

Improvements in the exploration of the CD-EoR in the era of AI



Stockholm
University

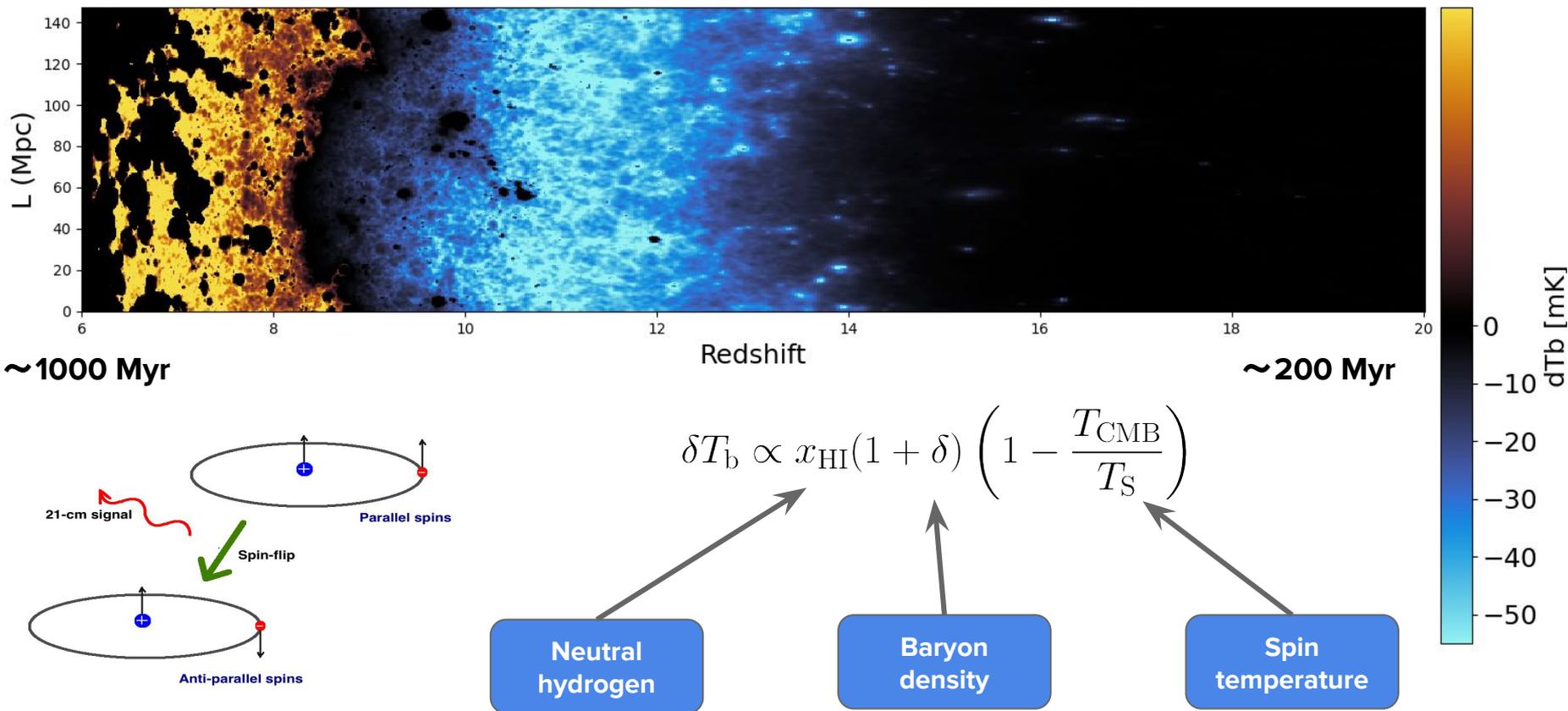
Sambit Giri
NORDITA fellow



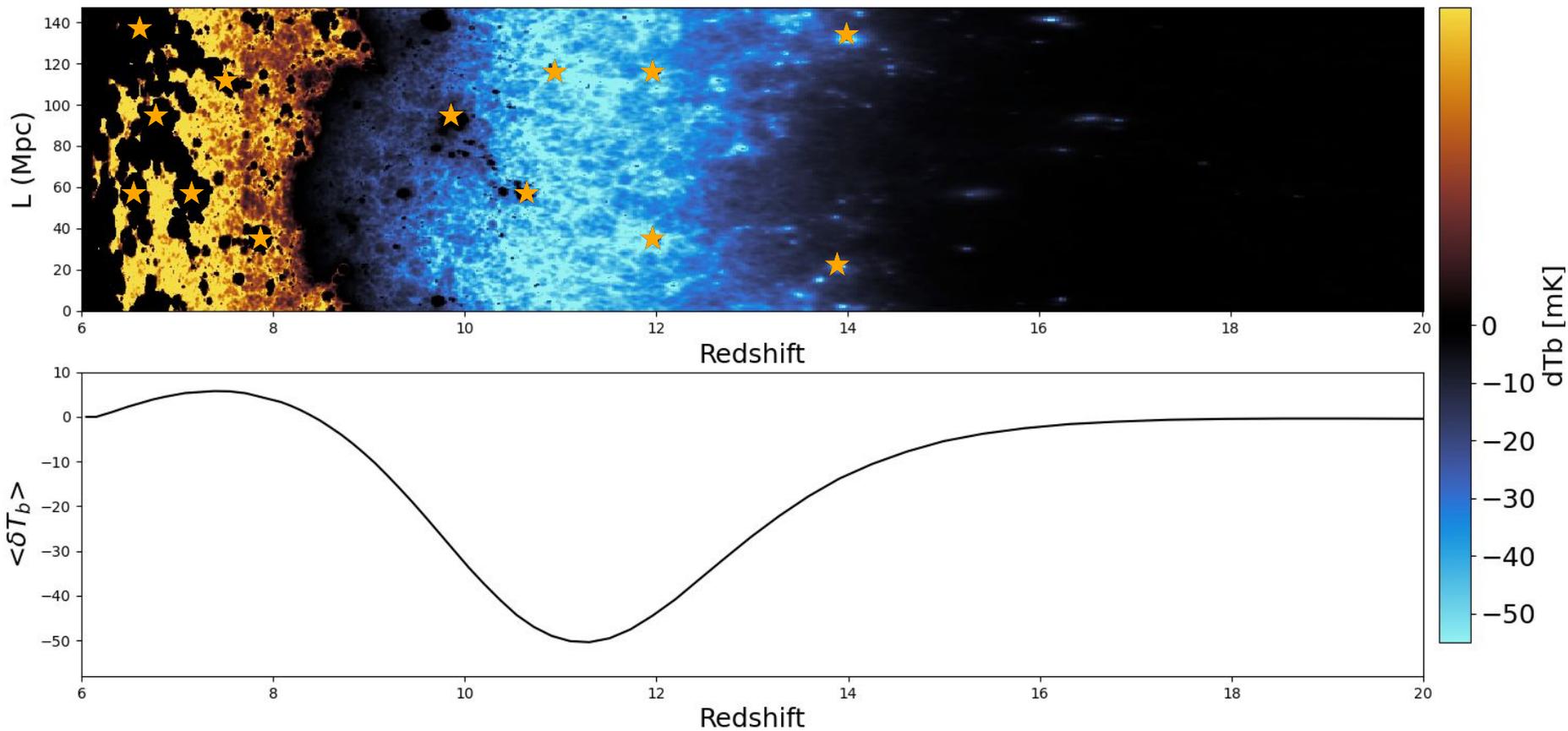
NORDITA

18 July 2024

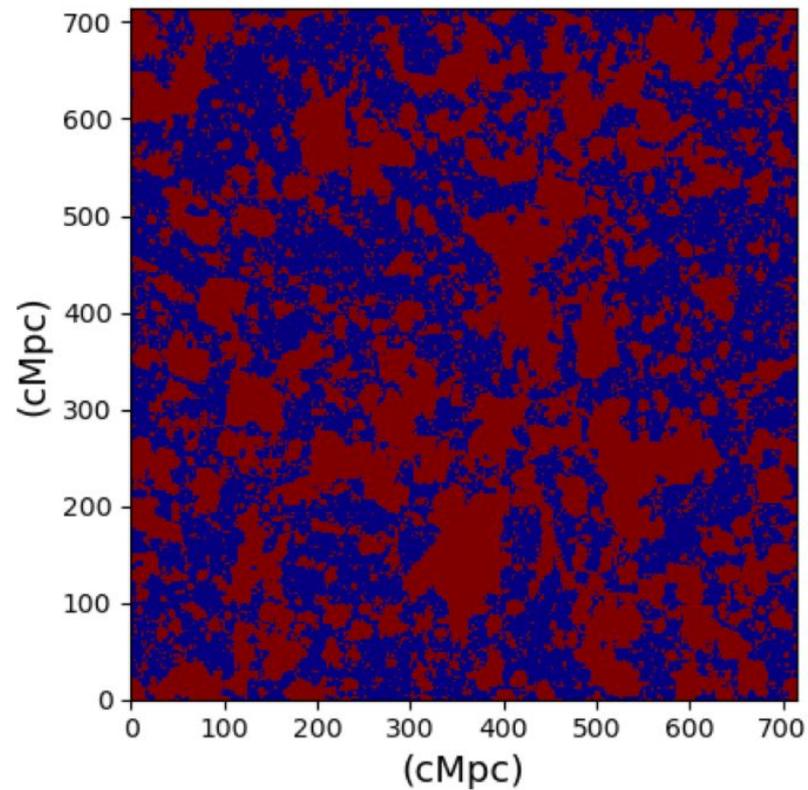
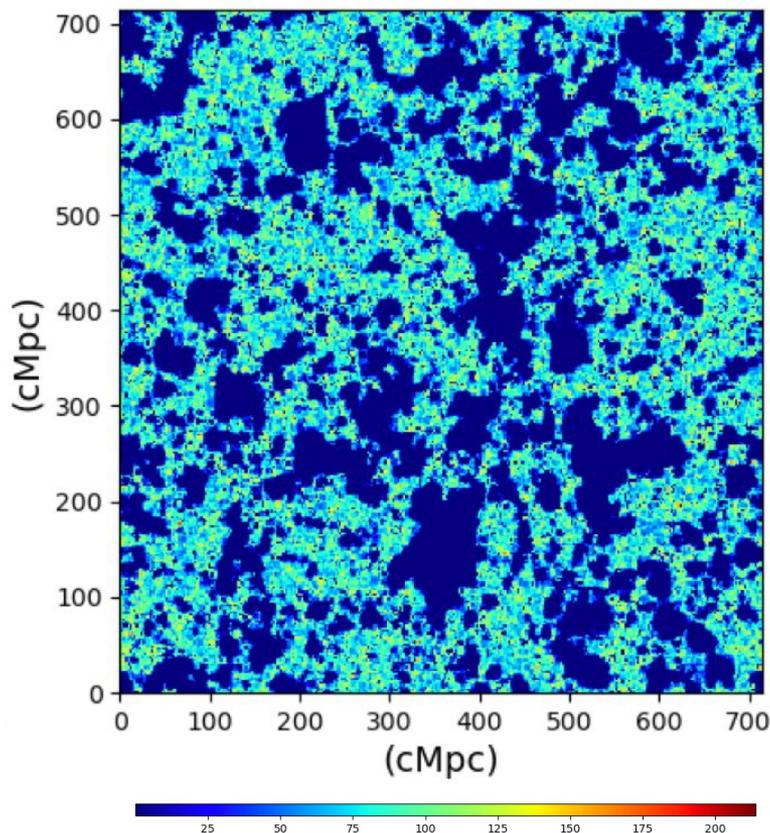
The 21-cm signal will probe the intergalactic medium



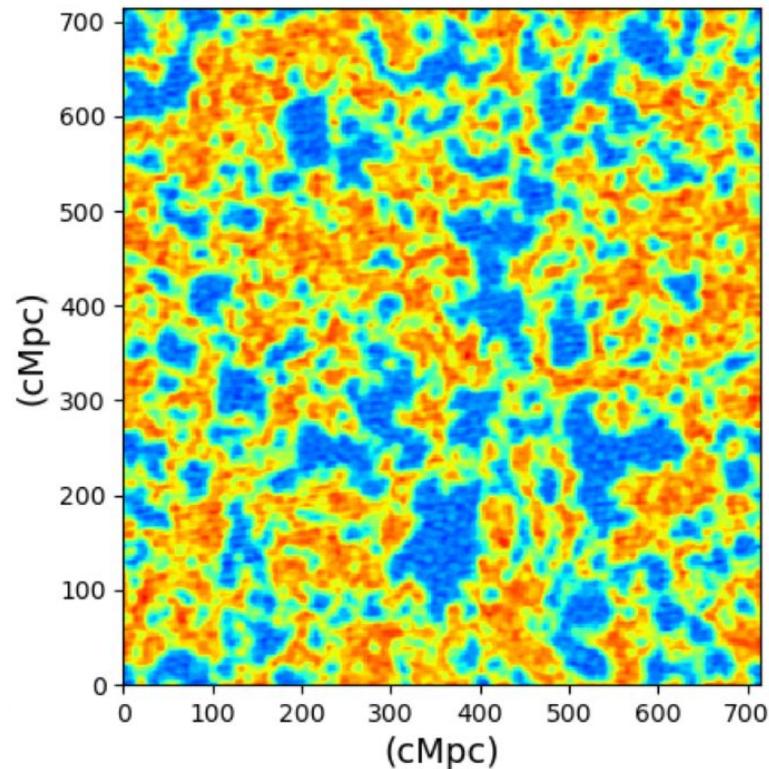
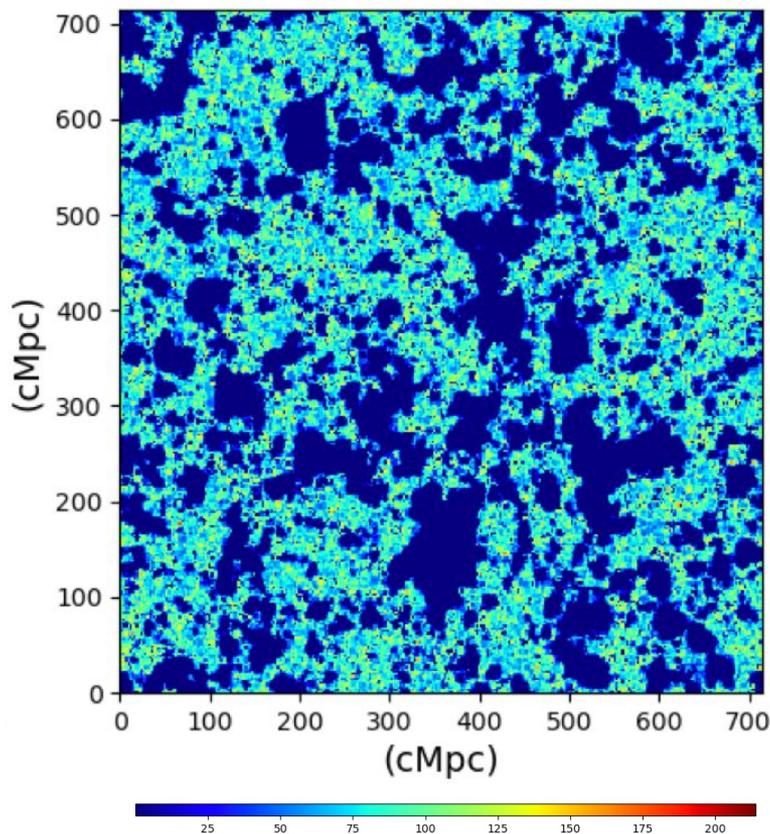
Evolution of the signal is driven by early luminous sources



Ionized regions are imprinted onto the 21-cm signal



Simulated SKA-Low image



How do we study ionized regions in SKA images?

Survey	Data per night/day	Galaxies	Cost	Scientists
DES	1 TeraB	~300 Million (all observed)	~\$70M	~400
DESI	40 GigaB	~35 Million (12M observed)	~\$70M	~600
Rubin-LSST	15 TeraB	~Billions	~\$1.0B	~1000
Euclid	850 GigaB	~Billions	~\$1.5B	~1500
SKA	1 PetaB	~Billions	~\$1.3B	~1000

Big Data!

(Lahav 2023)



Low frequency
component of Square
Kilometre Array
(SKA-Low)

Methods to automatically identify meaningful features

Survey	Data per night/day	Galaxies	Cost	Scientists
DES	1 TeraB	~300 Million (all observed)	~\$70M	~400
DESI	40 GigaB	~35 Million (12M observed)	~\$70M	~600
Rubin-LSST	15 TeraB	~Billions	~\$1.0B	~1000
Euclid	850 GigaB	~Billions	~\$1.5B	~1500
SKA	1 PetaB	~Billions	~\$1.3B	~1000

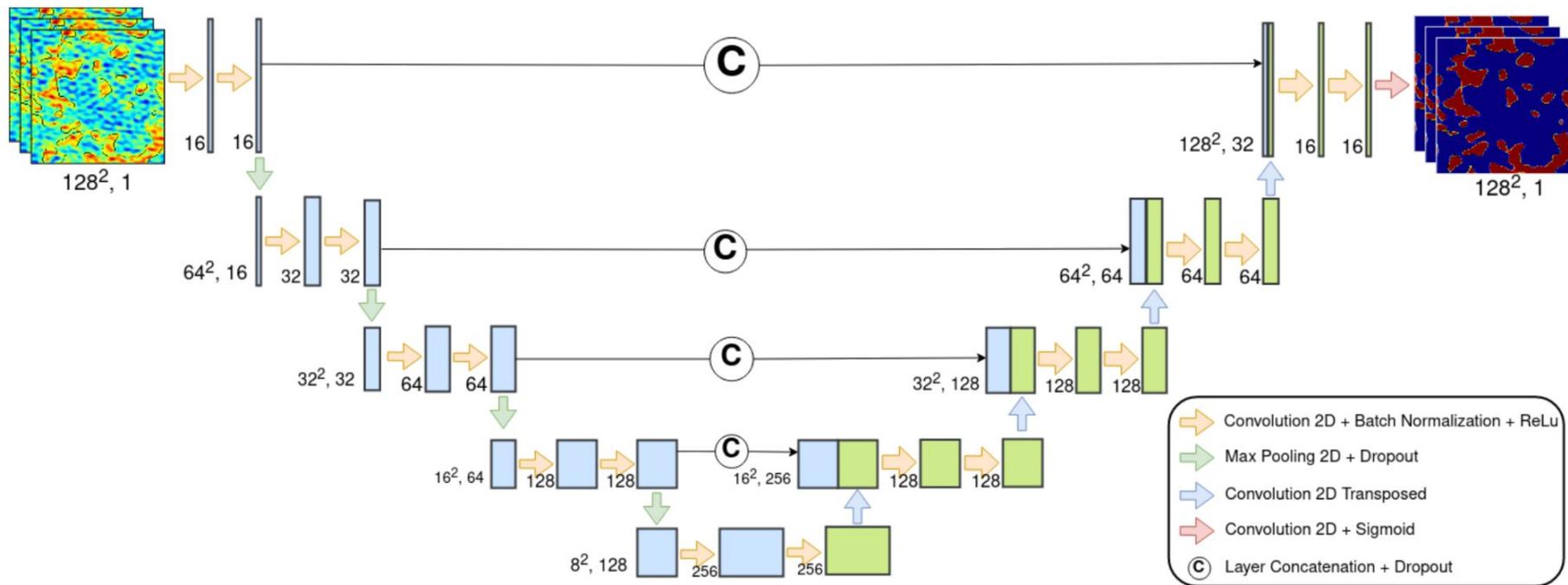
(Lahav 2023)

- **Superpixel method**
Image processing (**Giri+2018b**)
- **SegU-Net**
Deep learning (**Bianco, Giri+2021;2024**)



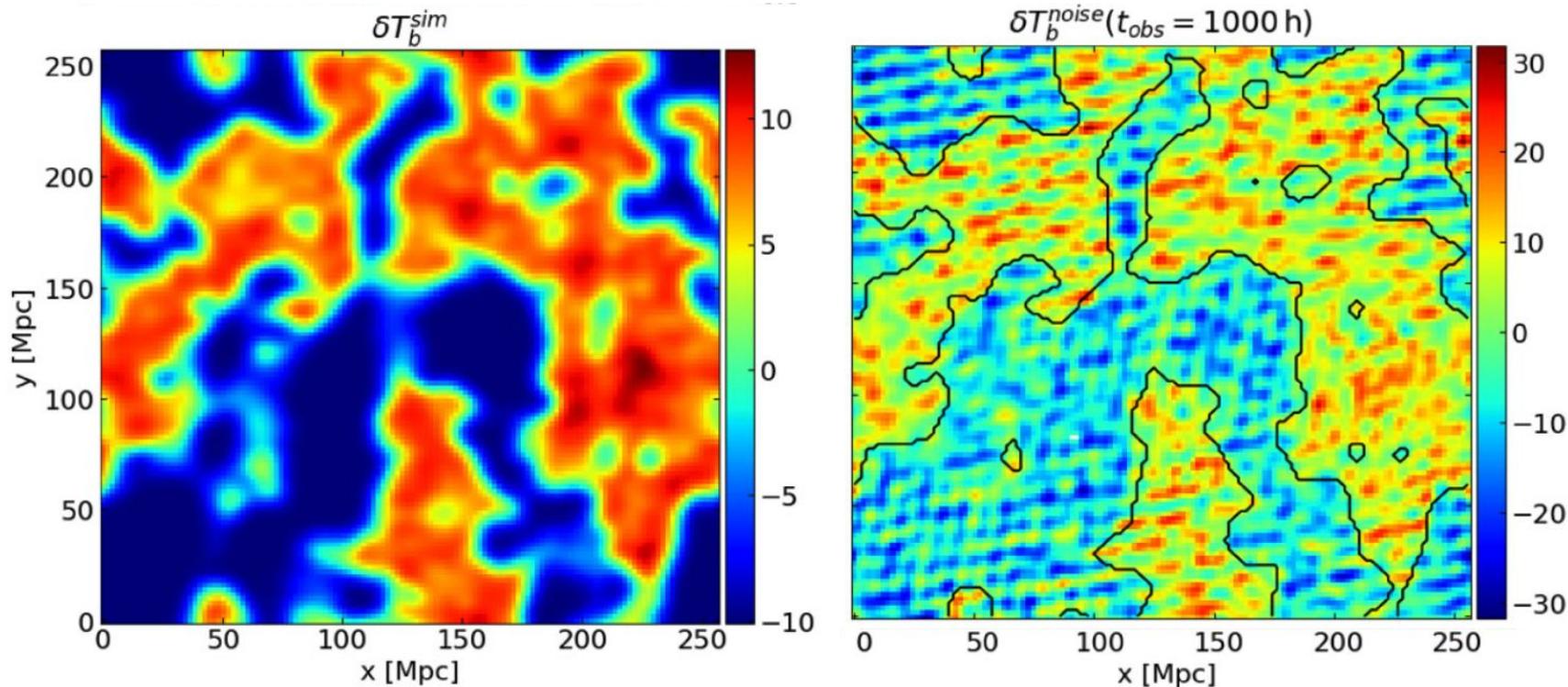
Low frequency
component of Square
Kilometre Array
(SKA-Low)

SegU-Net: deep learning segmentation model

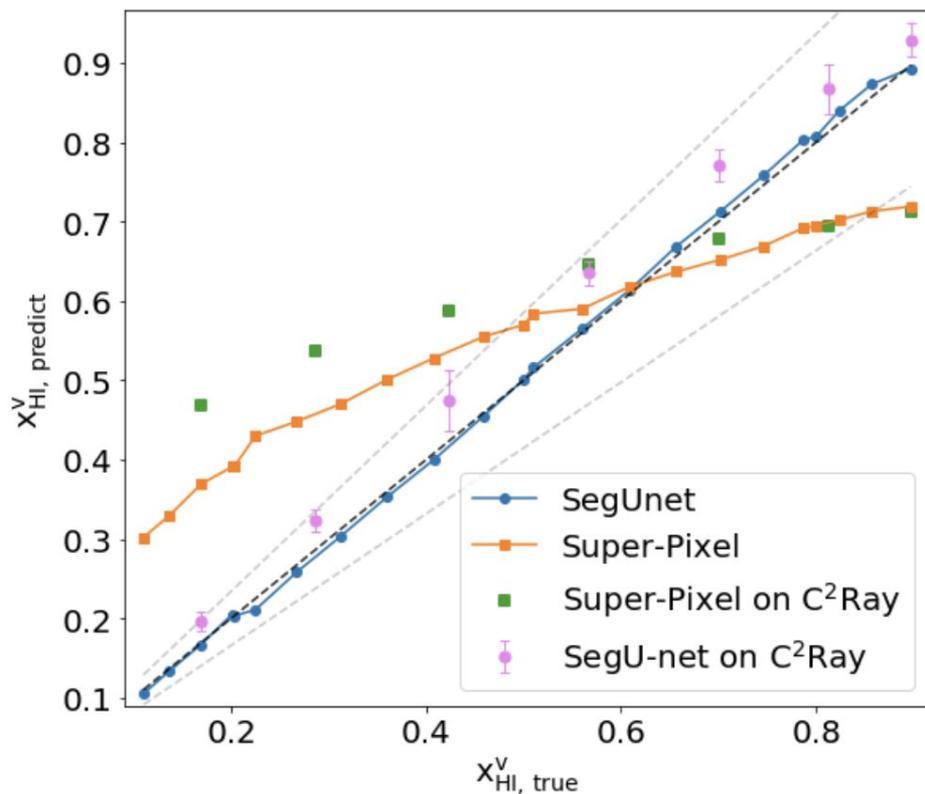


(Bianco, Giri+2021;2024)

SegU-Net identifying ionized regions

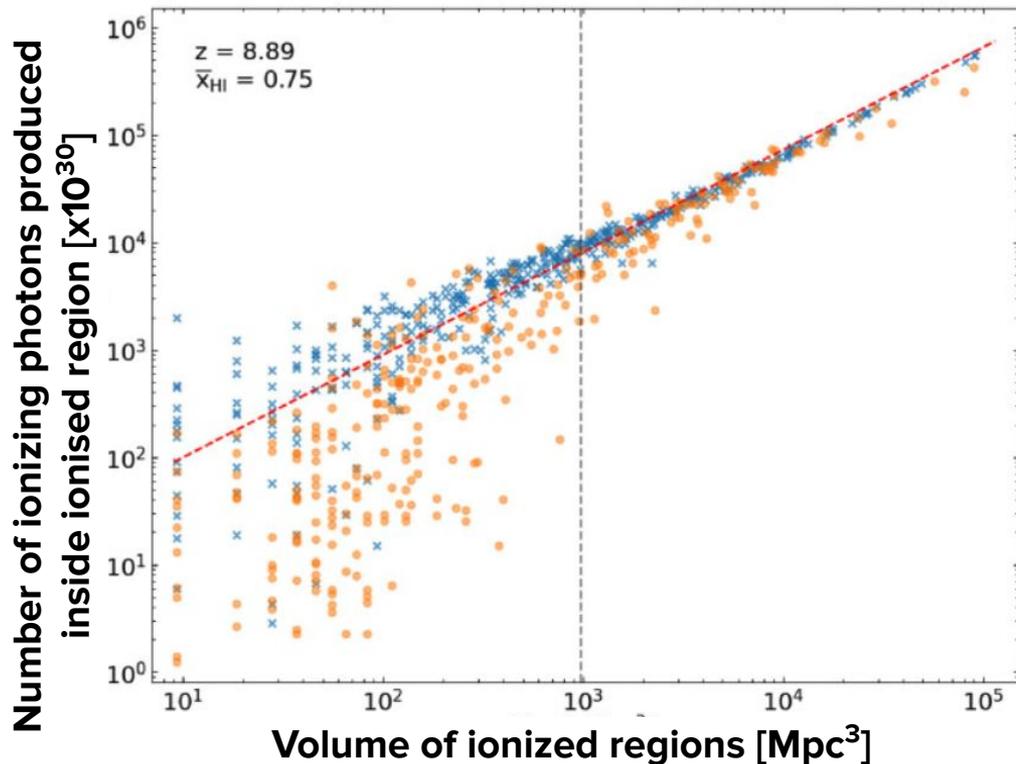


Estimate the mean neutral fraction



(Bianco, Giri+2021)

Estimate number of galaxies inside ionized regions

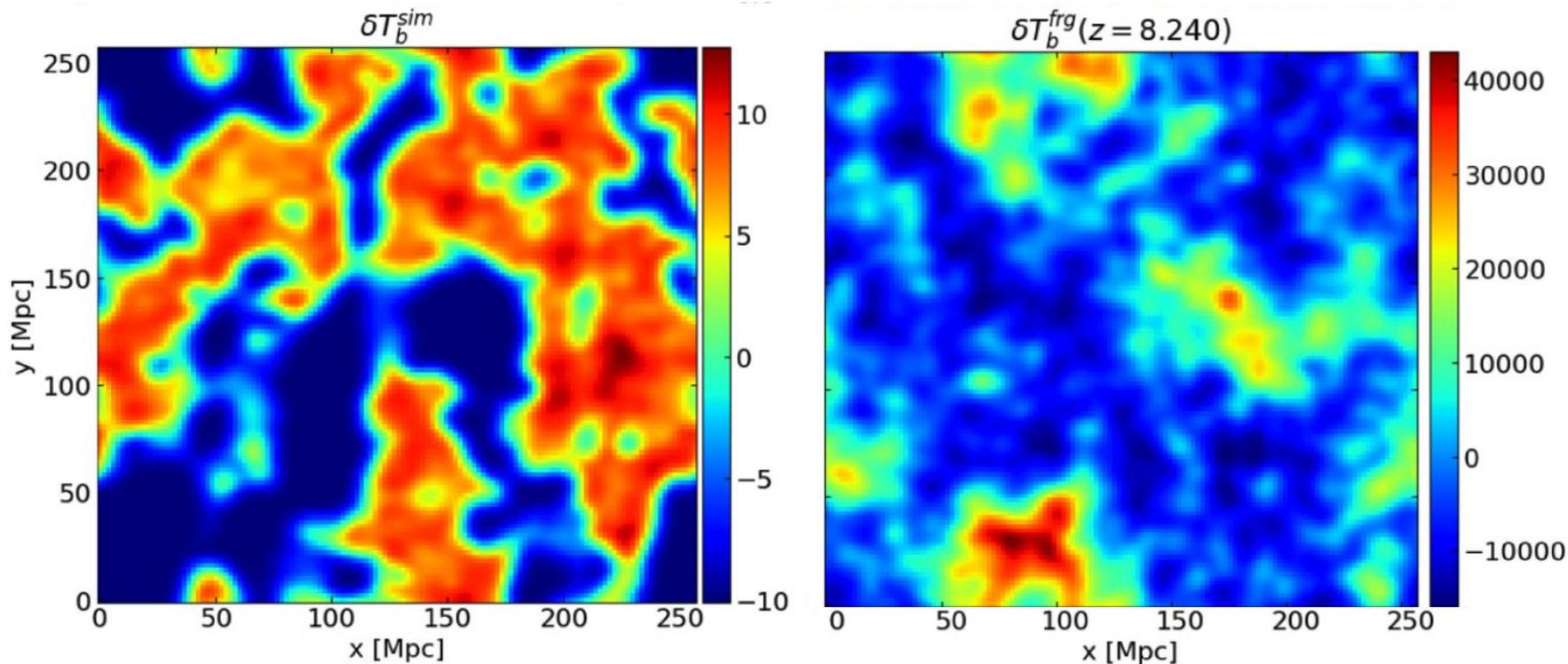


$$M_{\text{star,tot}} \approx 8 \times 10^{10} M_{\odot} \left(\frac{V_{\text{ion}}}{10^3 \text{cMpc}^3} \right) \left(\frac{\langle f_{\text{esc}} \rangle}{0.1} \right)^{-1}$$

(Zackrisson, ..., Giri+2020)

(Bianco, Giri+2024)

21-cm signal contaminated with Foreground



(Bianco, Giri+2021;2024)



Science Data Challenge 3

Foregrounds

[Scoring](#)[Teams](#)[Leaderboard](#)[Rules](#)[Data](#)[Test Data](#)

Our 'Foregrounds' challenge asks participants to remove obscuring sources of emission which prevent analysis of the underlying hydrogen-21cm signal from the Epoch of Reionisation. This foreground emission stems from both Galactic and extragalactic sources, both of which have previously observed, and unobserved components.

Given the lack of a model for the finer structure of Galactic emission at SKA-LOW frequencies, the removal of Galactic emission from the dataset represents a significant challenge. By similar reasoning, source confusion from previously unknown extragalactic sources, especially at the coarser resolution at metre-wavelengths, complicates the matter further.

From our synthetic [datasets](#), participants are asked to extract the cylindrically-averaged power spectrum of the EoR signal, clean from foregrounds contamination.

To assess resulting submissions, our [scoring](#) ('figure-of-merit') algorithms will take resulting cylindrical power spectra, and return a score. Ancillary analytical data products can be assessed, however, the cylindrical power spectra results will be the only ones that affect the 'leaderboard'. Given that this challenge runs over 6 months, we expect this to be an iterative process and encourage competing teams to make full use of the [computing facilities](#) assigned to them.



Science Data Challenge 3

Foregrounds

Astrophysicist

Michele Bianco,^{1,2,3}★ Sambit. K. Giri,^{4,5} Rohit Sharma,⁶ Tianyue Chen,¹ Shreyam Parth Krishna,¹

Chris Finlay⁷, Viraj Nistane⁷, Philipp Denzel,⁶ Massimo De Sanctis,⁸ Hatem Ghorbel⁸

¹ Laboratoire d'Astrophysique, Ecole Polytechnique Federale de Lausanne (EPFL), Observatoire de Sauverny, Versoix

² Astronomy Centre, Department of Physics & Astronomy, Pevensey III Building, University of Sussex, Falmer, Brighton

³ Institute for Particle Physics and Astrophysics, ETH Zurich, Wolfgang-Pauli-Str 27, 8093 Zurich, Switzerland

⁴ Nordita, KTH Royal Institute of Technology and Stockholm University, Hannes Alfvéns väg 12, SE-106 91 Stockholm, Sweden

⁵ Institute for Computational Science, University of Zurich, Winterthurerstrasse 190, 8057 Zurich, Switzerland

⁶ Fachhochschule Nordwestschweiz, Bahnhofstrasse 6, 5210 Windisch, Switzerland

⁷ Département de Physique Théorique and Center for Astroparticle Physics, Université de Geneve, 24 quai Ernest Ansermet, 1211 Geneve 4, Switzerland

⁸ Data Analytics Group, Haute Ecole Arc Ingénierie HES-SO, University of Applied Sciences Western Switzerland

From our synthetic [subsets](#), participants are asked to extract the cylindrically averaged power spectrum of the EoR signal, clean from foregrounds contamination.

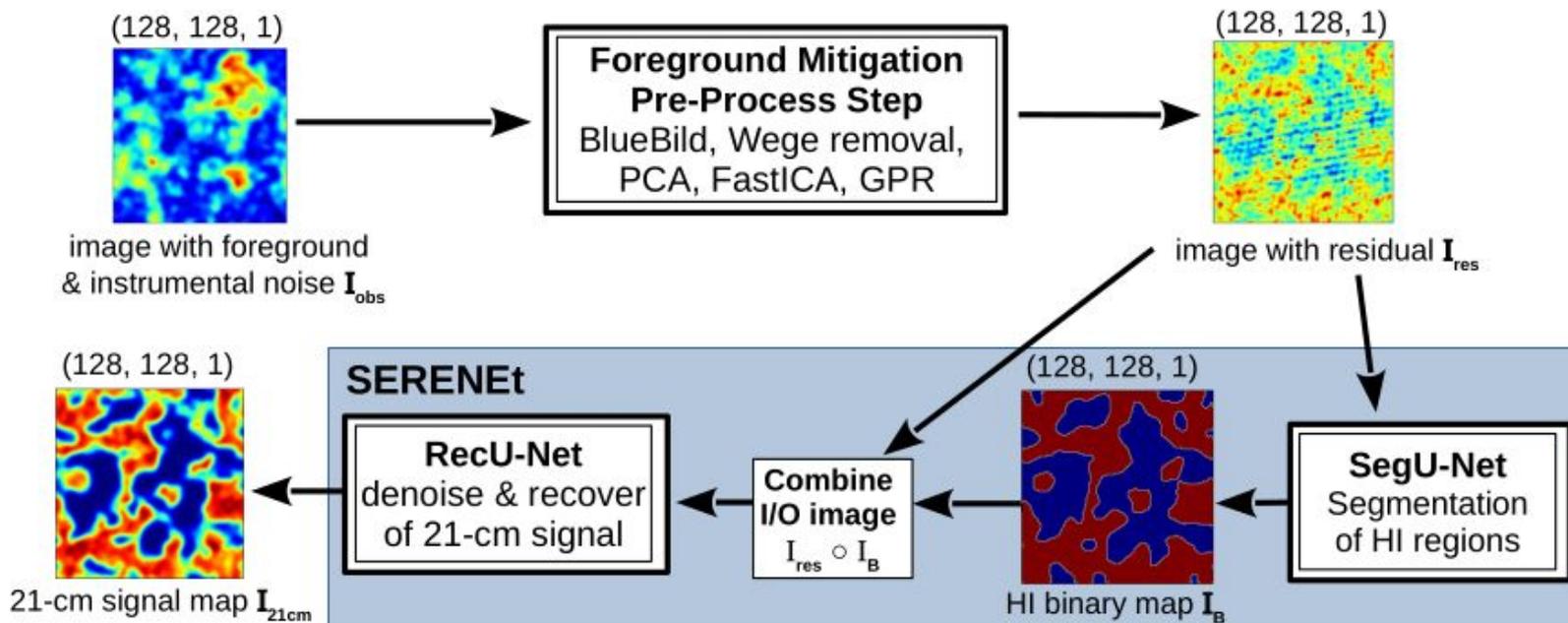
To assess resulting submissions, our [scoring](#) ('figure-of-merit') algorithms will take resulting cylindrical power spectra, and return a score. Ancillary analytical data products can be assessed, however, the cylindrical power spectra results will be the only ones that affect the 'leaderboard'. Given that this challenge runs over 6 months, we expect this to be an iterative process and encourage competing teams to make full use of the [computing facilities](#) assigned to them.

Data

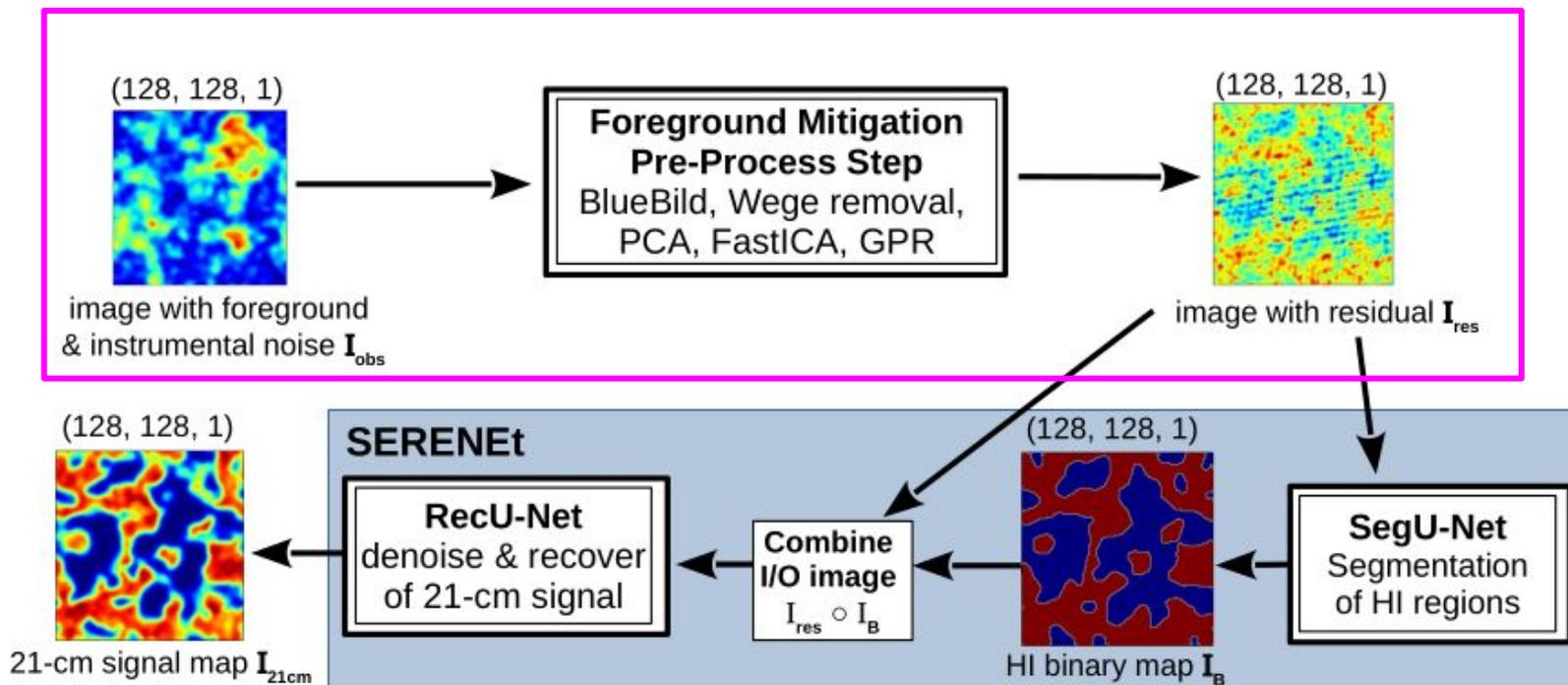
Test Data

Data Scientists

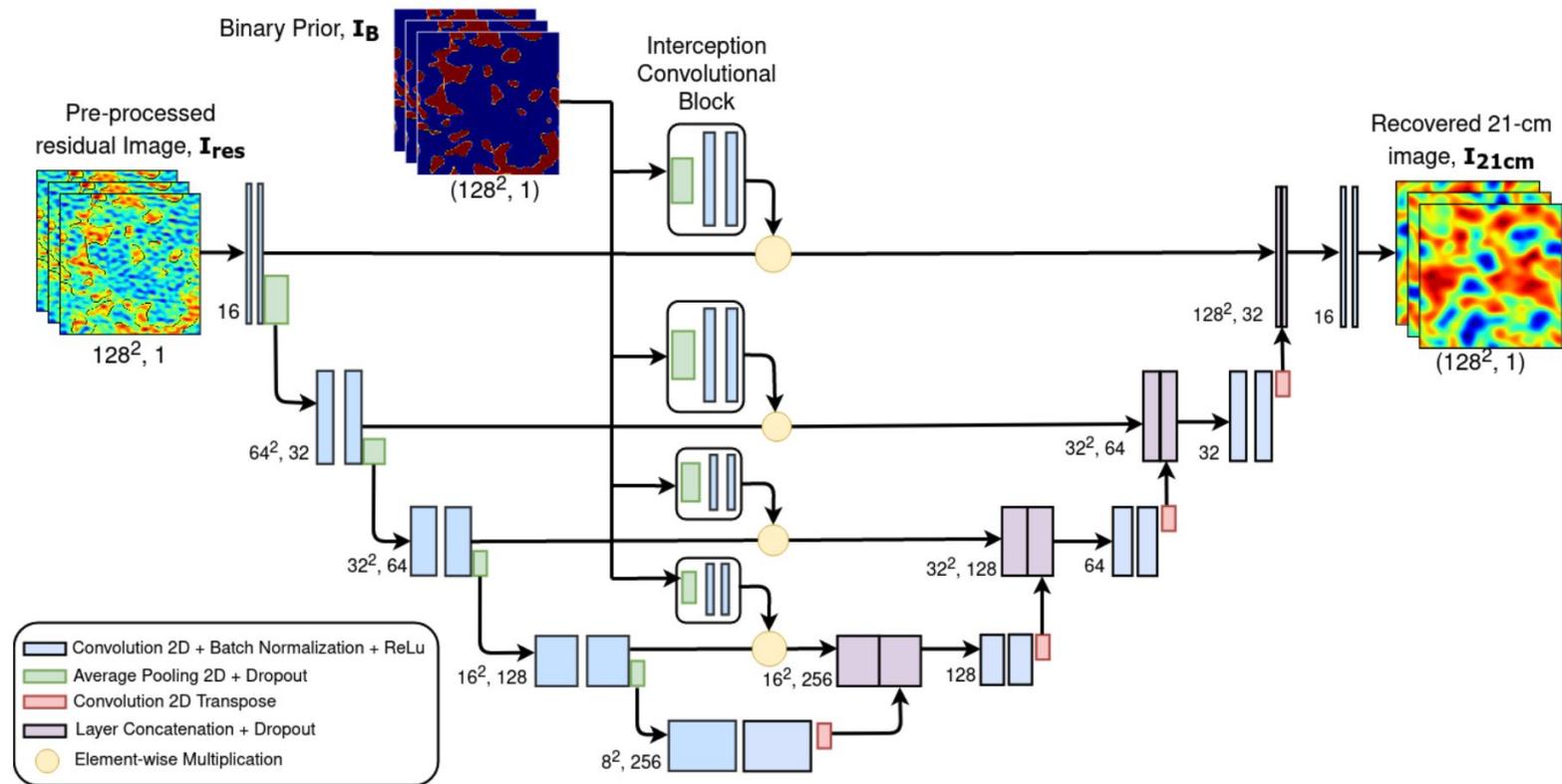
SERENet foreground mitigation pipeline



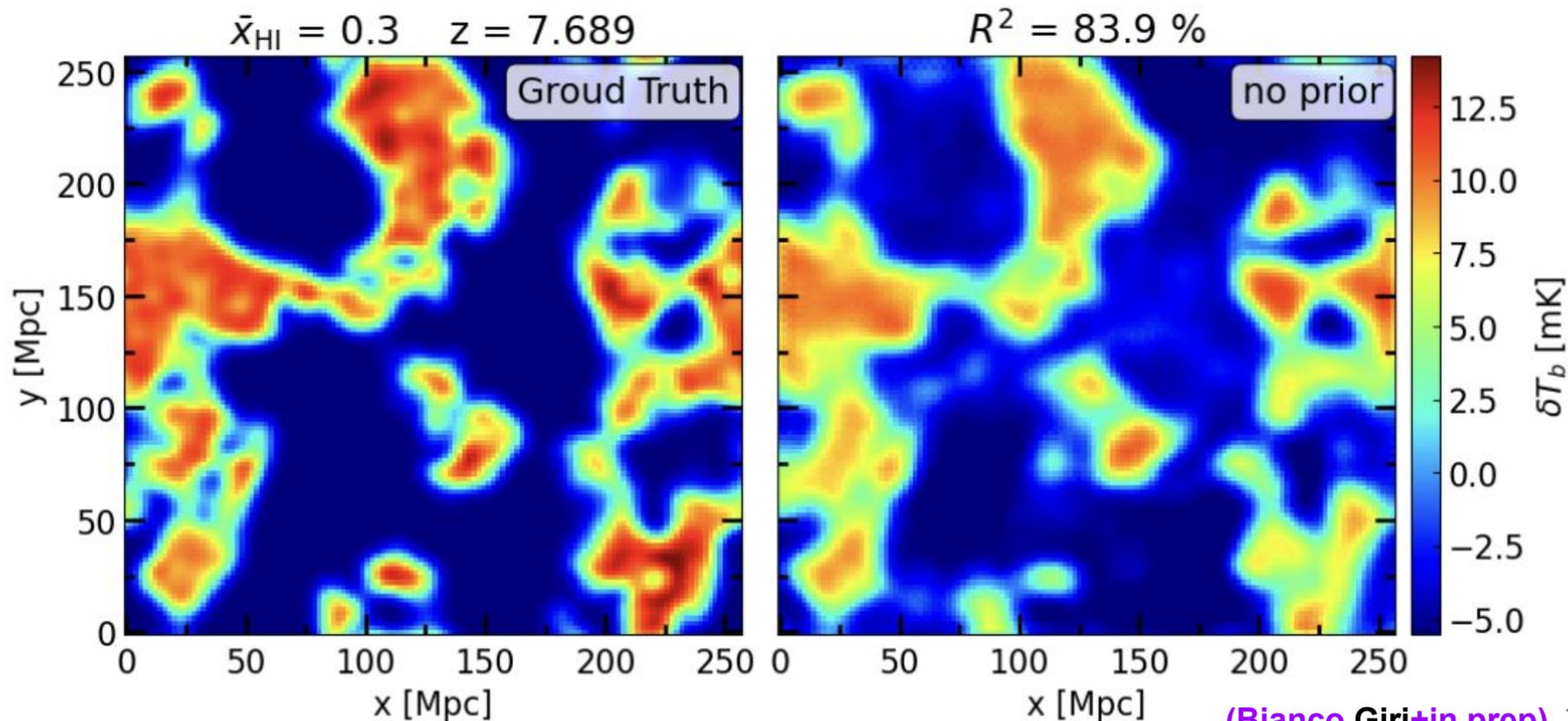
SERENet: Pre-processing with PCA



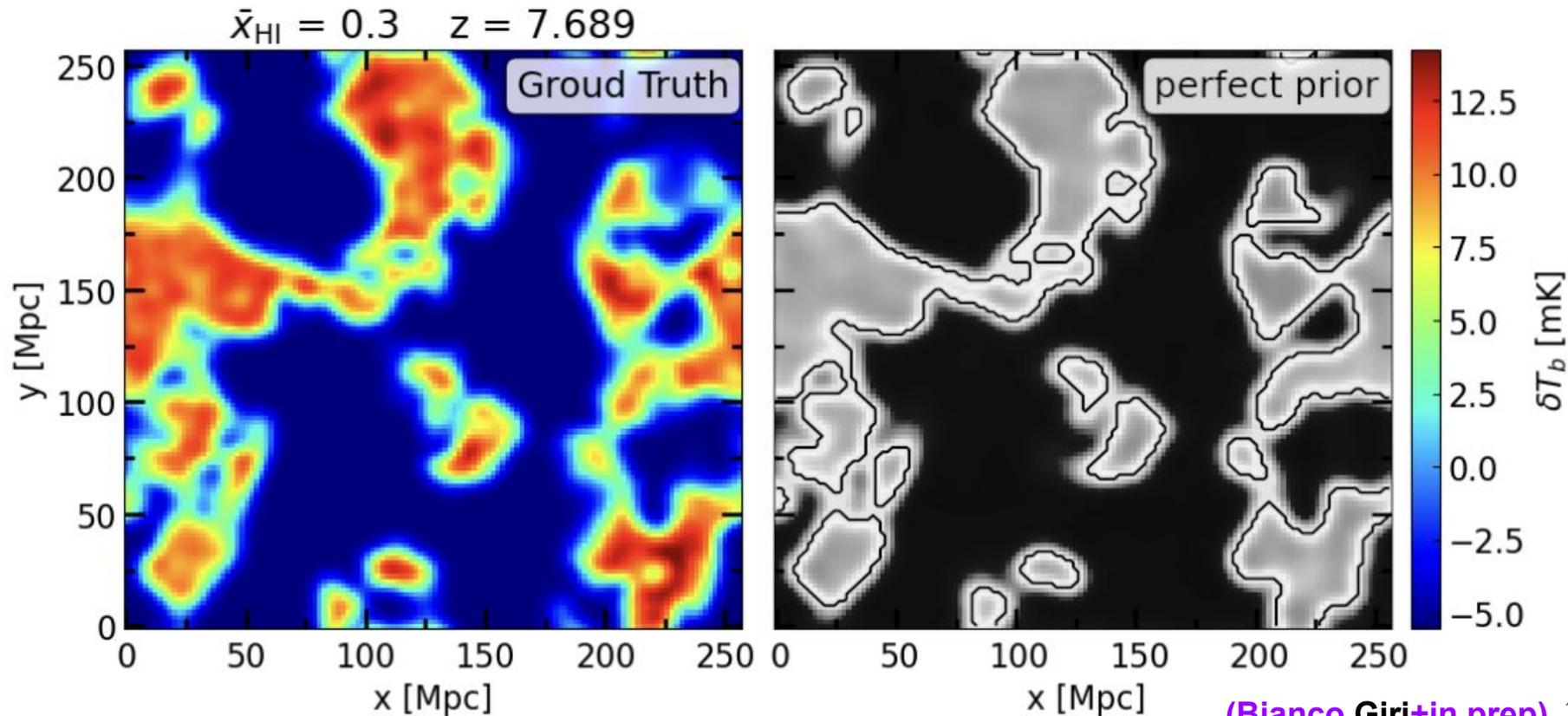
SERENet framework



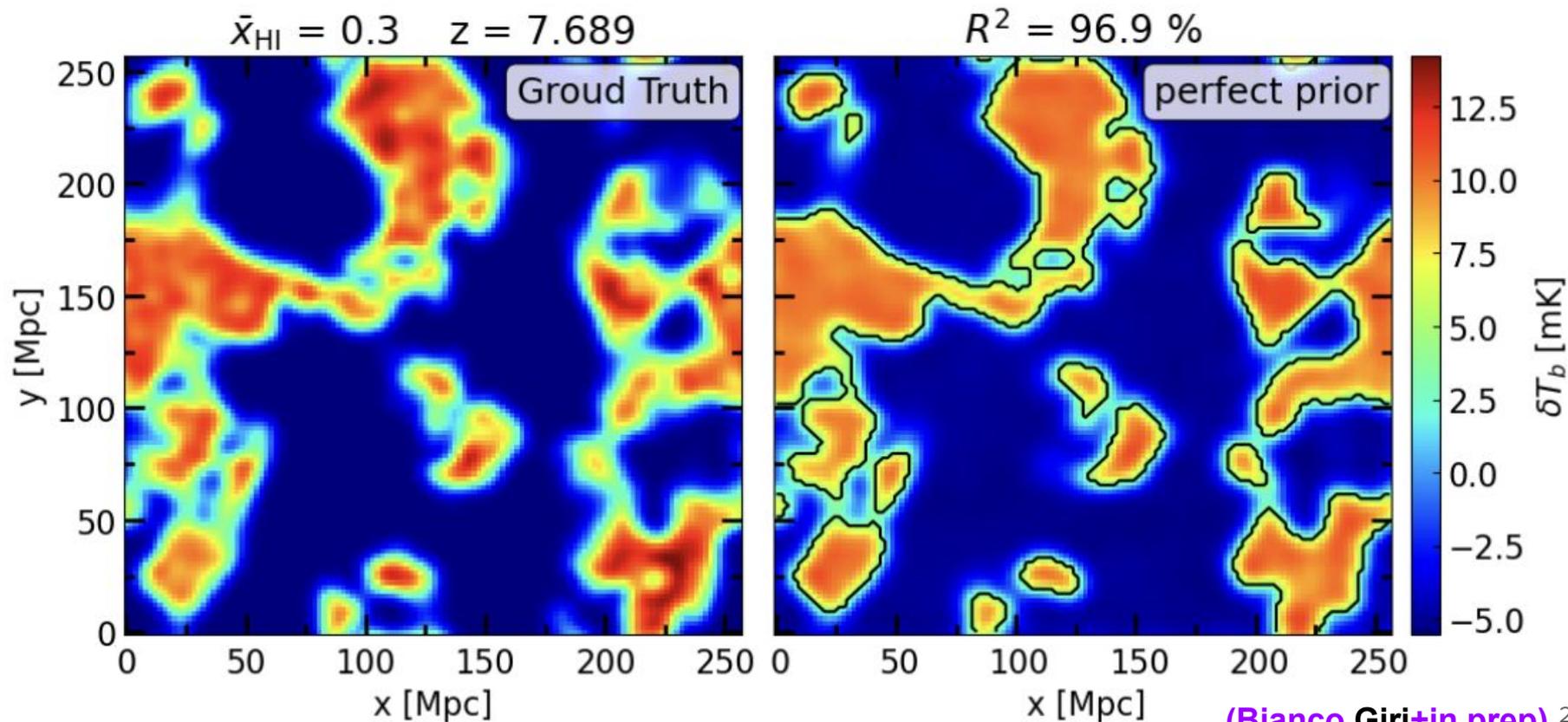
SERENet recovering the 21-cm image



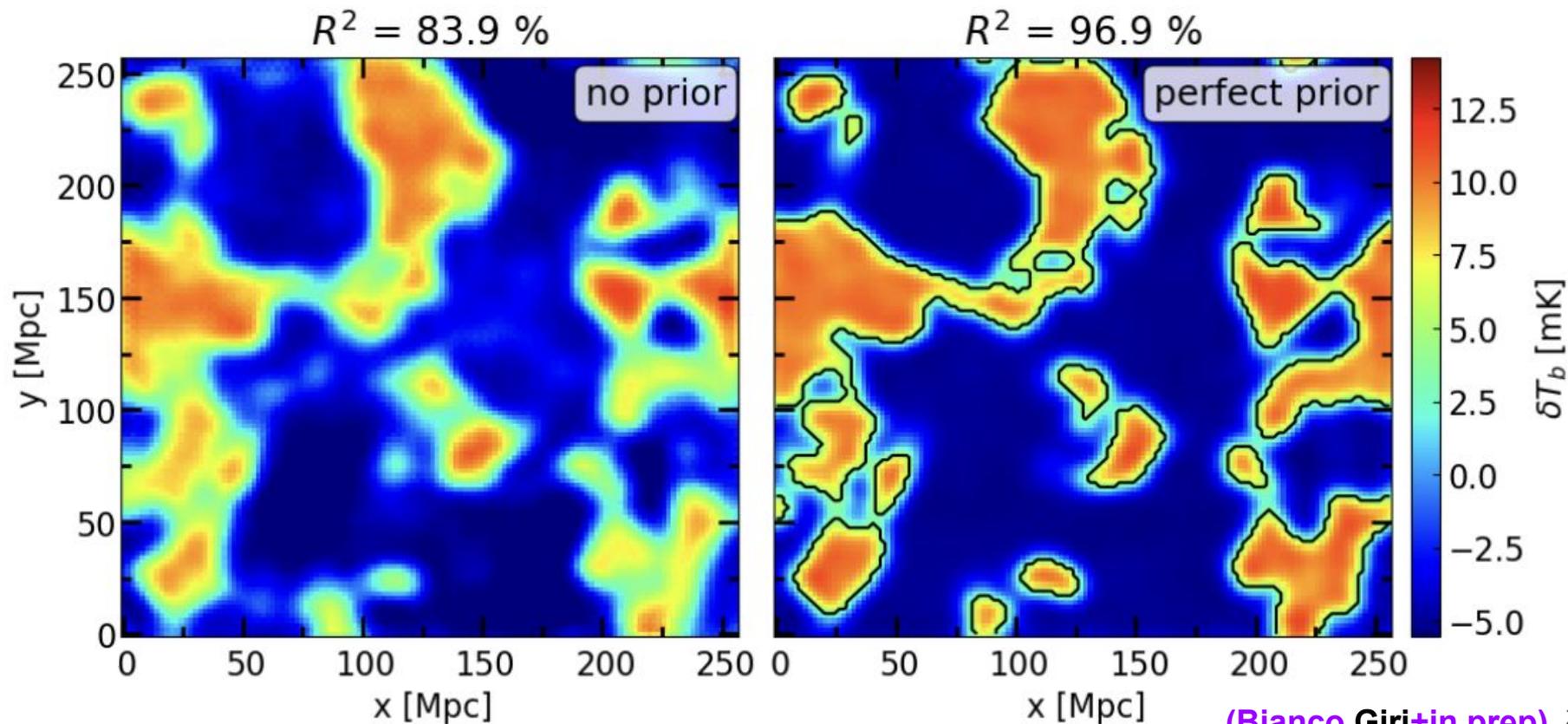
SERENet: binary maps of ionized regions as prior



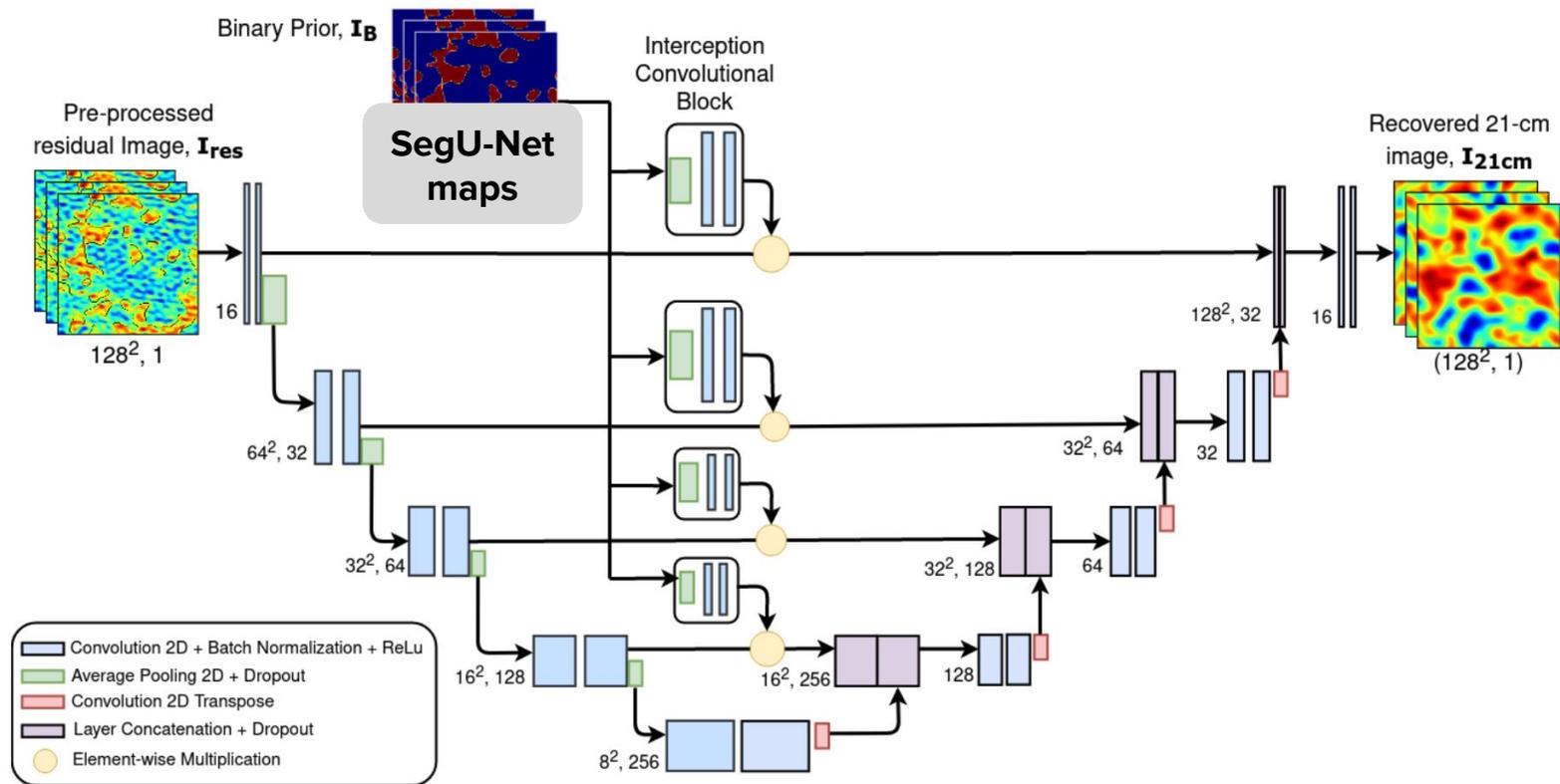
SERENet recovering the 21-cm image with prior



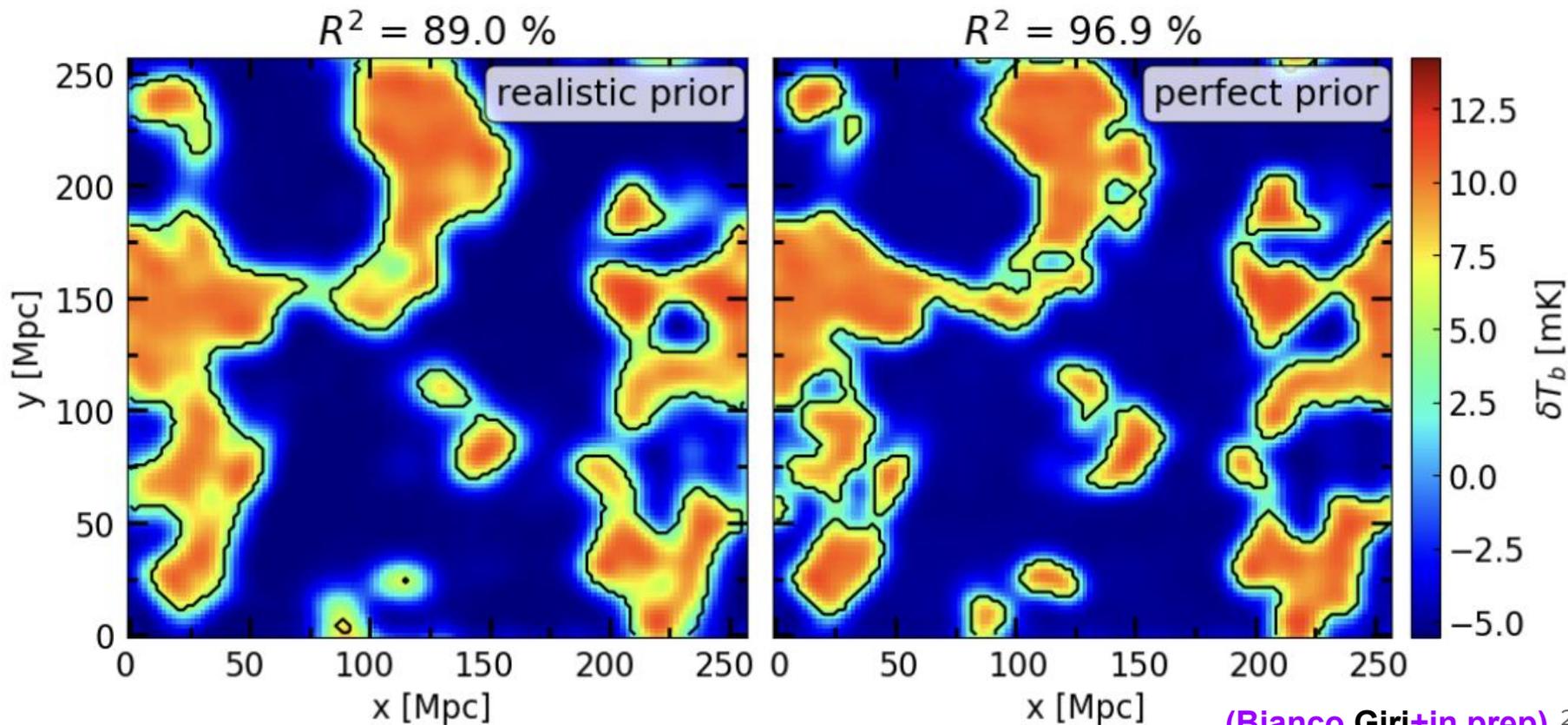
No prior vs Perfect prior



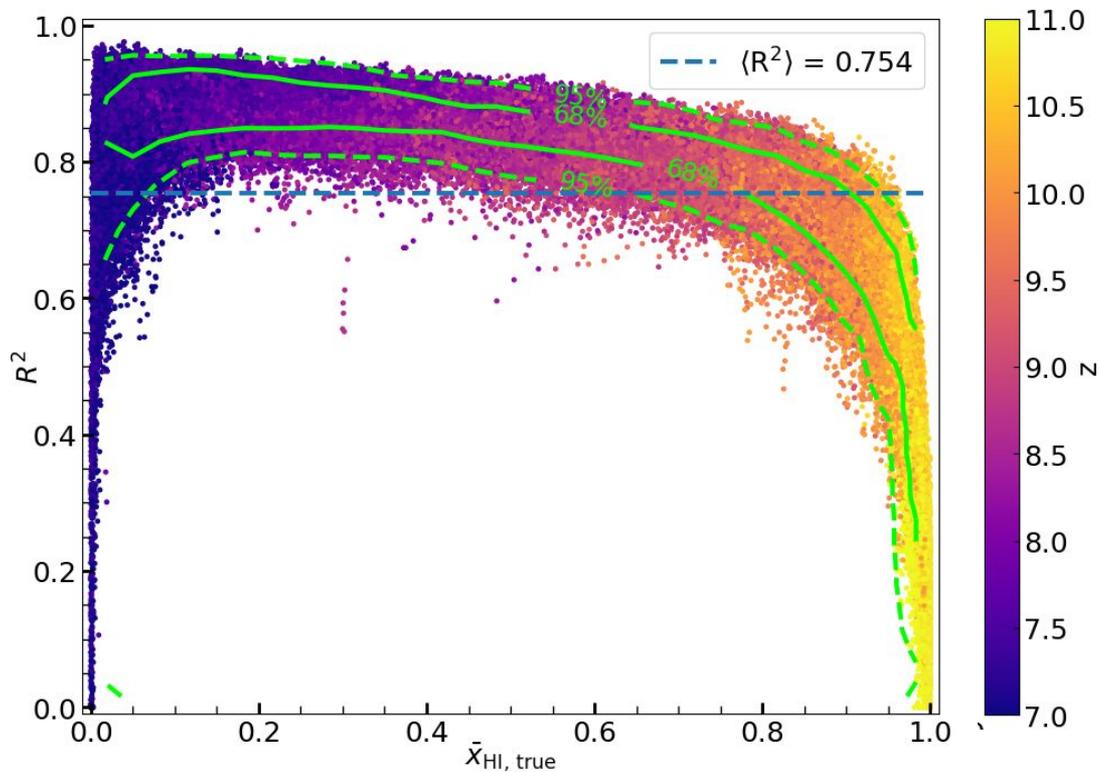
SERENet: realistic prior with SegU-Net



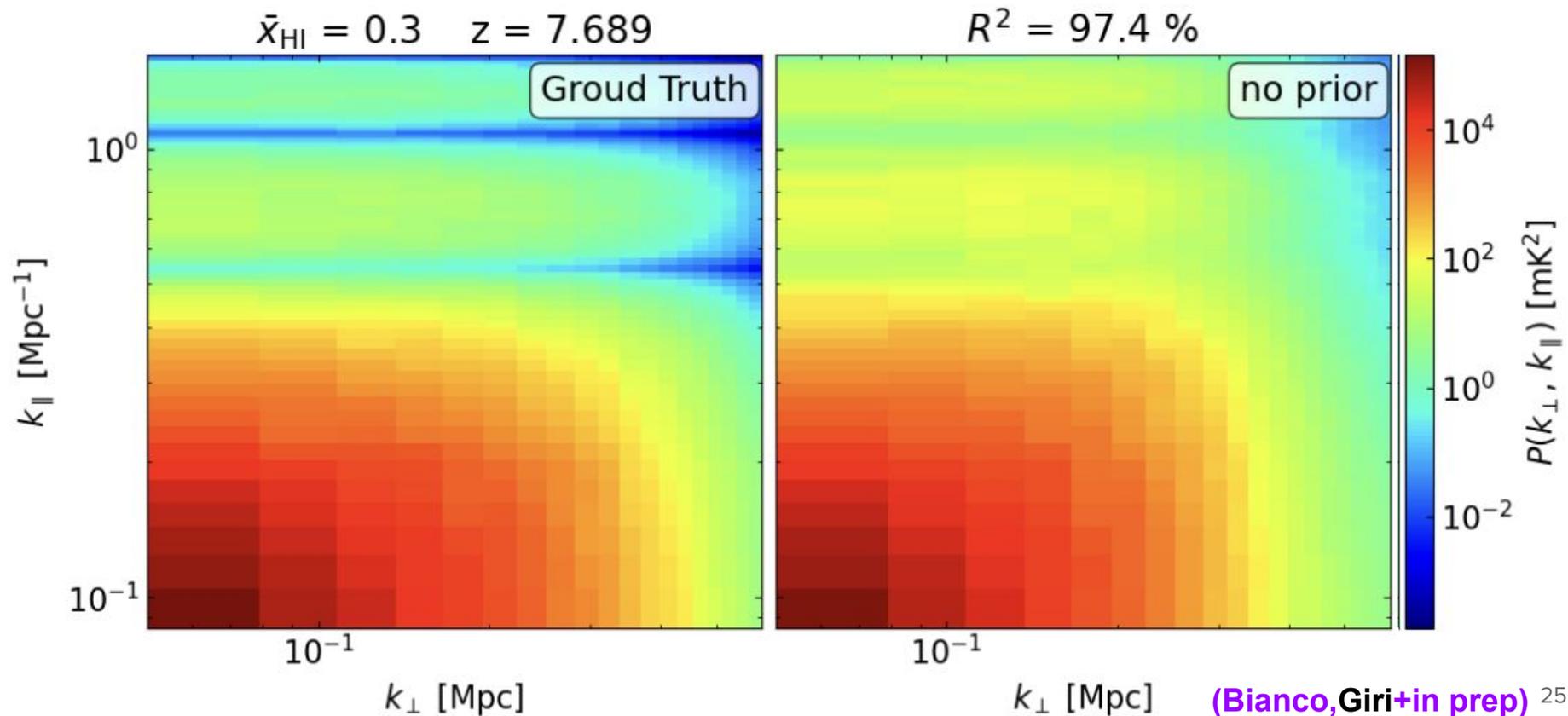
SERENet recovering the 21-cm image with prior



R^2 -score of the entire testing set



SERENet recovering the 21-cm power spectrum



Summary

- Data science methods will maximise the scientific outcomes of the huge amount of data that SKA will produce
- **SegU-Net**: deep learning-based framework to probe ionization state of the intergalactic medium
- **SERENet**: deep learning-based to mitigate foreground contamination in the 21-cm signal during the epoch of reionization