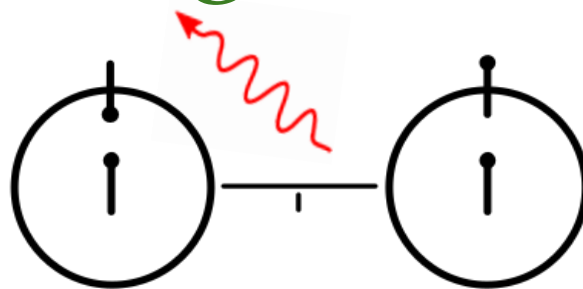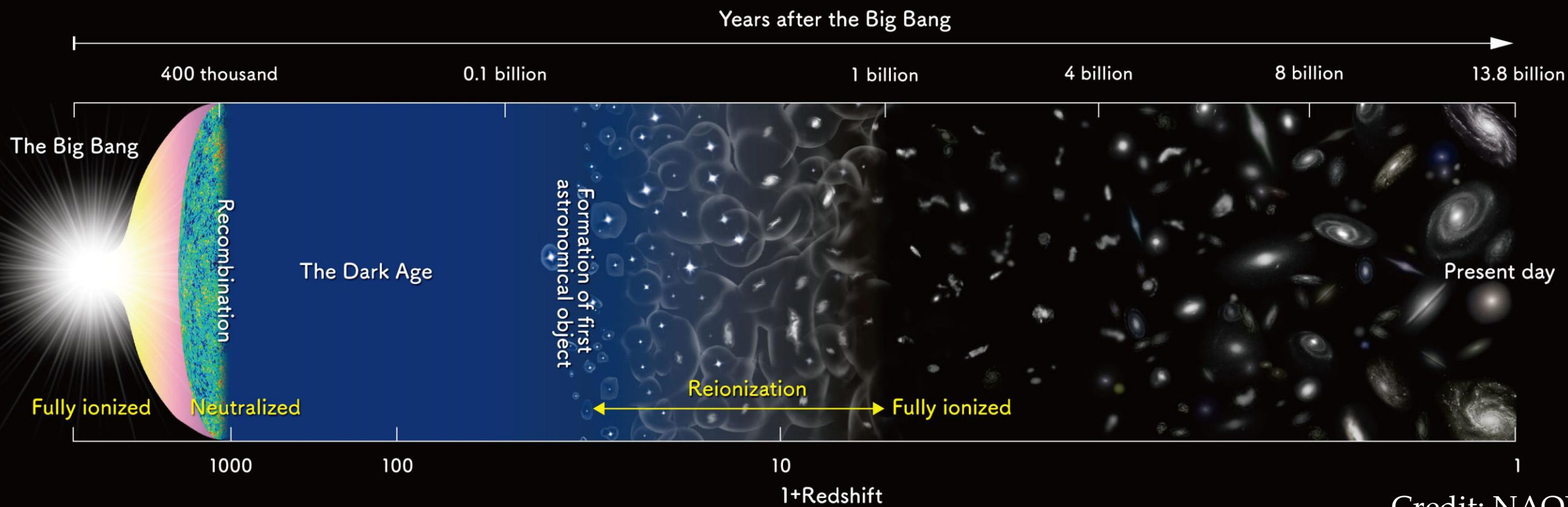# Machine Learning of the Cosmic 21-cm Signal

David Prelogović

PostDoc @ SISSA, Trieste, IT
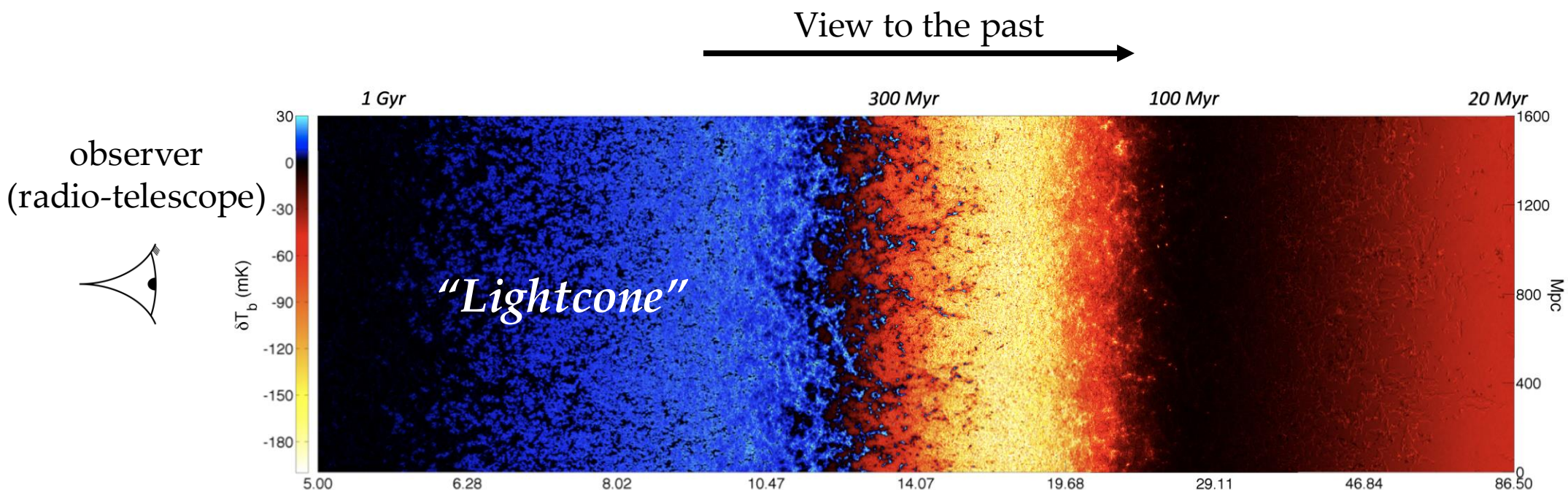
# 1. Cosmic 21-cm Signal

- Hydrogen atoms abundant throughout the Universe's evolution
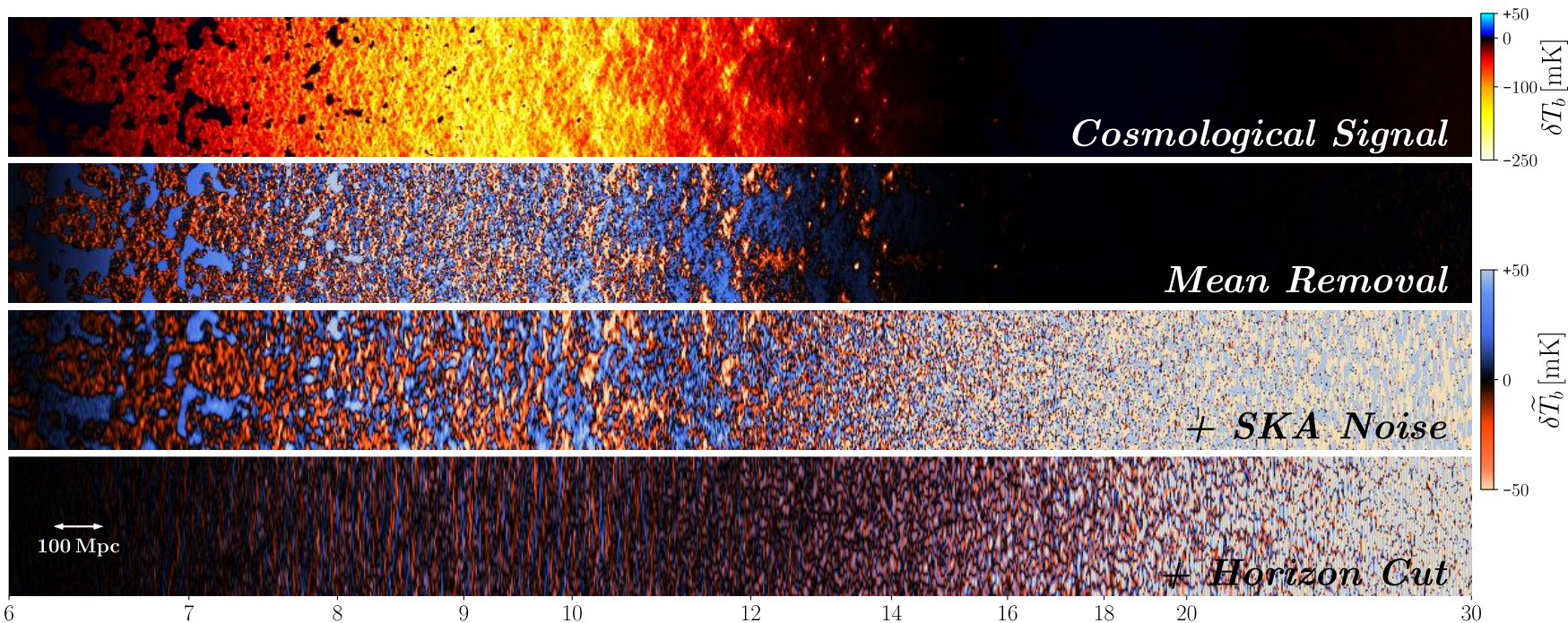- Encoding the first billion years



Credit: NAOJ

# 2. Cosmic 21-cm Signal

View to the past →



observer (radio-telescope)

*"Lightcone"*

- cosmo. + astro.

$$\delta T_b \approx 30\, x_{\mathrm{HI}}\, \Delta \left( \frac{H}{dv_r/dr + H} \right) \left( 1 - \frac{T_\gamma}{T_S} \right) \left( \frac{1+z}{10}\, \frac{0.15}{\Omega_{\mathrm{M}} h^2} \right)^{1/2} \left( \frac{\Omega_b h^2}{0.023} \right) \mathrm{mK}$$

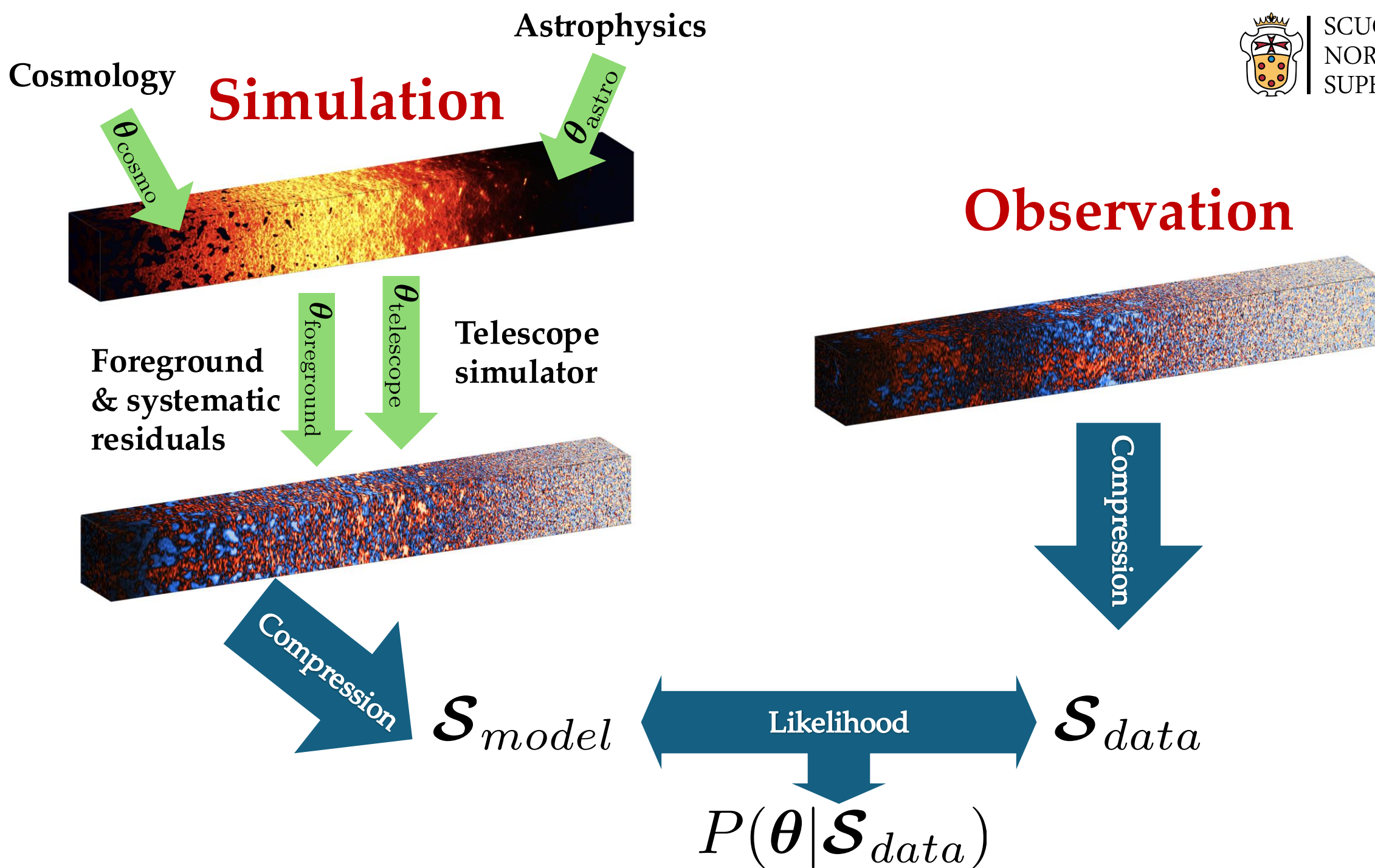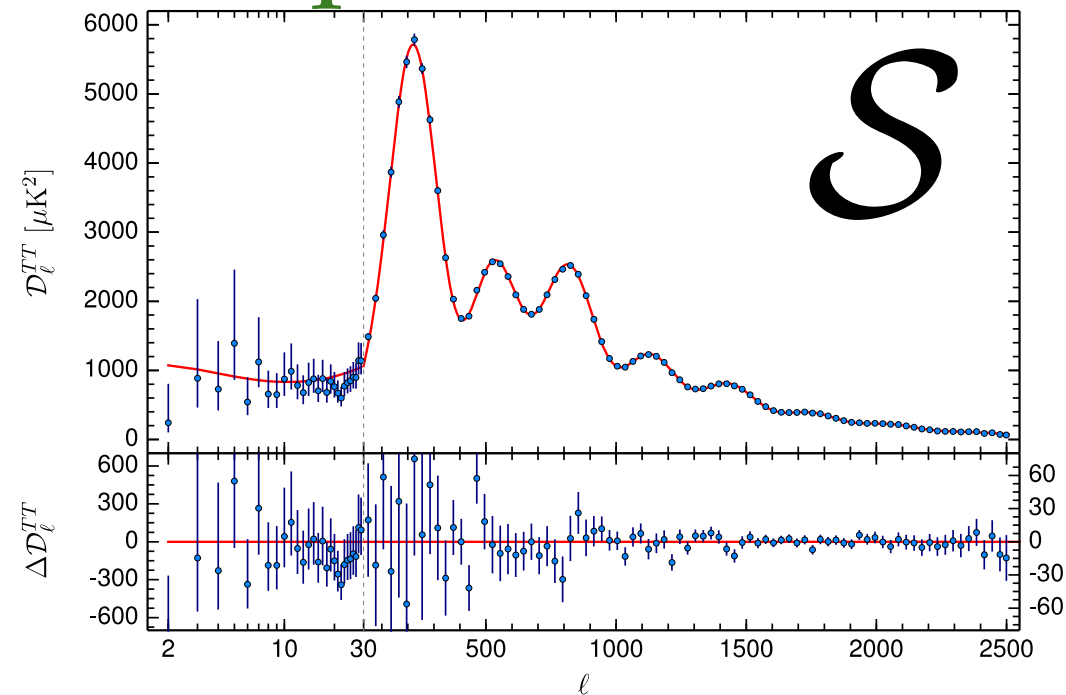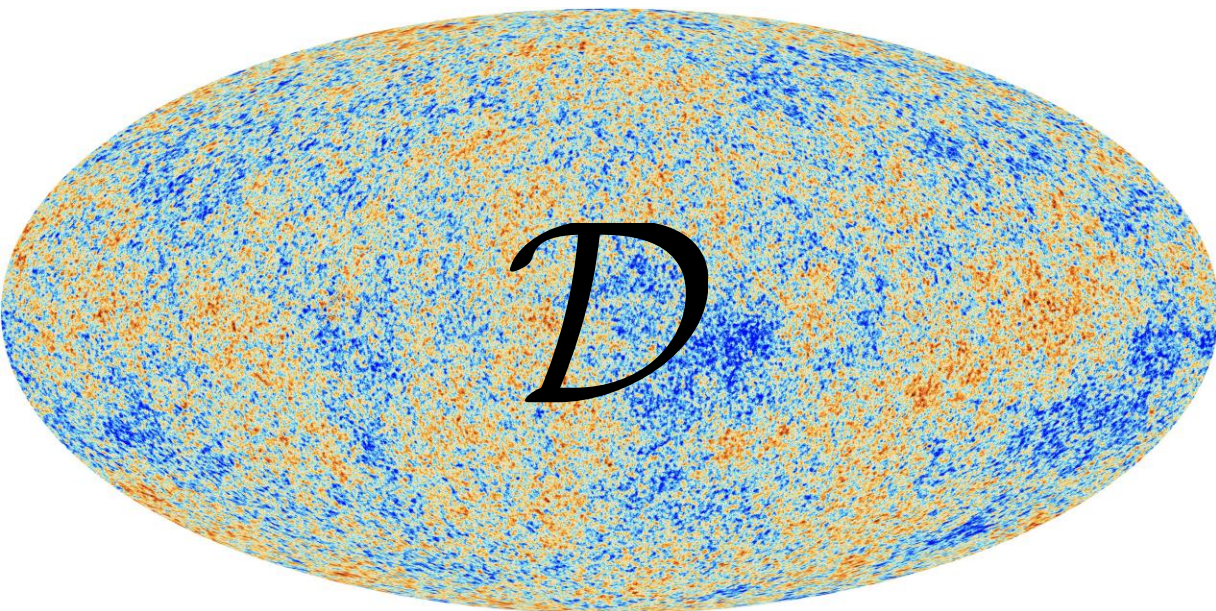# 3. Forward modeling pipeline



21cmFAST

SKA simulator

Foreground avoidance

Prelogović+2022

**Cosmology** **Simulation** **Astrophysics**

$\theta_{cosmo}$   $\theta_{astro}$

SCUOLA NORMALE SUPERIORE

**Observation**

$\theta_{foreground}$   $\theta_{telescope}$

**Foreground & systematic residuals**   **Telescope simulator**

Compression

Compression

$\mathcal{S}_{model}$ ⟷ Likelihood ⟷ $\mathcal{S}_{data}$
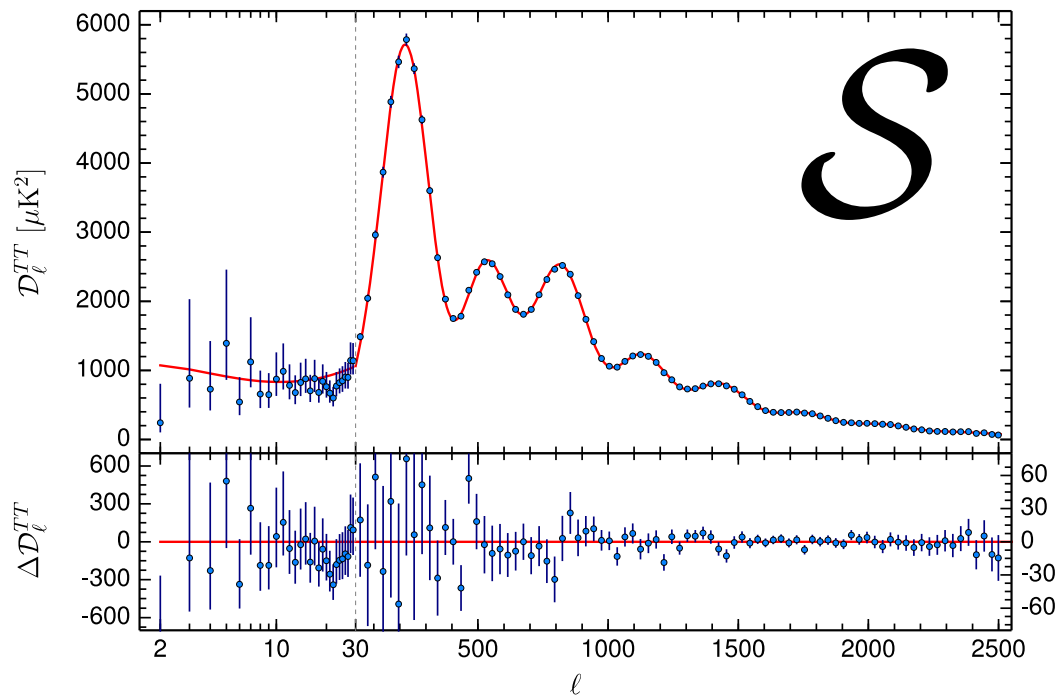
$$P(\boldsymbol{\theta}|\mathcal{S}_{data})$$

# 4.1 Classical Inference Example: CMB
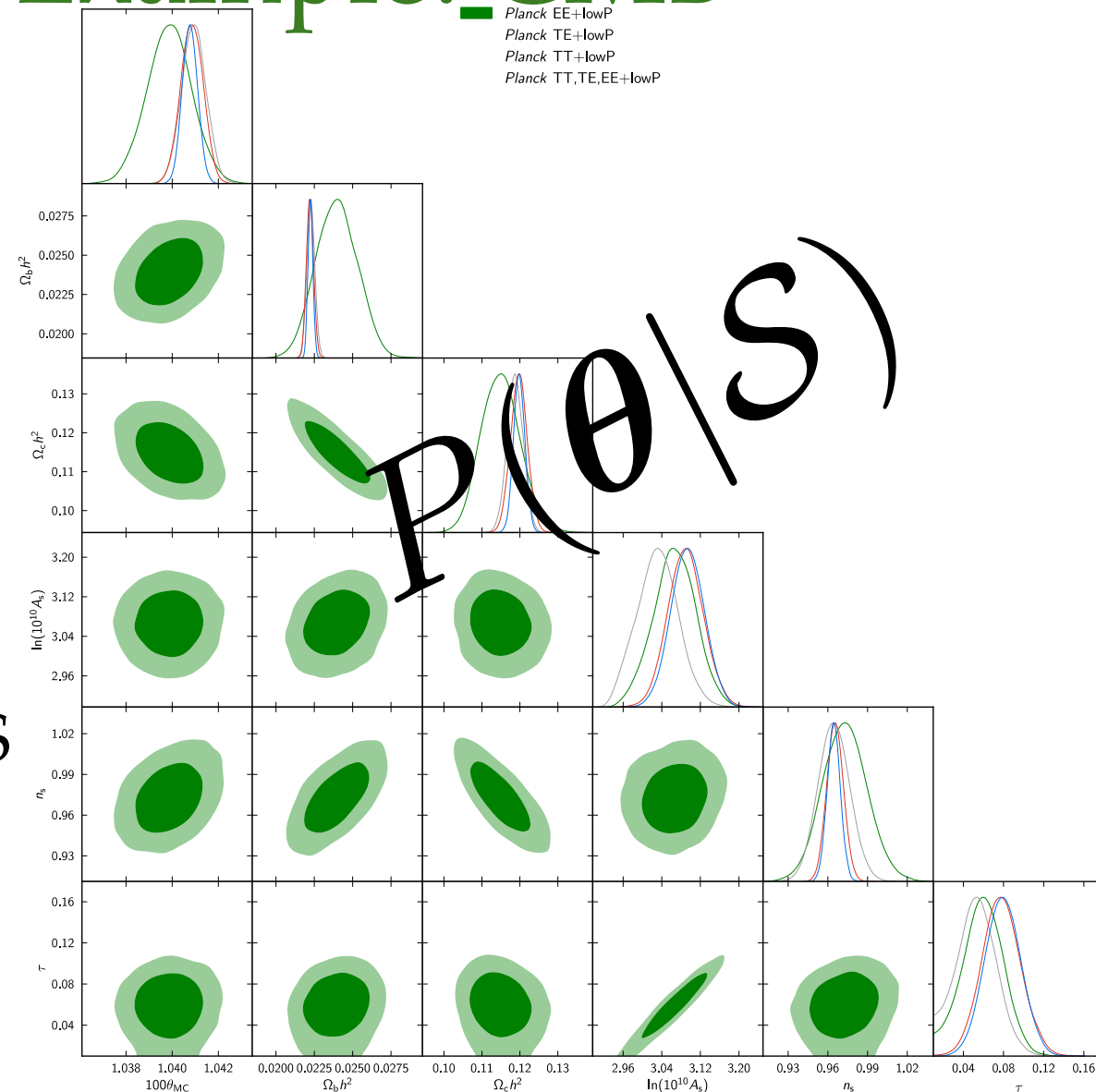


- Full sky map compressed to 1DPS
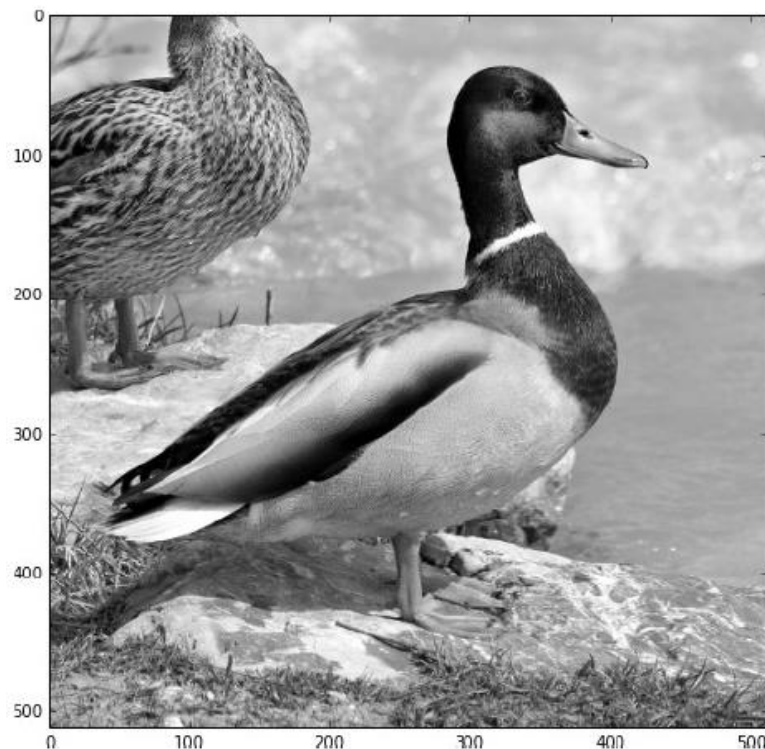  - Known, optimal compression
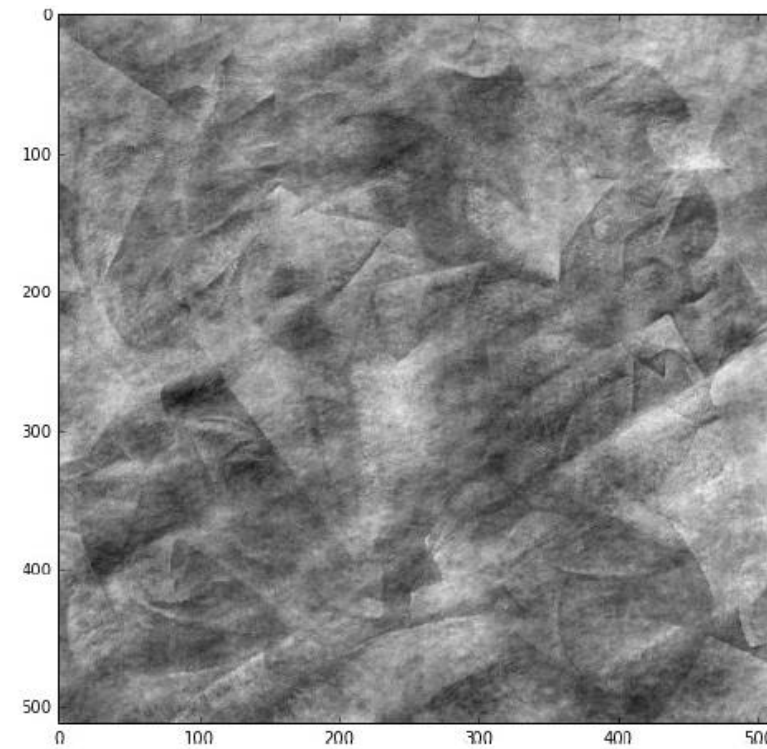
# 4.1 Classical Inference Example: CMB



- Full sky map compressed to 1DPS
  - Known, optimal compression
- From it we infer the cosmology
  - Known likelihood

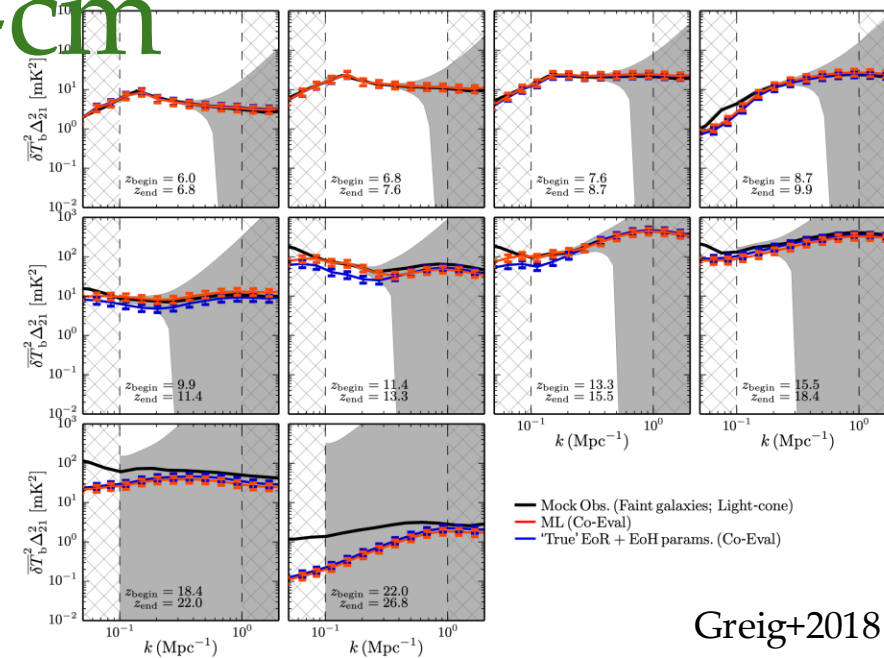Planck 2015

# 4.2 Compression for a Duck



change phases

- Same 2D PS
- Highly non-Gaussian

Credit: G. Bernardi

# 4.3 Compression for the 21-cm

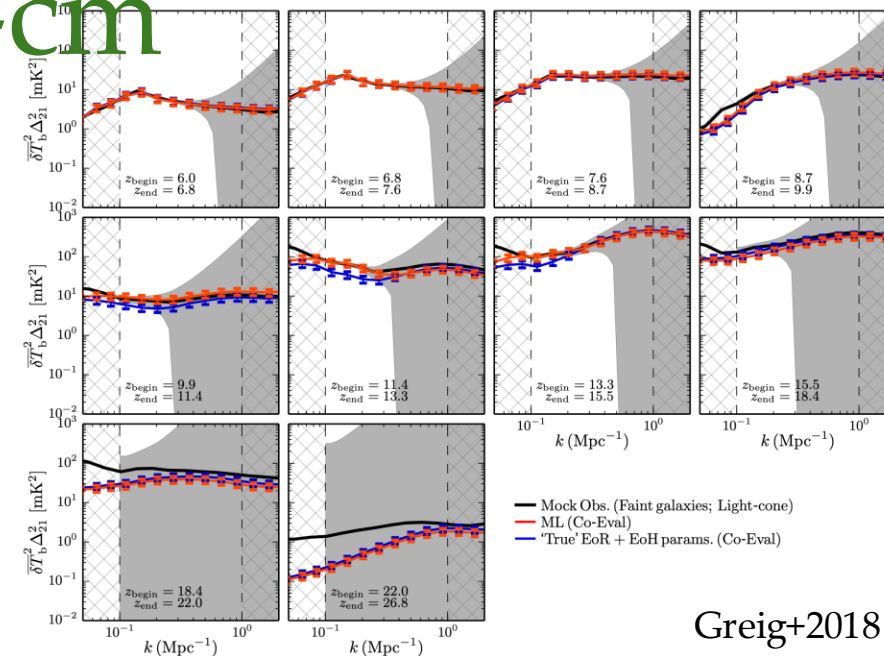- Simpler than a duck
  - Power spectrum



Greig+2018

# 4.3 Compression for the 21-cm

- ## Simpler than a duck
  - ### Power spectrum
  - ### Bispectrum



Greig+2018



Watkinson+2020

# 4.3 Compression for the 21-cm

- ## Simpler than a duck
  - ### Power spectrum
  - ### Bispectrum
  - ### Morphological spectra



Gazagnes+2020



Greig+2018



Watkinson+2020
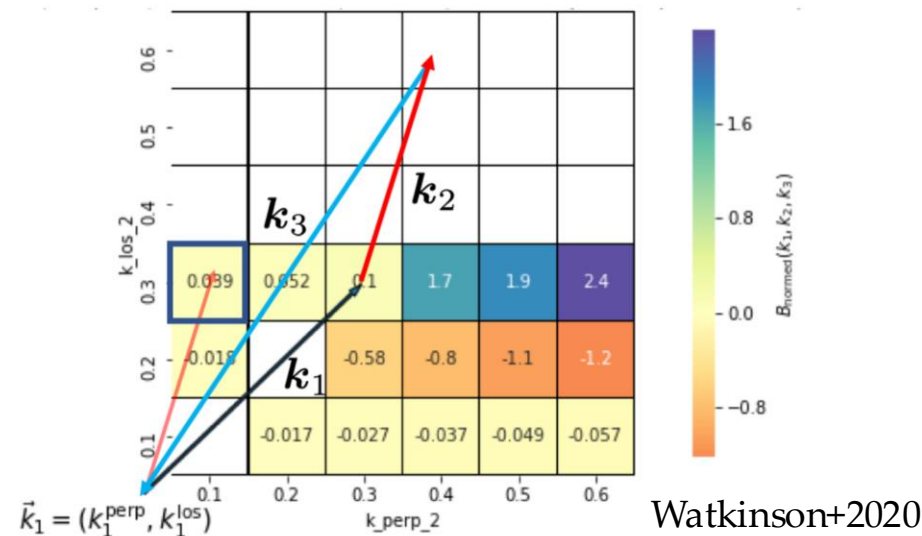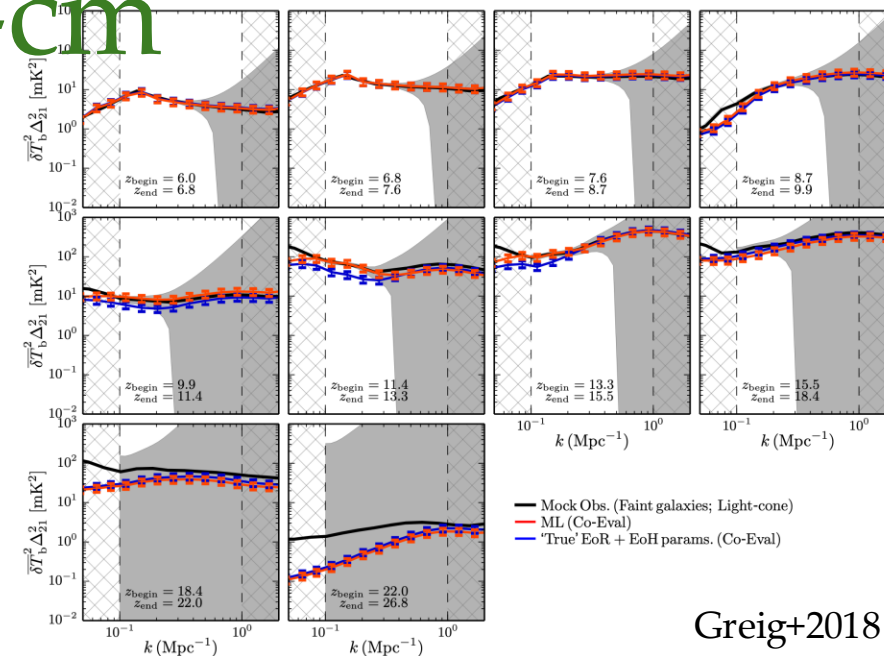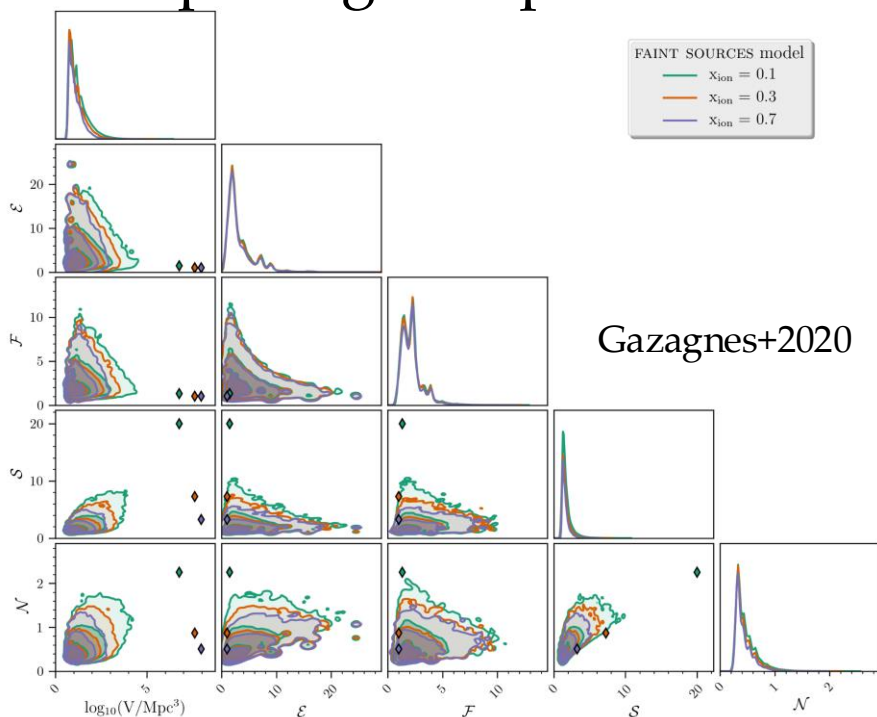
# 4.3 Compression for the 21-cm

- Simpler than a duck
  - Power spectrum
  - Bispectrum
  - Morphological spectra



Greig+2018

**Higher-order summaries**
**DO improve parameter inference!**

Gazagnes+2020

Watkinson+2020

# 5.1 ML role #1 - Compression

- 21-cm – no good a-priori physical motivation for a compression
- We cannot know THE optimal compression/summary

# 5.1 ML role #1 - Compression

- 21-cm – no good a-priori physical motivation for a compression
- We cannot know THE optimal compression/summary

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
- +++

- Joint space of

(parameters, summaries)

$$P(\mathcal{S}, \boldsymbol{\theta}) = P(\mathcal{S}|\boldsymbol{\theta}) \, P(\boldsymbol{\theta})$$

$P(\theta|S^*)$

$S$

$S^*$

$\theta^*$

$\theta$

$P(S|\theta^*)$

SCUOLA
NORMALE
SUPERIORE

$\theta^* \rightarrow$

$\mathcal{D}$

Initial Density Field

Lightcone

$\mathcal{S}$

Foreground wedge

Foreground wedge

z = 9.1

# 6. ML role #2 – Simulation Based Inference

- Joint space of

(parameters, summaries)

$$P(\mathcal{S}, \boldsymbol{\theta}) = P(\mathcal{S}|\boldsymbol{\theta}) \, P(\boldsymbol{\theta})$$

- Assumption – perfect data simulator
- Fitting the distribution with Neural Density Estimators (NDE)

$P(\theta|s^*)$

$S$

$s^*$

$\theta^*$

$\theta$

$P(s|\theta^*)$

# What is the likelihood of the 21-cm 1D power spectrum?

**Cosmology** **Simulation** **Astrophysics**

$\theta_{cosmo}$ $\theta_{astro}$

SCUOLA NORMALE SUPERIORE

**Observation**

$\theta_{foreground}$ $\theta_{telescope}$

**Foreground & systematic residuals**

**Telescope simulator**

Compression (1DPS)

Compression (1DPS)

$\mathcal{S}_{model}$ Likelihood $\mathcal{S}_{data}$

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

# 1. 1DPS has a non-Gaussian likelihood

**Gaussian data**
**=**
**Gaussian likelihood in the PS**

**Non-Gaussian data**
**=**
**Non-Gaussian likelihood, even in the PS**



Zhao+2022
Saxena+2023
Prelogović & Mesinger 2023

# 2. Classical inference (MCMC)

- Possible by approximating the PS likelihood with a Gaussian
  - Usually wrongly justified through the central limit theorem

$$P(\mathcal{S}|\boldsymbol{\theta}) = \mathcal{N}\left(\boldsymbol{\Sigma}_{\mathcal{S}}(\boldsymbol{\theta}), \boldsymbol{\mu}_{\mathcal{S}}(\boldsymbol{\theta})\right)$$

$$= \frac{1}{(2\pi)^{n/2}\sqrt{|\boldsymbol{\Sigma}_{\mathcal{S}}(\boldsymbol{\theta})|}} e^{-\frac{1}{2}(\mathcal{S}-\boldsymbol{\mu}_{\mathcal{S}}(\boldsymbol{\theta}))^T \boldsymbol{\Sigma}_{\mathcal{S}}^{-1}(\boldsymbol{\theta})(\mathcal{S}-\boldsymbol{\mu}_{\mathcal{S}}(\boldsymbol{\theta}))}$$

# 2. Classical inference (MCMC)

- Possible by approximating the PS likelihood with a Gaussian
  - Usually wrongly justified through the central limit theorem

$$P(\mathcal{S}|\boldsymbol{\theta}) = \mathcal{N}\left(\boldsymbol{\Sigma}_{\mathcal{S}}(\boldsymbol{\theta}), \boldsymbol{\mu}_{\mathcal{S}}(\boldsymbol{\theta})\right)$$

$$= \frac{1}{(2\pi)^{n/2}\sqrt{|\boldsymbol{\Sigma}_{\mathcal{S}}(\boldsymbol{\theta})|}} e^{-\frac{1}{2}(\mathcal{S}-\boldsymbol{\mu}_{\mathcal{S}}(\boldsymbol{\theta}))^T \boldsymbol{\Sigma}_{\mathcal{S}}^{-1}(\boldsymbol{\theta})(\mathcal{S}-\boldsymbol{\mu}_{\mathcal{S}}(\boldsymbol{\theta}))}$$

Prelogović&Mesinger 2023

- Common additional simplifications
  1) ignoring correlations by using diagonal Σ

# 2. Classical inference (MCMC)

- Possible by approximating the PS likelihood with a Gaussian
  - Usually wrongly justified through the central limit theorem

$$P(\mathcal{S}|\boldsymbol{\theta}) = \mathcal{N}\left(\boldsymbol{\Sigma}_{\mathcal{S}}(\boldsymbol{\theta}), \boldsymbol{\mu}_{\mathcal{S}}(\boldsymbol{\theta})\right)$$

$$= \frac{1}{(2\pi)^{n/2}\sqrt{|\boldsymbol{\Sigma}_{\mathcal{S}}(\boldsymbol{\theta})|}} e^{-\frac{1}{2}(\mathcal{S}-\boldsymbol{\mu}_{\mathcal{S}}(\boldsymbol{\theta}))^T \boldsymbol{\Sigma}_{\mathcal{S}}^{-1}(\boldsymbol{\theta})(\mathcal{S}-\boldsymbol{\mu}_{\mathcal{S}}(\boldsymbol{\theta}))}$$

- Common additional simplifications
  1) ignoring correlations by using diagonal $\Sigma$
  2) Fixing the covariance at fiducial parameters $\Sigma = \Sigma_{\theta\text{fid}}$

Greig&Mesinger 2018
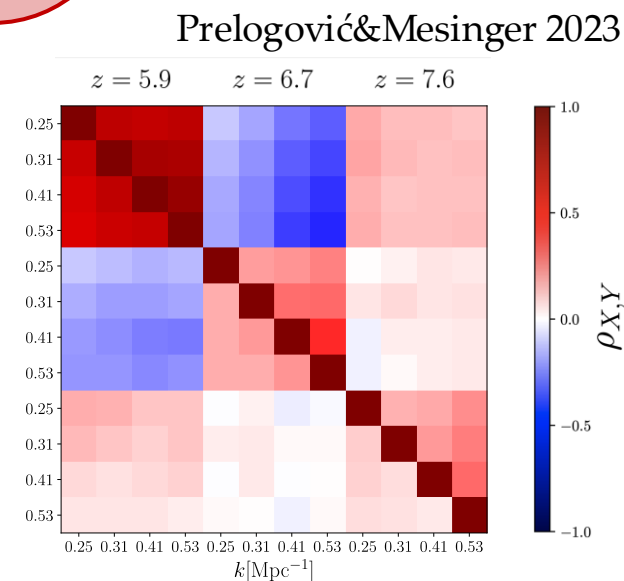Trott+2020
Mertens+2020
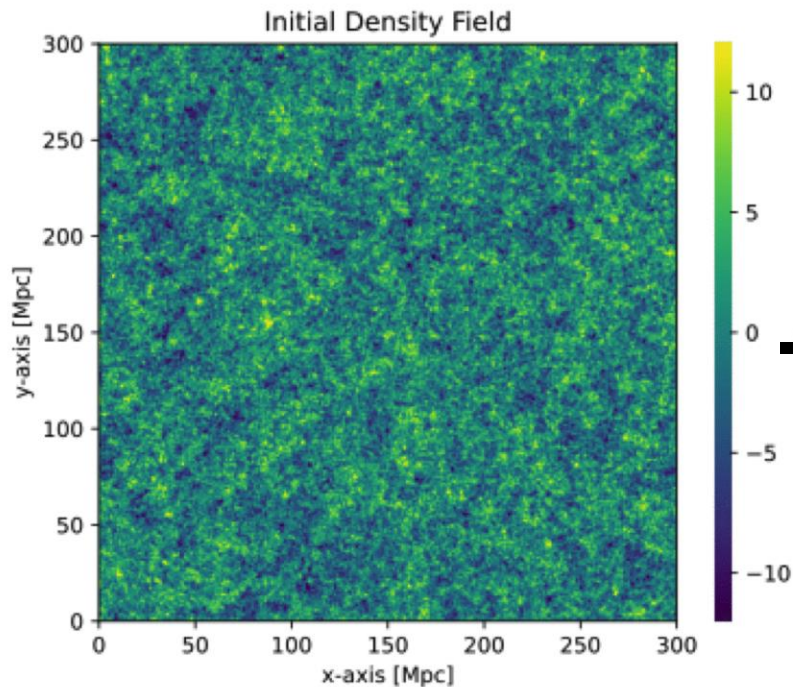HERA+2023

# 2. Classical inference (MCMC)

- Possible by approximating the PS likelihood with a Gaussian
  - Usually wrongly justified through the central limit theorem

$$P(\mathcal{S}|\boldsymbol{\theta}) = \mathcal{N}\left(\boldsymbol{\Sigma}_{\mathcal{S}}(\boldsymbol{\theta}), \boldsymbol{\mu}_{\mathcal{S}}(\boldsymbol{\theta})\right)$$

$$= \frac{1}{(2\pi)^{n/2}\sqrt{|\boldsymbol{\Sigma}_{\mathcal{S}}(\boldsymbol{\theta})|}} e^{-\frac{1}{2}(\mathcal{S}-\boldsymbol{\mu}_{\mathcal{S}}(\boldsymbol{\theta}))^T \boldsymbol{\Sigma}_{\mathcal{S}}^{-1}(\boldsymbol{\theta})(\mathcal{S}-\boldsymbol{\mu}_{\mathcal{S}}(\boldsymbol{\theta}))}$$

- Common additional simplifications
  1) ignoring correlations by using diagonal $\Sigma$
  2) Fixing the covariance at fiducial parameters $\Sigma = \Sigma_{\theta\text{fid}}$
  3) $\mu$ estimated from one simulation

Greig&Mesinger 2018
Trott+2020
Mertens+2020
HERA+2023

SCUOLA
NORMALE
SUPERIORE

$\theta^*$ →

Initial Density Field

$\mathcal{D}$

Lightcone

$\mathcal{S}$

Foreground wedge

Foreground wedge    z = 9.1

# 3. Simulation Based Inference

# 3. Simulation Based Inference

- Train a neural density estimator (NDE)
  - Gaussian mixture



$P(\mathcal{S}|\boldsymbol{\theta})$

Repin+2021

# 4. Results

**Including more
realistic likelihood
≠
more constraining
posterior**



—— NDE Gauss mixture

- - - CLASSIC fixed cov

—— CLASSIC fixed var

SCUOLA
NORMALE
SUPERIORE

# 4. Results

**Including more
realistic likelihood
≠
more constraining
posterior**



| | |
|---|---|
| —— NDE Gauss mixture | |
| – – CLASSIC fixed cov | |
| —— CLASSIC fixed var | |

**0.1-1 CPU h**

**100 000 CPU h**

- **Amortized inference**
- **x 2-3 smaller training DB vs.
one CLASSIC run**

# 4. Results

**Including more realistic likelihood**

**≠**

**more constraining posterior**

- **Amortized inference**
- **x 2-3 smaller training DB vs. one CLASSIC run**

**Varying covariance performs well**

Legend (left panel):
- NDE Gauss mixture
- CLASSIC fixed cov
- CLASSIC fixed var

Legend (right panel):
- NDE Gauss mixture
- NDE varying cov
- NDE varying var

Parameters: $f_{*,10}$, $f_{\mathrm{esc},10}$, $M_{\mathrm{turn}}$, $L_X/\mathrm{SFR}$, $E_0$

# 4. Results

***BUT:***
This is only qualitative description, and only for the <span style="color:red">mock observation</span>

- How does it perform for other points in the parameter space?
- Did the training converge?
- Can we quantify the best model?

**–> Simulation Based Calibration**

# 5. Simulation Based Calibration (SBC)

- *"prior" = "data averaged posterior"* $\quad P(\boldsymbol{\theta}) = \int P(\boldsymbol{\theta}|\tilde{y}) \, P(\tilde{y}|\tilde{\boldsymbol{\theta}}) P(\tilde{\boldsymbol{\theta}}) \, \mathrm{d}\tilde{y} \, \mathrm{d}\tilde{\boldsymbol{\theta}}$

# 5. Simulation Based Calibration (SBC)

- *"prior" = "data averaged posterior"* $\quad P(\boldsymbol{\theta}) = \int P(\boldsymbol{\theta}|\tilde{y}) \, P(\tilde{y}|\tilde{\boldsymbol{\theta}}) P(\tilde{\boldsymbol{\theta}}) \, \mathrm{d}\tilde{y} \, \mathrm{d}\tilde{\boldsymbol{\theta}}$

  1. Pull from prior $\qquad\qquad\qquad \tilde{\boldsymbol{\theta}} \sim P(\boldsymbol{\theta})$

# 5. Simulation Based Calibration (SBC)

- *"prior" = "data averaged posterior"*   $P(\boldsymbol{\theta}) = \int P(\boldsymbol{\theta}|\tilde{\boldsymbol{y}})\, P(\tilde{\boldsymbol{y}}|\tilde{\boldsymbol{\theta}})P(\tilde{\boldsymbol{\theta}})\,\mathrm{d}\tilde{\boldsymbol{y}}\,\mathrm{d}\tilde{\boldsymbol{\theta}}$

1. Pull from prior   $\tilde{\boldsymbol{\theta}} \sim P(\boldsymbol{\theta})$

2. Pull the data from the likelihood   $\tilde{\boldsymbol{y}} \sim P(\boldsymbol{y}|\tilde{\boldsymbol{\theta}}) \quad \Leftrightarrow \quad \tilde{\boldsymbol{y}} = \mathrm{simulator}(\tilde{\boldsymbol{\theta}})$

# 5. Simulation Based Calibration (SBC)

- **"prior" = "data averaged posterior"** $\quad P(\boldsymbol{\theta}) = \int P(\boldsymbol{\theta}|\tilde{\boldsymbol{y}}) \, P(\tilde{\boldsymbol{y}}|\tilde{\boldsymbol{\theta}}) P(\tilde{\boldsymbol{\theta}}) \, \mathrm{d}\tilde{\boldsymbol{y}} \, \mathrm{d}\tilde{\boldsymbol{\theta}}$

1. Pull from prior $\qquad\qquad\qquad\qquad\qquad\qquad \tilde{\boldsymbol{\theta}} \sim P(\boldsymbol{\theta})$

2. Pull the data from the likelihood $\qquad \tilde{\boldsymbol{y}} \sim P(\boldsymbol{y}|\tilde{\boldsymbol{\theta}}) \quad \Leftrightarrow \quad \tilde{\boldsymbol{y}} = \mathrm{simulator}(\tilde{\boldsymbol{\theta}})$

3. Calculate the posterior the sample $\quad P(\boldsymbol{\theta}|\tilde{\boldsymbol{y}})$

# 5. Simulation Based Calibration (SBC)

- *"prior" = "data averaged posterior"* $\quad P(\boldsymbol{\theta}) = \int P(\boldsymbol{\theta}|\tilde{\boldsymbol{y}})\, P(\tilde{\boldsymbol{y}}|\tilde{\boldsymbol{\theta}}) P(\tilde{\boldsymbol{\theta}})\, \mathrm{d}\tilde{\boldsymbol{y}}\, \mathrm{d}\tilde{\boldsymbol{\theta}}$

1. Pull from prior $\qquad\qquad\qquad\qquad\qquad \tilde{\boldsymbol{\theta}} \sim P(\boldsymbol{\theta})$

2. Pull the data from the likelihood $\qquad \tilde{\boldsymbol{y}} \sim P(\boldsymbol{y}|\tilde{\boldsymbol{\theta}}) \quad \Leftrightarrow \quad \tilde{\boldsymbol{y}} = \mathrm{simulator}(\tilde{\boldsymbol{\theta}})$

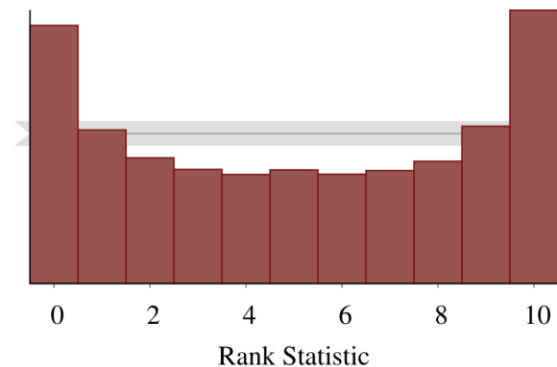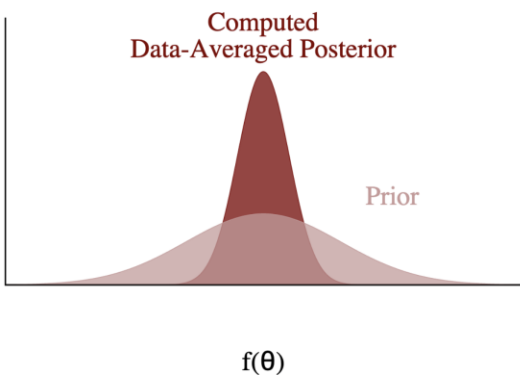3. Calculate the posterior the sample $\quad P(\boldsymbol{\theta}|\tilde{\boldsymbol{y}})$

4. Repeat and average posteriors $\qquad P(\boldsymbol{\theta}) \approx \dfrac{1}{N} \sum_{i=1}^{N} P_i(\boldsymbol{\theta}|\tilde{\boldsymbol{y}}_i)$
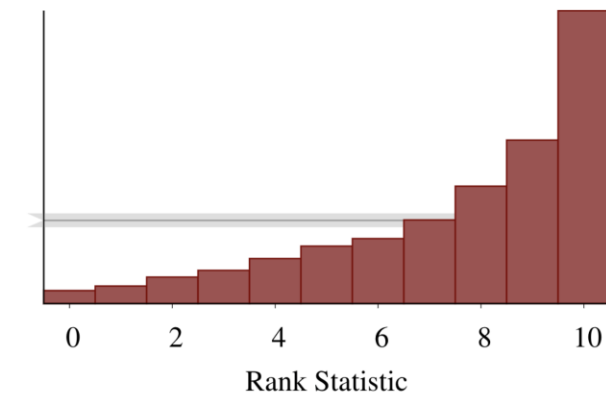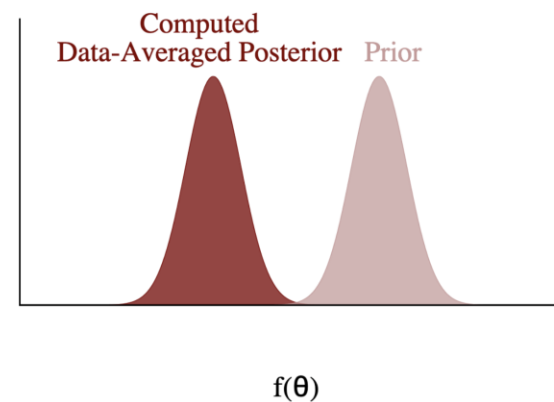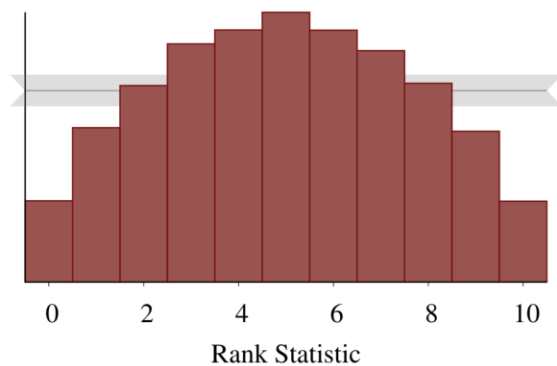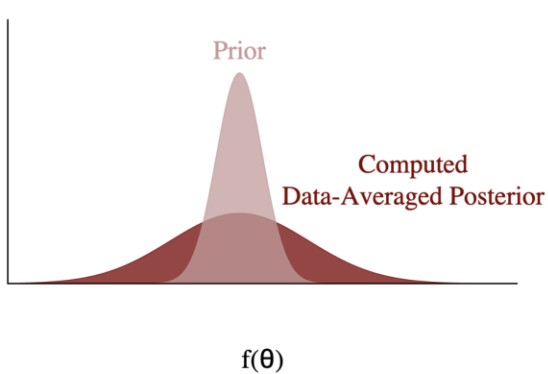
# 5. Simulation Based Calibration (SBC)

- **"prior" = "data averaged posterior"**   $P(\boldsymbol{\theta}) = \int P(\boldsymbol{\theta}|\tilde{y}) \, P(\tilde{y}|\tilde{\boldsymbol{\theta}}) P(\tilde{\boldsymbol{\theta}}) \, \mathrm{d}\tilde{y} \, \mathrm{d}\tilde{\boldsymbol{\theta}}$

- SBC – casting integral into 1D rank statistics distribution



Talts+2018

# 6. SBC for 21-cm PS

- 10 000 posteriors

- <span style="color:red">Would be useful for classic inference</span>, but is too expensive to compute

- NDE Gauss mixture – the best



NDE CMAF

$f_{*,10}$
$f_{\mathrm{esc},10}$
$M_{\mathrm{turn}}$
$L_X/\mathrm{SFR}$
$E_0$

Rank Statistics

NDE Gauss mixture

$f_{*,10}$
$f_{\mathrm{esc},10}$
$M_{\mathrm{turn}}$
$L_X/\mathrm{SFR}$
$E_0$

Rank Statistics

Prelogović & Mesinger 2023

# Conclusions

- SBI – current and future frontier in the 21-cm inference
    - Cheaper and more precise, by recovering a data-driven likelihood
    - Convergence / performance tests crucial!

# How informative are summaries of the 21-cm signal?

# 1. Fisher information matrix

- If we label data space as $\boldsymbol{d}$ and its likelihood as $P(\boldsymbol{d}|\boldsymbol{\theta})$

$$\boldsymbol{F}(\boldsymbol{\theta}^*)_{mn} = \mathrm{E}_{P(\boldsymbol{d}|\boldsymbol{\theta}^*)} \left[ \frac{\partial}{\partial \boldsymbol{\theta}_m} \ln P(\boldsymbol{d}|\boldsymbol{\theta}^*) \cdot \frac{\partial}{\partial \boldsymbol{\theta}_n} \ln P(\boldsymbol{d}|\boldsymbol{\theta}^*) \right]$$
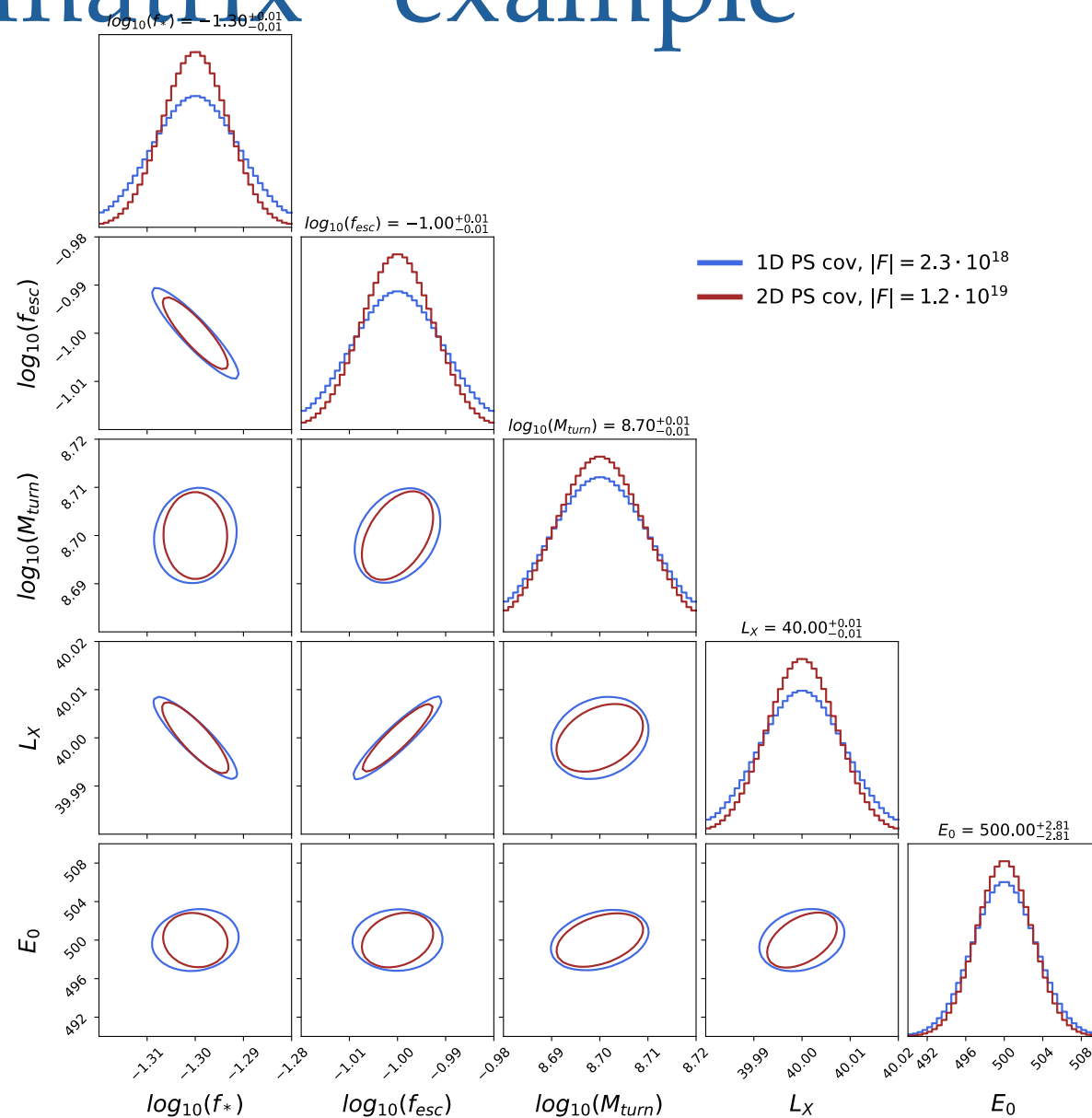
- The usefulness comes from

1D: $\mathrm{Var}\left(\hat{\boldsymbol{\theta}}_m\right) \geq (\boldsymbol{F}^{-1})_{mm}$    ND: $\det \mathrm{Cov}(\hat{\boldsymbol{\theta}}) \geq \det \boldsymbol{F}^{-1}$

*How well we can estimate a parameter is fundamentally limited by its Fisher information.*

*(i.e. one cannot go below it)*

Fisher 1935

# 1. Fisher information matrix - example

- We cannot perform better than the shown ellipse

- Different summary, different Fisher matrix

- det $\mathbf{F}^{-1}$ = volume of the ellipse
  - det $\mathbf{F}^{-1}$ smaller the better
  - det $\mathbf{F}$ bigger the better

# 3. Distribution of the Fisher information

- $\det \boldsymbol{F}(\boldsymbol{\theta}^*)$ is information measure just around one point
- Calculating around many different points is better
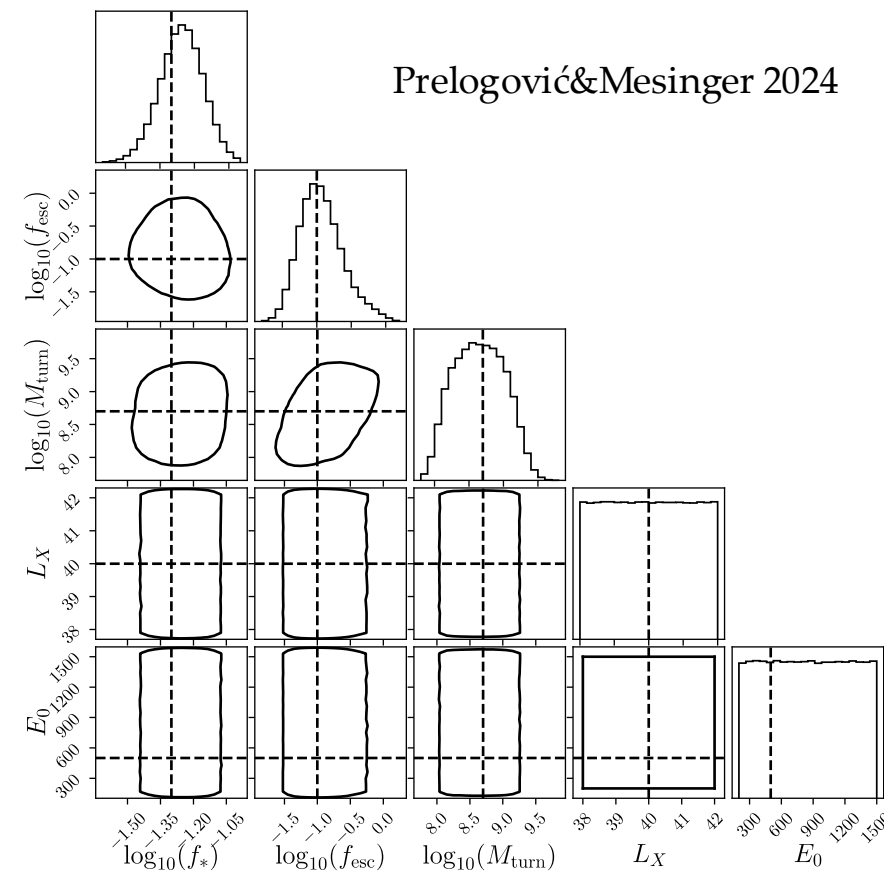
# 3. Distribution of the Fisher information

- $\det \boldsymbol{F}(\boldsymbol{\theta}^*)$ is information measure just around one point
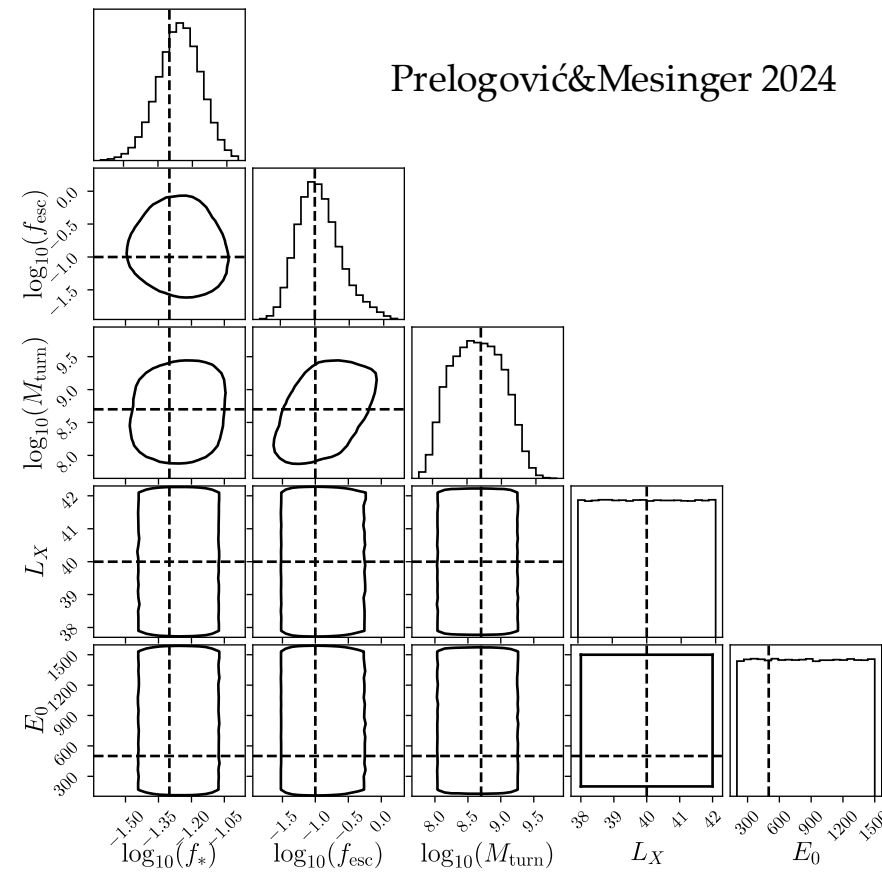
- Calculating around many different points is better

- Sample ~150 points from the prior

- Around each point construct simulations. needed to compute the Fisher matrix

Prelogović&Mesinger 2024

# 3. Distribution of the Fisher information

- $\det \boldsymbol{F}(\boldsymbol{\theta}^*)$ is information measure just around one point
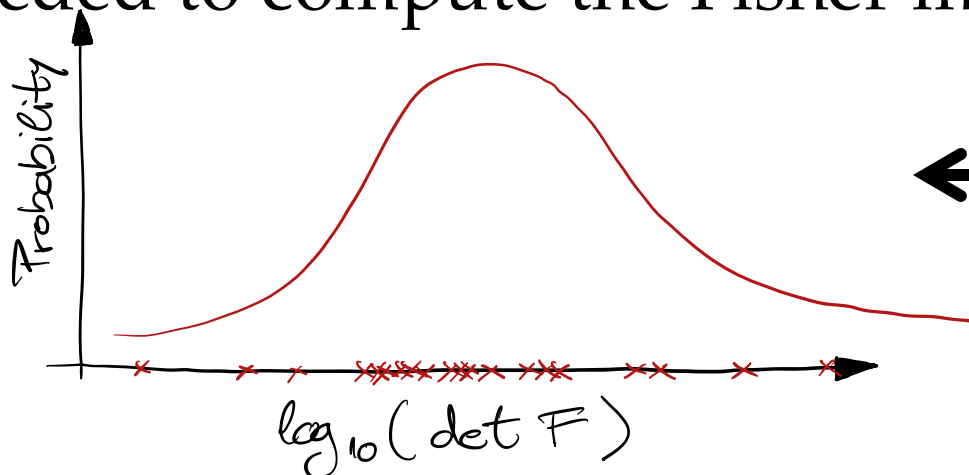- Calculating around many different points is better

- Sample ~150 points from the prior
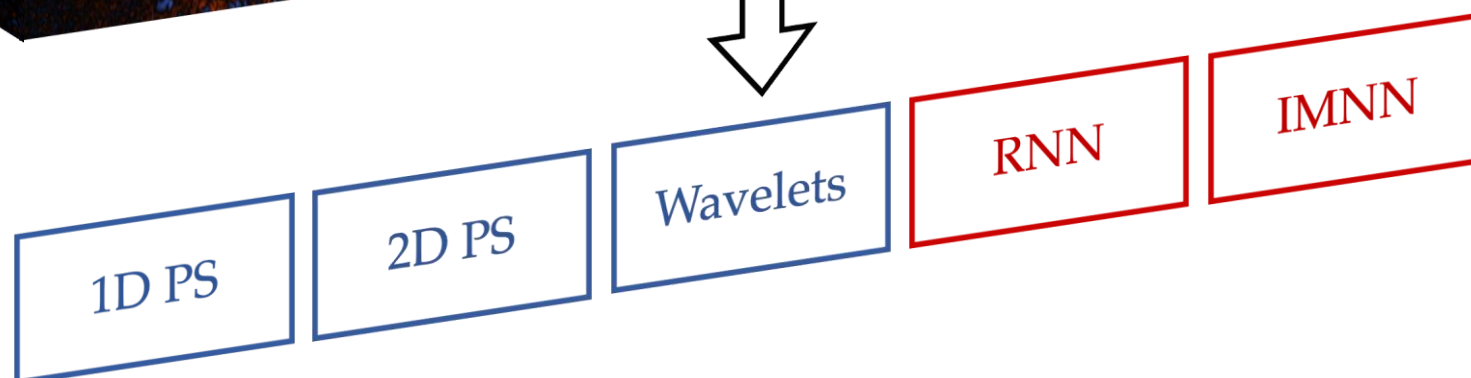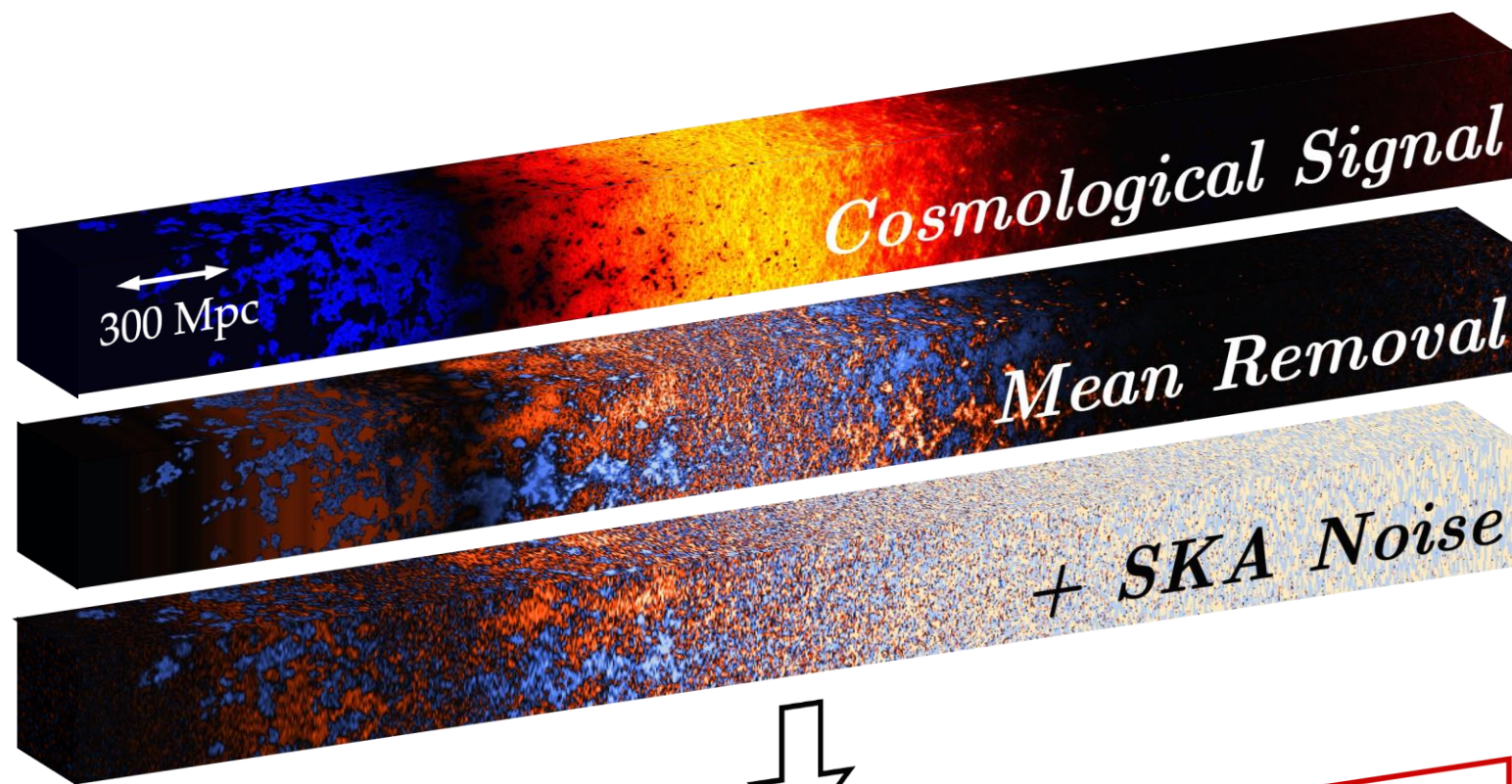- Around each point construct simulations. needed to compute the Fisher matrix

Prelogović&Mesinger 2024

Ce Sui+2023

**Later talk!**

# 4. Considered summaries



300 Mpc

Cosmological Signal

Mean Removal
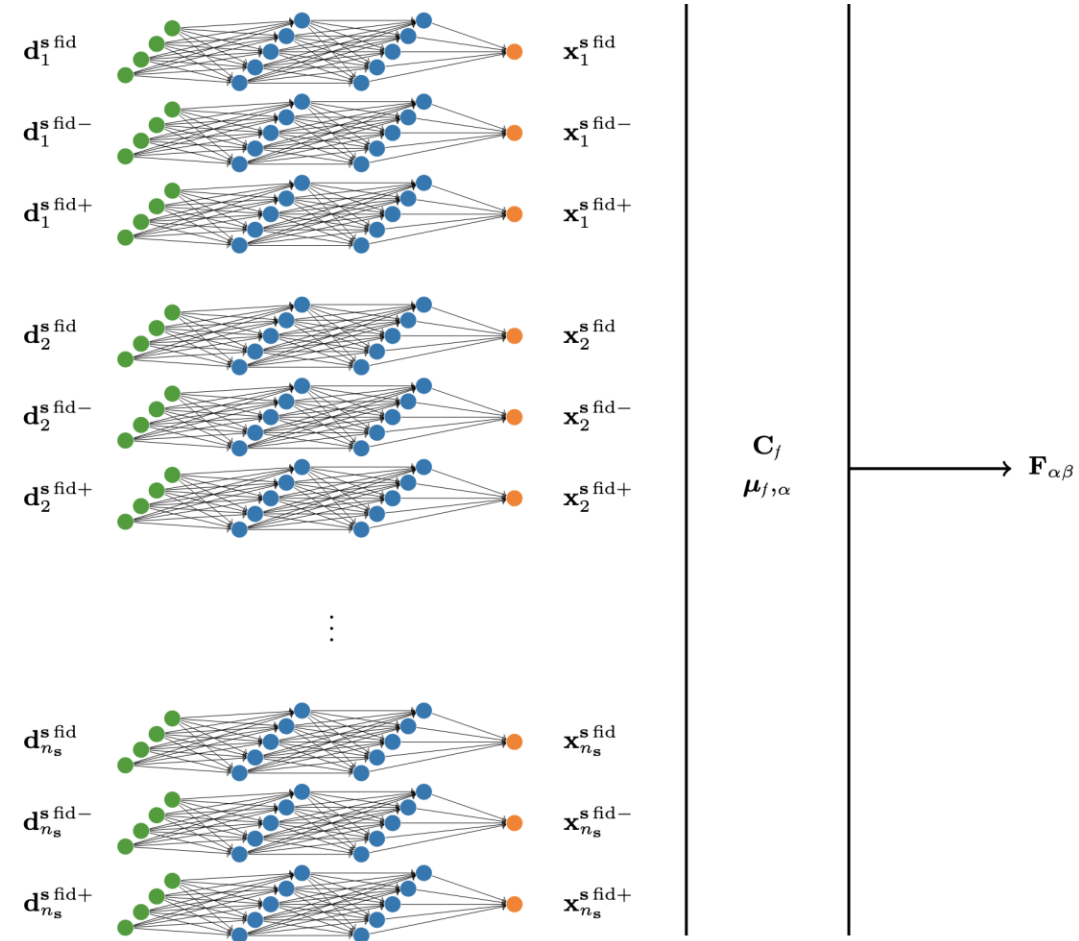
+ SKA Noise

1D PS | 2D PS | Wavelets | RNN | IMNN

# 4.1 Information Maximizing NN

- Unsupervised algorithm

- Simulate the data at a fiducial parameter set: $\mathbf{d}(\boldsymbol{\theta}_{\text{fid}})$

- Simulate around the fiducial parameters: $\mathbf{d}(\boldsymbol{\theta}_{\text{fid}}^{+}), \mathbf{d}(\boldsymbol{\theta}_{\text{fid}}^{-})$
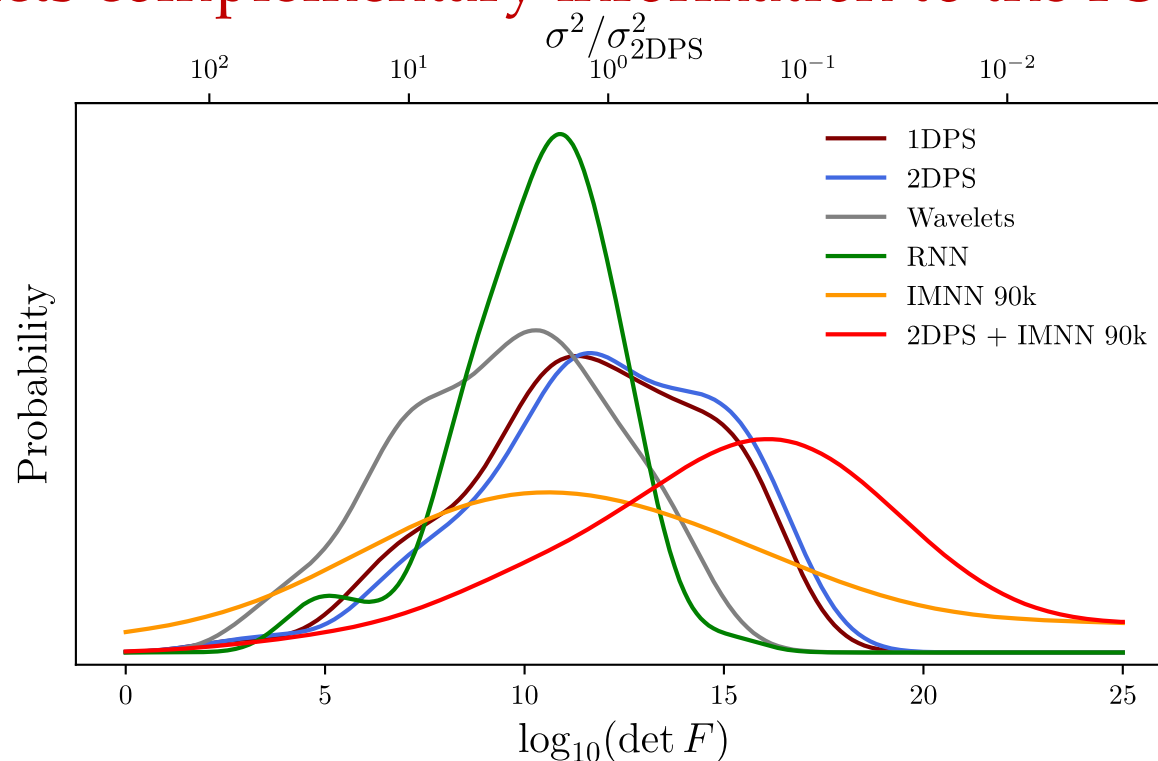
- Calculate compressed summary:

$$\boldsymbol{s} = NN(\boldsymbol{d})$$

- Maximize Fisher information:

$$\mathcal{L} = -\ln(\det \boldsymbol{F})$$



Charnock+2018
Makinen+2023
Prelogović+2024
Maitra+2024

# 5. Results

- 1DPS and 2DPS clear winners

- Combining 2DPS + IMNN
    - IMNN extracts complementary information to the PS

# Conclusions

- SBI – current and future frontier in the 21-cm inference
  - Cheaper and more precise, by recovering a data-driven likelihood
  - Convergence / performance tests crucial!


- Fisher distribution – information-based metric for a summary quality
  - Hard to beat the PS
  - Combination of classical + neural summaries as a powerful way forward

# Thank you!