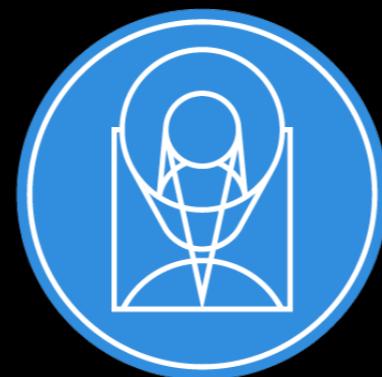


# Applications of Machine Learning in Astrophysics

Sultan Hassan (he/him)  
Space Telescope Science Institute

Online Machine Learning Workshop  
ICTP Physics Without Frontiers - Afghanistan  
October 1, 2025



Every dot is a galaxy!

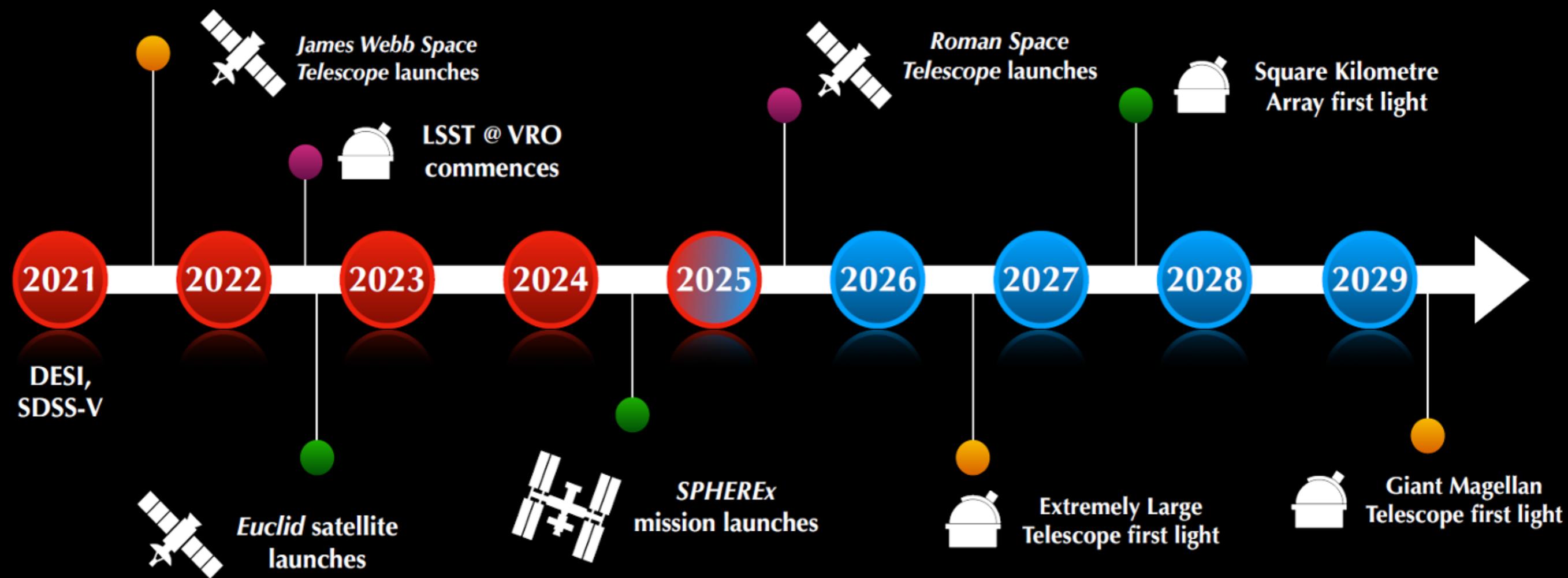
How do we extract hidden information within these images?



credit:<https://science.nasa.gov/asset/webb/webbs-first-deep-field-nircam-image/>

# Why Machine Learning?

New generation of surveys  
We are entering big data era



# What is Machine Learning?



?



# What is Machine Learning?



?



We show lots of data to learn from!



vs



# Same way for Astrophysics

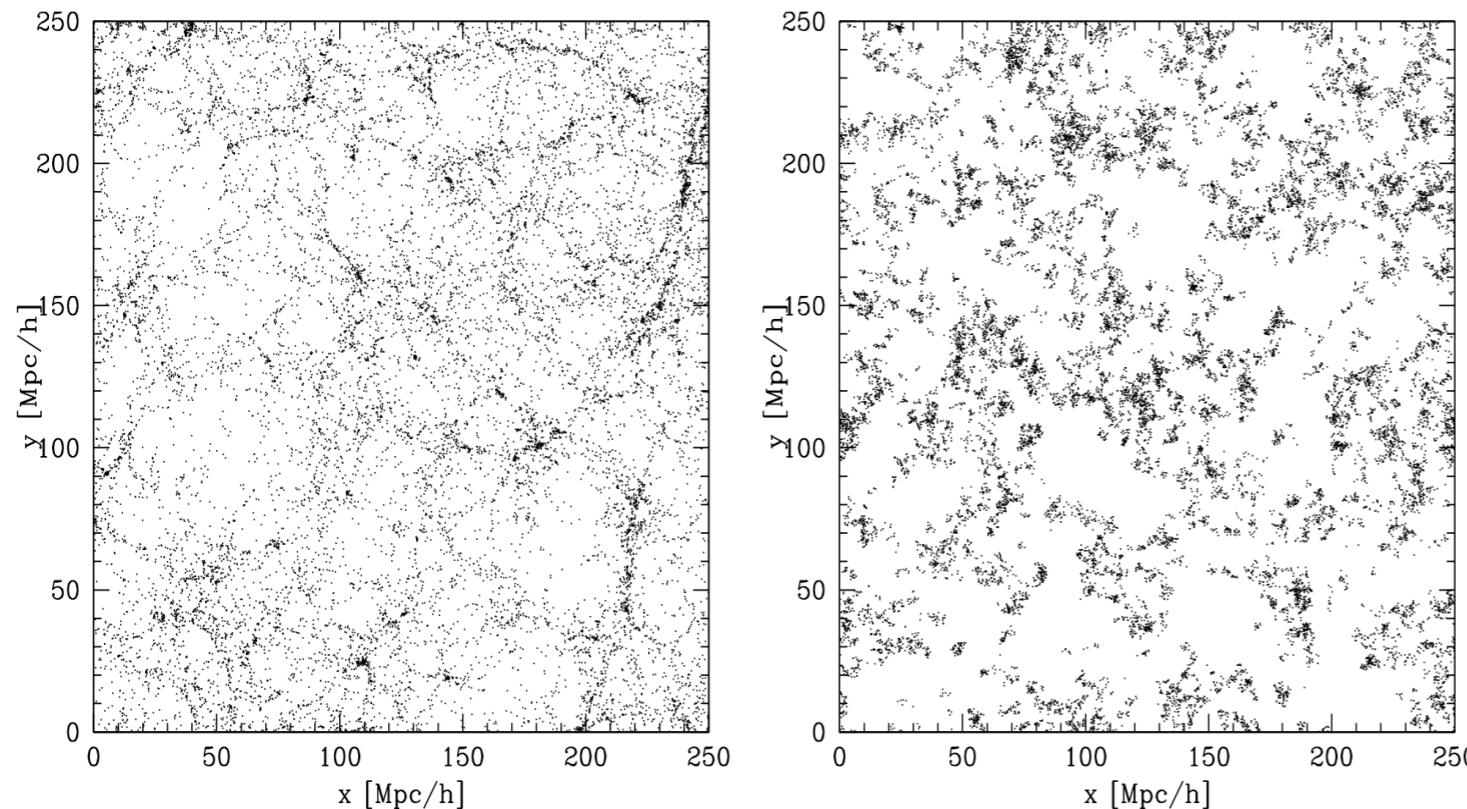
We show lots of galaxies or maps



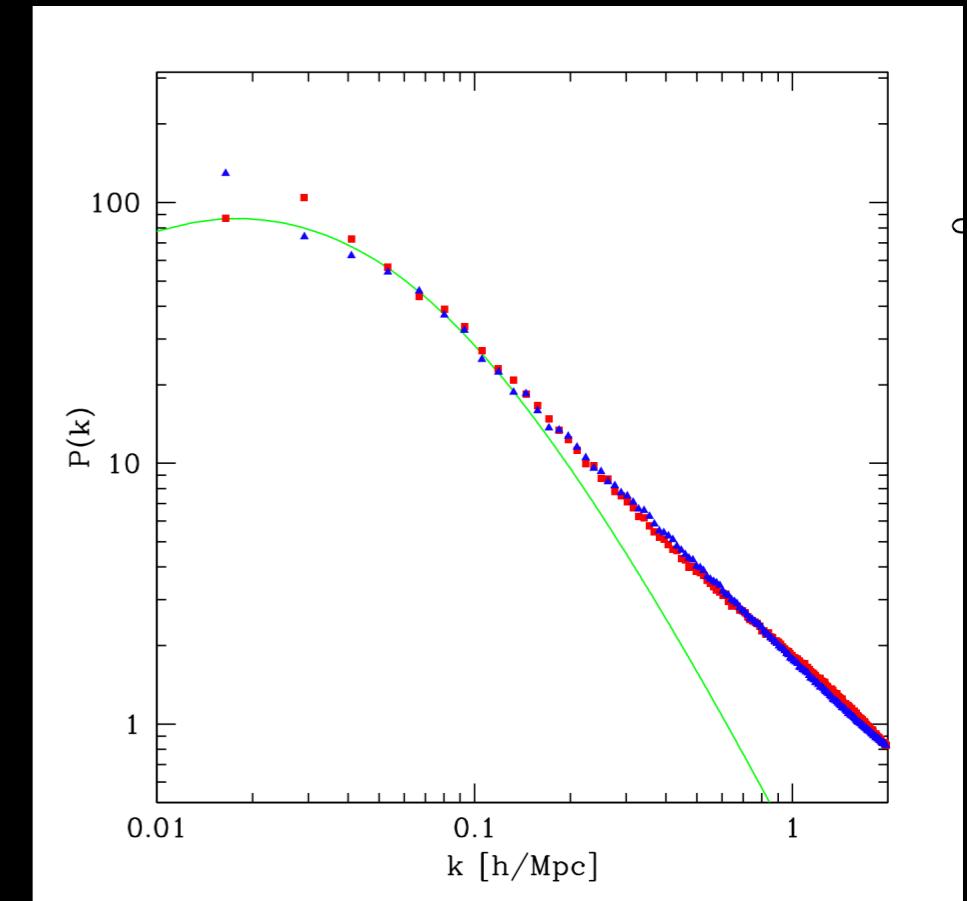
# Traditional methods: Data compression (summary statistics)

## Is this a good way to describe these fields?

Different fields

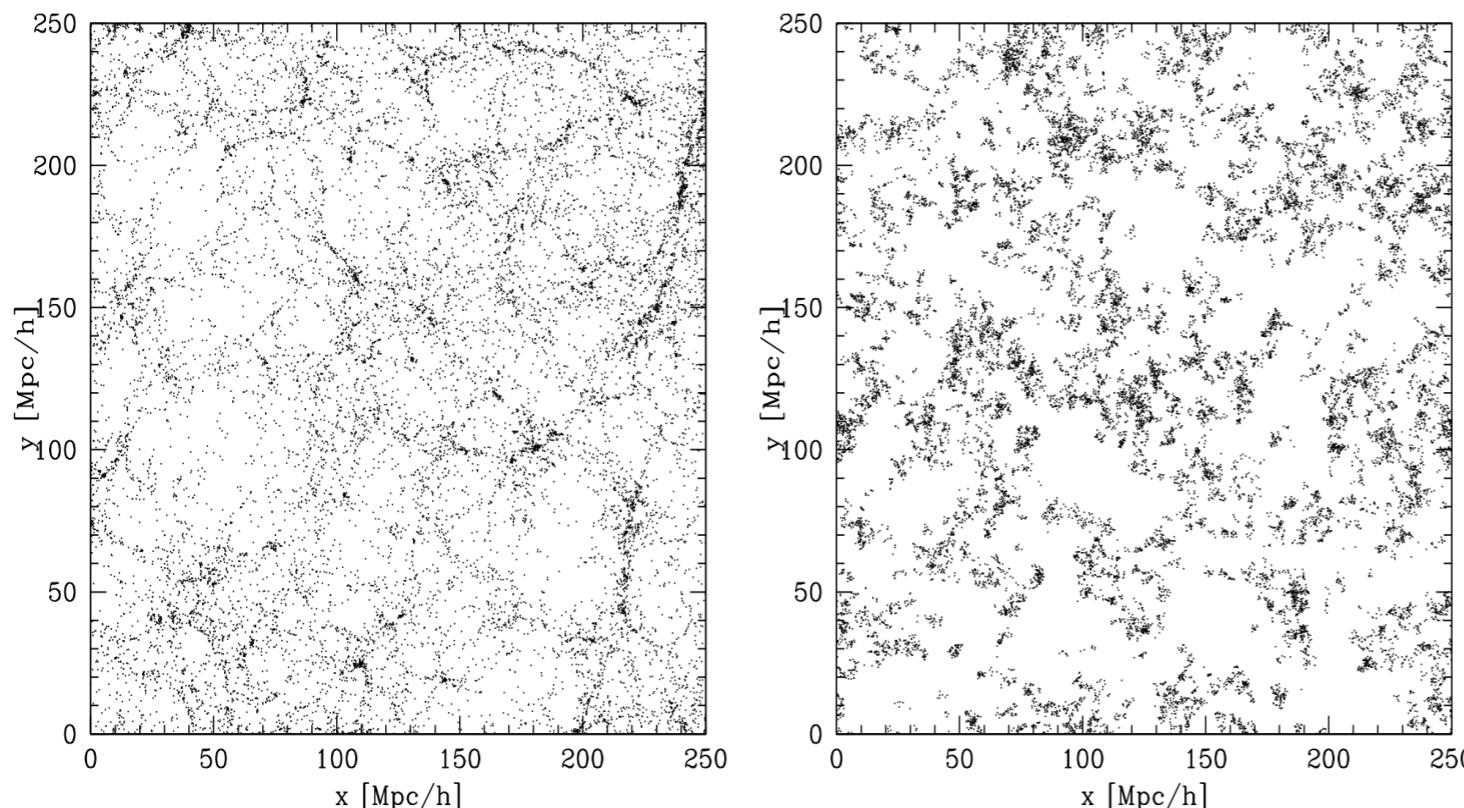


Same power spectrum

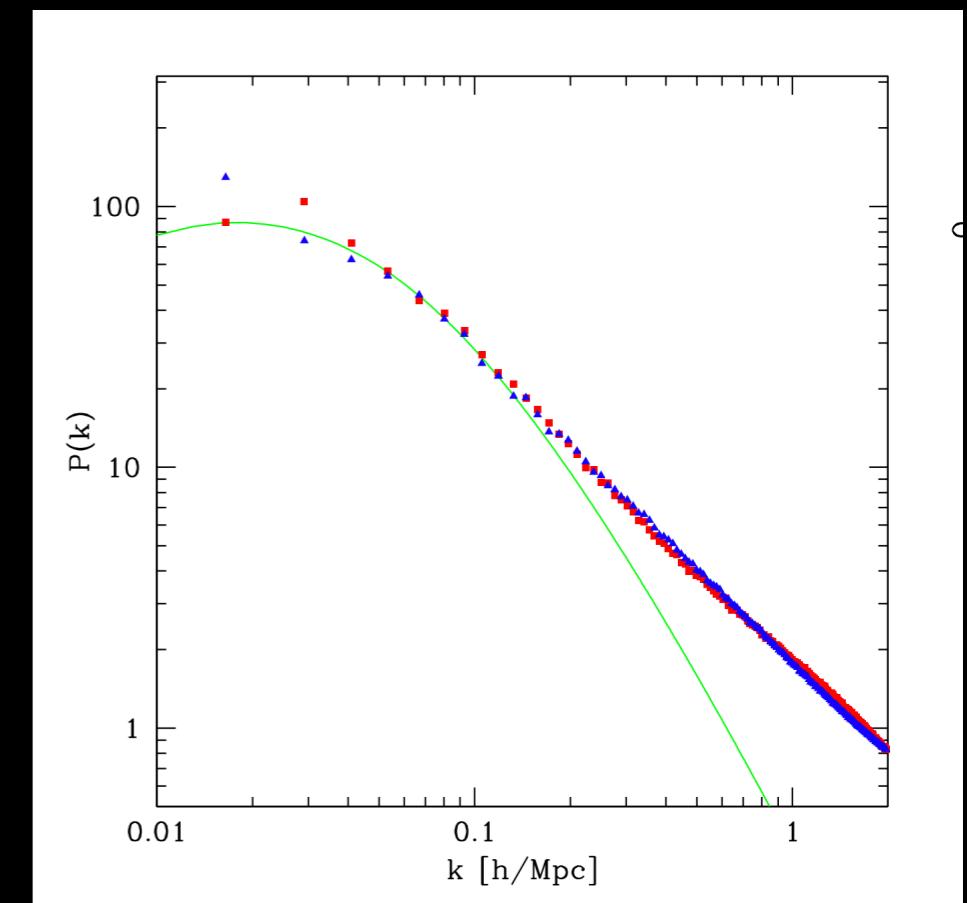


Traditional methods: Data compression (summary statistics)  
Is this a good way to describe these fields?

Different fields



Same power spectrum



We lose lots of information during data compression

Solution:

work directly with these large maps using machine learning

# Two main branches in Machine Learning..

$$p(\text{model} \mid \text{data}) \sim p(\text{data} \mid \text{model}) p(\text{model})$$

# Two main branches in Machine Learning..

Posterior

Likelihood

$$p(\text{model} \mid \text{data}) \sim p(\text{data} \mid \text{model}) p(\text{model})$$

# Two main branches in Machine Learning..

Discriminative



Posterior

Likelihood

$$p(\text{model} \mid \text{data}) \sim p(\text{data} \mid \text{model}) p(\text{model})$$

# Two main branches in Machine Learning..

Discriminative



Posterior

Generative



Likelihood

$$p(\text{model} \mid \text{data}) \sim p(\text{data} \mid \text{model}) p(\text{model})$$

# Two main branches in Machine Learning..

Discriminative



Posterior

Generative



Likelihood

$$p(\text{model} \mid \text{data}) \sim p(\text{data} \mid \text{model}) p(\text{model})$$

The current direction in ML is more towards generative models...

# Two main branches in Machine Learning..

Discriminative

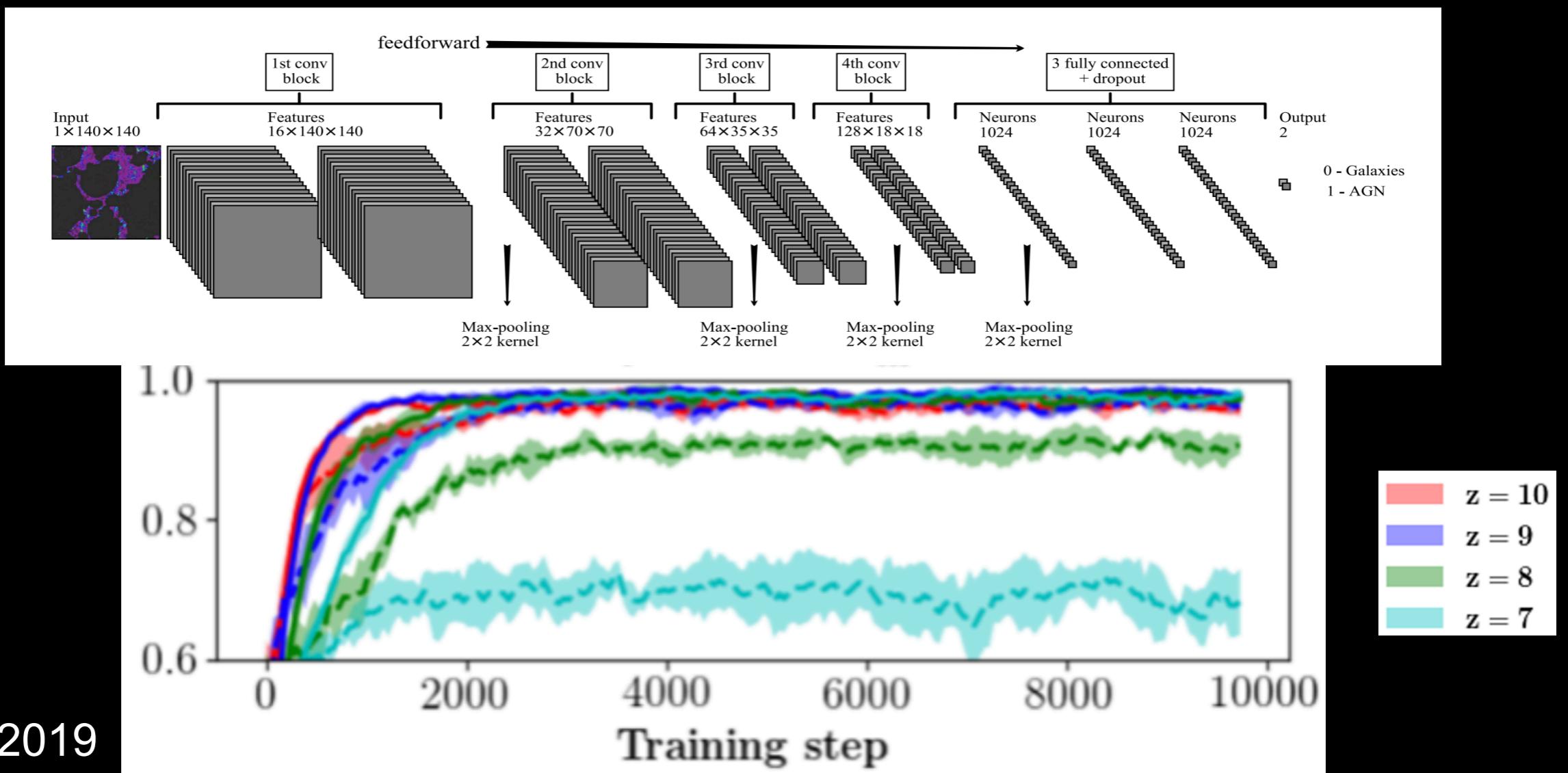
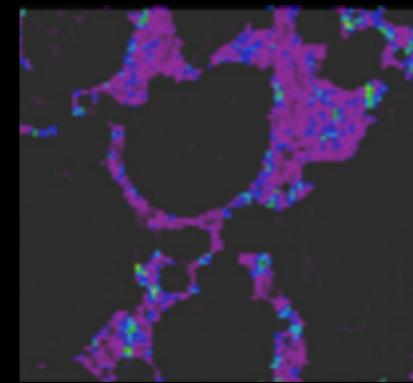


Posterior

$$p(\text{model} \mid \text{data}) \sim p(\text{data} \mid \text{model}) p(\text{model})$$

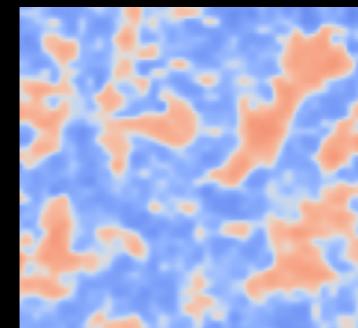
Convolutional Neural Networks (CNNs) achieve the state-of-the-art performance in extracting information (e.g. recovering parameters) from fields.

$$p(\text{AGN or Galaxies} \mid \text{Image})$$

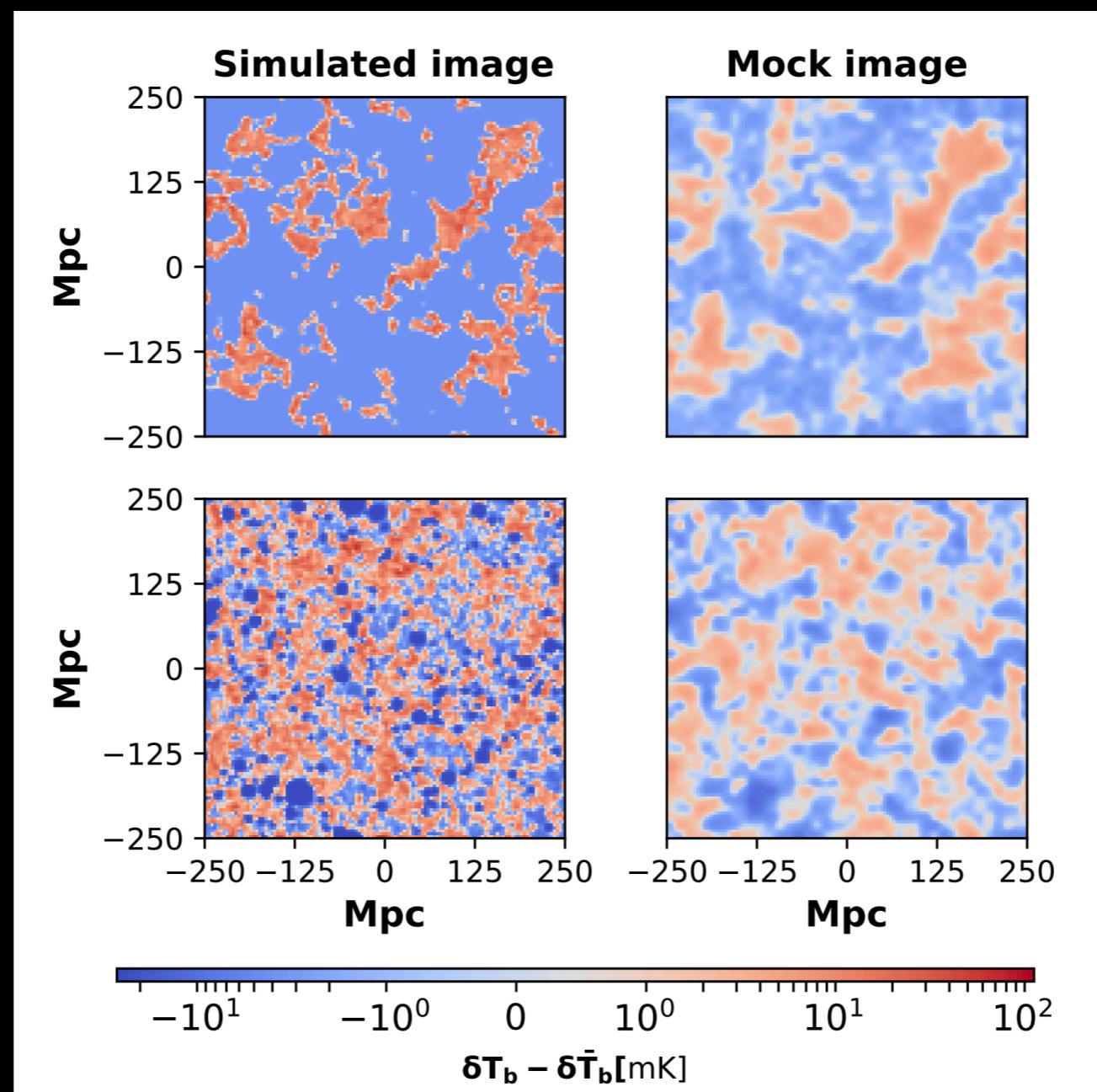


# Constraining the Reionization History with CNNs

$$p(x_{\text{HI}} | \quad )$$



Tumelo Mangena, SKA office  
South Africa

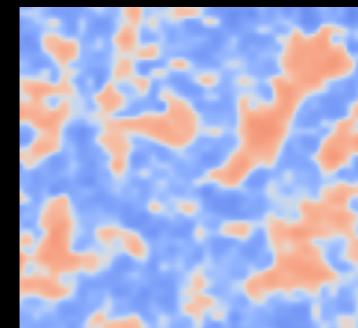


# Constraining the Reionization History with CNNs



Tumelo Mangena, SKA office  
South Africa

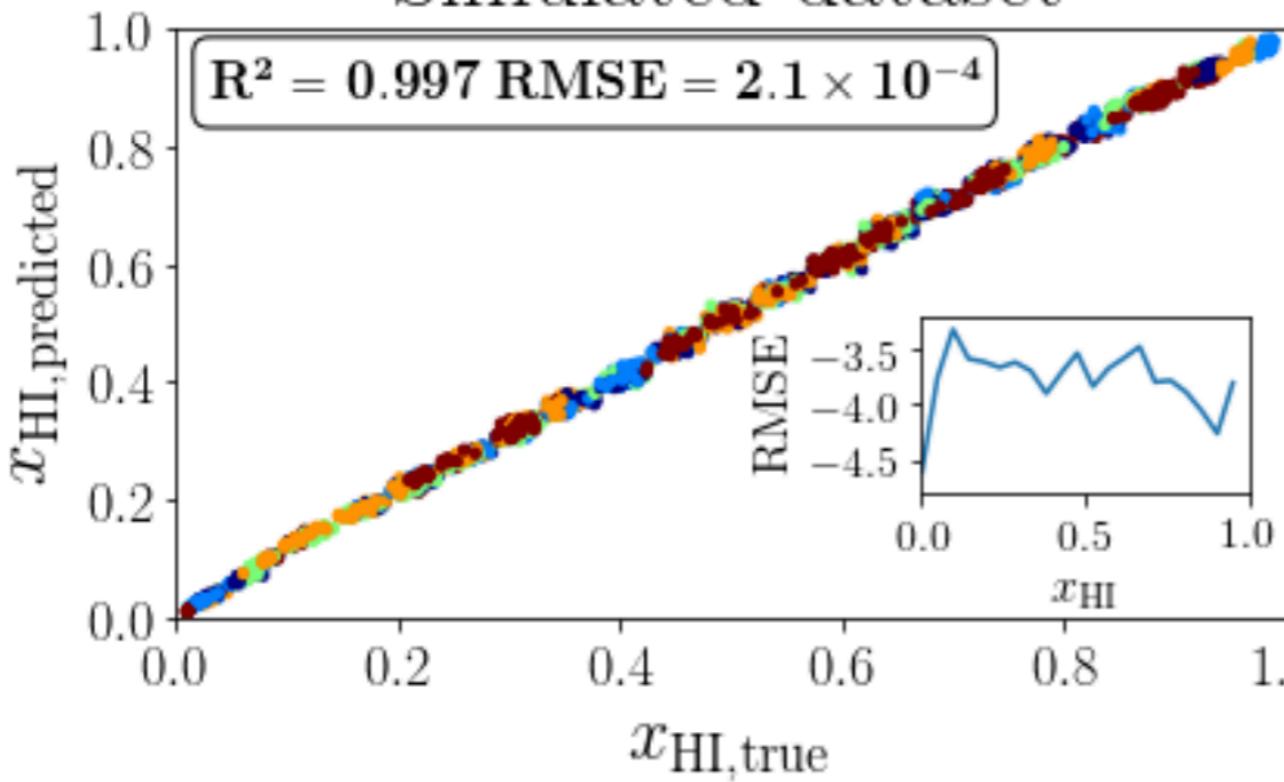
$$p(x_{\text{HI}} | \text{Simulated image} )$$



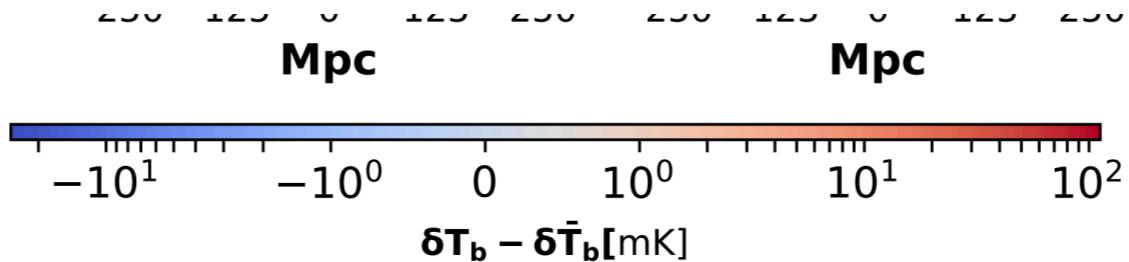
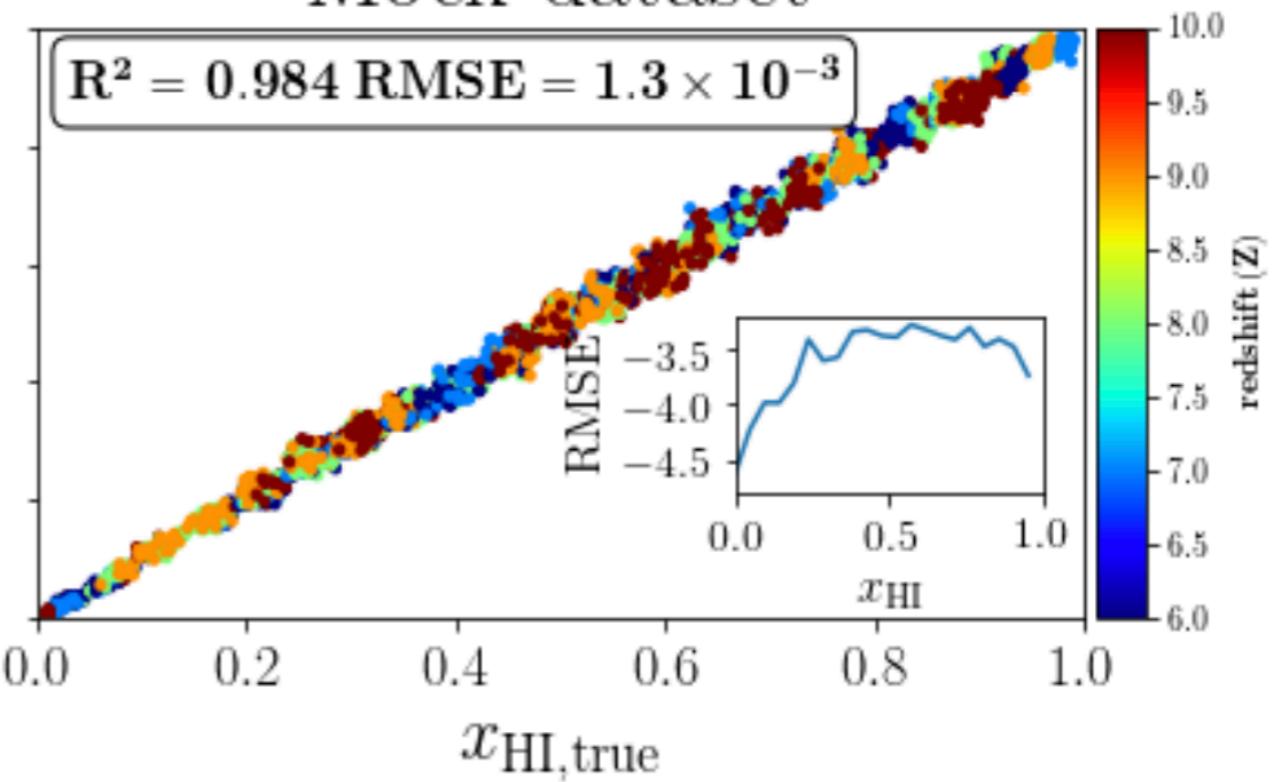
Simulated image

Mock image

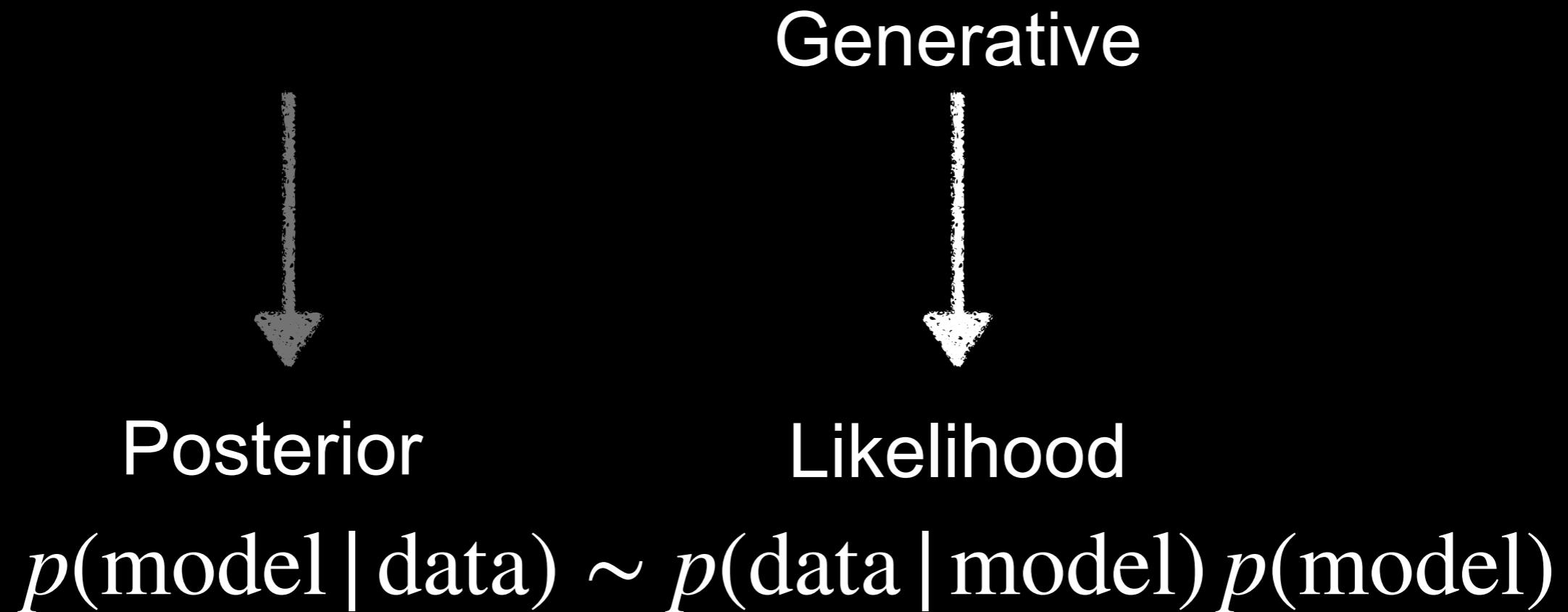
Simulated dataset



Mock dataset



# Two main branches in Machine Learning..



The current direction in ML is more towards generative models...

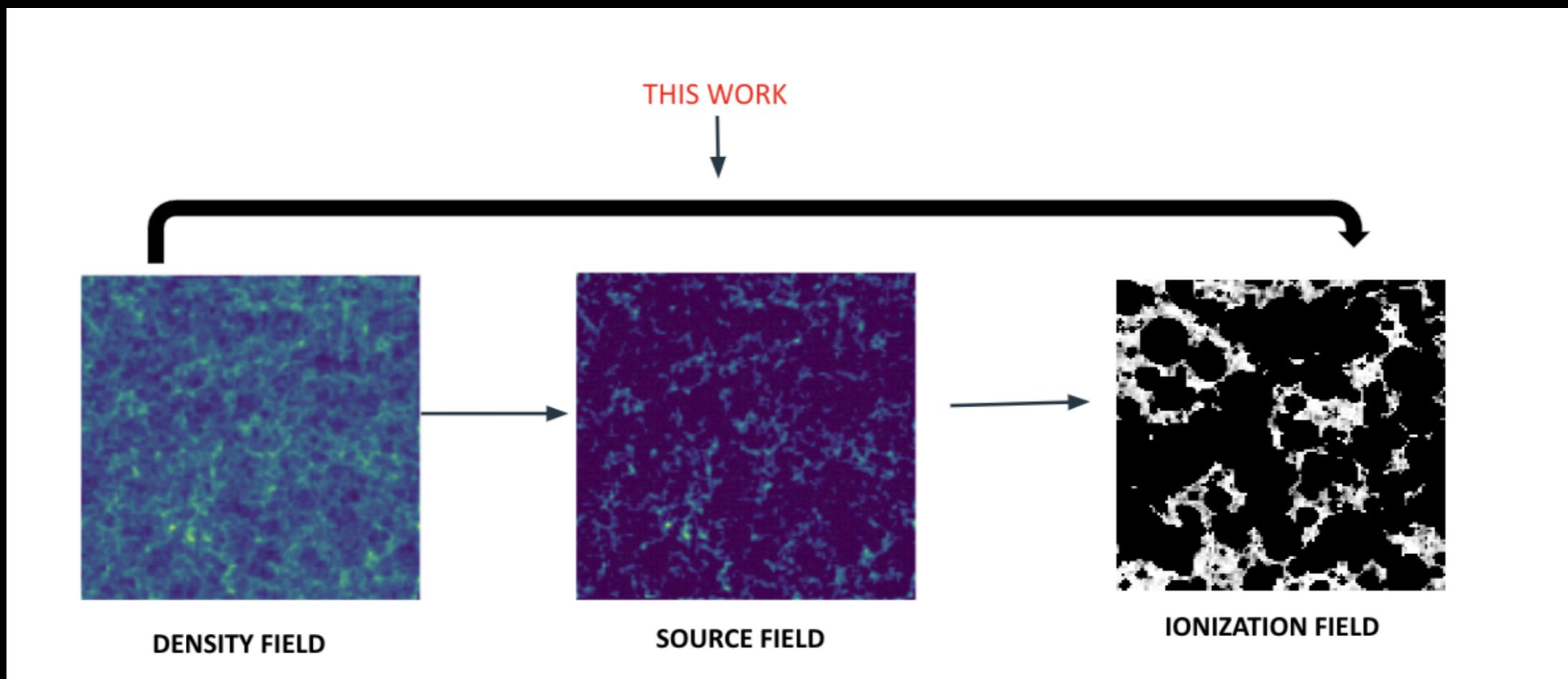
# Emulating radiation transport on large scales using denoising U-Nets

(Masipa, Hassan, Santos, Contardo, Cho 2023, accepted to ICLR 2023, arXiv:2303.12065 )



Mosima Masipa,  
PhD student  
University of the Western  
Cape, South Africa

$$p( \text{IONIZATION FIELD} \mid \text{DENSITY FIELD} )$$



# Emulating radiation transport on large scales using denoising U-Nets

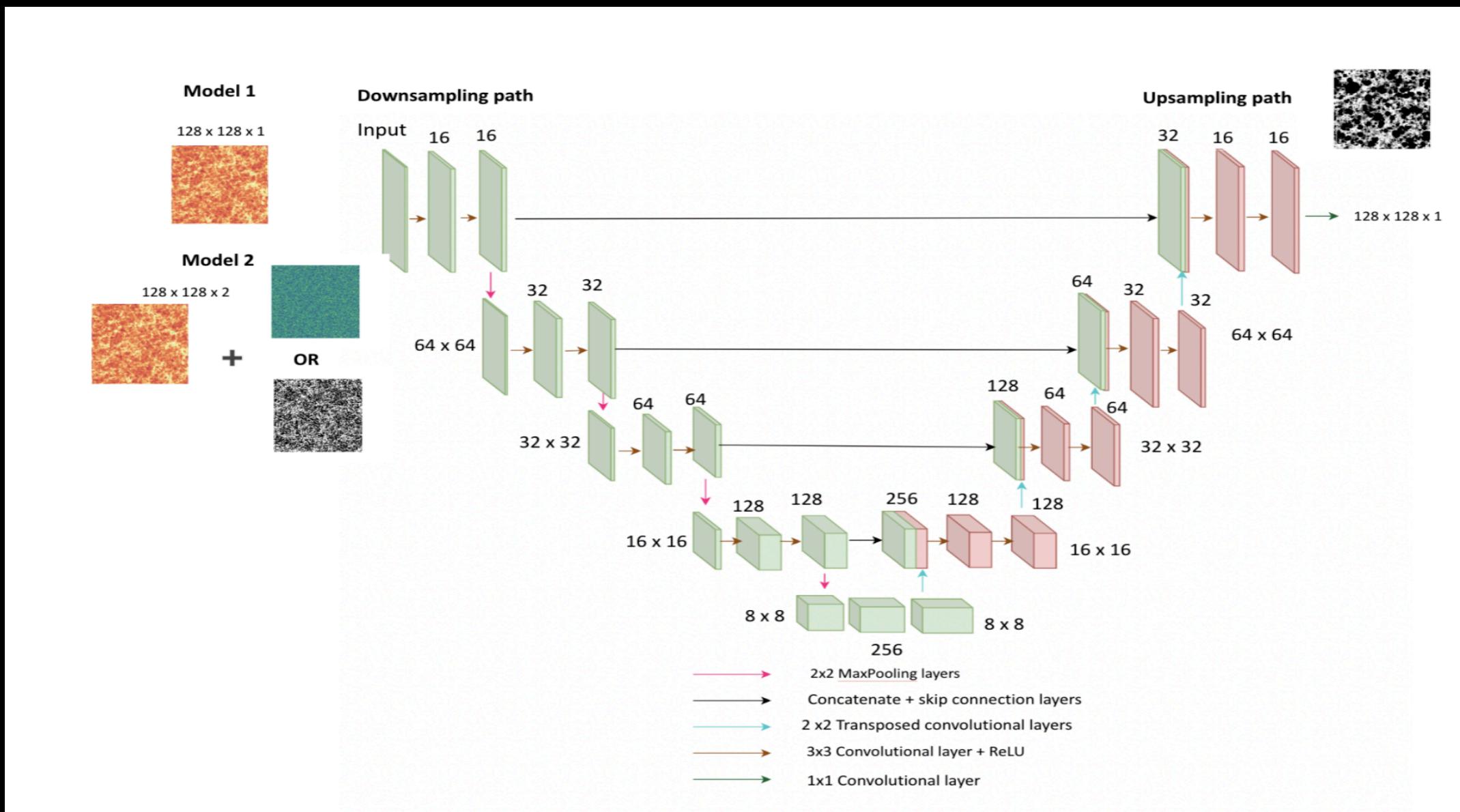
(Masipa, Hassan, Santos, Contardo, Cho 2023, accepted to ICLR 2023, arXiv:2303.12065 )



Mosima Masipa,

MSc student

University of the Western  
Cape, South Africa

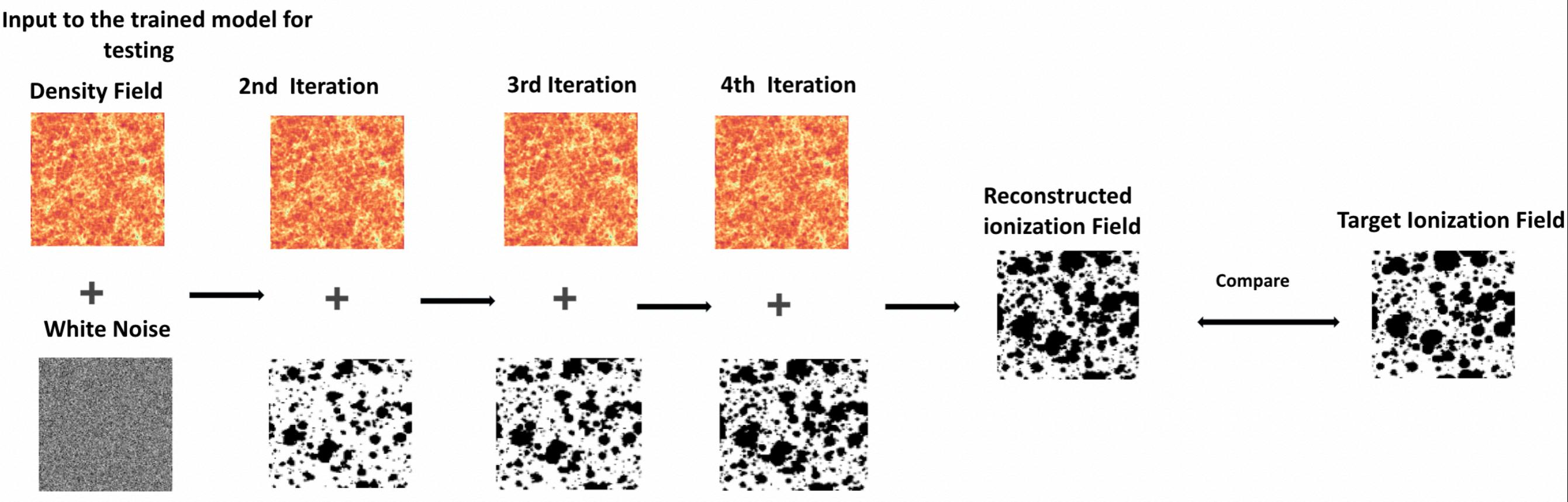


# Emulating radiation transport on large scales using denoising U-Nets

(Masipa, Hassan, Santos, Contardo, Cho 2023, accepted to ICLR 2023, arXiv:2303.12065 )



Mosima Masipa,  
MSc student  
University of the Western  
Cape, South Africa



# Emulating radiation transport on large scales using denoising U-Nets

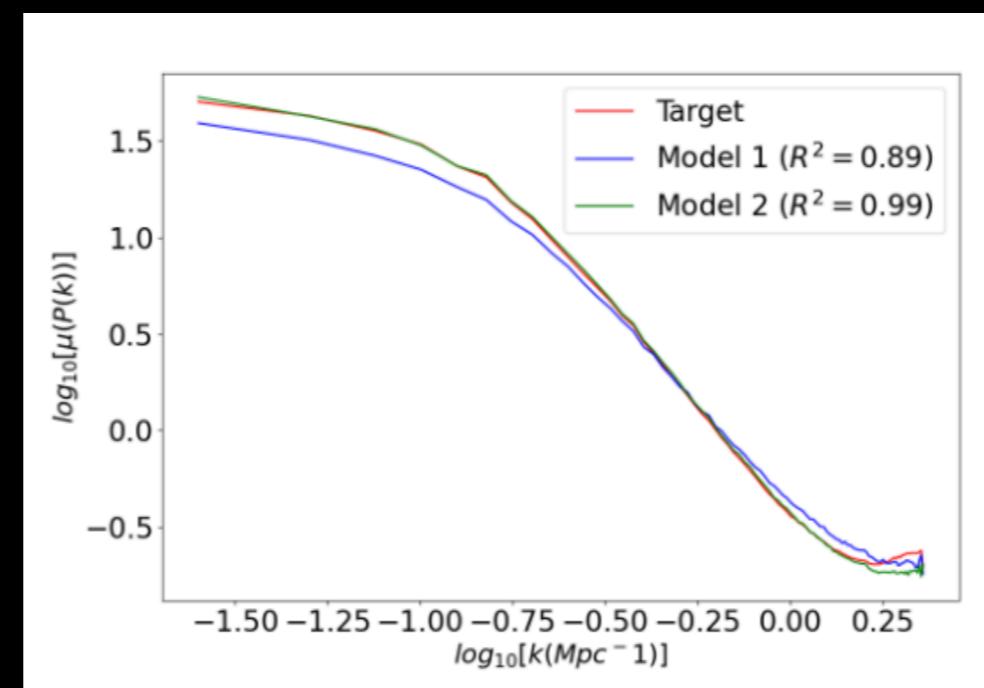
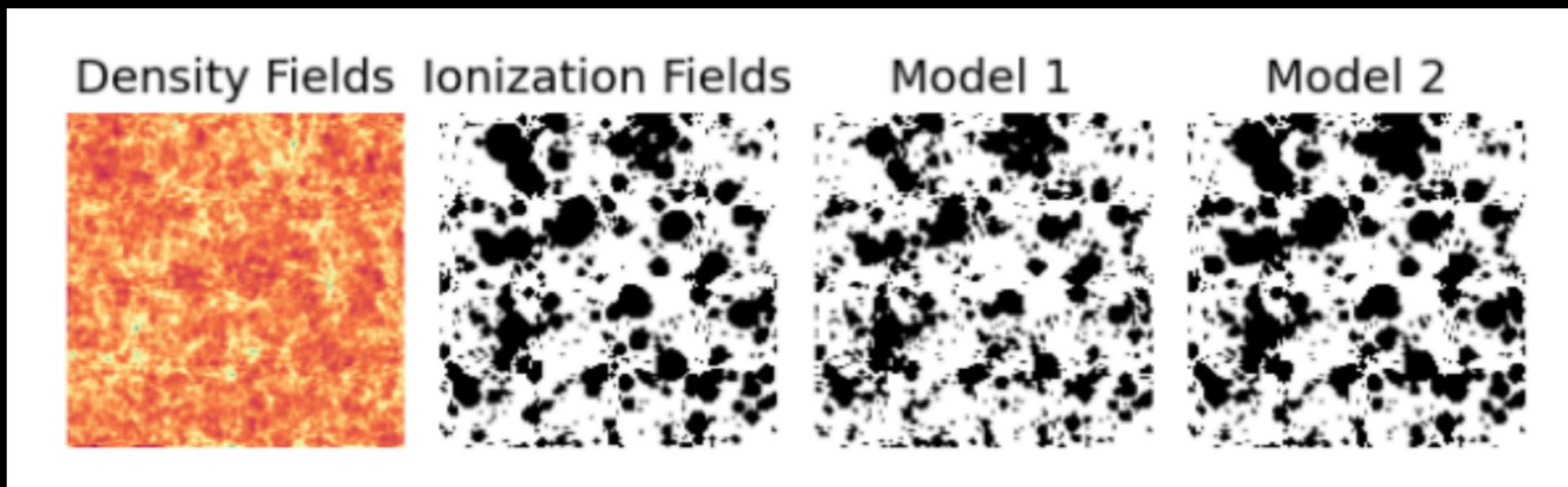
(Masipa, Hassan, Santos, Contardo, Cho 2023, accepted to ICLR 2023, arXiv:2303.12065 )



Mosima Masipa,

MSc student

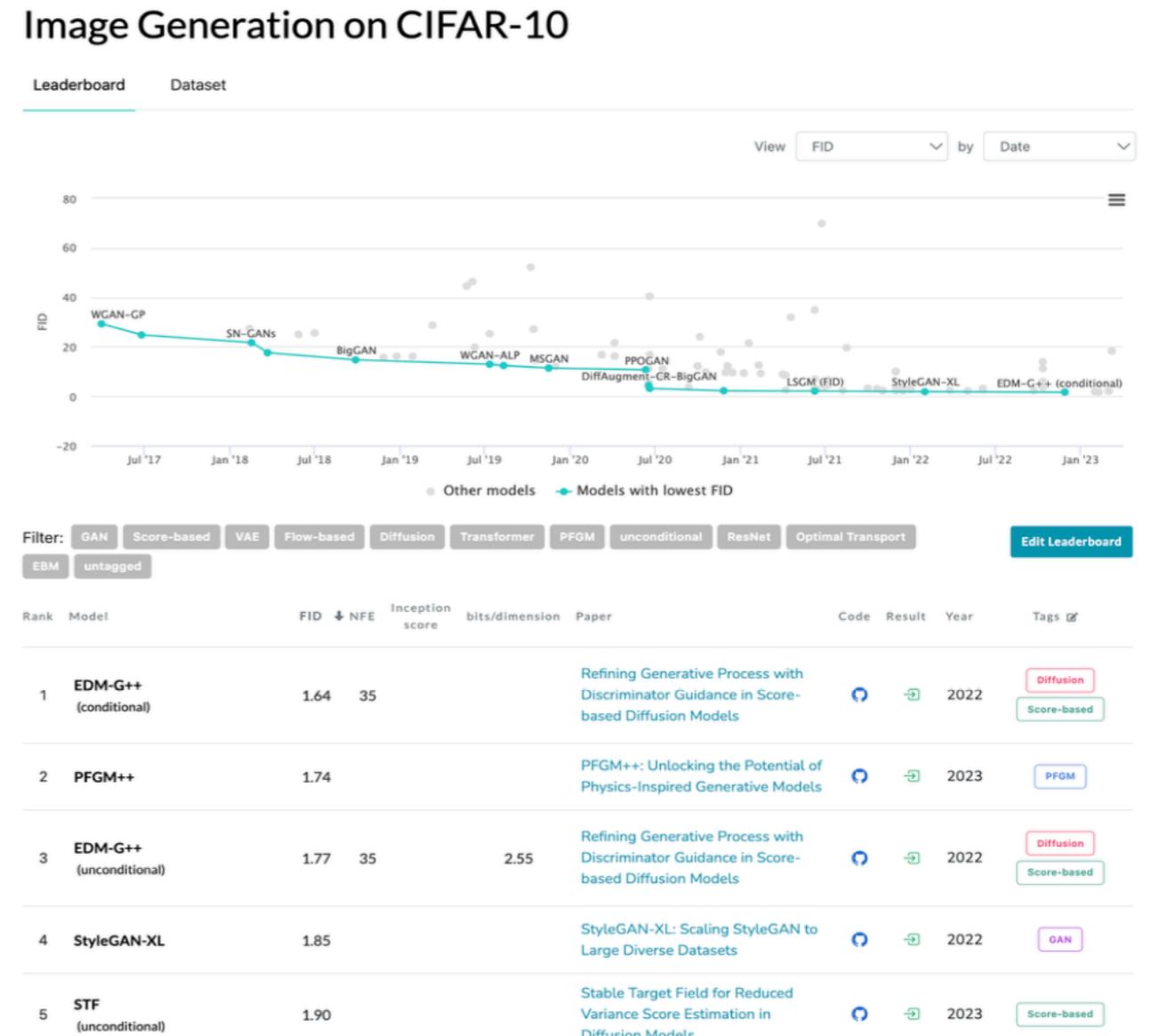
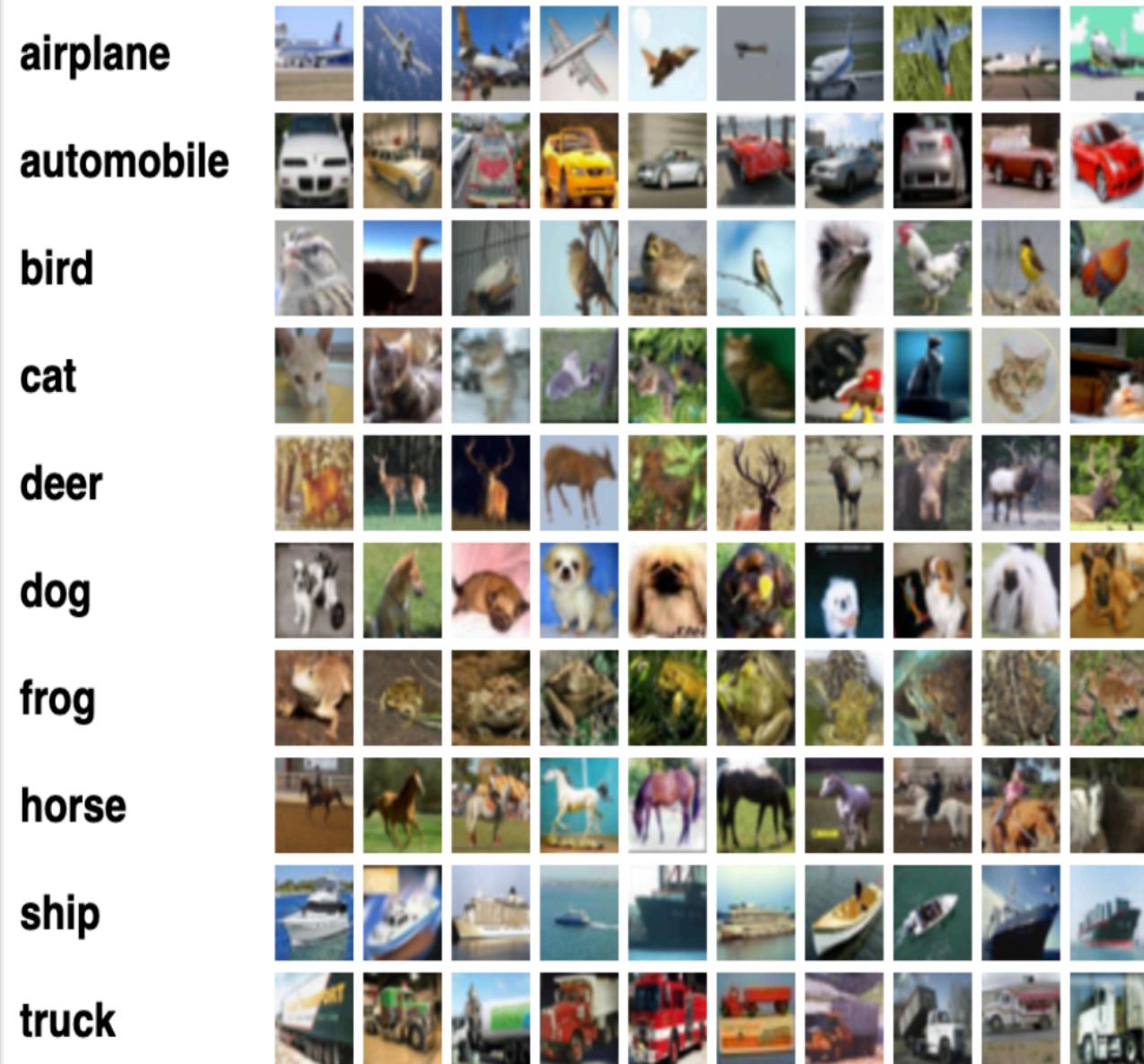
University of the Western  
Cape, South Africa



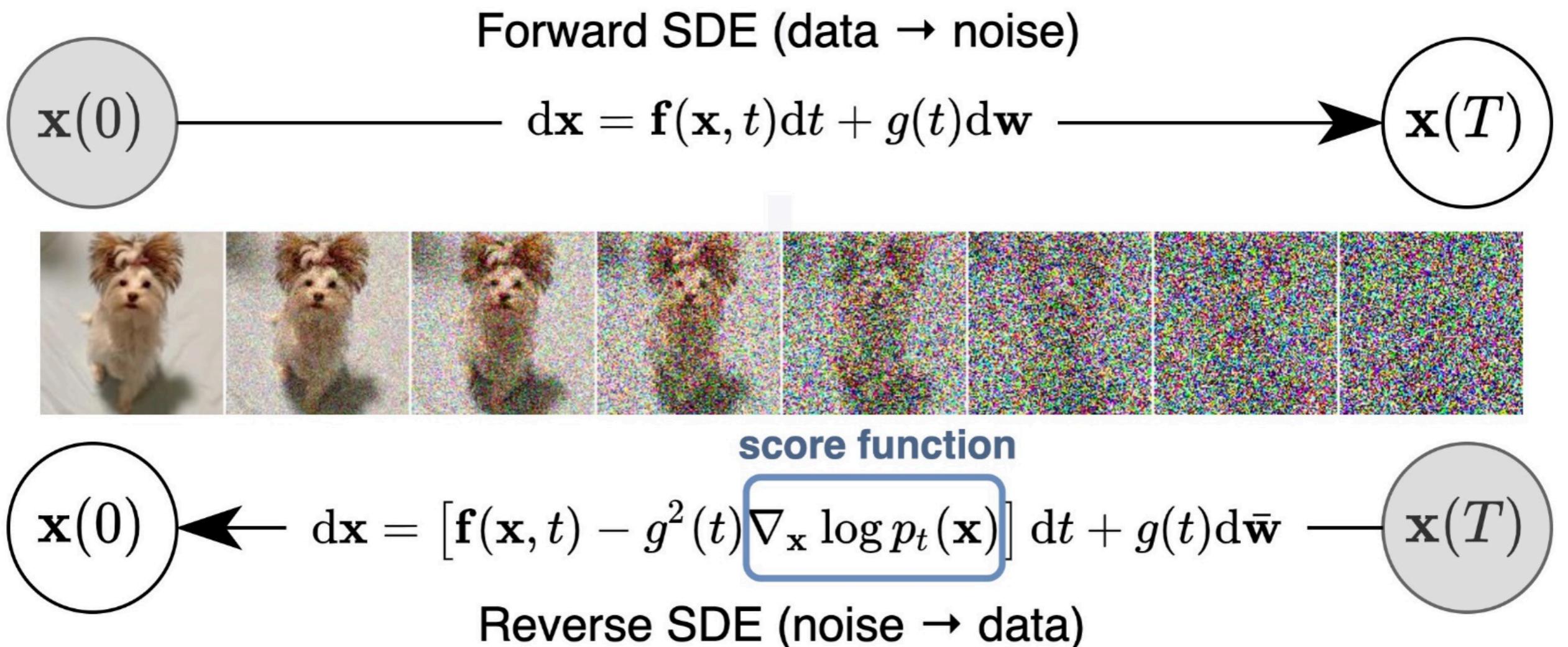
# Why diffusion models?

On CIFAR-10 alone:

- Best ever model in image generation.
- 3 out the first 5.
- 5 out of the first 10.

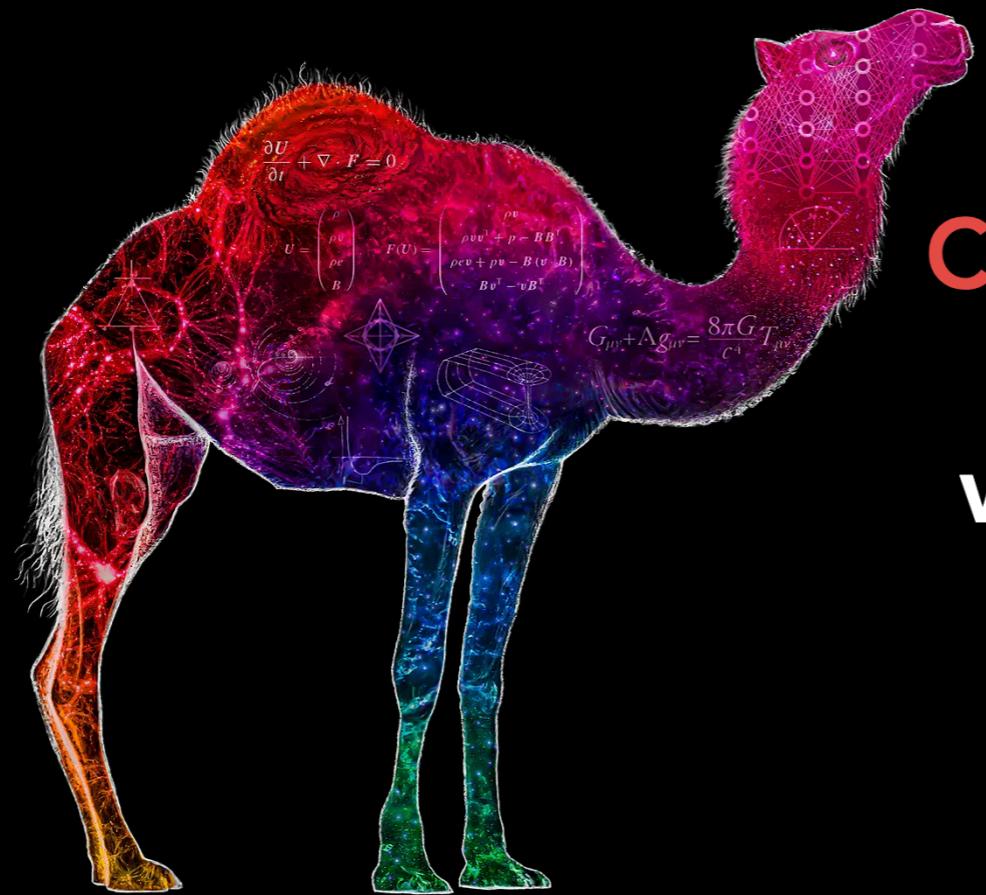


# Diffusion models



# Introducing the CAMELS dataset

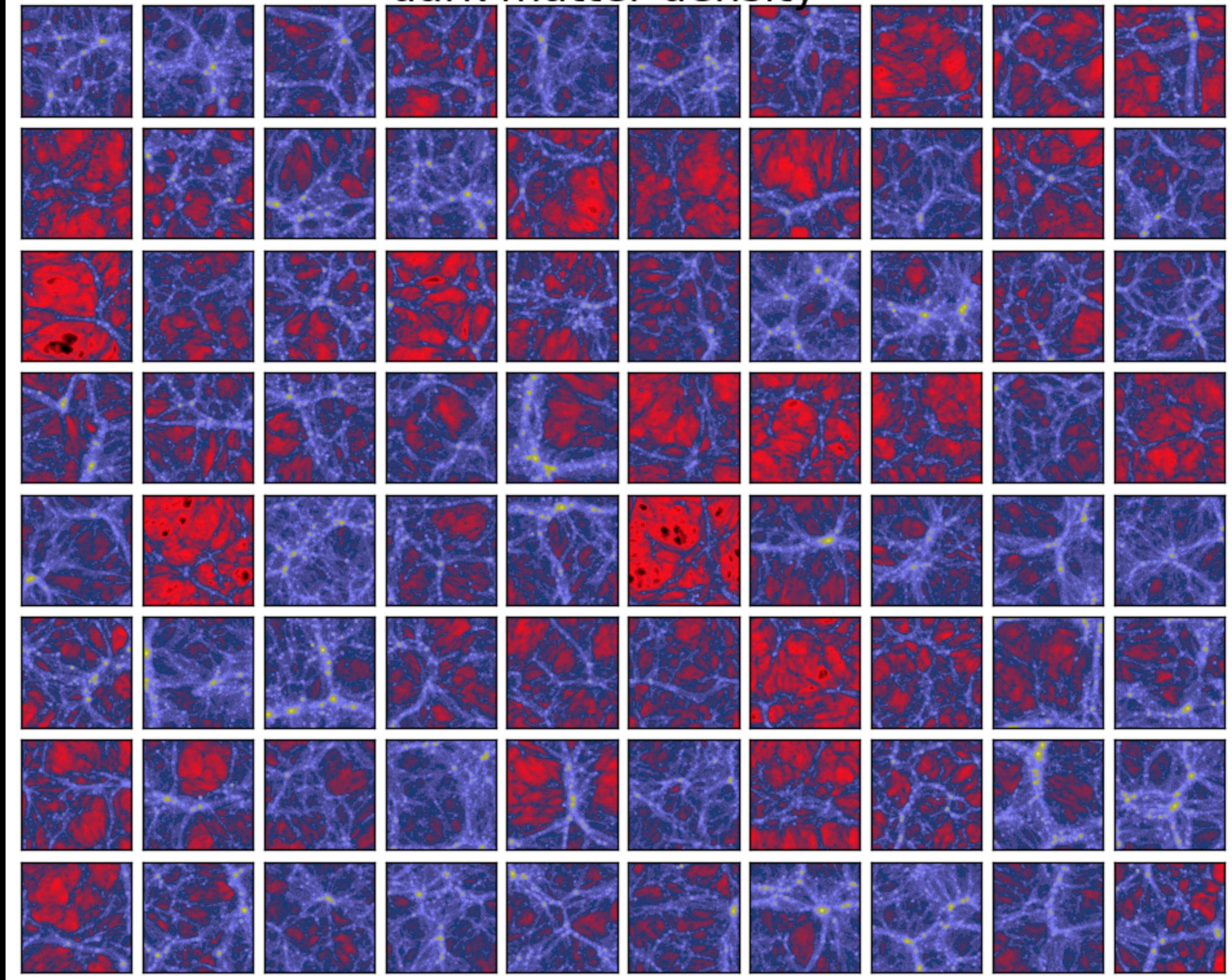
# CAMELS



**Cosmology and  
Astrophysics  
with MachinE  
Learning  
Simulations**

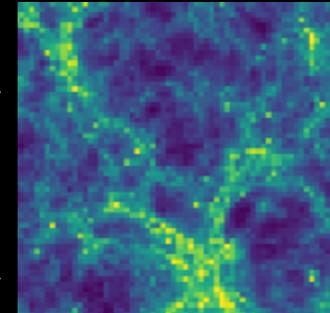
> 4,000 simulations run with *Simba* and *IllustrisTNG*  
designed to train machine learning algorithms

# dark matter density



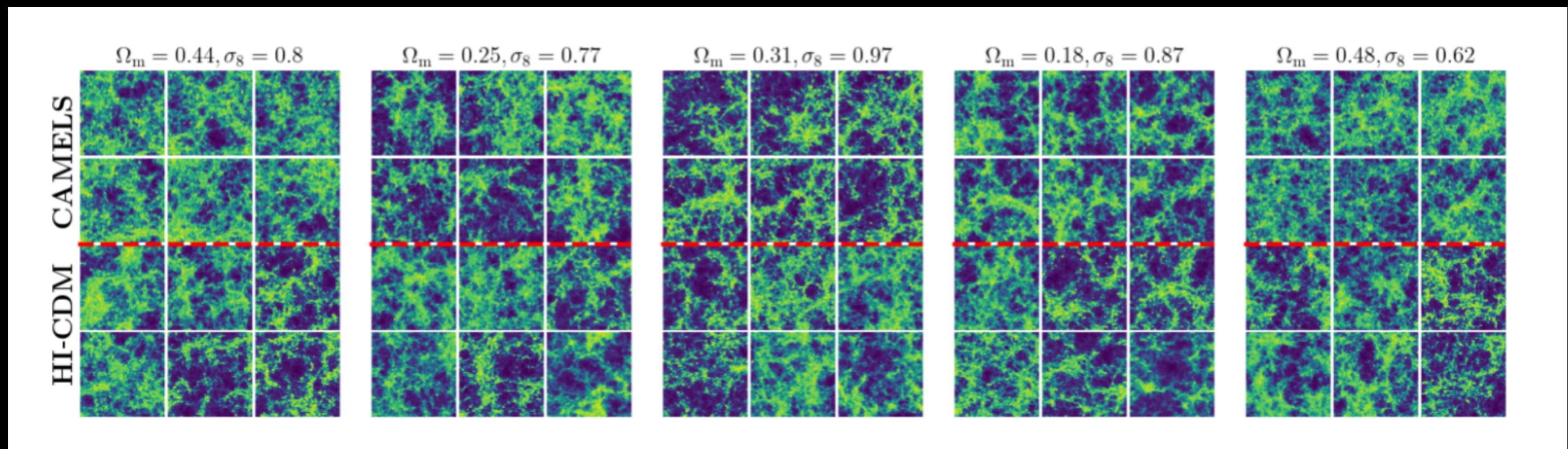
HI-CDM: Conditional diffusion  
model for high fidelity HI maps

$$p(\text{HI map} \mid \Omega_m, \sigma_8)$$



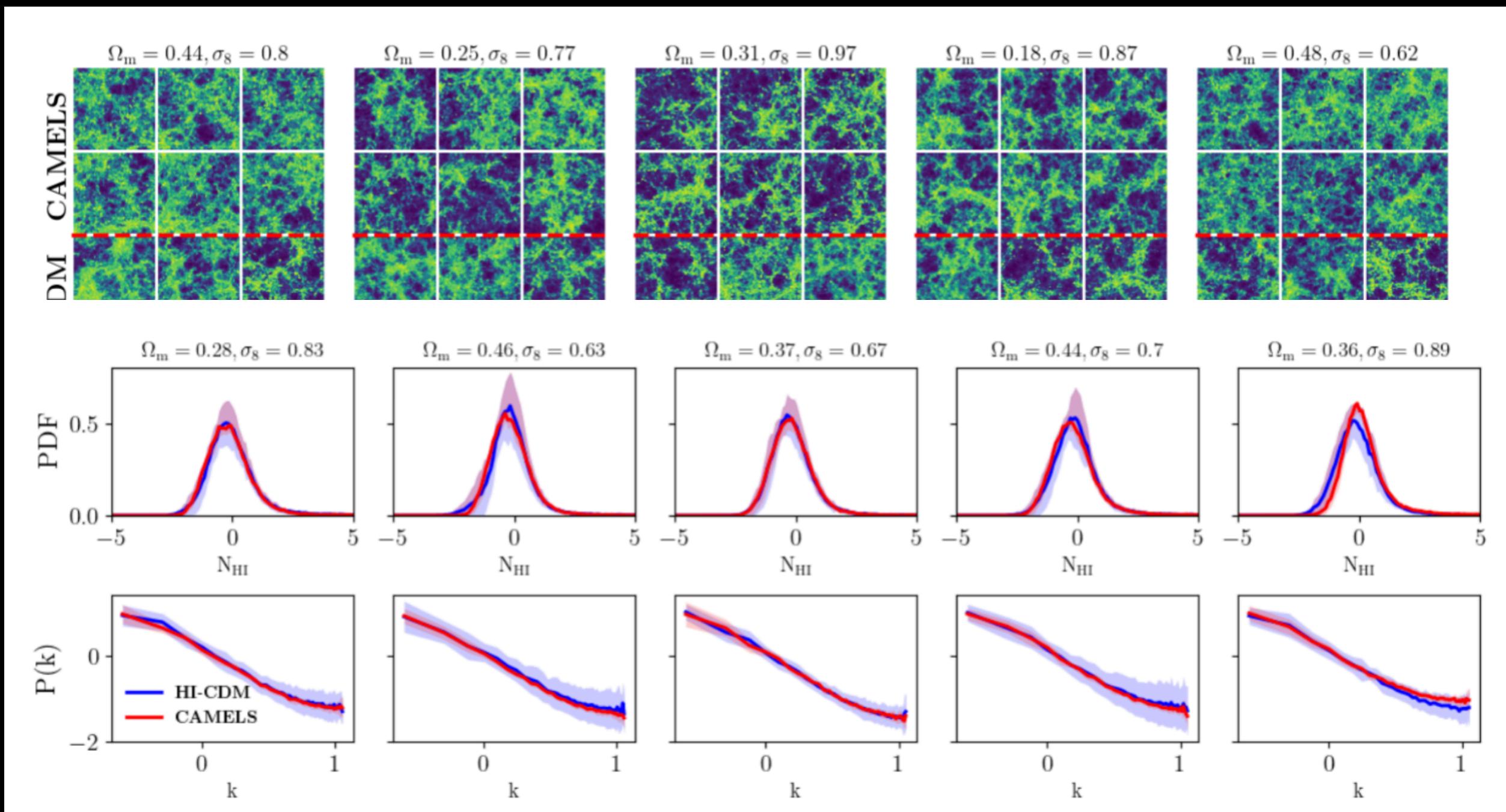
# HI-CDM: Conditional diffusion model for high fidelity HI maps

$$p(\text{HI map} \mid \Omega_m, \sigma_8)$$



# HI-CDM: Conditional diffusion model for high fidelity HI maps

$$p(\text{HI map} \mid \Omega_m, \sigma_8)$$

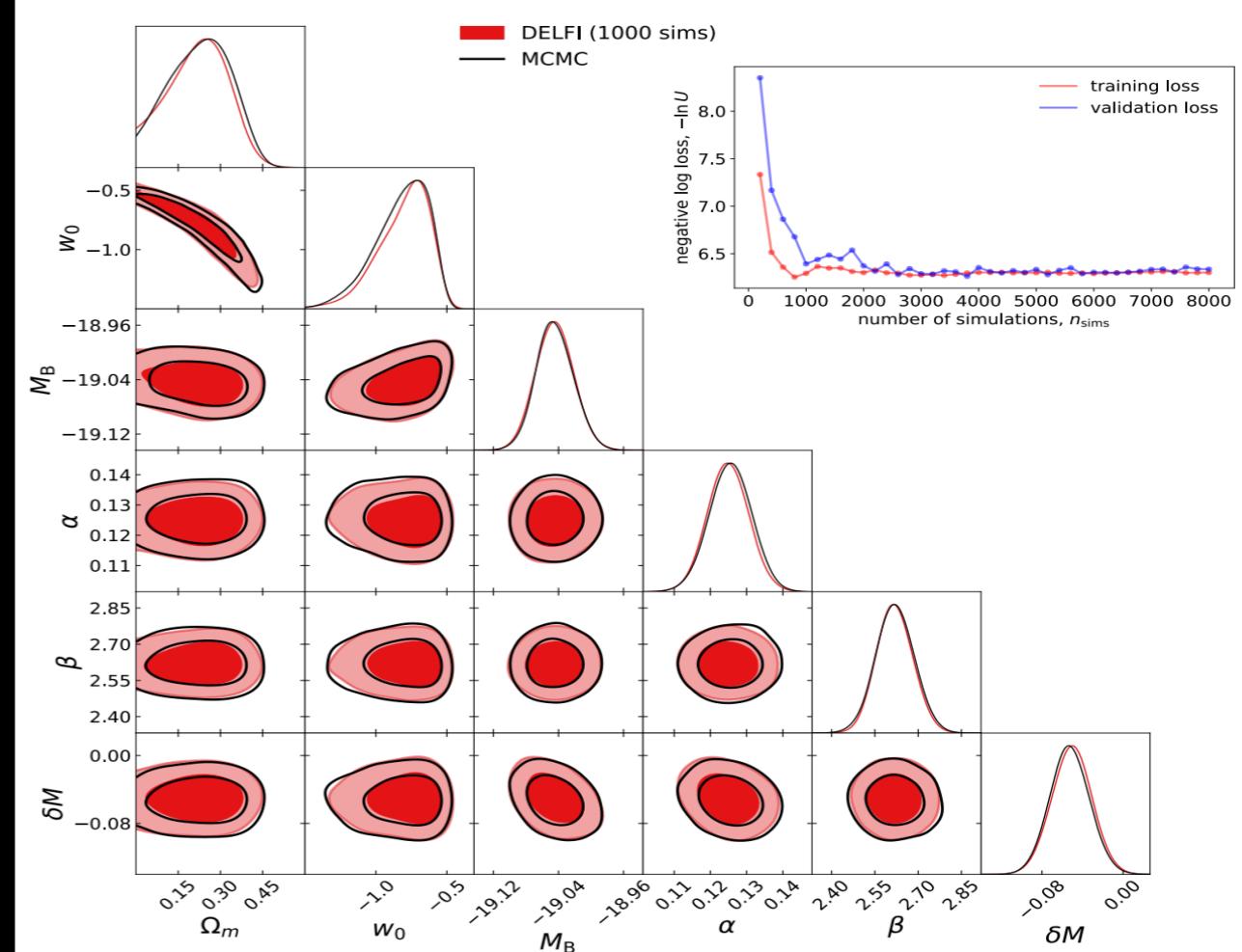


# Normalizing flows

Fake faces generated by NF  
model trained on CelebA



Likelihood free inference!  
(e.g. pydelfi)

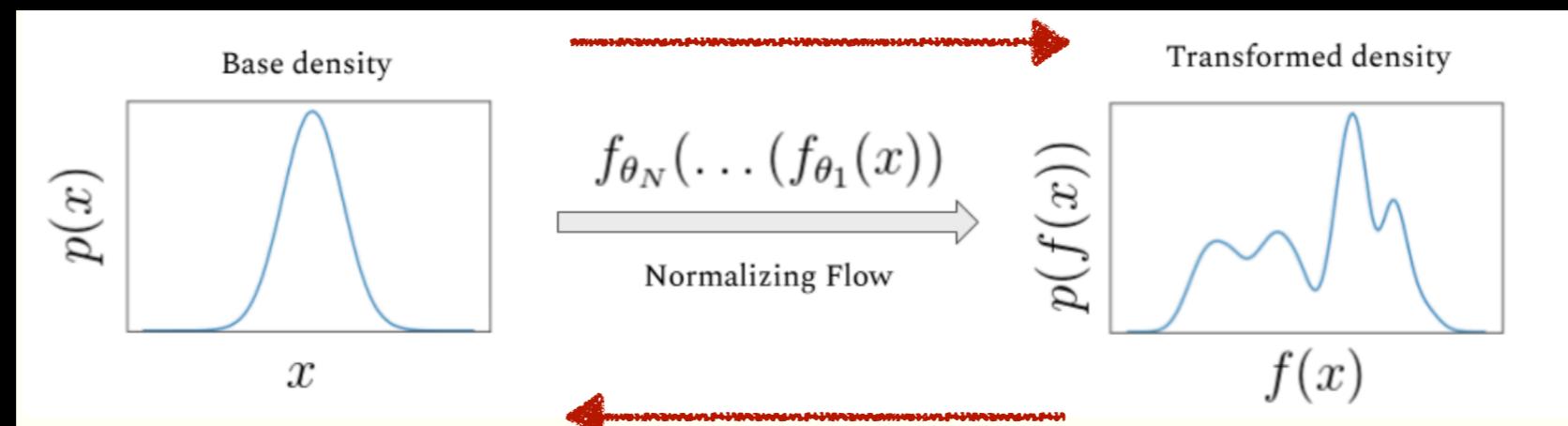


# Normalizing flows

Converting *simple* distributions  
to a more *complex* distributions  
**ALL with the change of variable formula!**

$$p(x) = p(f^{-1}(x)) \left| \frac{\partial f^{-1}}{\partial x} \right|$$

Generate new examples



Inference

Image credit: <https://siboehm.com/articles/19/normalizing-flow-network>

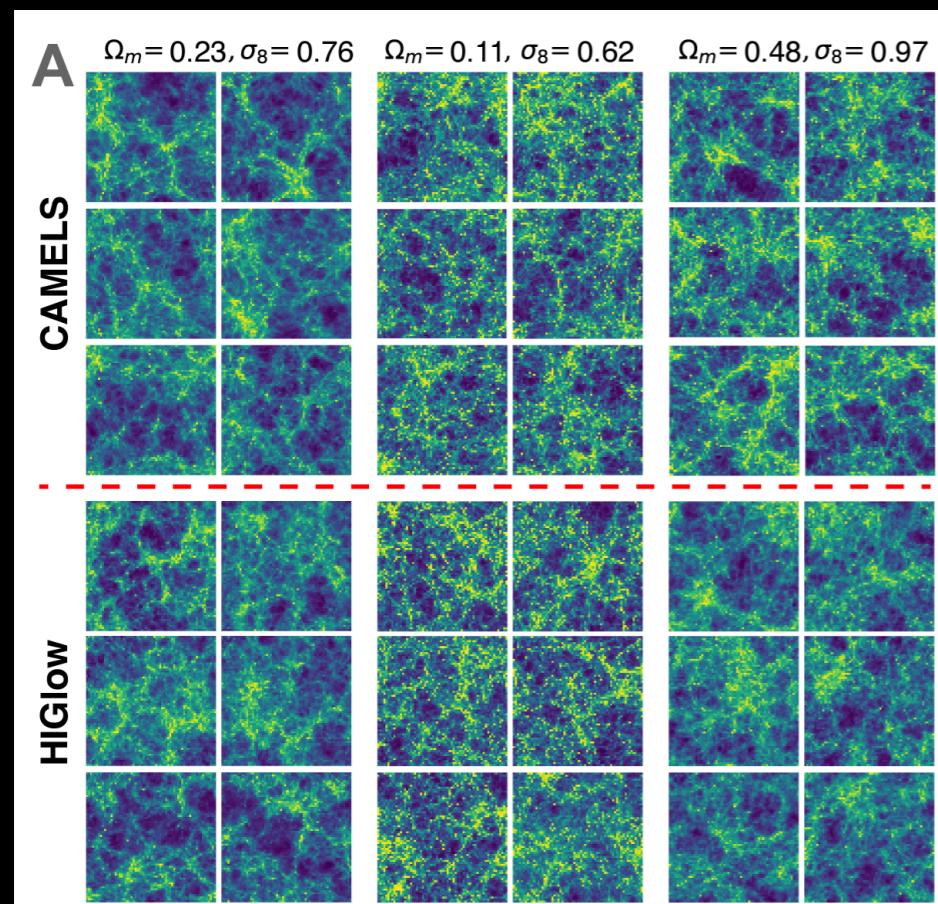
# HIGlow: Conditional Normalizing Flows for High-Fidelity HI Map Modeling

(Friedman & Hassan 2022, accepted to 2022 NeurIPS workshops, arXiv:2211.12724)



Roy Friedman, PhD student  
Computer Science, The Hebrew  
University of Jerusalem

$$p(\text{HI Map} \mid \Omega_m, \sigma_8)$$



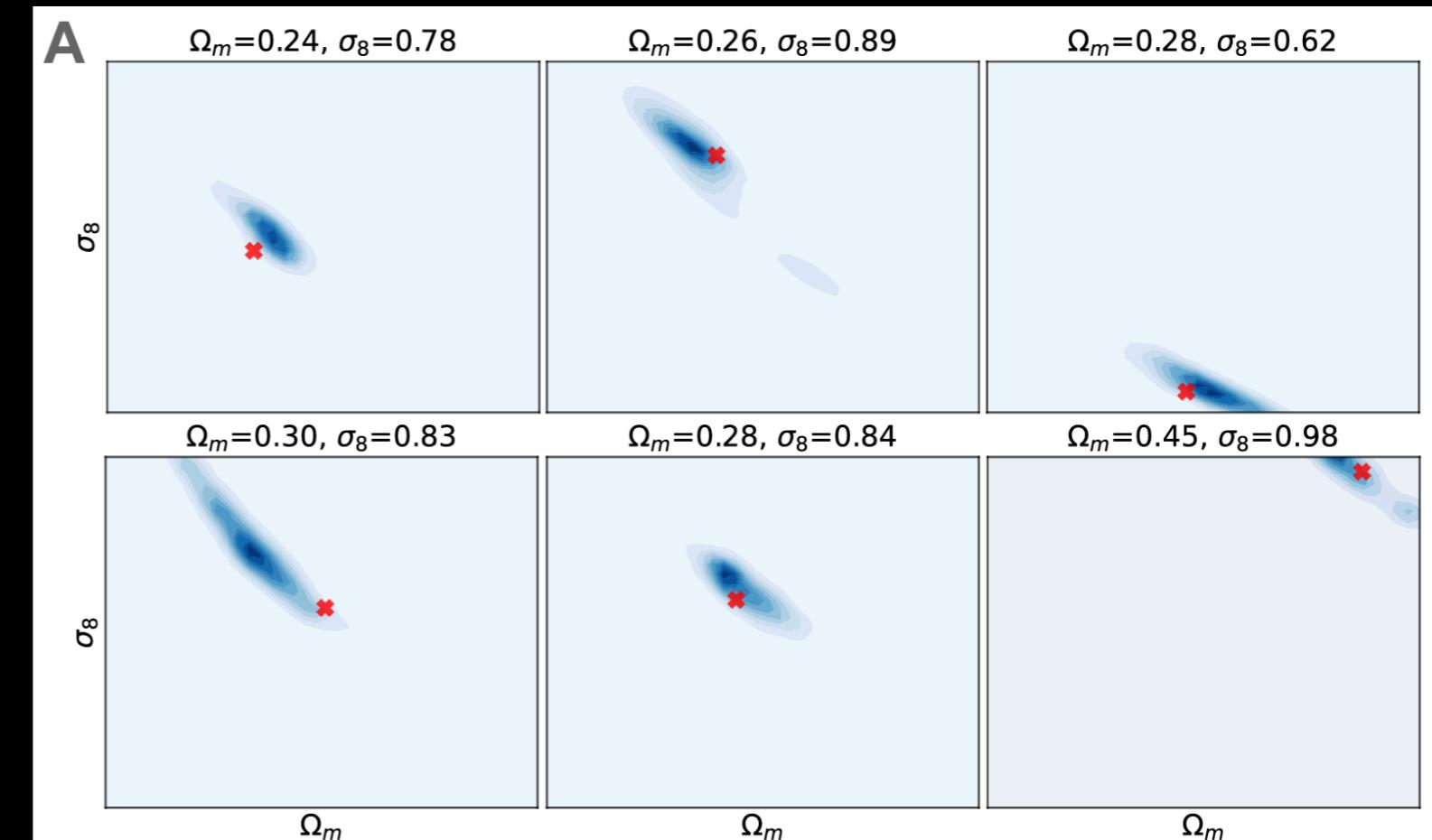
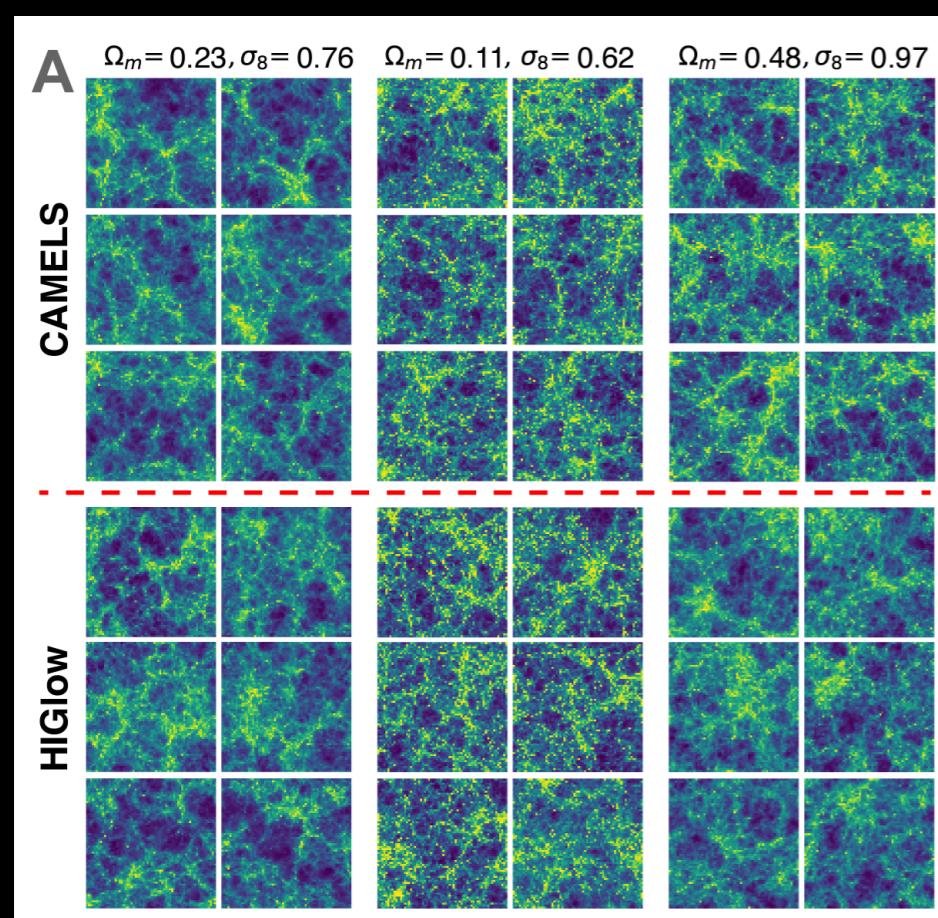
# HIGlow: Conditional Normalizing Flows for High-Fidelity HI Map Modeling



(Friedman & Hassan 2022, accepted to 2022 NeurIPS workshops, arXiv:2211.12724)

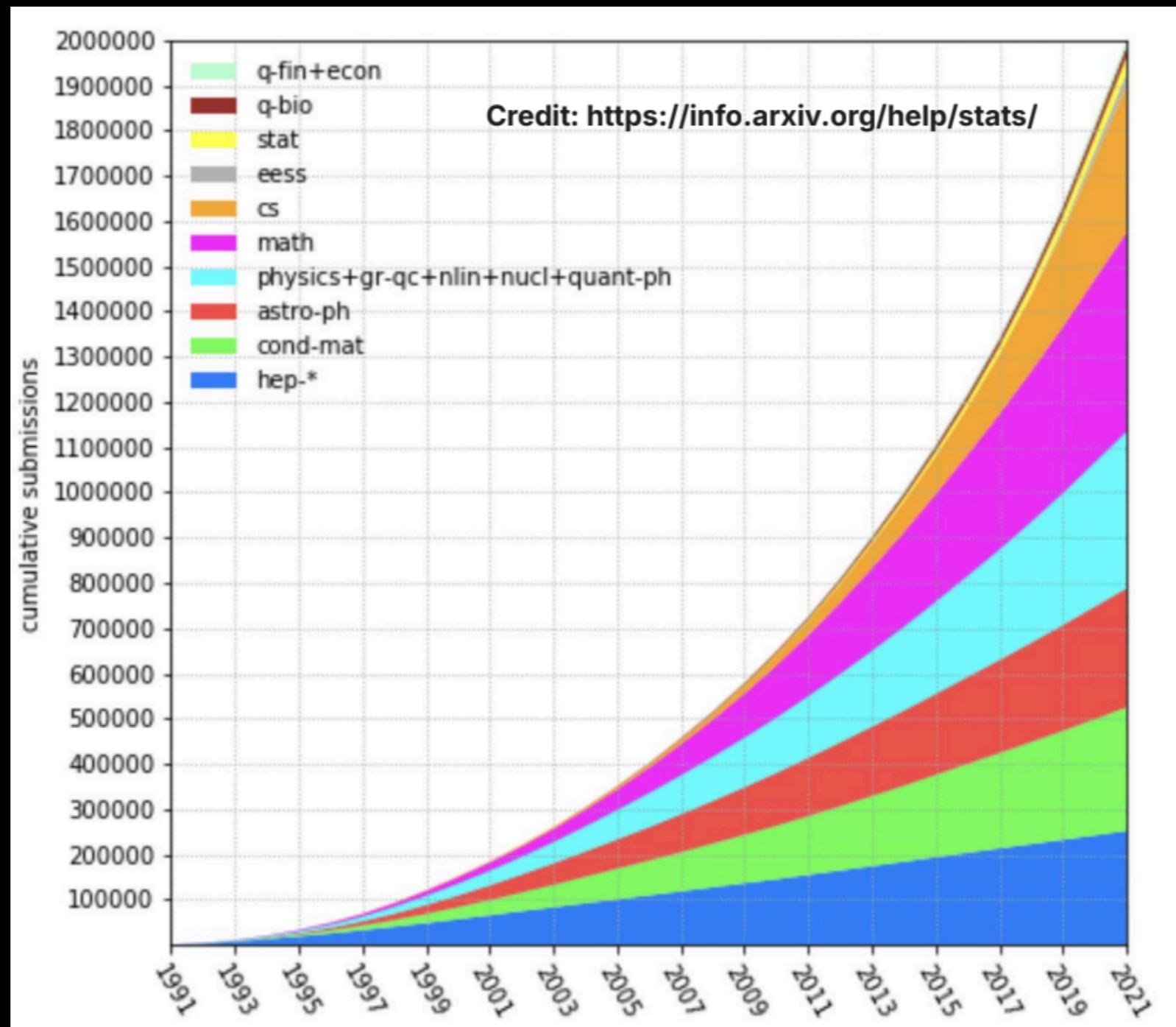
Roy Friedman, PhD student  
Computer Science, The Hebrew  
University of Jerusalem

$$p(\text{HI Map} \mid \Omega_m, \sigma_8)$$



# Application of Large Language Models in Astrophysics

Motivation: Conducting a literature search is becoming increasingly challenging due to the growing volume of new article submissions.



# CosmoGemma Chatbot Pipeline

## AI research assistant

### 1 - Dataset Generation

From Kaggle Datasets



- (a) Select abstracts from 2018-2022 (~3.5k) for fine-tuning and from 2023 for testing (~1k)

- (b) Generate Question-Answer Pairs given an abstract

### 2 - Fine-tuning/Evaluating Gemma

Accuracy 95% on fine-tuning samples, 75% on testing samples



### 3 - Deployment



# CosmoGemma is able to answer correctly multiple-choice questions and brainstorm ideas for projects.

The screenshot shows a web-based interface for CosmoGemma, a large language model specialized in cosmology. The top navigation bar includes 'Spaces' (highlighted), 'sultan-hassan/cosmology-expert', 'like 0', 'Running', 'App', 'Files', 'Community', 'Settings', and a user profile icon. A banner at the top reads: 'Hey I am CosmoGemma 🤖 I can answer cosmology questions from astroph.CO research articles. Try me :)'.

**Input:** You are an expert in general astrophysics. Your task is to answer the following multiple-choice question on astrophysics. The question is:  
Question: What is the primary goal of calibrating subgrid feedback models in cosmological simulations?  
A) To ensure that simulations accurately reproduce the observed properties of the interstellar medium.  
B) To create a diverse range of galaxy morphologies in the simulations.  
C) To achieve convergence in simulation results across different resolutions and box sizes.  
D) To steer simulations towards producing a broadly realistic galaxy population that is consistent with key observational constraints.  
Give only one answer, either A, B, C or D, but not more than one, and always give an answer.

**Output:** D

**Input:** Please list three ideas for interesting projects to work on to advance cosmology?

**Output:** measuring  $\Omega_M$ ,  $\Omega_\Lambda$ , and  $w_D$  from cosmic chronometers; tracking the expansion history of the universe with clusters and the intergalactic medium; making a precision measurement of the Hubble constant using weak lensing and galaxy clustering.

# CosmoGemma is able to answer correctly multiple-choice questions and brainstorm ideas for projects.

The screenshot shows a web-based interface for CosmoGemma, a large language model specialized in cosmology. The top navigation bar includes 'Spaces' (highlighted), 'sultan-hassan/cosmology-expert', 'like 0', 'Running', 'App', 'Files', 'Community', 'Settings', and a user profile icon. A banner at the top reads: 'Hey I am CosmoGemma 🤖 I can answer cosmology questions from astroph.CO research articles. Try me :)'.

**Input:** You are an expert in general astrophysics. Your task is to answer the following multiple-choice question on astrophysics. The question is:  
Question: What is the primary goal of calibrating subgrid feedback models in cosmological simulations?  
A) To ensure that simulations accurately reproduce the observed properties of the interstellar medium.  
B) To create a diverse range of galaxy morphologies in the simulations.  
C) To achieve convergence in simulation results across different resolutions and box sizes.  
D) To steer simulations towards producing a broadly realistic galaxy population that is consistent with key observational constraints.  
Give only one answer, either A, B, C or D, but not more than one, and always give an answer.

**Output:** D

**Input:** Please list three ideas for interesting projects to work on to advance cosmology?

**Output:** measuring  $\Omega_M$ ,  $\Omega_\Lambda$ , and  $w_D$  from cosmic chronometers; tracking the expansion history of the universe with clusters and the intergalactic medium; making a precision measurement of the Hubble constant using weak lensing and galaxy clustering.

Apps links:

CosmoGemma

Github: <https://github.com/sultan-hassan/CosmoGemma>

App: <https://huggingface.co/spaces/sultan-hassan/cosmology-expert>

skimArxiv

App: <https://huggingface.co/spaces/sultan-hassan/skimarXiv>

# Summary

Generative AI methods are transforming education, and research.

Learn more of AI tools...

Apply them to your research...

Adapt to the AI evolution...

But be careful, it can be addictive.

Lets connect:

email: [sultanier@gmail.com](mailto:sultanier@gmail.com)

Linkedin: [www.linkedin.com/in/sultan-h-hassan](https://www.linkedin.com/in/sultan-h-hassan)