



Finanziato  
dall'Unione europea  
NextGenerationEU



Ministero  
dell'Università  
e della Ricerca



Italiadomani

PIANO NAZIONALE  
DI RIPRESA E RESILIENZA



Centro Nazionale di Ricerca in HPC,  
Big Data and Quantum Computing

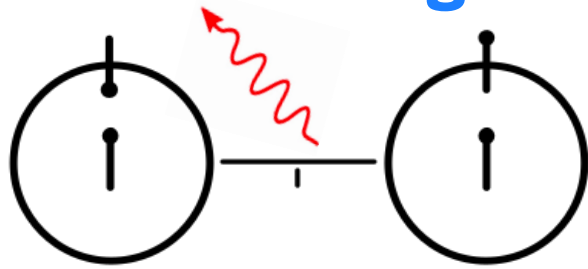
# *Optimal compression and Simulation-Based Inference of the cosmic 21-cm signal*

*David Prelogović, Andrei Mesinger*

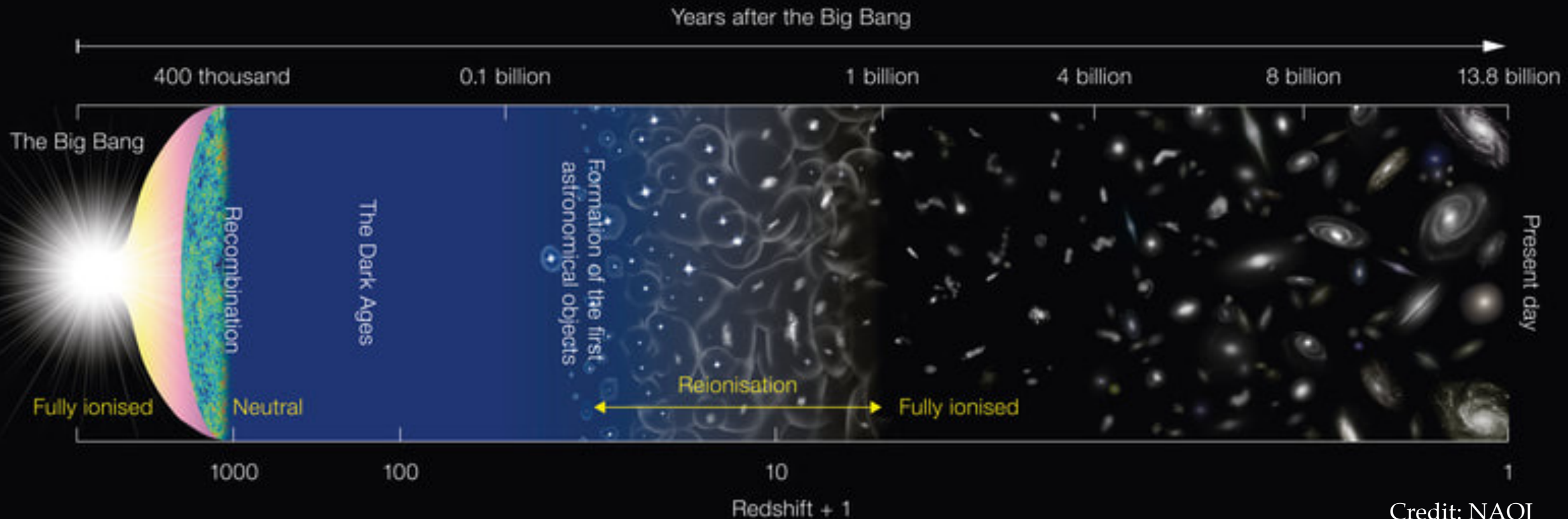
*Scuola Normale Superiore*

Spoke 3 Technical Workshop, Trieste October 9 / 11, 2023

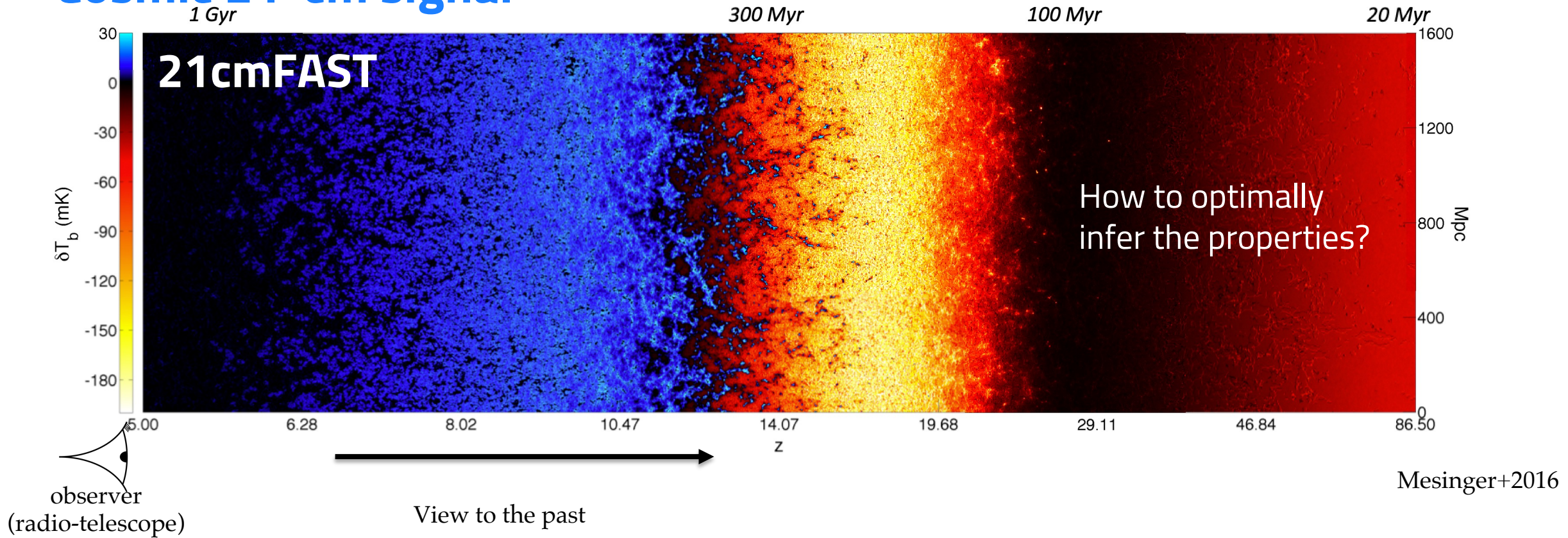
# Cosmic 21-cm signal



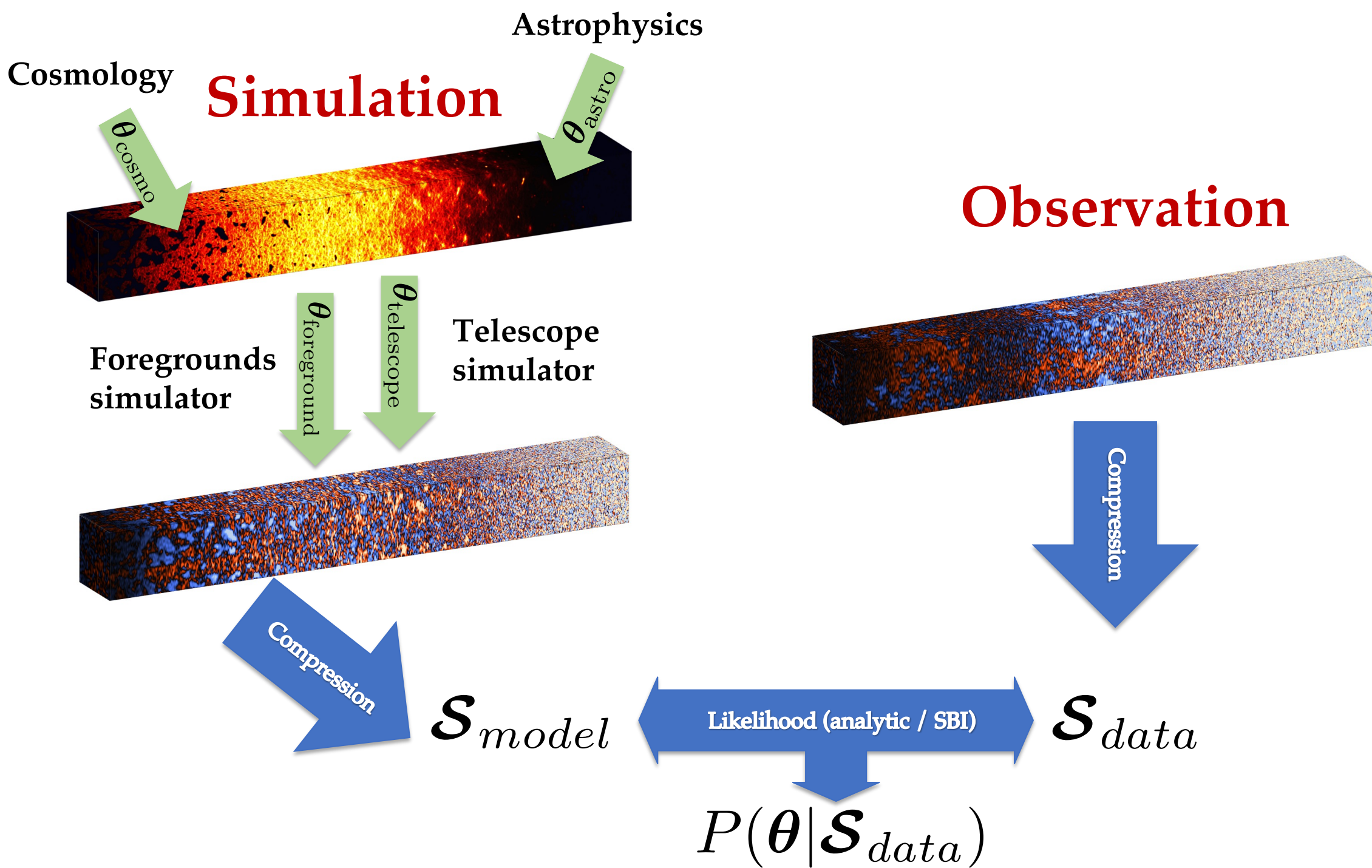
- Over 90% of the "normal" matter in the Universe is **hydrogen**



# Cosmic 21-cm signal



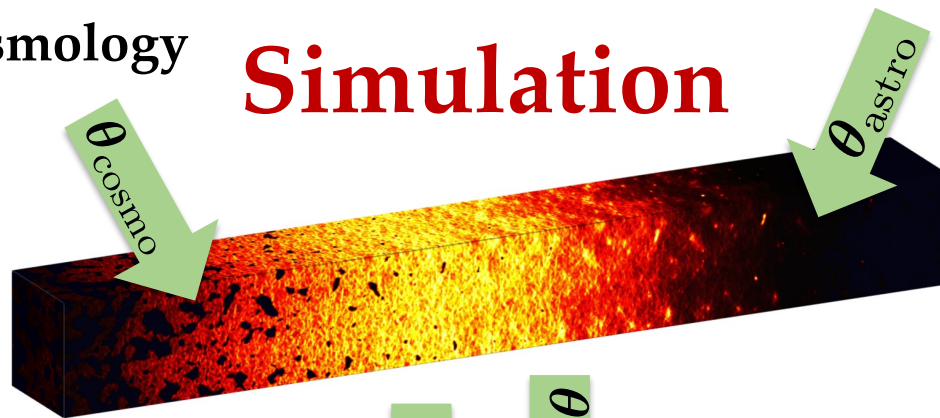
- The properties of the ***unseen first galaxies*** are encoded in the timings and patterns of this signal!



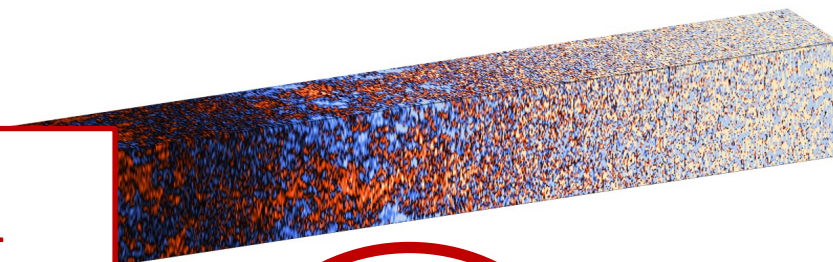
Astrophysics

Cosmology

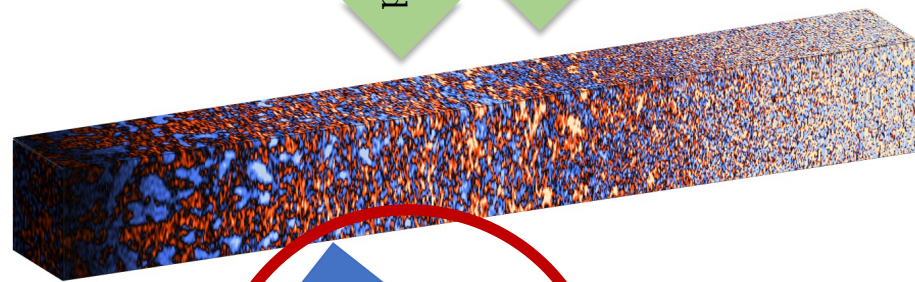
Simulation



Observation



# 1 – optimal compression



$\mathcal{S}_{model}$



$\mathcal{S}_{data}$

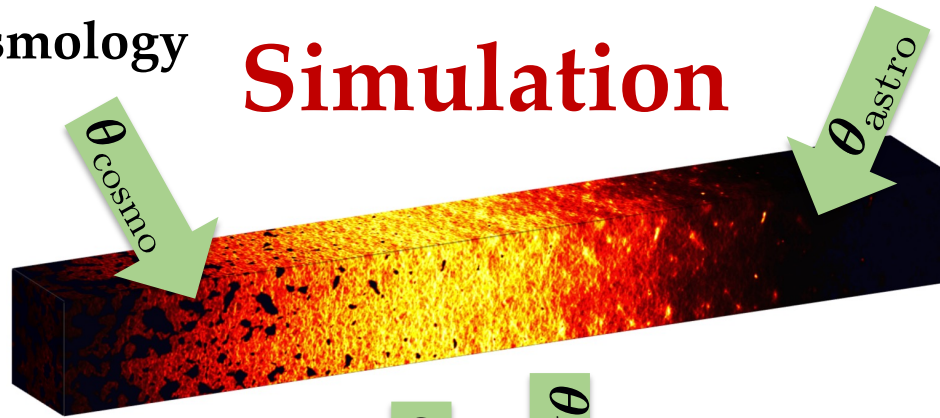


$$P(\theta | \mathcal{S}_{data})$$

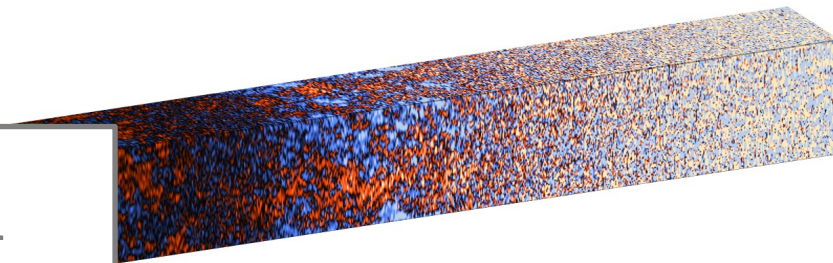
Cosmology

Astrophysics

**Simulation**



**Observation**



# 1 – optimal compression

# 2 – “optimal” likelihood



$\mathcal{S}_{model}$



$\mathcal{S}_{data}$



$$P(\theta | \mathcal{S}_{data})$$

# Scientific Rationale

- **The 21-cm signal has non-Gaussian formation with unknown likelihood**
- **Future telescopes such as SKA will be able to record images of the signal**
- **Advanced Inference techniques are necessary:**
  - I. to extract as much information as possible form the signal**
  - II. to reduce computational cost of classical inference techniques**

## Main objectives:

**#1 – finding the optimal compression of the signal**

**#2 - “writing” the likelihood of the signal – Simulation-Based Inference**

The background is a deep blue gradient. On the left side, there are numerous bright blue light trails and dots that appear to be moving towards the center, creating a sense of depth and motion. The trails are composed of many small, overlapping lines and points of light, some of which are blurred, suggesting speed. The overall effect is reminiscent of a data stream or a futuristic digital environment.

# **Simulation-Based Inference for the 21-cm PS**



## Likelihood of the 21-cm

- Analytically not available
- Numerically non-tractable

## Likelihood of the 21-cm

- Analytically not available
- Numerically non-tractable
- Eg. – power spectrum

$$\mathcal{L}(\boldsymbol{\sigma}, \boldsymbol{\mu}; \mathbf{x}) = \frac{1}{(2\pi)^{n/2} \sqrt{\prod_i \sigma_i^2}} e^{-\frac{1}{2} \sum_i (x_i - \mu_i)^2 / \sigma_i^2}$$

$$\mathcal{L}(\boldsymbol{\Sigma}, \boldsymbol{\mu}; \mathbf{x}) = \frac{1}{(2\pi)^{n/2} \sqrt{|\boldsymbol{\Sigma}|}} e^{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}$$

## Likelihood of the 21-cm

- Analytically not available
- Numerically non-tractable
- Eg. – power spectrum

$$\mathcal{L}(\boldsymbol{\sigma}, \boldsymbol{\mu}; \mathbf{x}) = \frac{1}{(2\pi)^{n/2} \sqrt{\prod_i \sigma_i^2}} e^{-\frac{1}{2} \sum_i (x_i - \mu_i)^2 / \sigma_i^2}$$

$$\mathcal{L}(\boldsymbol{\Sigma}, \boldsymbol{\mu}; \mathbf{x}) = \frac{1}{(2\pi)^{n/2} \sqrt{|\boldsymbol{\Sigma}|}} e^{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}$$

### **Classical inference simplifications**

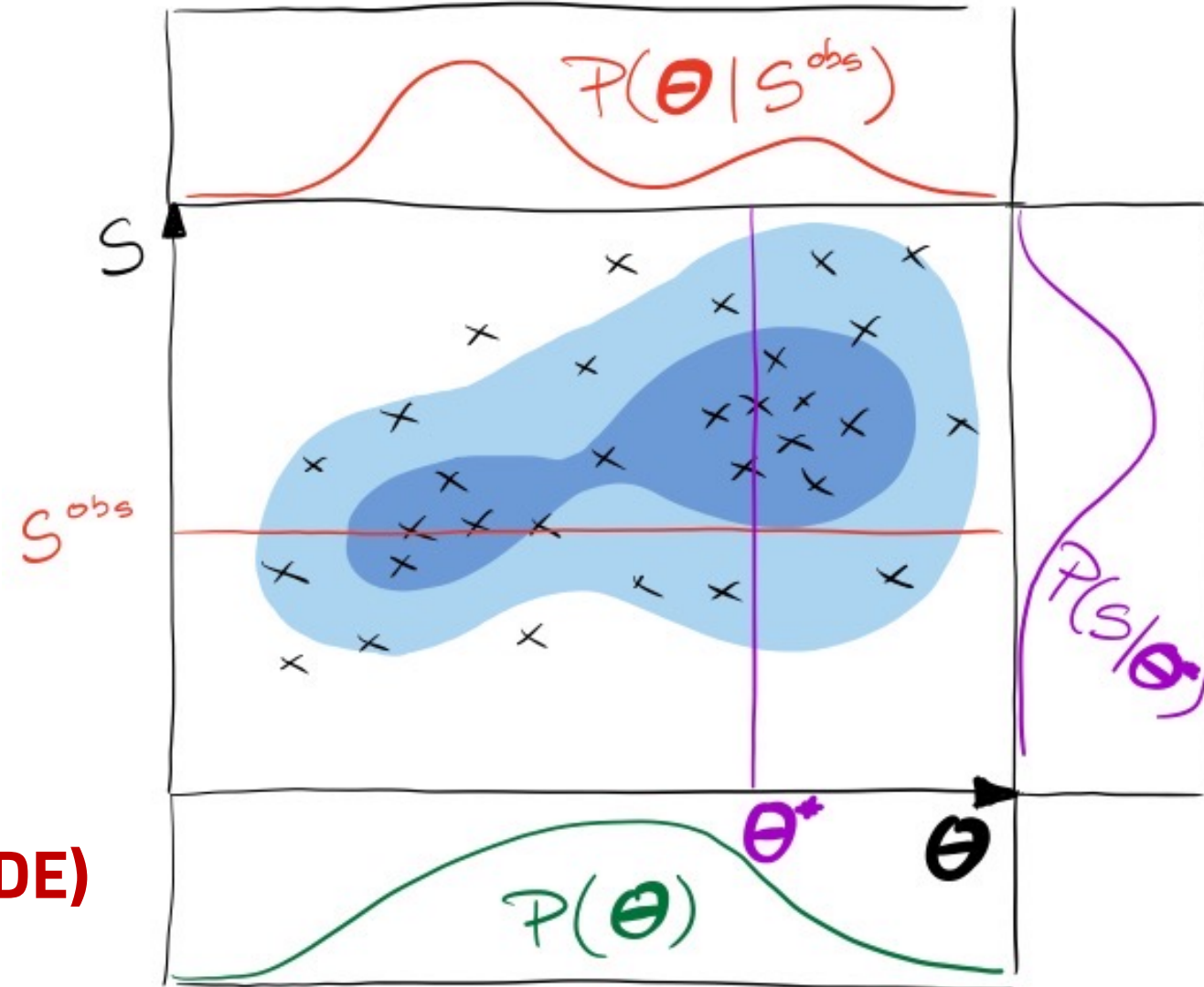
- 1)  $\boldsymbol{\mu}$  estimated from one sim.
- 2)  $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_{\text{fid}}$
- 3) No higher order terms

→ Simulation-Based Inference

## Simulation-Based Inference

- Pull from the prior  $\theta^* \sim P(\theta)$
- Pull from the likelihood  
 $S^* \sim P(S|\theta^*)$   
i.e. simulate and compress  
 $S^* = \text{compress}(\text{simulate}(\theta^*))$
- Repeat many times
- Fit the distribution with

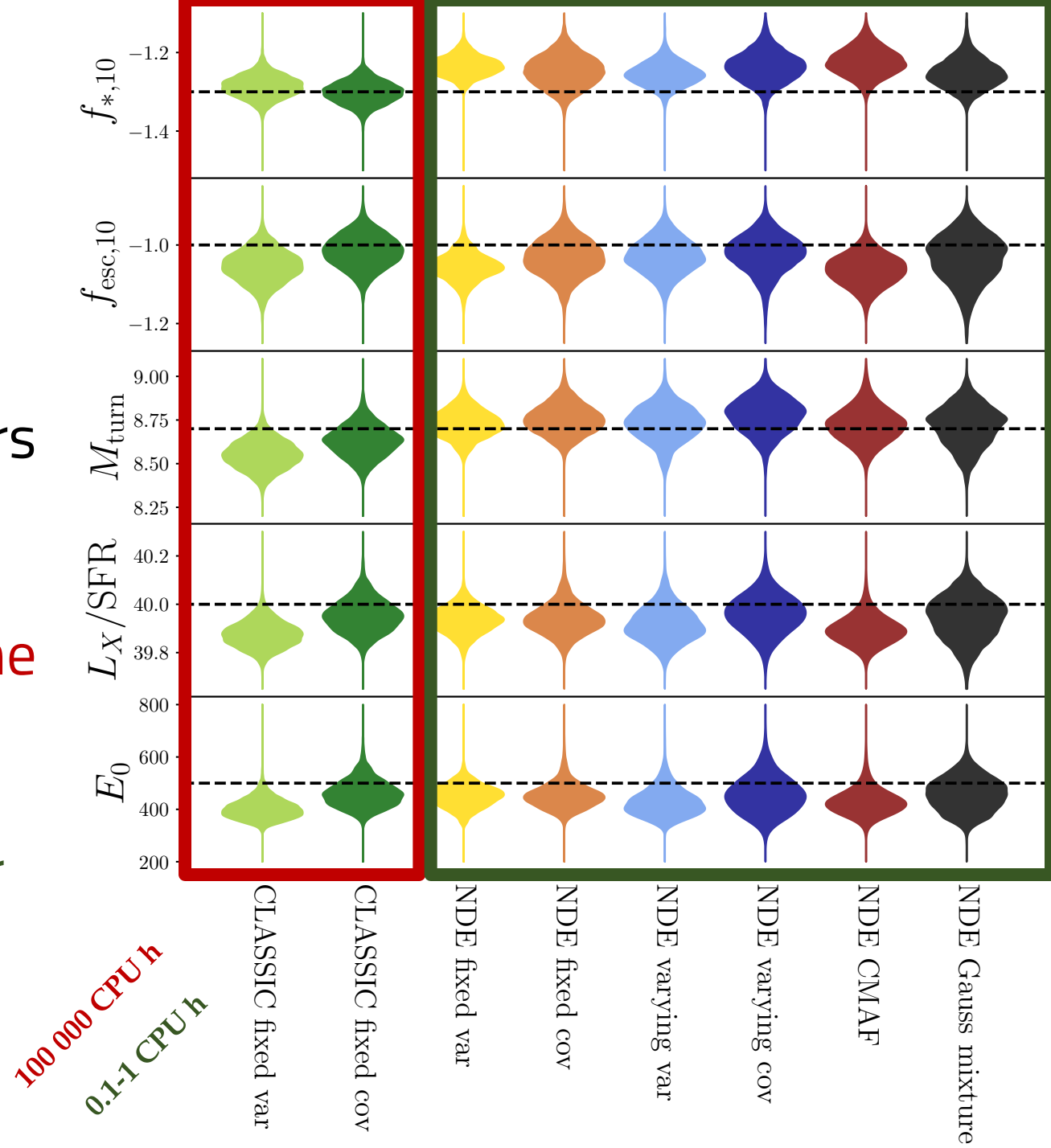
## Neural Density Estimators (NDE)



Credit: Tom Charnock

## SBI for the 21-cm PS

- Able to infer the parameter posteriors
- At a fraction of computational cost
- **Caveat: large database needed for the NDE training!**
- Once constructed, reusable for other inferences, other summaries



The background is a deep blue gradient. On the left side, there are numerous bright blue light trails that appear to be moving towards the center, creating a sense of depth and motion. These trails are composed of many small, bright blue dots or particles. The overall effect is reminiscent of a digital or scientific visualization, such as a light cone or a data stream.

# **Optimal Compression of the 21-cm lightcones**

## Ad-hoc compressions

- Without good a-priori physical motivation, we cannot know what is THE optimal compression/summary, i.e. providing tightest recovery of astrophysics and cosmology

## Ad-hoc compressions

- Without good a-priori physical motivation, we cannot know what is THE optimal compression/summary, i.e. providing tightest recovery of astrophysics and cosmology

Solution:

Let the machines figure it out for us!

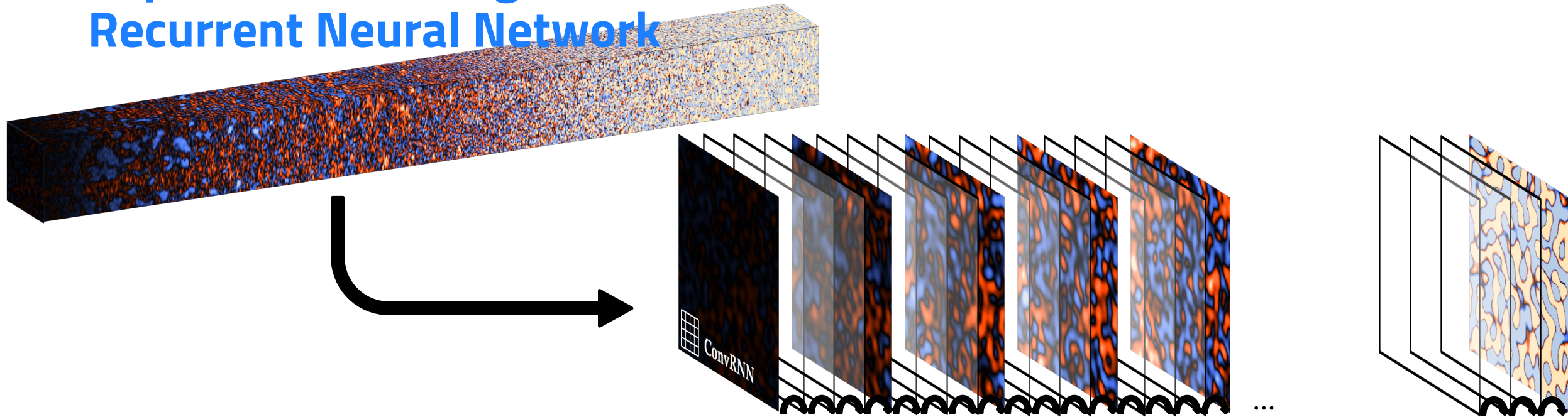
(Neural Network)

- Gillet+2018
- La Plante & Ntampaka 2019
- Makinen+2020
- Mangena+2020
- Hortúa+2020
- Prelogović+2021
- +++





# Supervised learning: Recurrent Neural Network

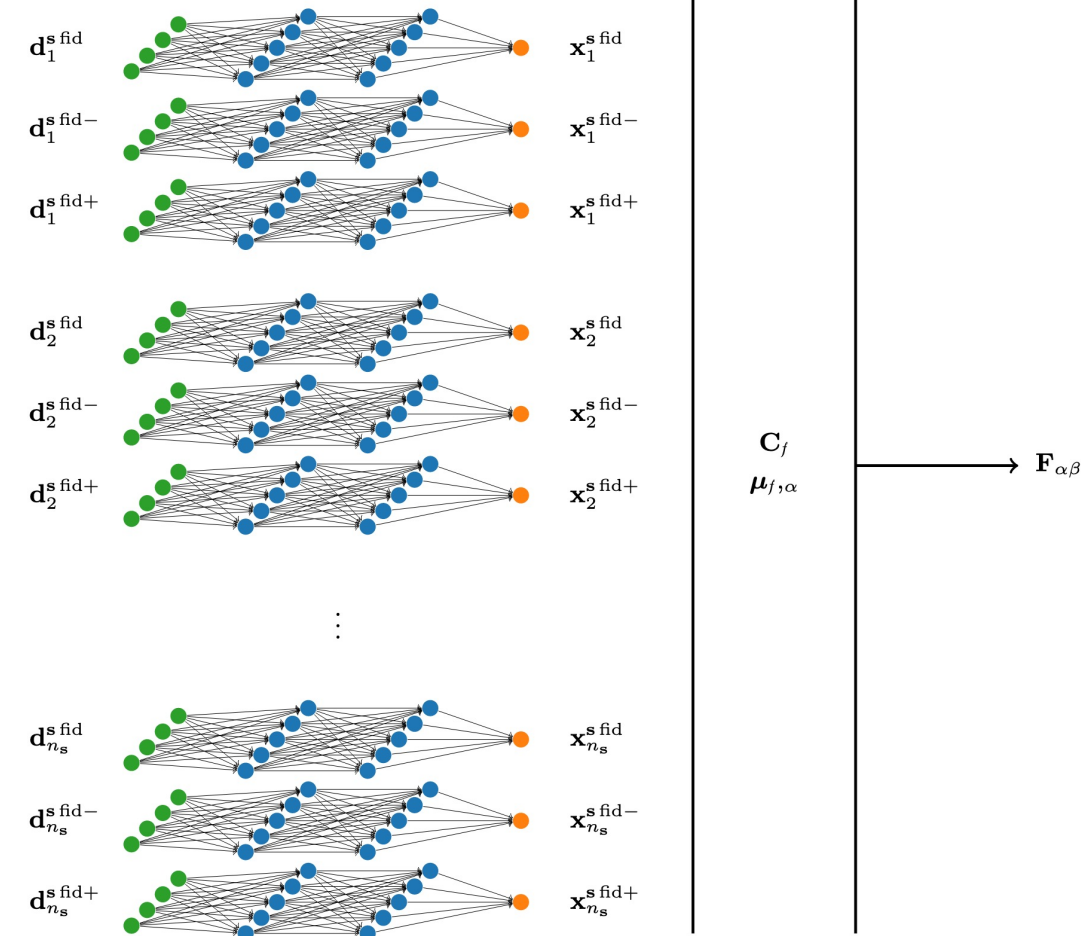


- RNN – encoding correlations across all frequency bins at once
- 2D Convolutional NN– local correlations in sky-plane

# Unsupervised learning: Information Maximizing Neural Network

- Simulate the data at a fiducial parameter set:  $\mathbf{d}(\boldsymbol{\theta}_{\text{fid}})$
- Simulate around the fiducial parameters:  $\mathbf{d}(\boldsymbol{\theta}_{\text{fid}}^+)$ ,  $\mathbf{d}(\boldsymbol{\theta}_{\text{fid}}^-)$
- Calculate compressed summary:  
$$\mathbf{s}(\boldsymbol{\theta}) = \text{NN}(\mathbf{d}(\boldsymbol{\theta}))$$

- Maximize Fisher information:  
$$L = -\ln |F|$$



# Methods and results

## Database:

- for ~150 points pulled from the prior, construct finite—difference “Fisher database”  
→ 50 000 samples

## Methodology:

- distribution of Fisher information ( $\ln|F|$ ) as a proxy for summary quality

## Summaries:

- RNN & IMNN for deep learning summaries
- Wavelets, 1DPS and 2DPS for classical summaries

# Methods and results

## Database:

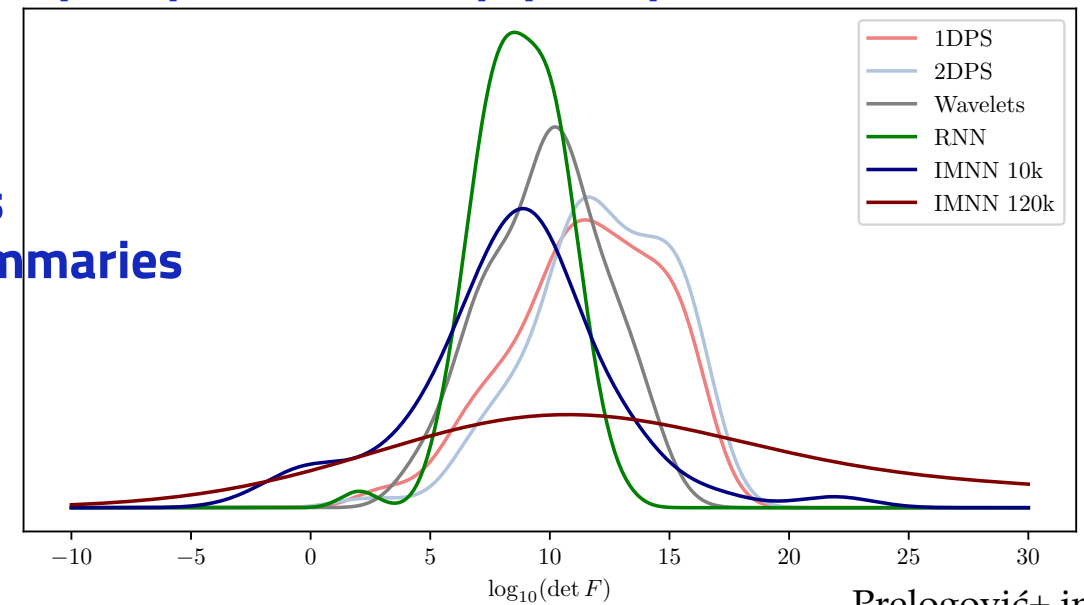
- for ~150 points pulled from the prior, construct finite—difference “Fisher database”  
→ 50 000 samples

## Methodology:

- distribution of Fisher information ( $\ln|F|$ ) as a proxy for summary quality

## Summaries:

- RNN & IMNN for deep learning summaries
- Wavelets, 1DPS and 2DPS for classical summaries



Prelogović+ in.prep.

The background is a deep blue gradient. On the left side, there are numerous light trails and dots in shades of cyan and white, creating a sense of depth and movement, similar to a data visualization or a futuristic tunnel. The word "Summary" is centered in the middle of the image.

# Summary

# Milestones and KPIs

## Construction of databases:

- SBI database ~100 000 21cmFAST simulations
- IMNN database ~ 20 000
- Fisher database ~ 50 000

## NN trainings:

- IMNN – 5 000 GPU h per model
- RNN – 1 000 GPU h per model
- SBI NDEs – relatively cheap (~few GPUh per model)

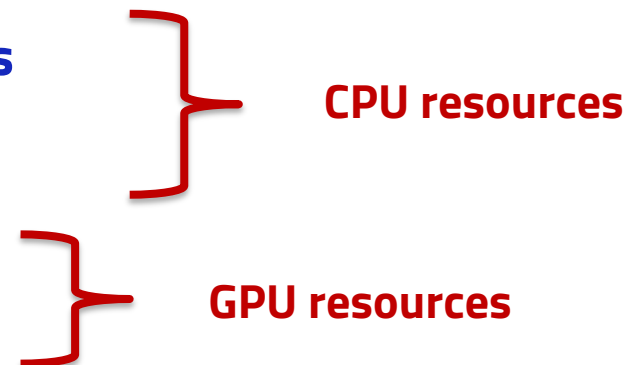
## Code release:

- two publicly available Python libraries:
  - > 21cmLikelihoods – SBI library for exploring 21cmPS likelihoods  
<https://github.com/dprelogo/21cmLikelihoods>
  - > conditionalKDE – helper library for prior sampling  
[https://github.com/dprelogo/conditional\\_kde](https://github.com/dprelogo/conditional_kde)

## Publications:

Exploring the likelihood of the 21-cm power spectrum with simulation-based inference  
[arXiv:2305.03074](https://arxiv.org/abs/2305.03074)

## Bottlenecks:



## Next Steps and Expected Results

### Code release:

**summaries21cm: library for exploration of different summaries of the 21cm signal**

### Publications:

**How informative are summaries of the cosmic 21-cm signal?**

### Database:

**start a construction of a 9-dim database, for multi-wavelength SBI (talk tomorrow)  
(granted a fraction of CPU resources)**

### Projects:

**constrained realizations of the 21-cm signal from galaxy probes  
(granted small Cineca C call – Leonardo, more resources needed)**