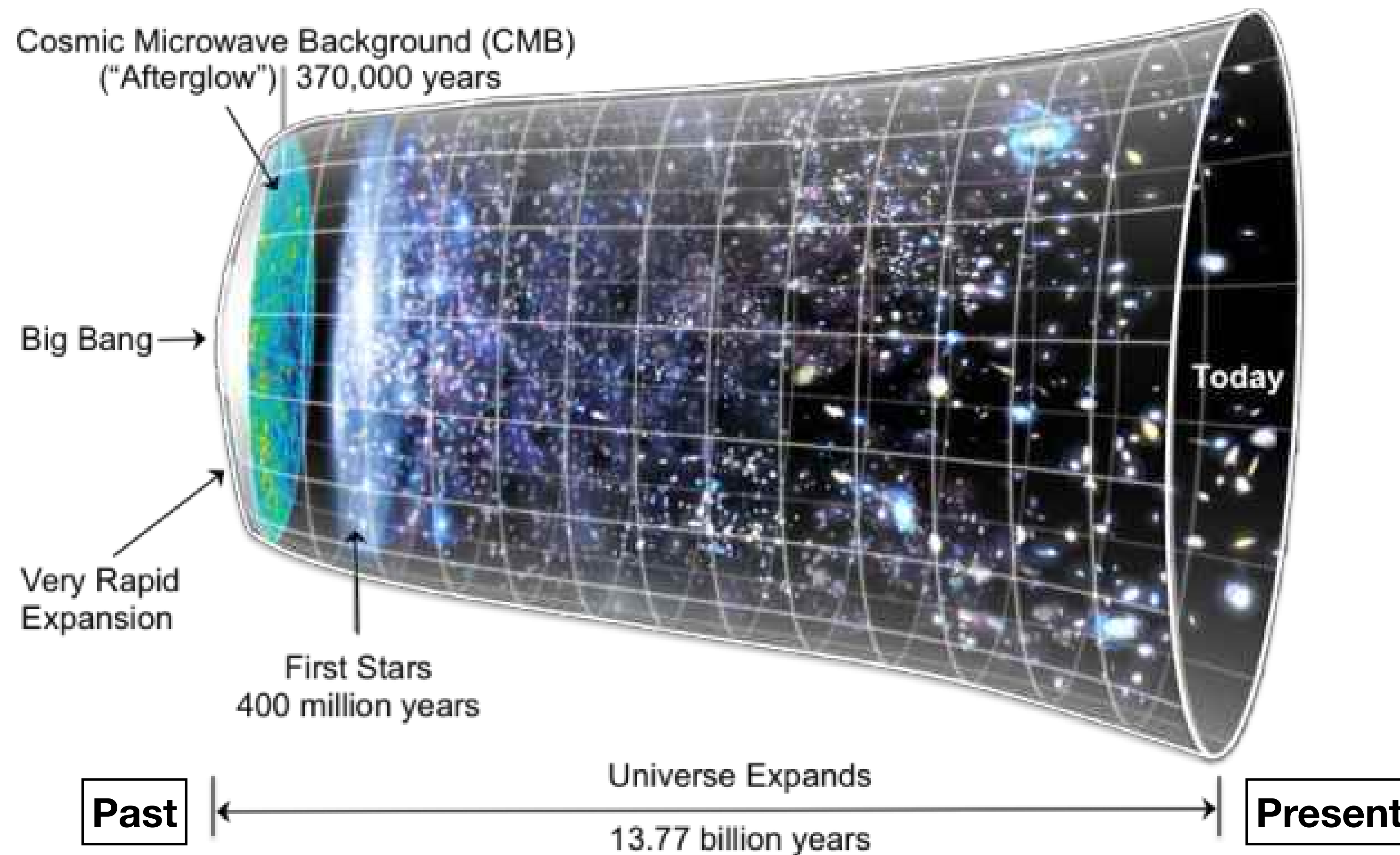


Deep Learning Application for Reconstruction of the Large-Scale Structure of the Universe

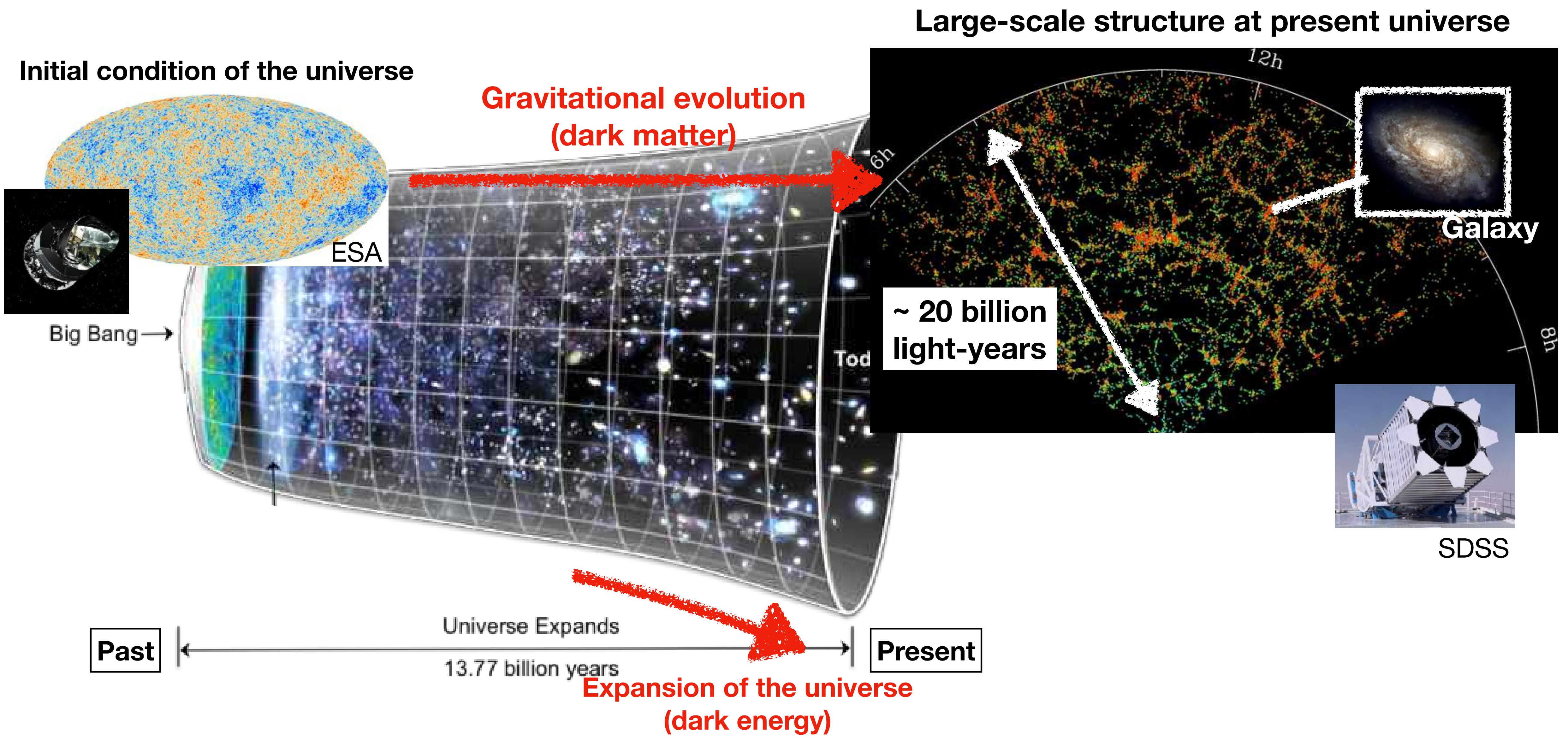
Kana Moriwaki (UTokyo UTAP/RESCEU)

ipi seminar
2022/10/19

Evolution of the Universe and the Large-Scale Structure of the Universe



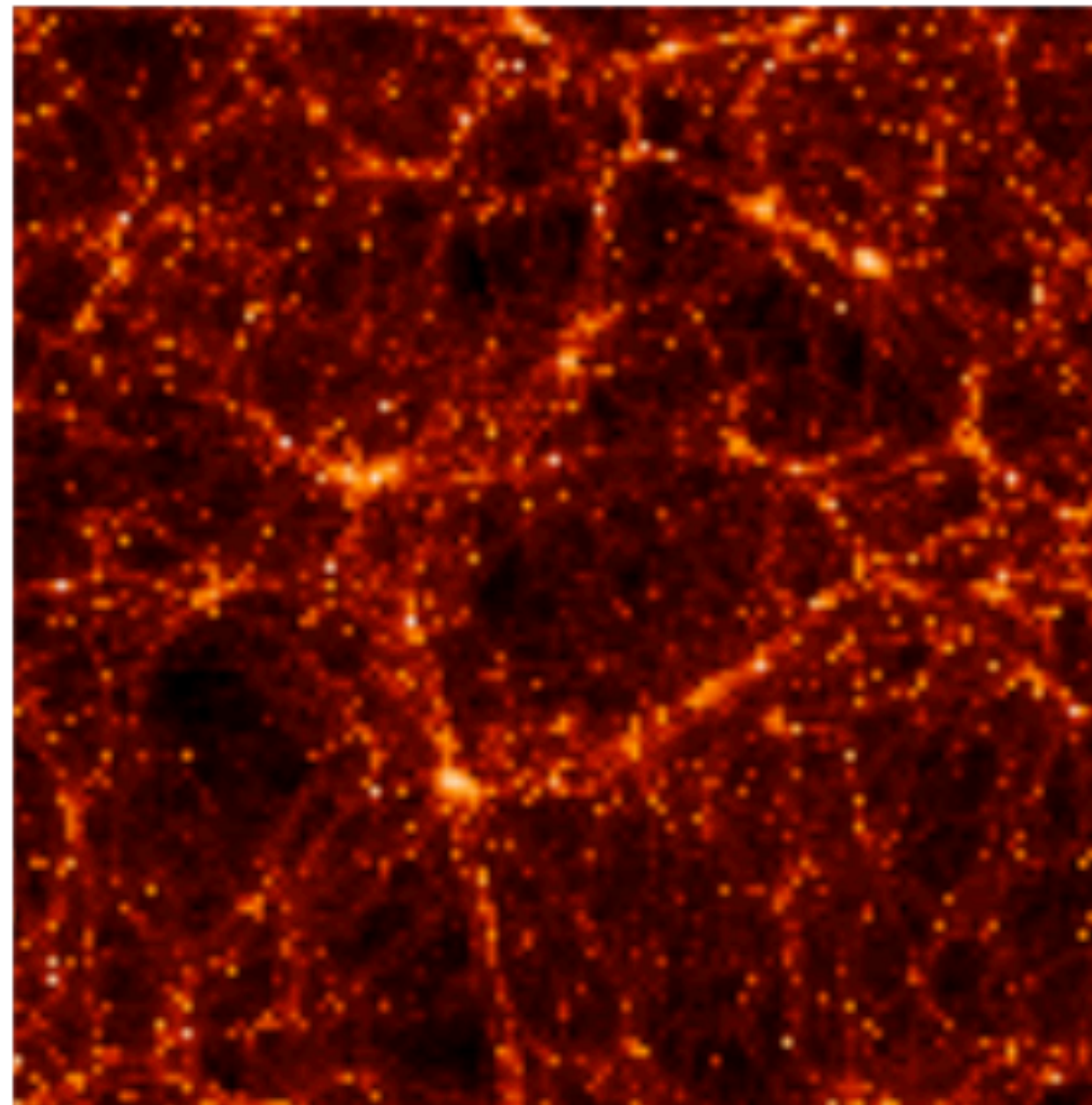
Evolution of the Universe and the Large-Scale Structure of the Universe



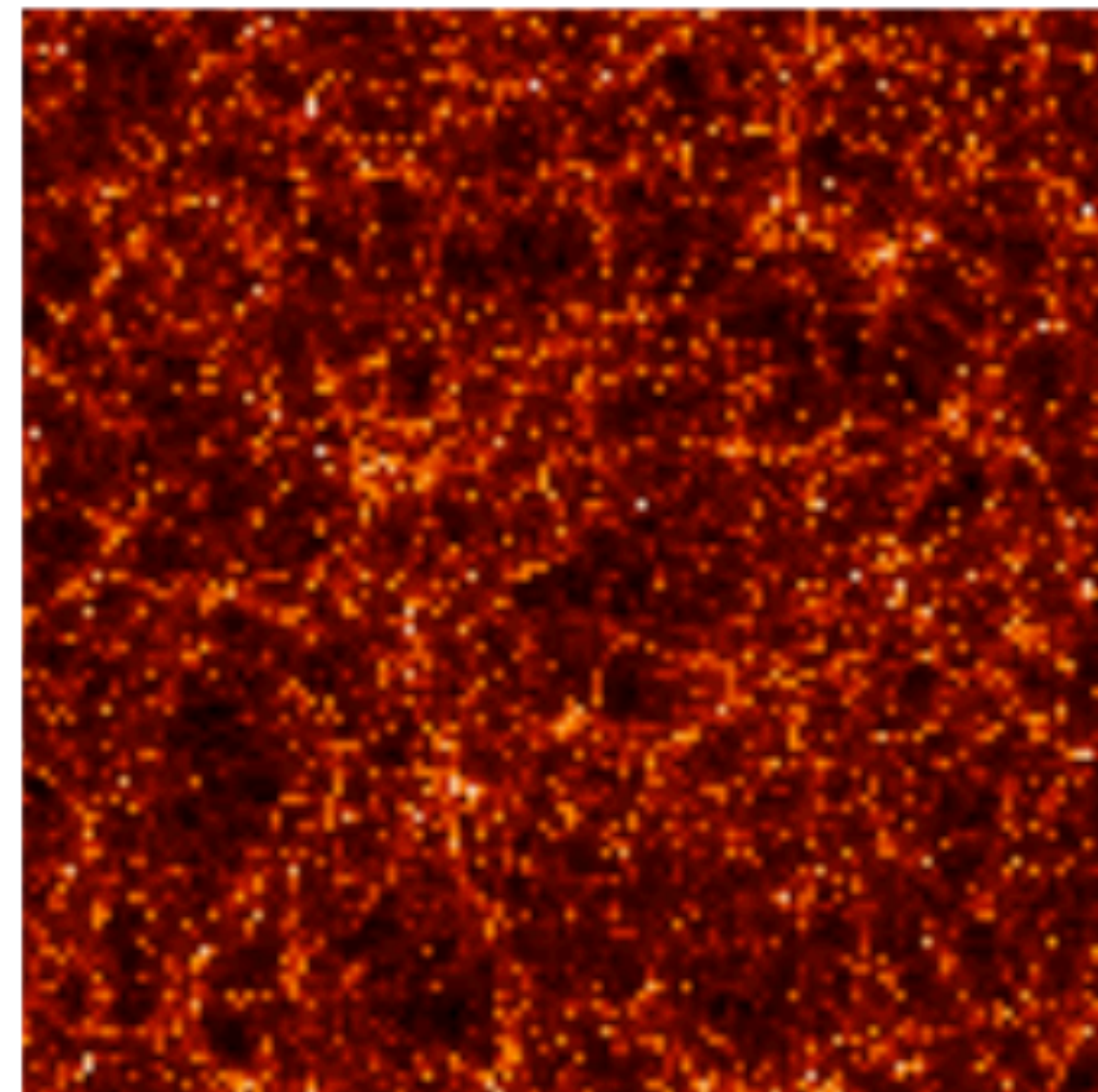
Evolution of the Universe and the Large-Scale Structure of the Universe

- **The large-scale structure tells us about:**
 - **Contents of the universe including dark matter and dark energy**
 - **Initial condition of the universe**
 - **etc.**

**30% matter
70% dark energy**

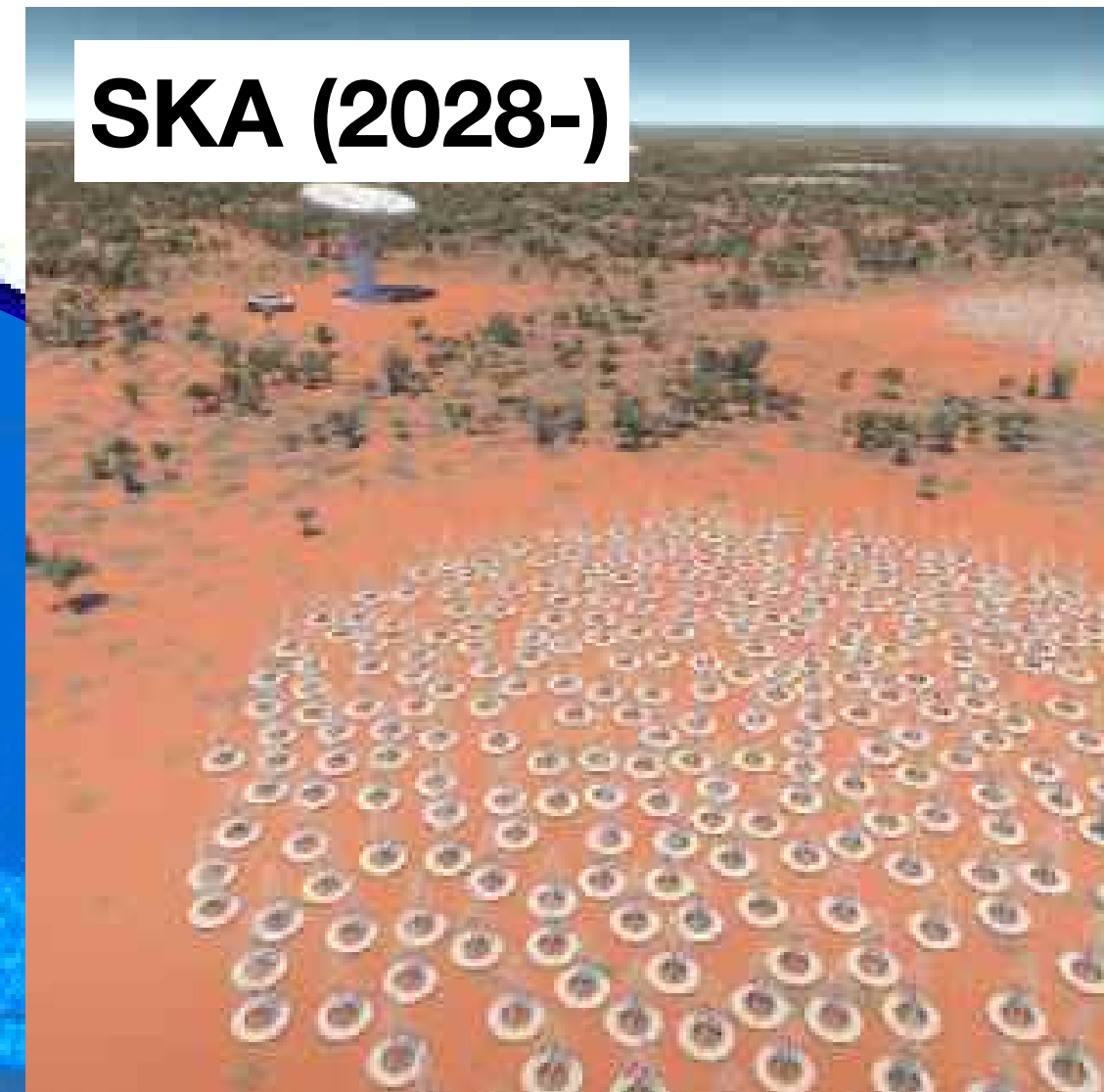
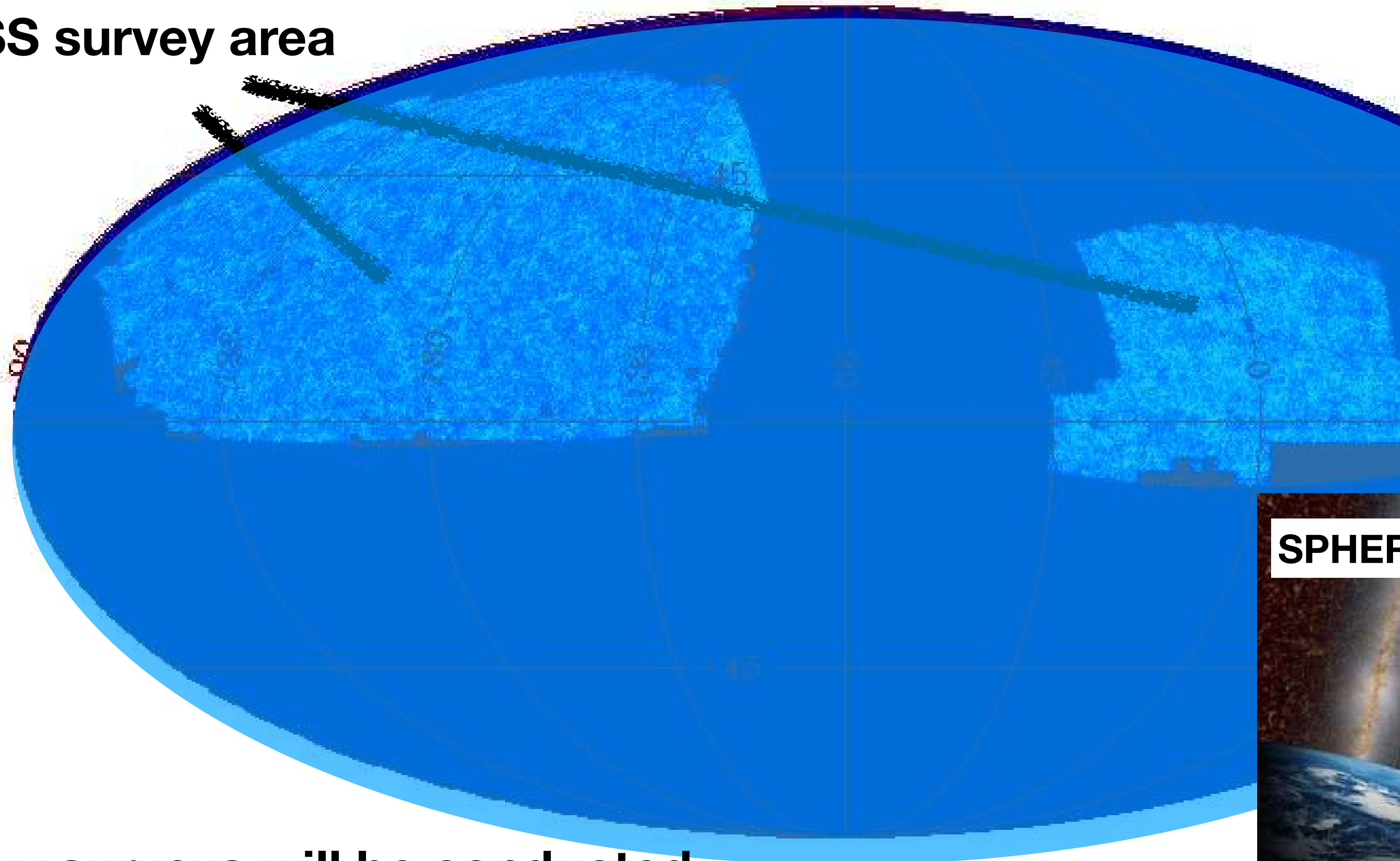


100% matter



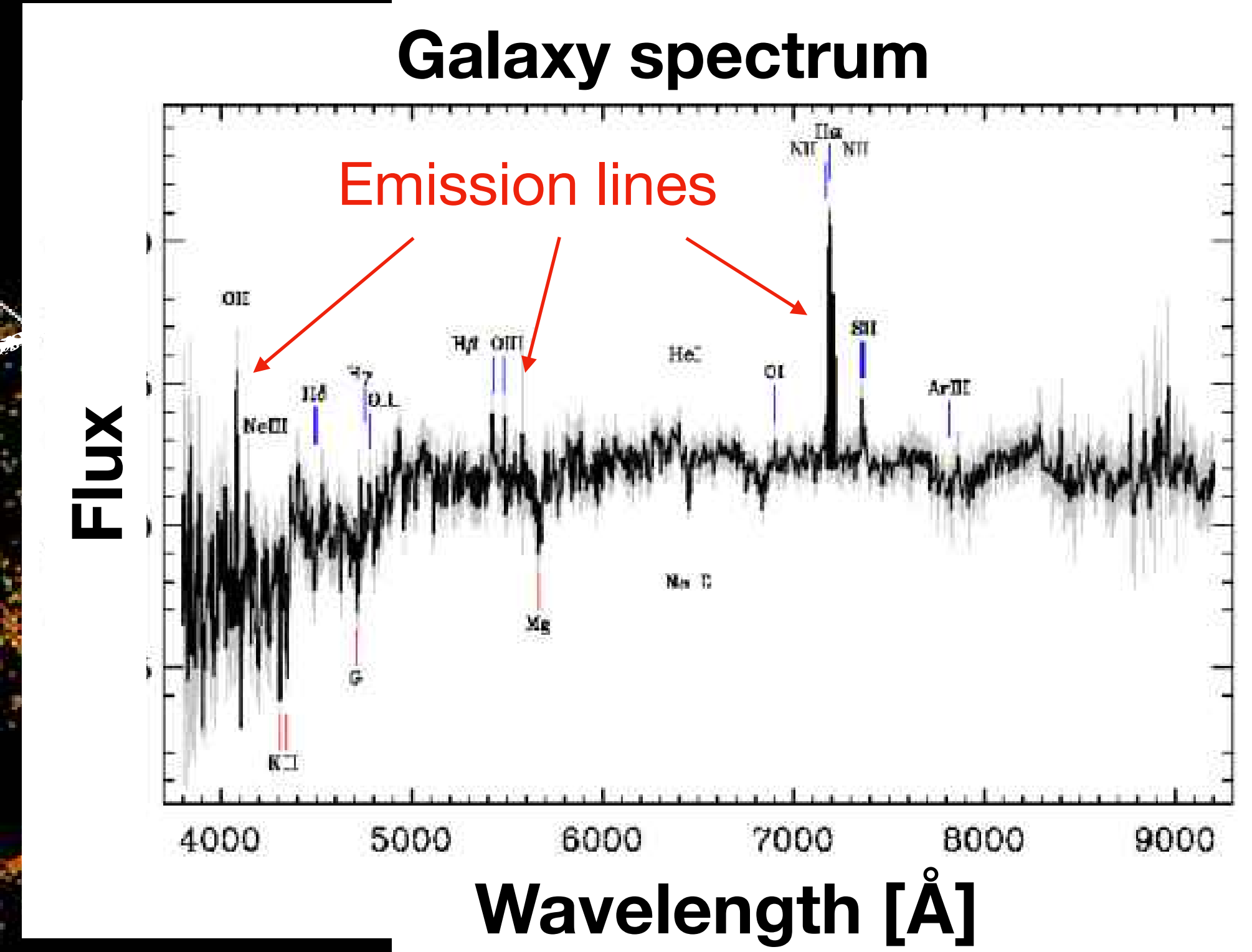
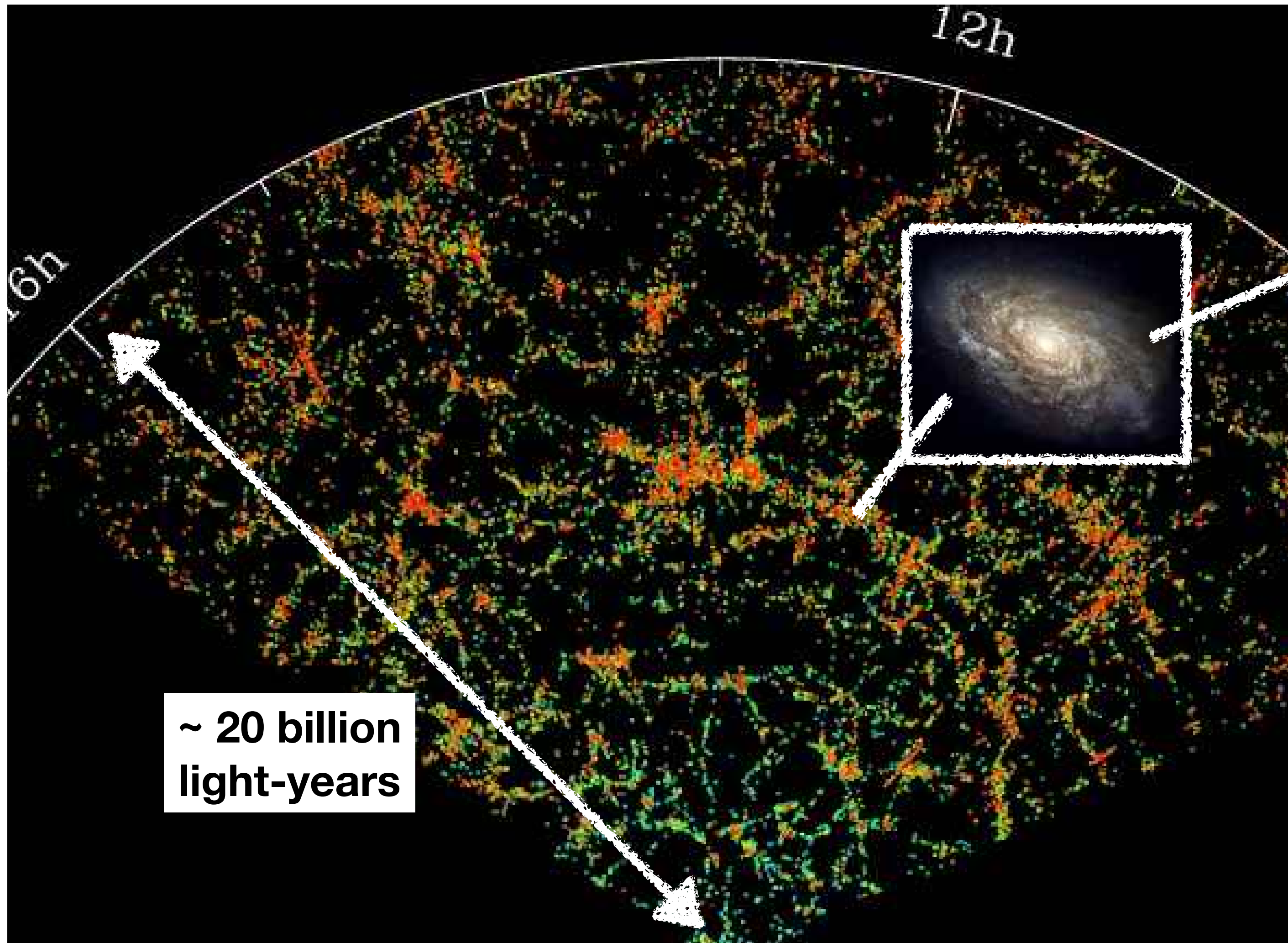
Largest-volume data will be available soon!

SDSS survey area



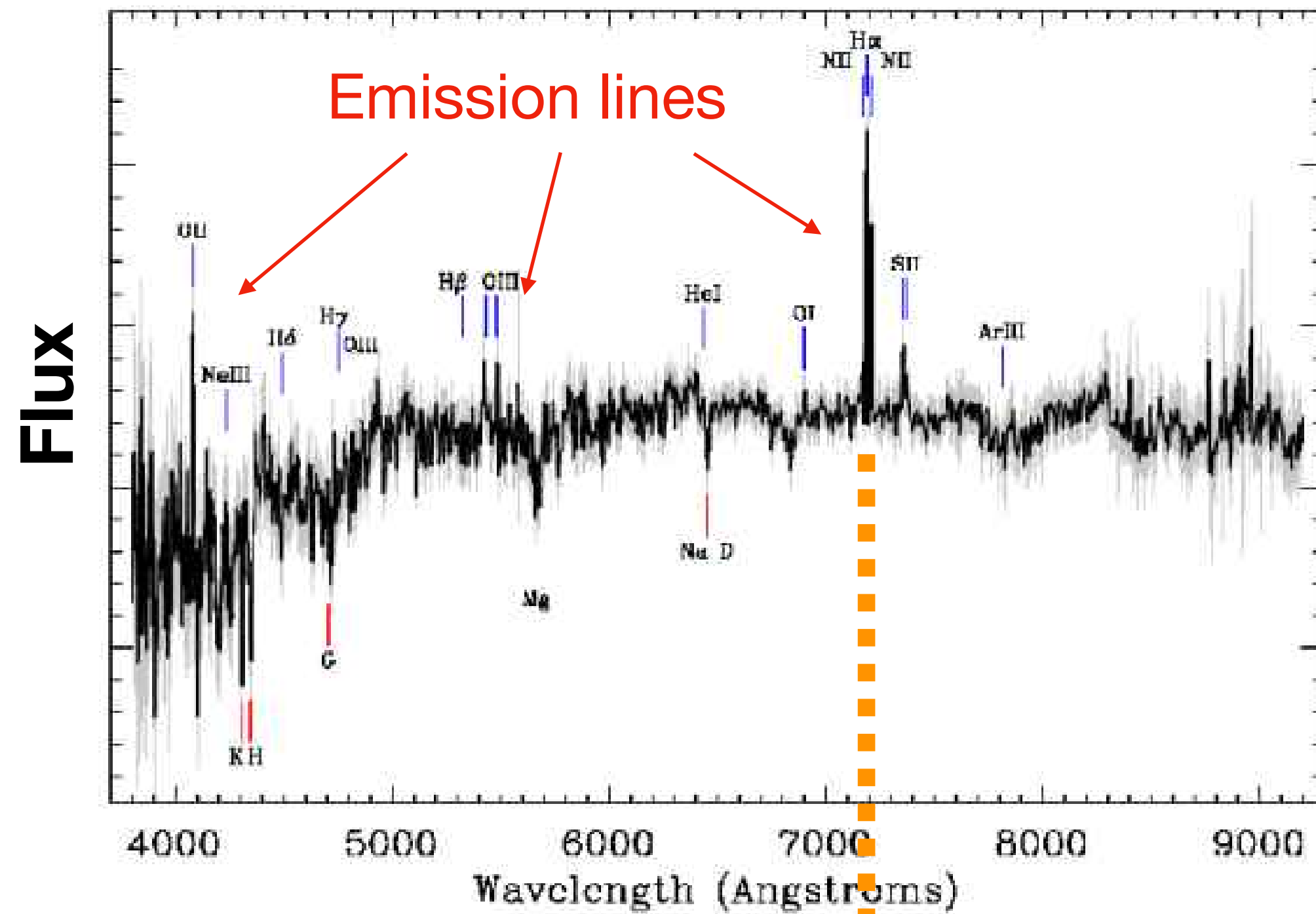
All-sky surveys will be conducted.

Spectroscopic Observations to Measure the Large-scale Distributions



Emission Line: a Key to Measure 3D Distributions

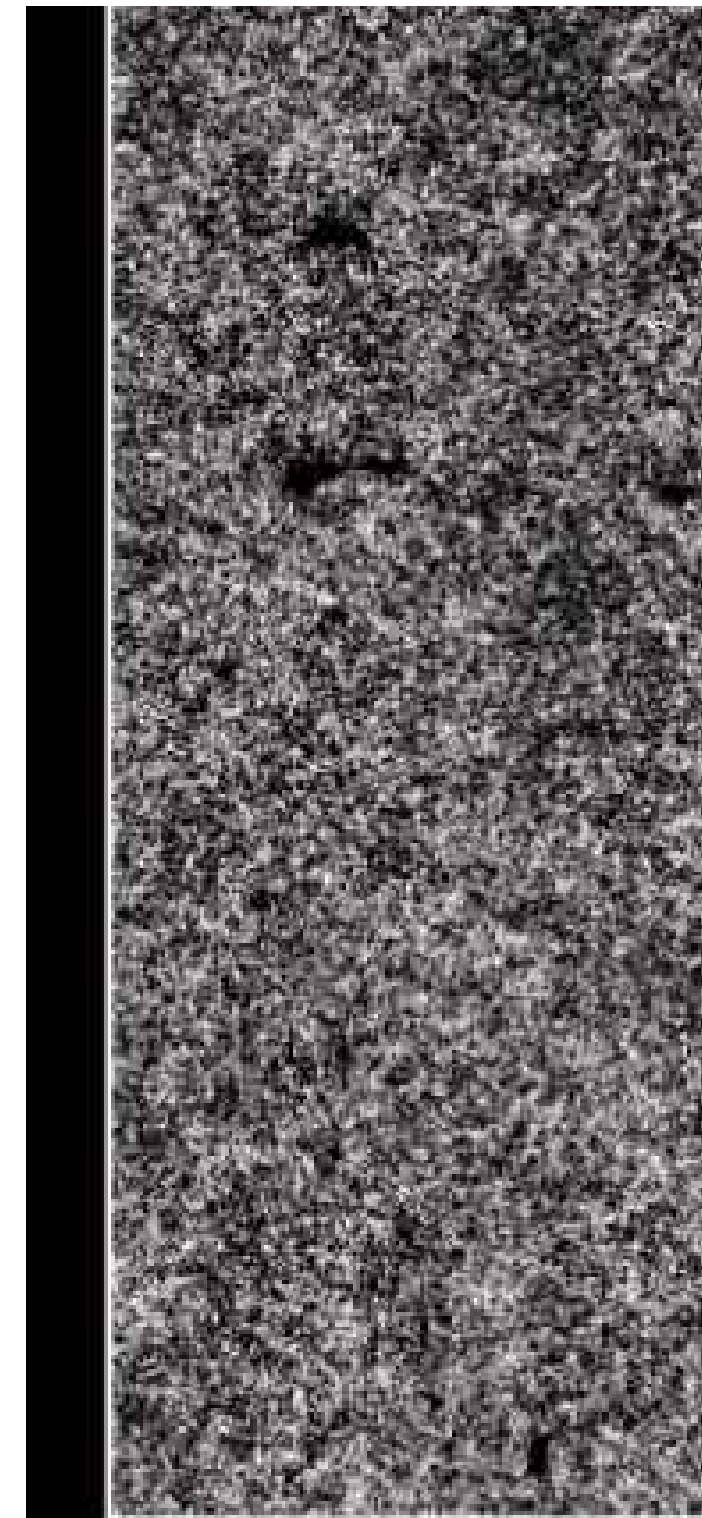
The observed wavelengths of emission lines are a measure of the distance.



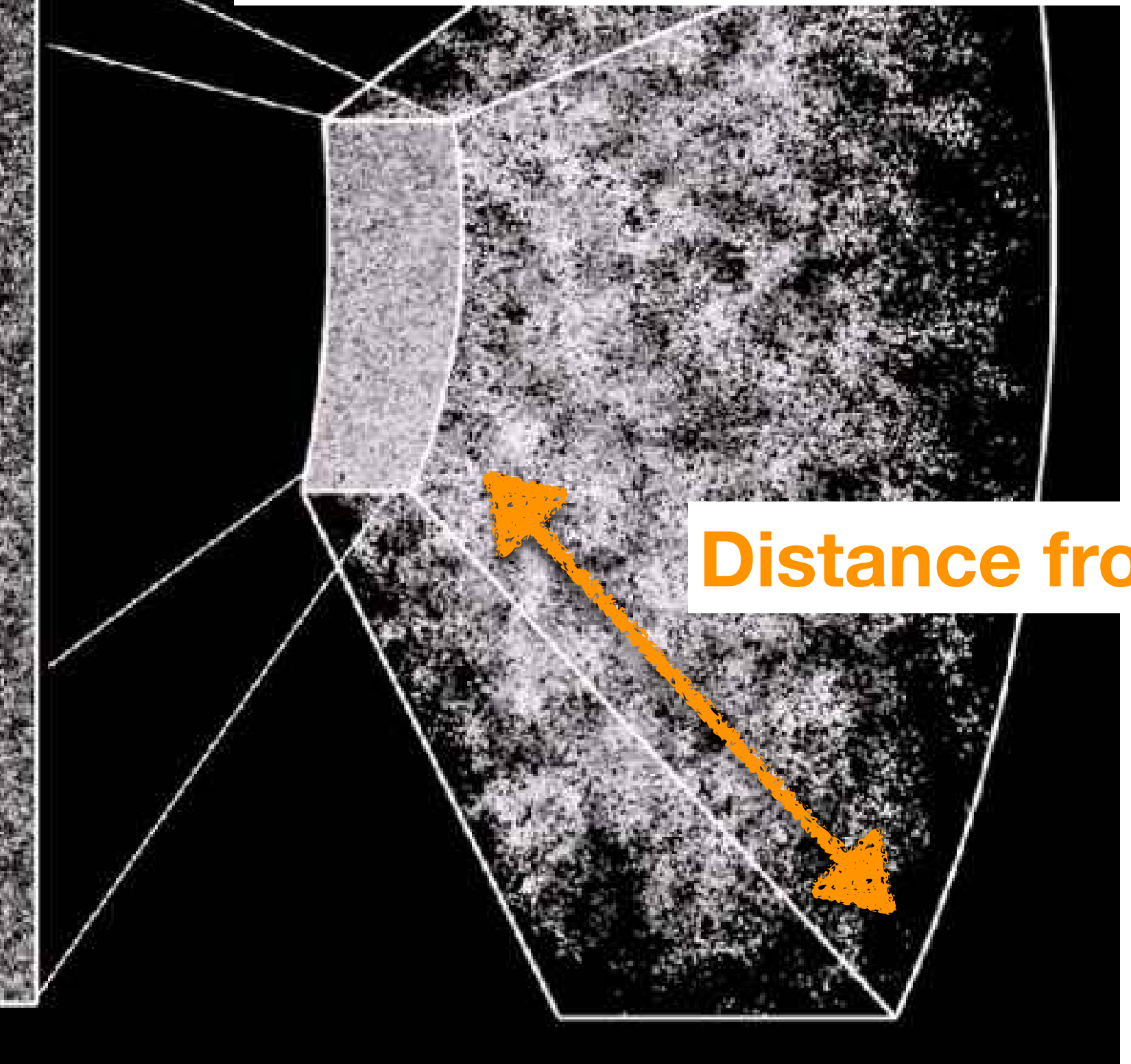
Observed wavelength

$$\lambda_{\text{obs}} = \lambda_{\text{rest}} (1 + z)$$

2D distribution
on the sky



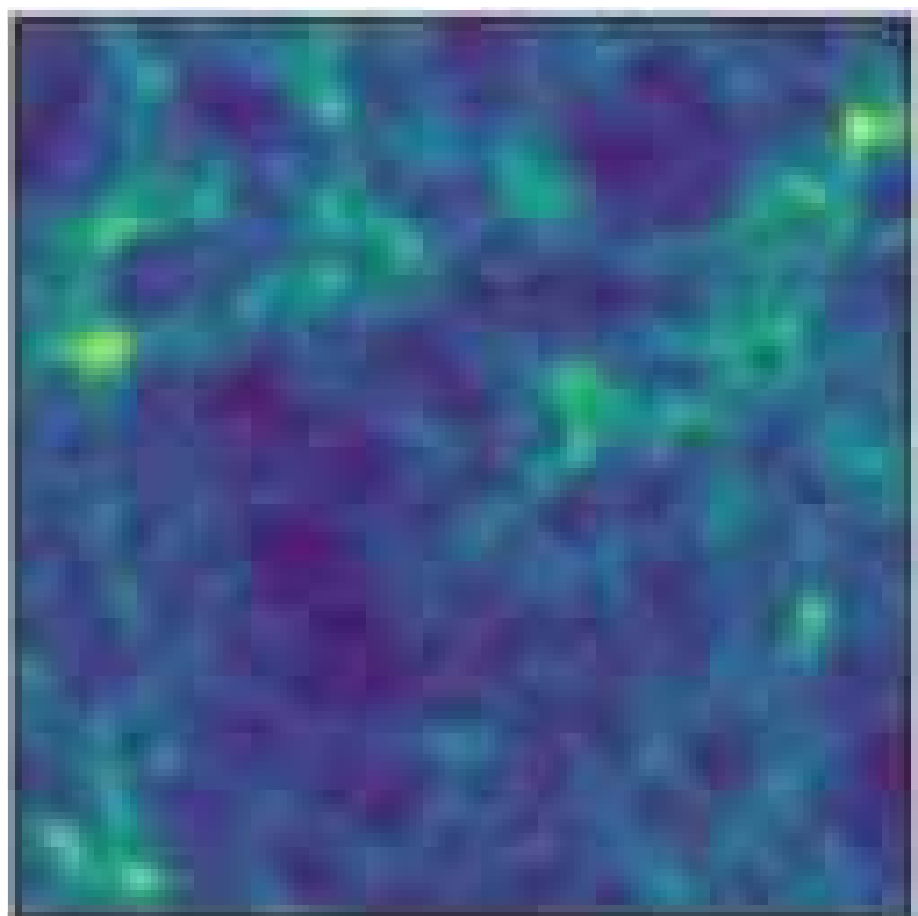
3D intensity distribution
(= 3D galaxy distributions)
is measured



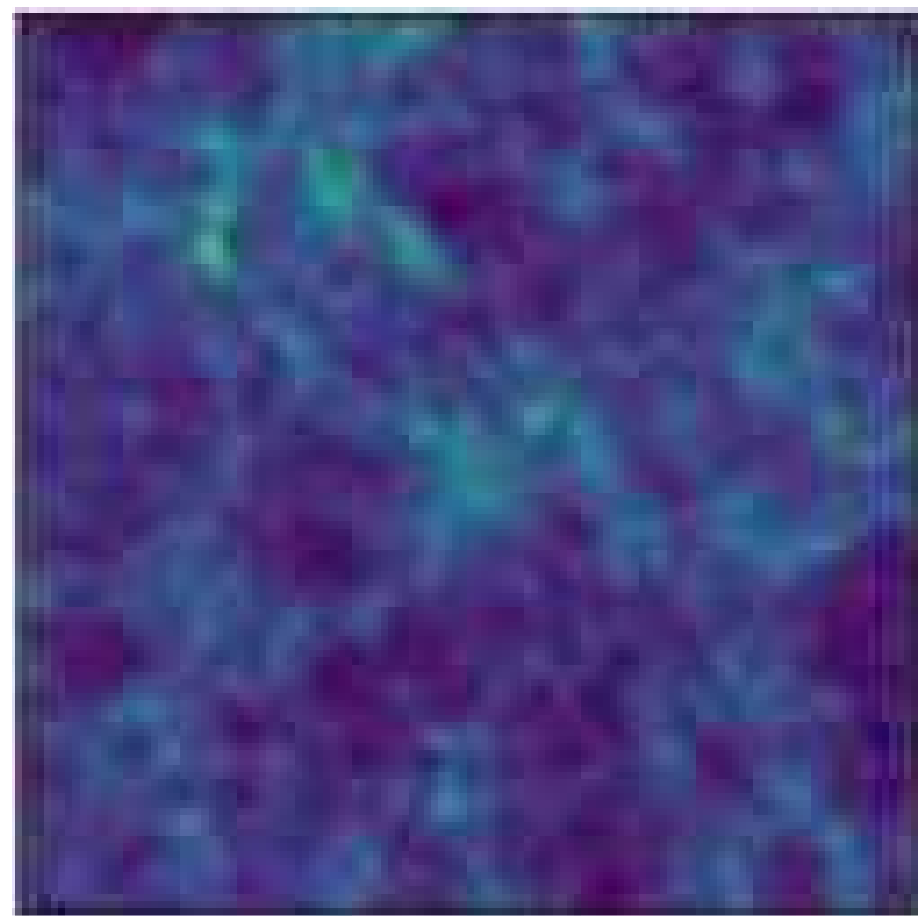
Distance from us

A Serious Problem: Contaminations and Noises

Hydrogen line signals
from near galaxies



Oxygen line signals
from distant galaxies



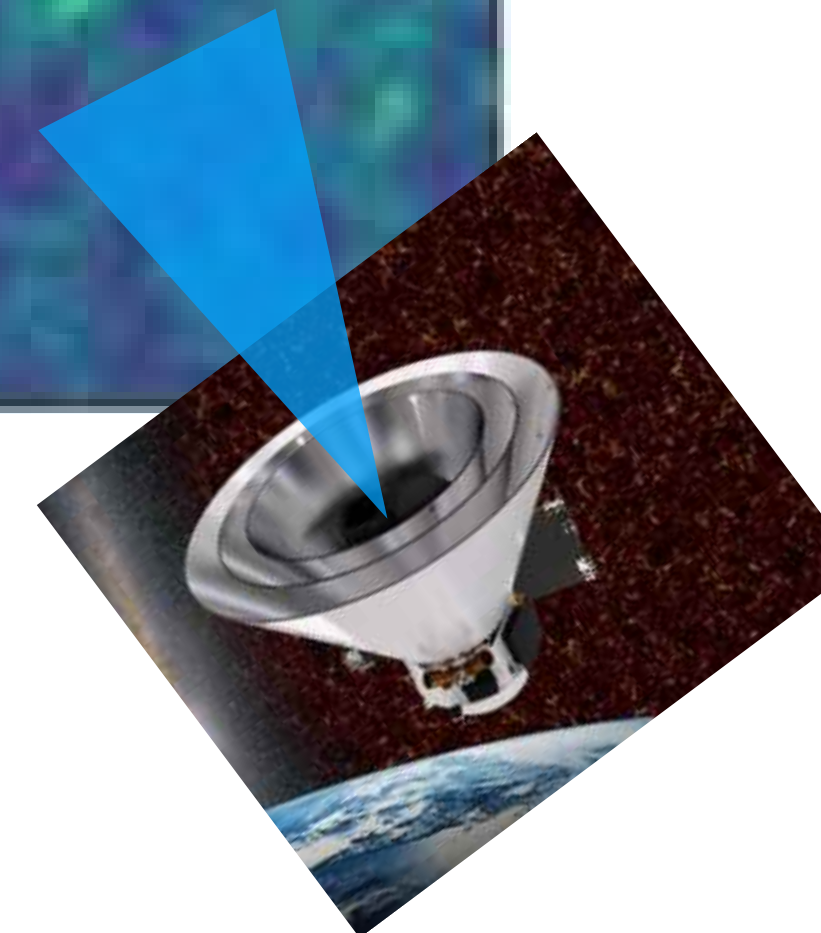
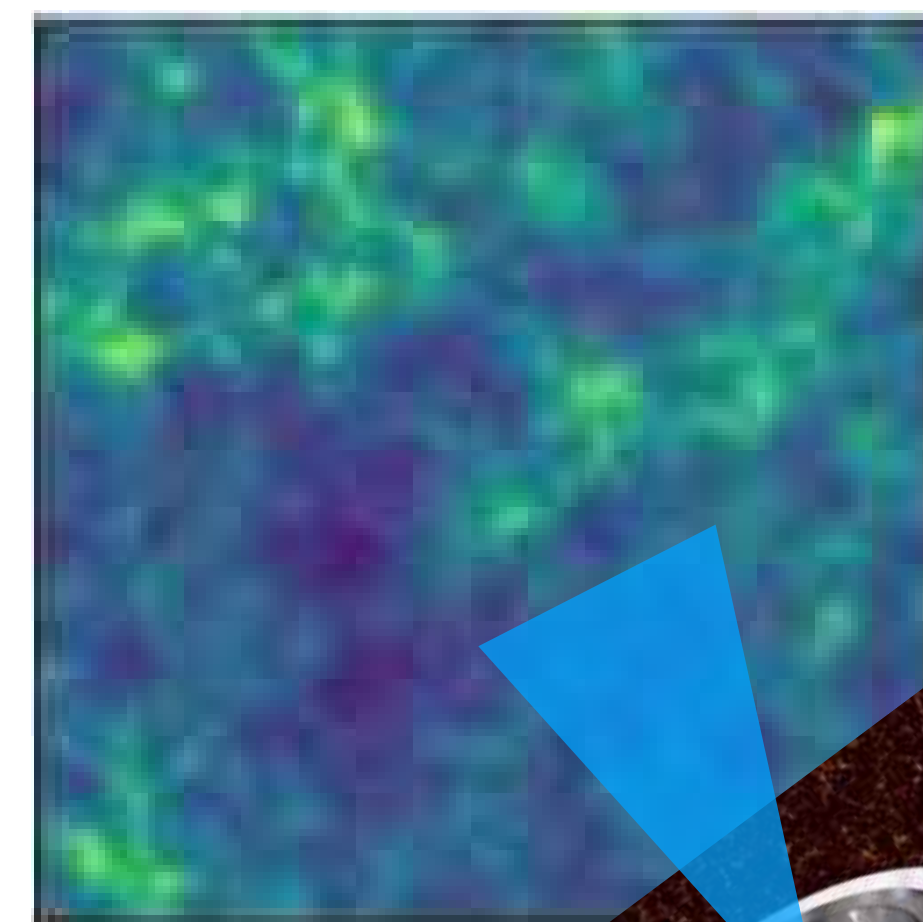
+

+

obs.
noise



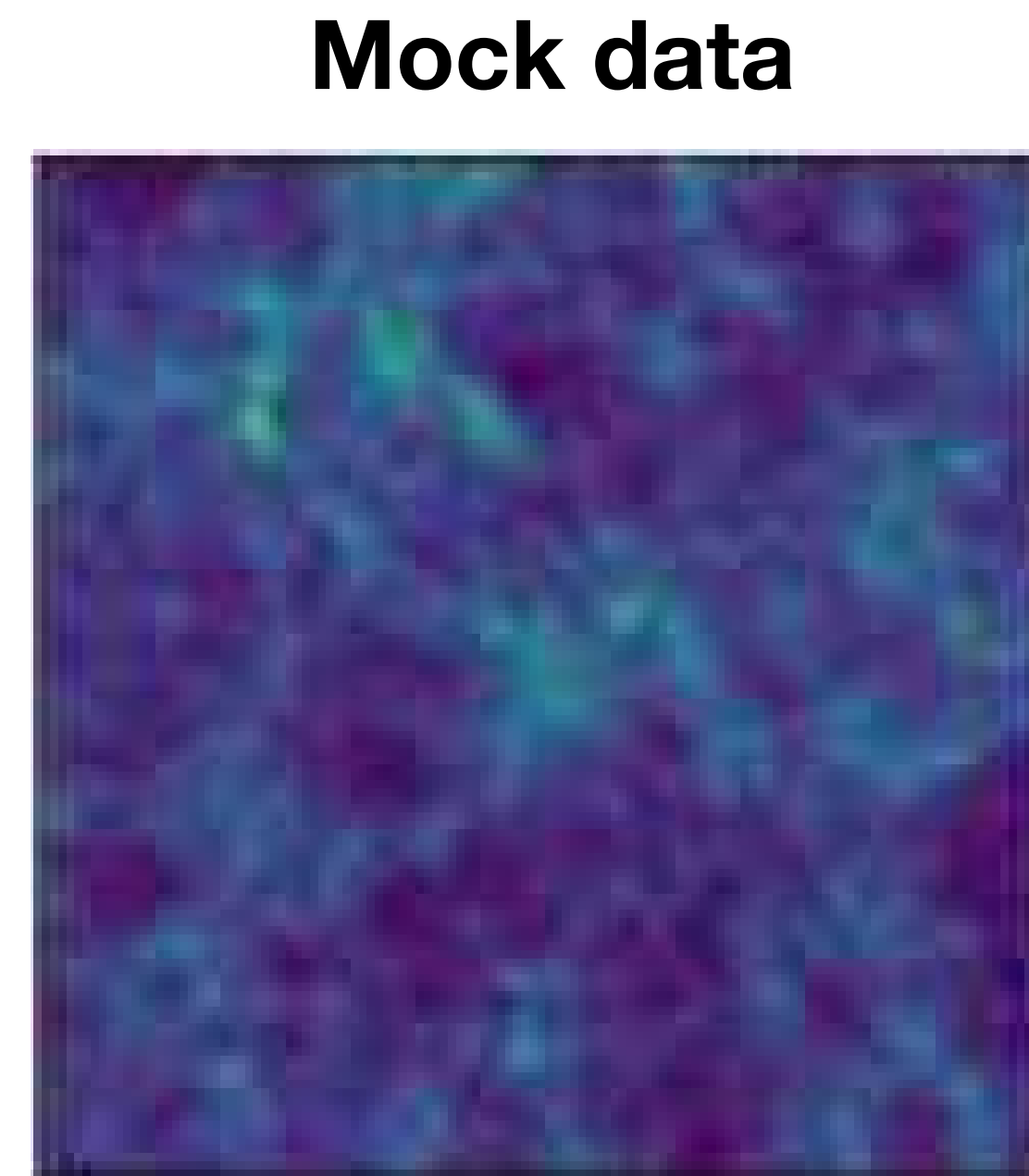
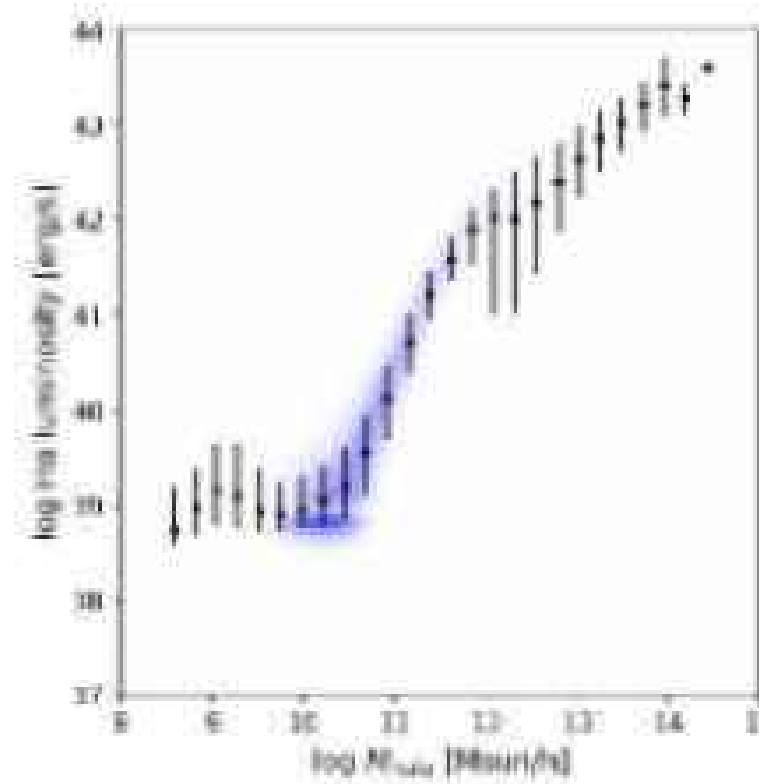
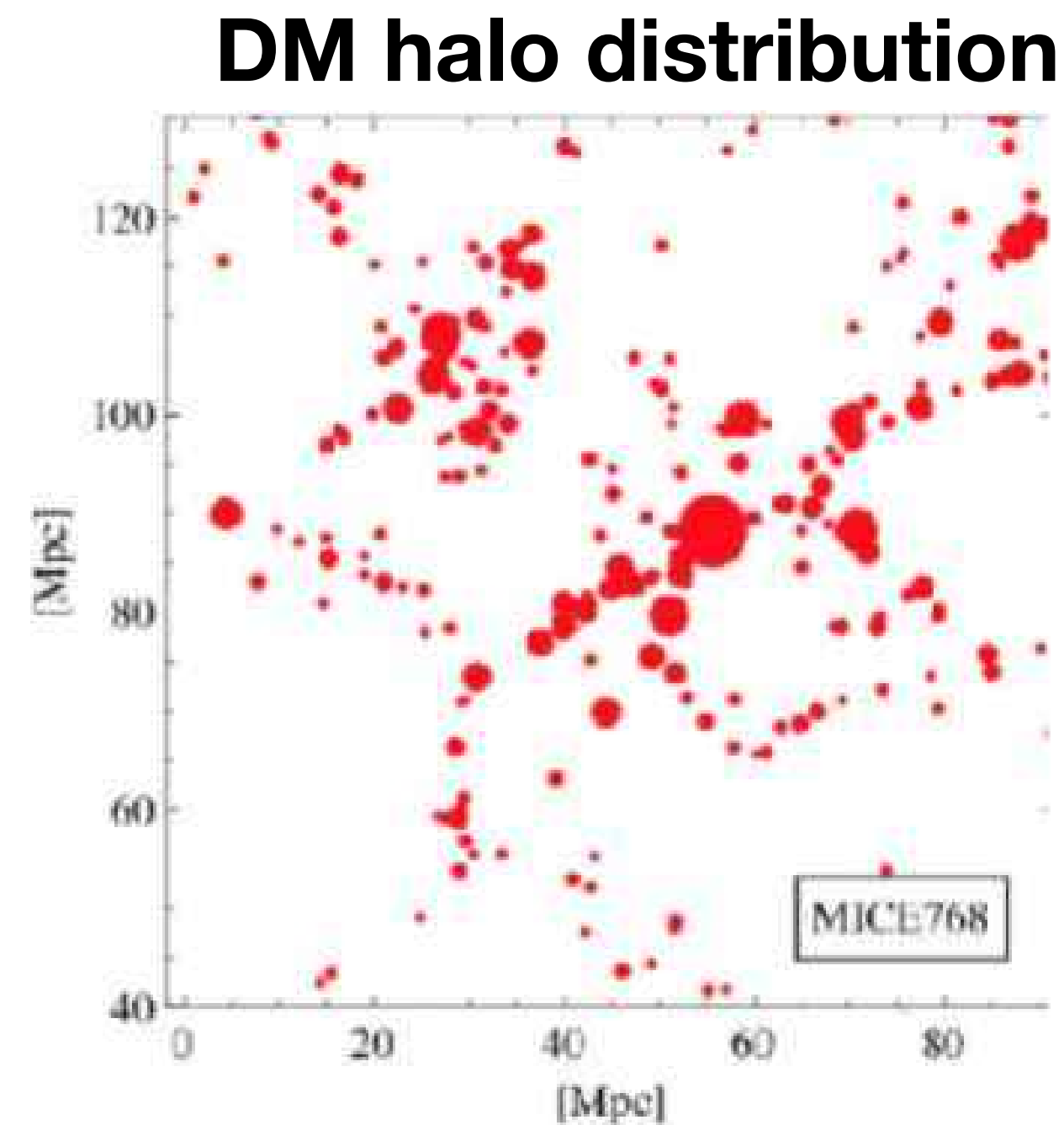
Observed data
at wavelength λ_{obs}



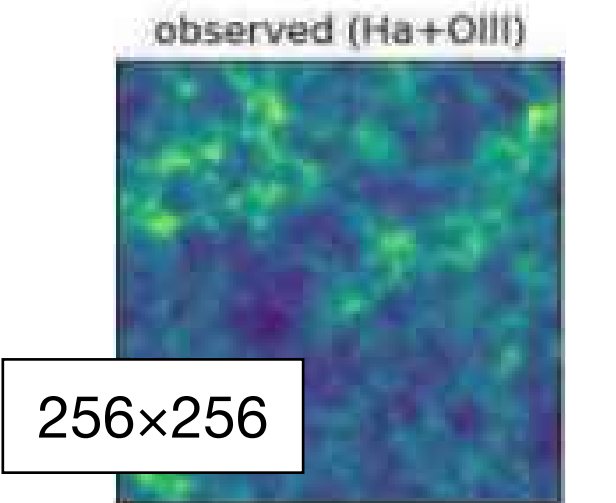
Train a Deep Learning Model with Mock Observational Data

Generate ~30,000 realistic mock observational maps using fast DM simulatino code + emission line model

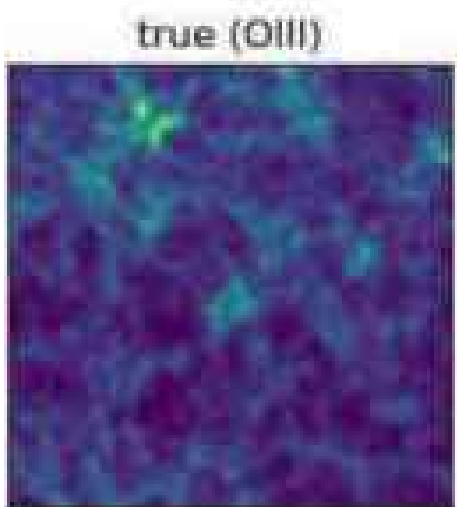
**Emission line model
(mass-to-luminosity relation)**



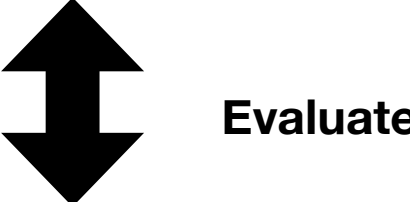
**× 2 × 30,000
+ noise maps**



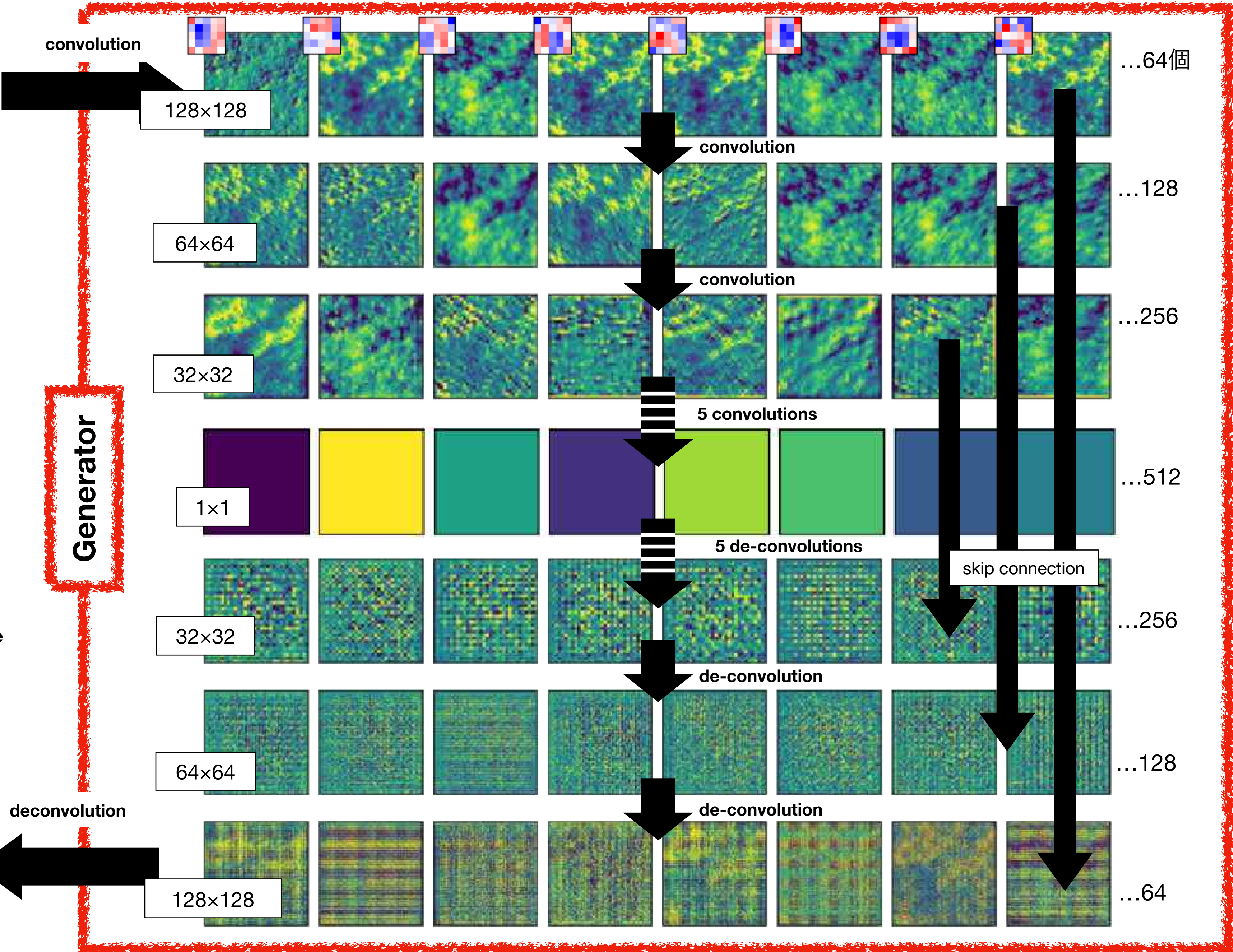
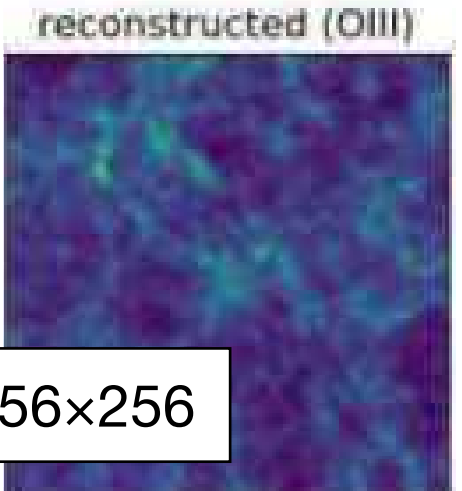
input: Observed map



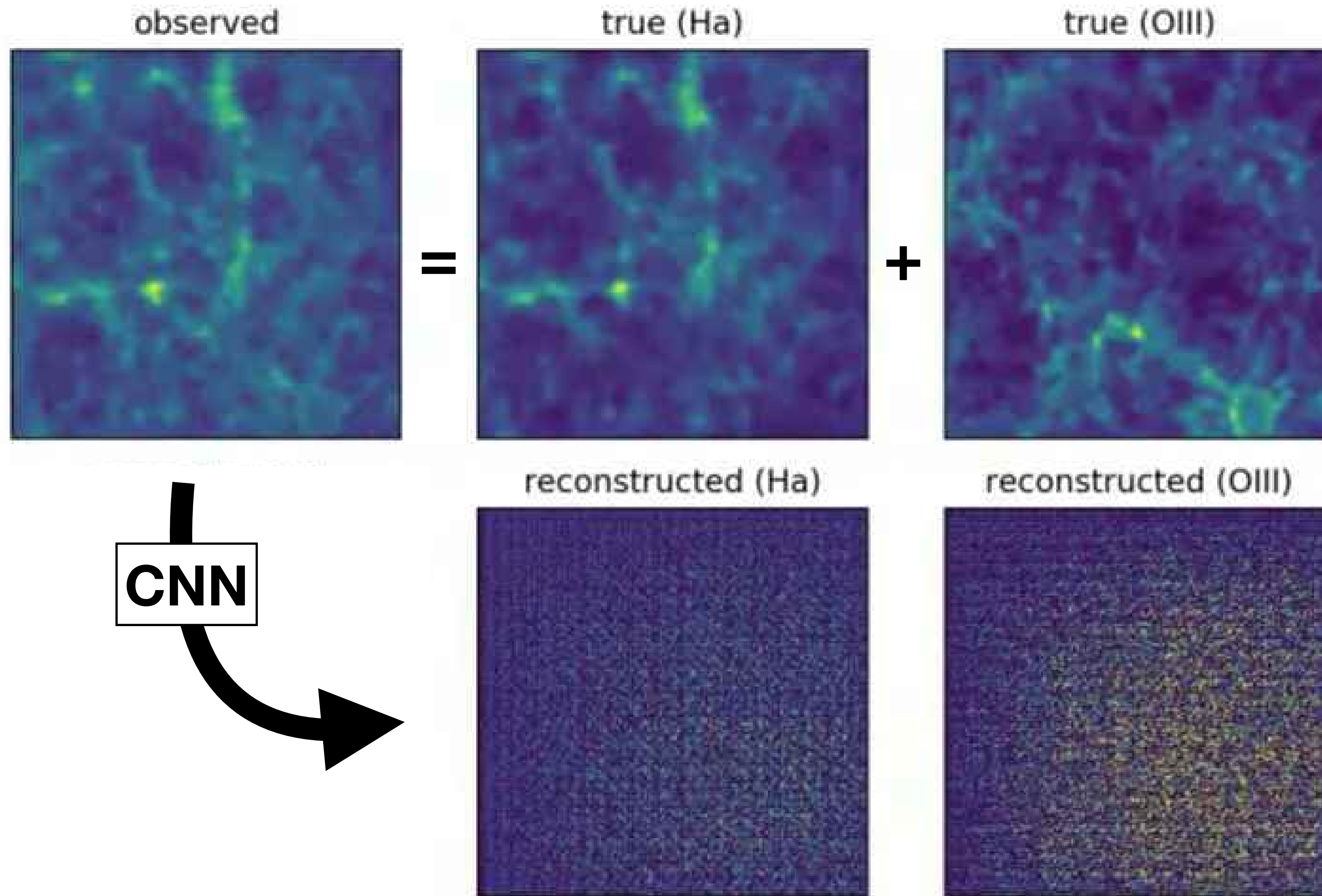
Ground truth X_{true}



output: Signal map

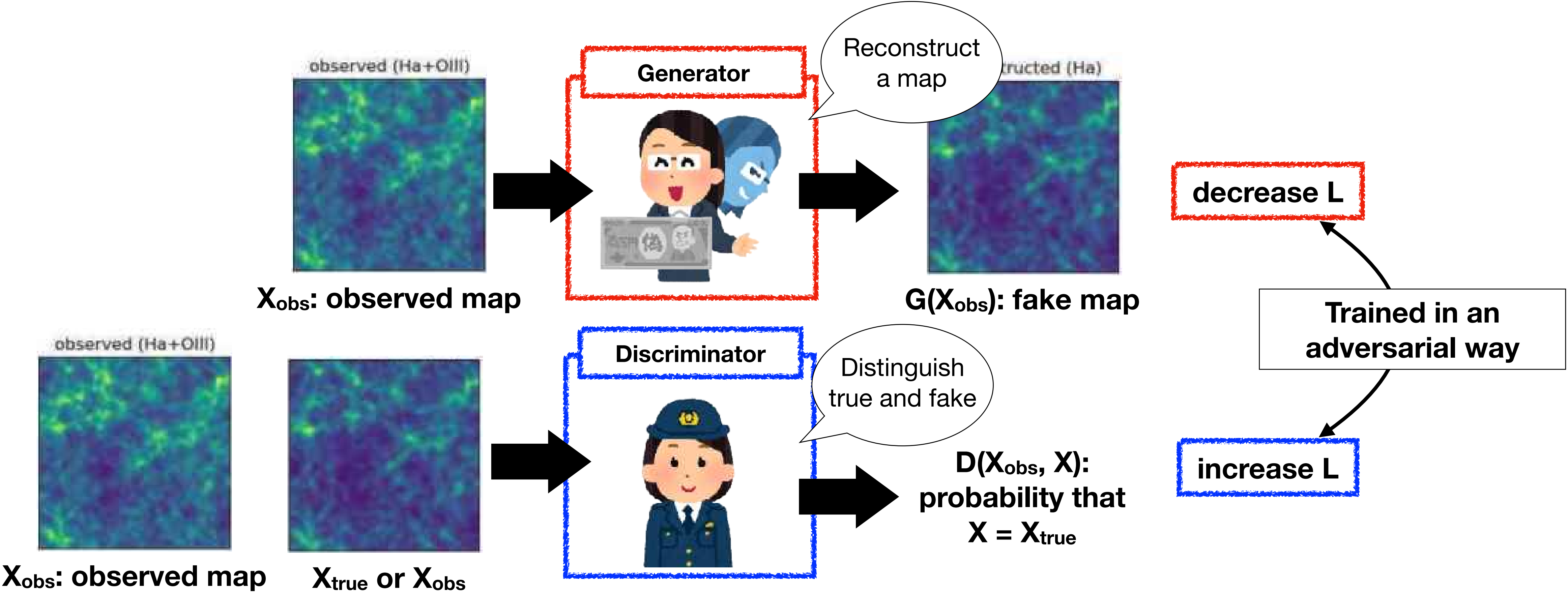


Machine Learns the Large-scale Structure...



Conditional Generative Adversarial Network

- GAN: **Generator** and **Discriminator** are updated in an adversarial way.

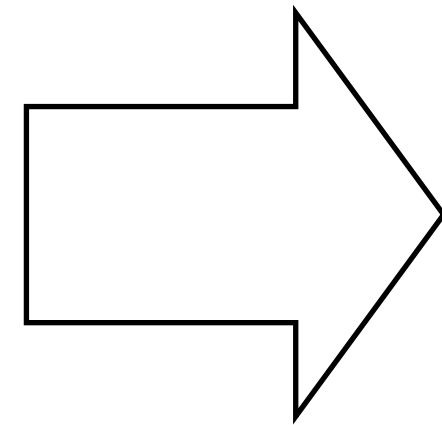
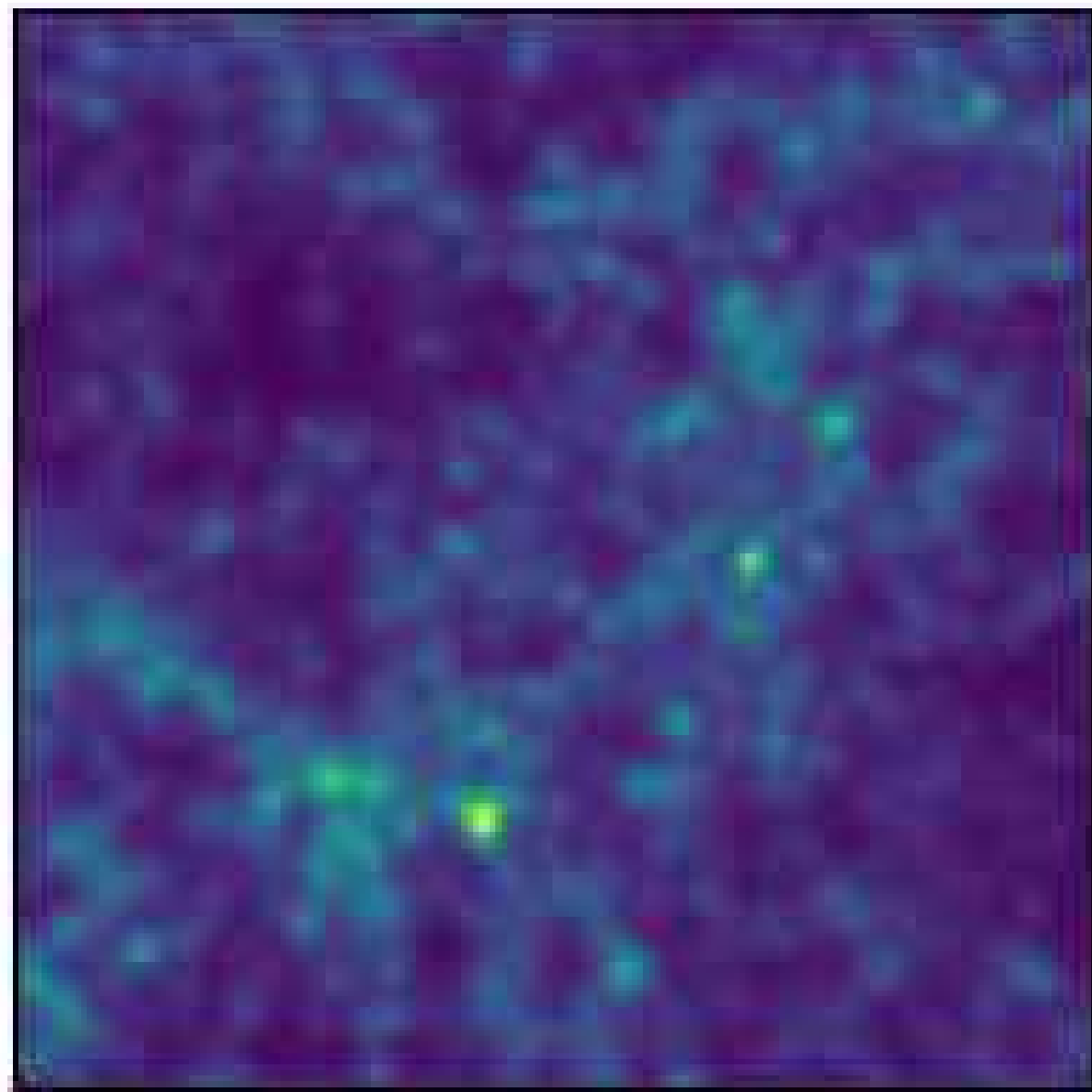


Loss function:
$$L[G, D] = \log D(X_{\text{obs}}, X_{\text{true}}) + \log[1 - D(X_{\text{obs}}, G(X_{\text{obs}}))] + \lambda \langle \|X_{\text{true}} - G(X_{\text{obs}})\| \rangle$$

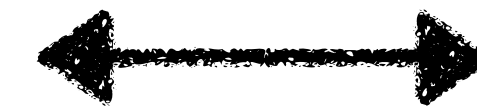
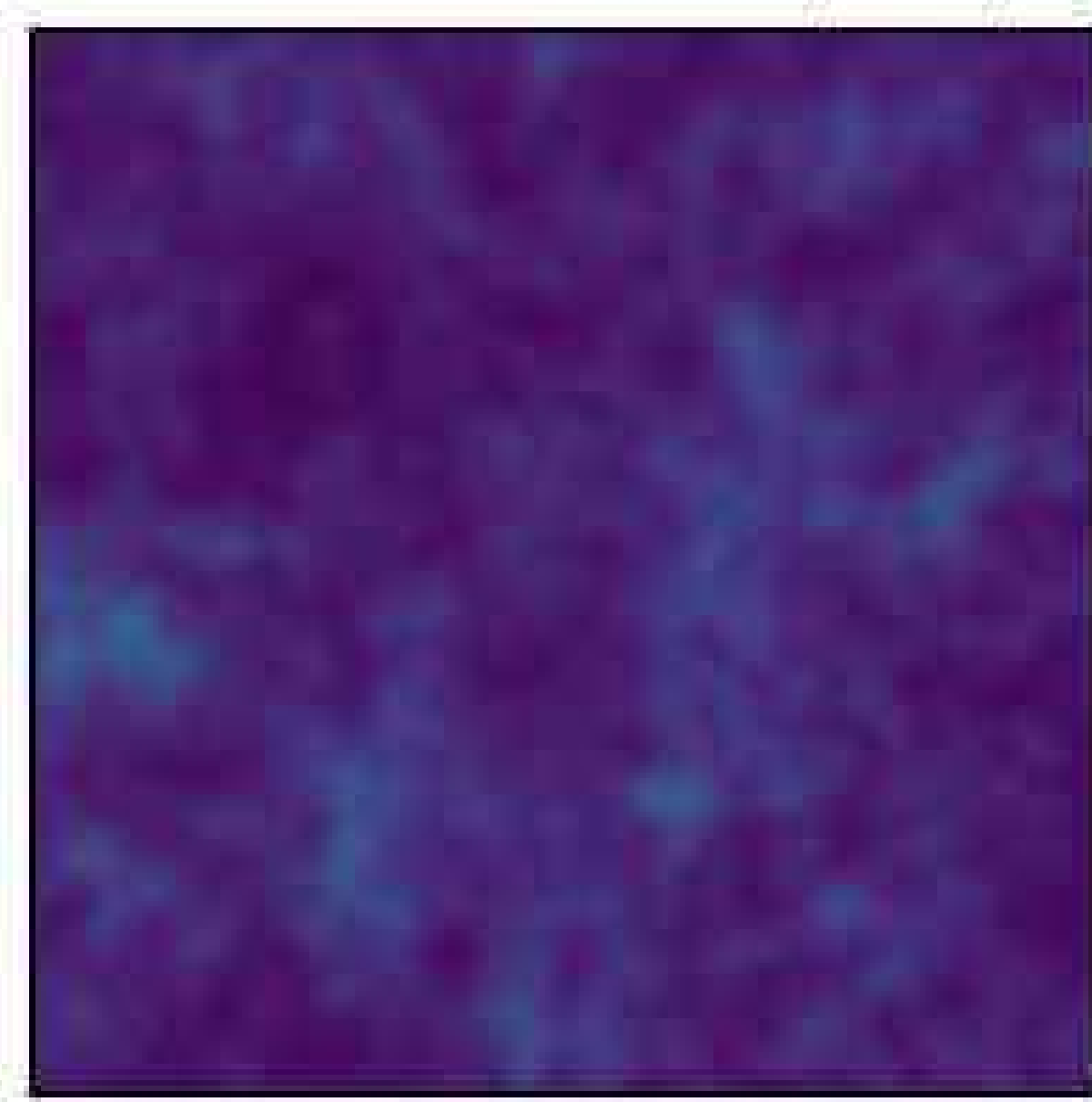
What if we do not use GAN?

The network tends to reproduce obscured images

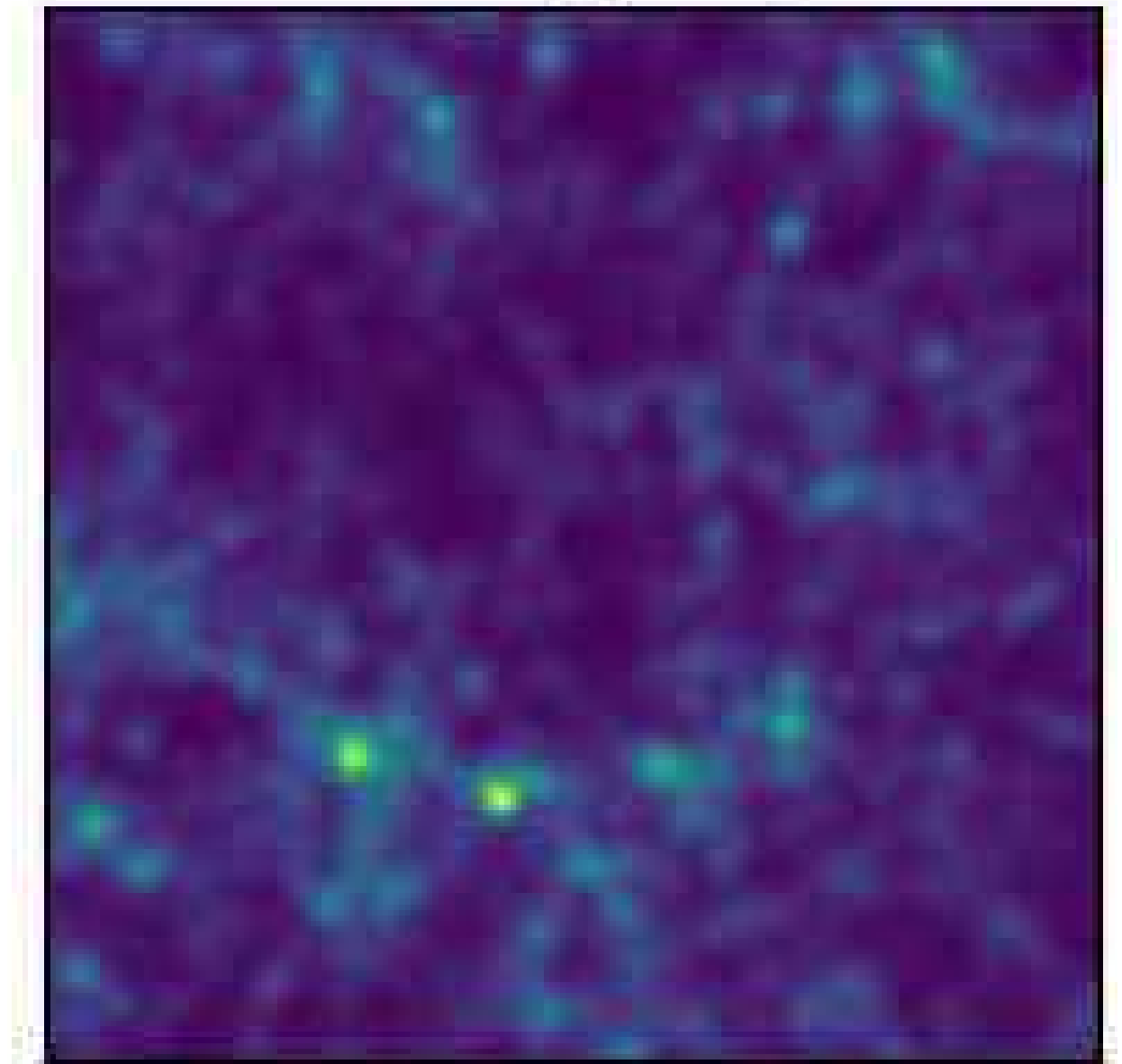
Observed (Line1+Line2)



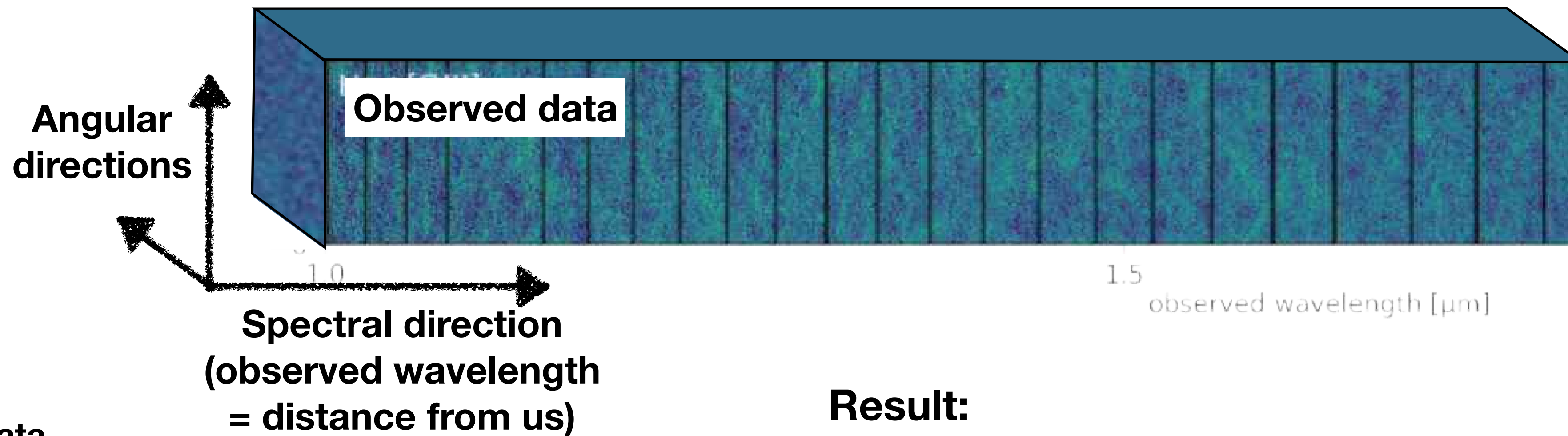
Reconstruct (Line2)



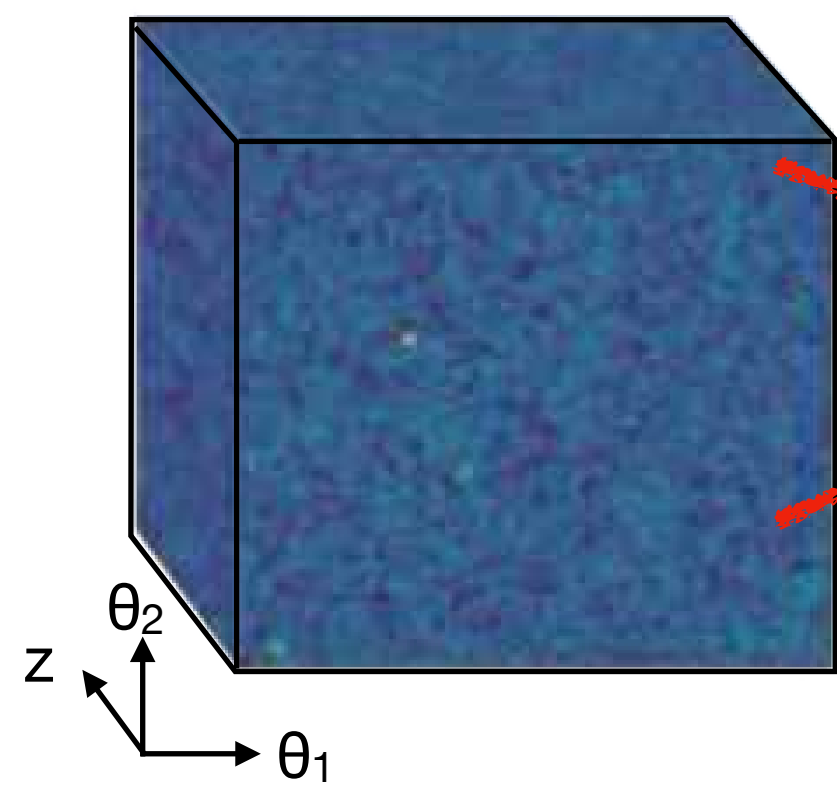
True (Line2)



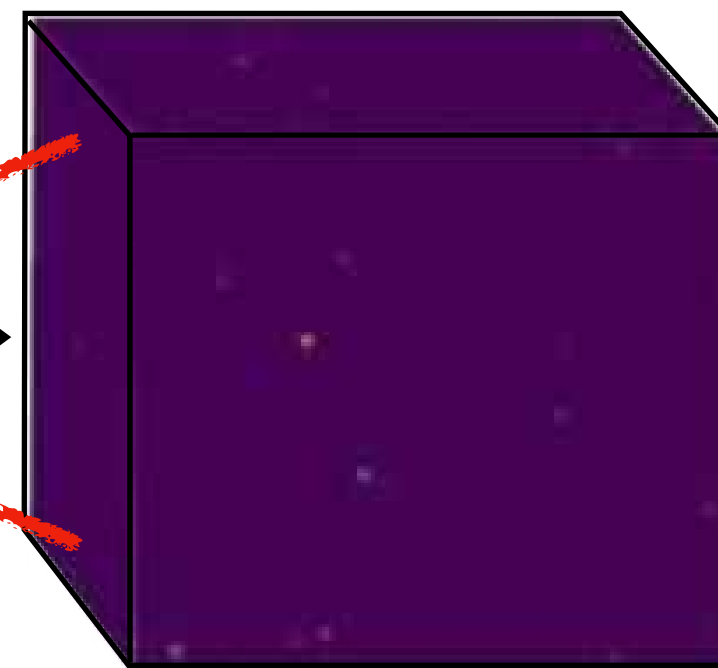
Reconstruction of 3D Maps



Observed data
(Line1 + Line 2 + noise)

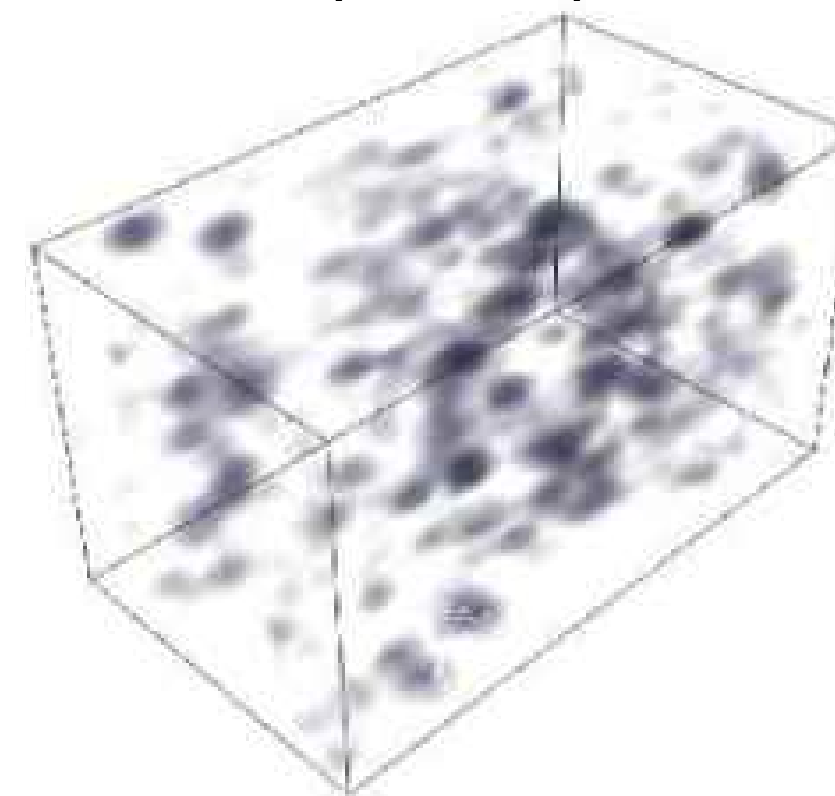


Line1

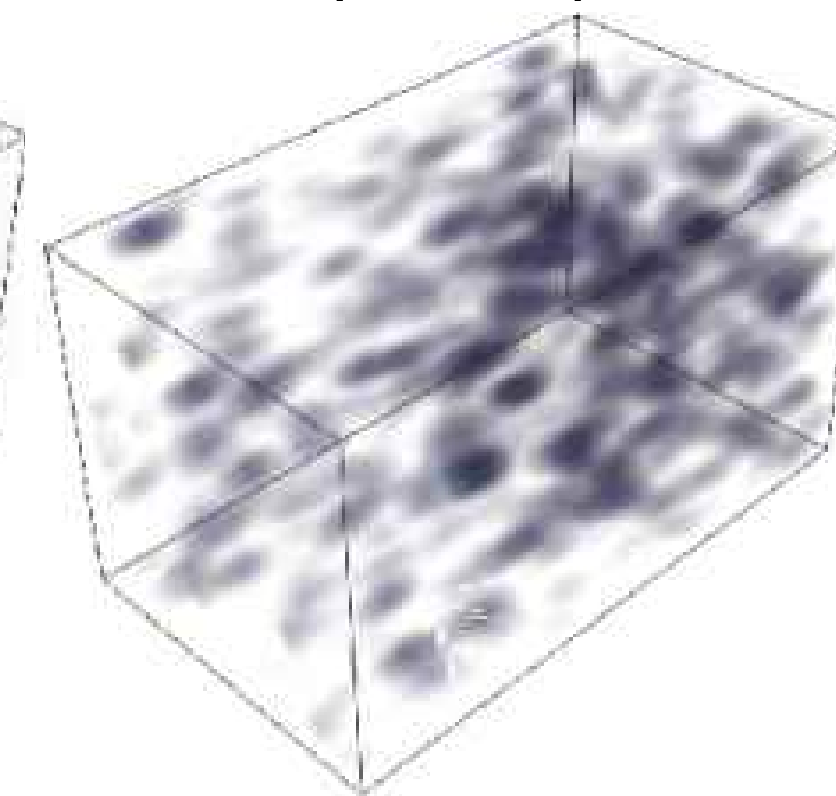


Result:

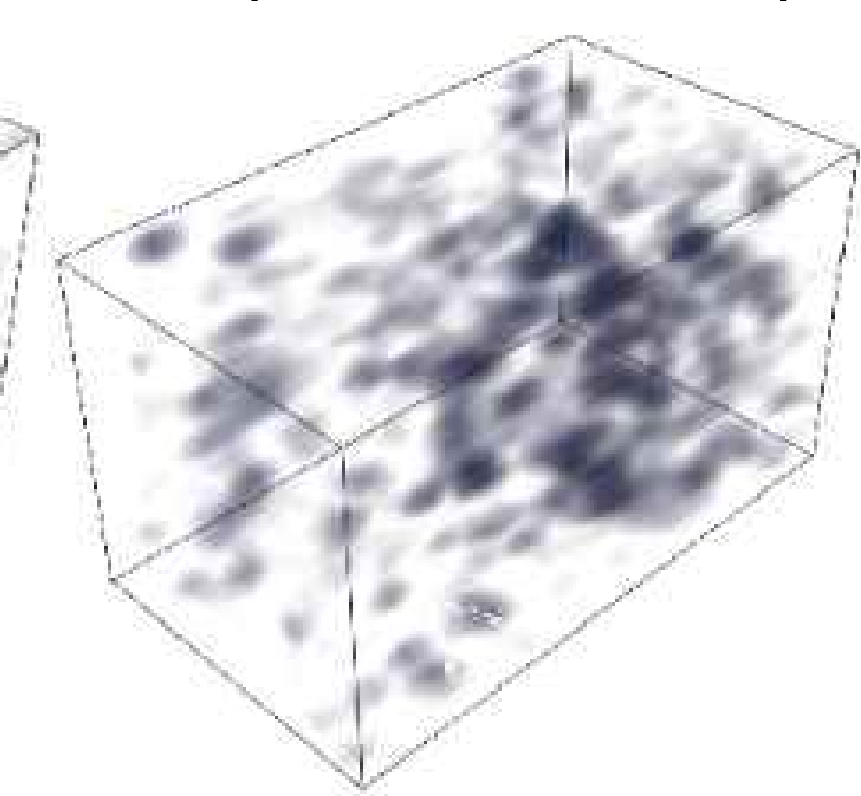
True
(Line2)



Reconstruct
(Line2)

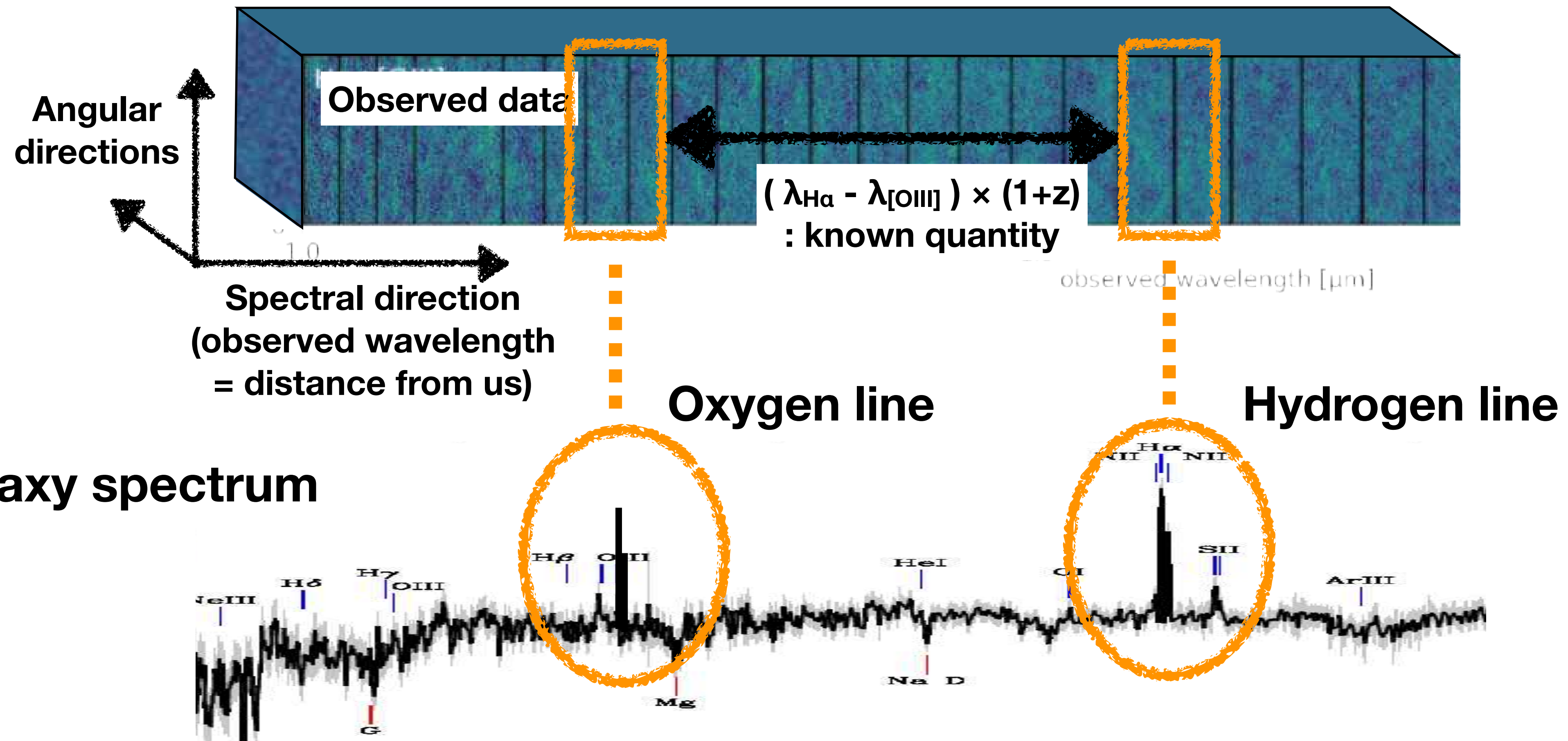


cf.) True
(Line1+Line2)



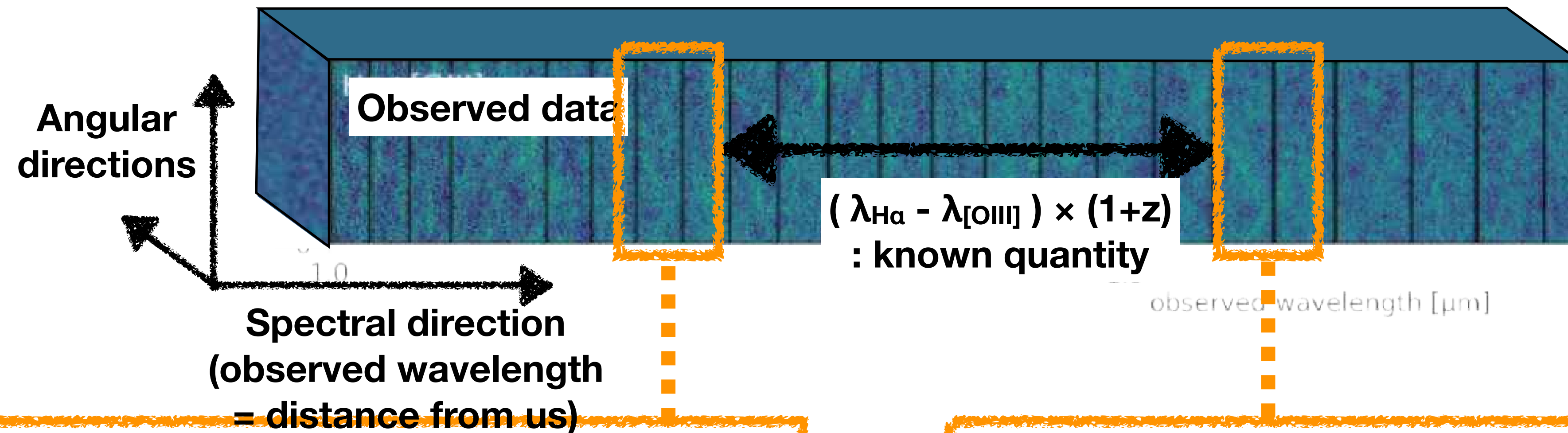
A simple network does not work well...

Pre-processing Input Data with *Physical information*

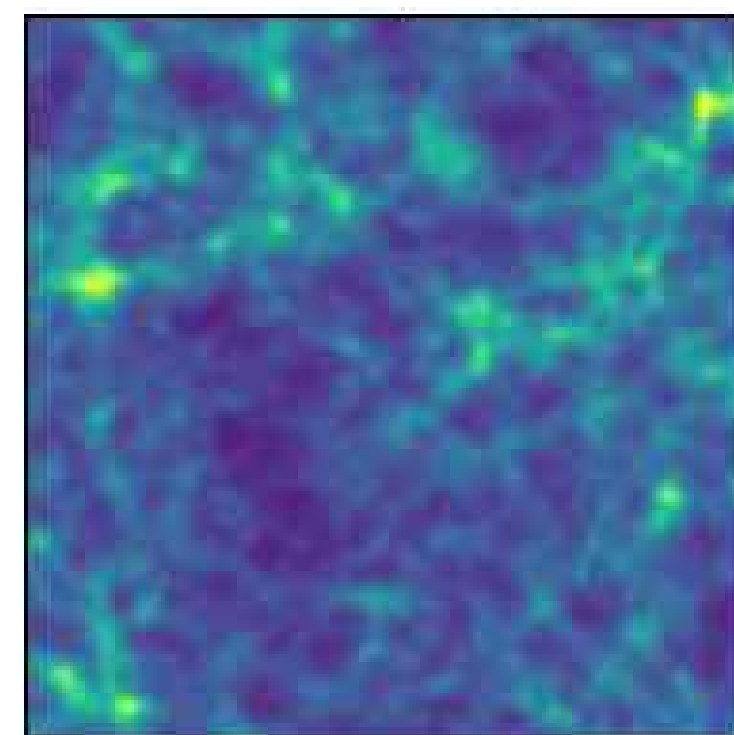


The contribution from the galaxy appears in two wavelengths

Pre-processing Input Data with *Physical information*



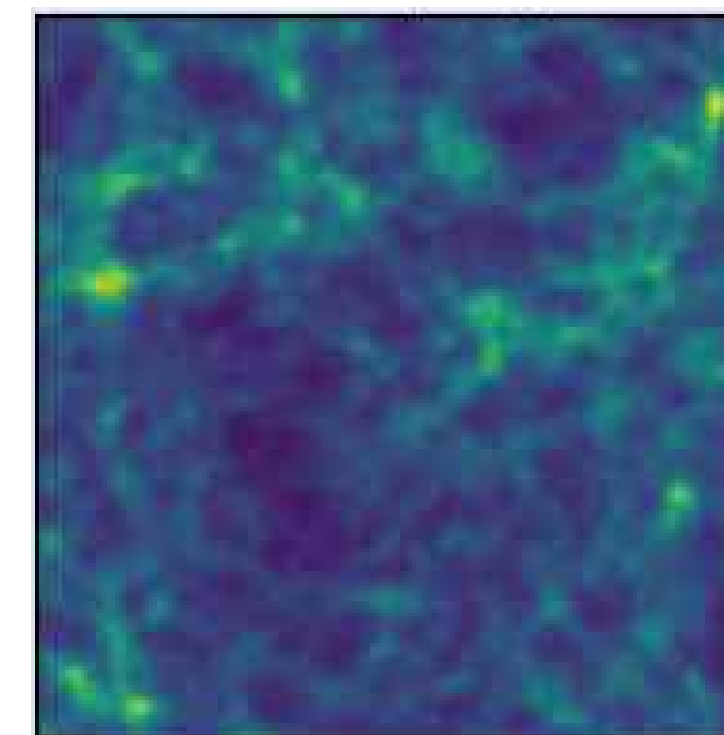
The structure traced
by oxygen line



+ contaminations
& obs. noise

(Almost) the same signals are
buried in the two data cubes

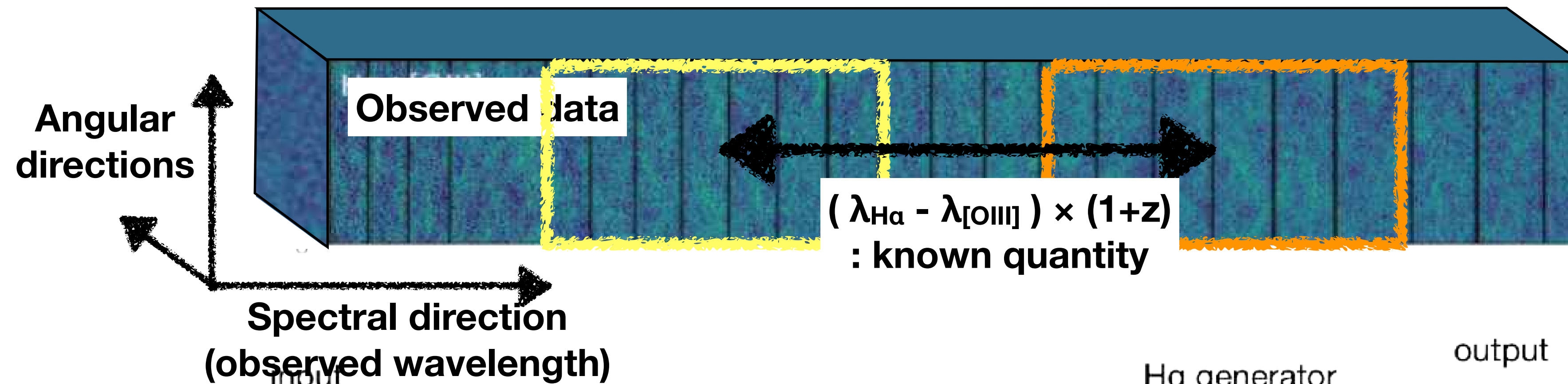
The structure traced
by hydrogen line



+ contaminations
& obs. noise

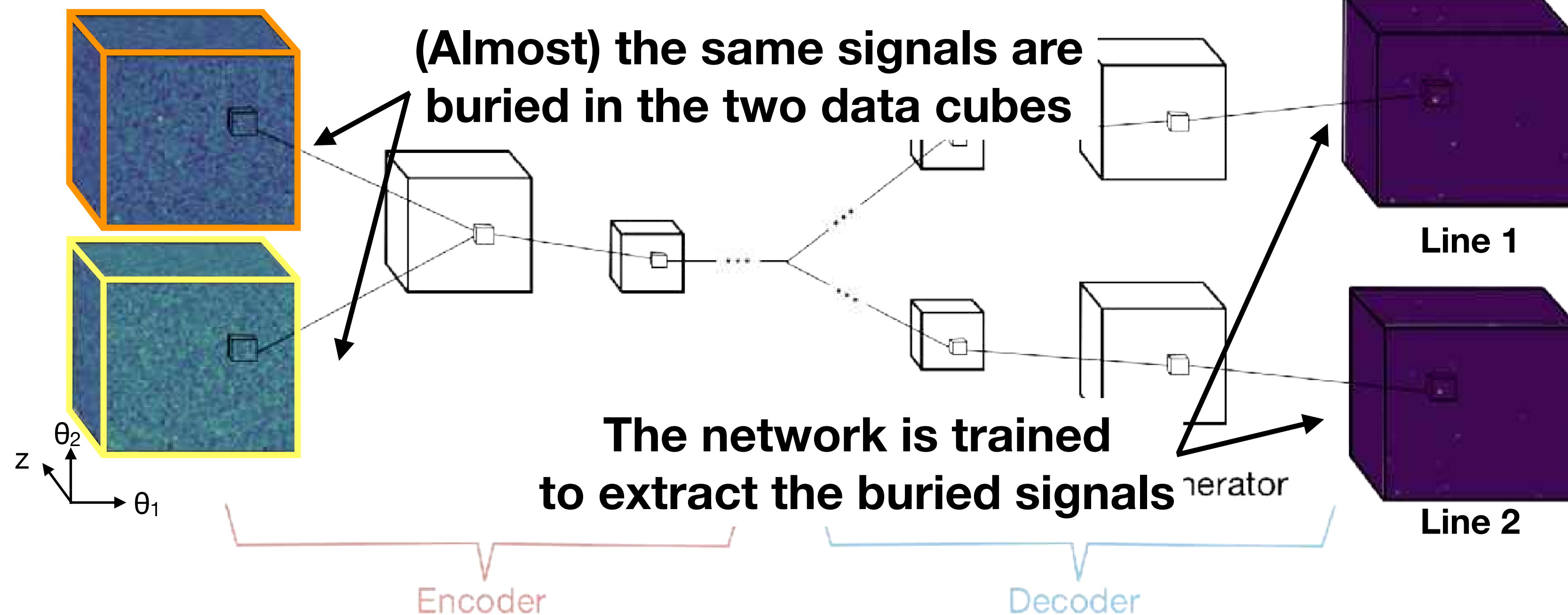
Pre-processing Input Data with *Physical information*

KM & Yoshida 2021



Input 1

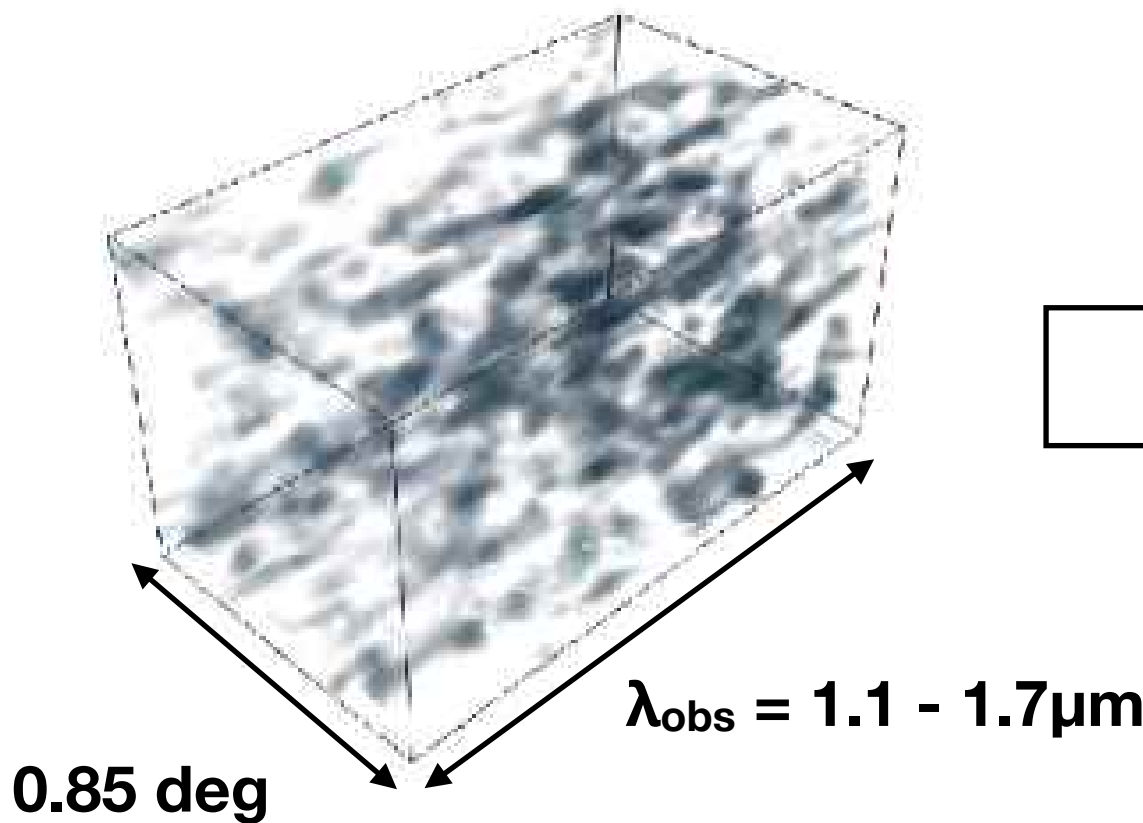
Input 2



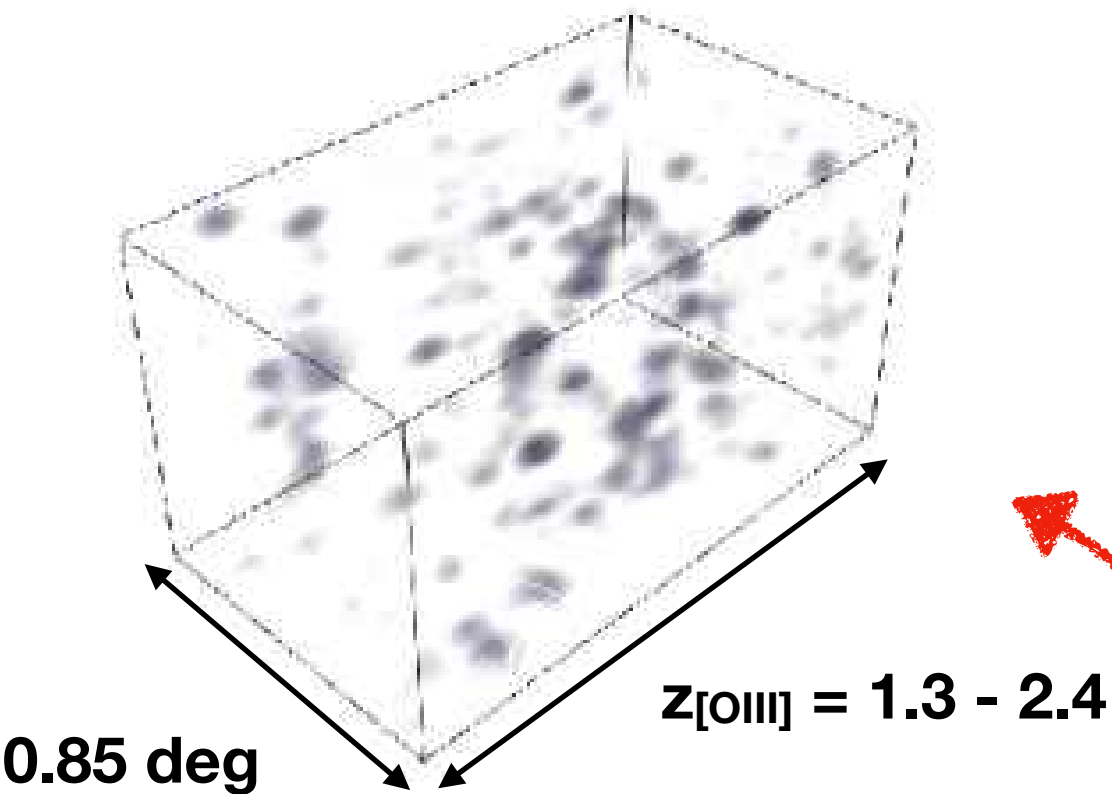
Reconstruction Result

KM & Yoshida 2021

Observed (H α + [OIII]+noise)

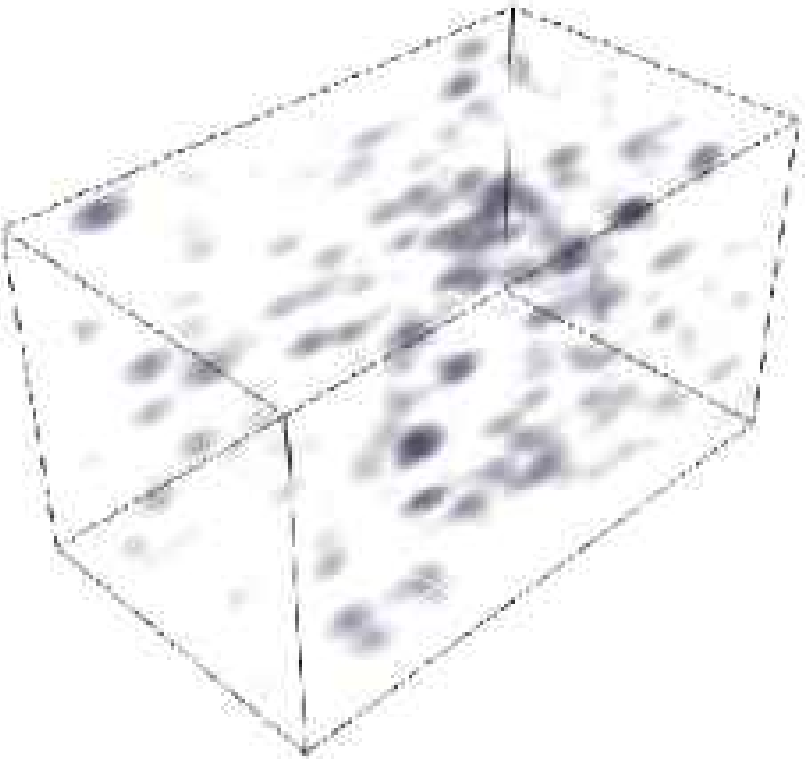


Reconstruct [OIII]

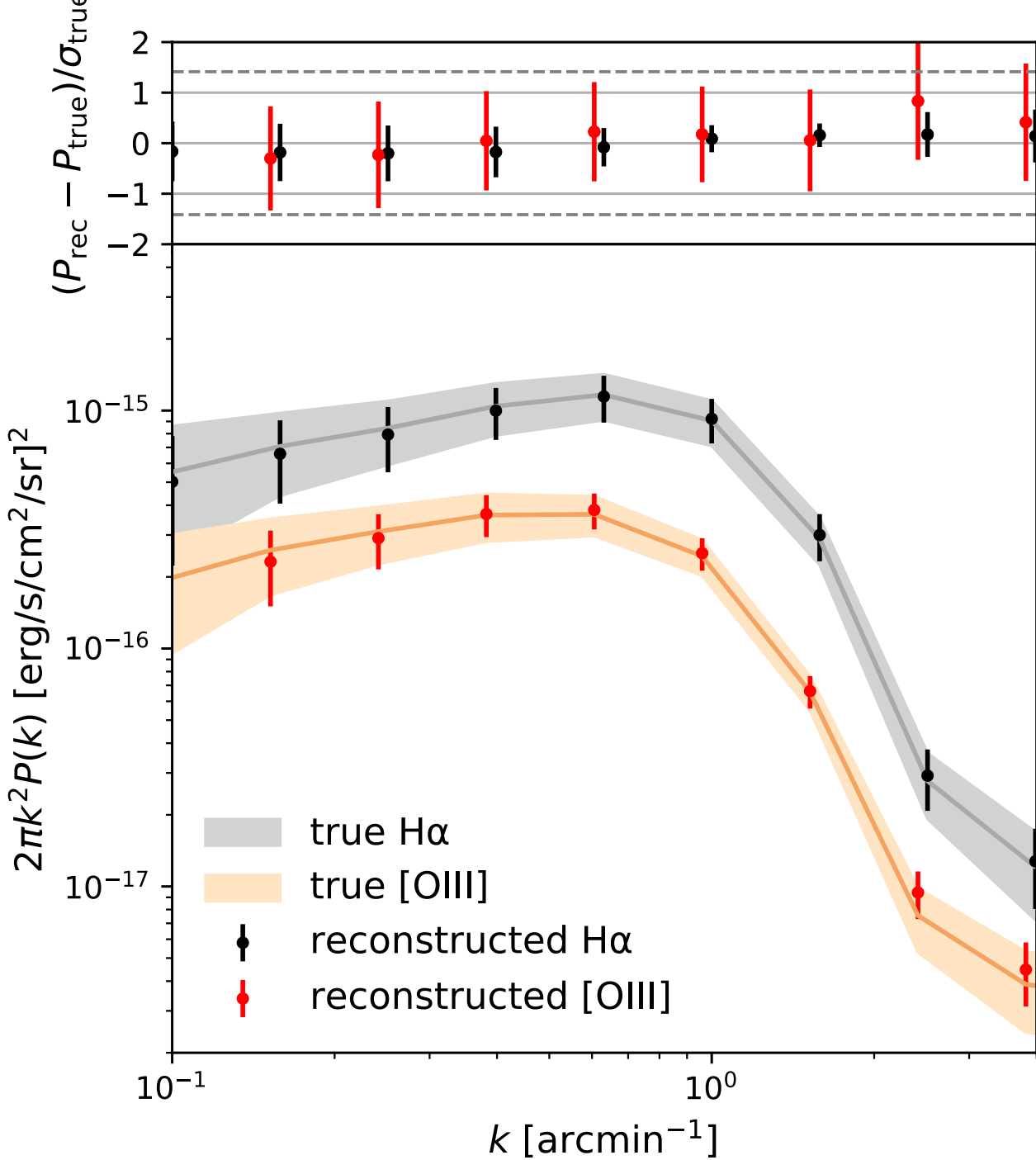


Bright sources are properly reconstructed

True [OIII]



Statistics are also reproduced e.g., Power spectrum

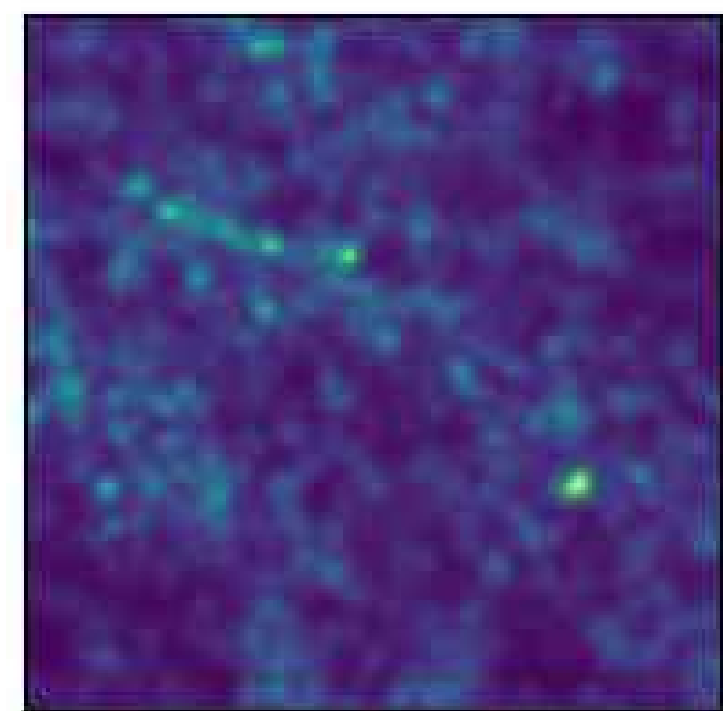


Reconstructed maps can be used for studying cosmology and galaxy formation and evolution

Reproducibility of bright peaks
 Peak detectability of H α and [OIII]
 precision = 82%, 68%
 recall = 80%, 77%

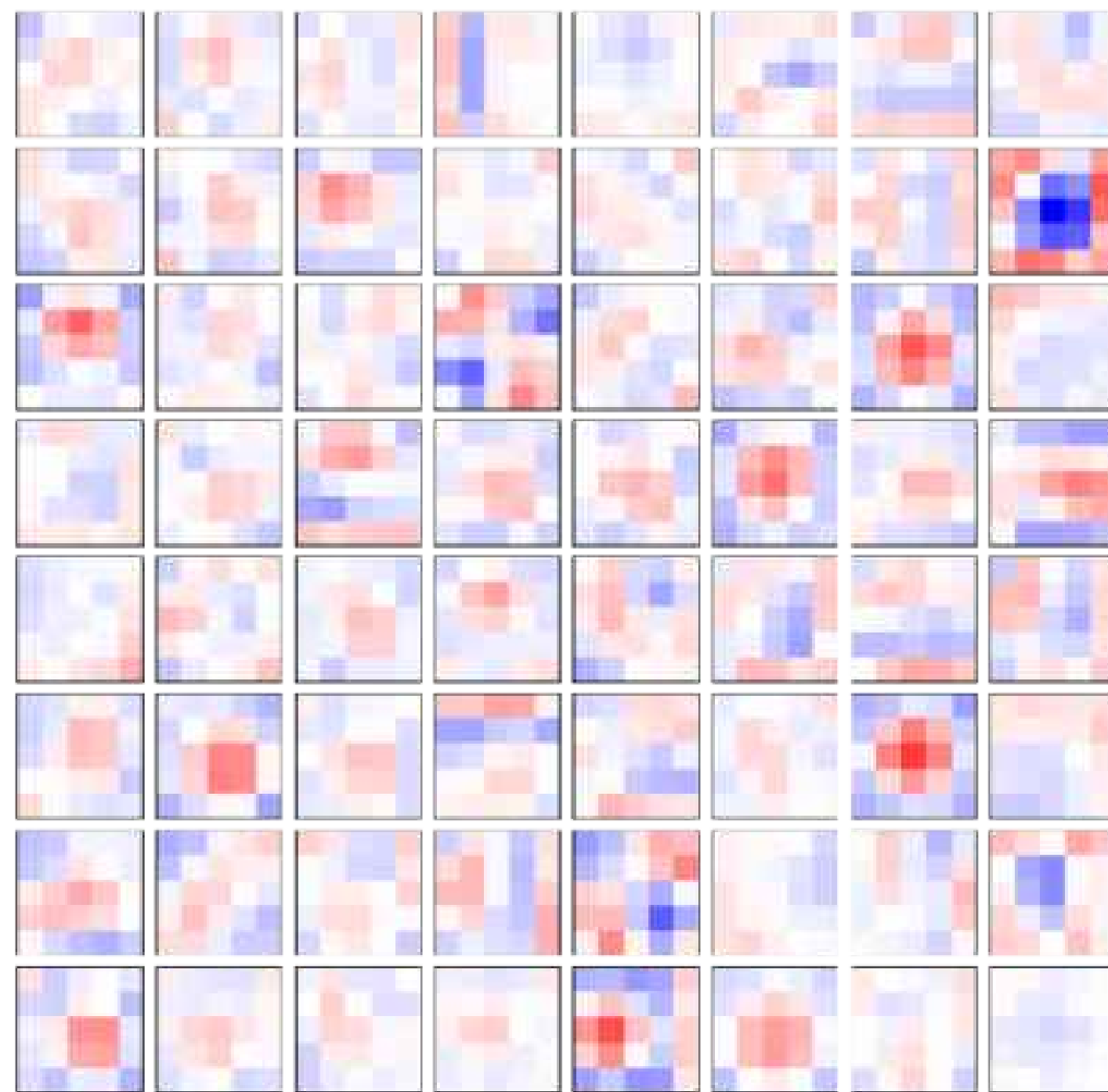
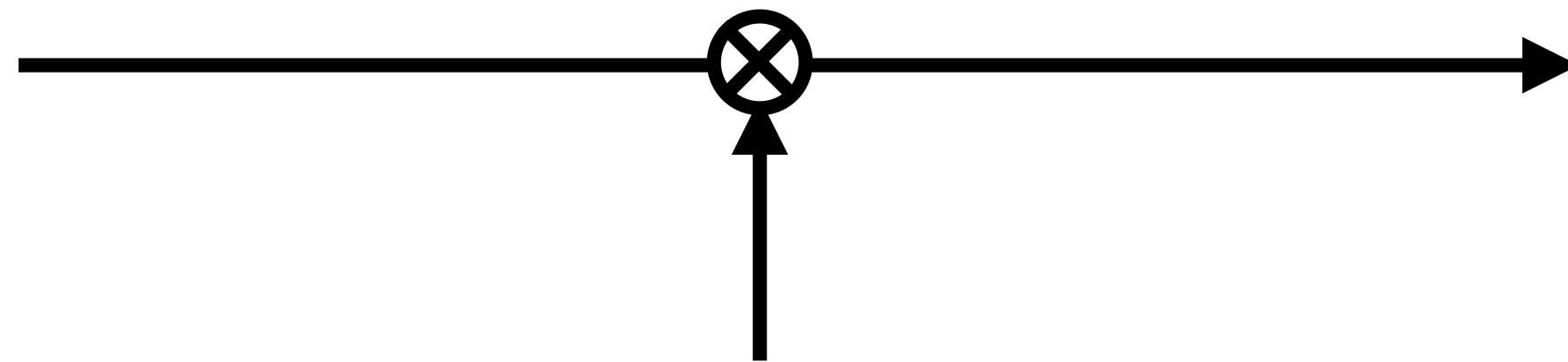
What Does the Machine Learn to Separate the Signals?

Let's have a look at the convolutional filters.

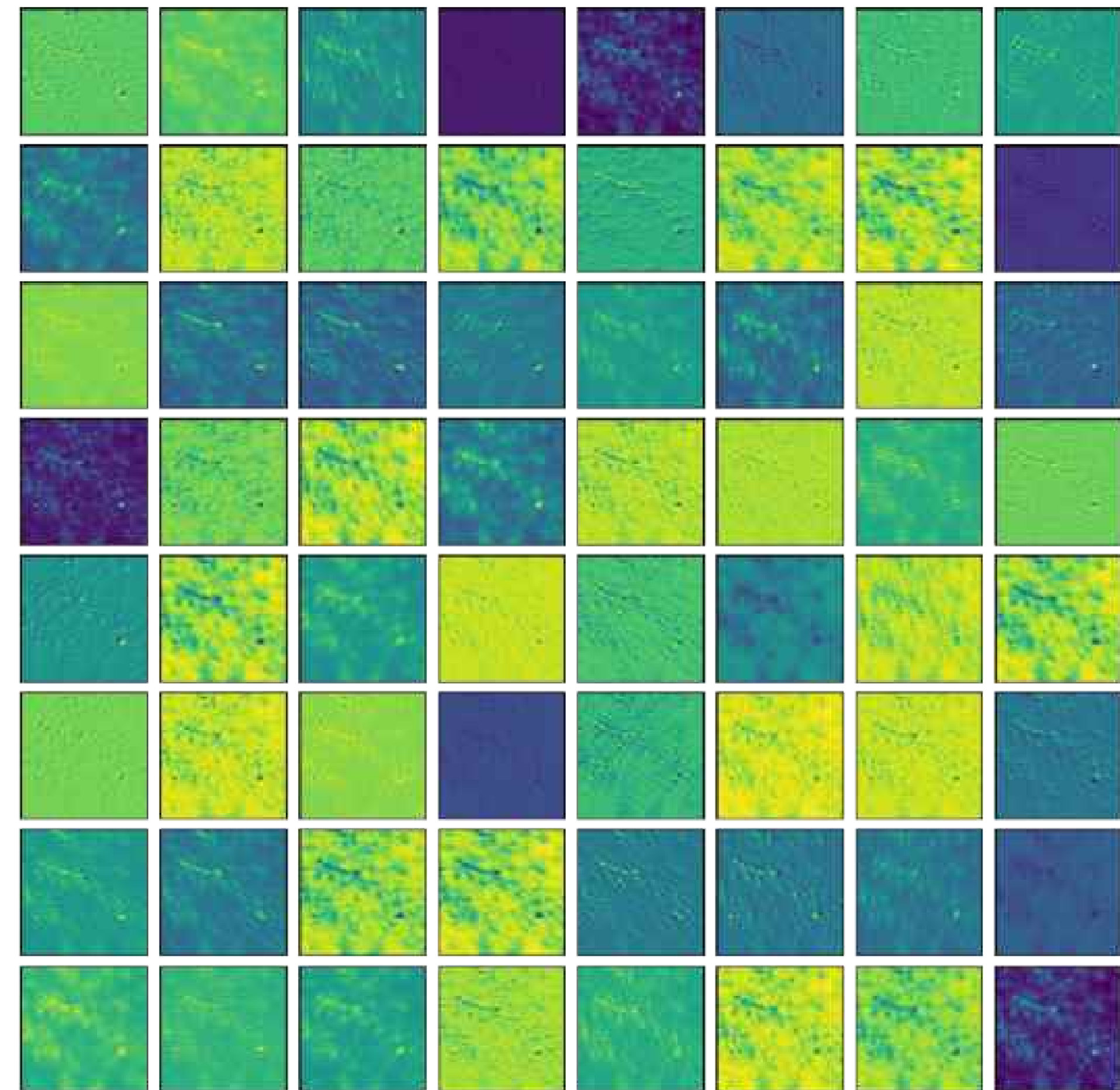


Input (observed)

Convolution



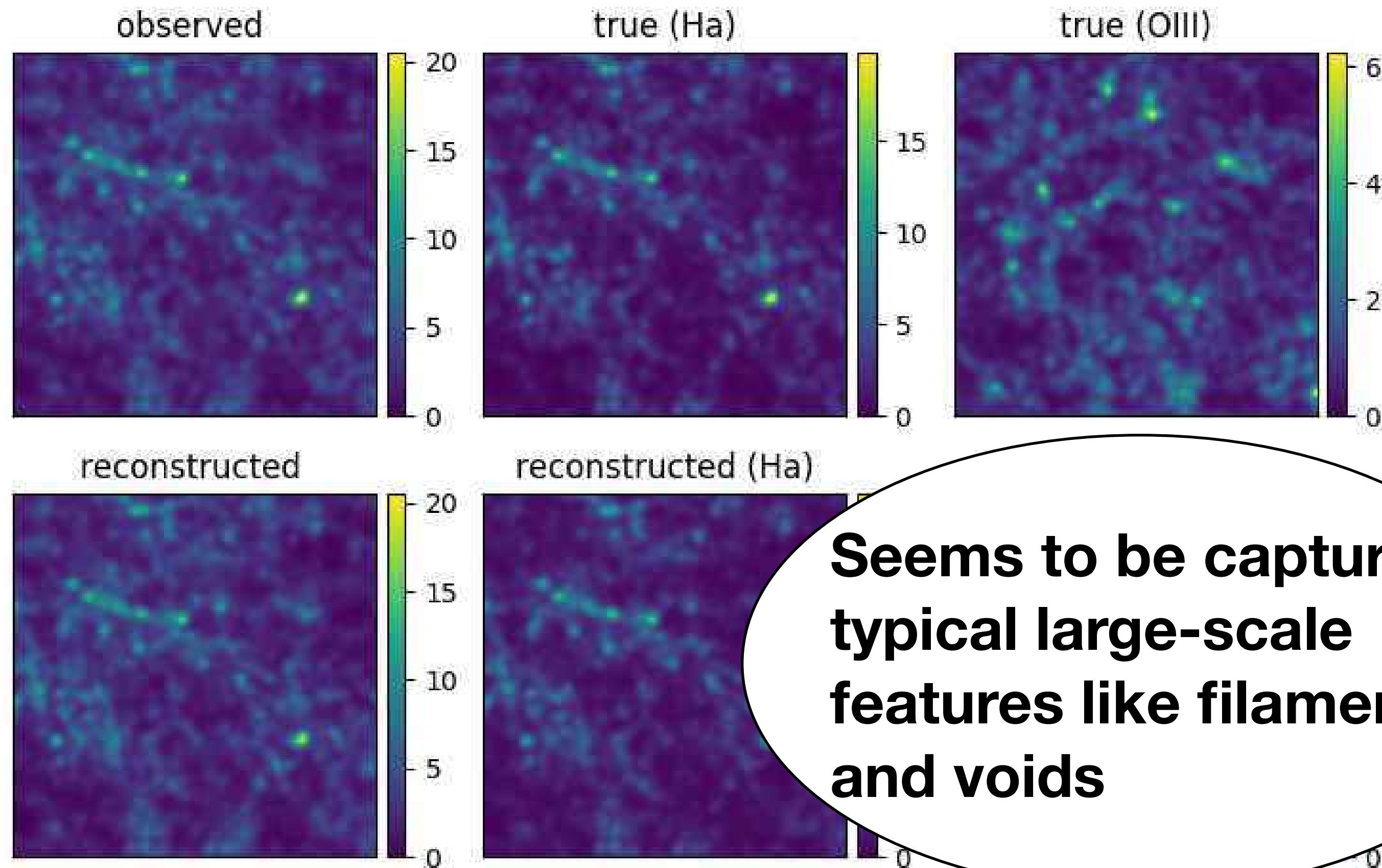
Convolutional filters



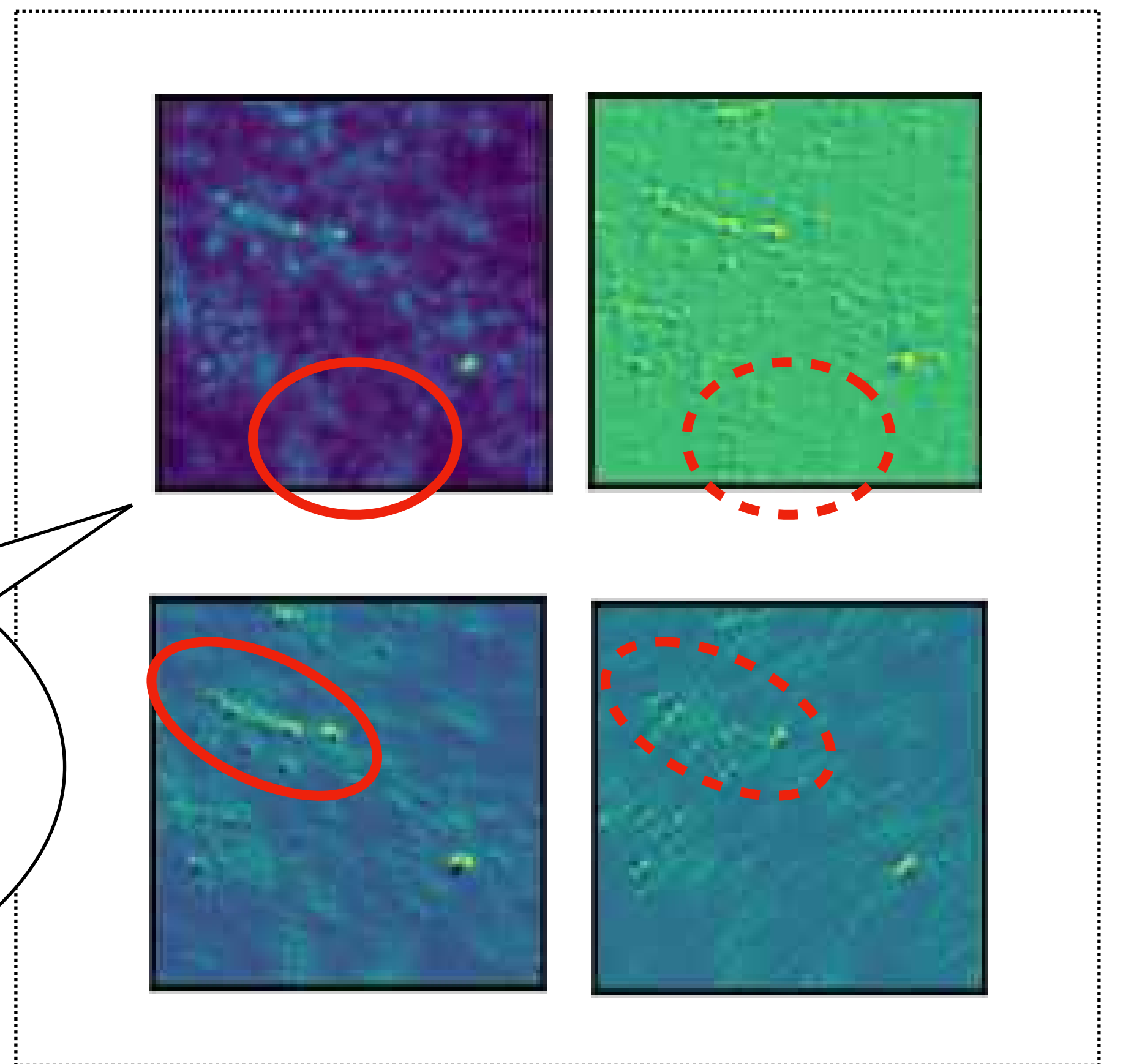
1st layer outputs



Filters in 2D Separation Models

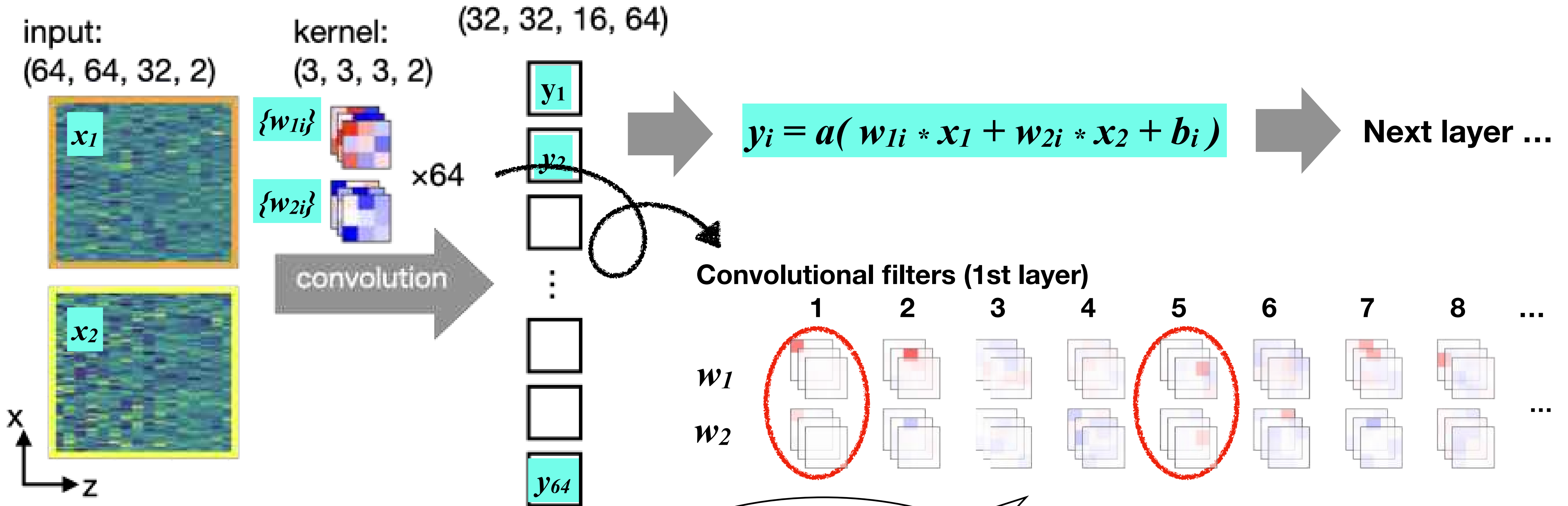


Examples of the filters in the 1st layer



**Structures at different distances have different features (e.g., scale length, bias).
→ The machine might distinguish signals from different distances by learning them.**

Filters in 3D Network with Pre-processing



Some filters pick up synchronizing signals in two inputs

→ Our method with physical information is working effectively as expected!

Challenges in Application for Actual Observational Data

Can we trust the reconstructed maps?

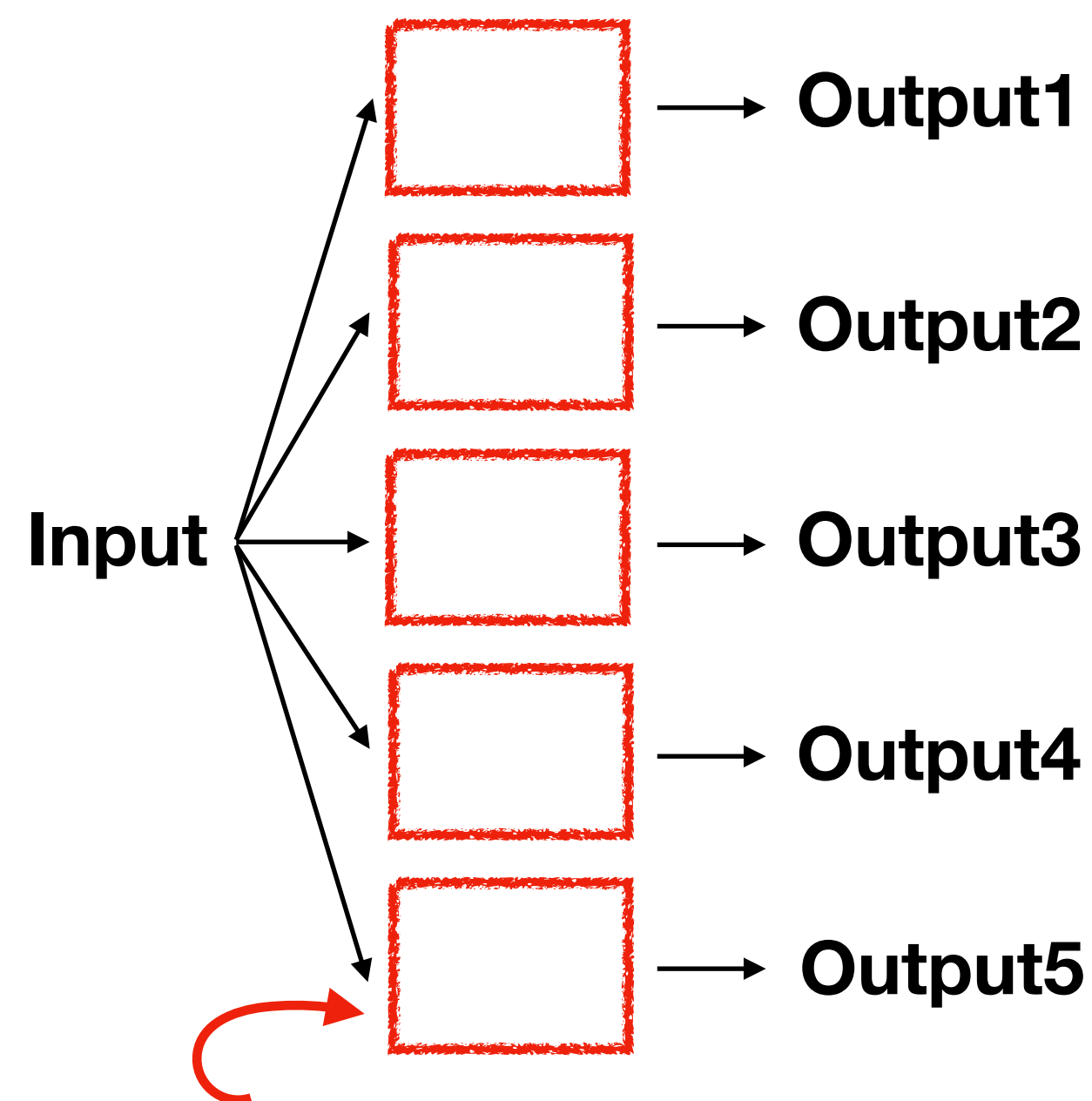
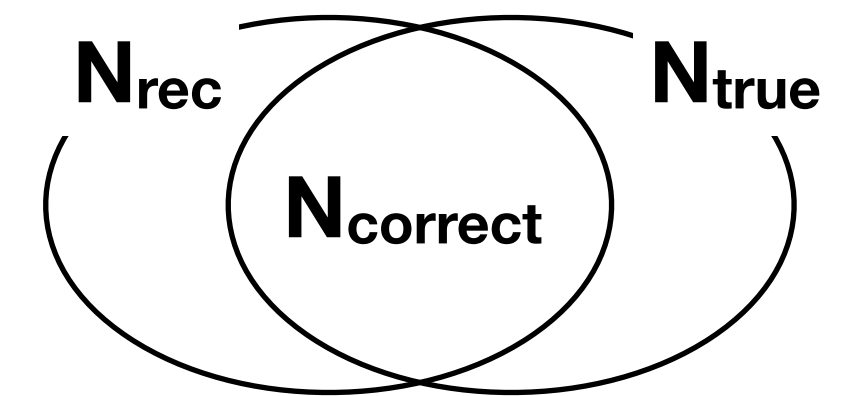
- **How precise is the reconstructed map? Is there a generation error?**
- **Is the model dependent on the assumption in the training model?**

→ Evaluation of the generation error and the effects of the assumed model is important to extract cosmological information from future observational data.

How Precise Is the Reconstructed Map?

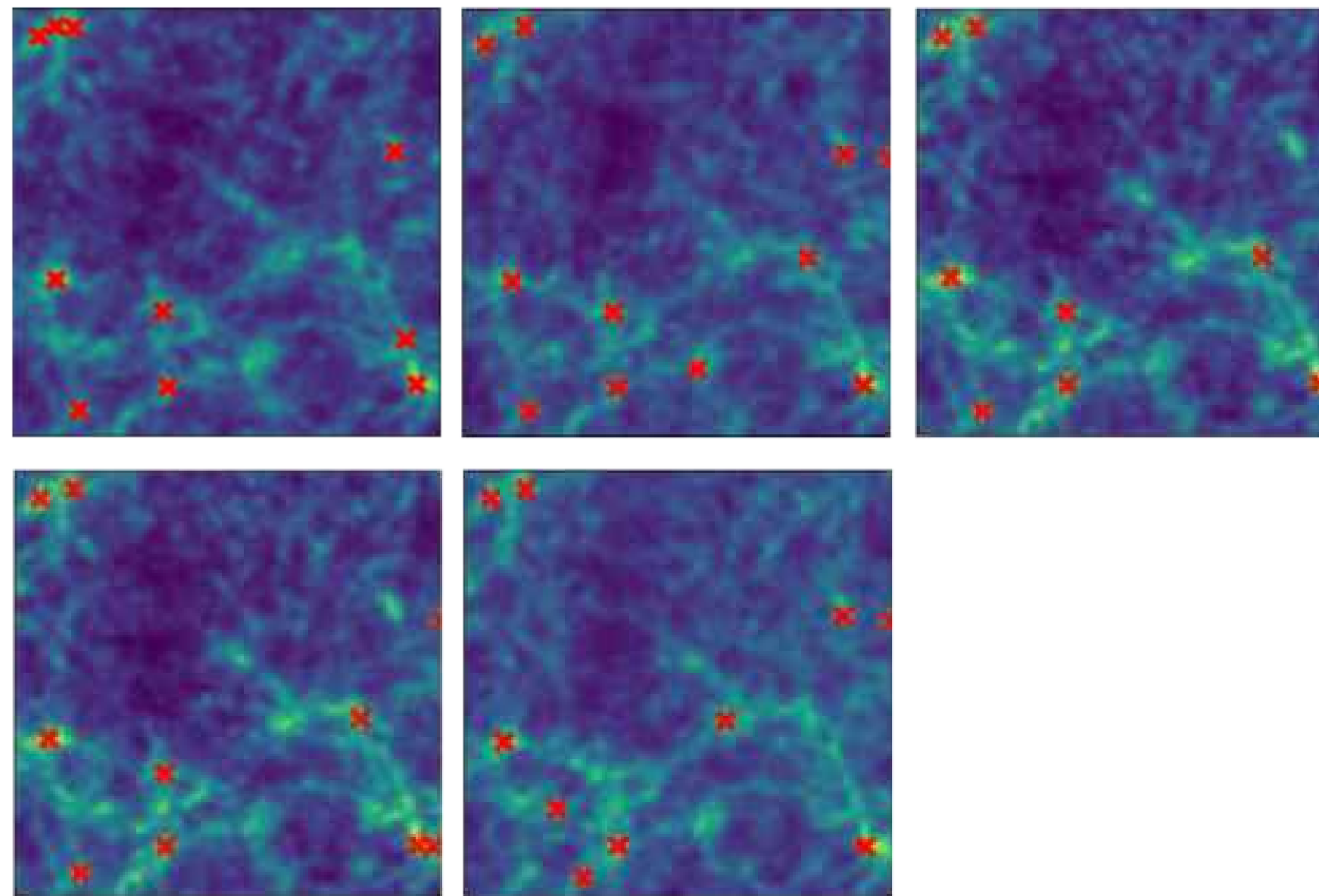
Detectability of $> 3\sigma$ peaks

- Precision ($N_{\text{correct}}/N_{\text{rec}}$) of a machine: **76%**
- Precision when we *combine* five networks (bugging): **91%**

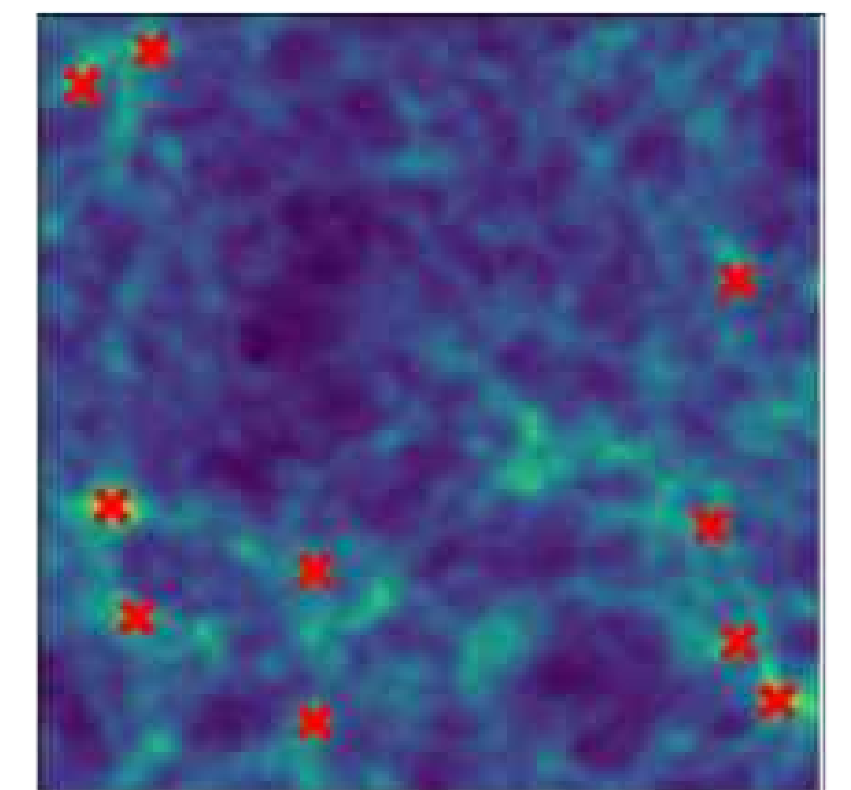


**Machines trained on the different data
and with different random feeds**

outputs of 5 networks



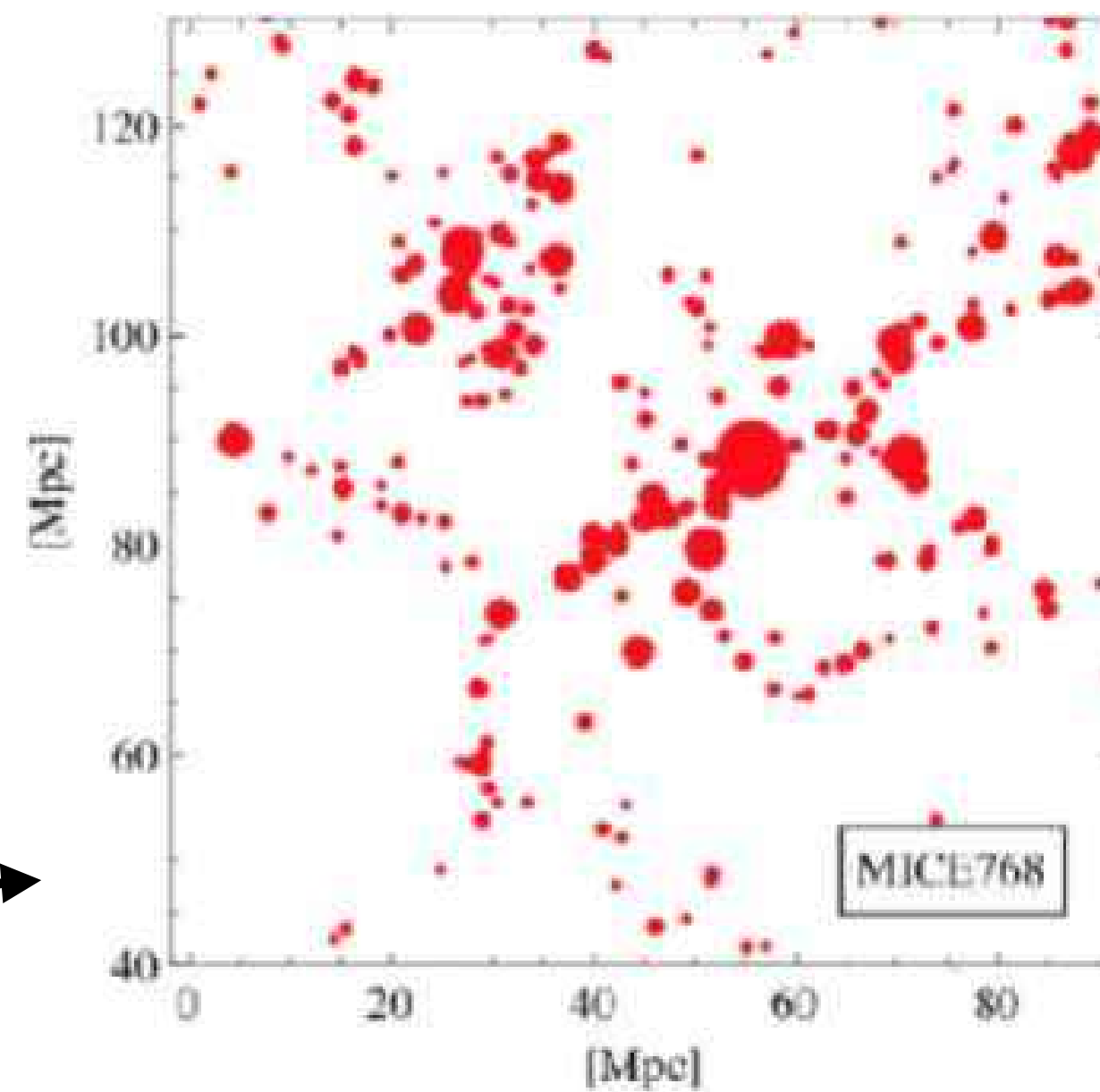
true distributions



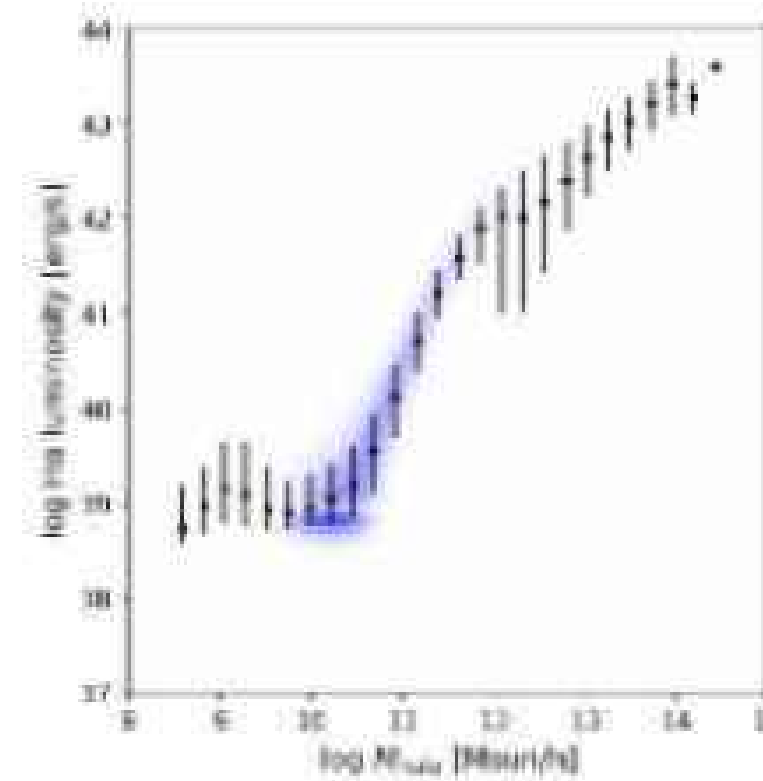
It is also possible to evaluate the generation error by combining multiple networks (e.g., taking the variance).

Does the Reconstruction Depend on the Assumed Model?

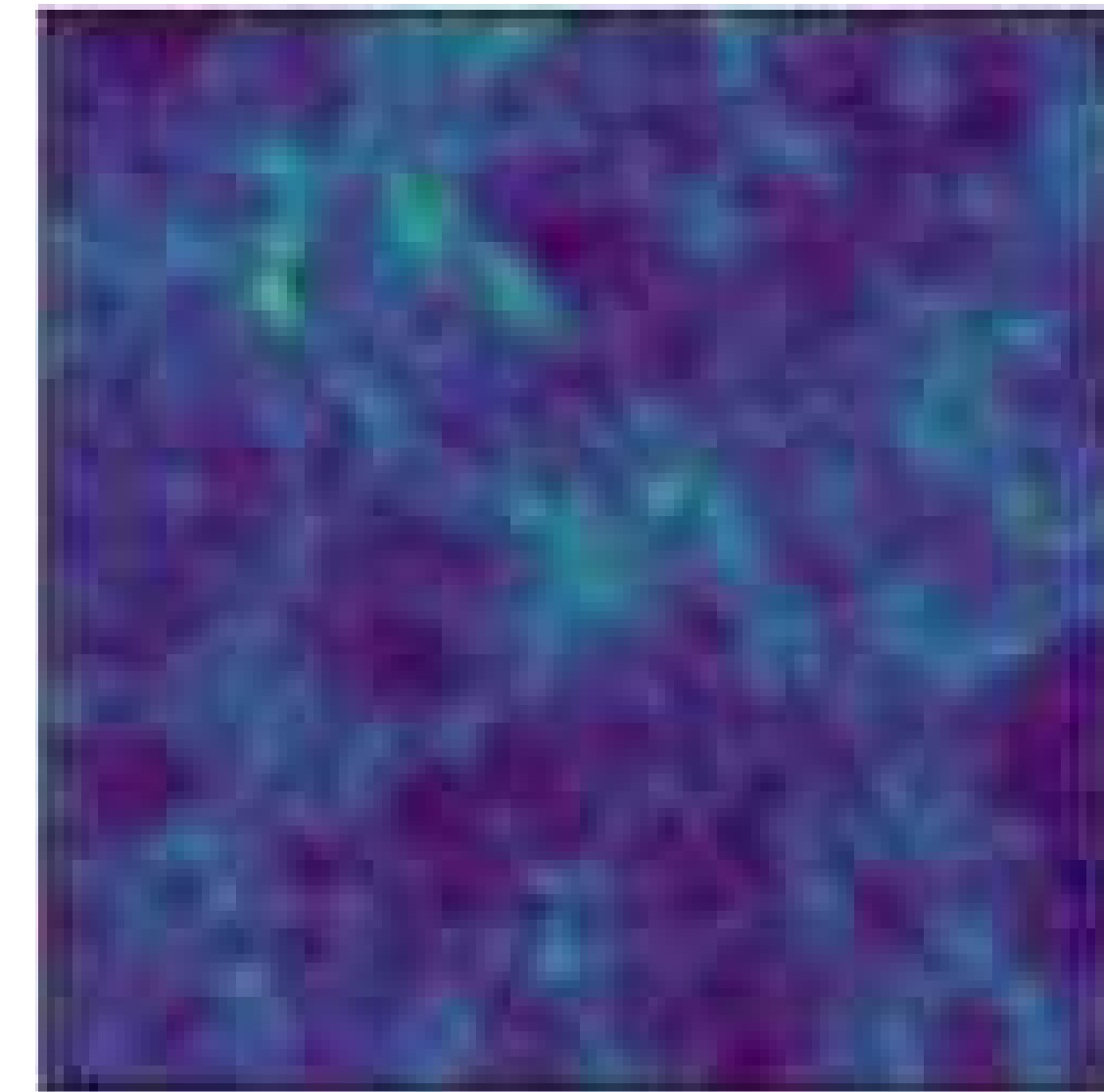
DM halo distribution



Emission line model:
(what we assume)



Mock data

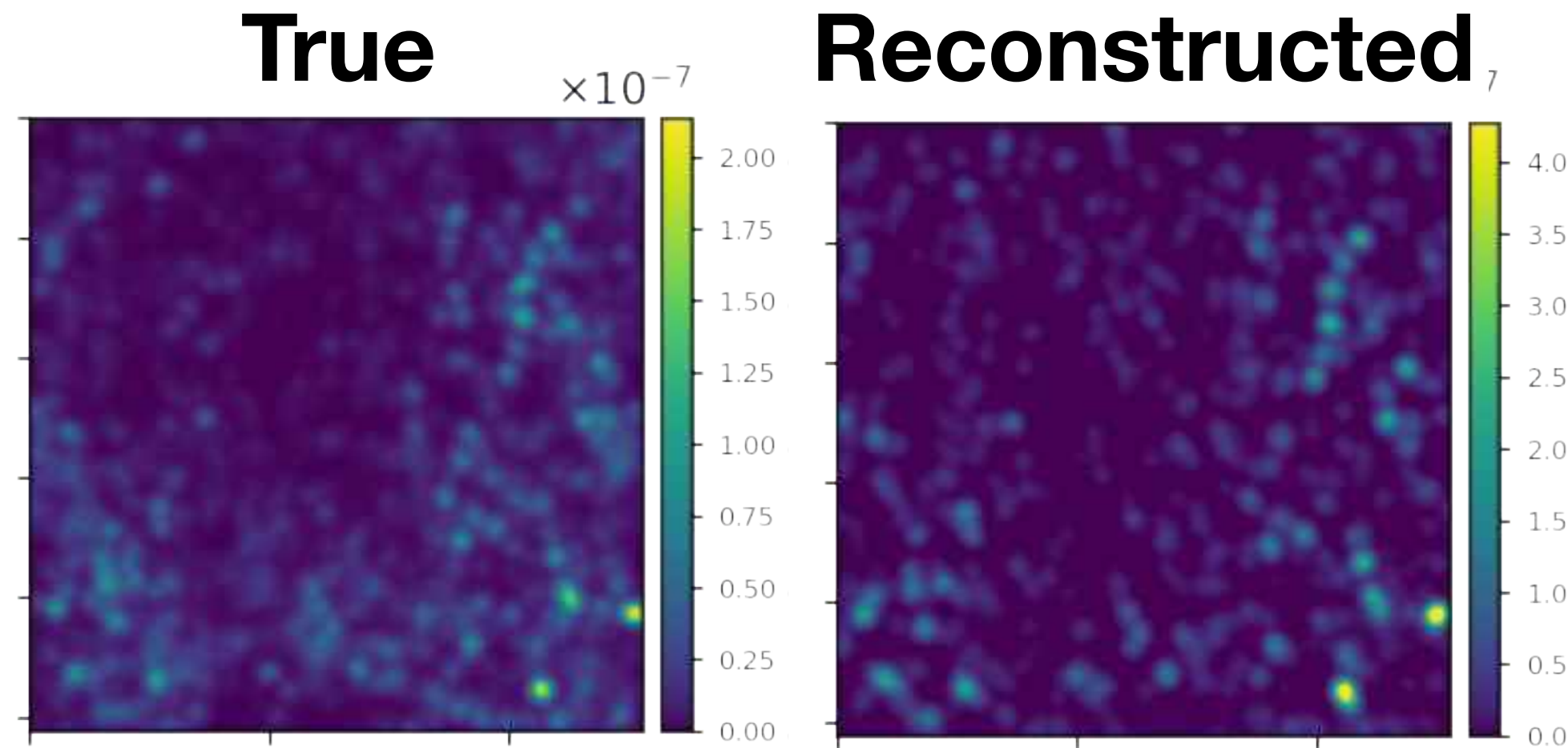


Cosmological information
(what we want to know from the observation)

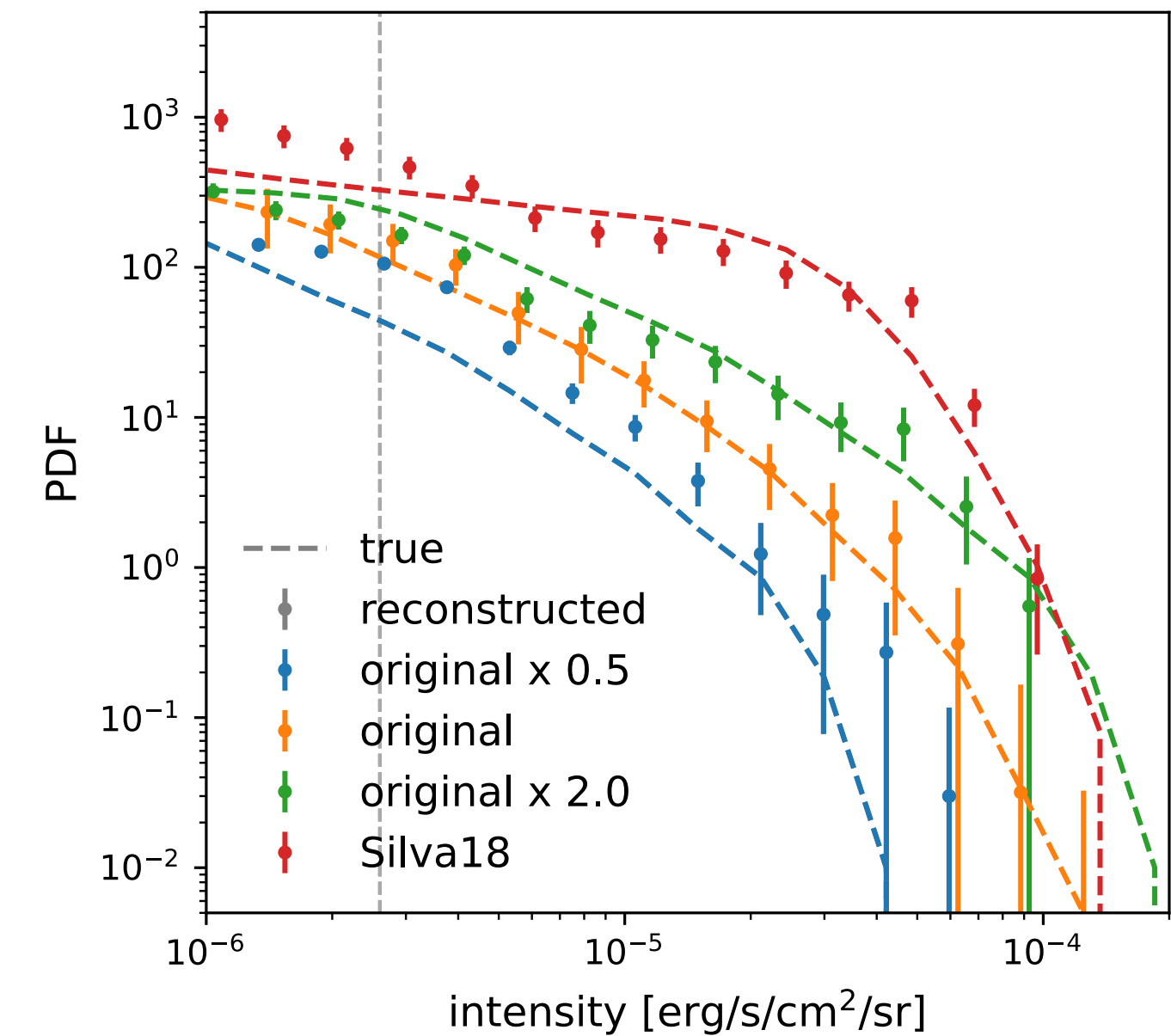
**What if the assumed line emission model
in training data is wrong?**

Test with Different Line Emission Models

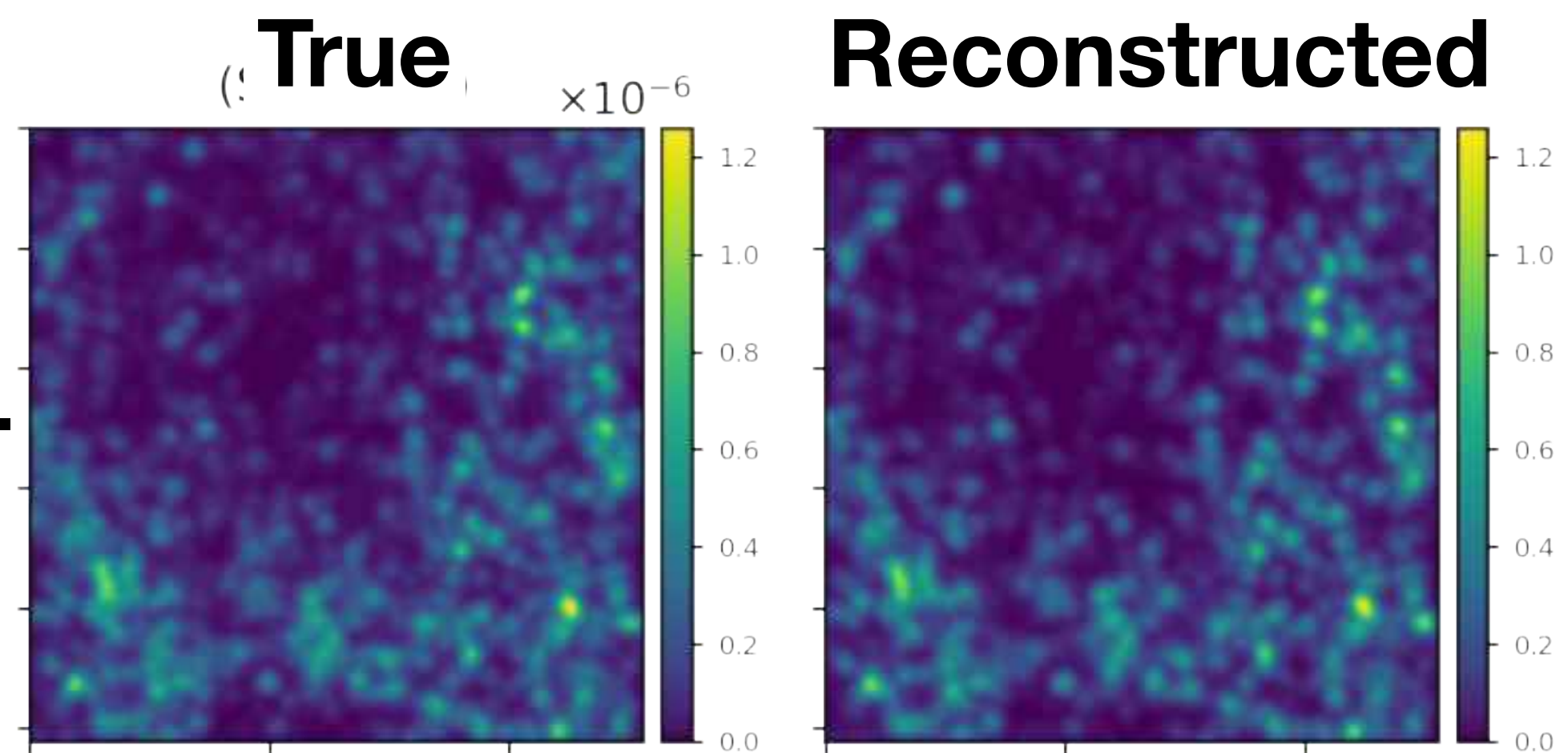
Model 1
($\times 2$ brighter
intensity model)



**Statistics as well as bright pixel
positions are reproduced
properly irrespective to the
assumed models in test data**



Model 2
(different mass-to-
luminosity model)



**What about noise model? More
different models? \rightarrow Future study**

Summary

- **A generative adversarial network can be used to reconstruct the large-scale distributions of the universe from noisy observational maps.**
- **We can get good reproducibility by pre-processing the input data based on physical information.**
- **The machine learns the typical features in the large-scale structure as well as the synchronizing signals in two input data.**
- **Generation errors and the uncertainties in assumed models should be carefully evaluated in future actual use – combining multiple machines would be an important strategy.**