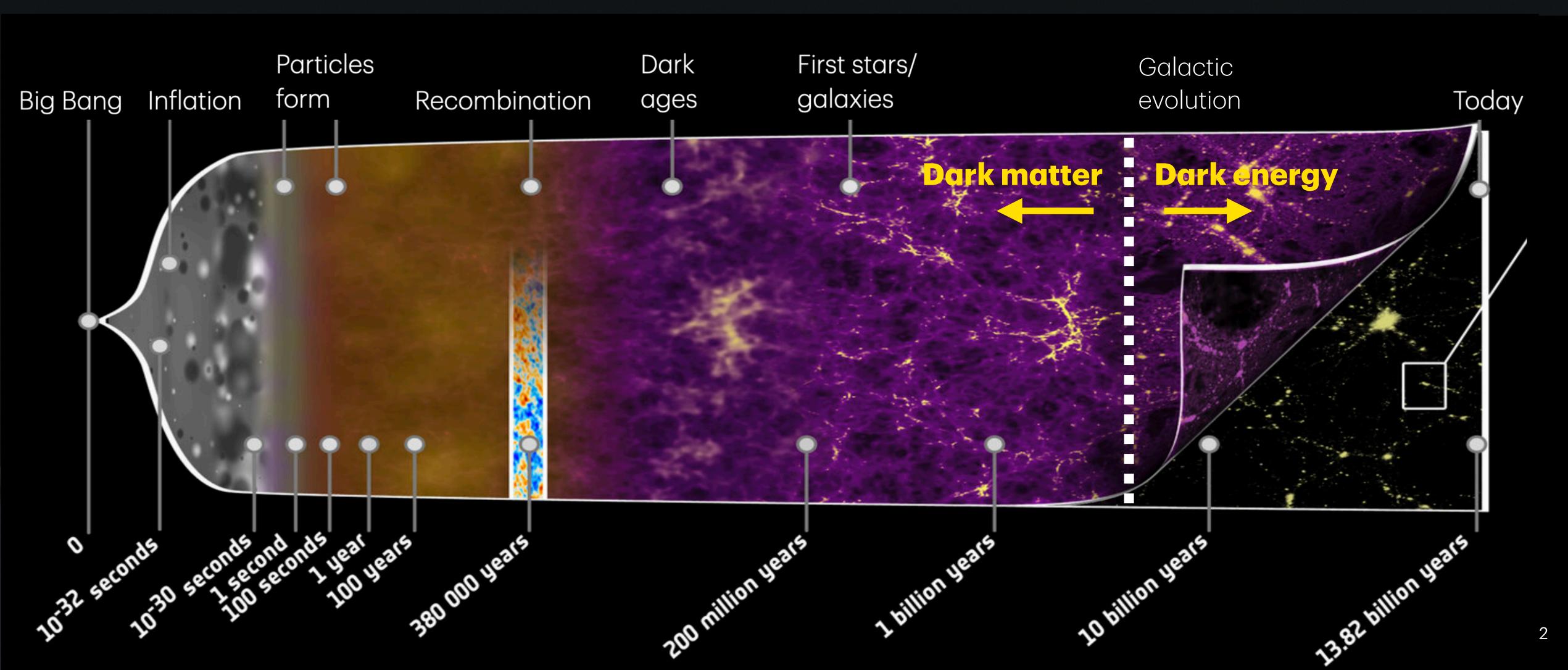


Shivam Pandey (Johns Hopkins University)

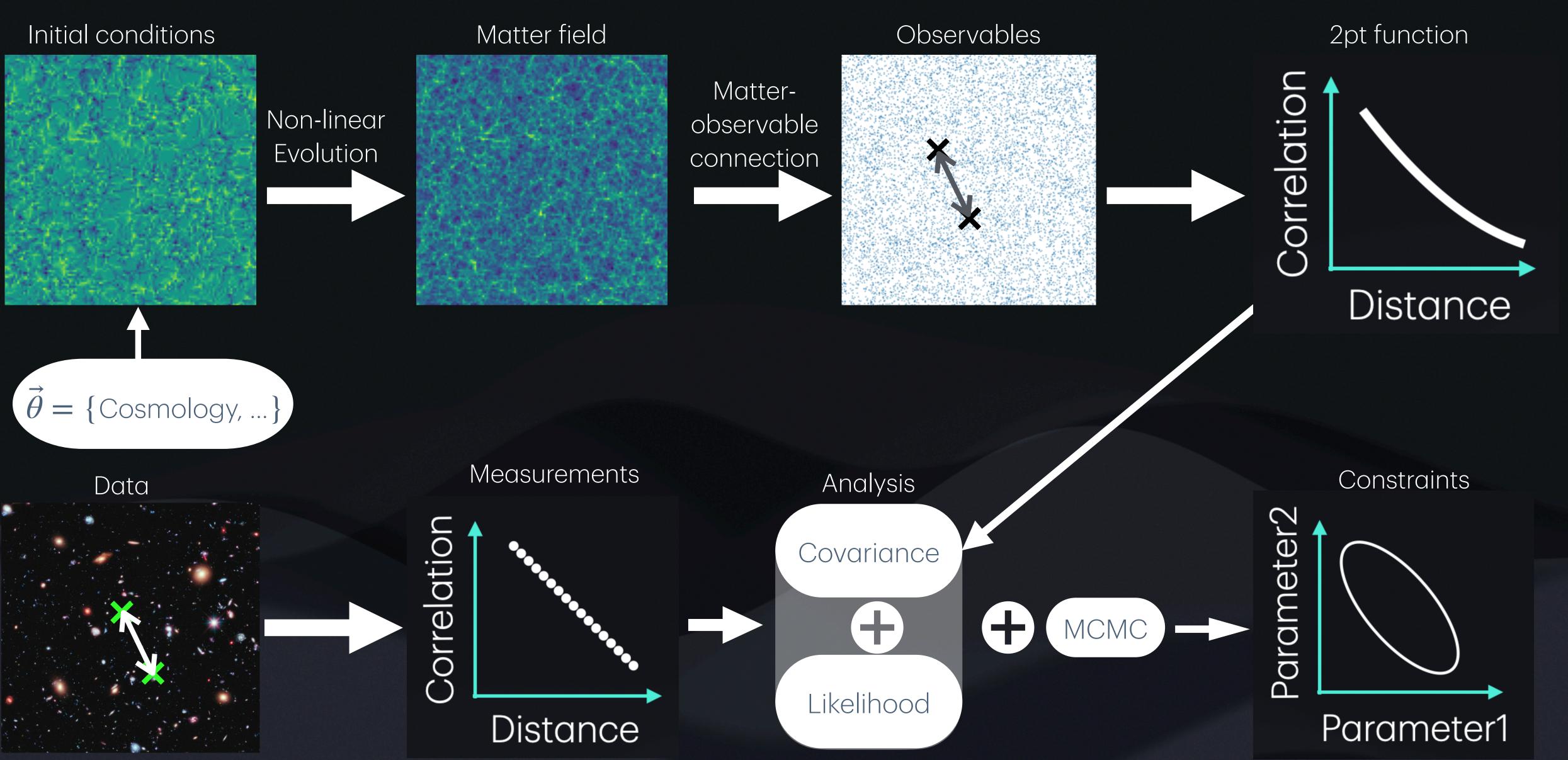
works with Ben Wandelt, Chirag Modi, Francois Lanusse, Chris Lovell, Learning the Universe Collaboration

What we want to understand...

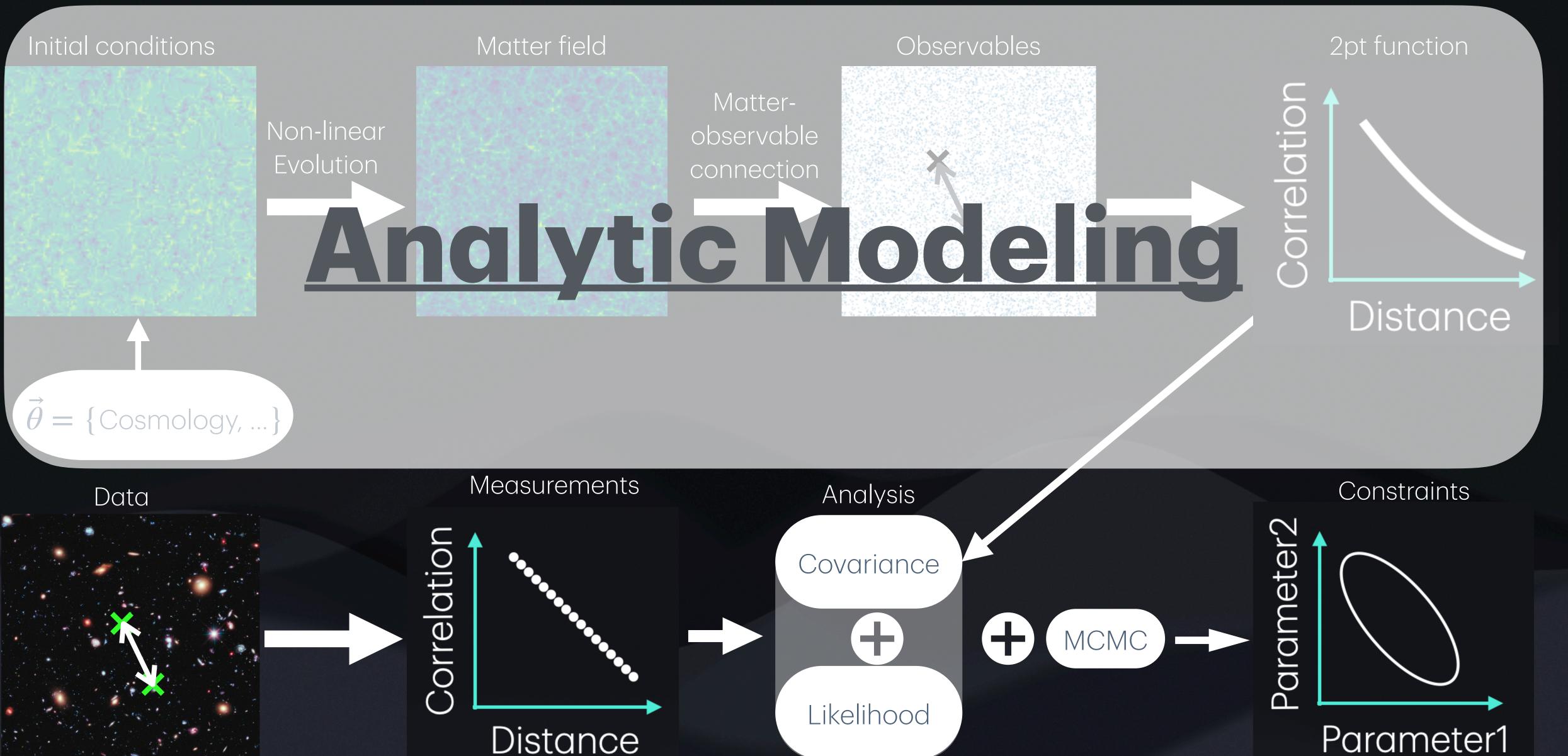
Evolution of the Large Scale Structure



Traditional Inference pipeline

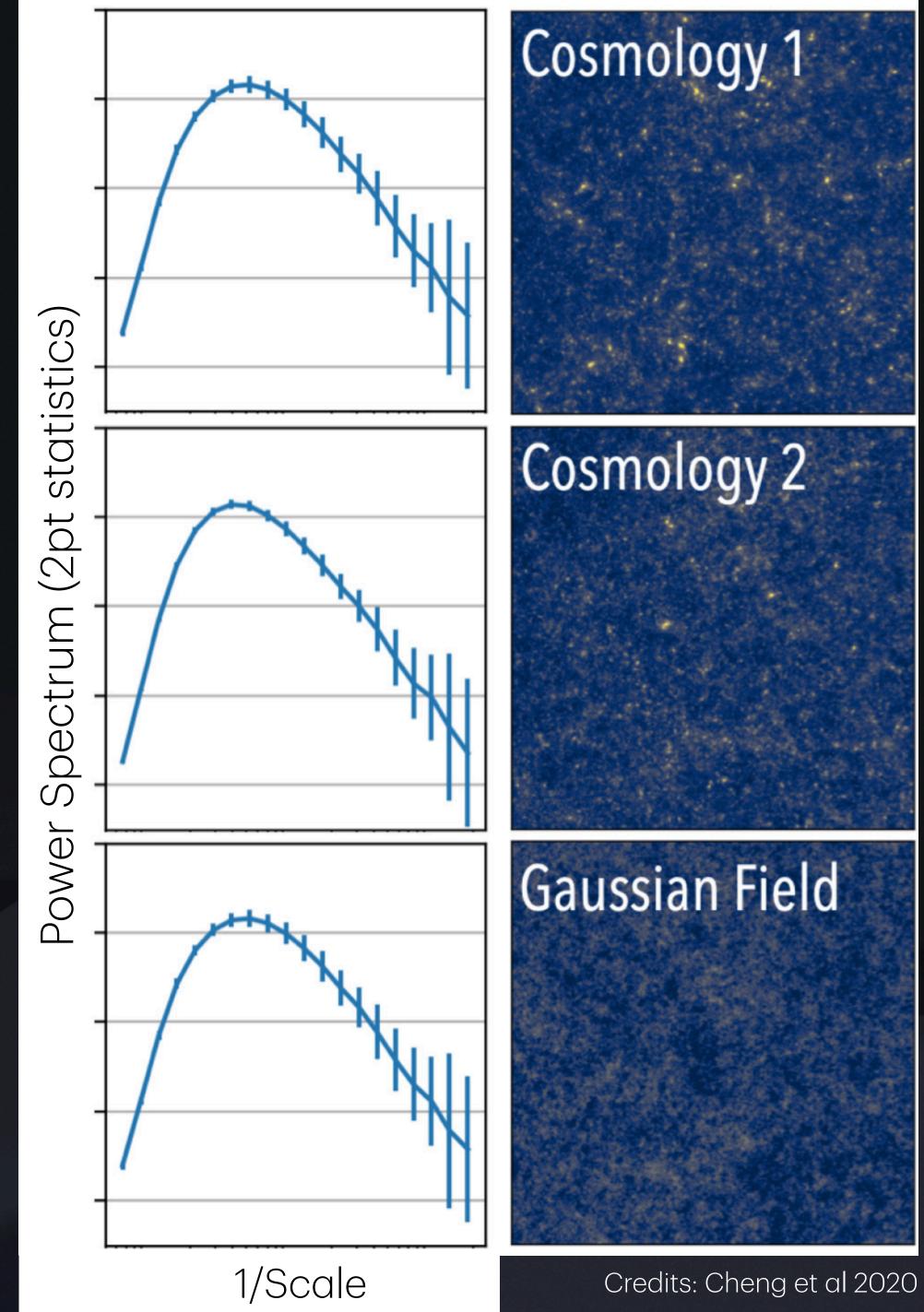


Traditional Inference pipeline



There is much more information in the observation!

2pt correlation alone does not capture all the information in a non-gaussian field.



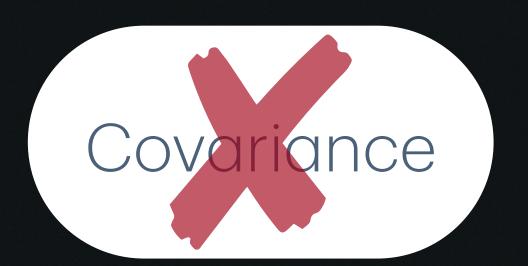
How do we do better?

- Go beyond 2pt statistics!
- Multiple options:
 - Marked correlation function
 - -Wavelets transforms
 - -Graph neural networks
 - -Convolutional neural networks
 - -.... (many more)

How do we do better?

- Go beyond 2pt statistics!
- Multiple options:
 - Marked correlation function
 - -Wavelets transforms
 - -Graph neural networks
 - -Convolutional neural networks
 - -.... (many more)



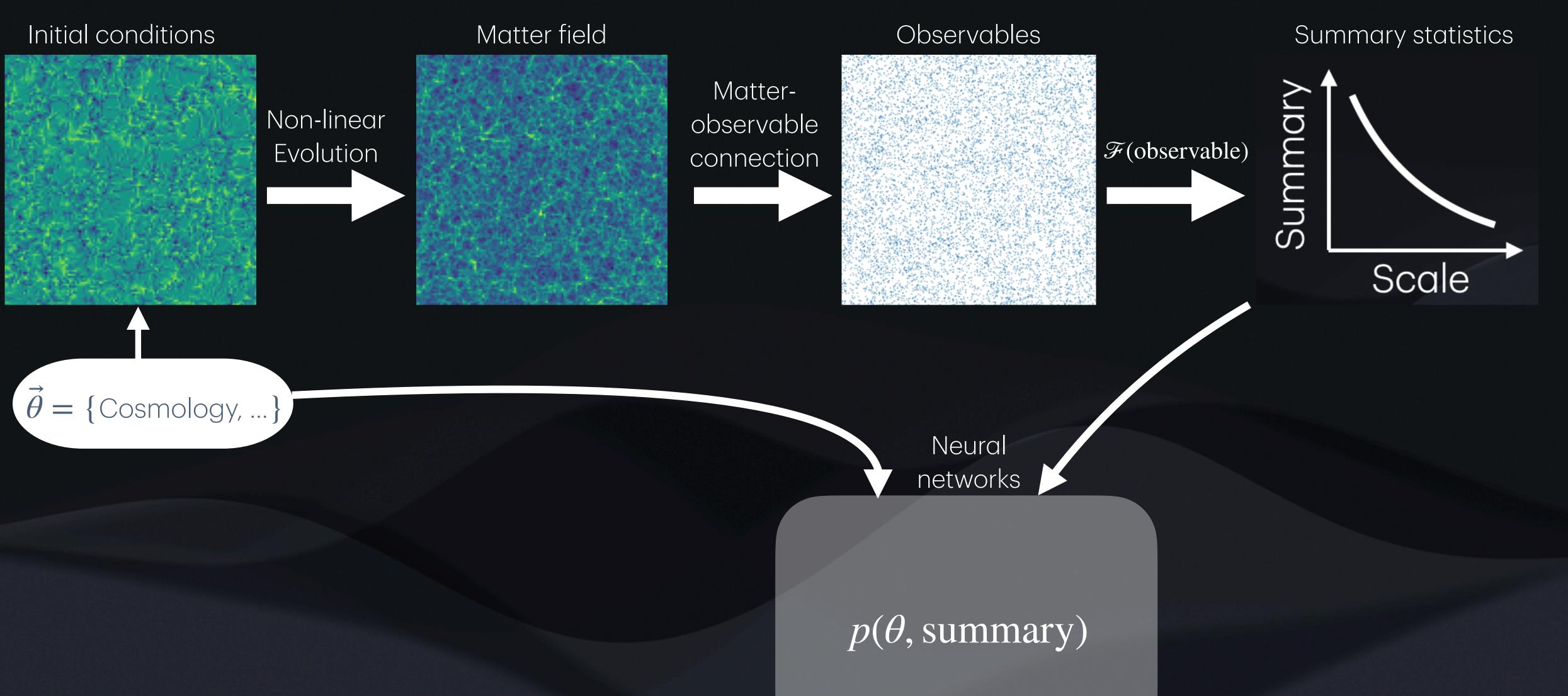


(Too non-linear)
(Need too many simulations)

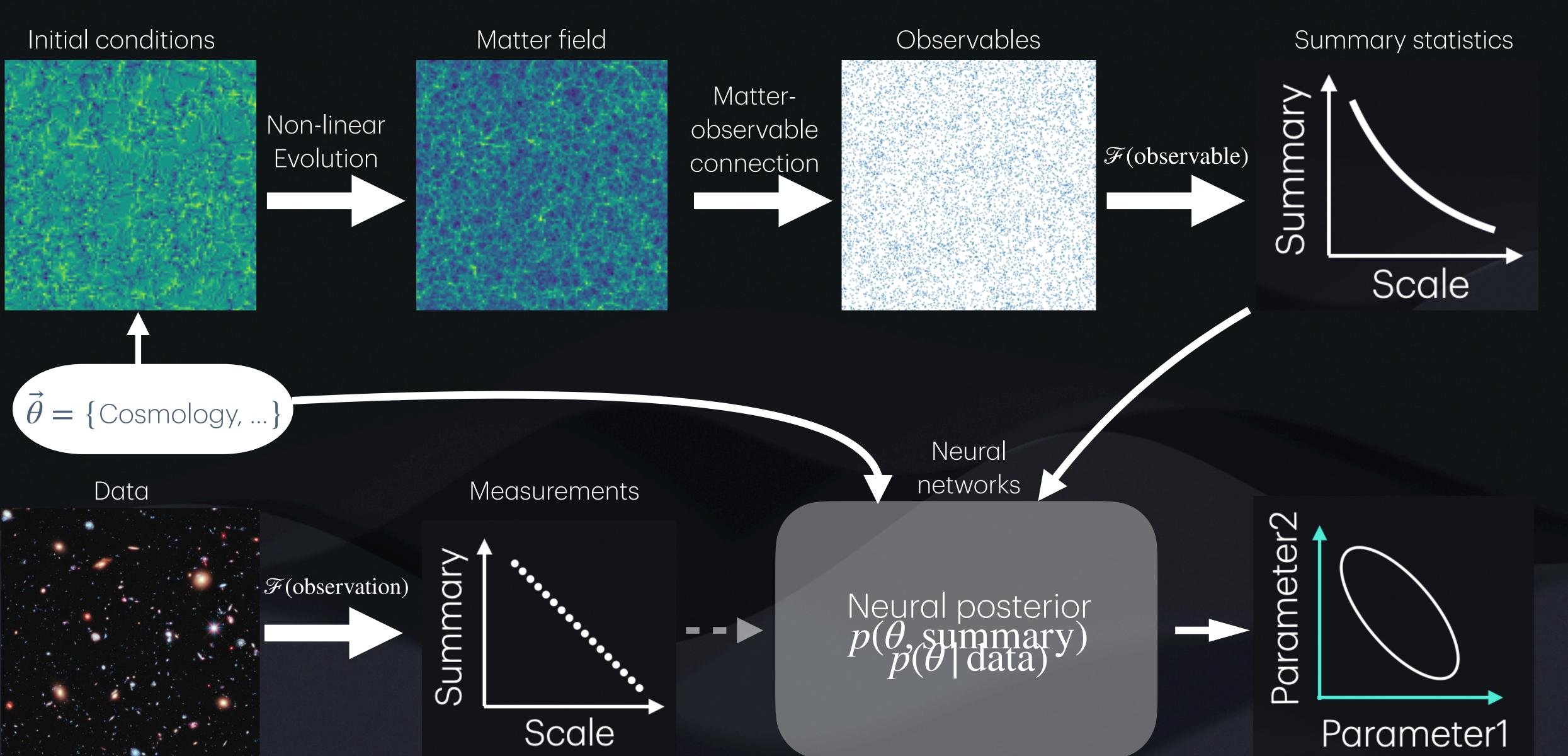


(Away from central limit theorem)

Simulation Based Inference (SBI) pipeline

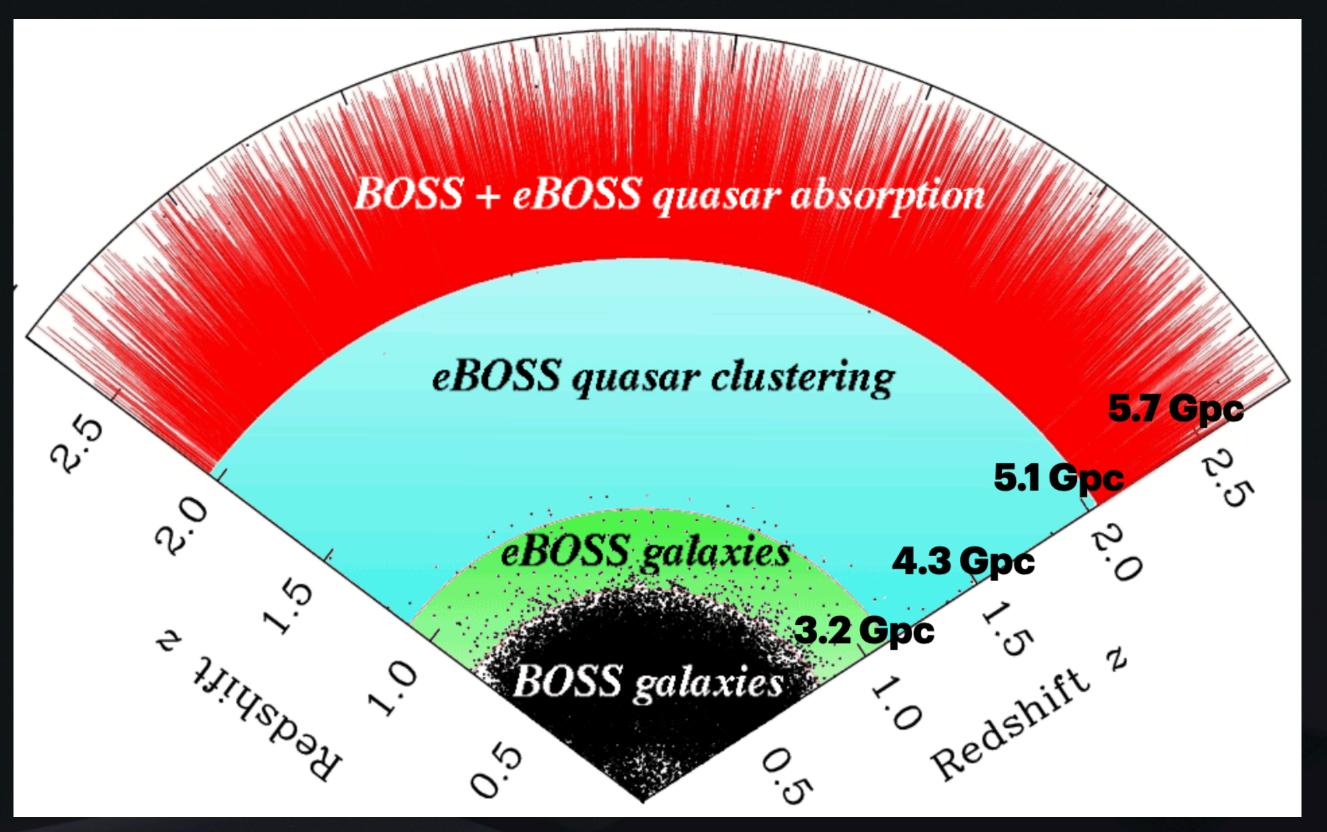


Simulation Based Inference (SBI) pipeline

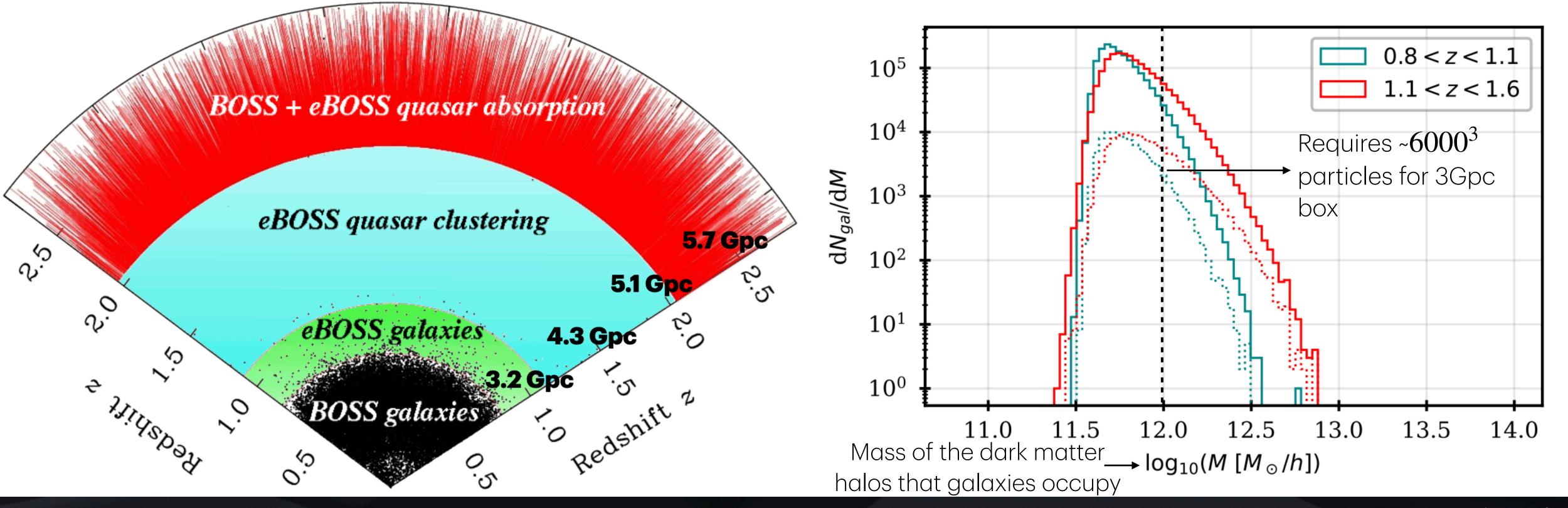


What are the main bottlenecks in this pipeline?

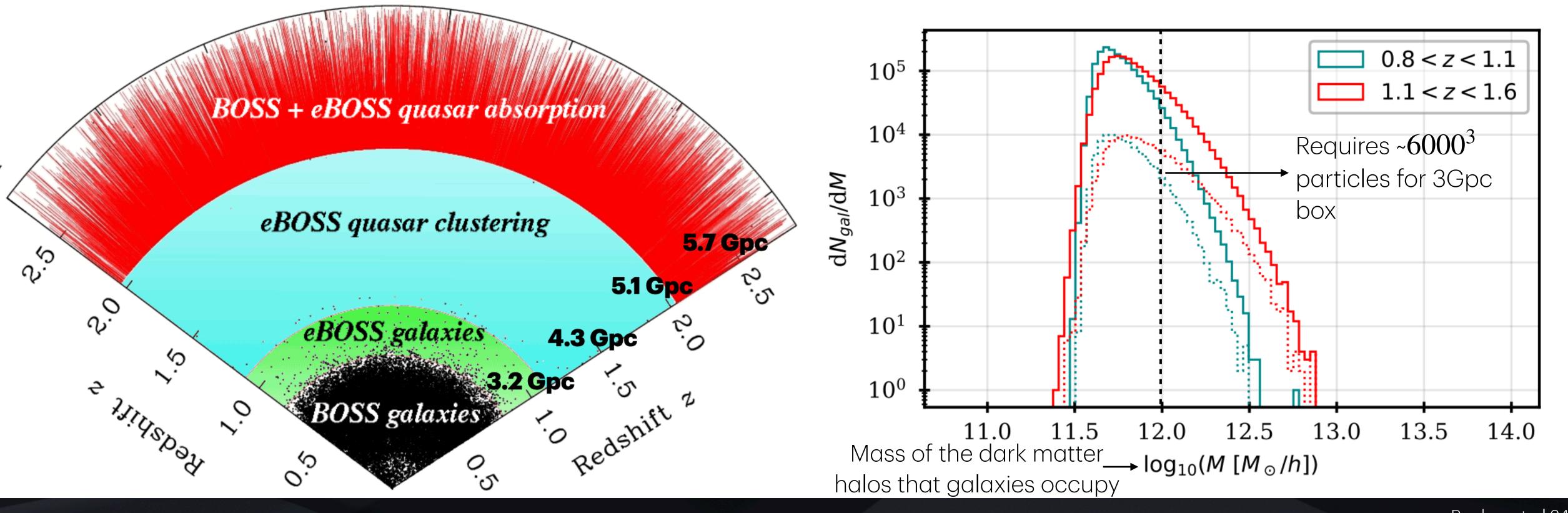
Simulation volume



Simulation volume and resolution



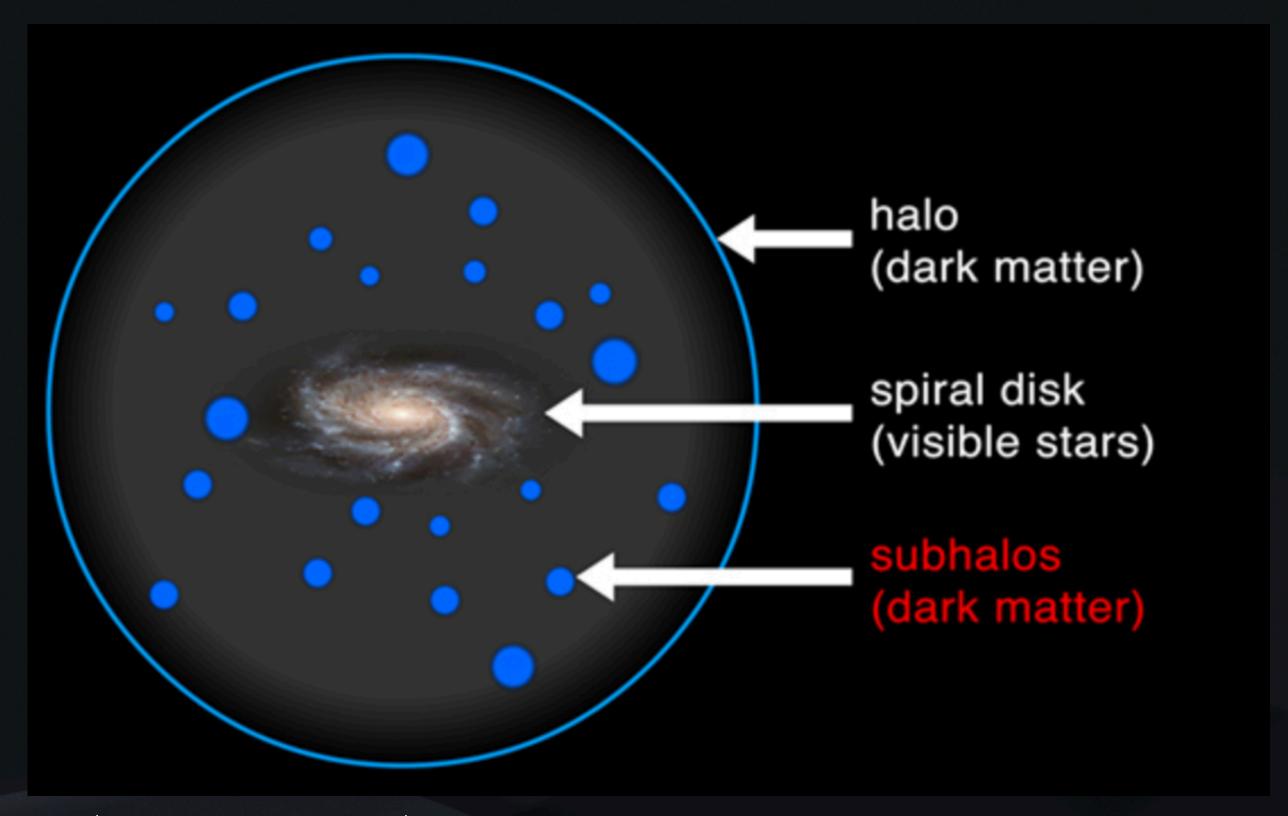
Simulation volume and resolution



Rocher et al 24

Typically <u>need to run a lot of them</u> to compare against the observations using techniques like SBI as mentioned before

Let's start with accelerating halo finding in the most extensive simulation suite we have

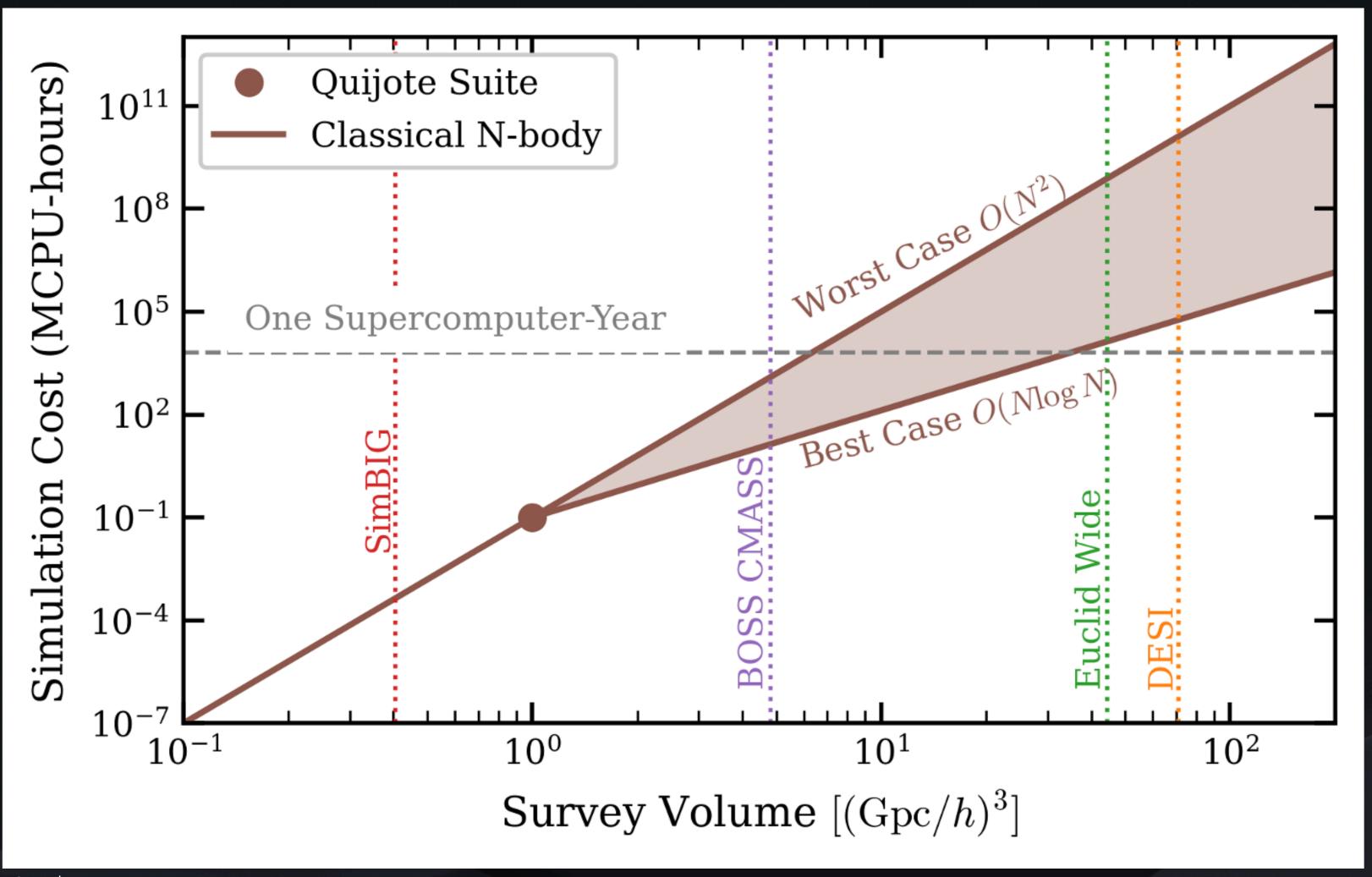


Credit: Luna Zagorac + Astrobites

Quijote N-body



- •1 $(Gpc/h)^3$ volume, 1024^3 particles
- ~2000 sims varying cosmologies
- ~5000 CPU hours/sim
- $10^{13}\,M_{\odot}$ halo mass resolution



Credit: Matt Ho + LtU-ILI team

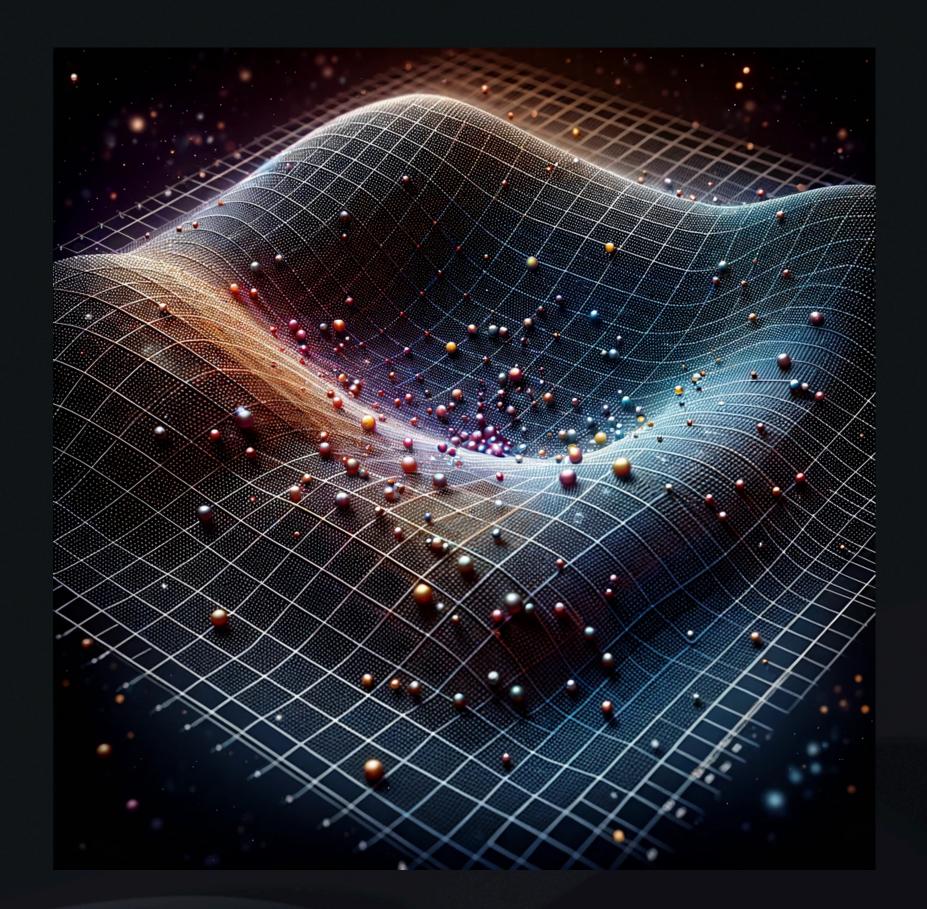
Naively scaling Quijote simulations to the needed volume will not work!

Quijote N-body



- •1 $(Gpc/h)^3$ volume, 1024^3 particles
- ~2000 sims varying cosmologies
- ~5000 CPU hours/sim
- $10^{13} M_{\odot}$ halo mass resolution
- >270 million CPU hours needed for full SDSS-like volume LH set

Particle mesh



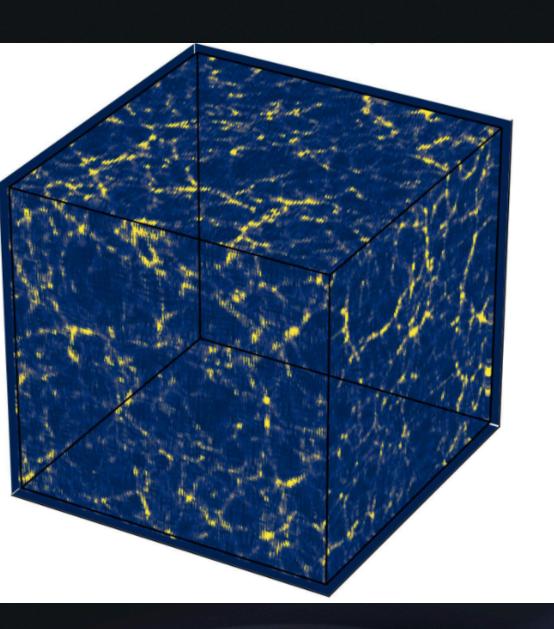
- $1(Gpc/h)^3$ volume
- 384³ particles
- ~5 CPU hours/sim or 10 GPU sec/sim
- IC —> Matter can be differentiable and much faster with GPUs

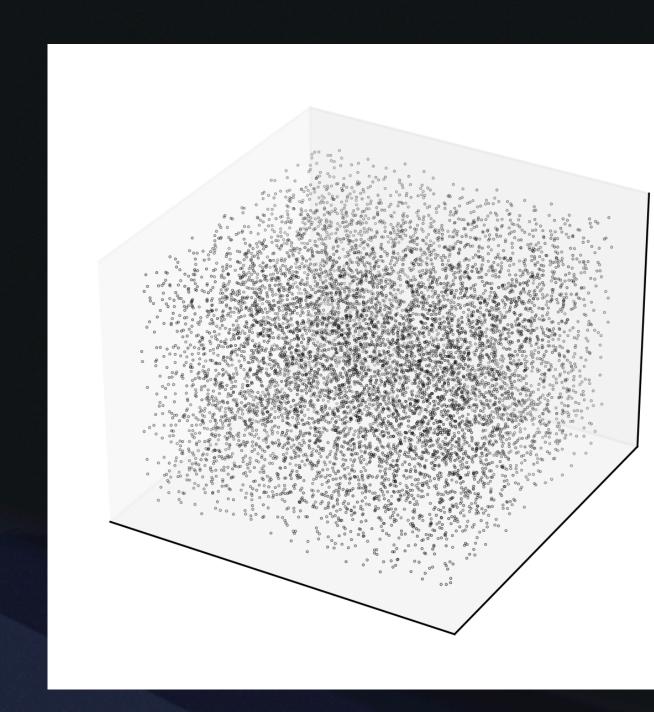
Quijote N-body



- $1 (Gpc/h)^3$ volume
- 1024³ particles
- ~5000 CPU hours/sim
- Not-differentiable
- >270 million CPU hours needed for full SDSS-like volume LH set

Connect PM to N-body halos

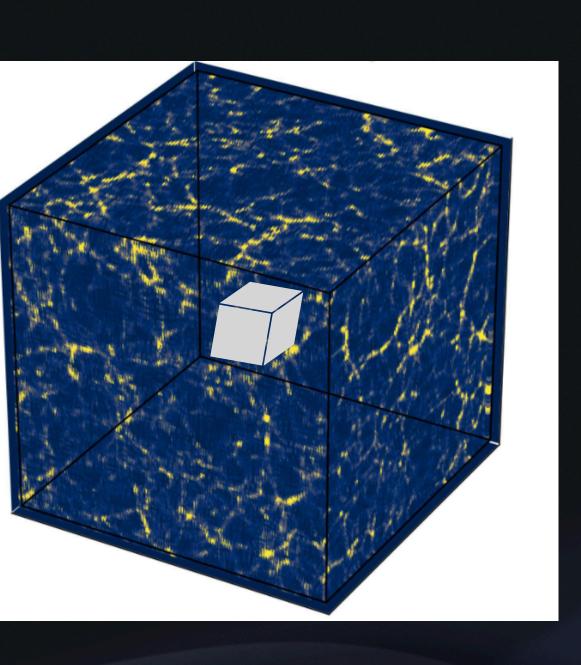




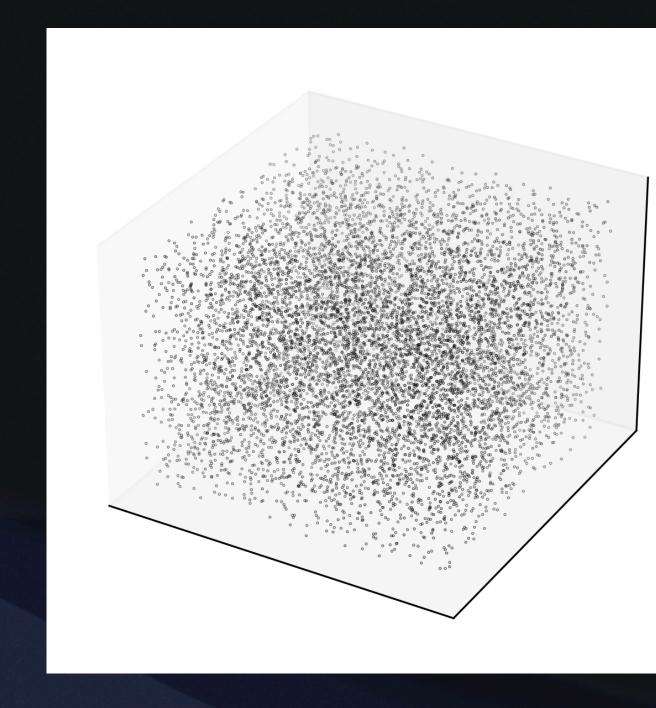
3D
Matter density field
(PM Simulation)
Input

3D
Halo distribution
(N-body Simulation)
Target

Voxelize the 3D volume (~8Mpc voxels)



 $p(\{M_1, M_2, \ldots, M_{N_j}\} \mid \vec{\delta}_j)$

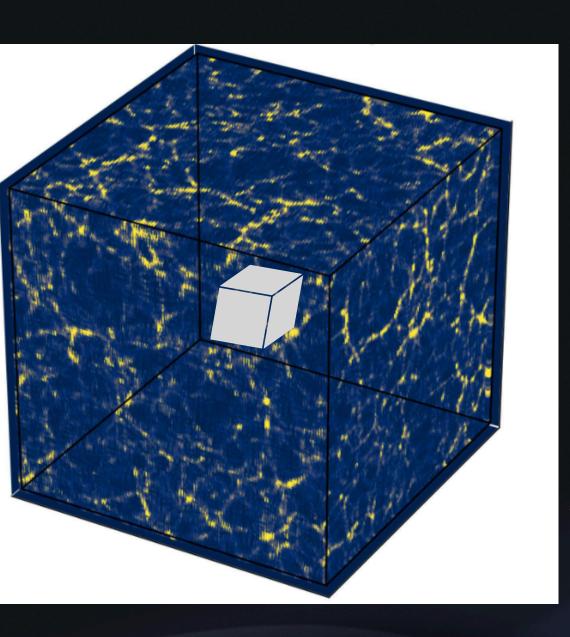


3D
Matter density field
(PM Simulation)
Input

3D Halo distribution (N-body Simulation) Target

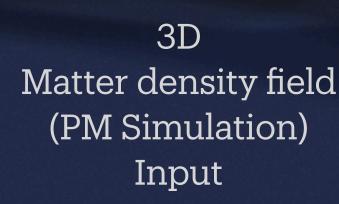
Creating Halos with Auto-Regressive Multi-stage networks

(CHARM, arxiv:2409.09124)



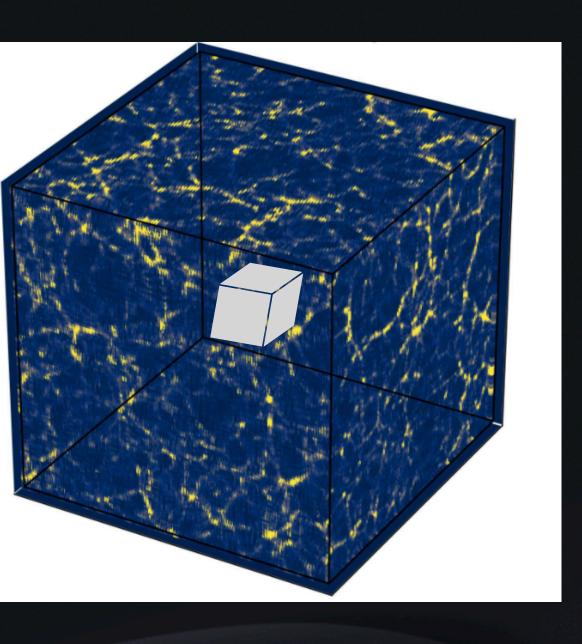
$$p(\{M_{1}, M_{2}, \dots, M_{N_{j}}\} | \vec{\delta}_{j}) =$$

$$p(N_{j} | \vec{\delta}_{j}) \times p(M_{1} | \vec{\delta}_{j}, N_{j}) \times p(M_{2} | M_{1}, \vec{\delta}_{j}, N_{j}) \times p(M_{N_{j}} | \{M_{1}, \dots, N_{j-1}\}, \vec{\delta}_{j}, N_{j})$$

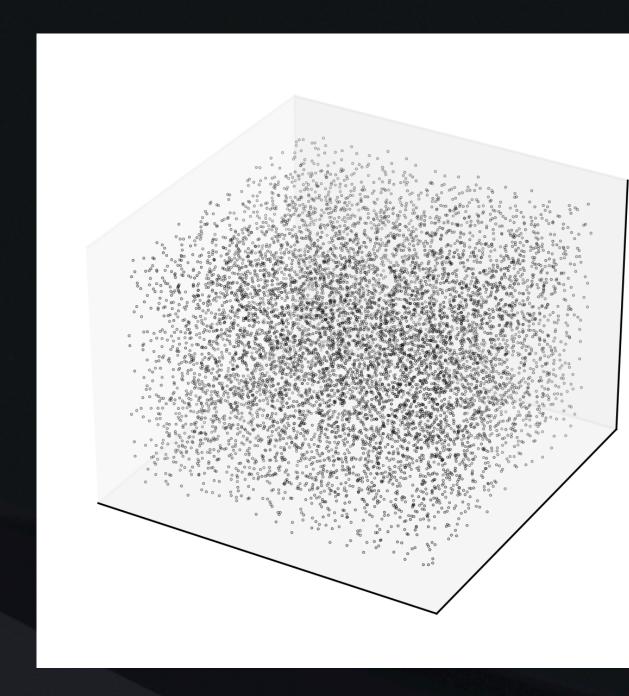


3D
Halo distribution
(N-body Simulation)
Target

CHARM OVEIVIEW arxiv:2409.09124



3D
Matter density field
(PM Simulation)
Input

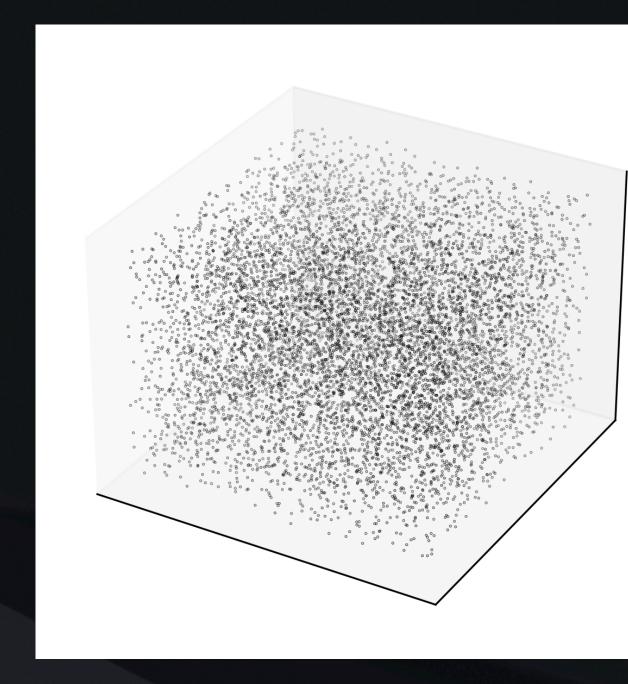


3D
Halo distribution
(N-body Simulation)
Target

CHARM OVEIVIEW arxiv:2409.09124

ResNet x2
Learned Features extraction

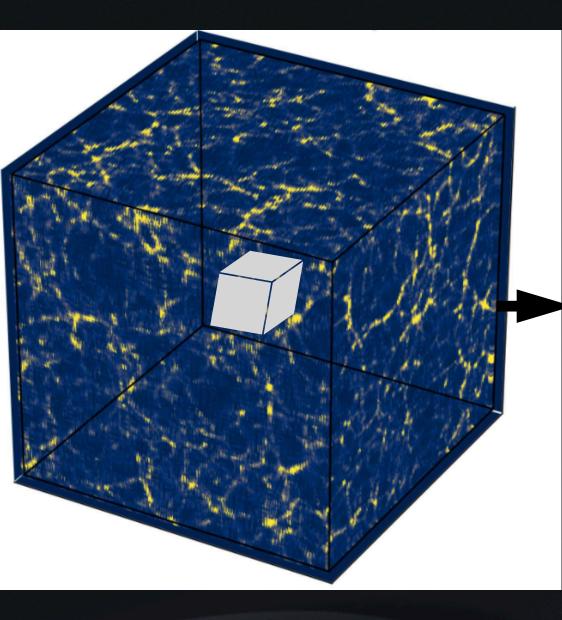
3D
Matter density field
(PM Simulation)
Input



3D
Halo distribution
(N-body Simulation)
Target

CHARM Overview

arxiv:2409.09124

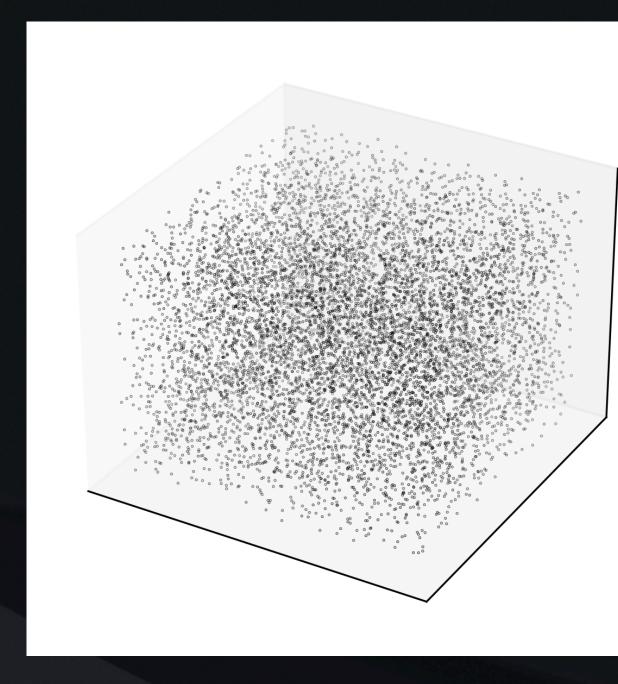


ResNet x2

Learned Features extraction

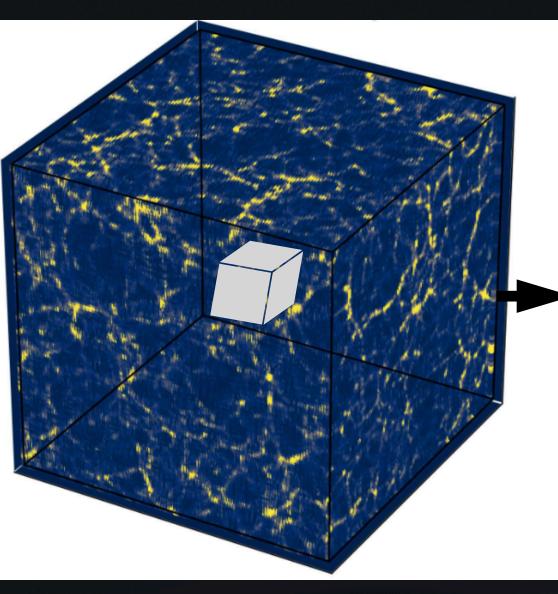
3D
Matter density field
(PM Simulation)
Input

Cosmological Parameters

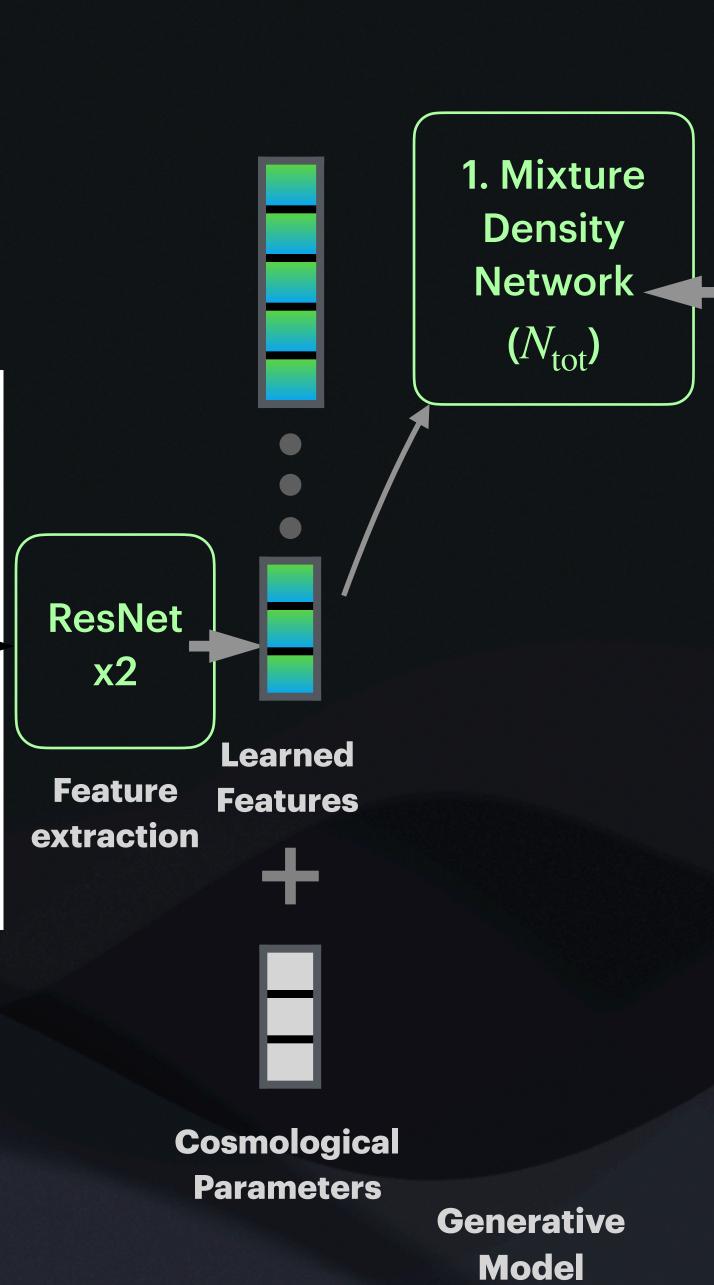


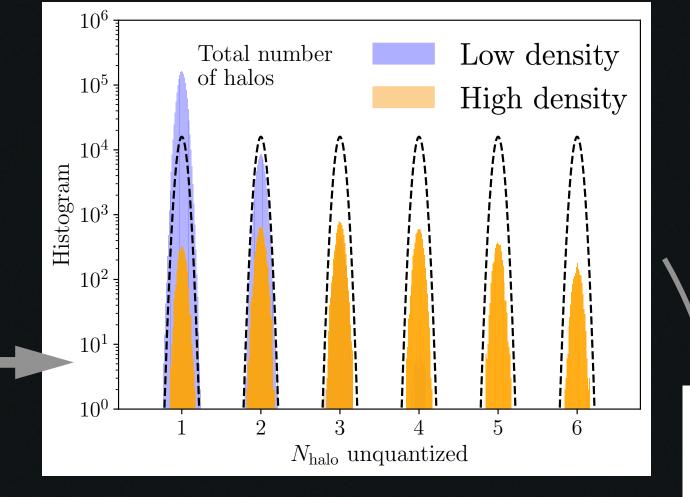
3D
Halo distribution
(N-body Simulation)
Target

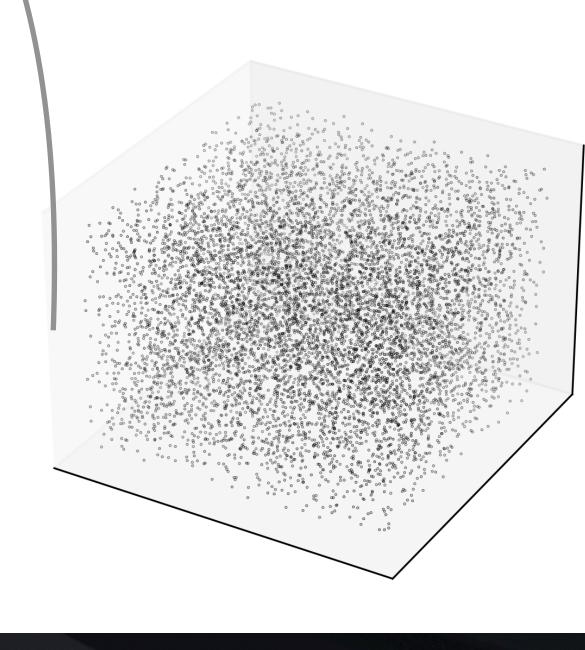
CHARM OVEIVIEW arxiv:2409.09124



3D
Matter density field
(PM Simulation)
Input

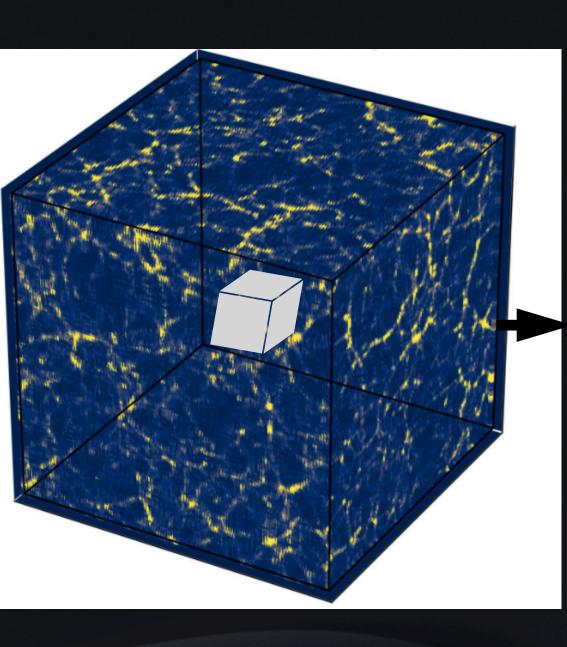




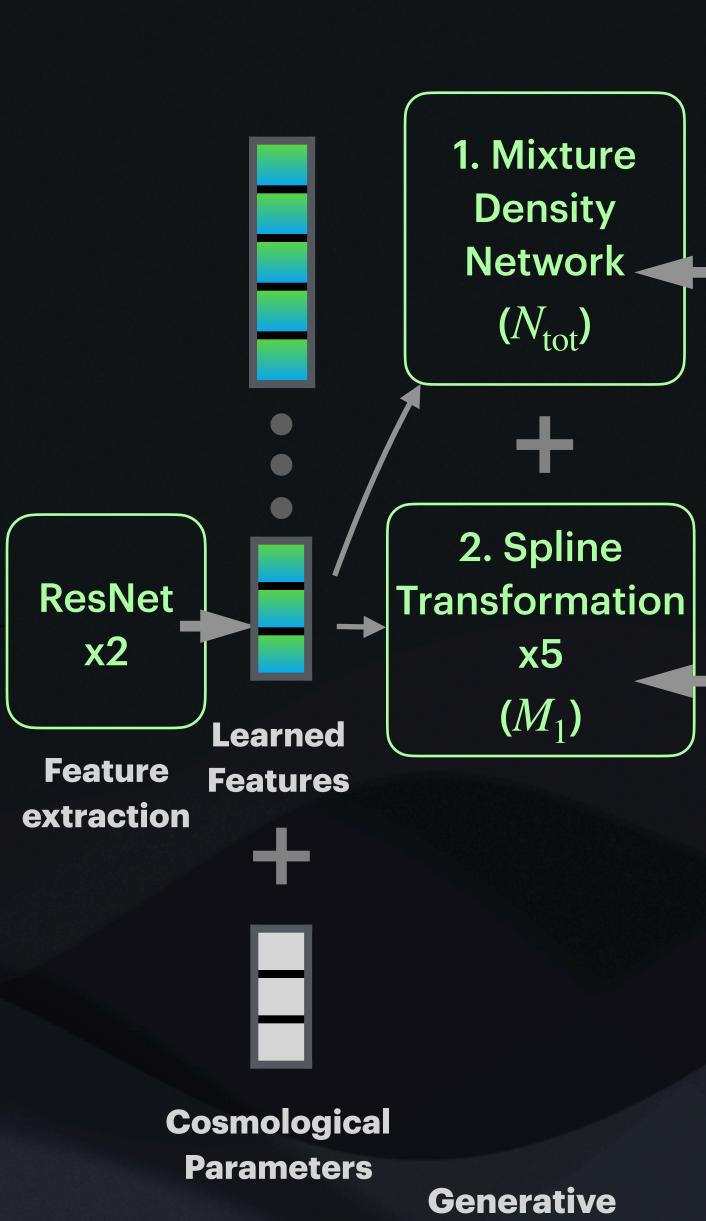


3D
Halo distribution
(N-body Simulation)
Target

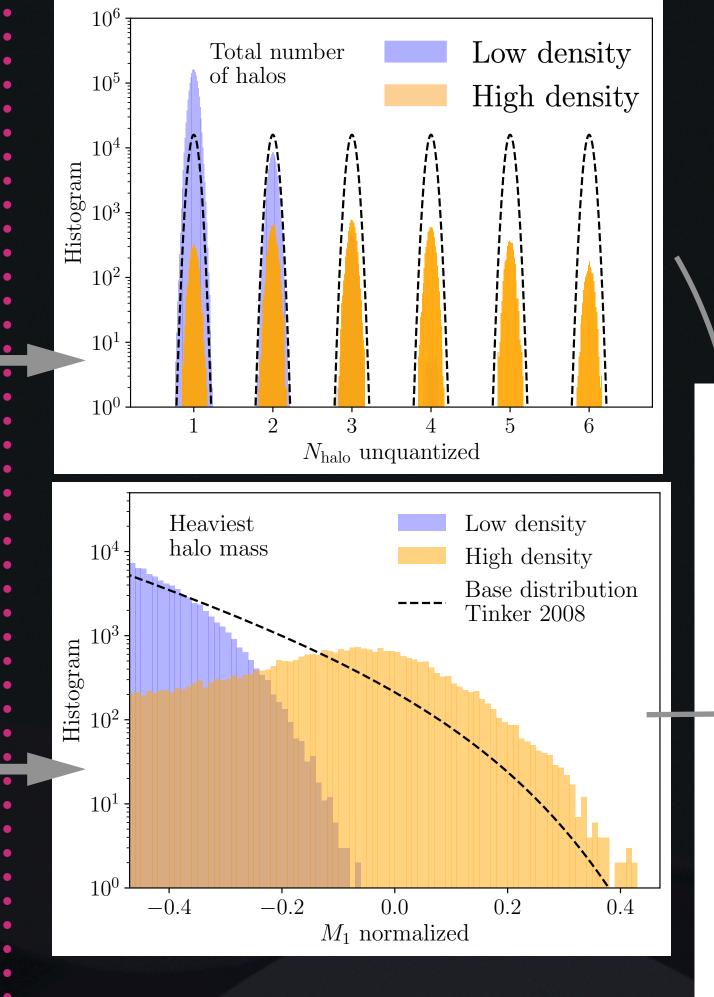
CHARM OVERVIEW arxiv:2409.09124



3D
Matter density field
(PM Simulation)
Input

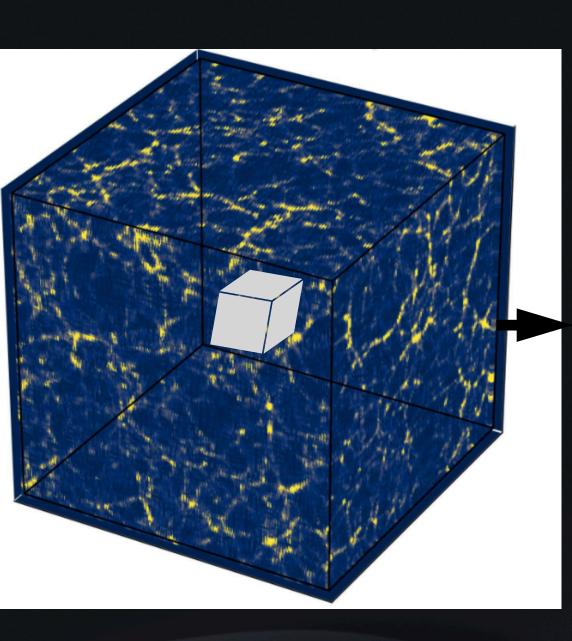


Model

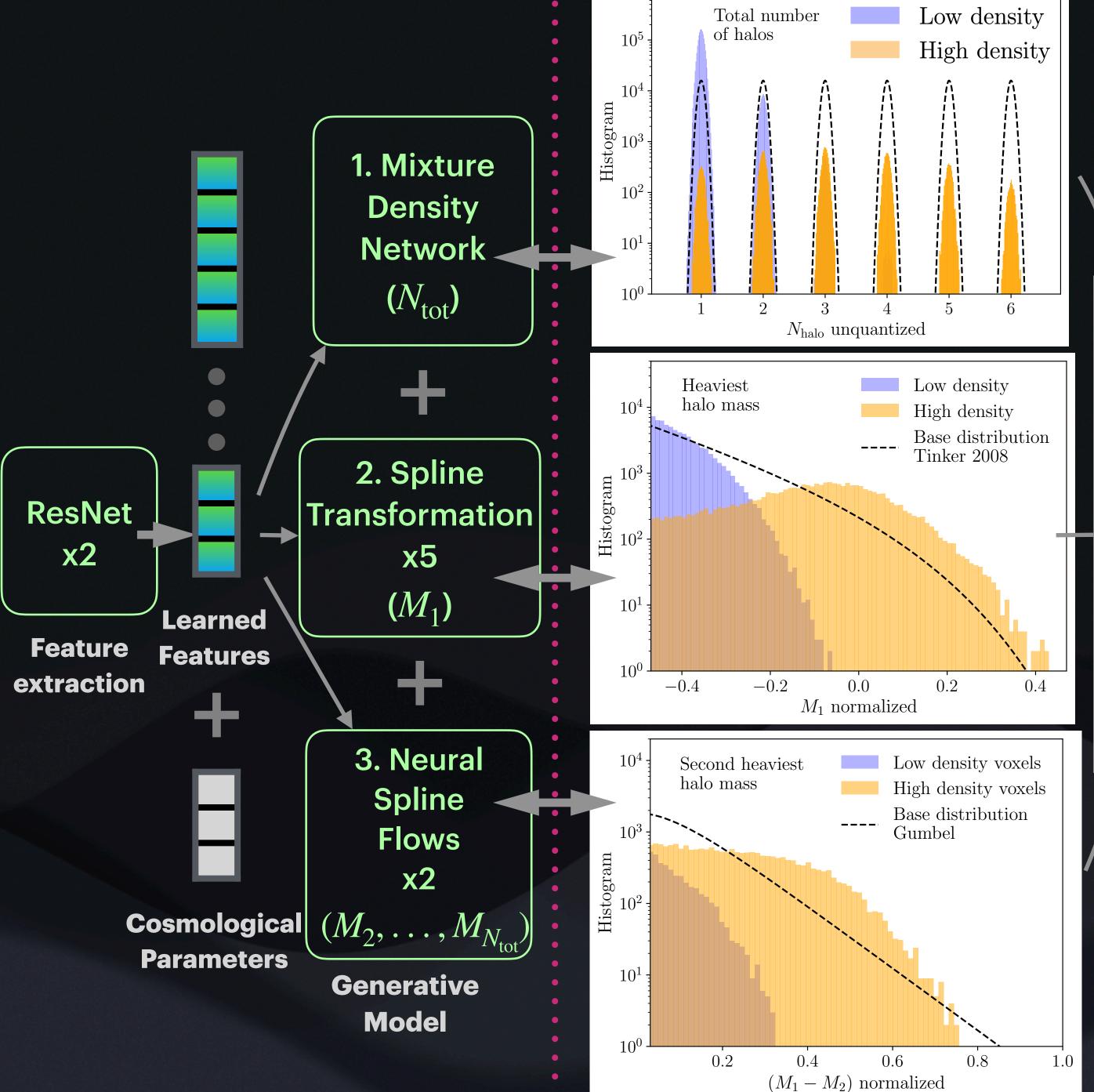


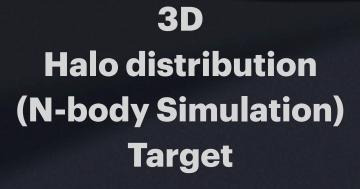


CHARM OVEIVIEW arxiv:2409.09124

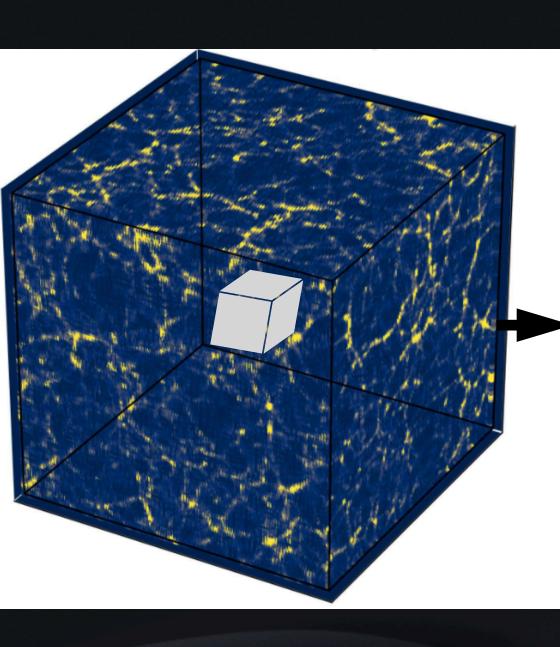


3D
Matter density field
(PM Simulation)
Input

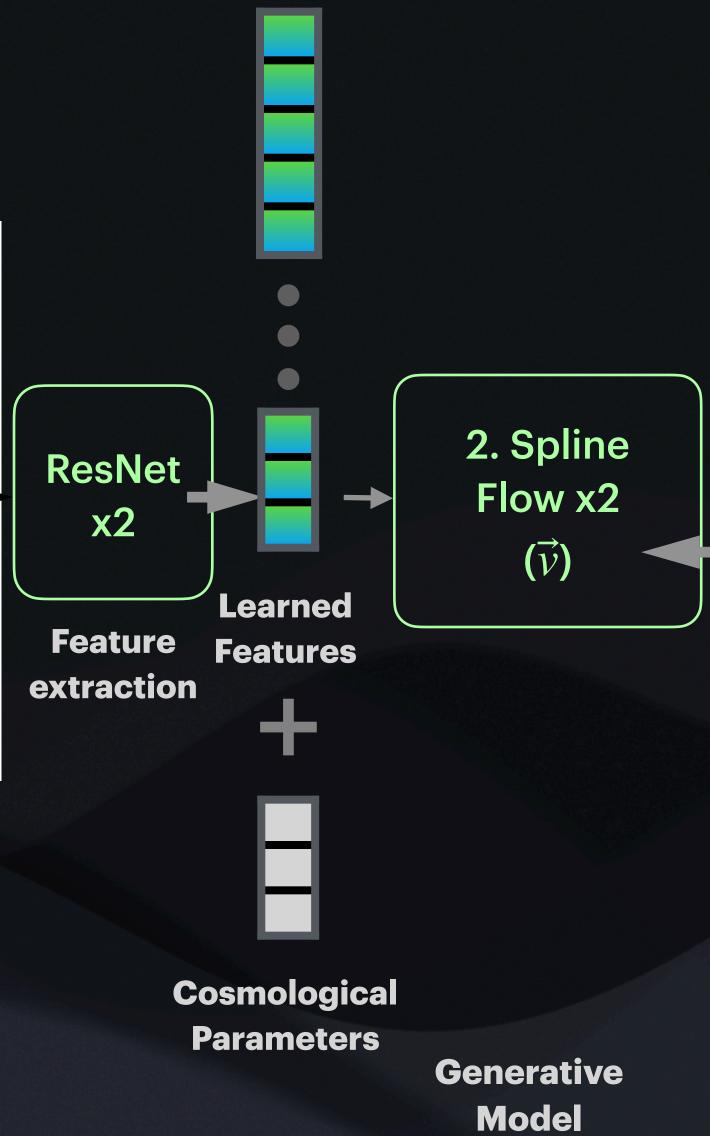


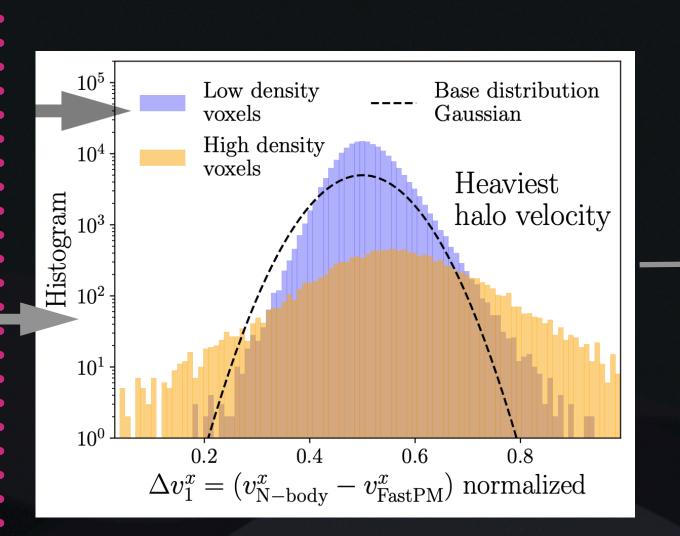


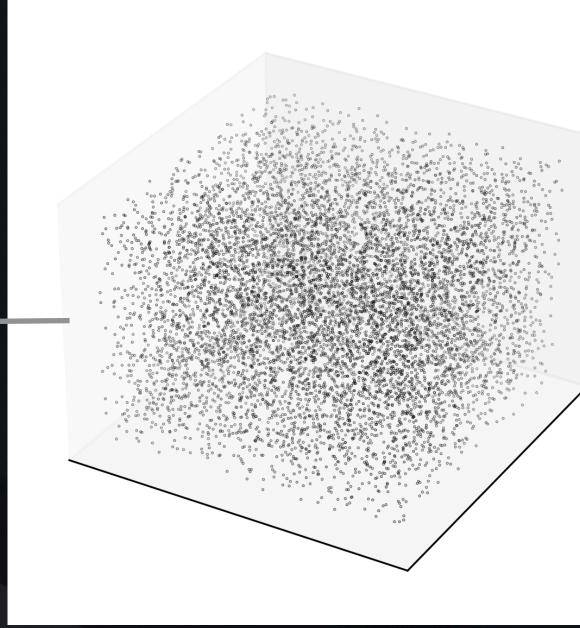
CHARM OVEIVIEW arxiv:2409.09124



3D
Matter density field
(PM Simulation)
Input

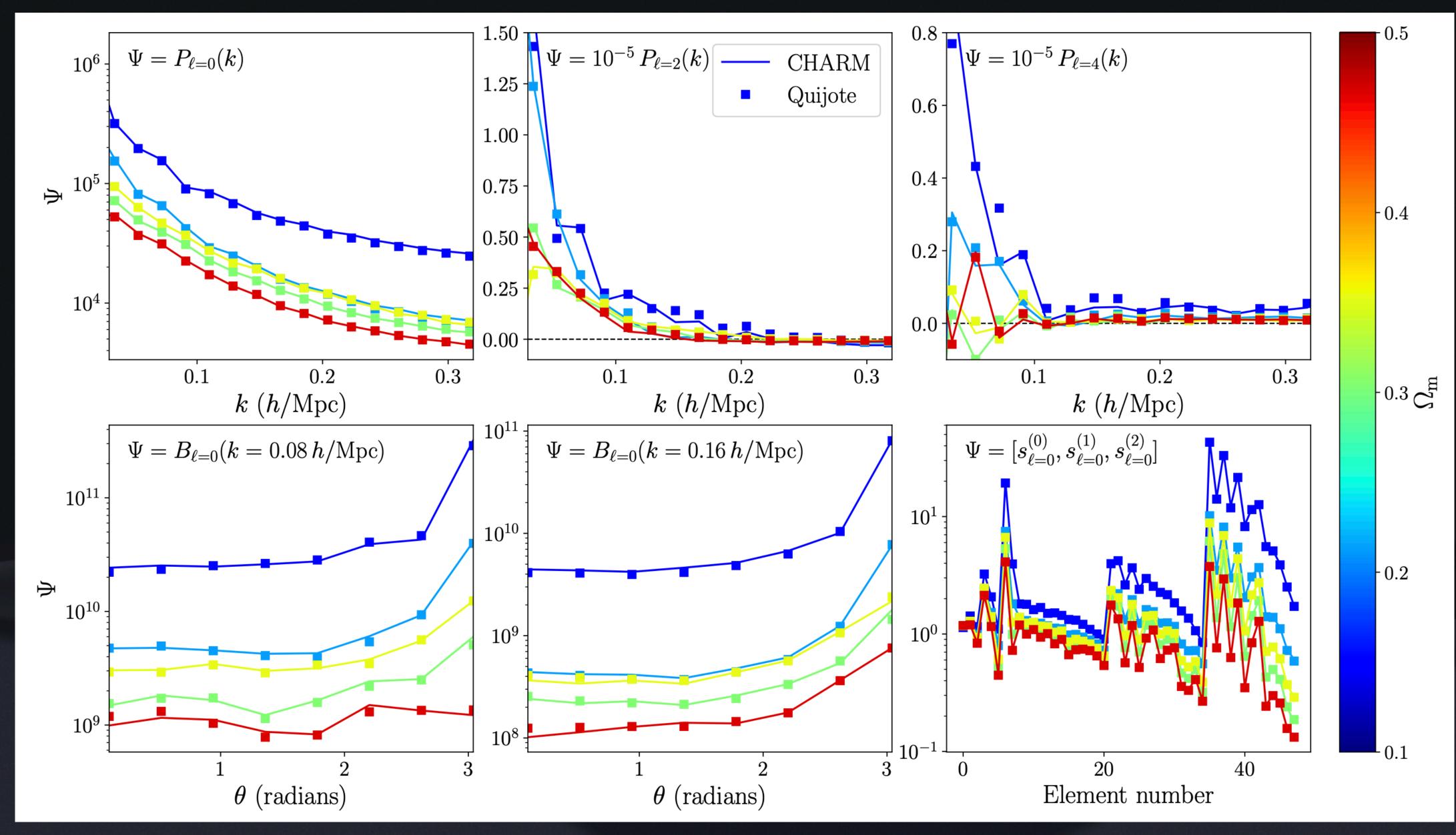




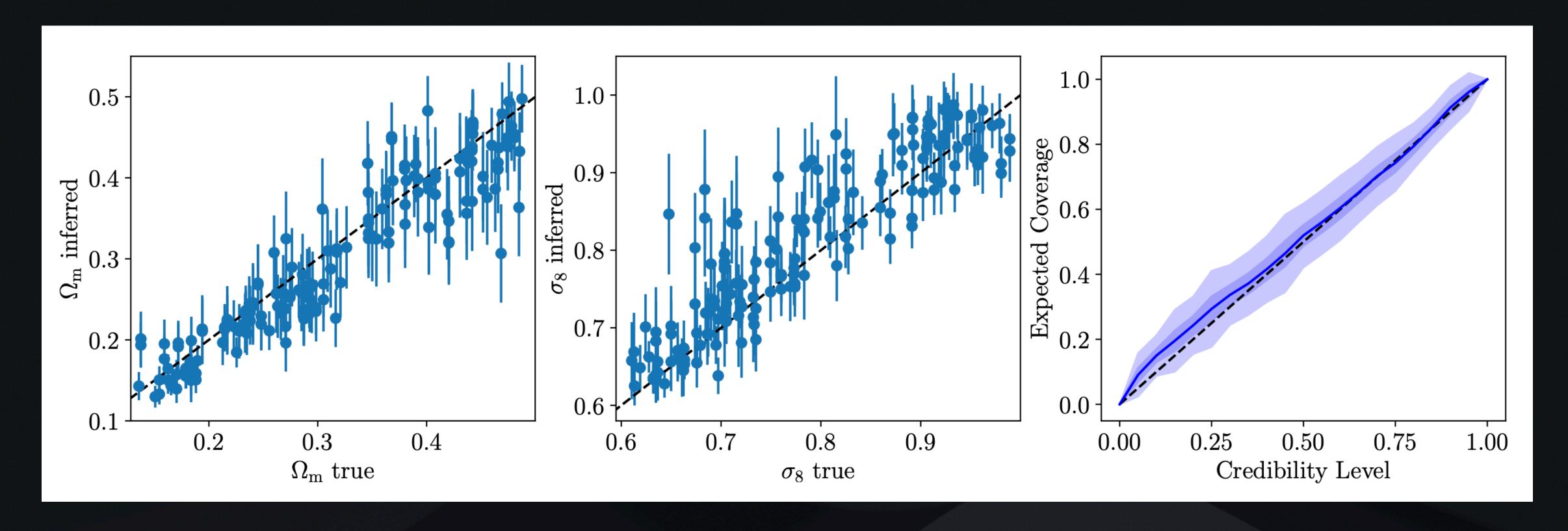


3D
Halo distribution
(N-body Simulation)
Target

N-pt performance — Redshift space



Cosmology inference performance with SBI

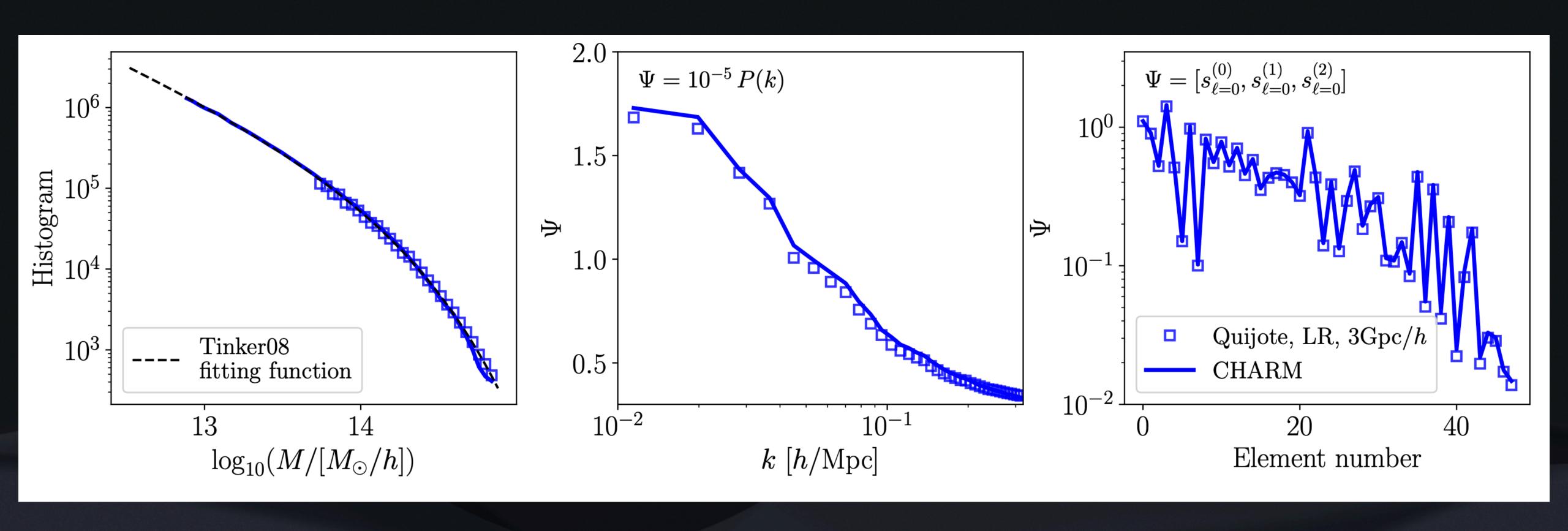


Accounts for halo-galaxy connection with analytical frameworks (eHOD)

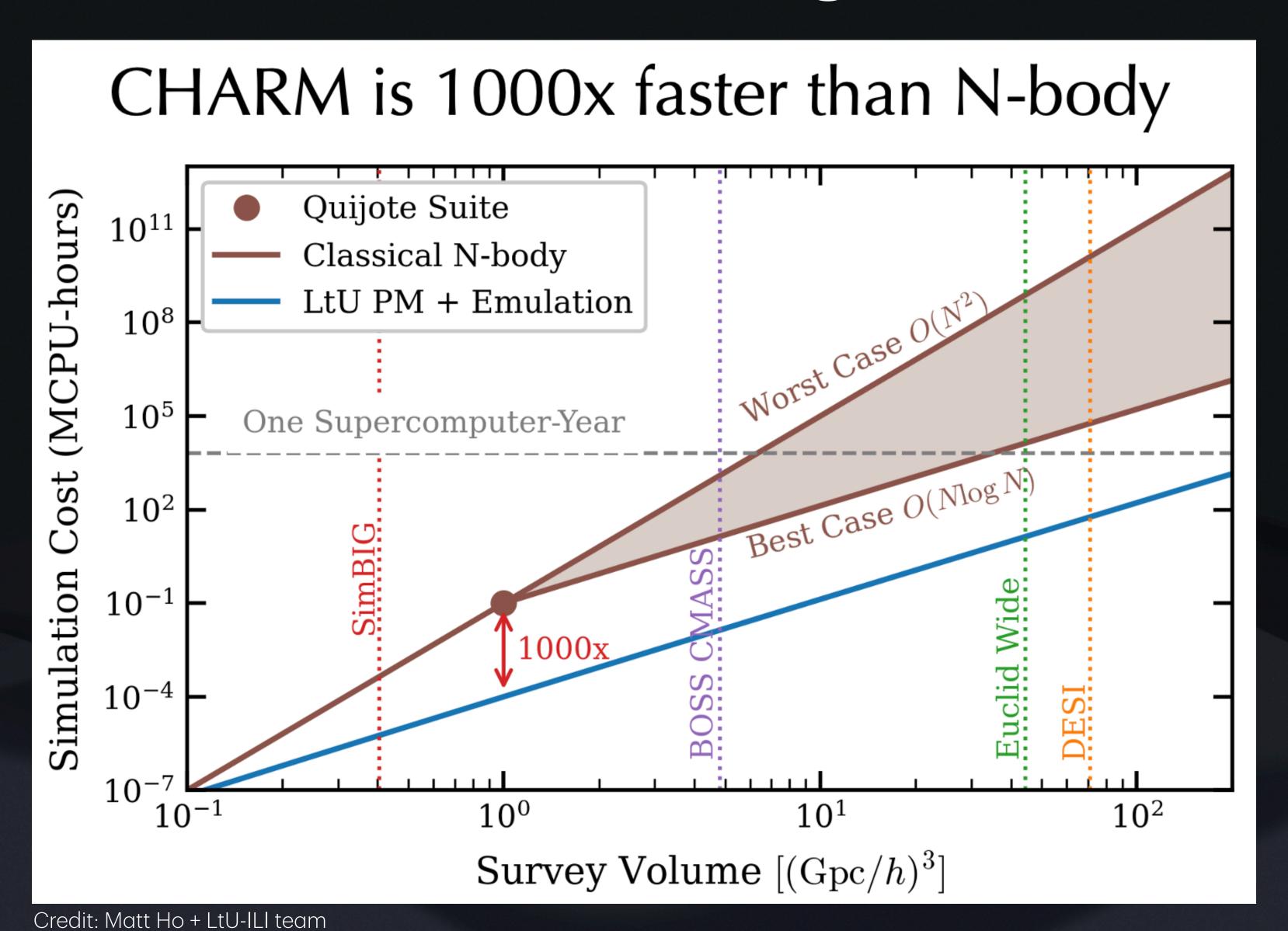
Currently being developed further to analyze the SDSS-NGC data within the LtU.

Being further developed to generalize to redshifts to get halo lightcone

Generalizes to large volumes



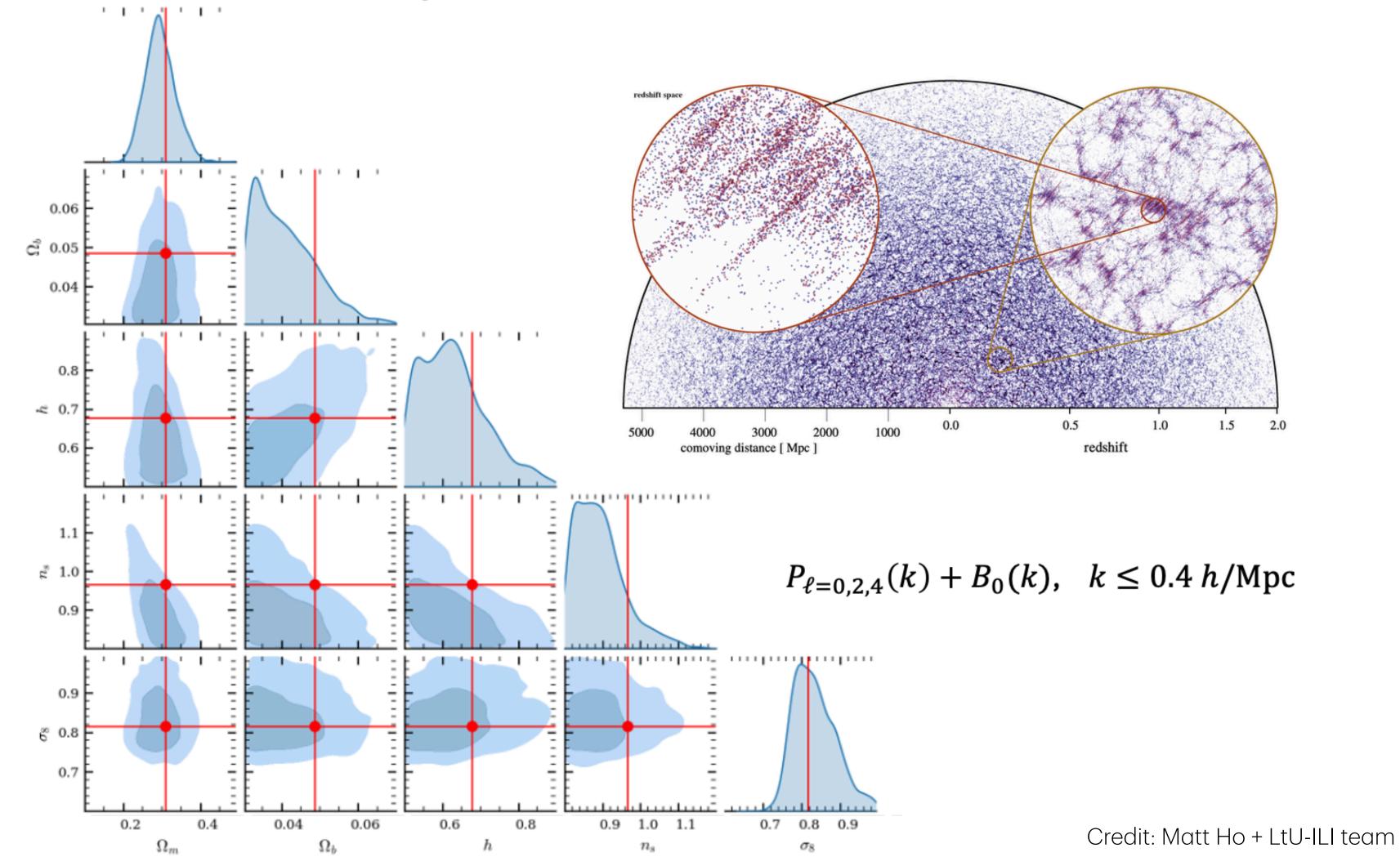
Can now run suite of large simulations



Robustness

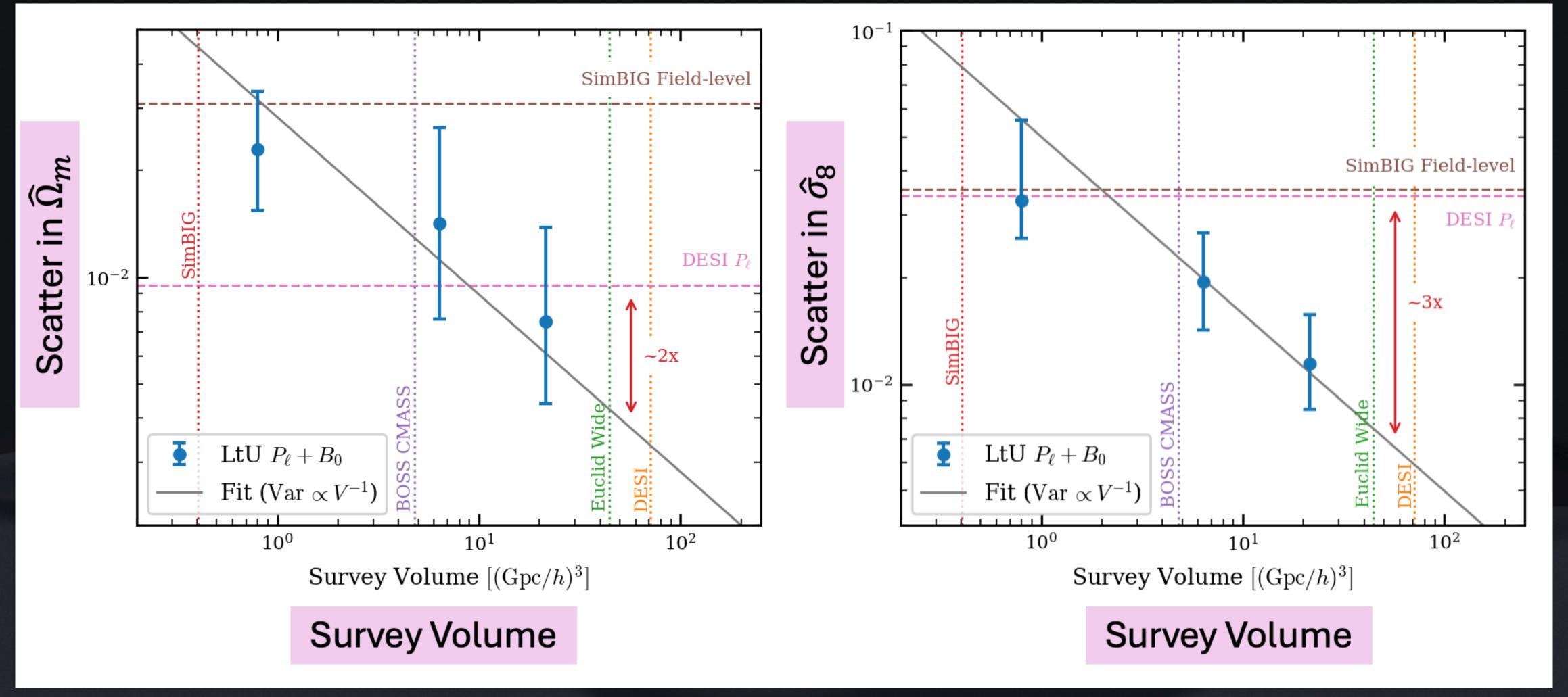
- Validating on completely independent simulation:
 - Different halo finder
 - Using a completely different way of populating the galaxies
- Also similarly validated on the Abacus simualtions:
 - Different gravity solver, different mass definitions, multiple cosmologies





Forecasted constraints

SDSS data analysis ongoing within LtU



Here halo locations are localized to ~8Mpc. But we want to go to smaller scales!

Can treat positions as properties as well.

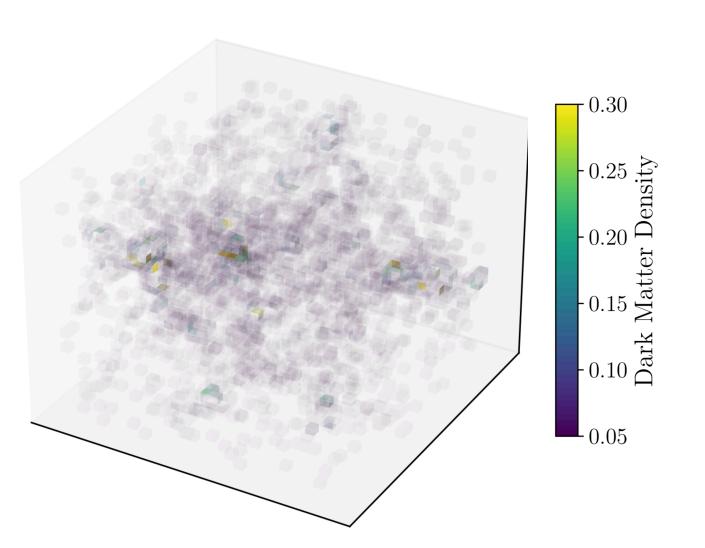
Here halo locations are localized to ~8Mpc. But we want to go to smaller scales!

Can treat positions as properties as well.

We want an architecture that is very efficient and scalable for long-length auto-regressive dependencies.

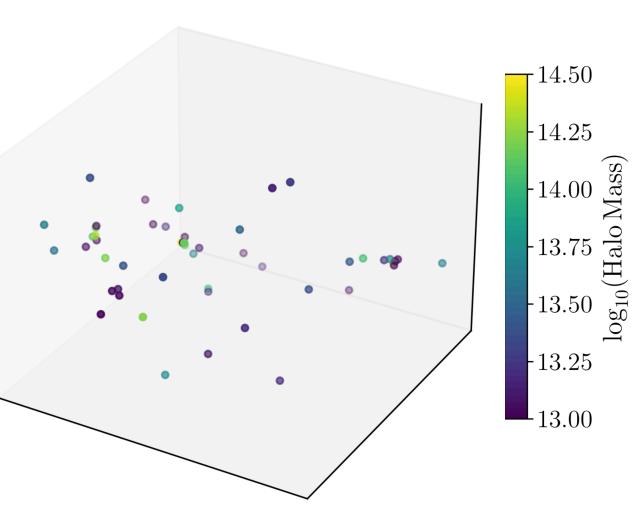
How else to go from left to right distribution?

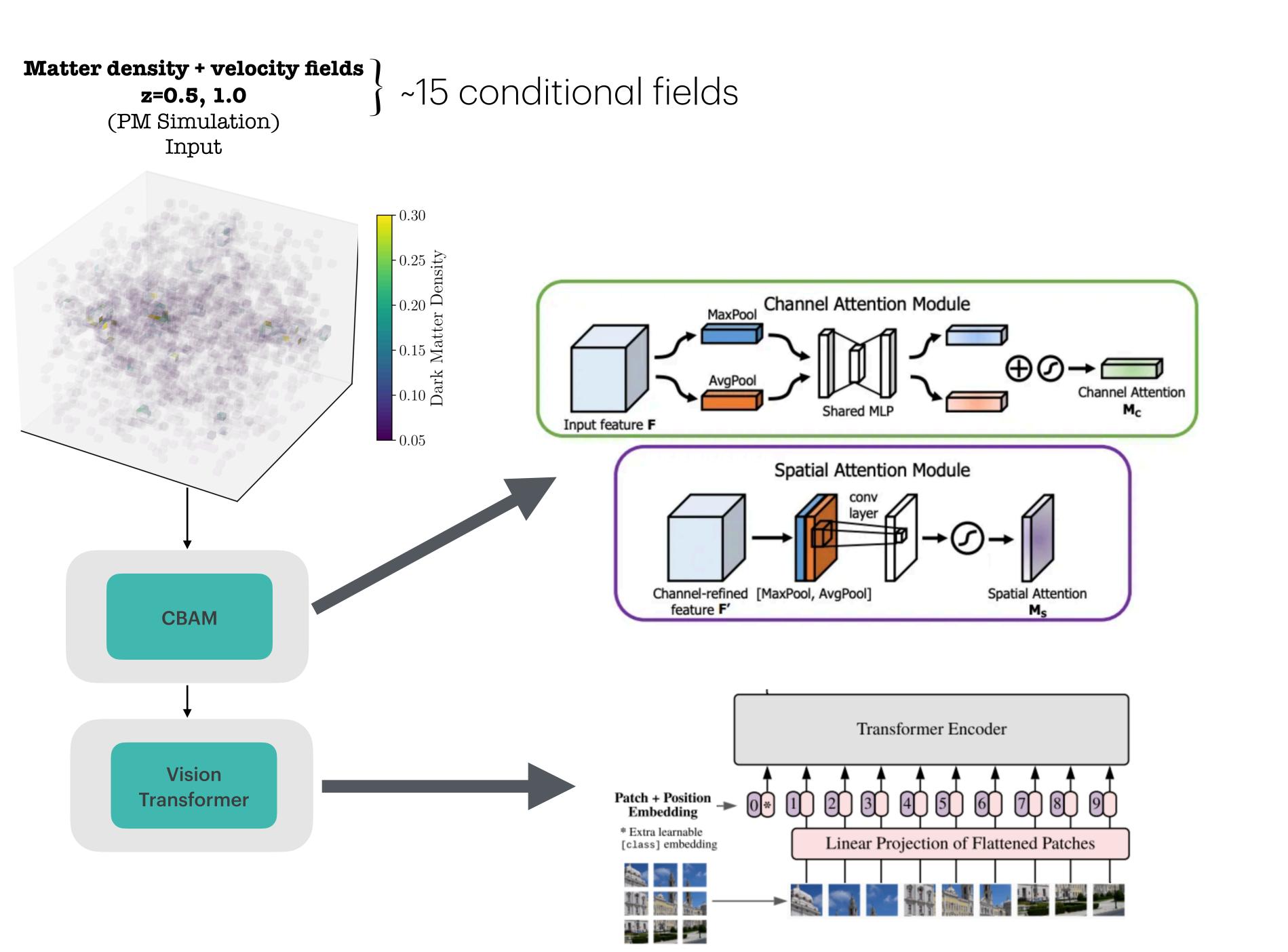
(PM Simulation) Input

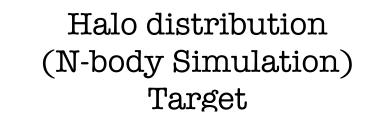


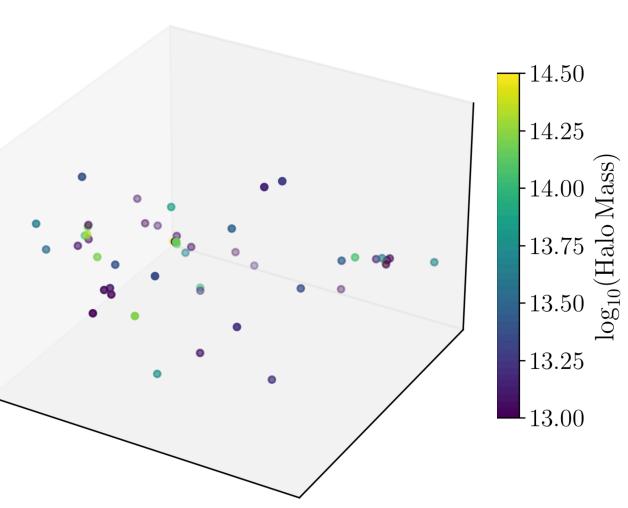


Halo distribution (N-body Simulation) Target

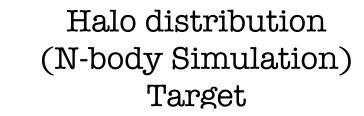


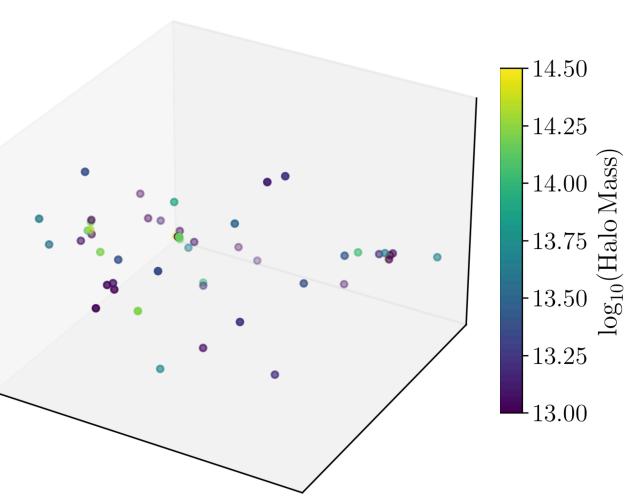


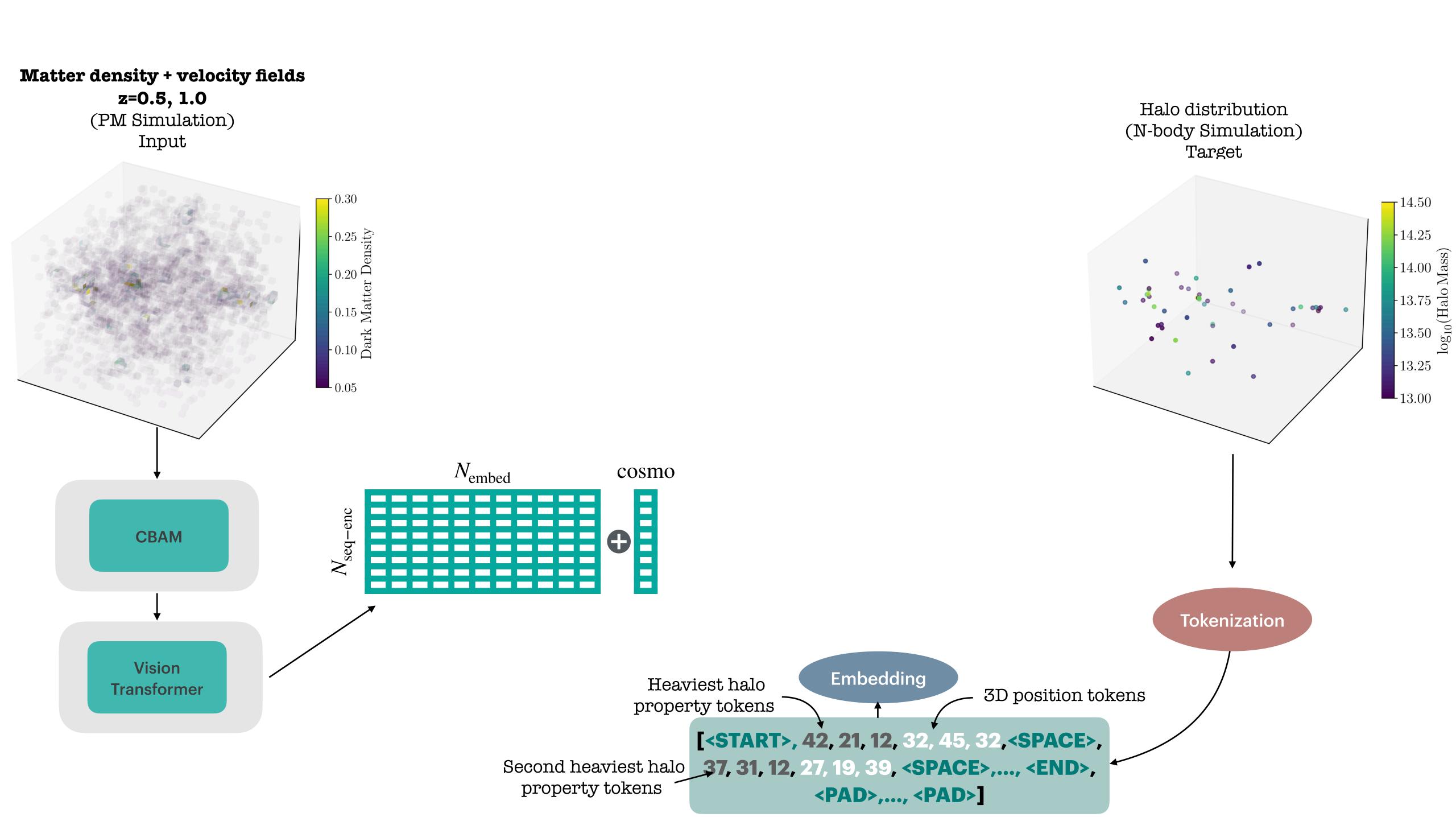


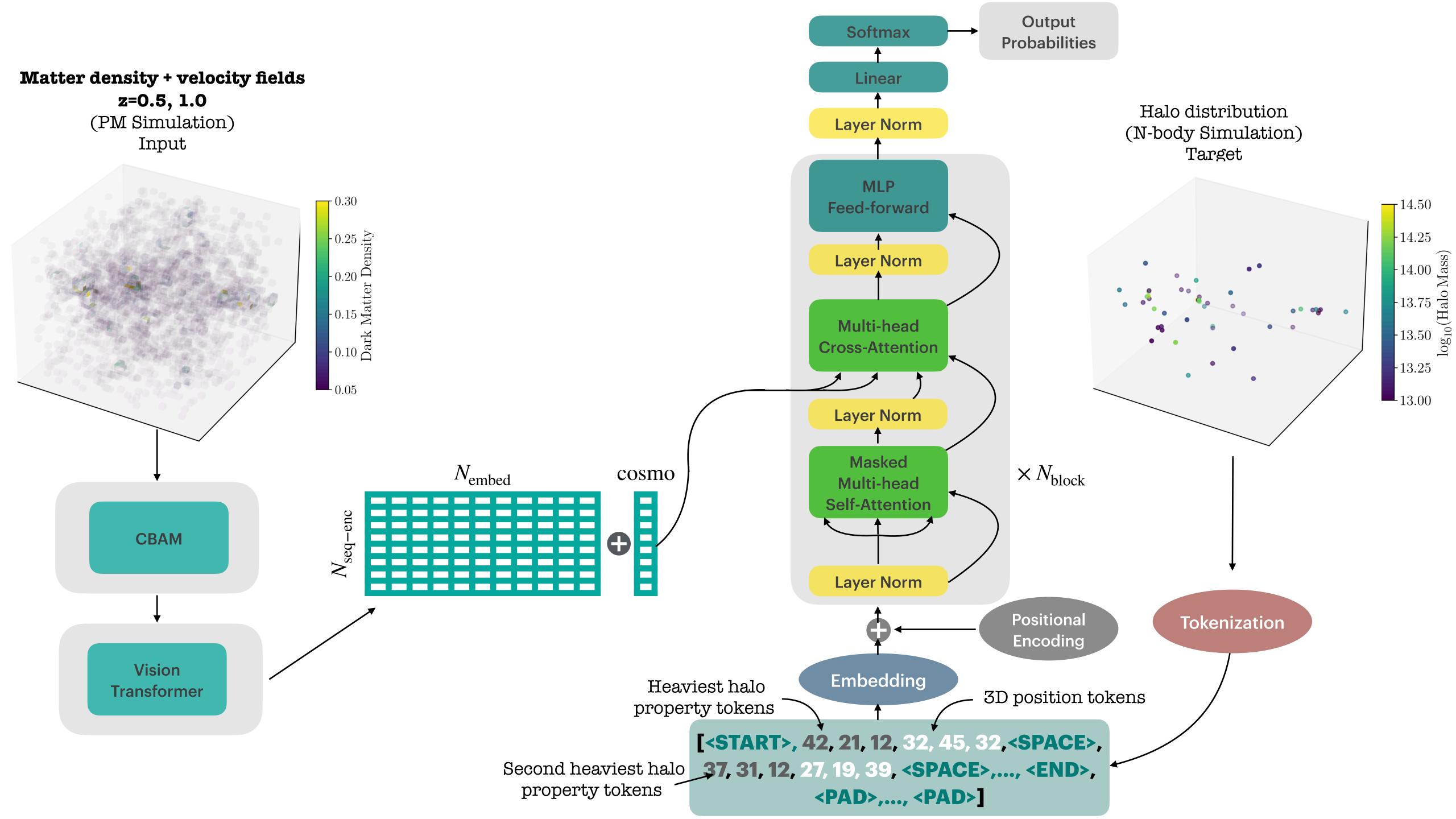


Matter density + velocity fields **z=0.5, 1.0** (PM Simulation) Input $\top 0.30$ - 0.25 - 0.20 - 1.00 - Dark $N_{\rm embed}$ cosmo • CBAM Vision Transformer

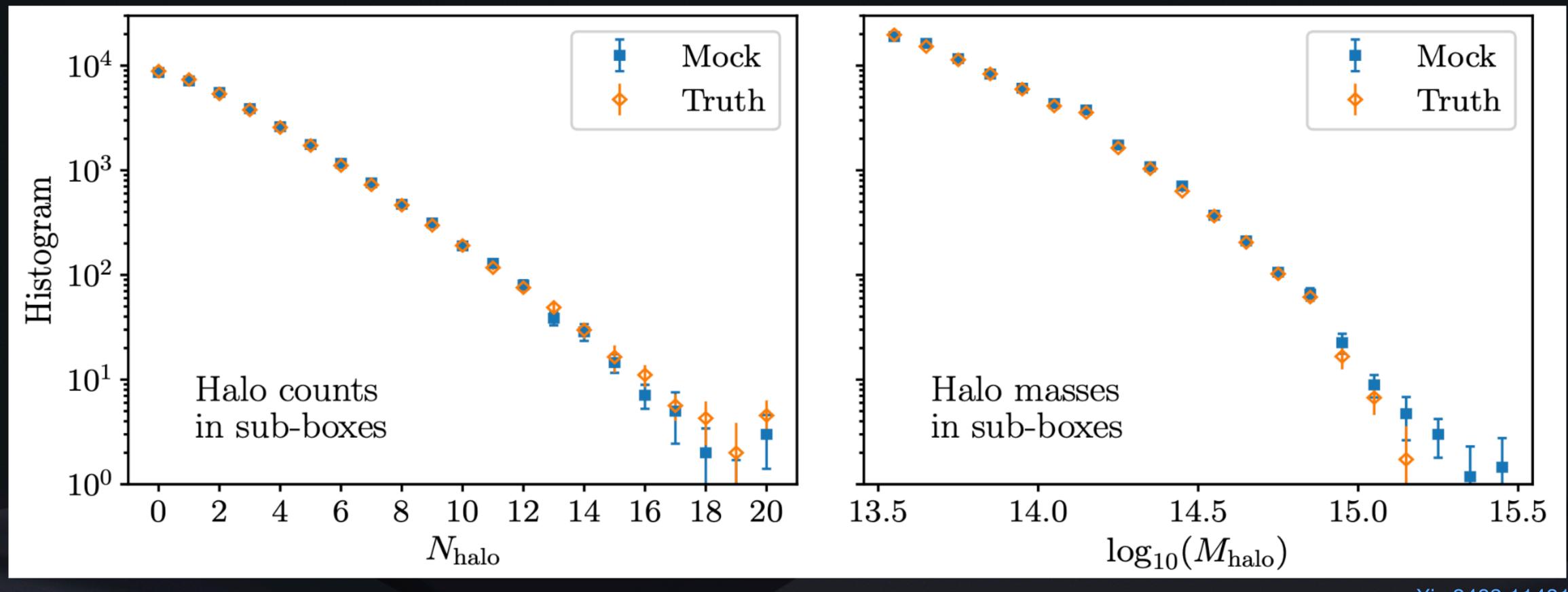




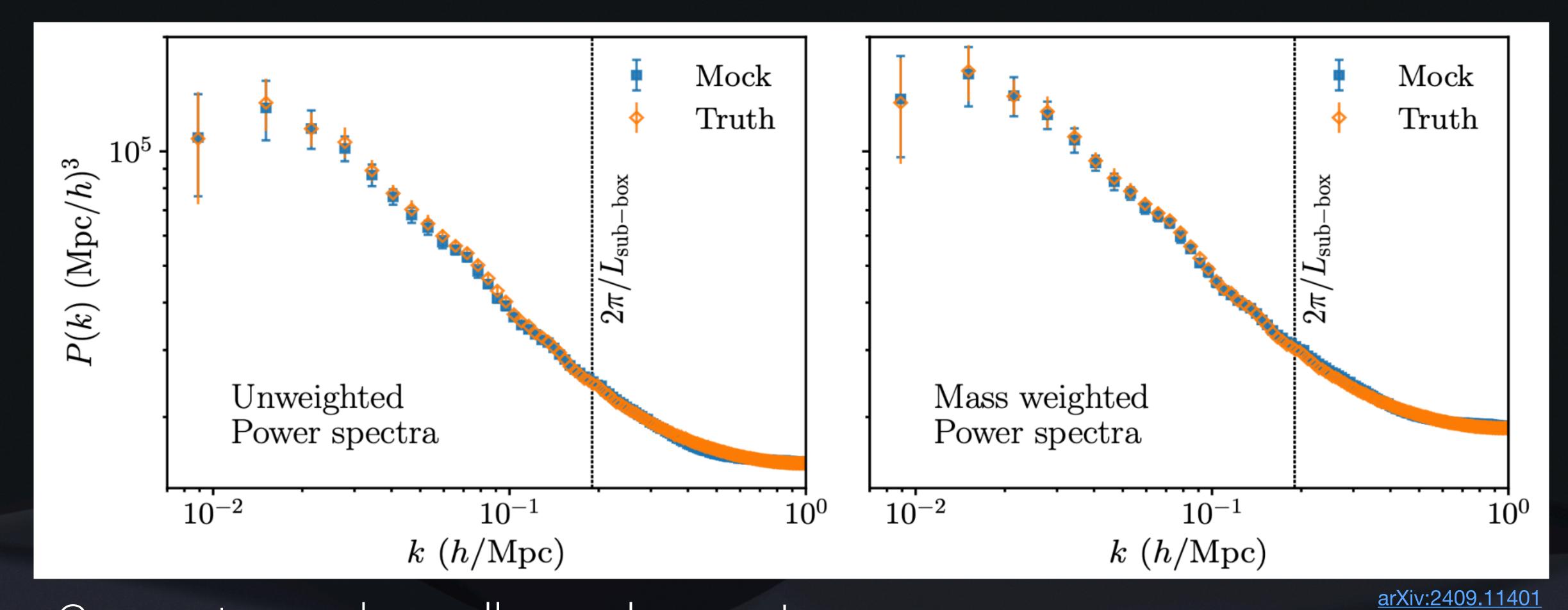




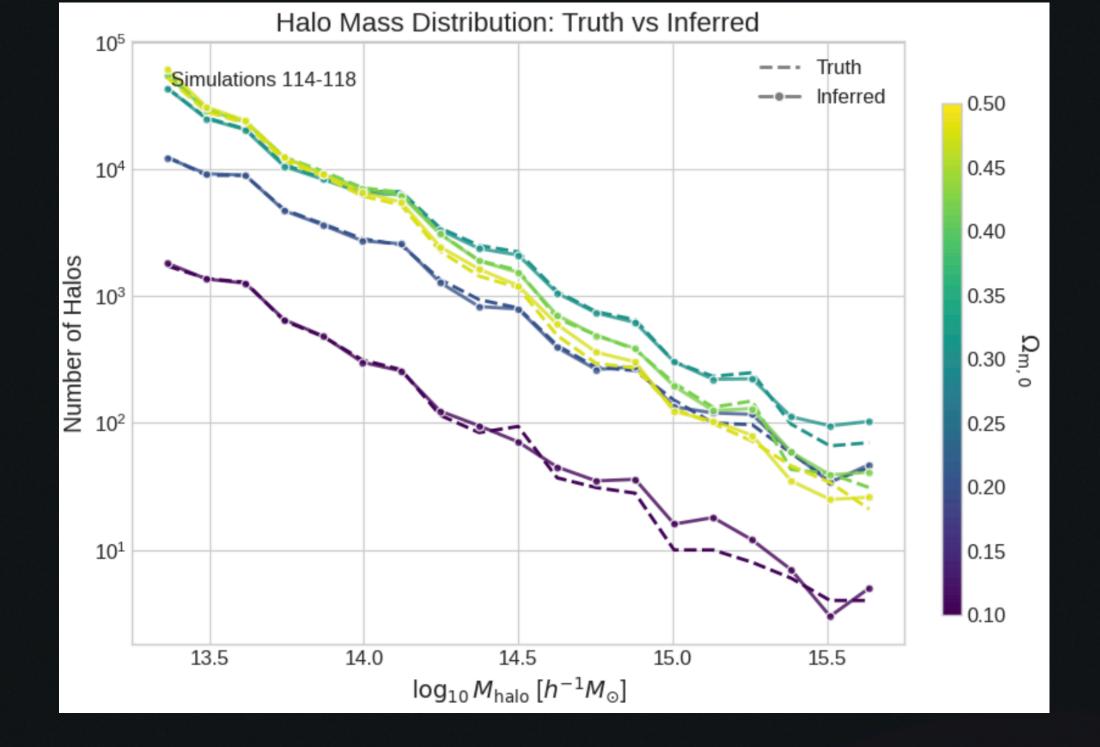
Fixed cosmology: 1-pt performance



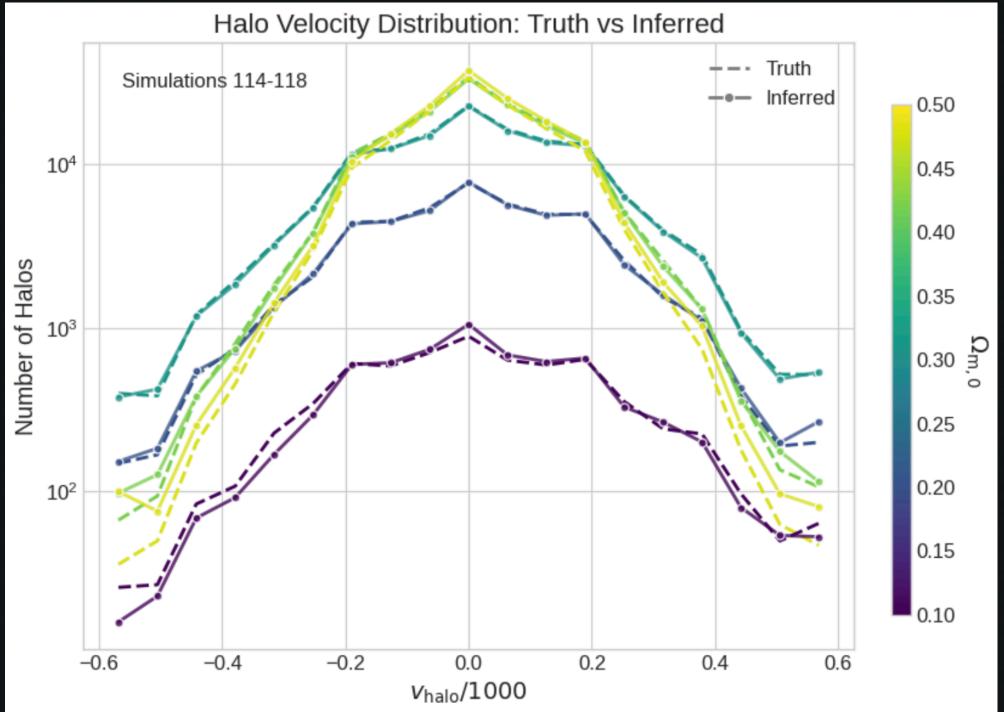
Fixed cosmology: 2-pt performance

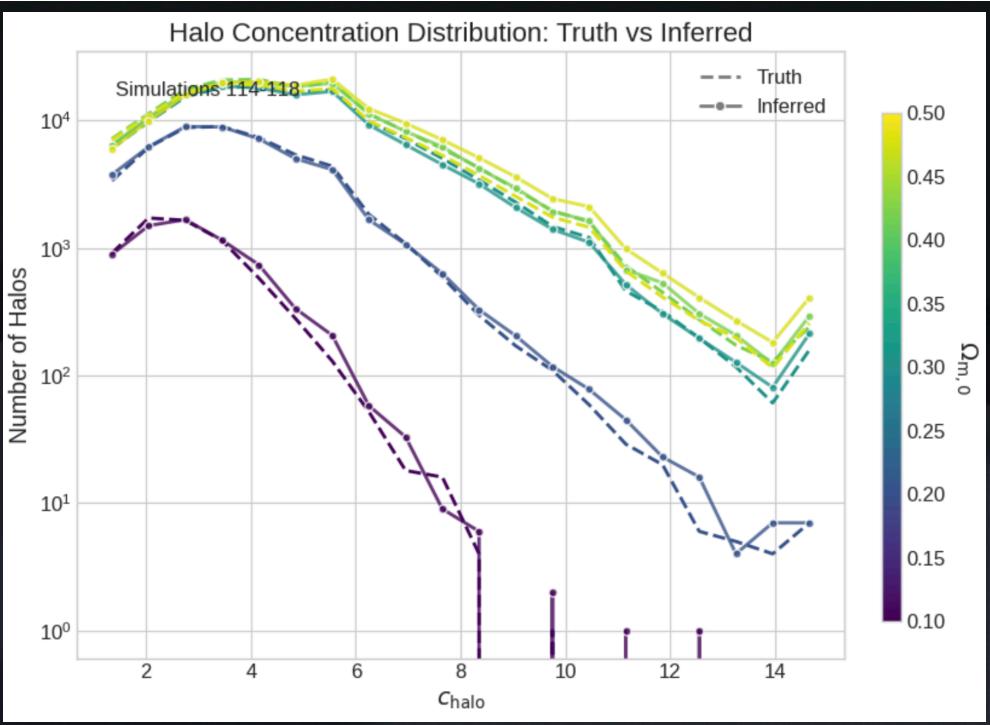


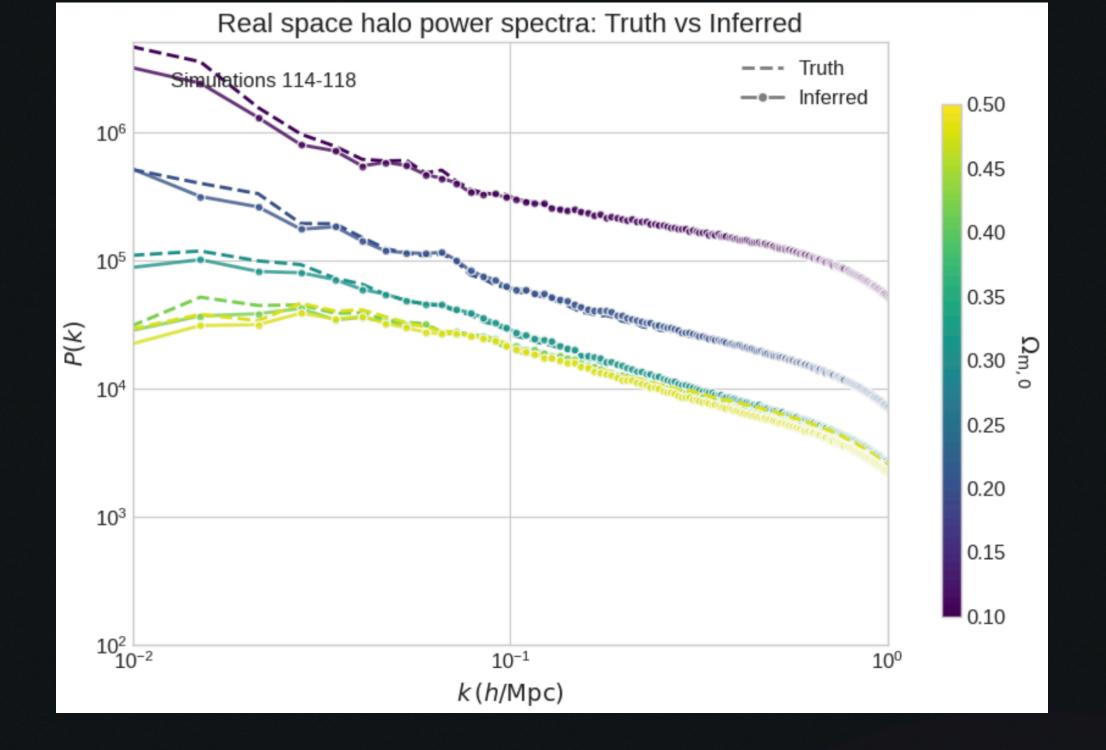
Can go to much smaller scales now!



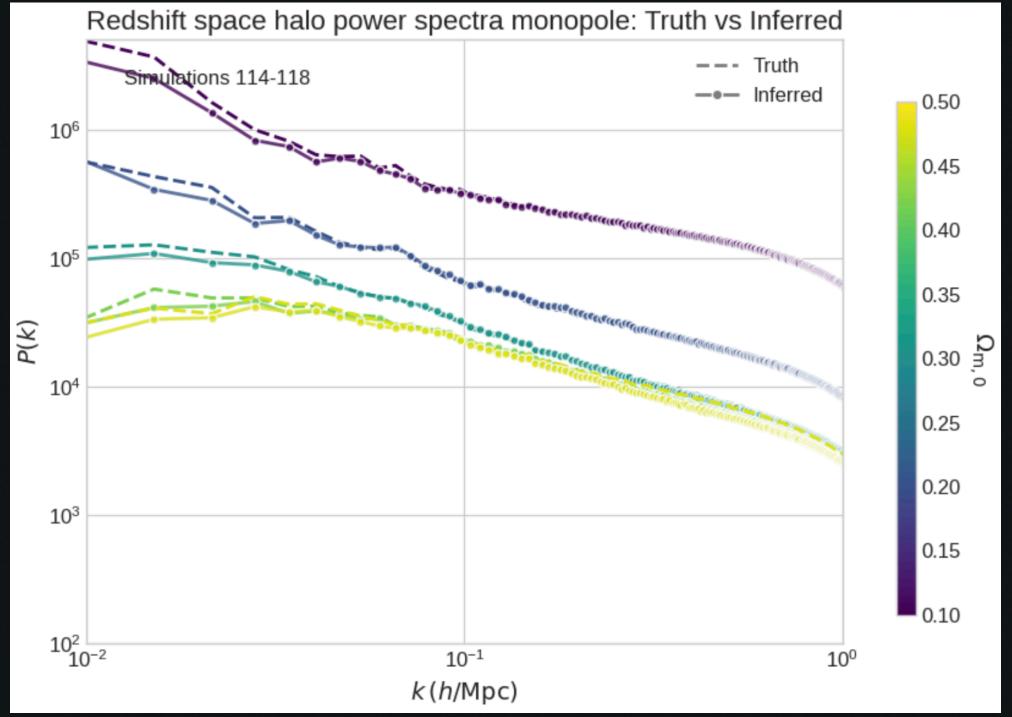
- Preliminary results
- 1-pt inference on test simulations
- Varying cosmologies have orders-of-magnitude different number of halos

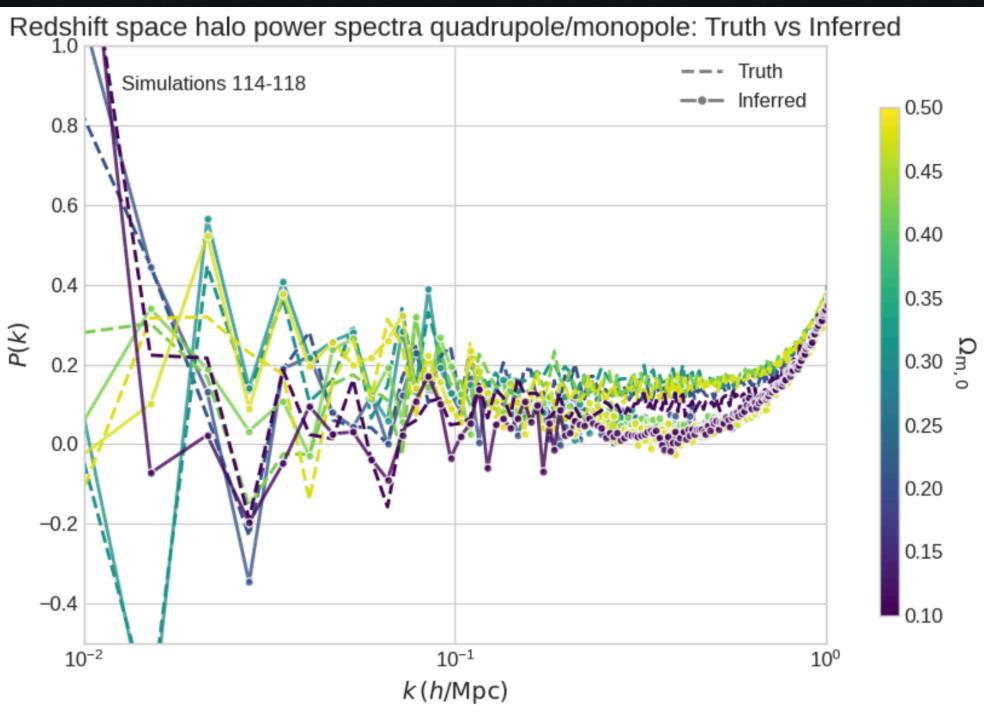






- Preliminary results
- 2-pt inference on test simulations
- Showing both real-space and redshift-space power spectra multipoles

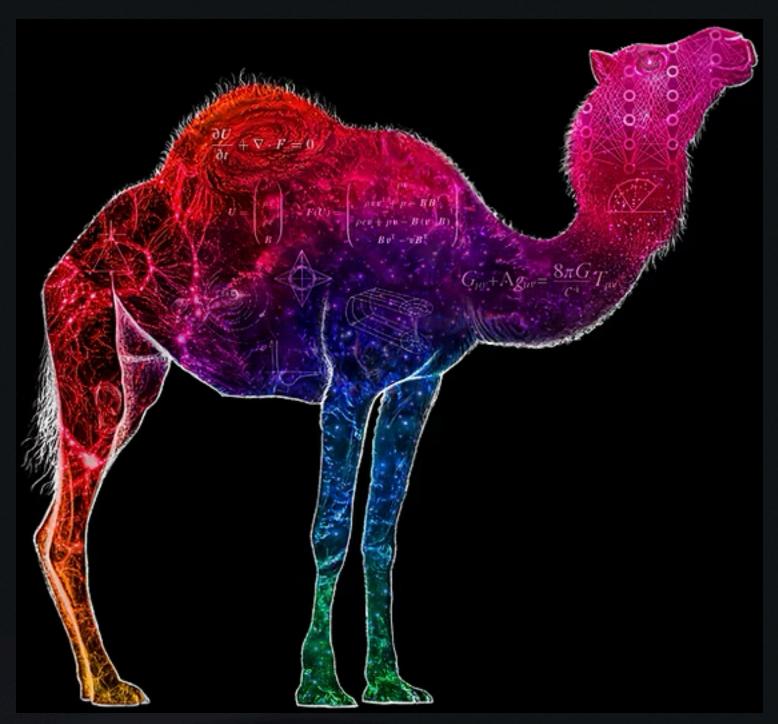




This was for halos and their properties only, requires halo-galaxy connection!

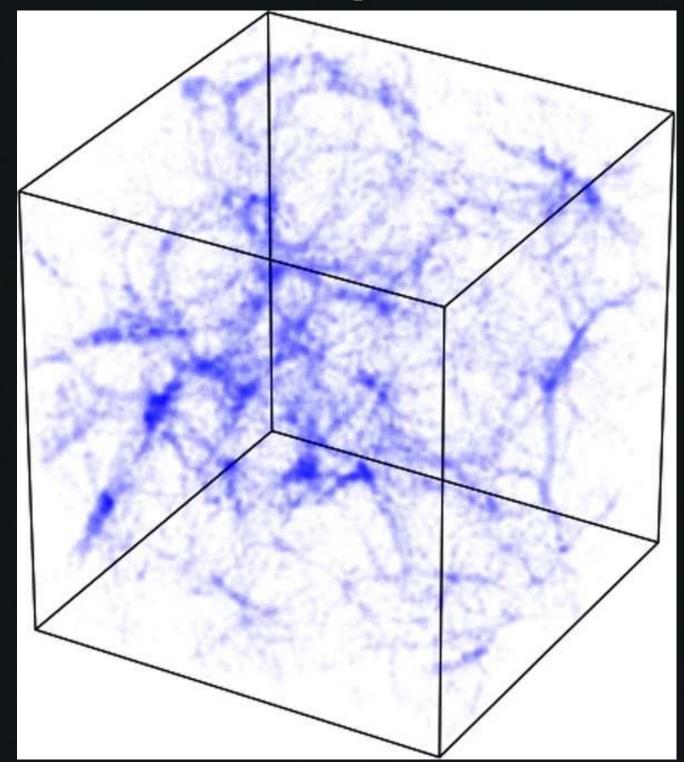
Can we go to galaxies and their observed properties directly while marginalizing over cosmology?

CAMELS Hydro-sims



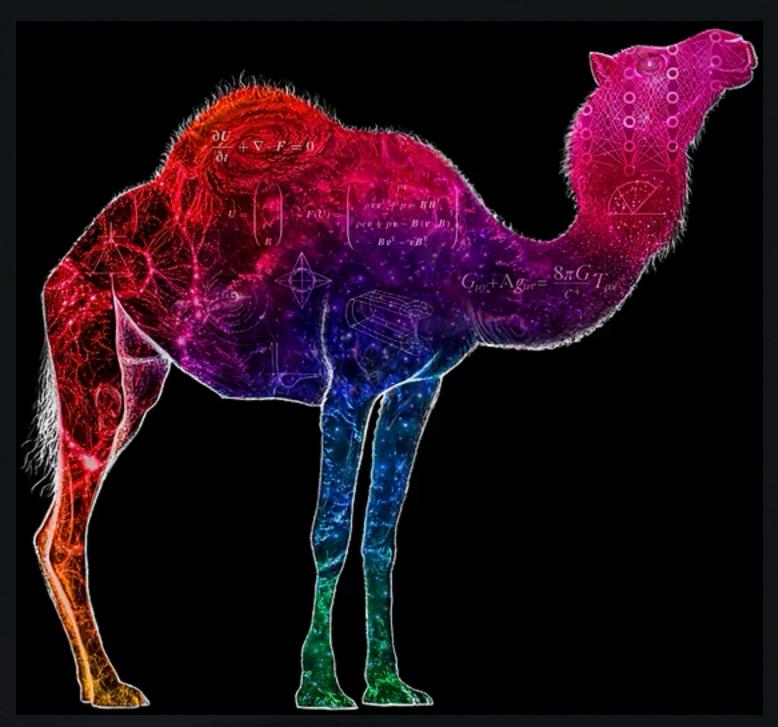
- $(25 \,\mathrm{Mpc}/h)^3 \,\mathrm{volume}$
- ~6000 CPU hours/sim
- Four astrophysical and two cosmological parameters varied
- 1000+ simulations in a latinhypercube space of parameters

N-body sims



- Dark matter only
- Matched IC and cosmology
- ~100x cheaper/sim

CAMELS Hydro-sims

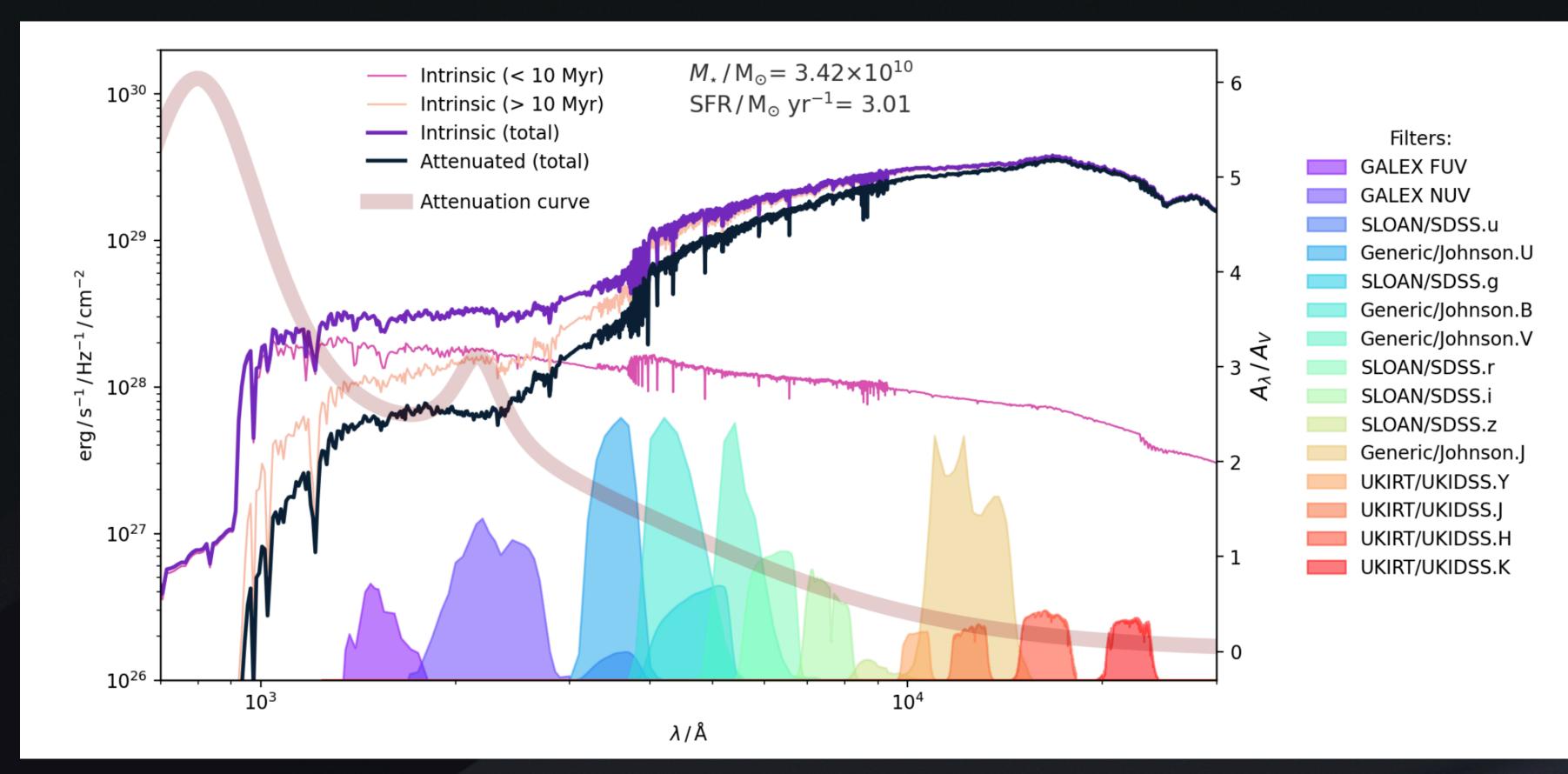


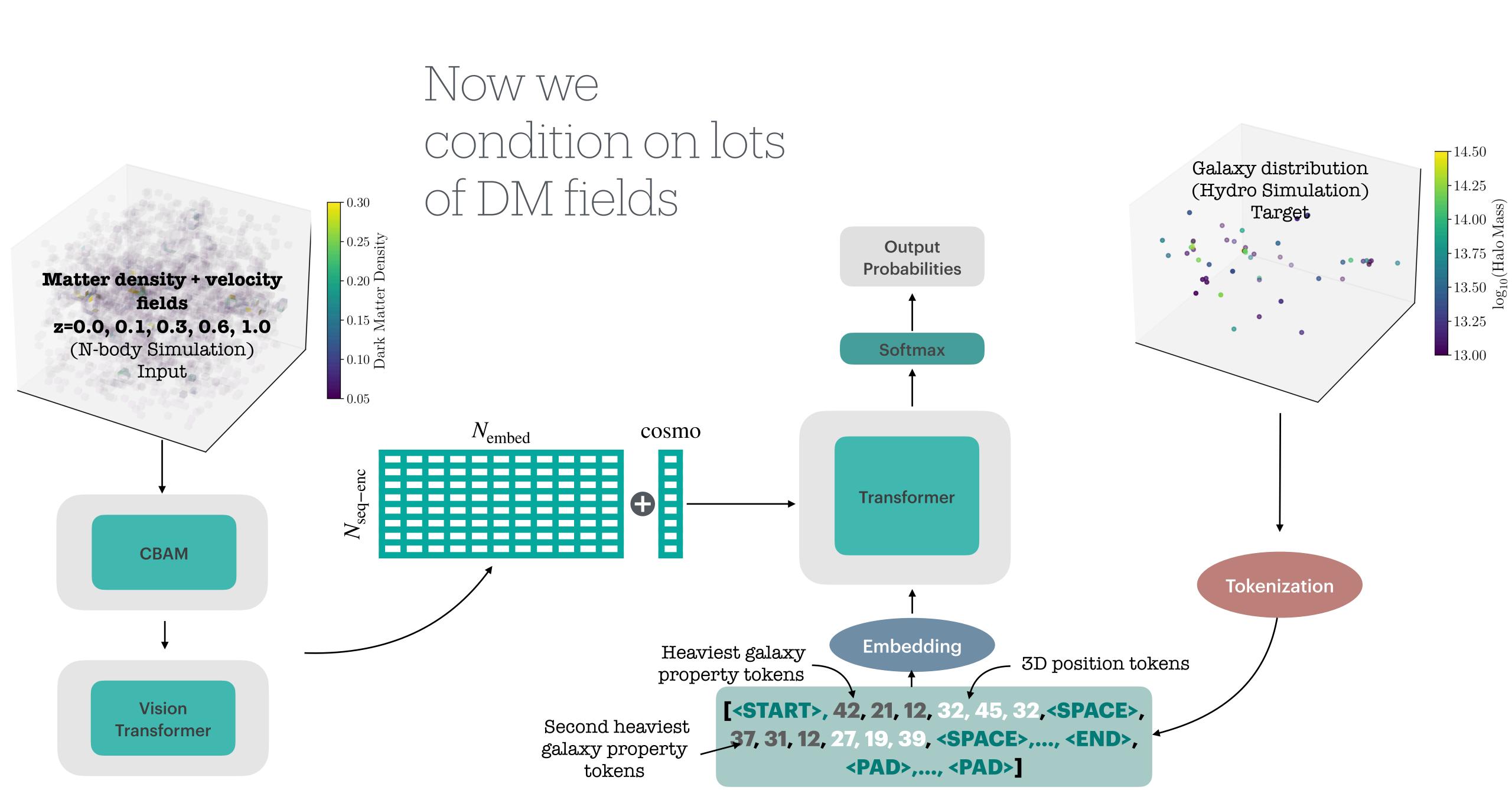
- $(25 \,\mathrm{Mpc}/h)^3 \,\mathrm{volume}$
- ~6000 CPU hours/sim
- Four astrophysical and two cosmological parameters varied
- 1000+ simulations in a latinhypercube space of parameters

Galaxy properties

IllustrisTNG-only

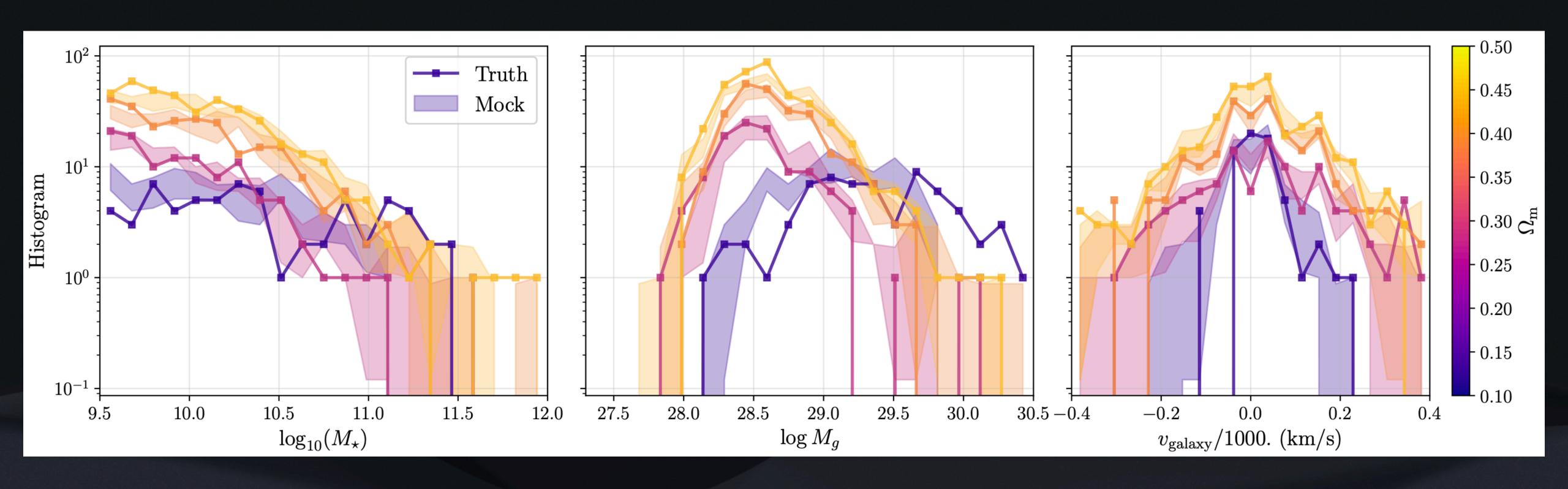
- Stellar mass
- Galaxy line-of-sight velocity
- SDSS photometry (dust attenuated)

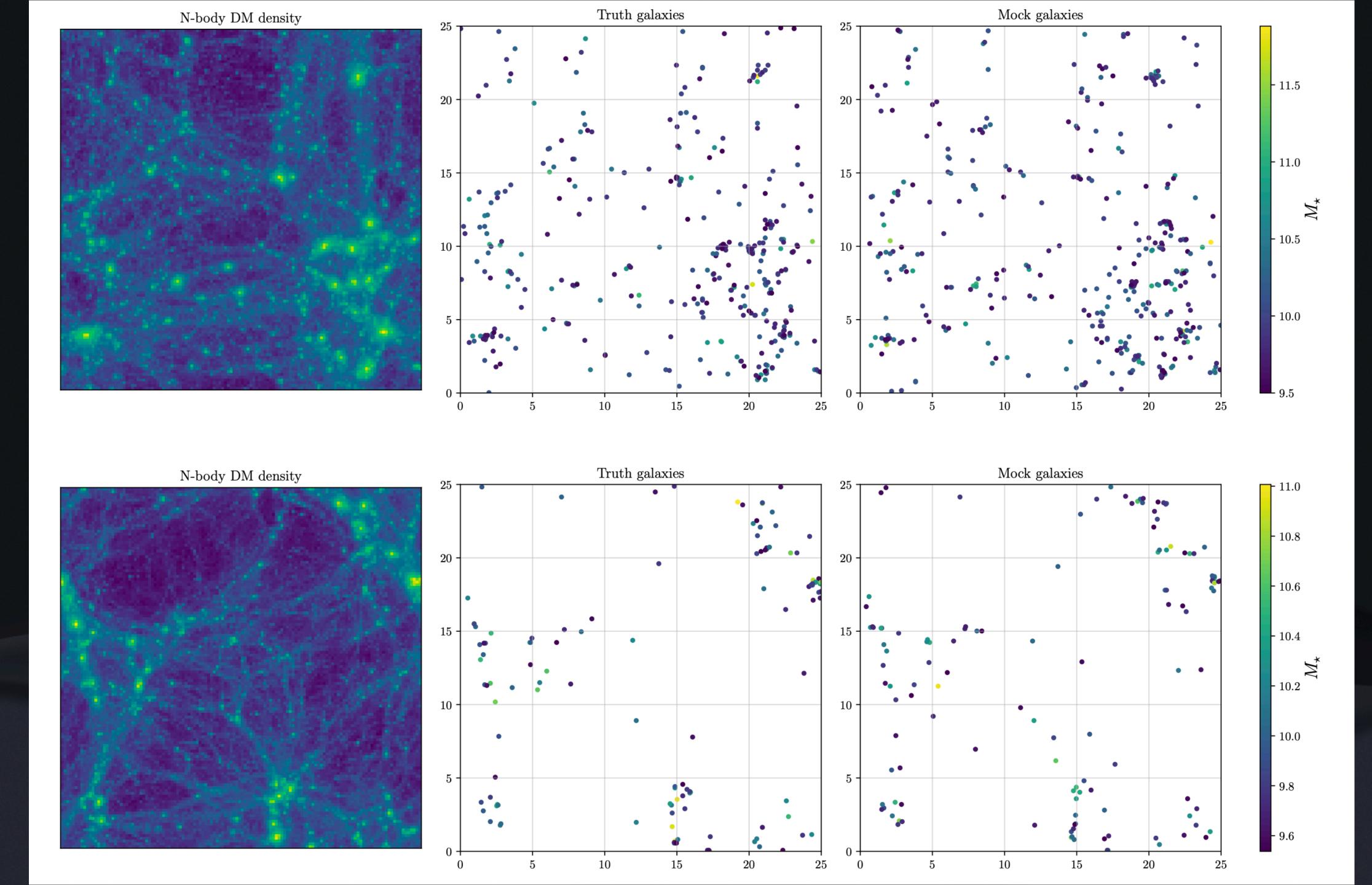




1pt PDFs

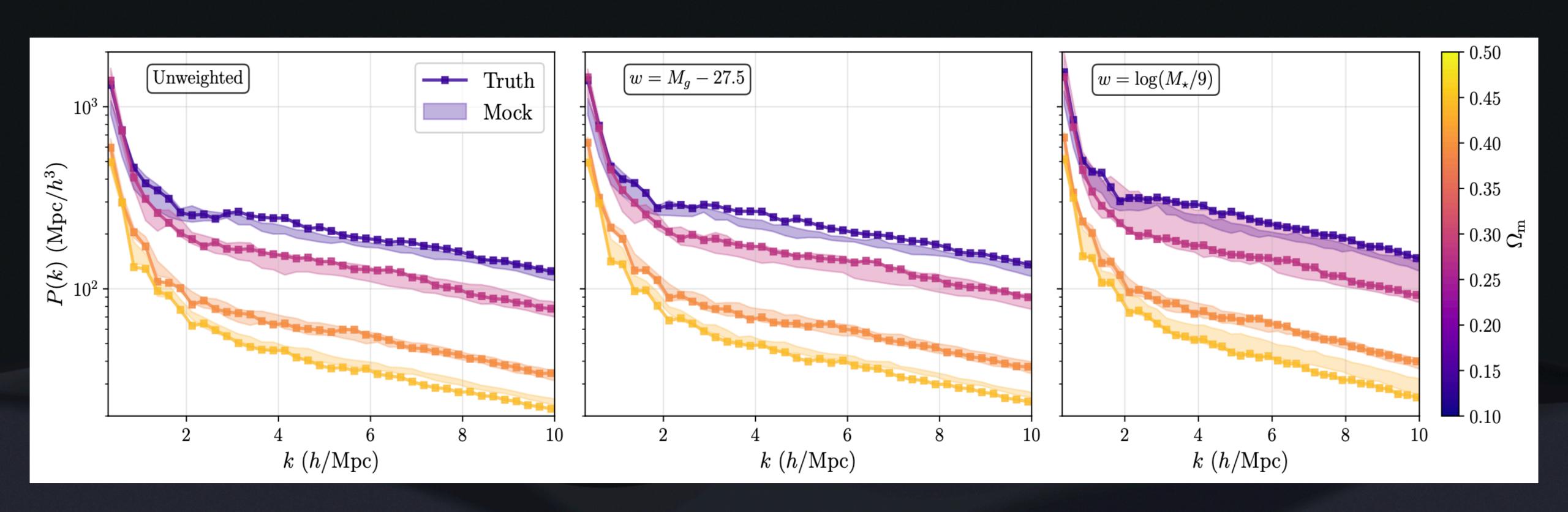
• Inference of galaxy properties on test LH simulations.





2pt (and beyond) correlations

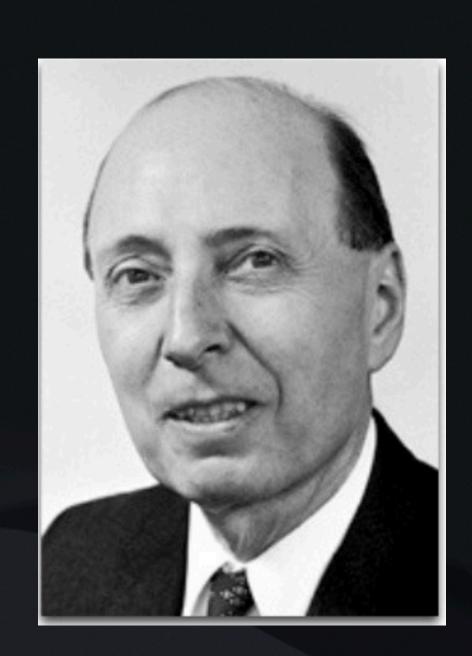
• Inference of galaxy properties on test LH simulations.



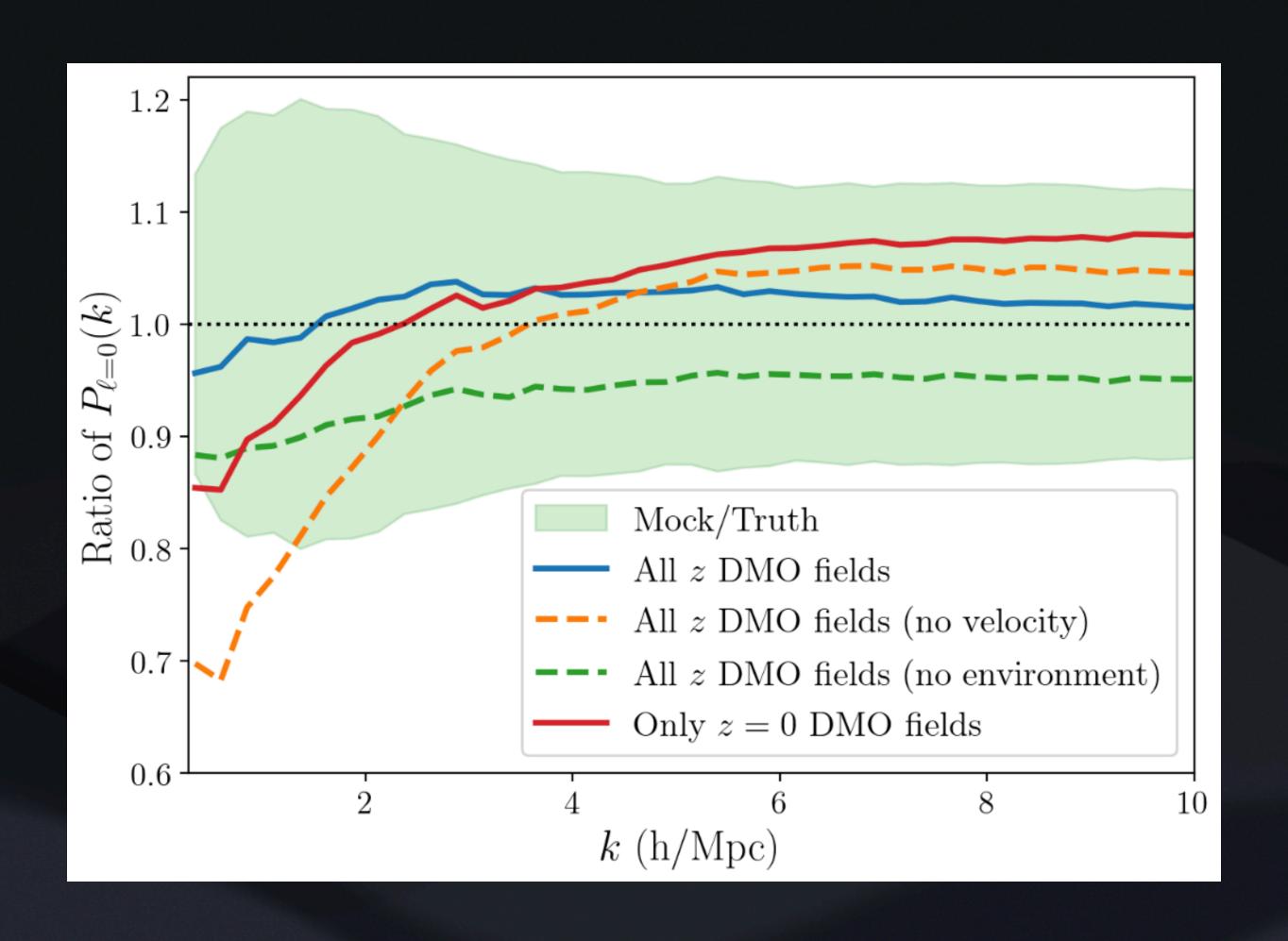
Interpreting the learning

It is nice to know that the computer understands the problem. But I would like to understand it too.

- Eugene Wigner



Interpreting the learning



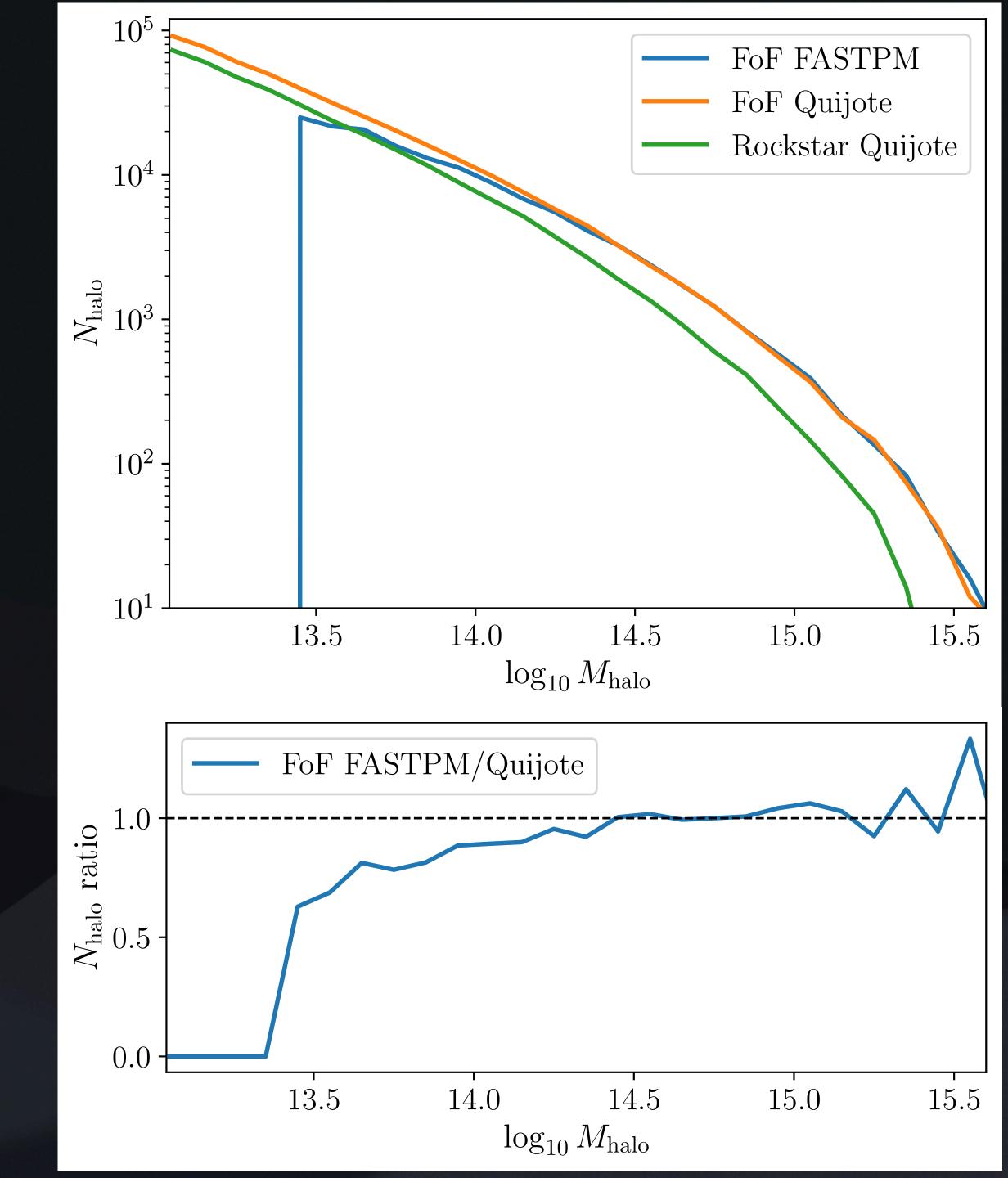
Conclusions

- Accelerated forward models are needed if we wish to extract full cosmological information from current/future surveys.
- The generalizability of transformers to multi-modal inputs can be extended to cosmological simulations as well!
- More work is needed to integrate it into final data pipeline and to interpret the performance of the model.

Extra slides

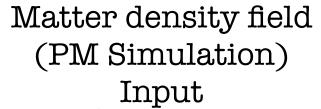
Halo finding in PM

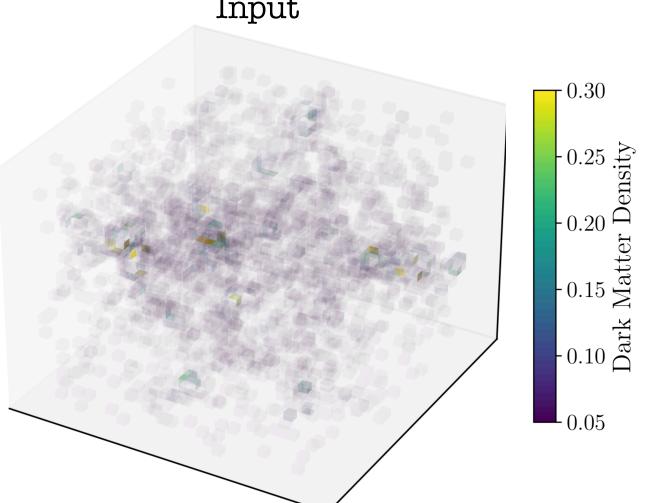
- PM can only reliably find friends-offriends halos of high masses
 - Rockstar halos much more reliable.
 - SDSS/DESI LRG galaxies occupy halos with $M \sim 10^{13} \, M_{\odot}/h$



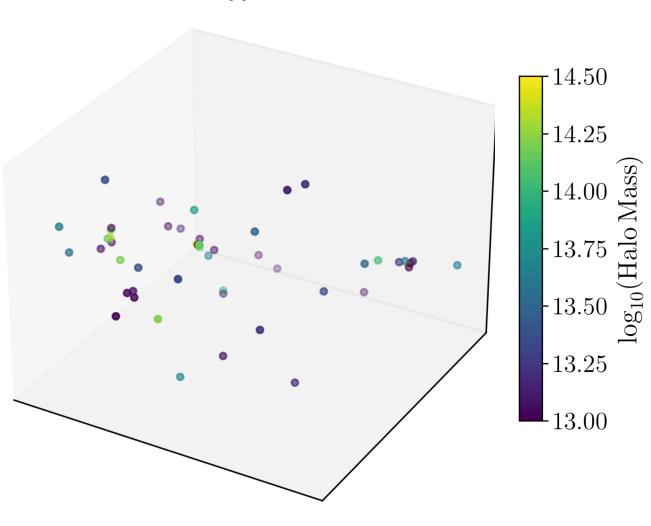
Can we go to small scales?

How else to go from left to right distribution?

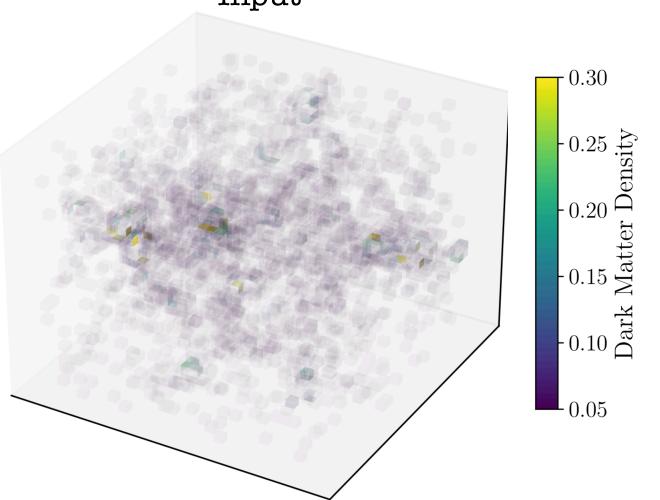




Halo distribution (N-body Simulation) Target



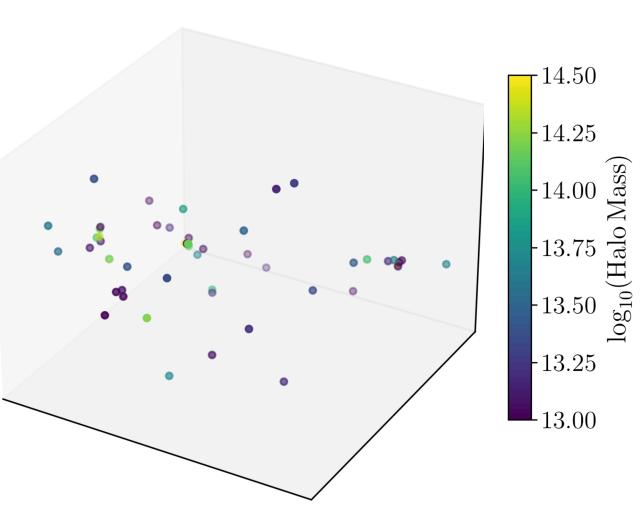
Matter density field (PM Simulation) Input

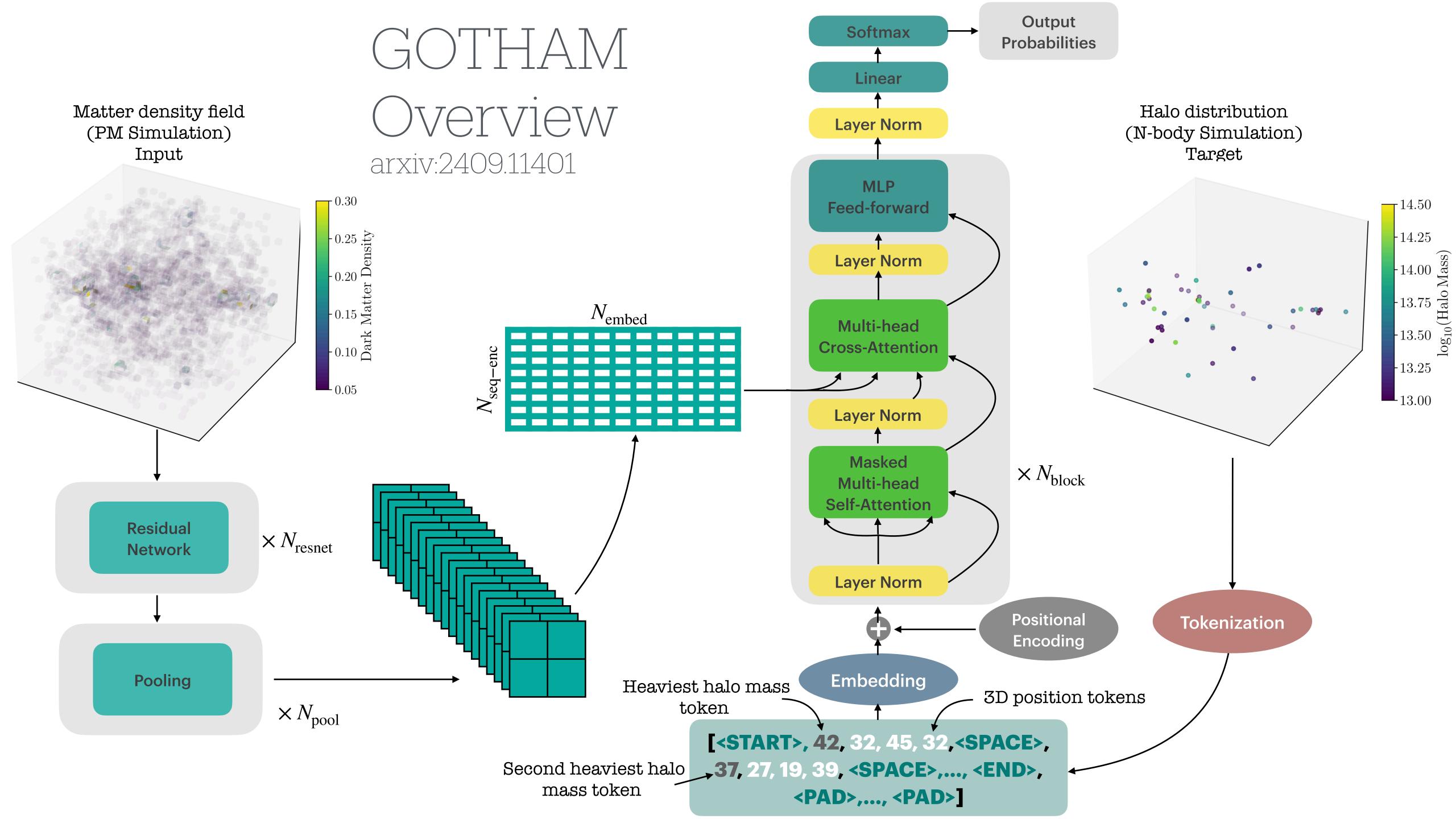


Can treat it as a language translation problem

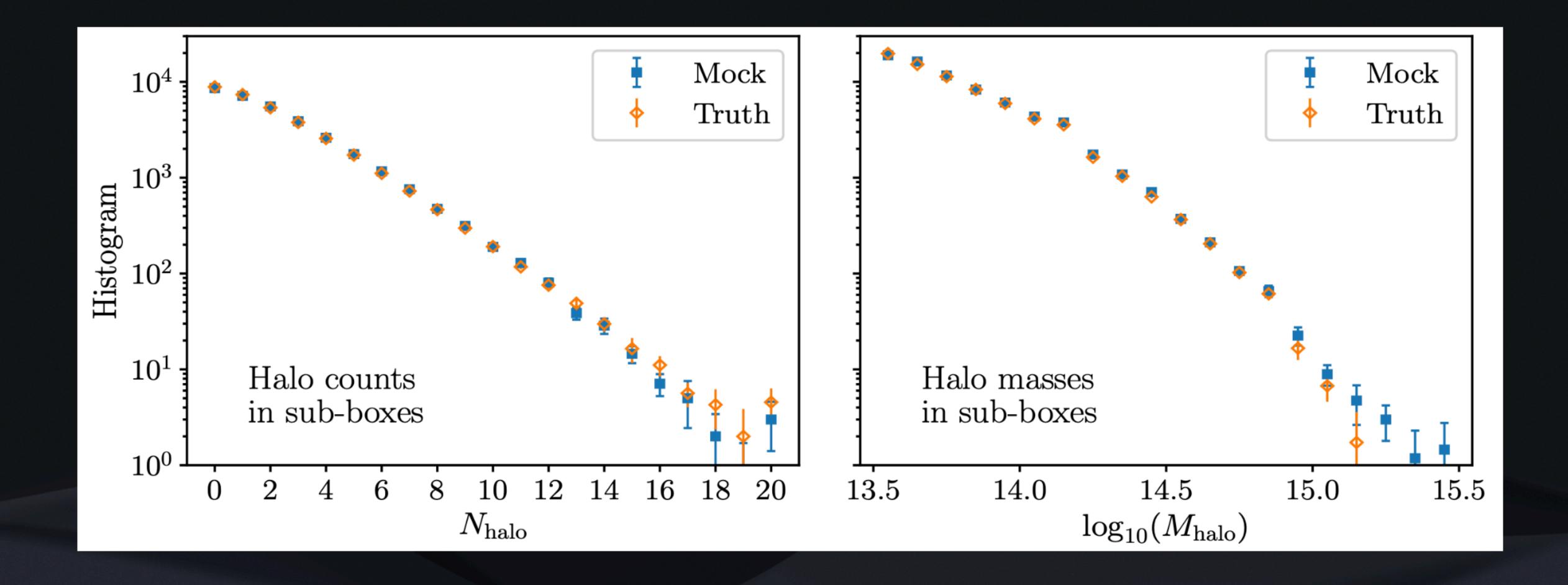


Halo distribution (N-body Simulation) Target

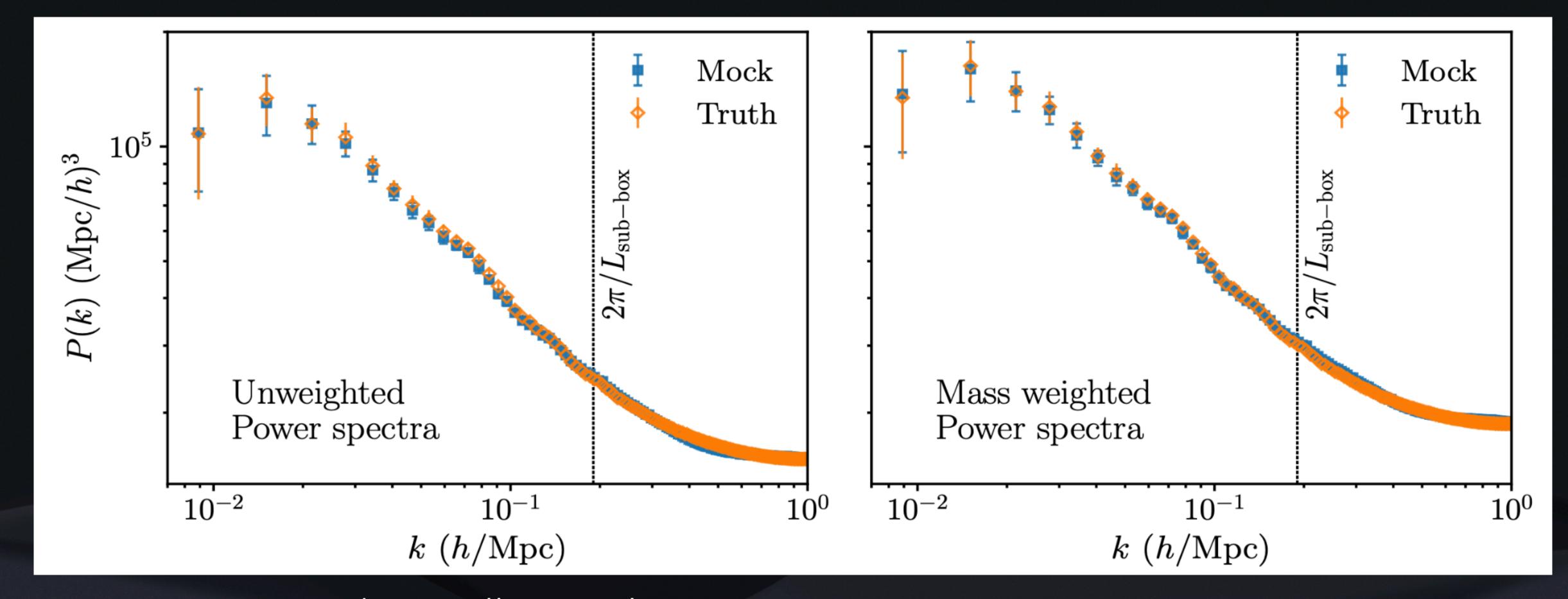




1-pt performance



2-pt performance



Can go to much smaller scales now!