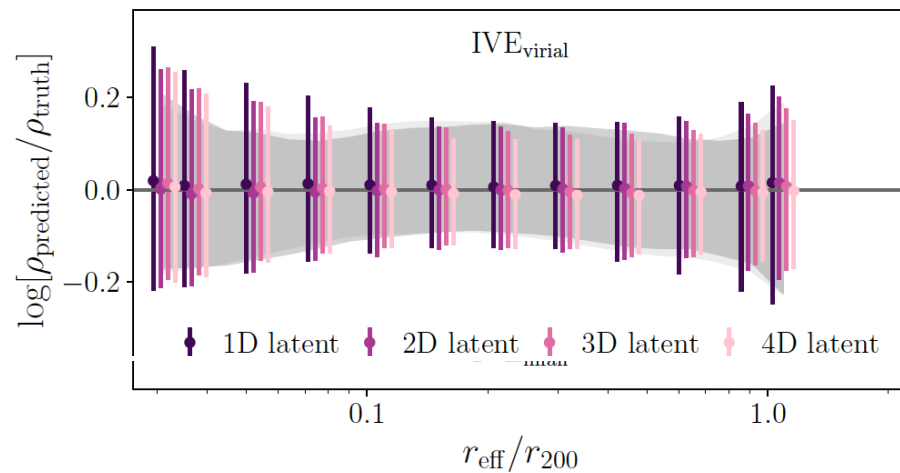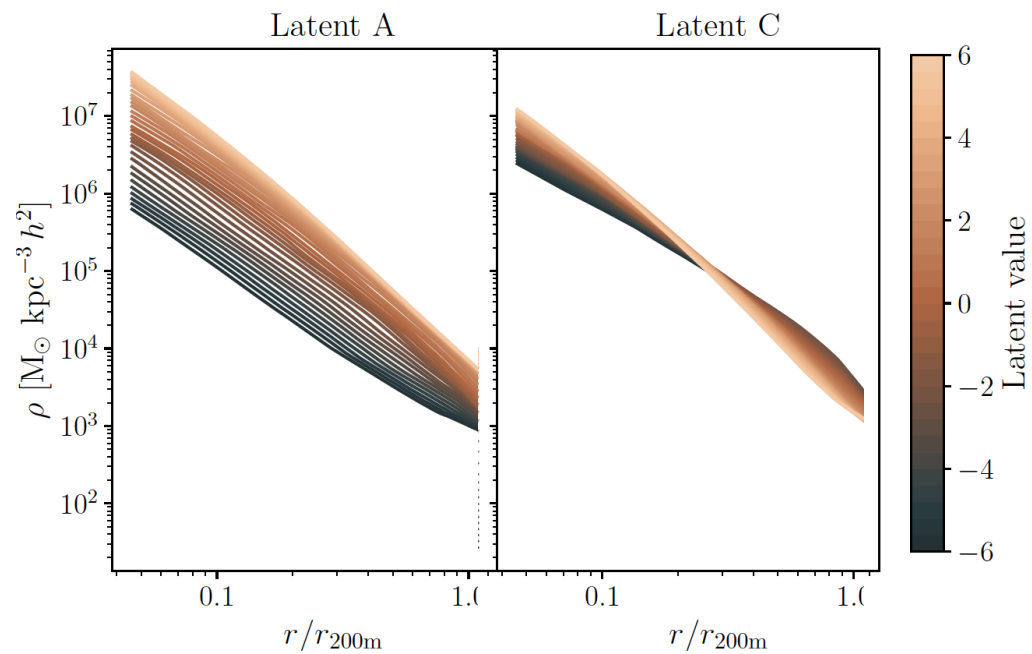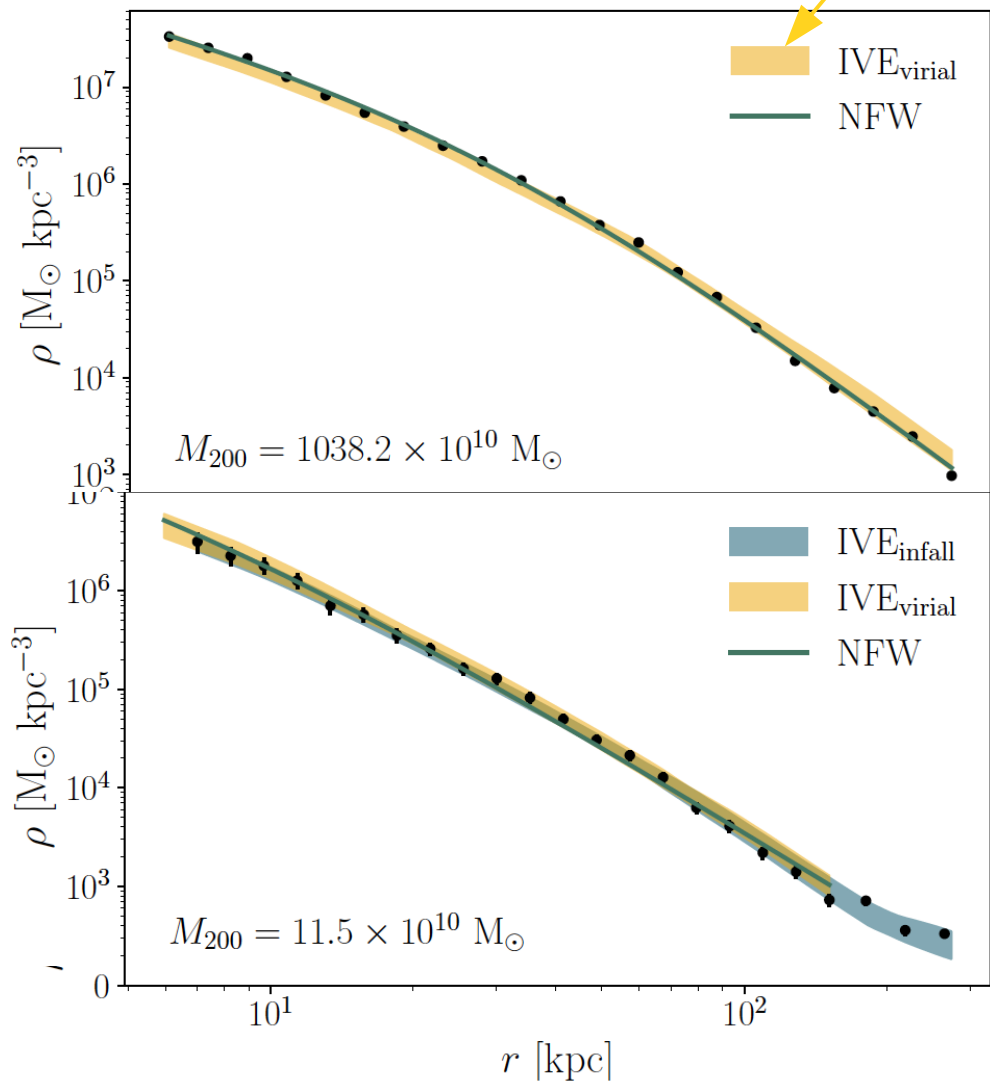
Goodness of fit

Disentanglement of latents

$$\mathcal{L}_{\text{pred}} = \frac{1}{N} \sum_{i=1}^{N} \left( \log_{10} \rho_{i,\text{true}} - \log_{10} \rho_{i,\text{pred}} \right)^2, \qquad \mathcal{D}_{\text{KL}}(\mathcal{N}(\mu_{\boldsymbol{z}}, \sigma_{\boldsymbol{z}}); \mathcal{N}(0, 1)) = -\frac{1}{2} \sum_{i=1}^{L} \left[ 1 + 2 \log \sigma_i - \mu_i^2 - \sigma_i^2 \right].$$

# NFW96 insights into z=0 halo density profiles

I. Spherically averaged halo density profiles within R200 can be fit over the resolved radial range to within the noise due to substructure and counting statistics by a smooth function of just two variables.

II. Profiles are homologous: they can be fit by a "universal" curve, with the two parameters corresponding to a characteristic radius and a characteristic density, hence to offsets of the universal curve parallel to the *x*- and *y*-axes in a log-log plot.

III. The characteristic densities and radii are correlated: bigger halos are less dense.

DL,  r < R200

$M_{200} = 1038.2 \times 10^{10}$ M$_\odot$

IVE$_{\text{virial}}$
NFW

IVE$_{\text{infall}}$
IVE$_{\text{virial}}$
NFW

$M_{200} = 11.5 \times 10^{10}$ M$_\odot$

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

$r$ [kpc]

Latent A

Latent C

$\rho$ [M$_\odot$ kpc$^{-3}$ $h^2$]

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

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

Latent value

IVE$_{\text{virial}}$

$\log[\rho_{\text{predicted}}/\rho_{\text{truth}}]$

1D latent    2D latent    3D latent    4D latent

$r_{\text{eff}}/r_{200}$

# NFW96 insights into z=0 halo density profiles

I. Spherically averaged halo density profiles within R200 can be fit over the resolved radial range to within the noise due to substructure and counting statistics by a smooth function of just two variables. ✔

II. Profiles are homologous; they can be fit by a "universal" curve, with the two parameters corresponding to a characteristic radius and a characteristic density, hence to offsets of the universal curve parallel to the *x*- and *y*-axes in a log-log plot.

III. The characteristic densities and radii are correlated: bigger halos are less dense.

# NFW96 insights into z=0 halo density profiles

I. Spherically averaged halo density profiles within R200 can be fit over the resolved radial range to within the noise due to substructure and counting statistics by a smooth function of just two variables. ✔

II. Profiles are homologous; they can be fit by a "universal" curve, with the two parameters corresponding to a characteristic radius and a characteristic density, hence to offsets of the universal curve parallel to the *x*- and *y*-axes in a log-log plot. ✘

III. The characteristic densities and radii are correlated: bigger halos are less dense.

# NFW96 insights into z=0 halo density profiles

**I.** Spherically averaged halo density profiles within R200 can be fit over the resolved radial range to within the noise due to substructure and counting statistics by a smooth function of just two variables. ✓

**II.** Profiles are homologous; they can be fit by a "universal" curve, with the two parameters corresponding to a characteristic radius and a characteristic density, hence to offsets of the universal curve parallel to the *x*- and *y*-axes in a log-log plot. ✗

**III.** The characteristic densities and radii are correlated: bigger halos are less dense. ✗

"Failure" is, in part, a consequence of the disentanglement requirement