

Hydrogen Intensity Mapping: **the ultimate signal is the weakest of all**

Isabella Paola Carucci

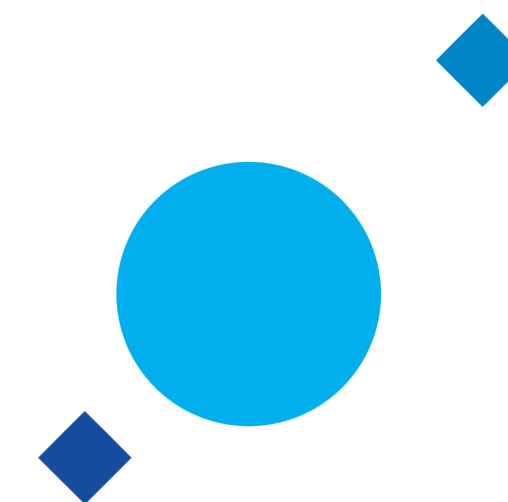
INAF - OATs

Cosmology 2023 in Miramare



Funded by
the European Union

NextGenerationEU



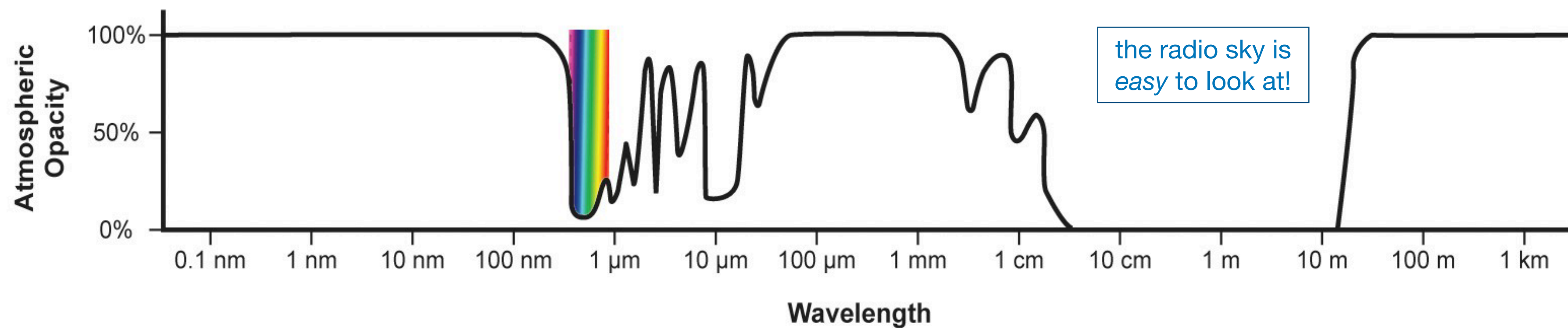
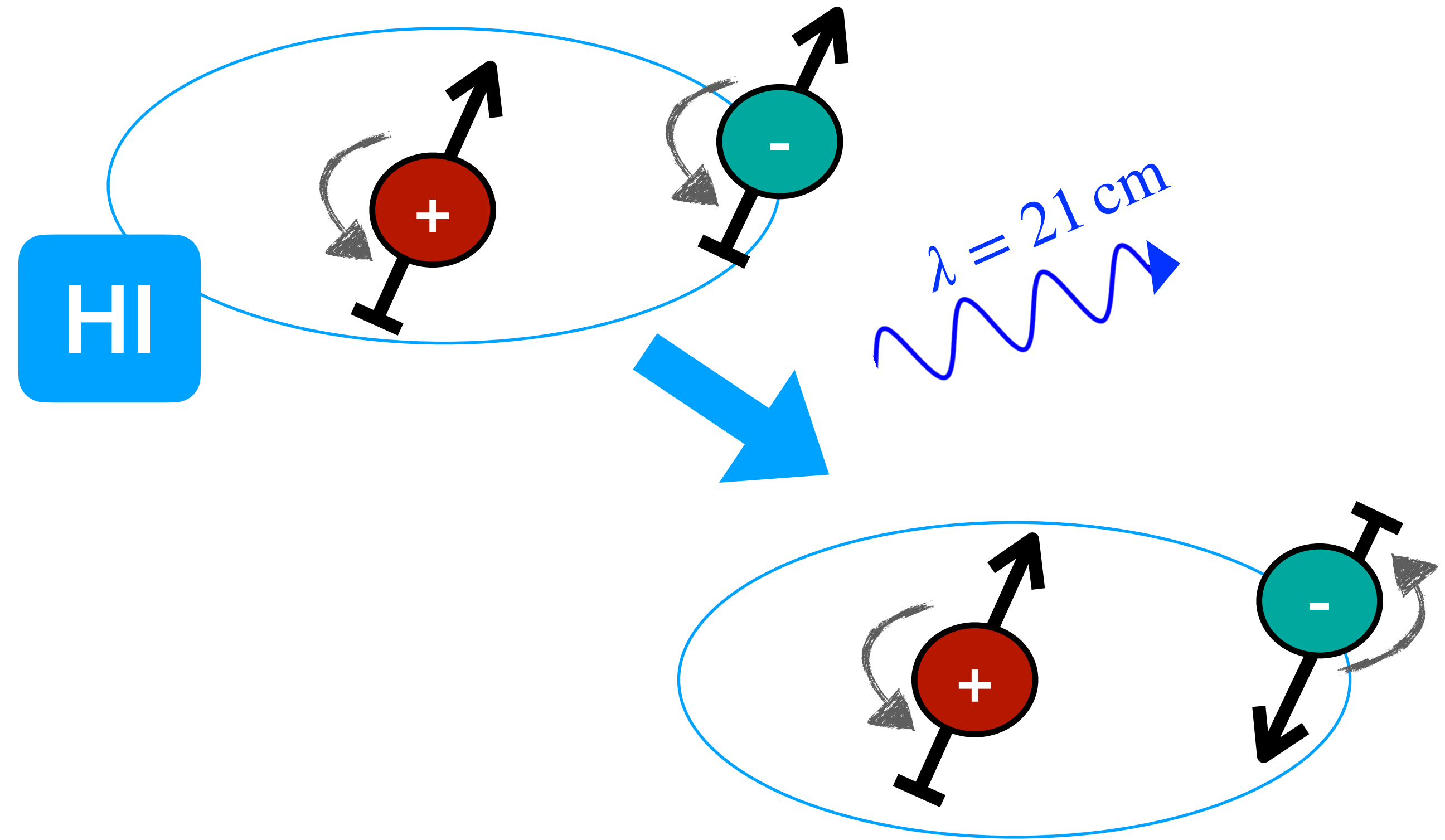
INAF

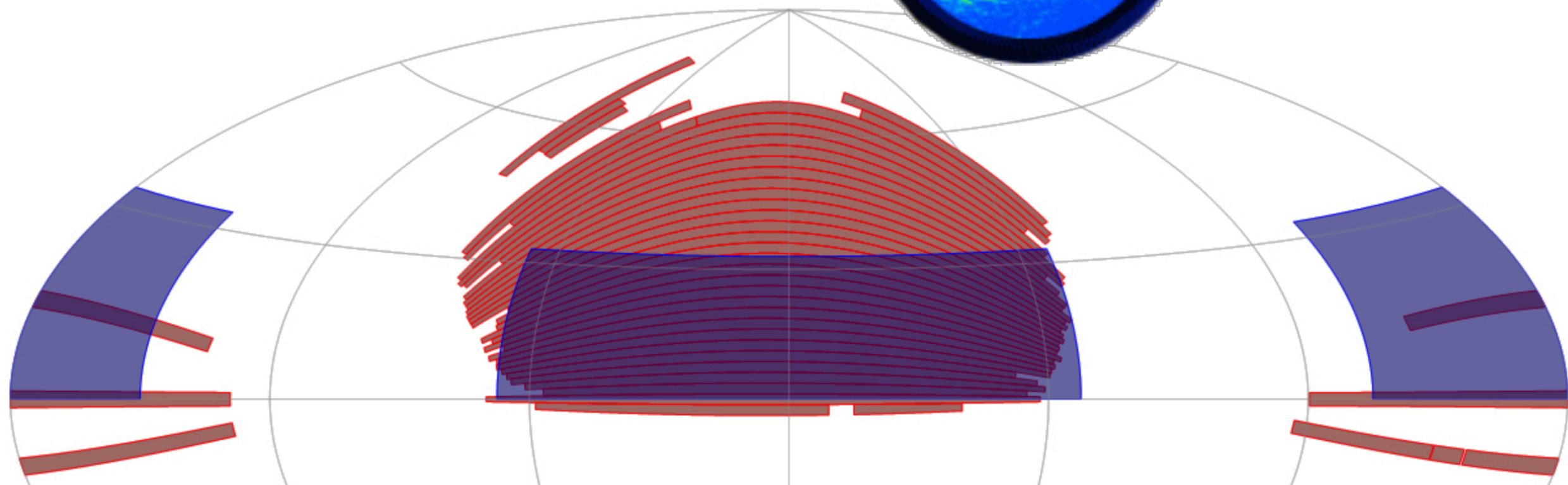
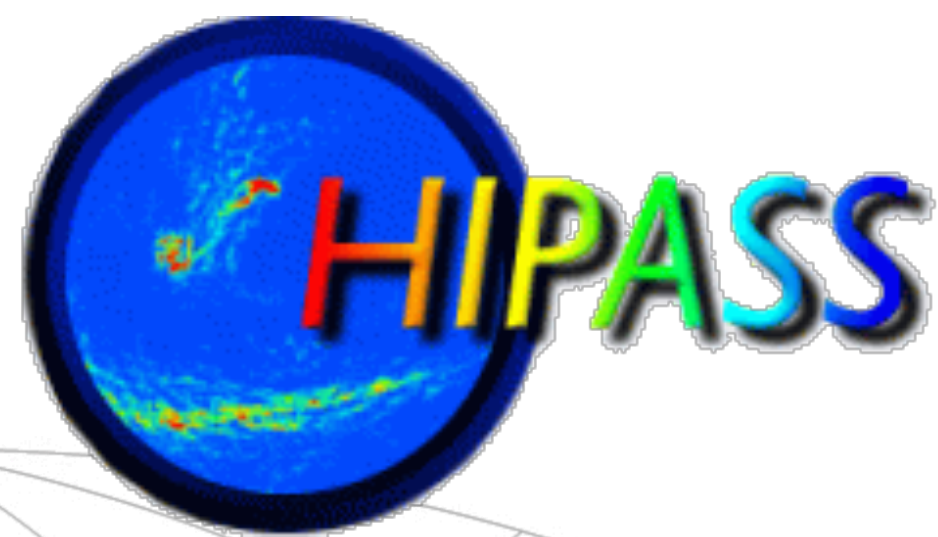
ISTITUTO NAZIONALE
DI ASTROFISICA

1. **Hydrogen Intensity Mapping (IM):**
what is it and why to do it
2. **IM is hard!** Biggest challenge:
weakness of the IM signal compared
to contaminants
3. We are getting there. **MeerKLASS**



21-cm radiation

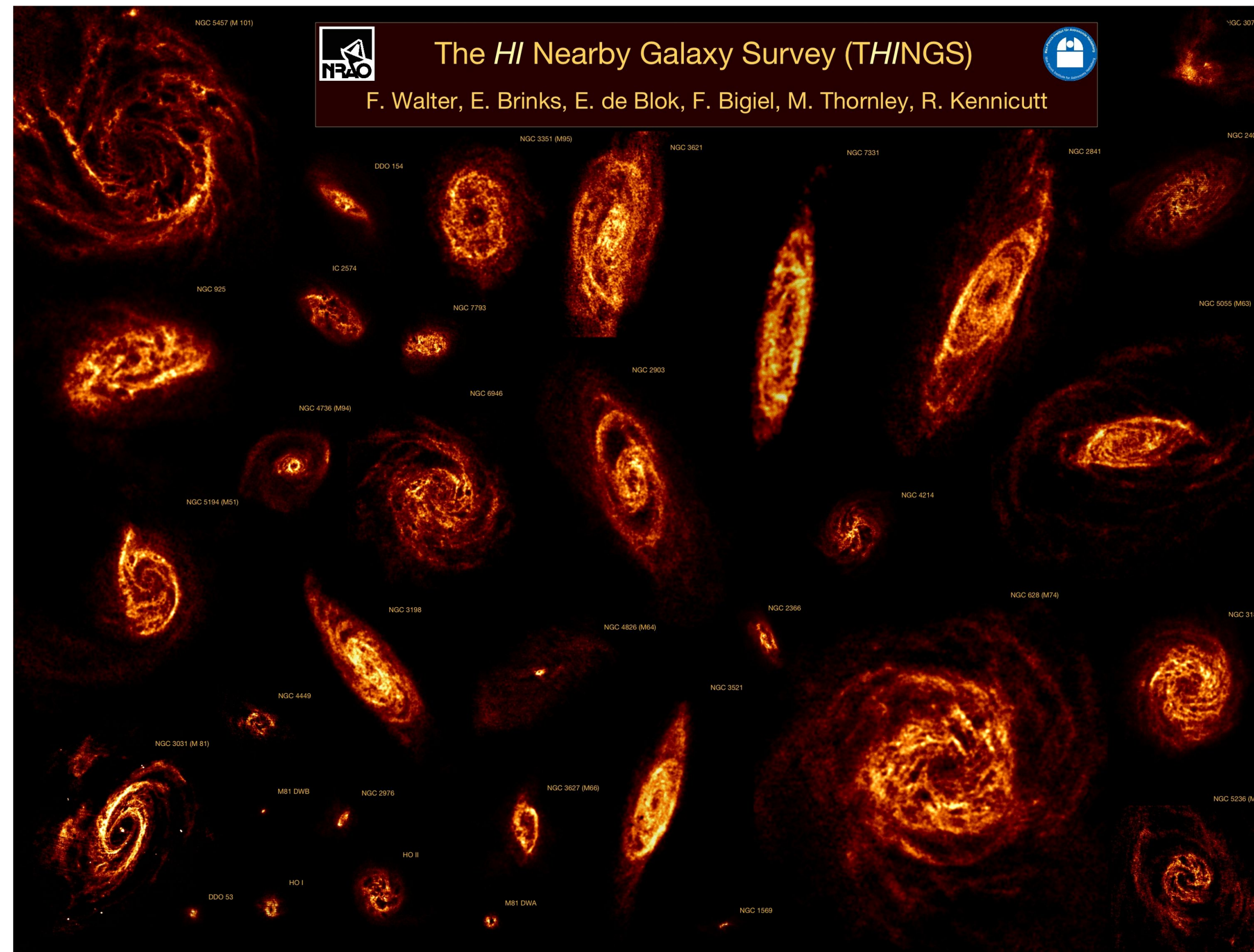
- strongly forbidden:
 $t_{1/2} \sim 10^7$ years
- VERY abundant
- Spectrally isolated
- Small obscuration



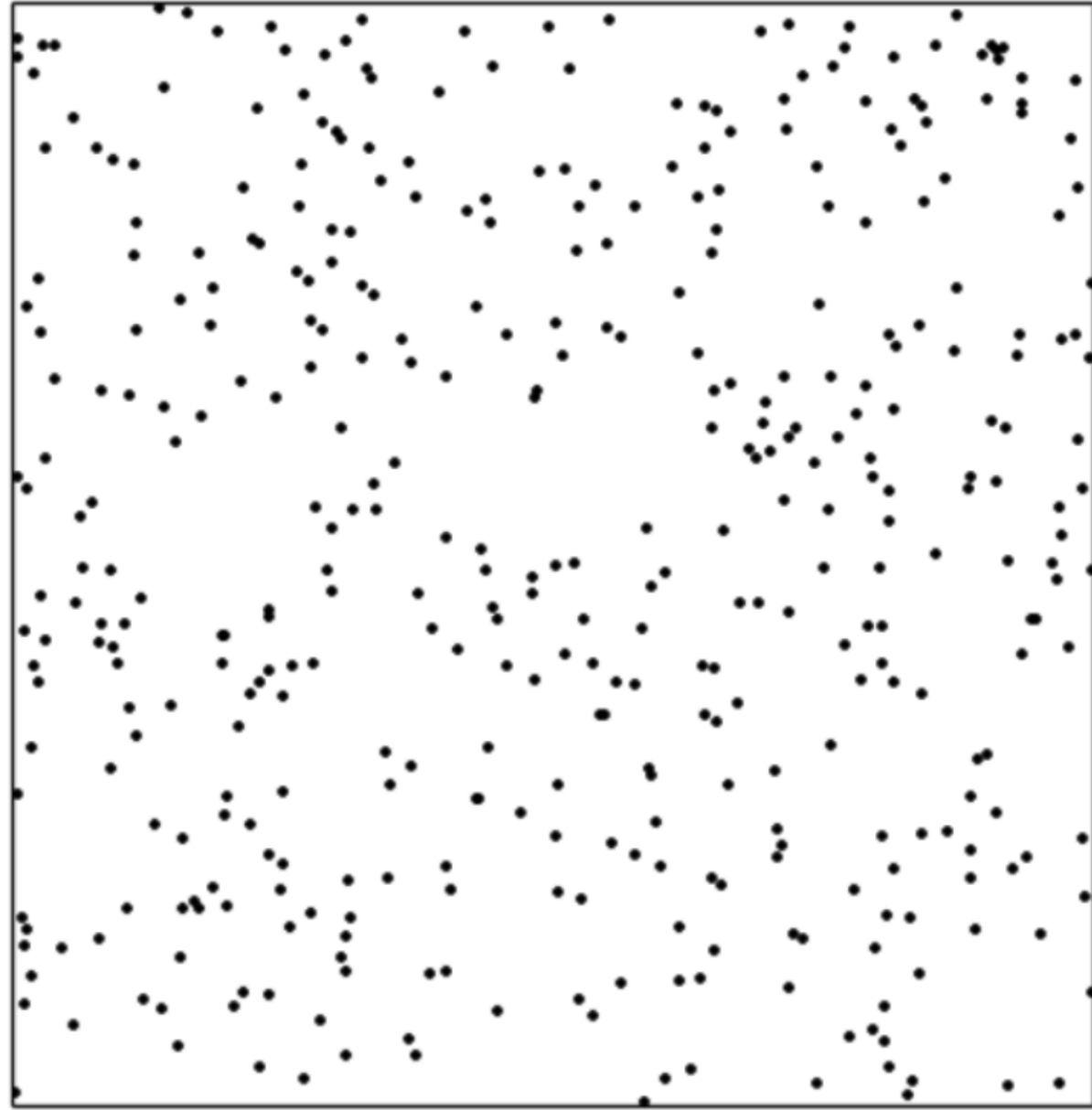


The Arecibo Legacy Fast ALFA Survey

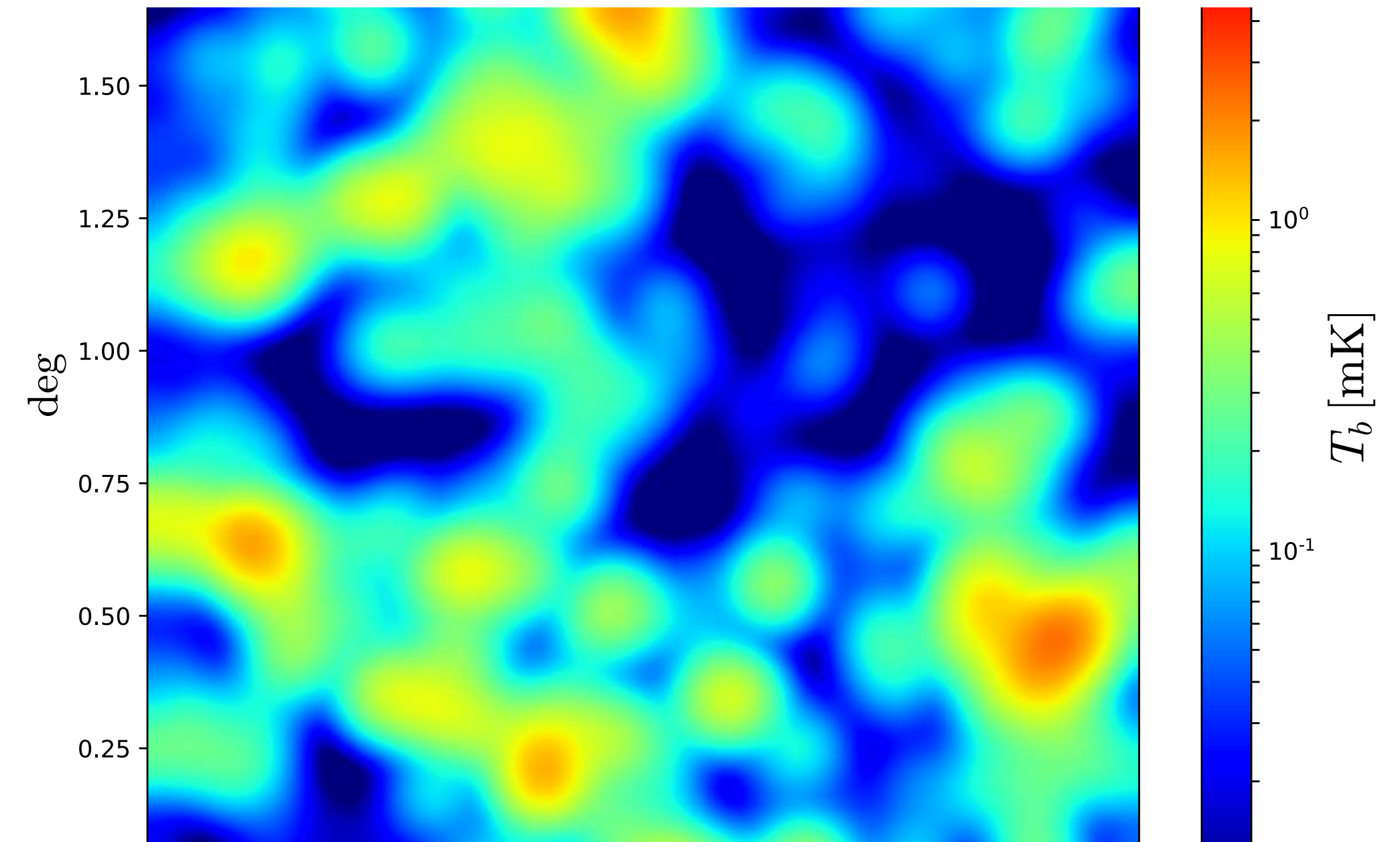
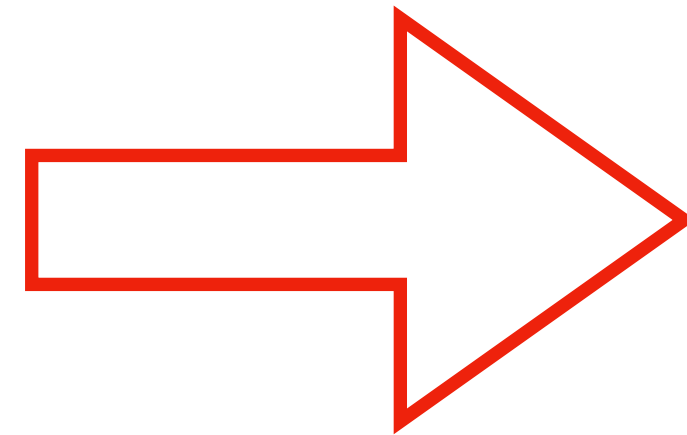
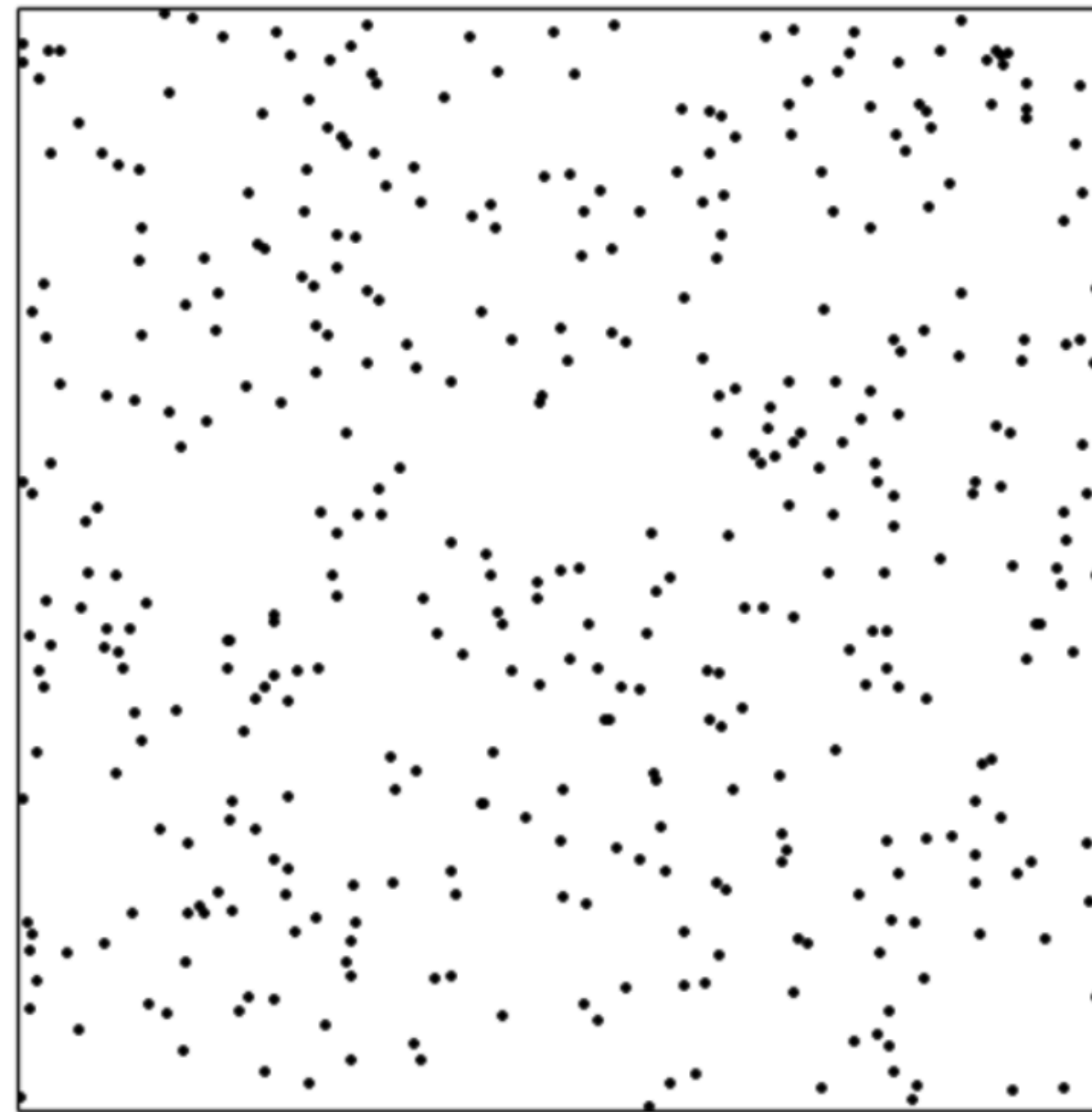
 The *HI* Nearby Galaxy Survey (*THINGS*) 
F. Walter, E. Brinks, E. de Blok, F. Bigiel, M. Thornley, R. Kennicutt



Record $z = 0.376$ detection of 21 cm emitting galaxy with 178 hours from VLA [Fernández et al, 2016]

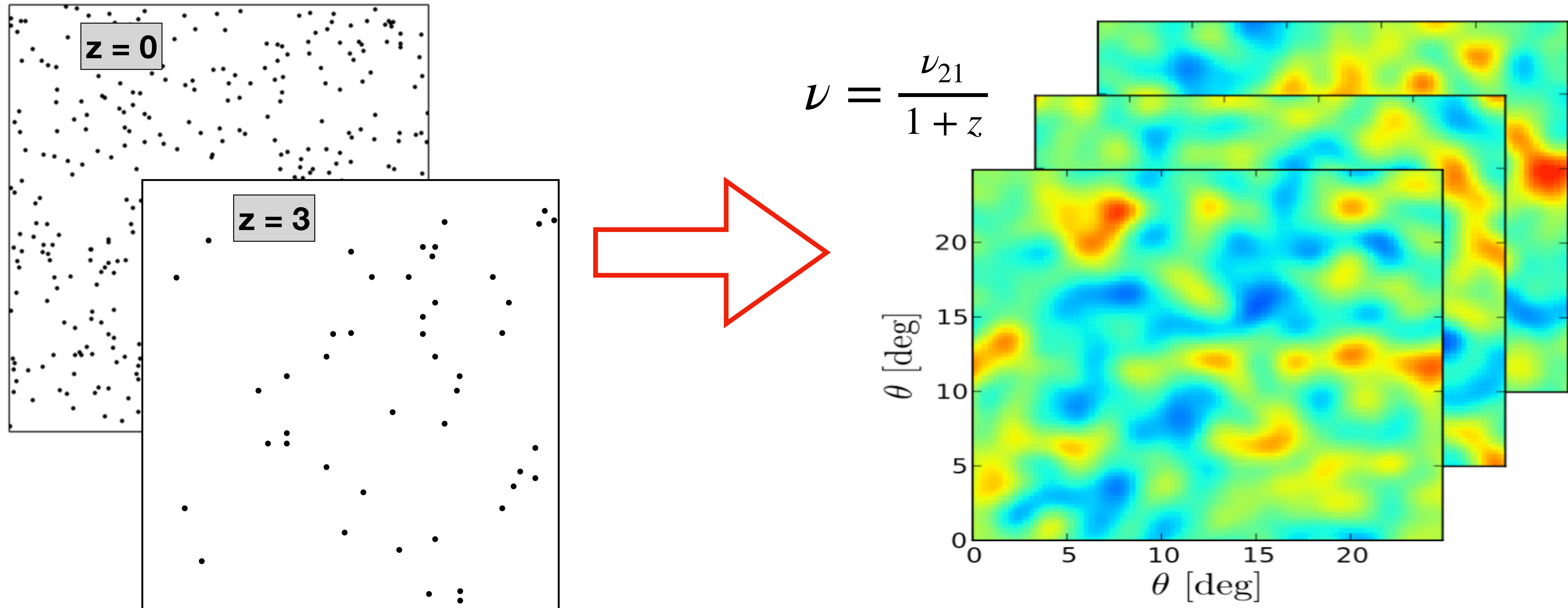


21-cm intensity mapping






Put signal-to-noise where you really need it: **linear large scale modes**

21-cm intensity mapping

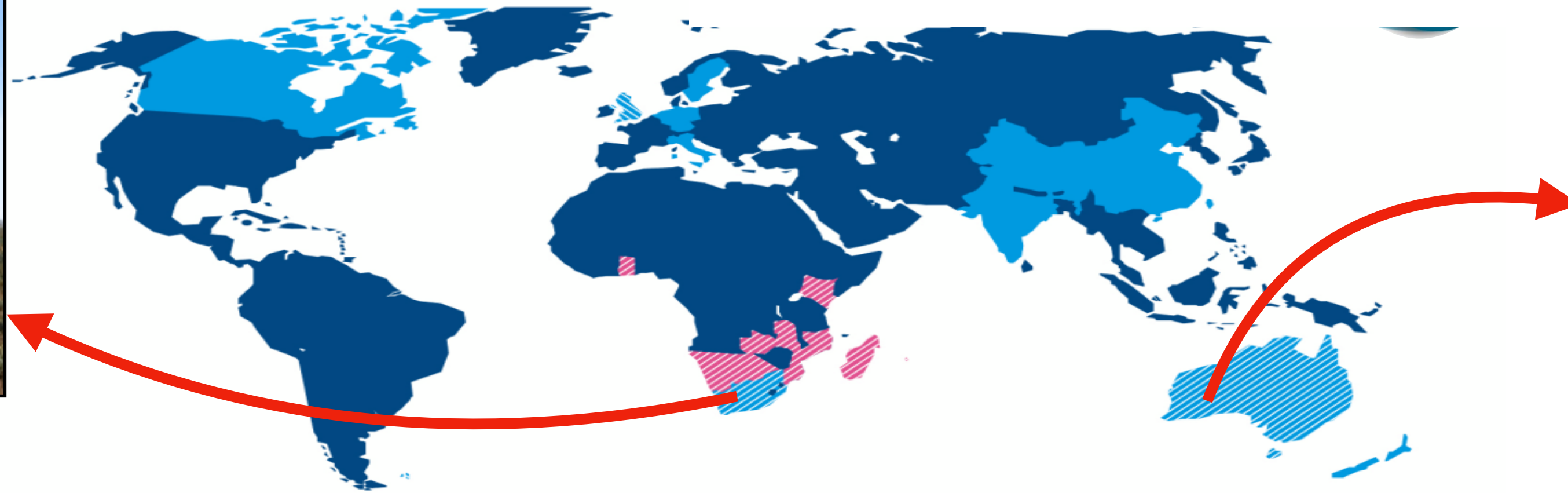


Big volumes (for cheap) and high redshift resolution

A large-scale structure scientist's wish list:

-  1. Large areas
-  2. Deep and accurate redshifts (distances)
-  3. Better coverage of the Universe epochs

SKAO



- Full members
- SKA Headquarters host country
- SKA Phase 1 and Phase 2 host countries

- African partner countries (non-member SKA Phase 2 host countries)

This map is intended for reference only and is not meant to represent legal borders

SKA1-mid

the SKA's mid-frequency instrument

$0 < z < 3$



Location:
South Africa



Frequency range:
350 MHz
to
15.3 GHz
with a goal of 24 GHz



197 dishes
(including 64 MeerKAT dishes)



Maximum baseline:
150km

SKA1-low

the SKA's low-frequency instrument

$3 < z < 27$



Location: Australia



Frequency range:
50 MHz
to
350 MHz



~131,000
antennas spread between
512 stations



Maximum baseline:
~65km

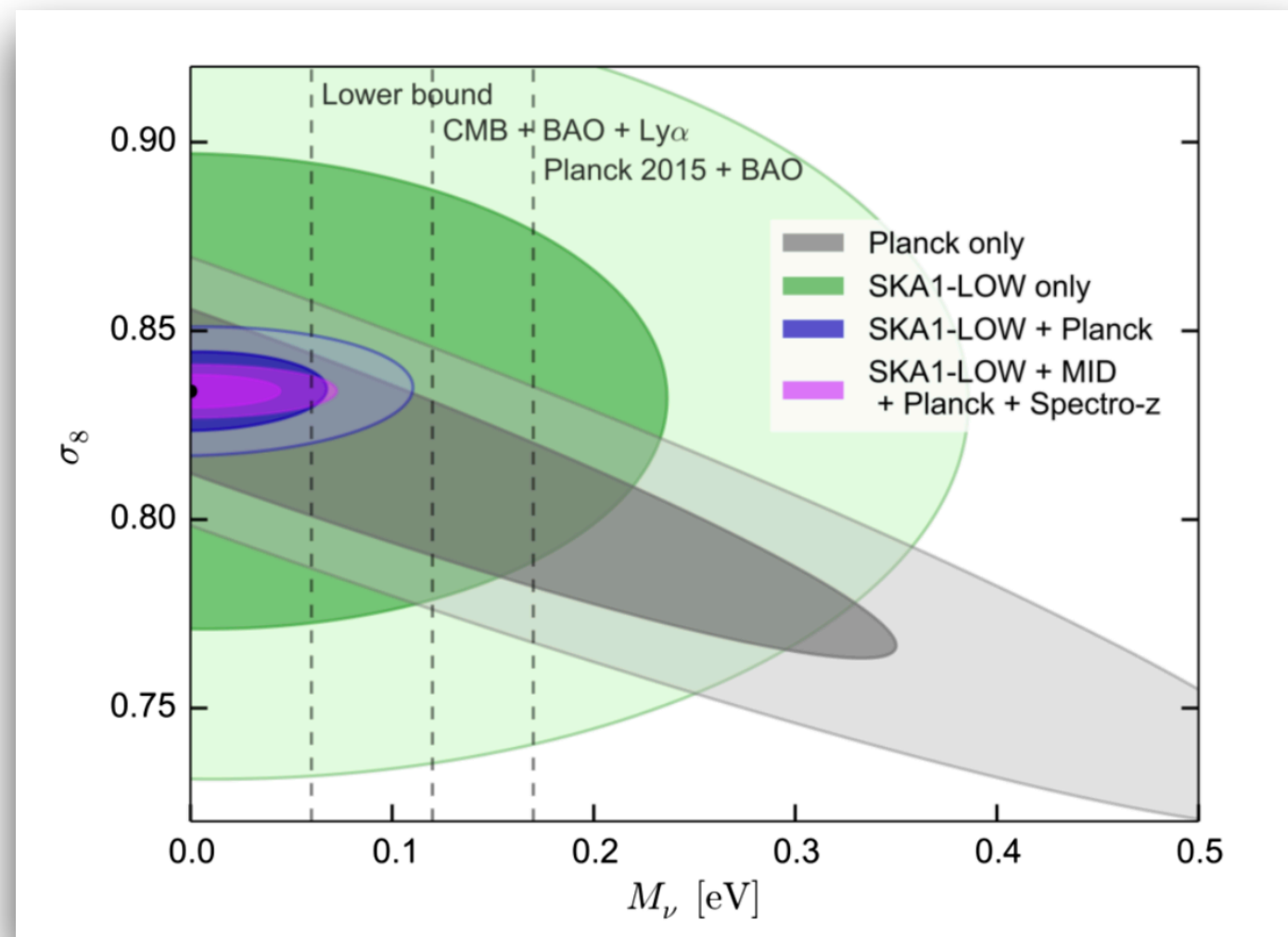
HI intensity mapping with the SKAO

Proposed SKA1 Cosmology Surveys

- a) Medium-Deep Survey of 5,000 deg² at 0.95-1.4 GHz for
 - HI galaxy redshift survey with 3.5 million objects
 - Weak Lensing shape measurements with ~50 million objects
 - Continuum galaxy survey with ~60 million objects
- b) Wide Survey of 20,000 deg² at 0.35-1.05 GHz for
 - Continuum galaxy survey with ~100 million objects
 - • HI intensity maps for $0.35 < z < 3$
- c) Deep Survey 100 deg² at 200-350 MHz for
 - • HI intensity maps for $3 < z < 6$

Cosmology with Phase 1 of the Square Kilometre Array **Red Book** 2018:
Technical specifications and performance forecasts

Intensity mapping with the SKAO



SKA1 Cosmology Surveys

• Deep Survey of 5,000 deg² at 0.95-1.4 GHz for galaxy redshift survey with 3.5 million objects

• Weak Lensing shape measurements with ~50 million objects

• Continuum galaxy survey with ~60 million objects

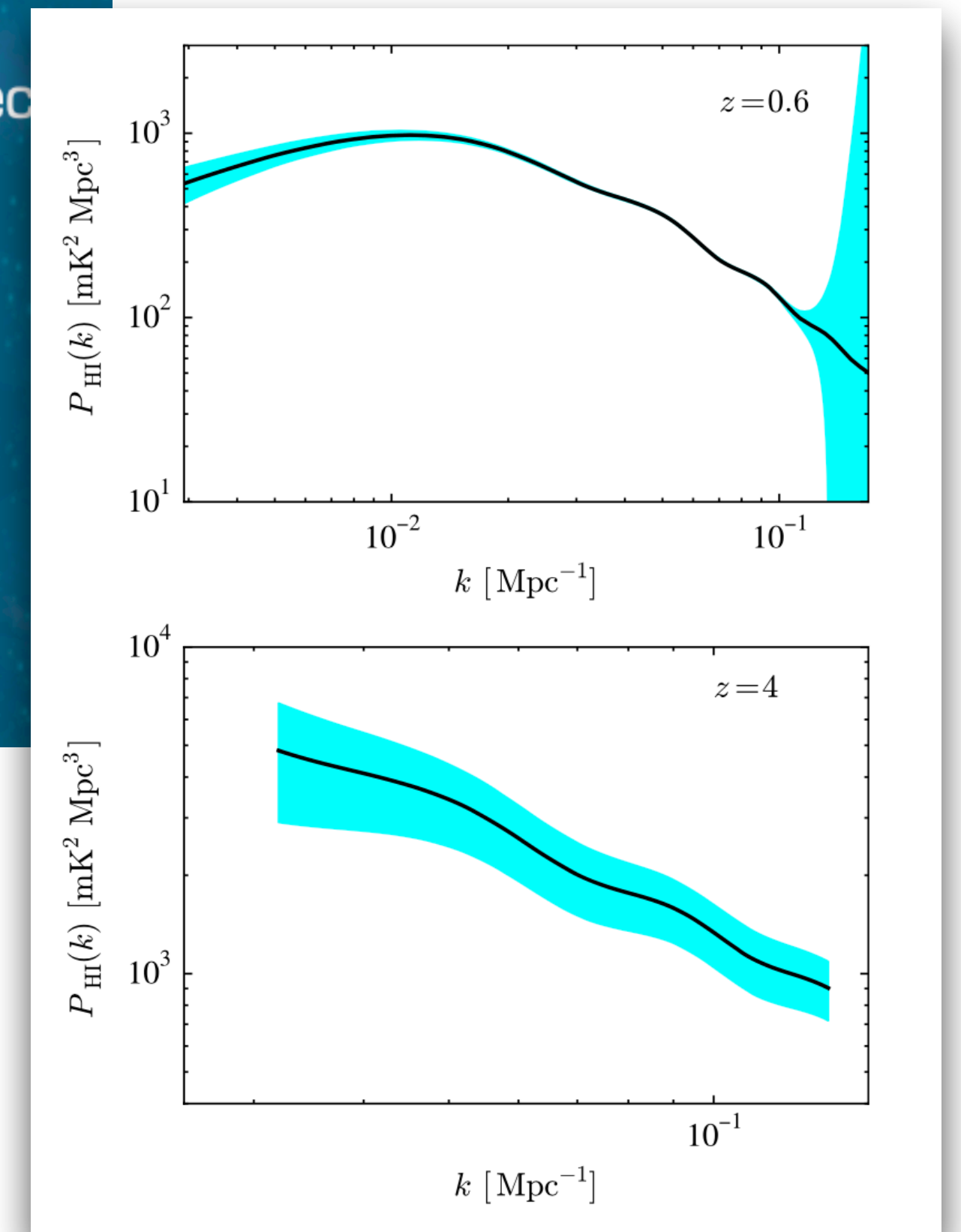
b) Wide Survey of 20,000 deg² at 0.35-1.05 GHz for

- Continuum galaxy survey with ~100 million objects
- HI intensity maps for 0.35 < z < 3



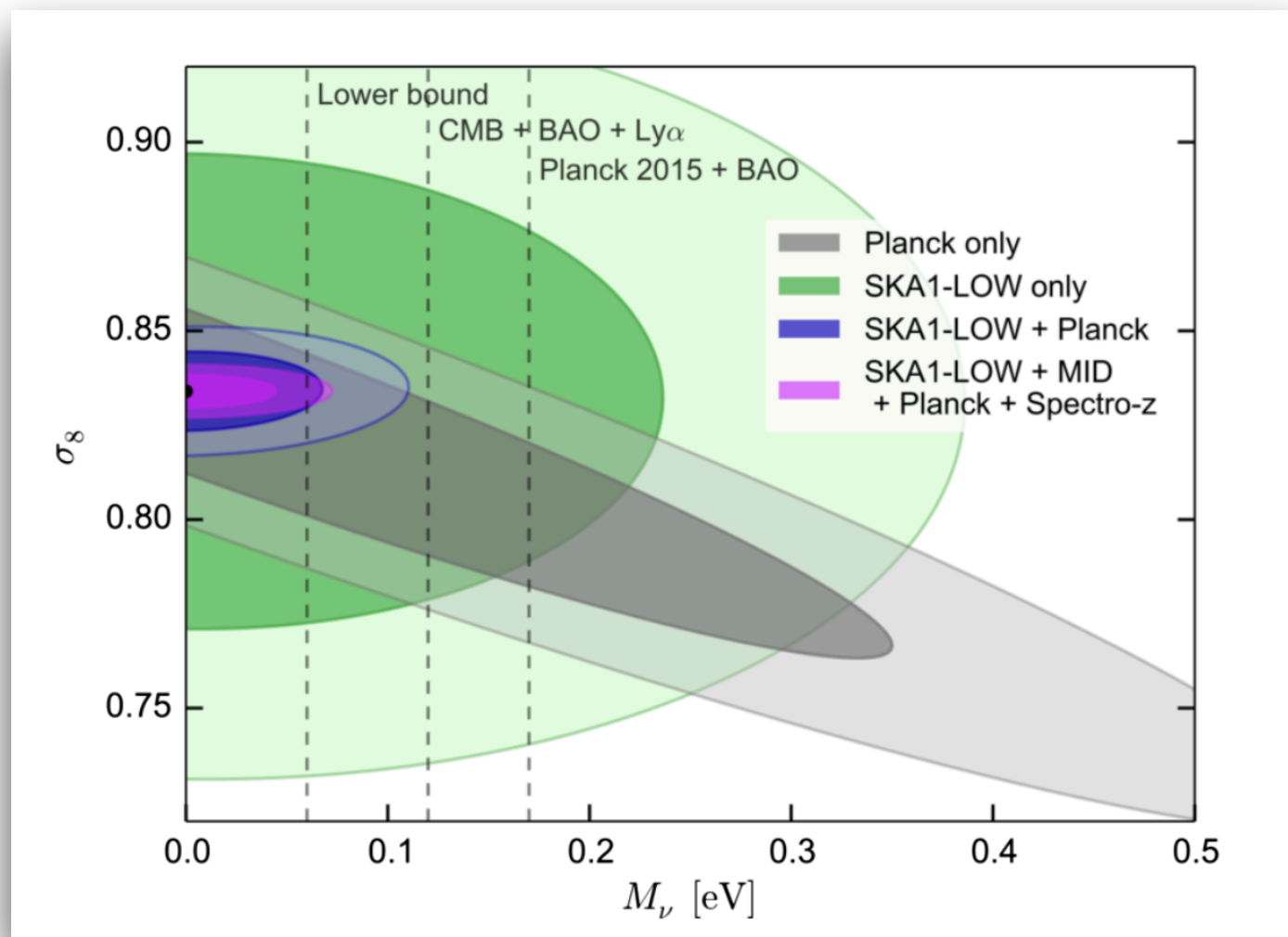
c) Deep Survey 100 deg² at 200-350 MHz for

- HI intensity maps for 3 < z < 6



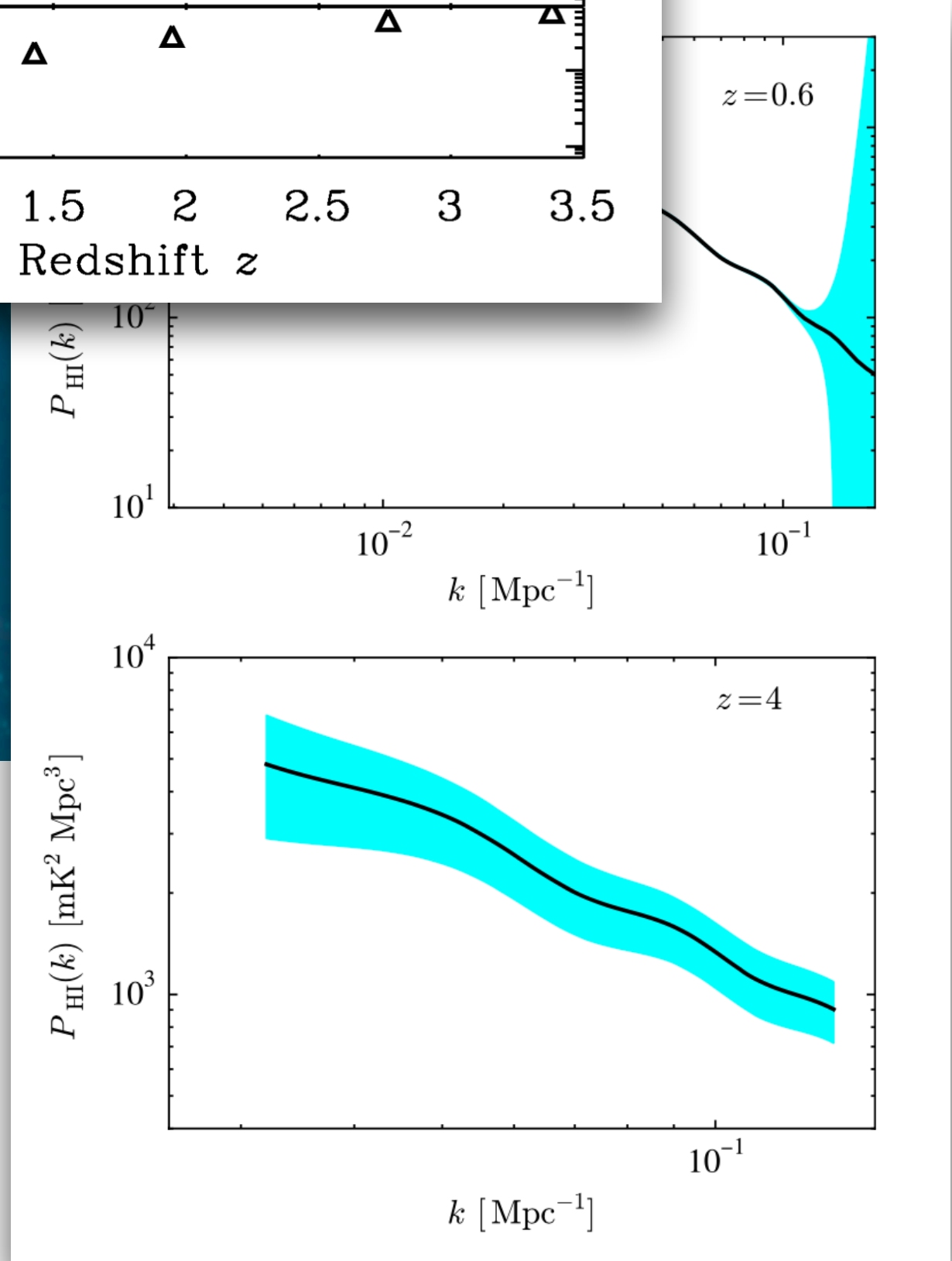
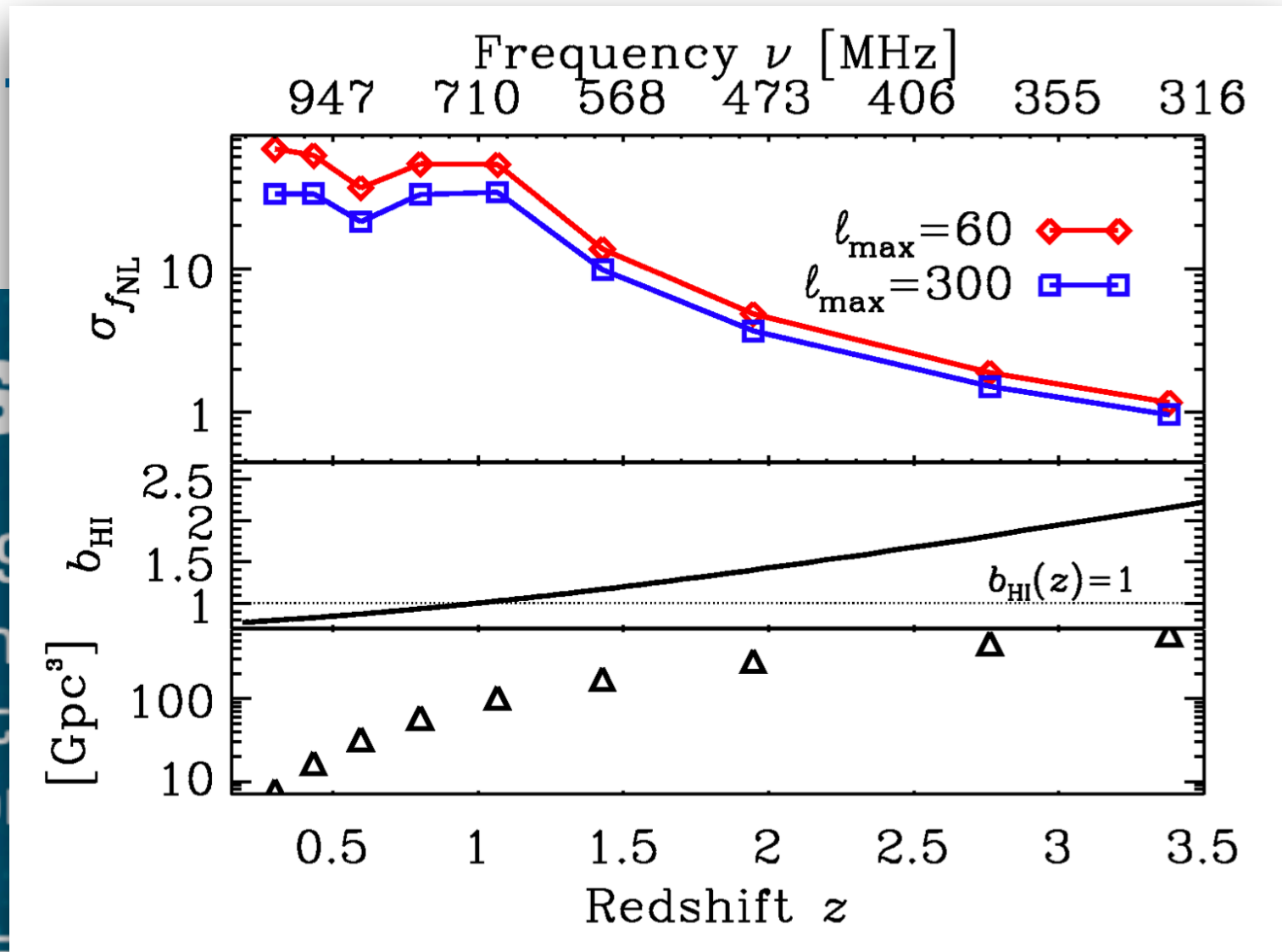
Cosmology with Phase 1 of the Square Kilometre Array **Red Book 2018:**
 Technical specifications and performance forecasts

Intensity mapping with



SKA1 Cosmology S

- Continuum galaxy survey with ~60 million objects
- Deep Survey of 5,000 deg² at 0.9-1.6 GHz for galaxy redshift survey with 3.5 million objects
- Weak Lensing shape measurements with 100 million objects
- Continuum galaxy survey with ~60 million objects
- HI intensity maps for 0.35 < z < 3
- HI intensity maps for 3 < z < 6



Cosmology with Phase 1 of the Square Kilometre Array **Red Book 2018:**
Technical specifications and performance forecasts

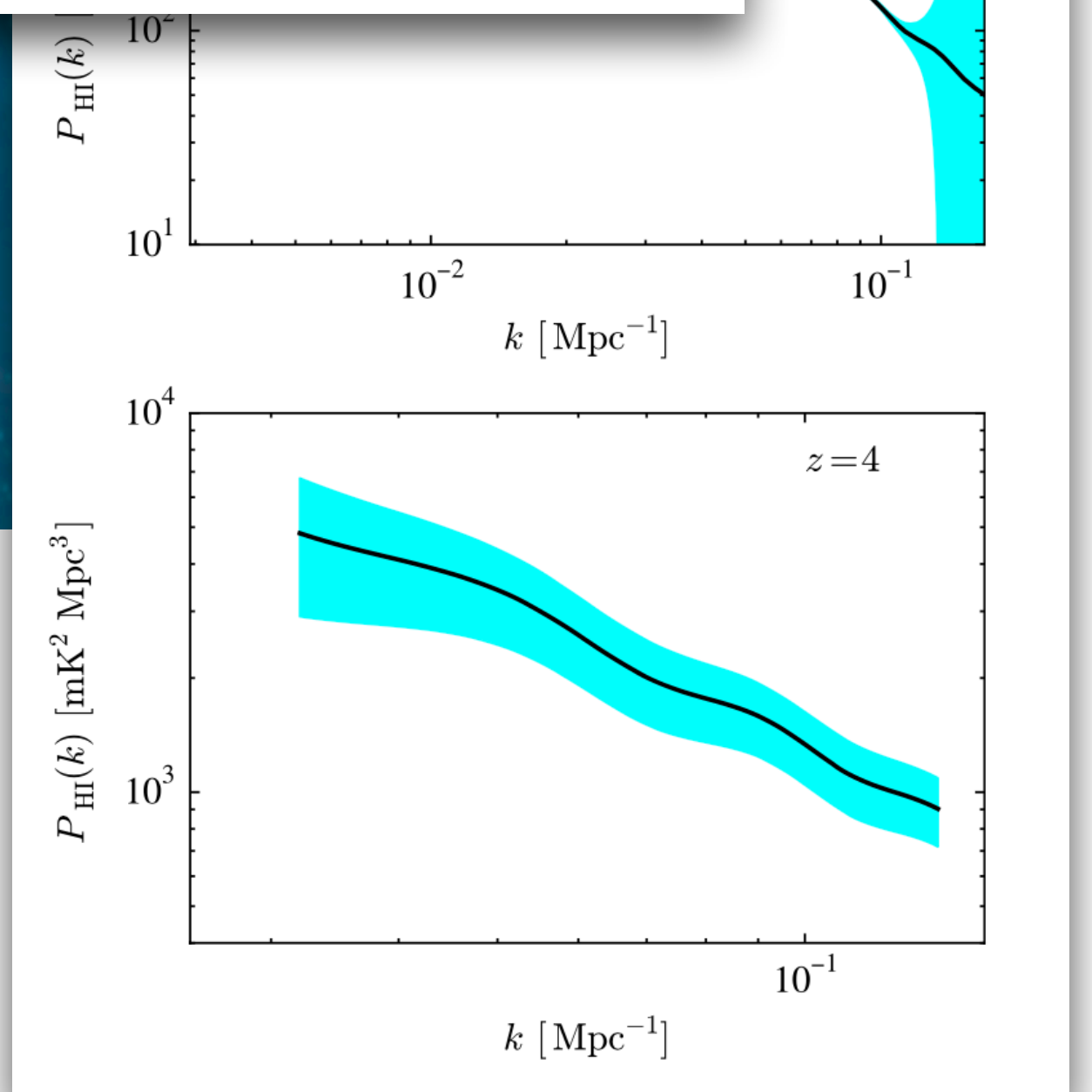
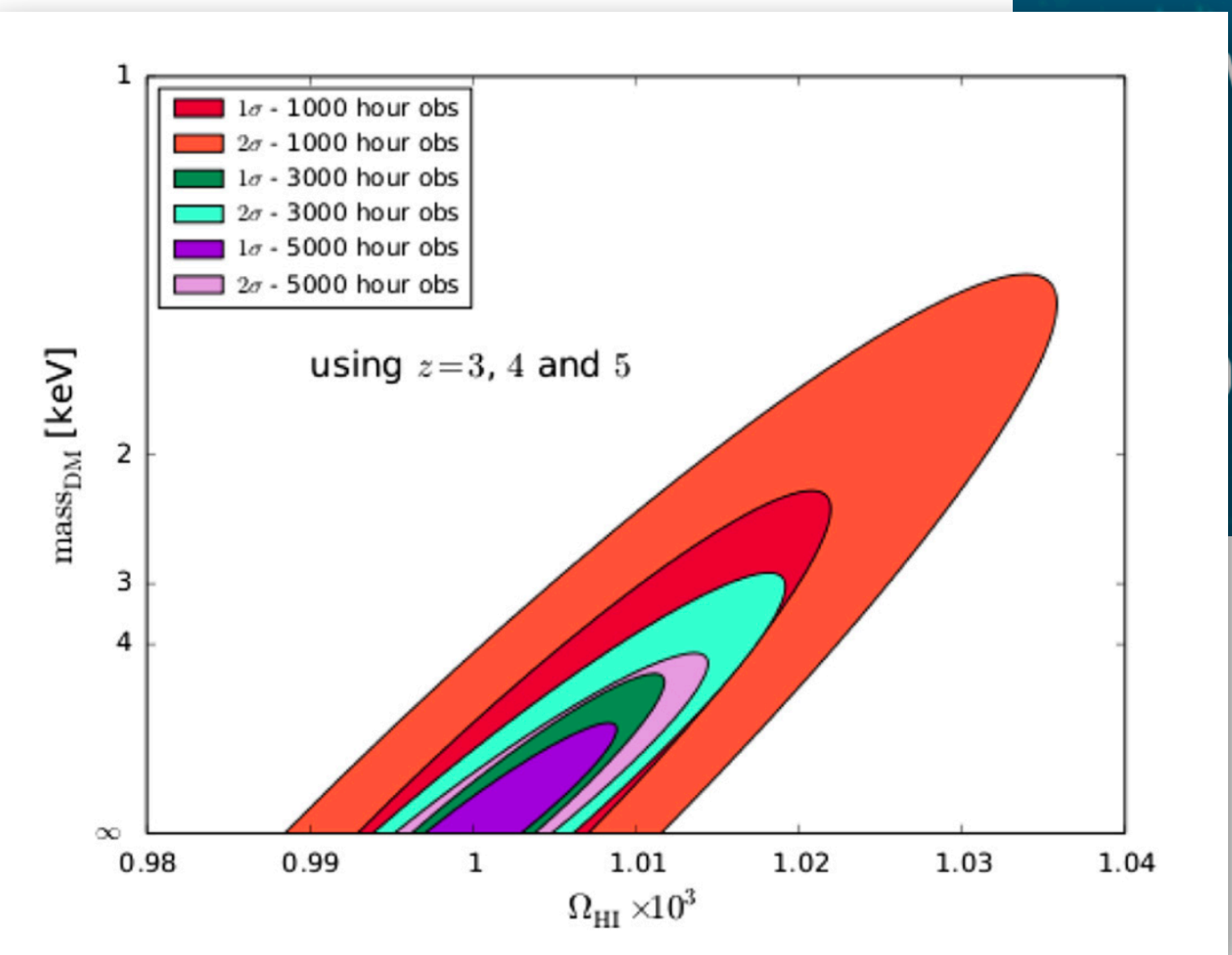
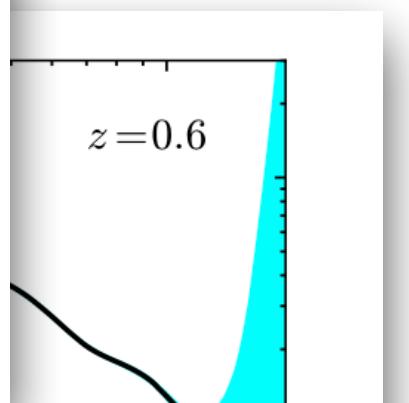
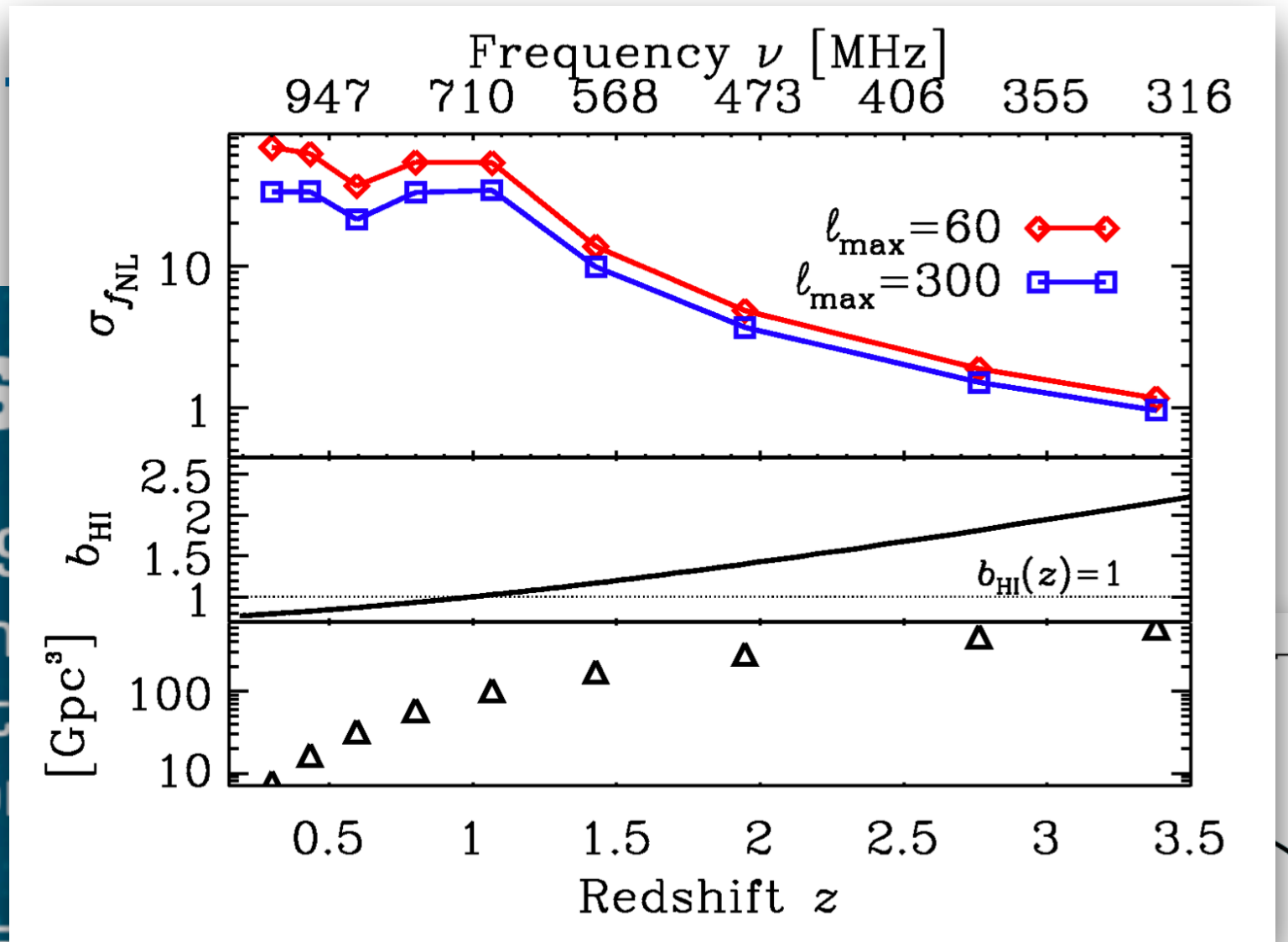
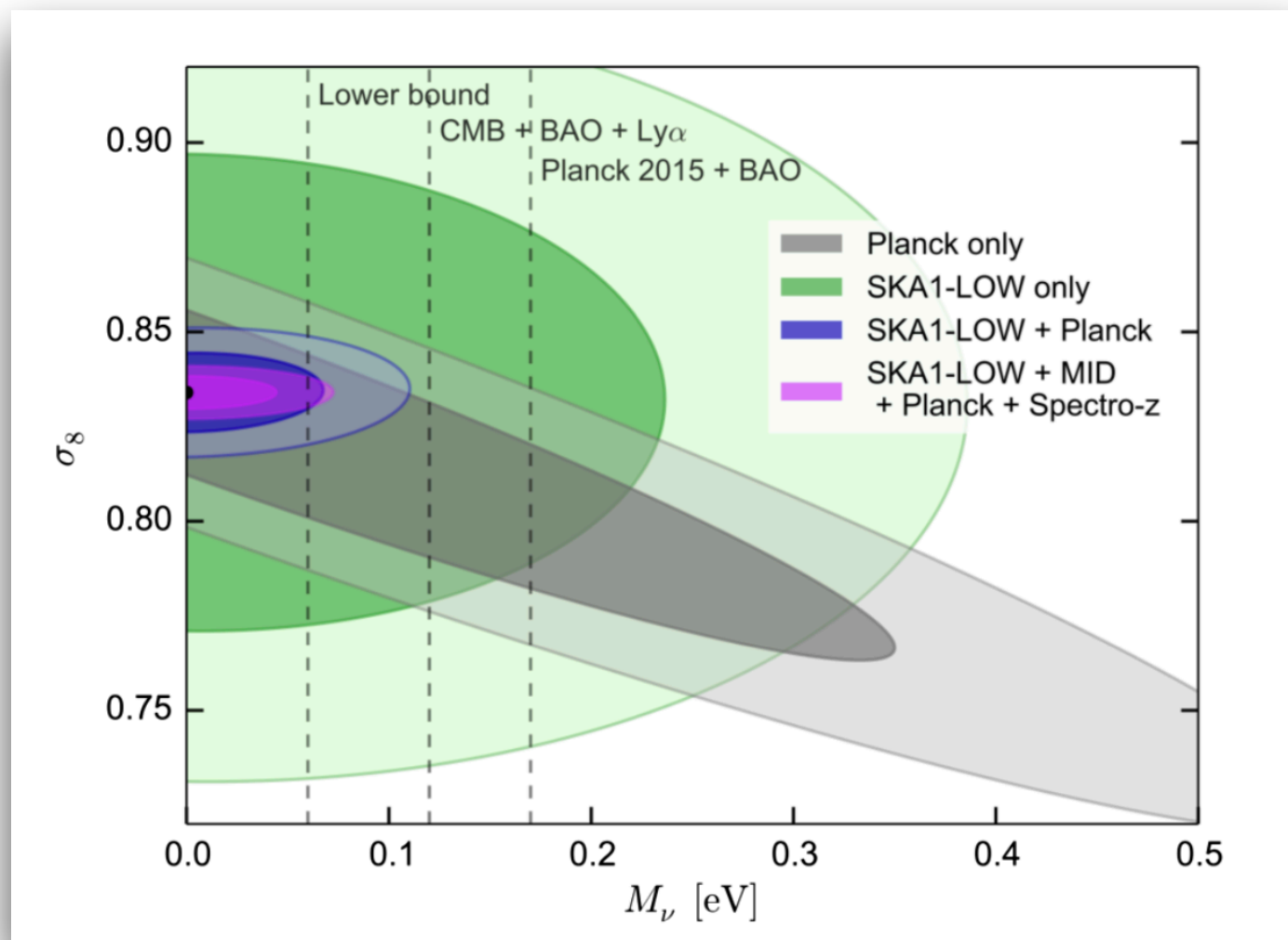
Intensity mapping with

SKA1 Cosmology S

Deep Survey of 5,000 deg² at 0.9
 galaxy redshift survey with 3.5 million
 Lensing shape measurements with
 Continuum galaxy survey with ~60 millio

Wide Survey of 20,000 deg² at 0.35-1.05 GHz for
 Continuum galaxy survey with ~100 million objects
 HI intensity maps for 0.35 < z < 3

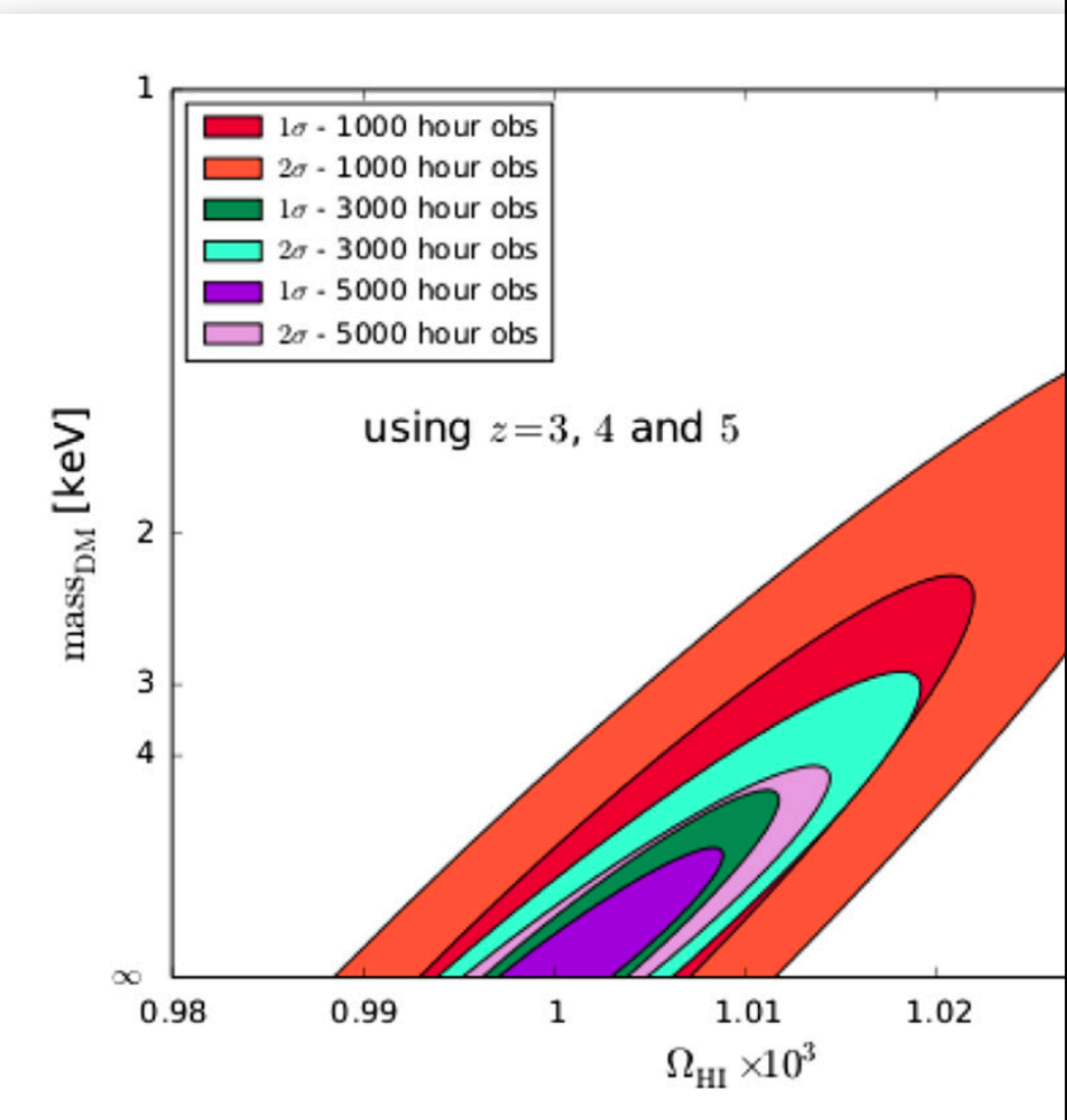
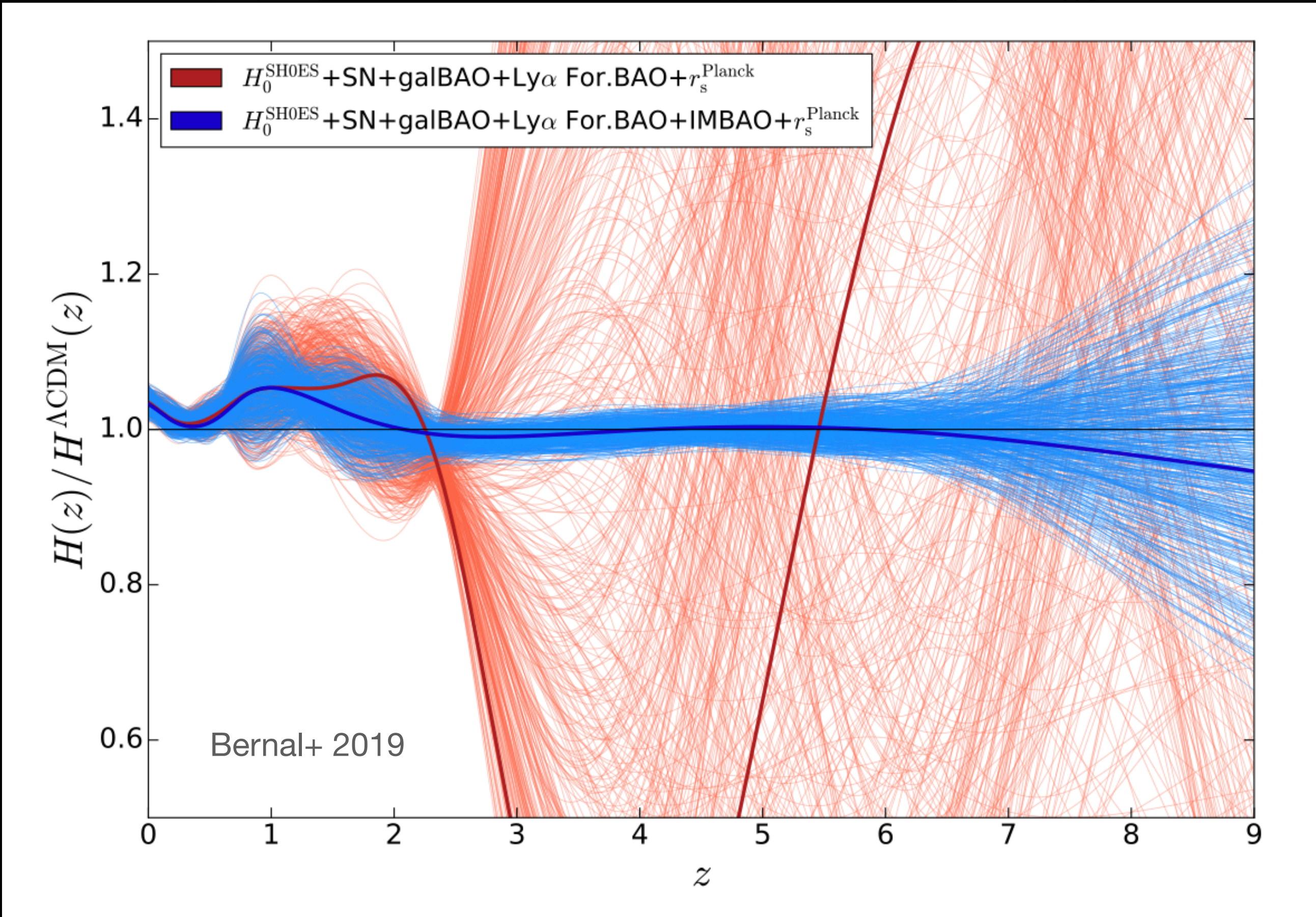
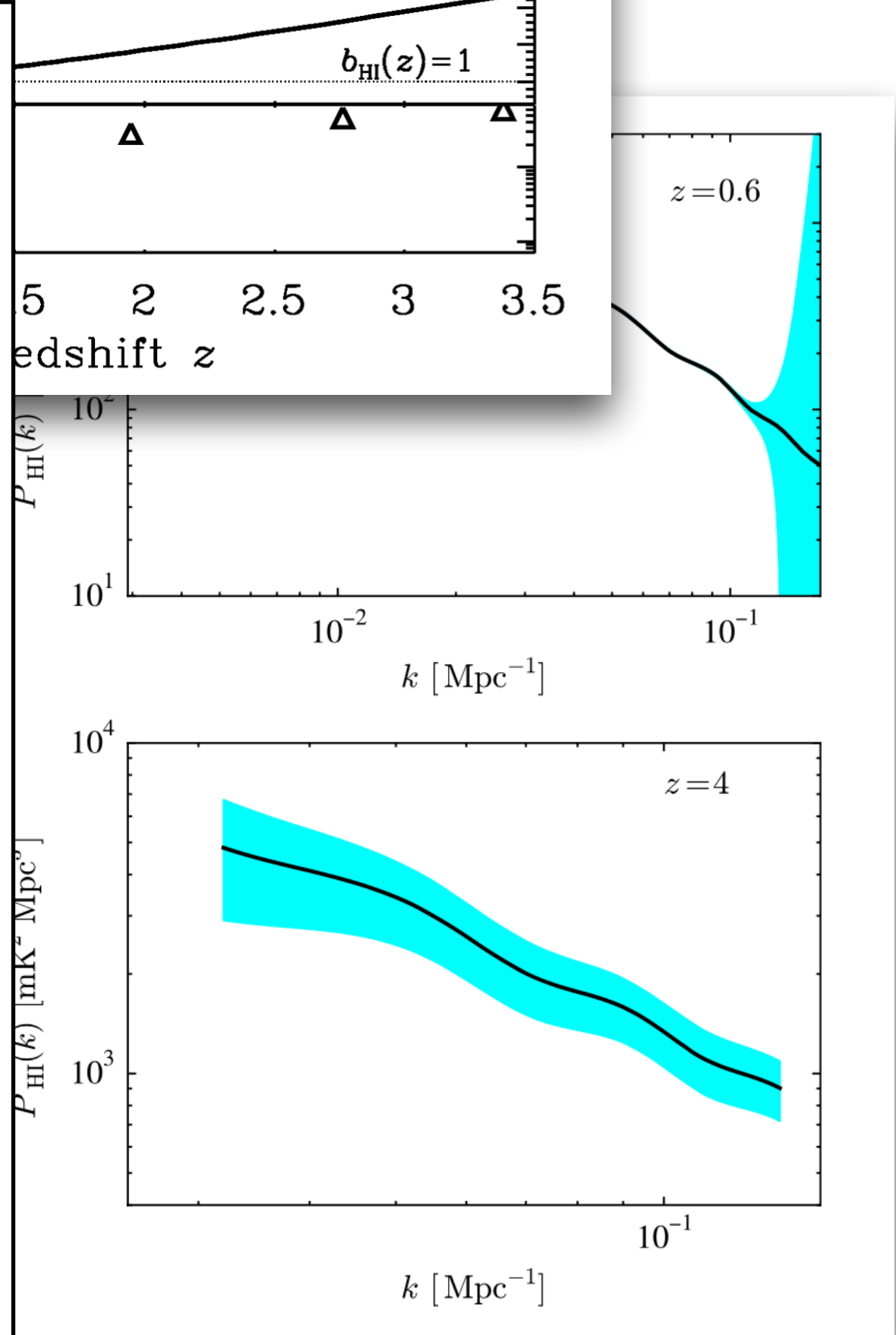
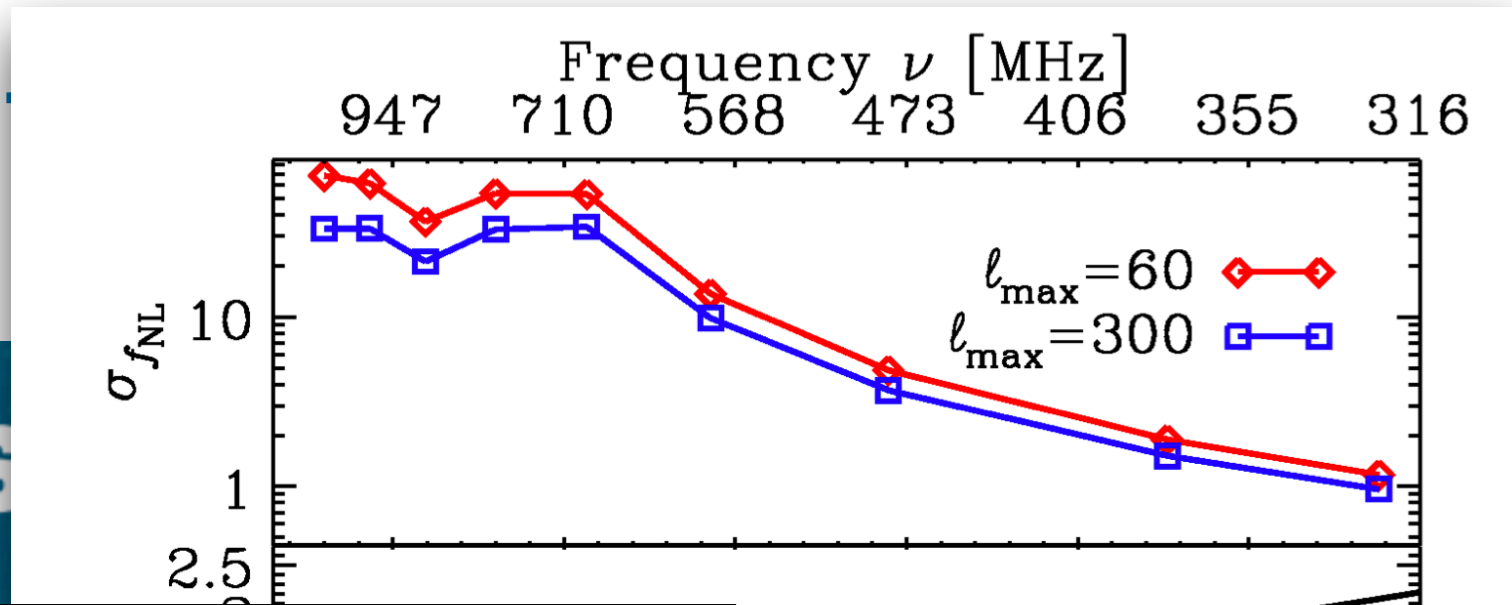
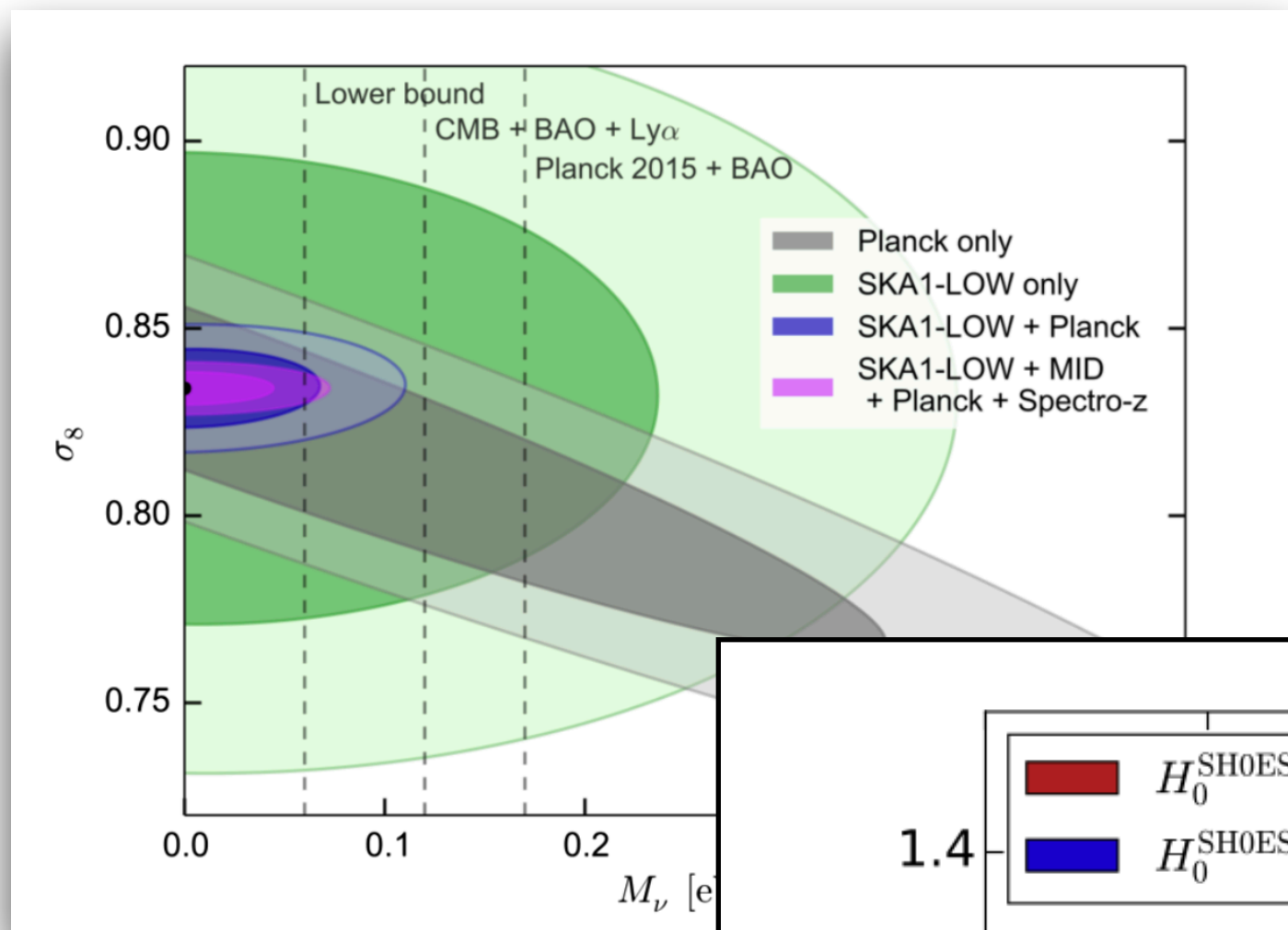
Deep Survey 100 deg² at 200-350 MHz for
 HI intensity maps for 3 < z < 6



Cosmology with Phase 1 of the Square Kilometre Array **Red Book 2018:**
 Technical specifications and performance forecasts

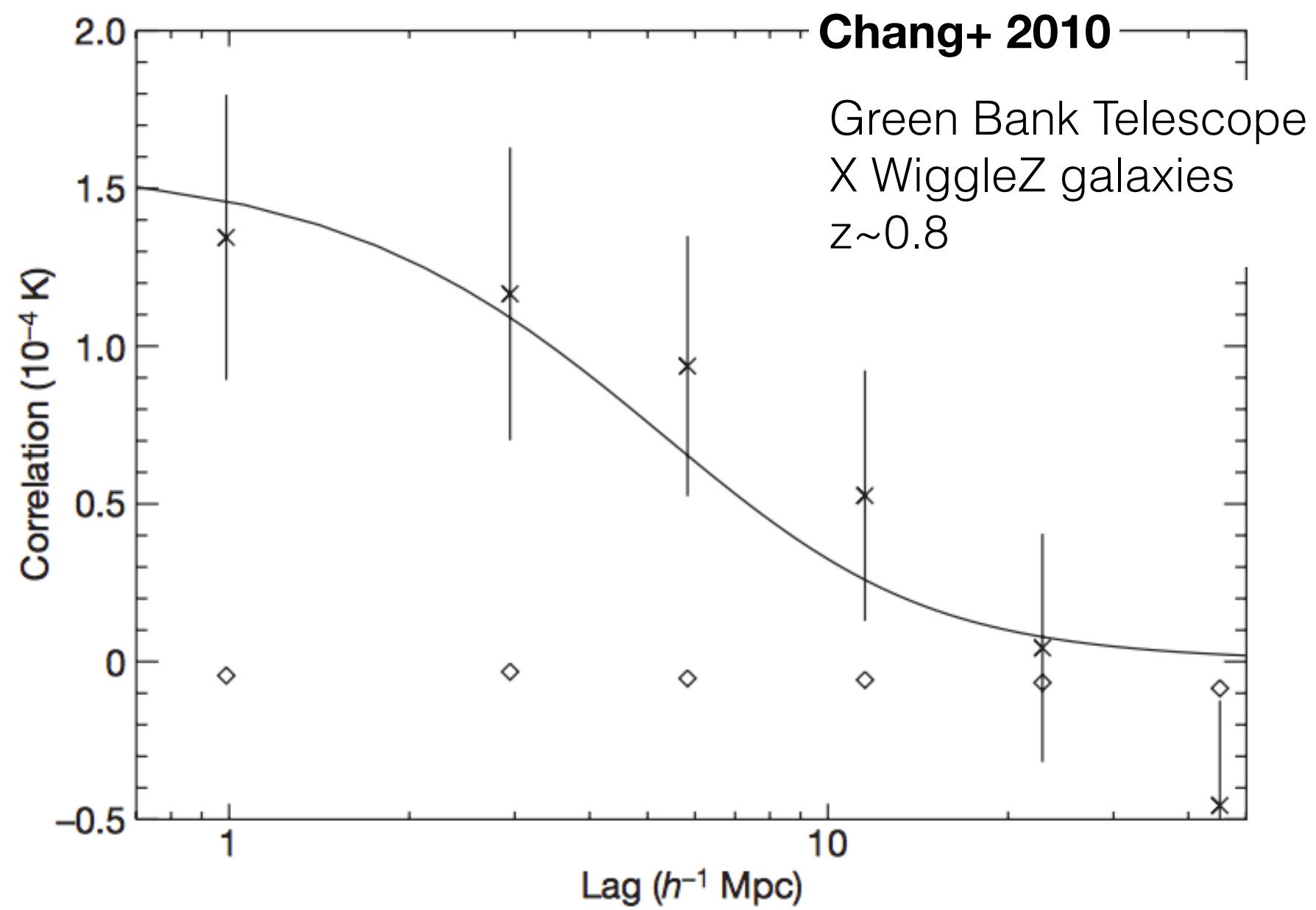
Density mapping with

red SKA1 Cosmology S

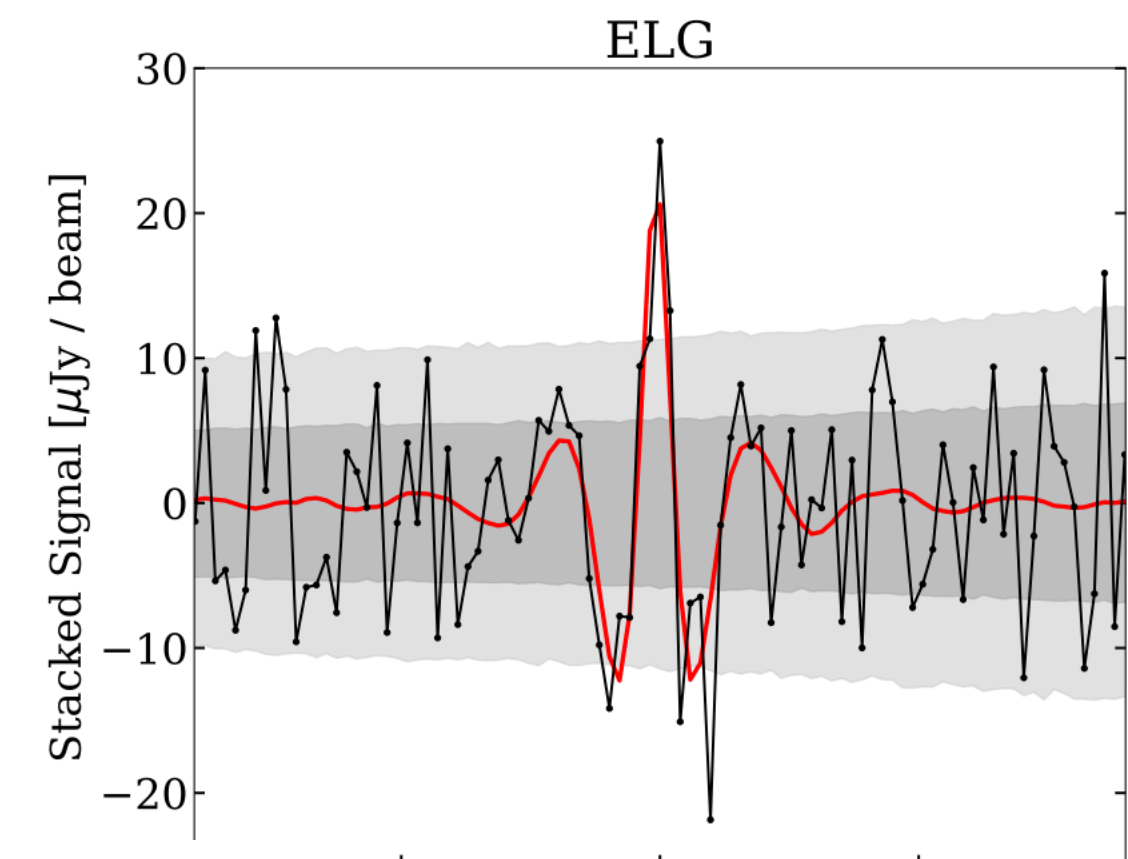
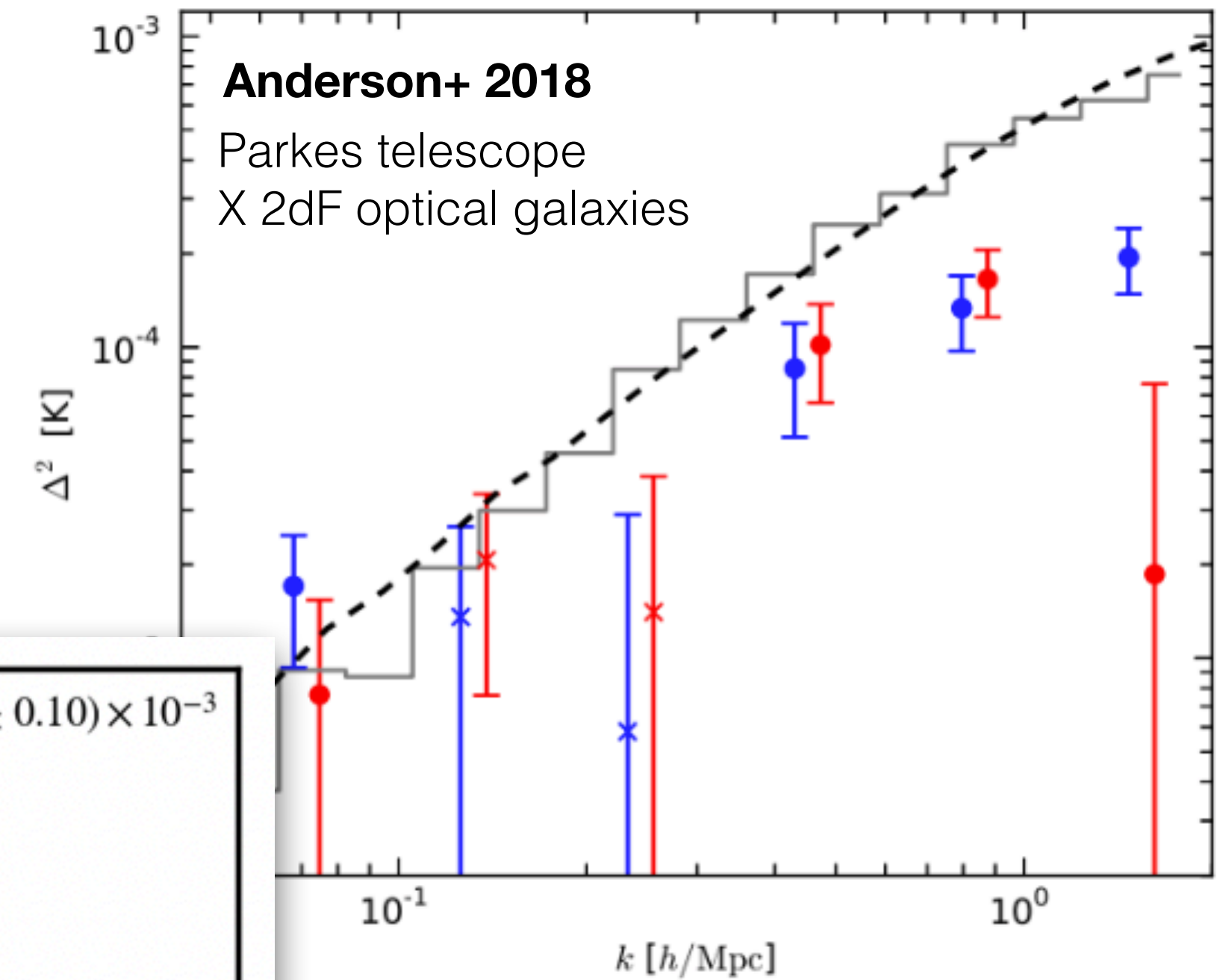
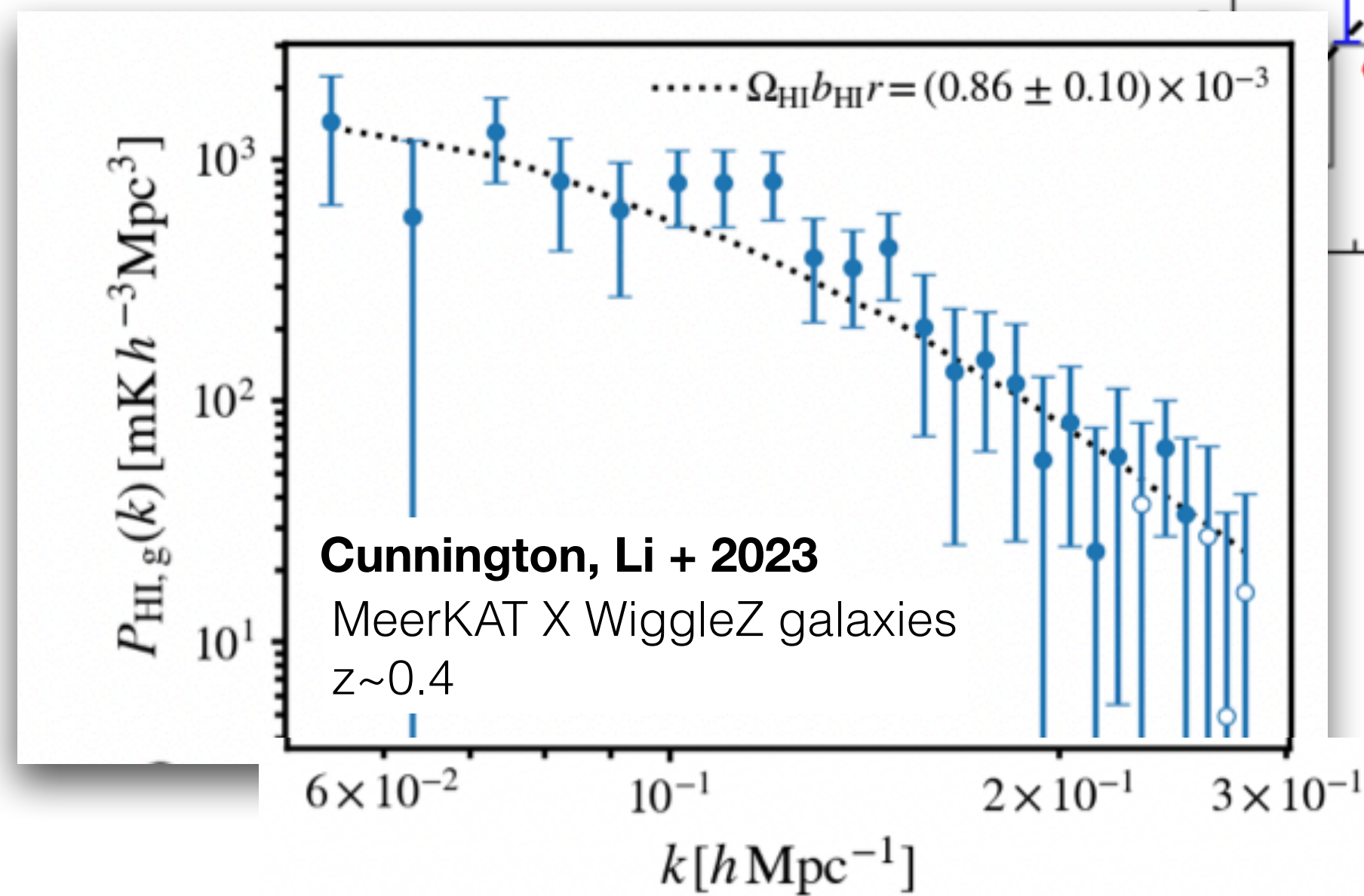


HI intensity mapping

State-of-the-art



also Masui+ 2013, Switzer+ 2013,
Wolz+ 2017,2022



**Contaminants are THE challenge
to overcome with HI intensity
mapping**

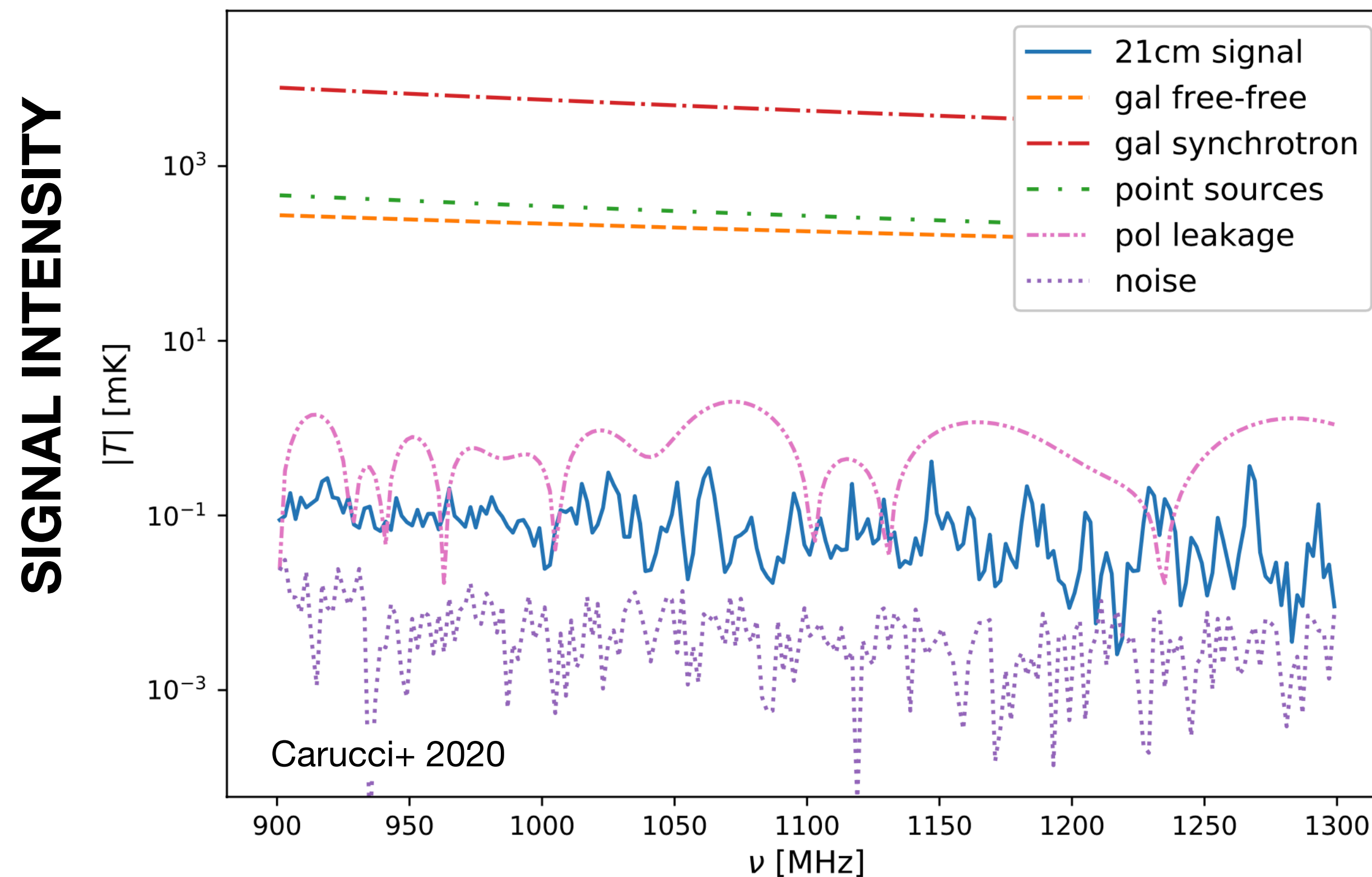


HI intensity mapping

buried under the contaminants

CHALLENGES:

- Foregrounds
- Systematics

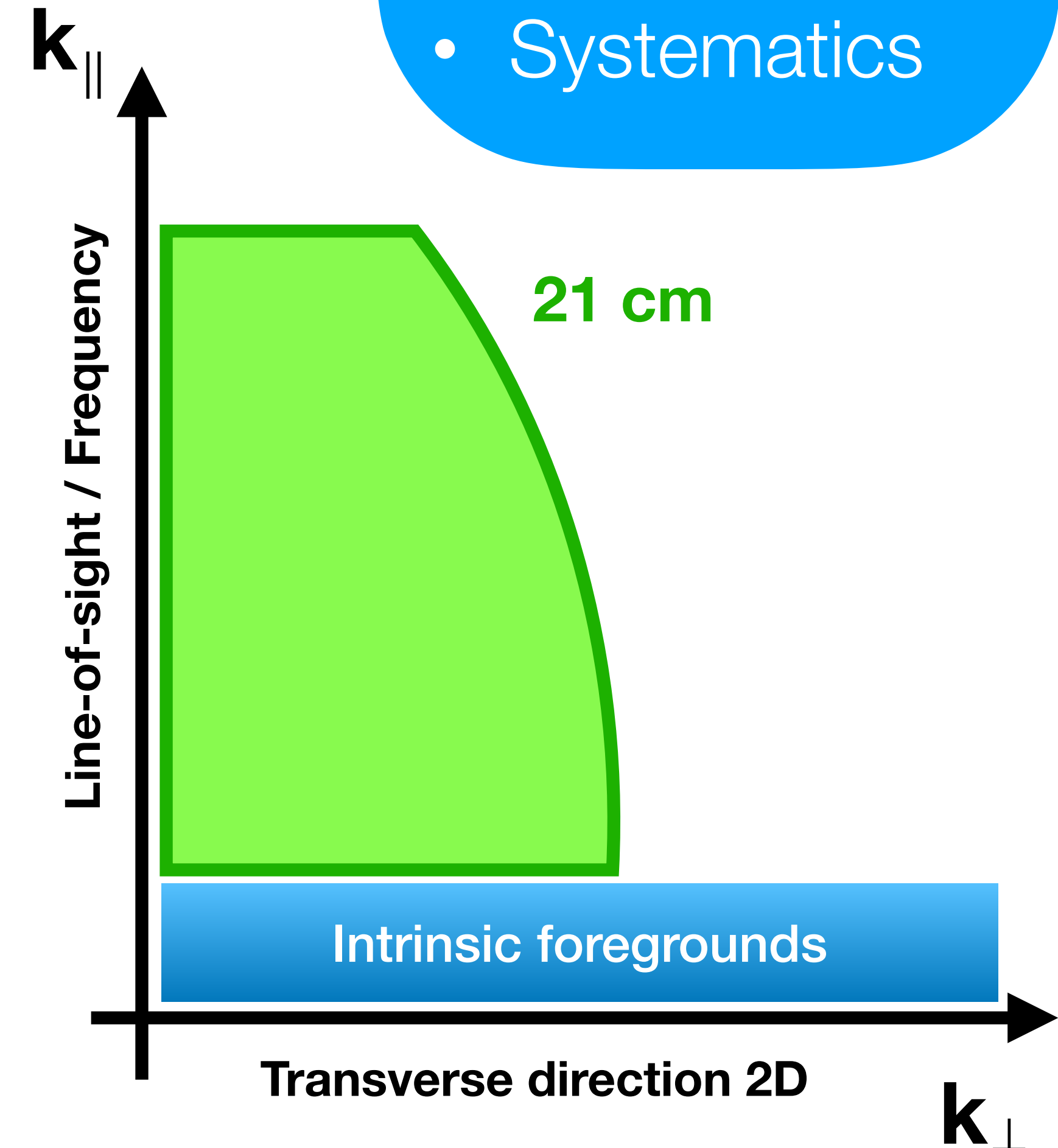
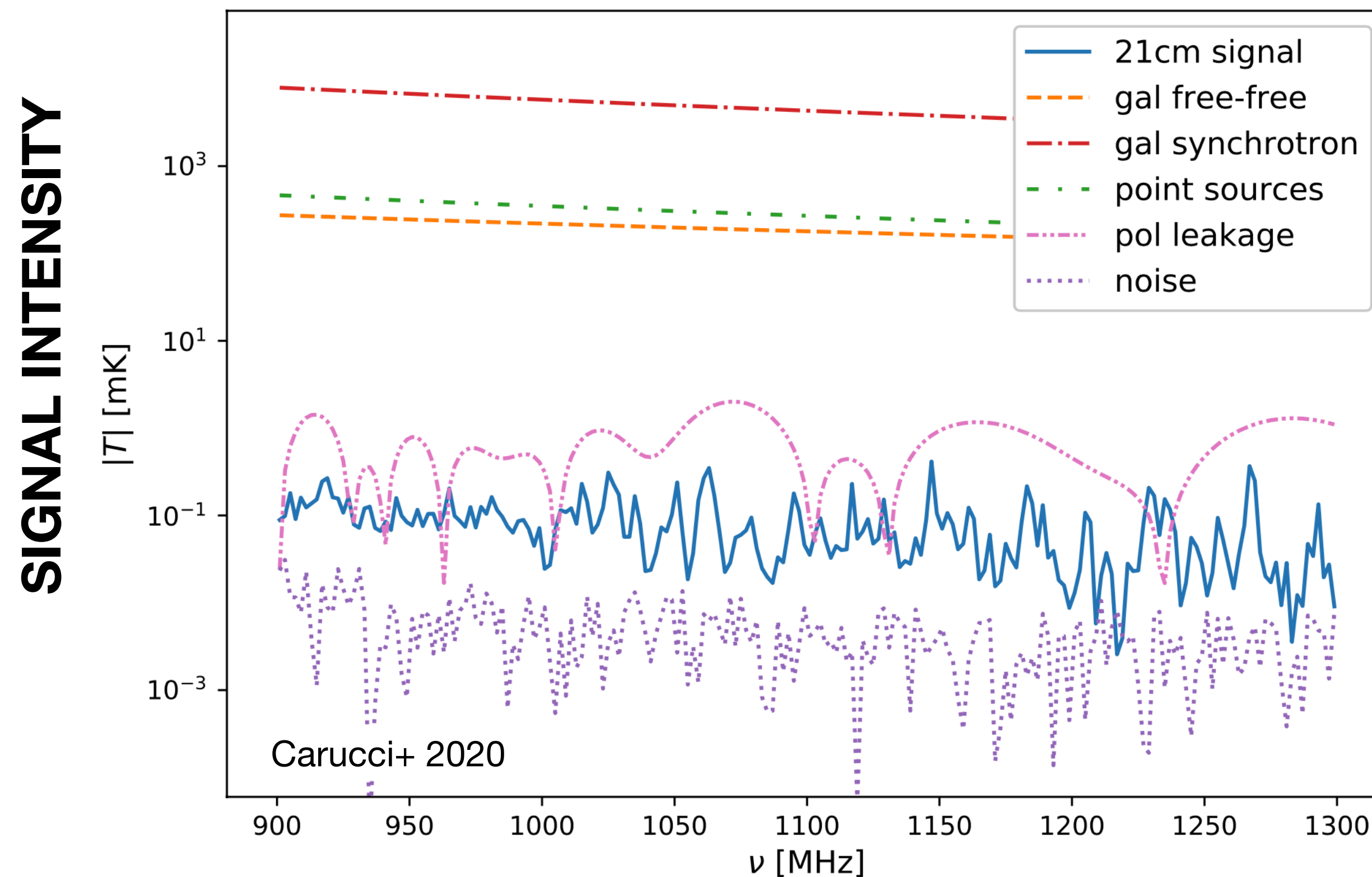


HI intensity mapping

buried under the contaminants

CHALLENGES:

- Foregrounds
- Systematics

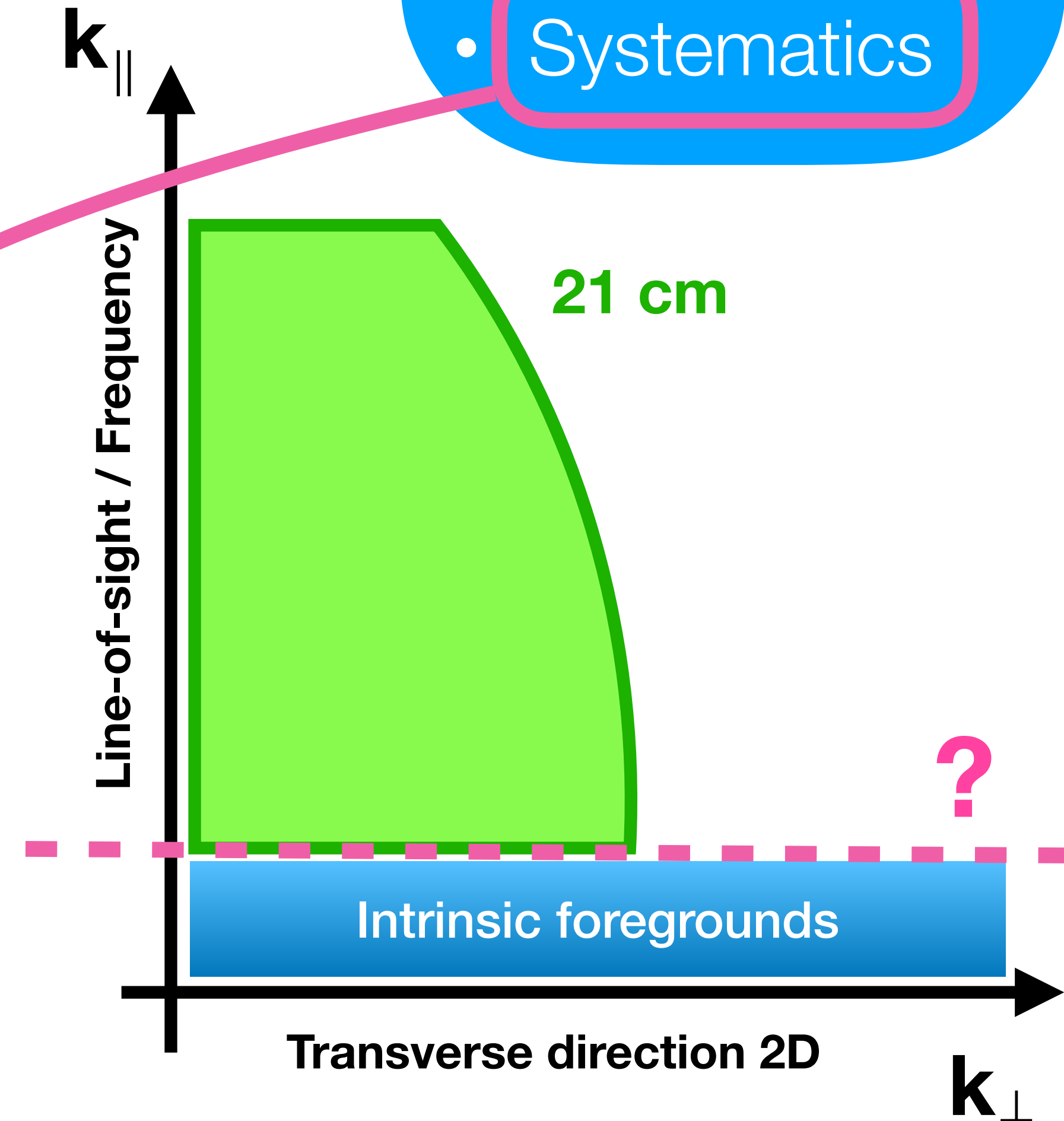
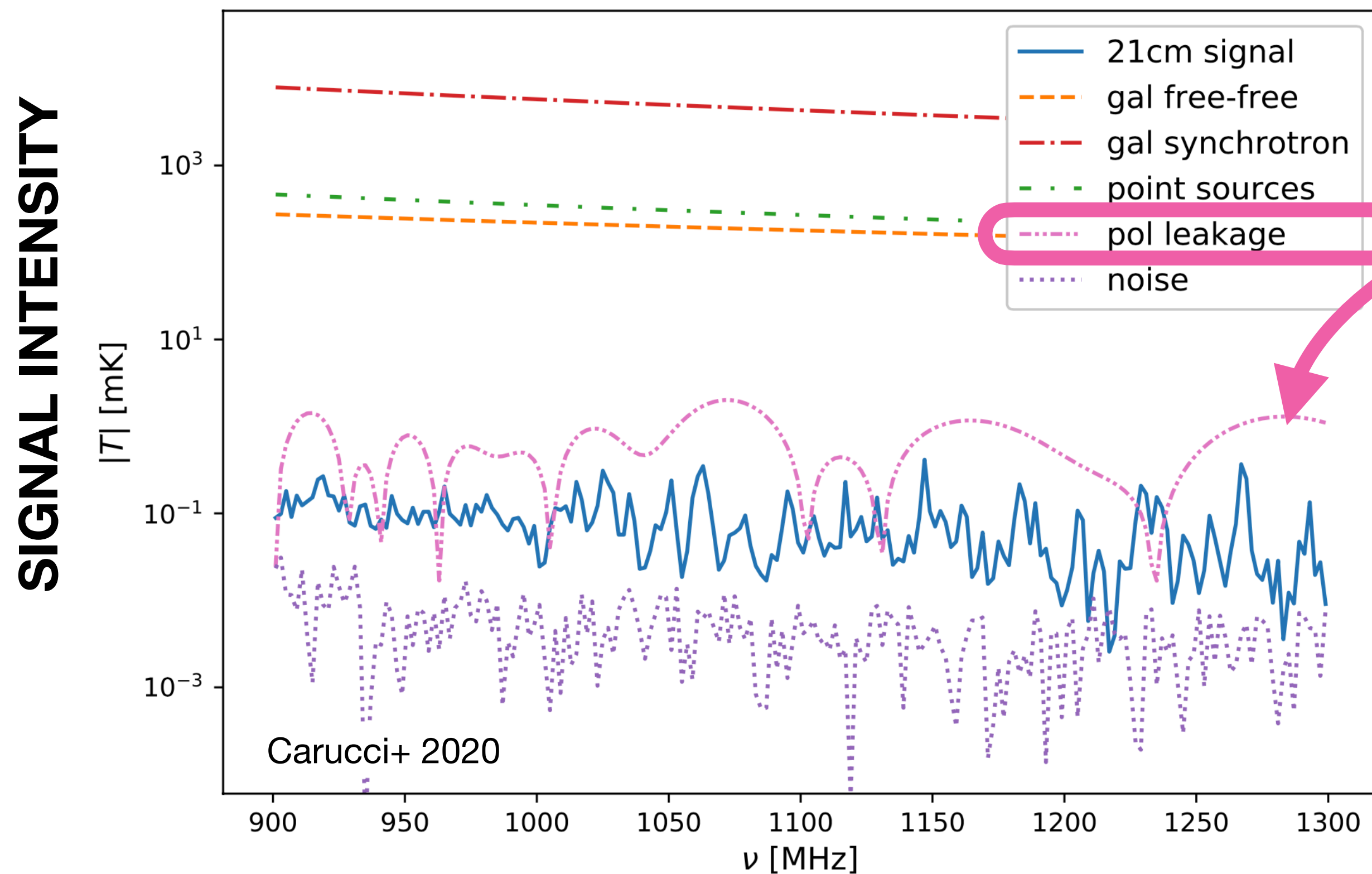


HI intensity mapping

buried under the contaminants

CHALLENGES:

- Foregrounds
- Systematics



Blind Source Separation algorithms

The separation of a set of source signals (contaminants) from a set of mixed signals (the maps), with little or no info about the source signal or the mixing process.

$$\begin{array}{c} \text{signal} \\ \mathbf{X} \\ (f,p) \end{array} = \begin{array}{c} \text{mixing} \\ \text{matrix } (f,n) \\ \mathbf{A} \end{array} \begin{array}{c} \mathbf{S} \\ \text{sources} \\ (n,p) \end{array} + \begin{array}{c} \mathbf{N} \\ \text{HI signal!} \end{array}$$

- **Decorrelation** \rightarrow diagonalise the covariance matrix
- **Independence** \rightarrow as more independent sources are mixed the signal becomes more Gaussian (central limit theorem). So, let's maximise the non-gaussianity of the sources to *unmix* them.

Principal Component Analysis (**PCA**)

Independent Component Analysis (**ICA**)

Hyvarinen + 1999

Blind Source Separation algorithms

The separation of a set of source signals (contaminants) from a set of mixed signals (the maps), with little or no info about the source signal or the mixing process.

Need to set number n of sources!

$$\begin{array}{c} \mathbf{X} \\ \text{signal} \\ (f,p) \end{array} = \begin{array}{c} \text{mixing} \\ \text{matrix } (f,n) \end{array} \begin{array}{c} \mathbf{S} \\ \text{sources} \\ (n,p) \end{array} + \begin{array}{c} \mathbf{N} \\ \text{HI signal!} \end{array}$$

- **Decorrelation** \rightarrow diagonalise the covariance matrix
- **Independence** \rightarrow as more independent sources are mixed the signal becomes more Gaussian (central limit theorem). So, let's maximise the non-gaussianity of the sources to *unmix* them.

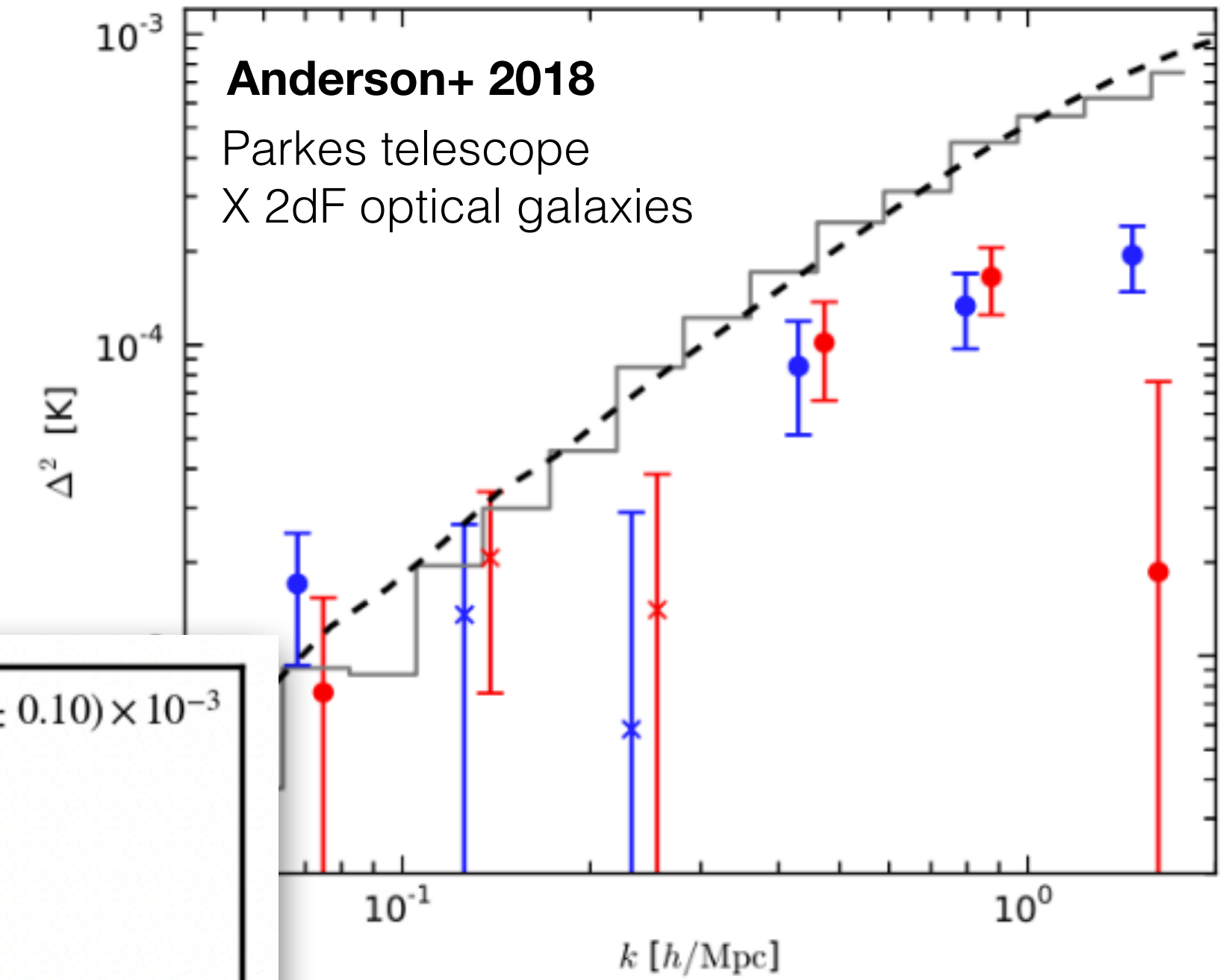
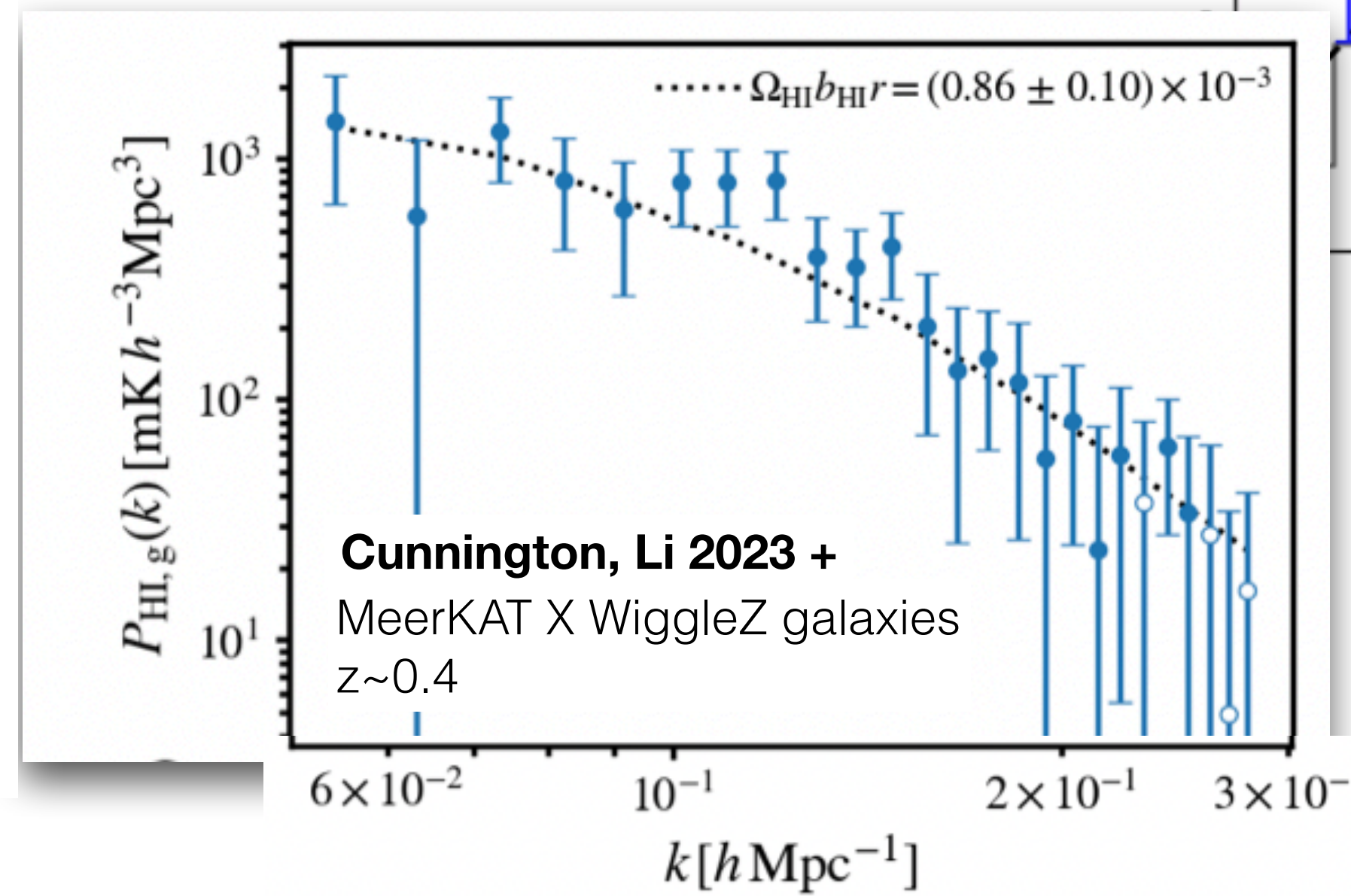
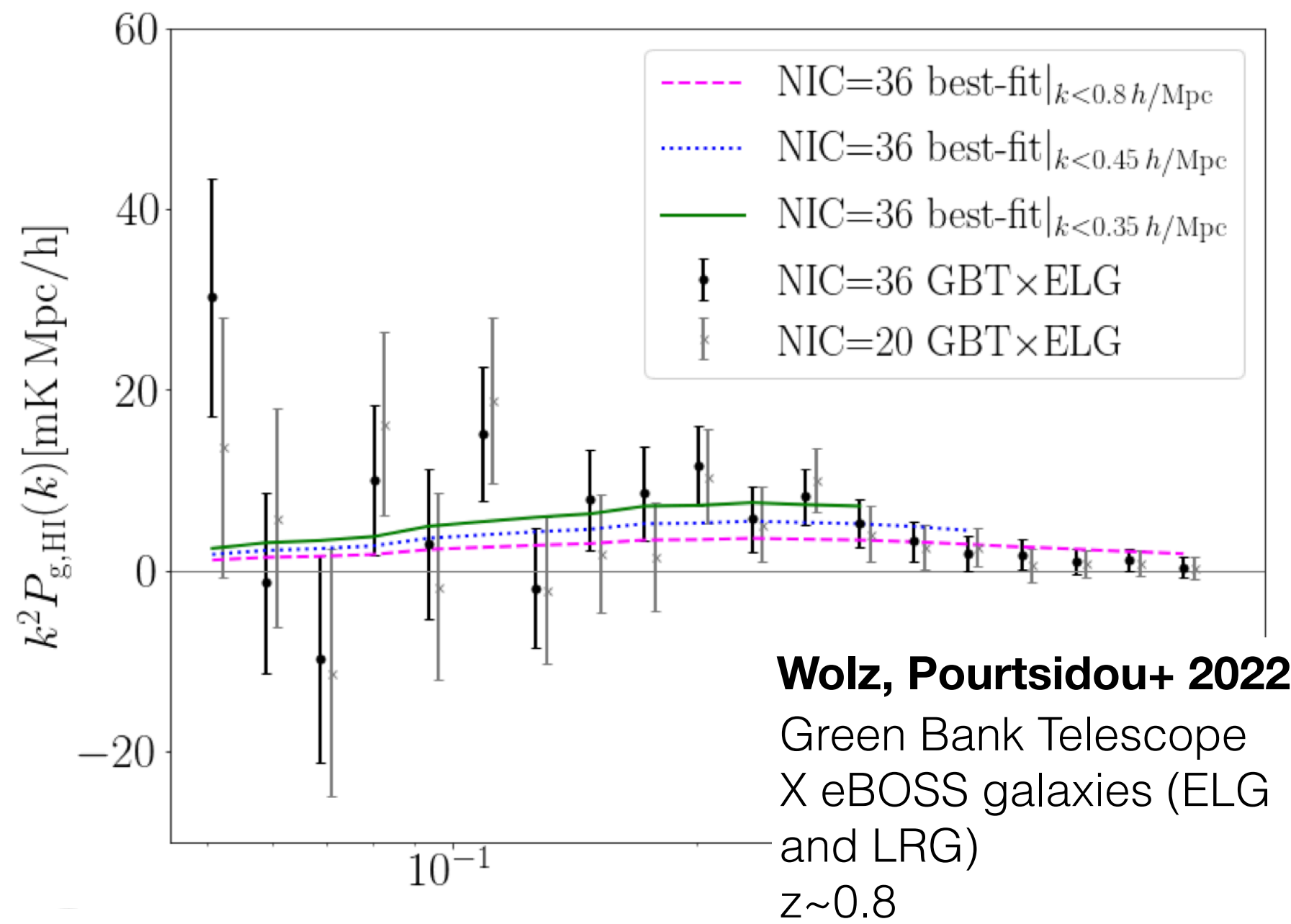
Principal Component Analysis (**PCA**)

Independent Component Analysis (**ICA**)

Hyvarinen + 1999

HI intensity mapping

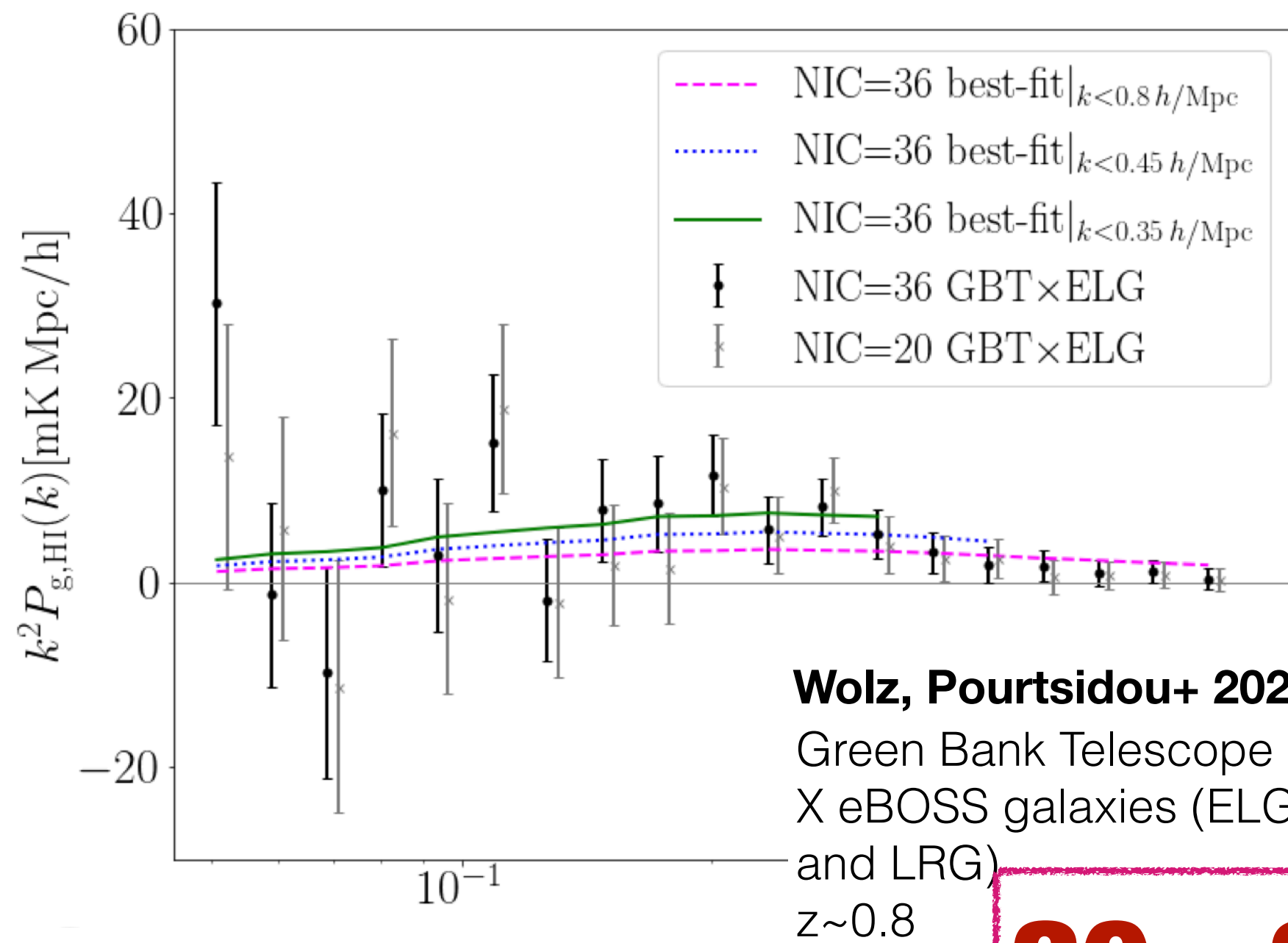
State-of-the-art



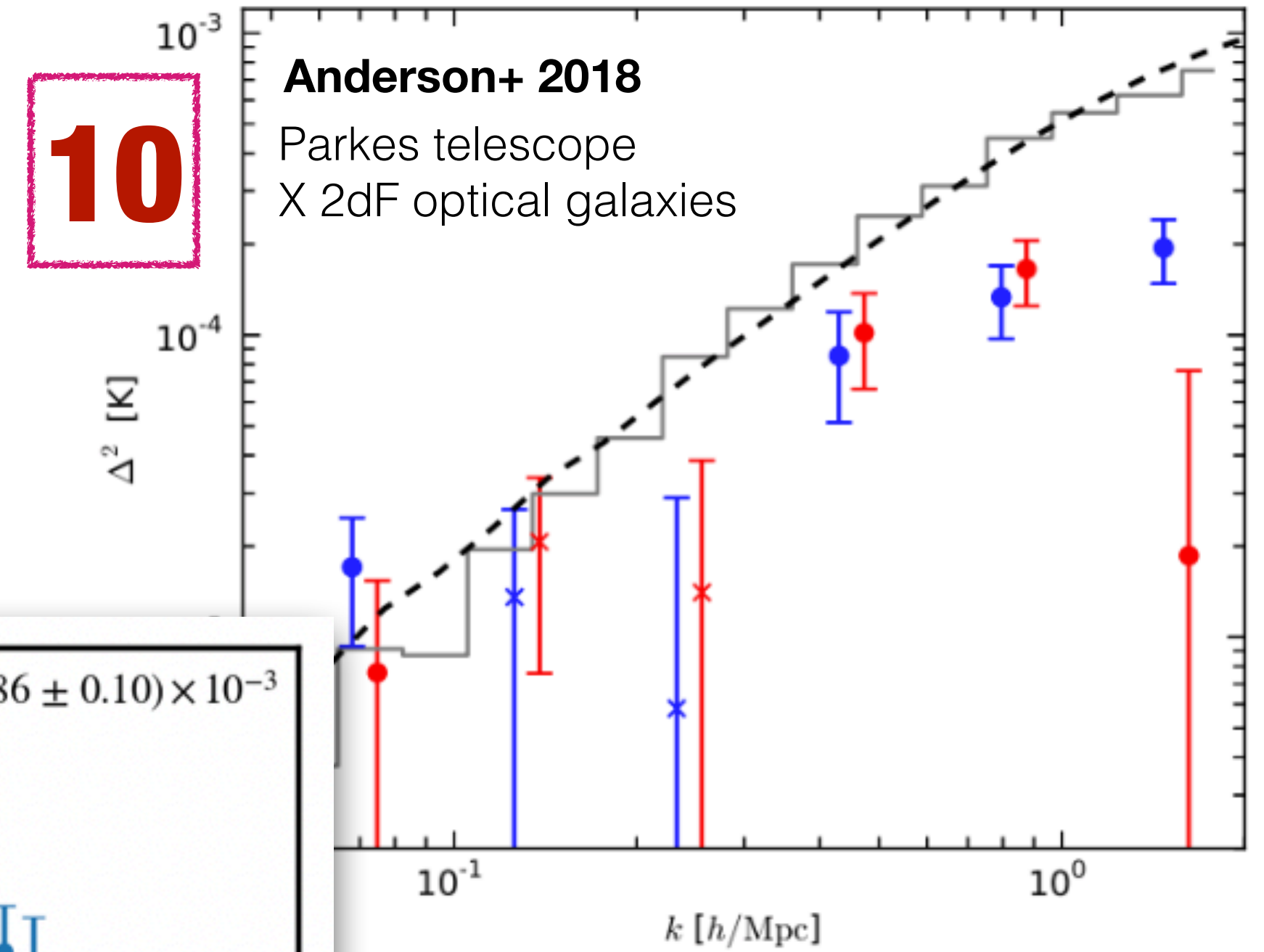
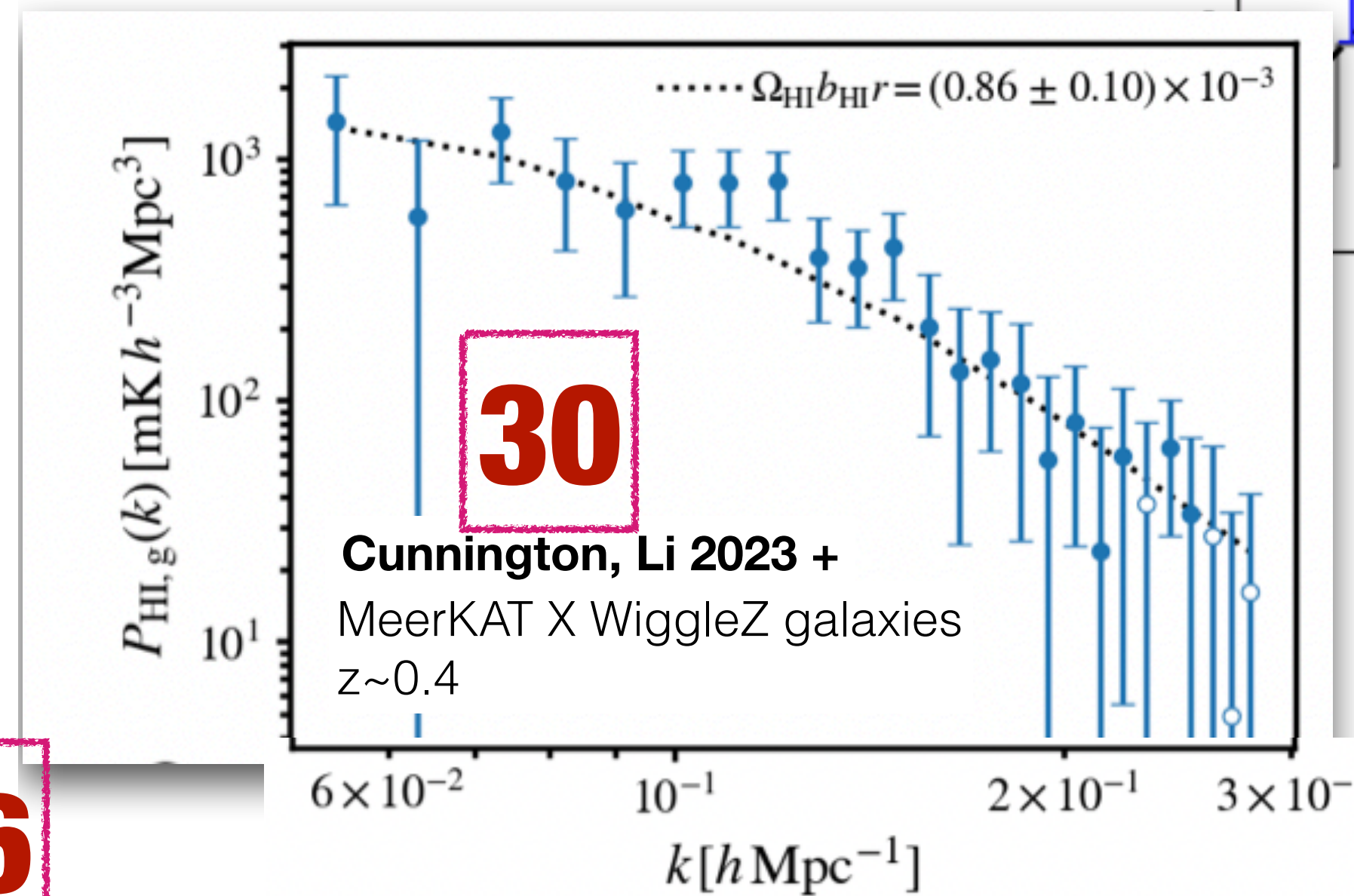
See also Masui+ 2013, Switzer+ 2013, Wolz+ 2017

HI intensity mapping

State-of-the-art



20 - 36



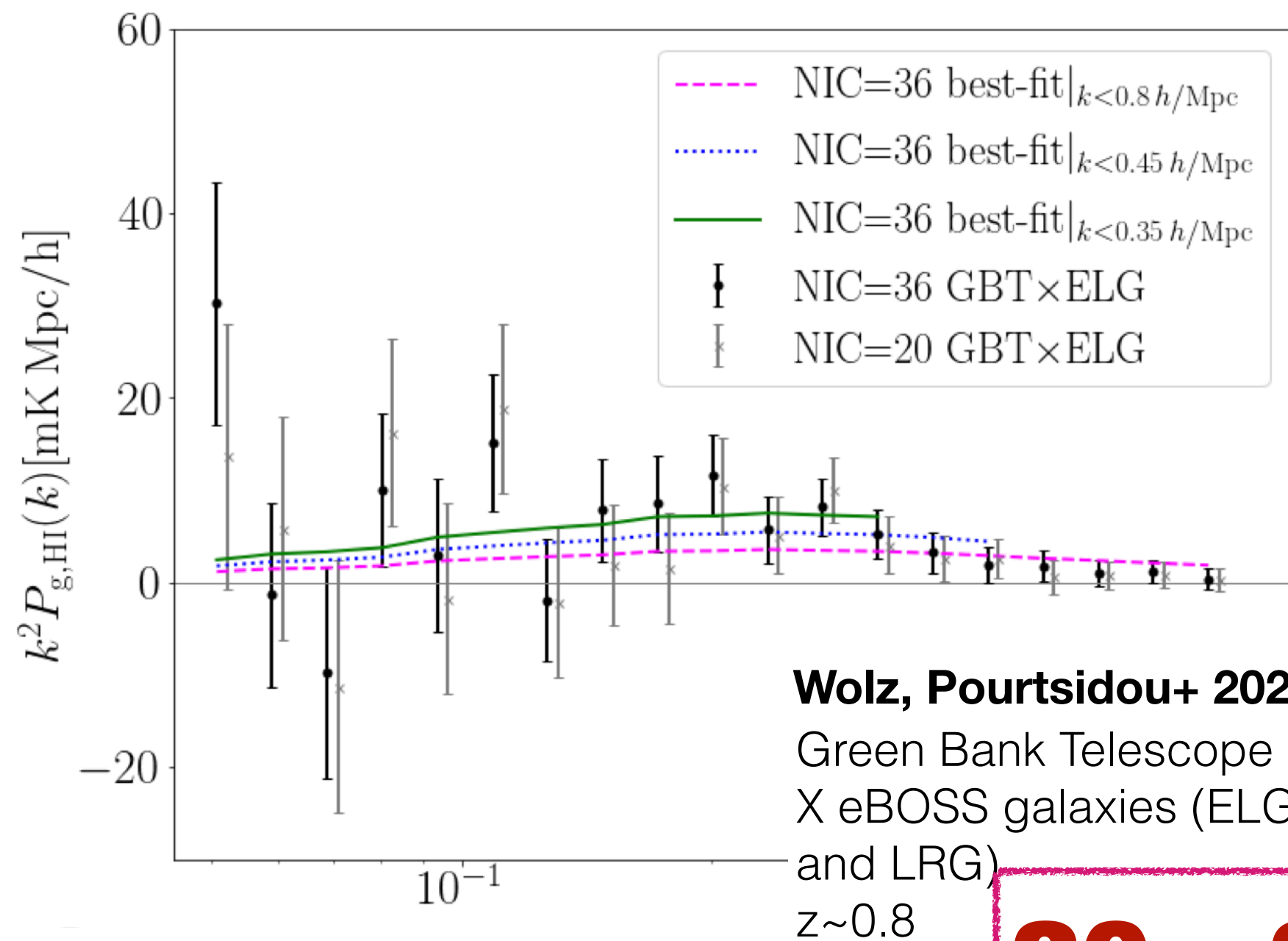
10

See also Masui+ 2013, Switzer+ 2013, Wolz+ 2017

10 - 20

HI intensity mapping

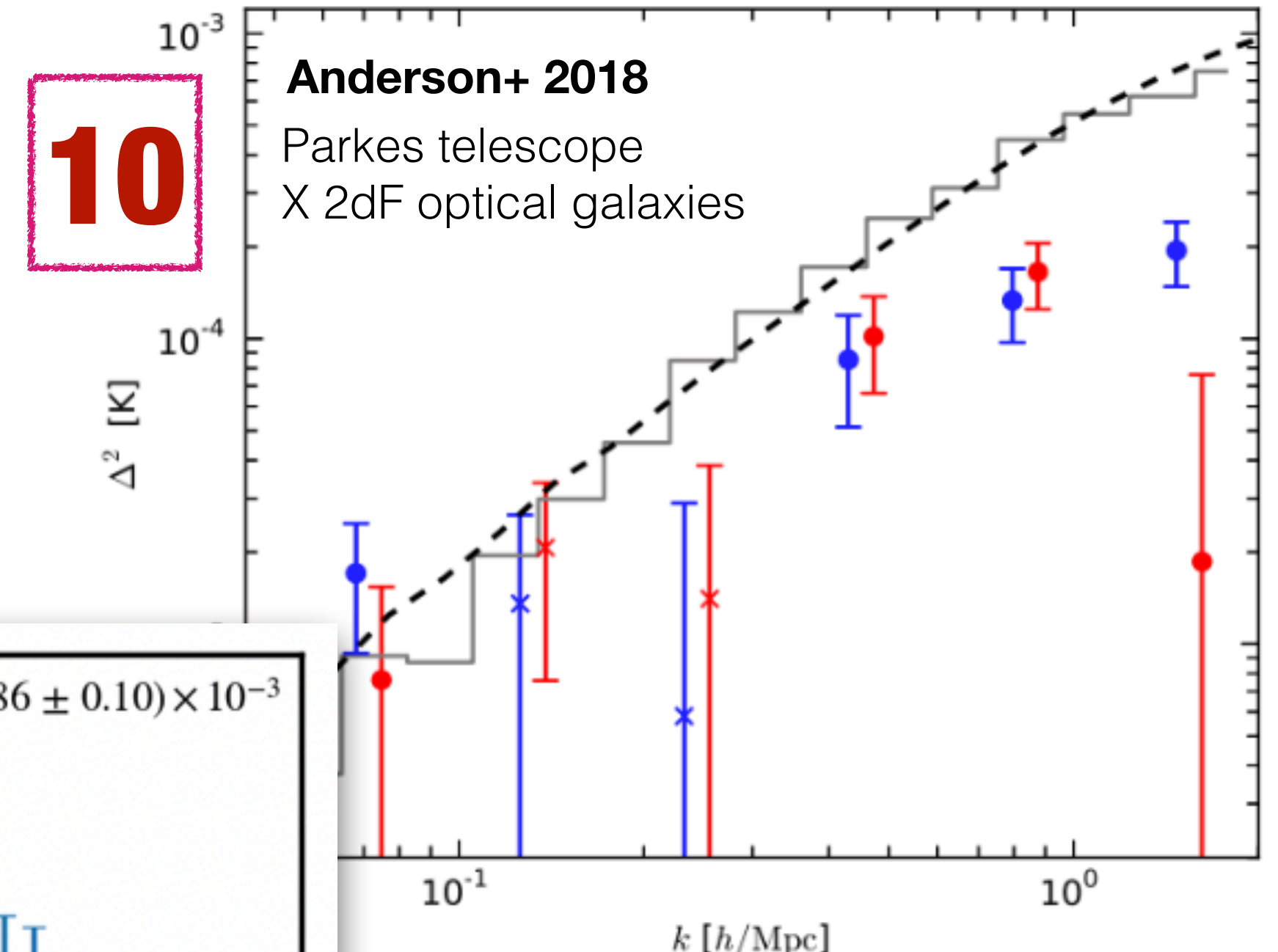
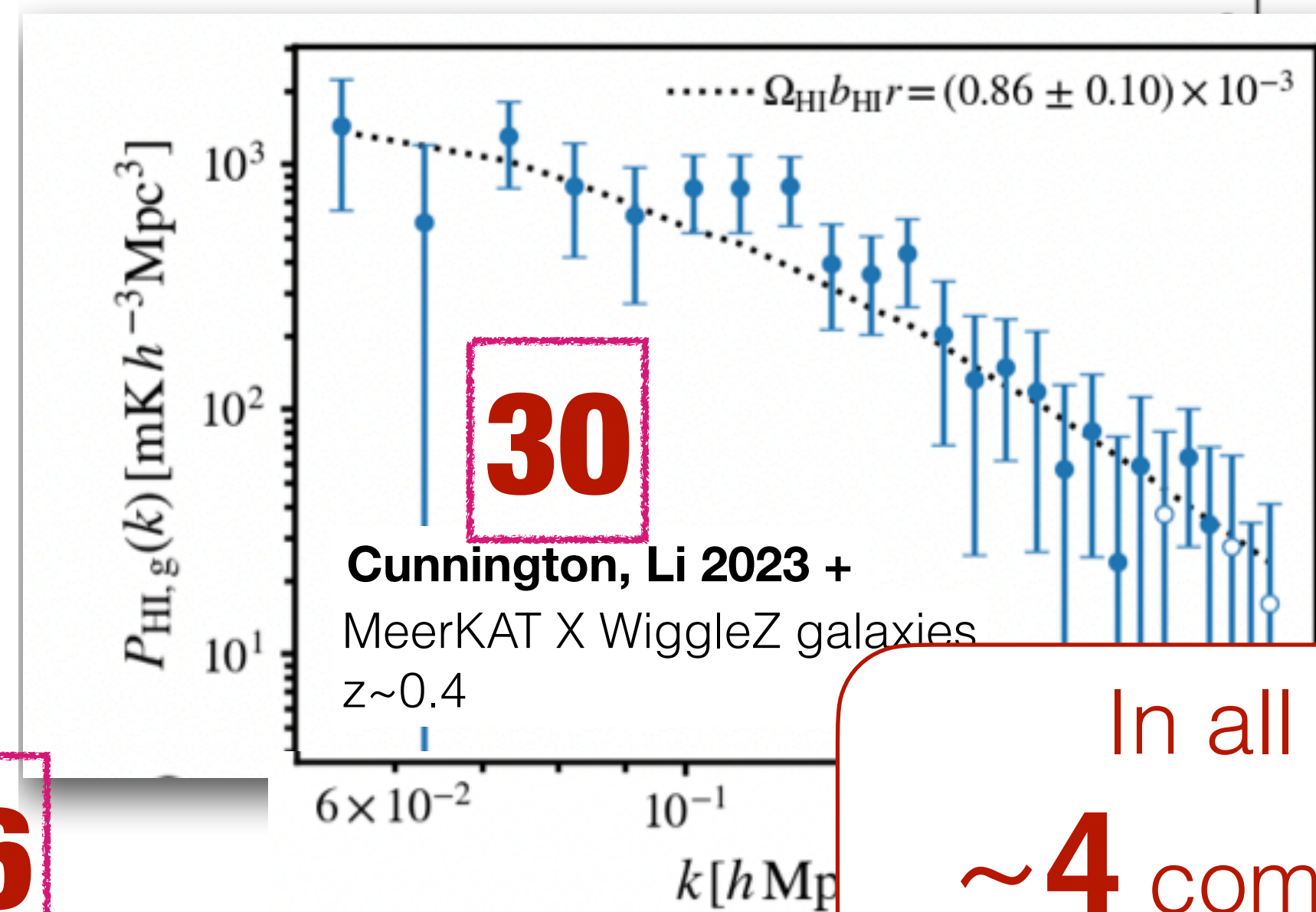
State-of-the-art



20 - 36

See also Masui+ 2013, Switzer+ 2013, Wolz+ 2017

10 - 20



In all theoretical works:
~4 components removed are
 enough
 (e.g., Wolz+ 2014, Alonso+ 2015, Cunnington+ 2019, ...)

1. **Hydrogen Intensity Mapping (IM):**
what is it and why to do it
2. **IM is hard!** Biggest challenge:
weakness of the IM signal compared
to contaminants

3. We are getting there. **MeerKLASS**

Blind Source Separation algorithms

The separation of a set of source signals (contaminants) from a set of mixed signals (the maps), with little or no info about the source signal or the mixing process.

$$\begin{array}{c} \mathbf{X} \\ \text{signal} \\ (f,p) \end{array} = \begin{array}{c} \text{mixing} \\ \text{matrix } (f,n) \\ \mathbf{A} \end{array} \begin{array}{c} \mathbf{S} \\ \text{sources} \\ (n,p) \end{array} + \begin{array}{c} \mathbf{N} \\ \text{HI signal!} \end{array}$$

- **Decorrelation** —> diagonalise the covariance matrix
- **Independence** —> as more independent sources are mixed the signal becomes more Gaussian (central limit theorem). So, let's maximise the non-gaussianity of the sources to *unmix* them.
- **Sparsity** —> mixtures are less sparse than sources!

Principal Component Analysis (**PCA**)

Independent Component Analysis (**ICA**)

Generalised Morphological Component Analysis (**GMCA**)

Bobin + 2007, 2008, 2012

Blind Source Separation algorithms

The separation of a set of source signals (contaminants) from a set of mixed signals (the maps), with little or no info ab

$$\begin{array}{c} \mathbf{X} \\ \text{signal} \\ (f,p) \end{array} = \begin{array}{c} \text{mixing} \\ \text{matrix } (f,n) \\ \mathbf{A} \end{array} \begin{array}{c} \mathbf{S} \\ \text{sources} \\ (n,p) \end{array} + \begin{array}{c} \mathbf{N} \\ \text{HI signal} \end{array}$$

- **Decorrelation** —> diagonalise the covariance matrix
- **Independence** —> as more independent sources are mixed the signal becomes more Gaussian (central limit theorem). So, let's maximise the non-gaussianity of the sources to *unmix* them.
- **Sparsity** —> mixtures are less sparse than sources!

Tested on **CMB** data
(e.g. Bobin+ 2016)
on **EoR** signal (Patil+ 2017)
on **X-ray** images of **Supernova**
remnants (Picquenot+ 2019)

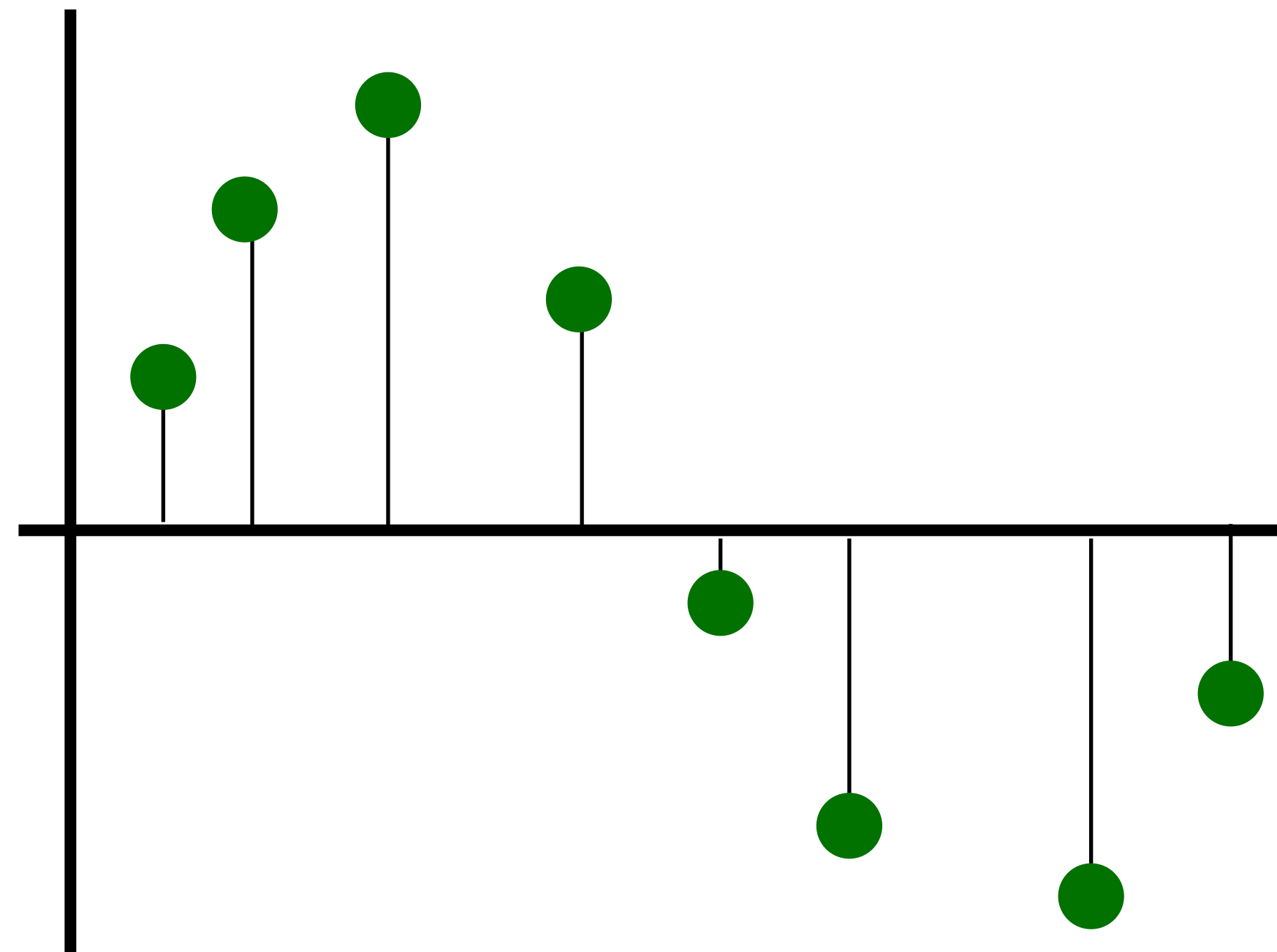
Principal Component
Analysis (**PCA**)

Independent
Component Analysis
(**ICA**)

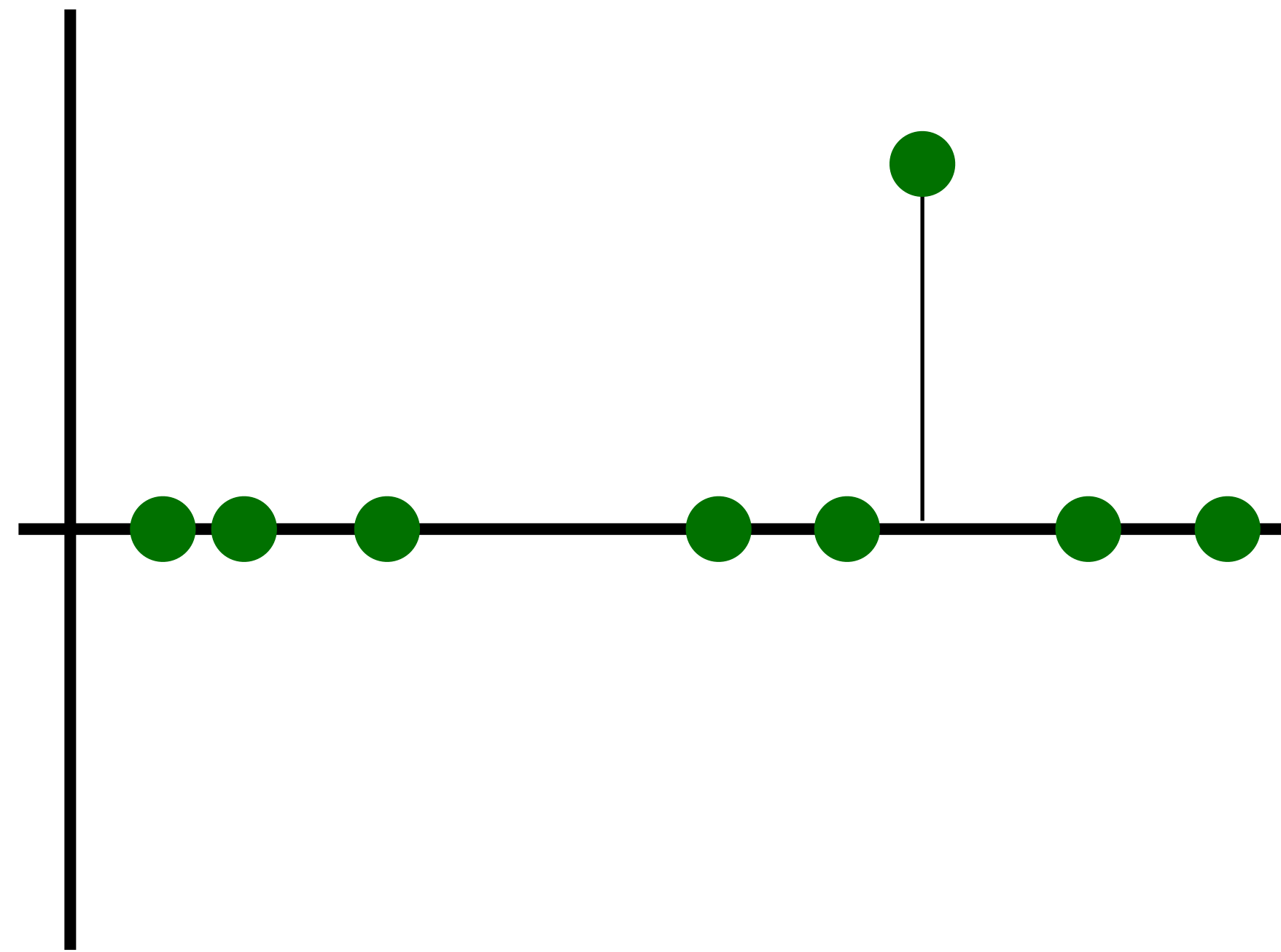
Generalised Morphological
Component Analysis
(**GMCA**)

Bobin + 2007, 2008, 2012

sparsity

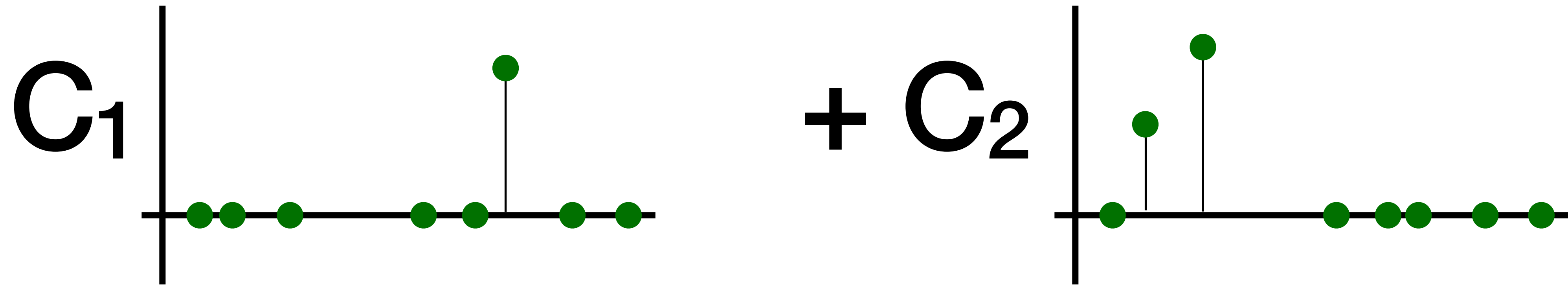


sparsity



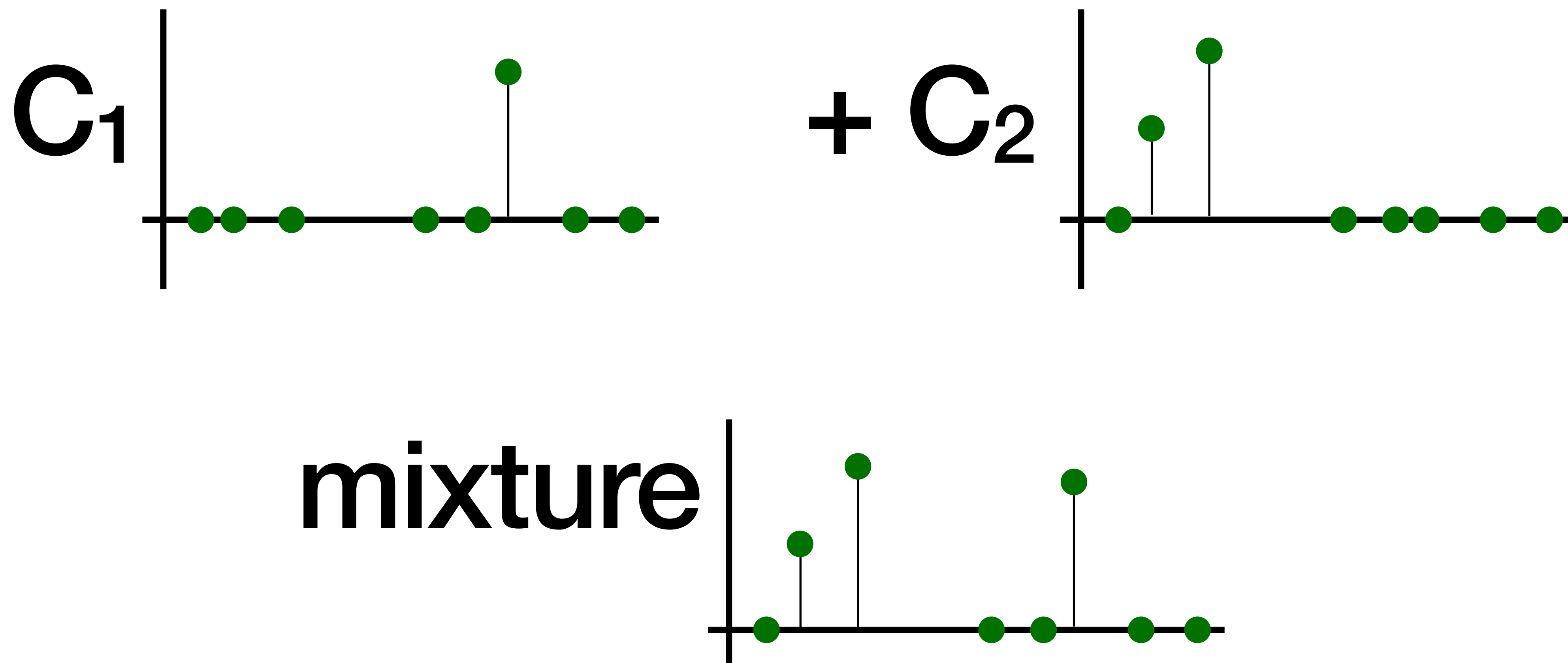
why sparsity?

mixtures are less sparse than components

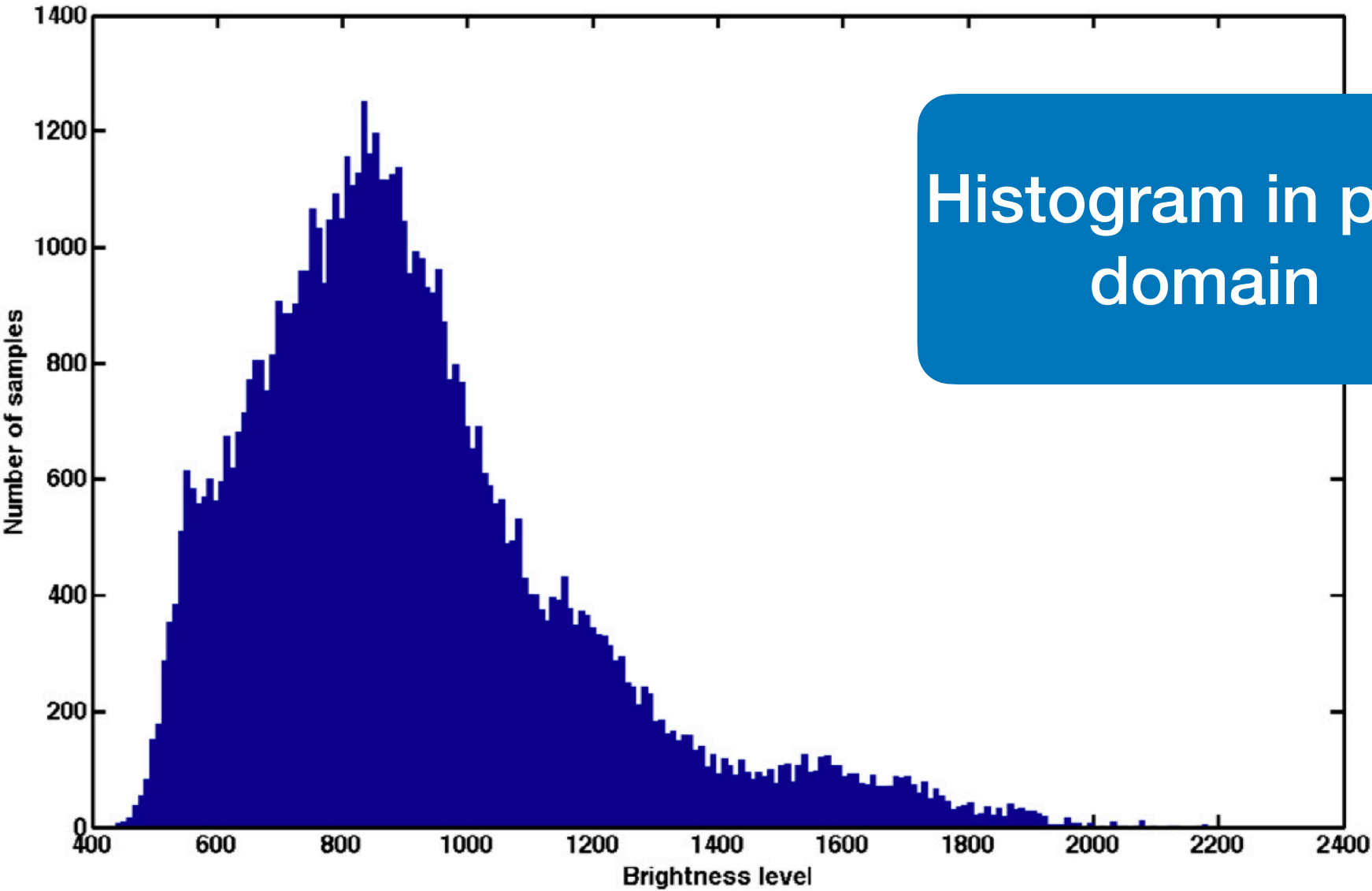
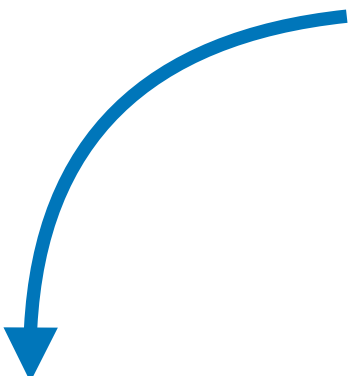
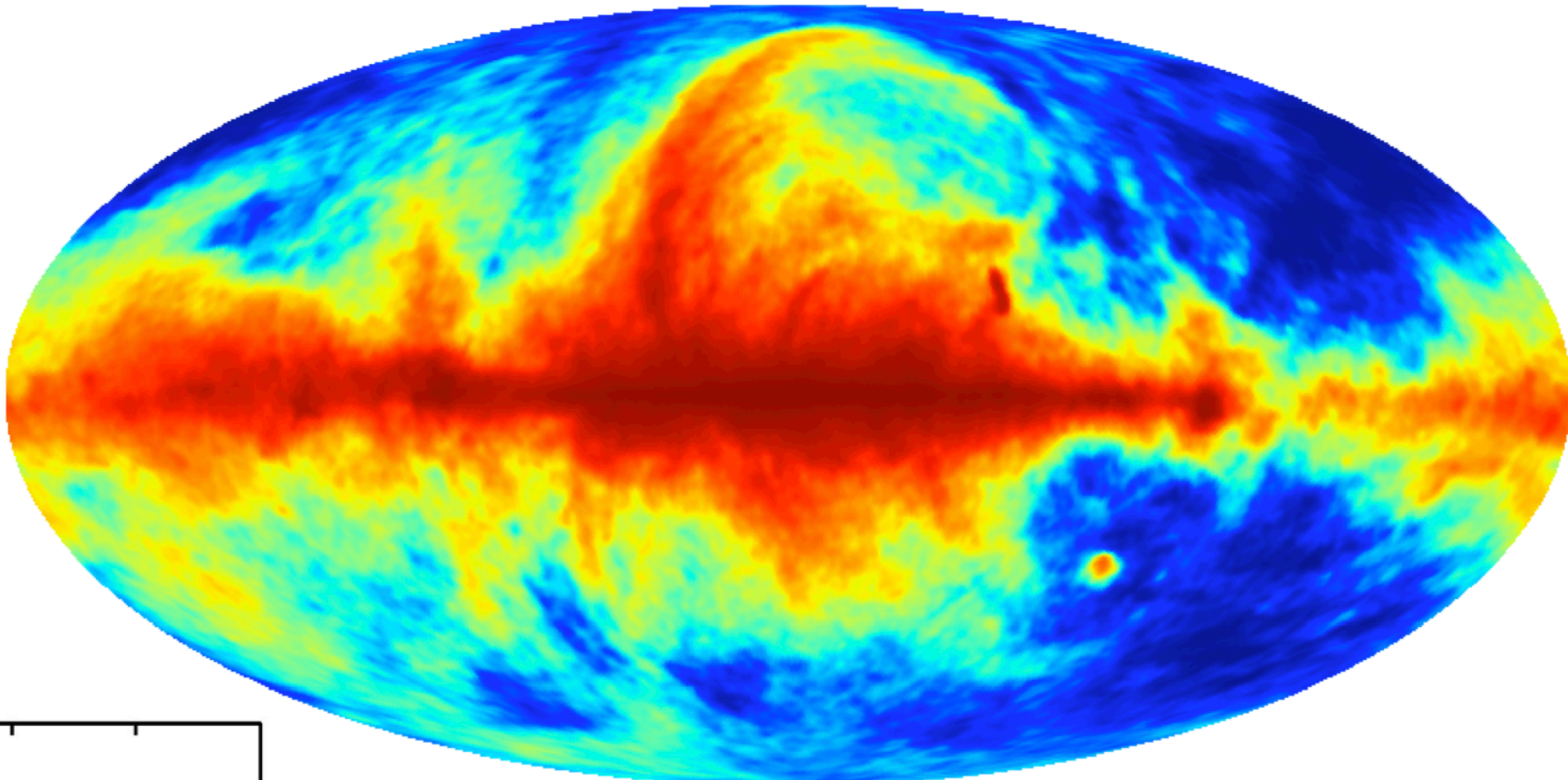


why sparsity?

mixtures are less sparse than components

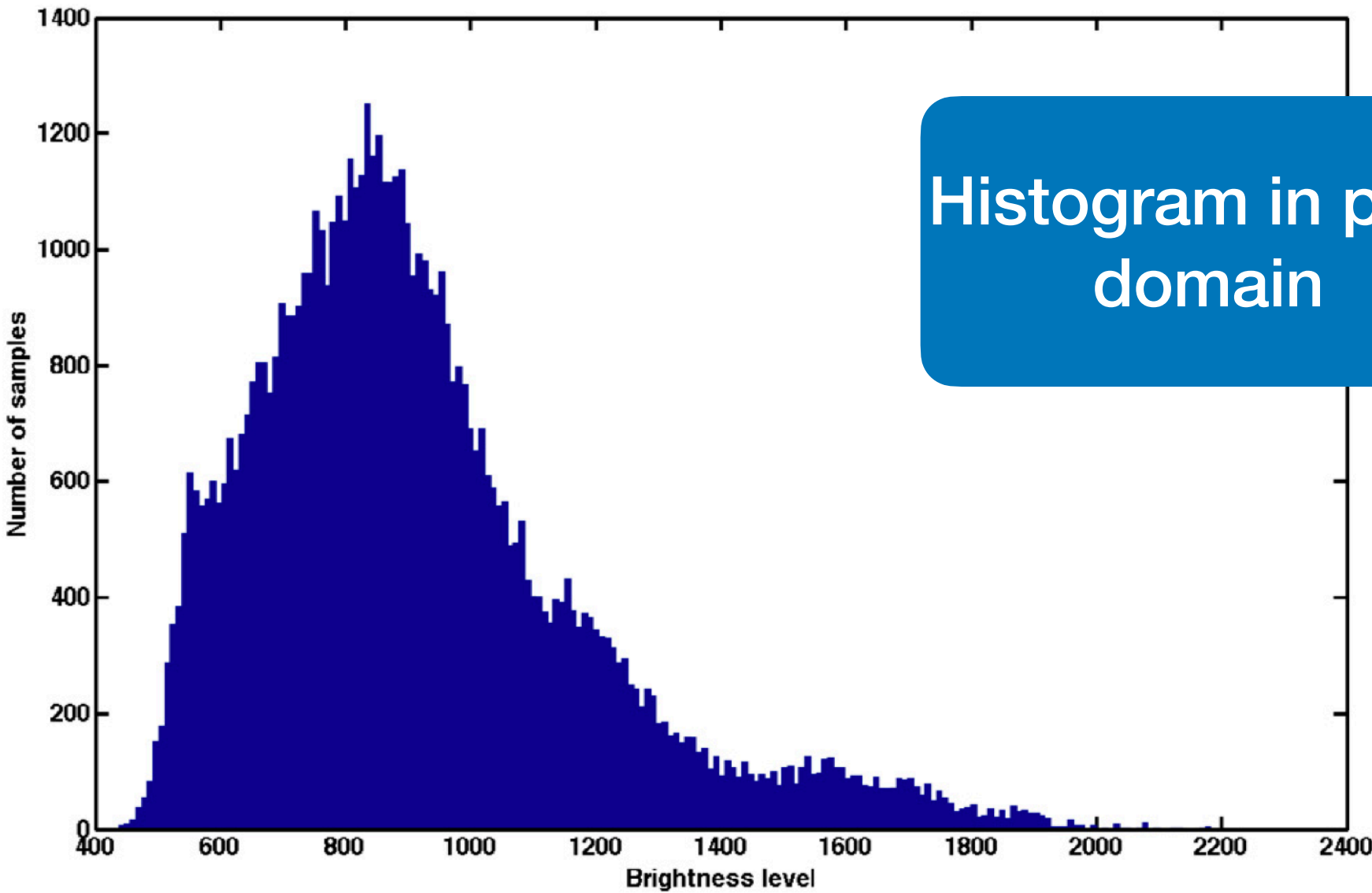
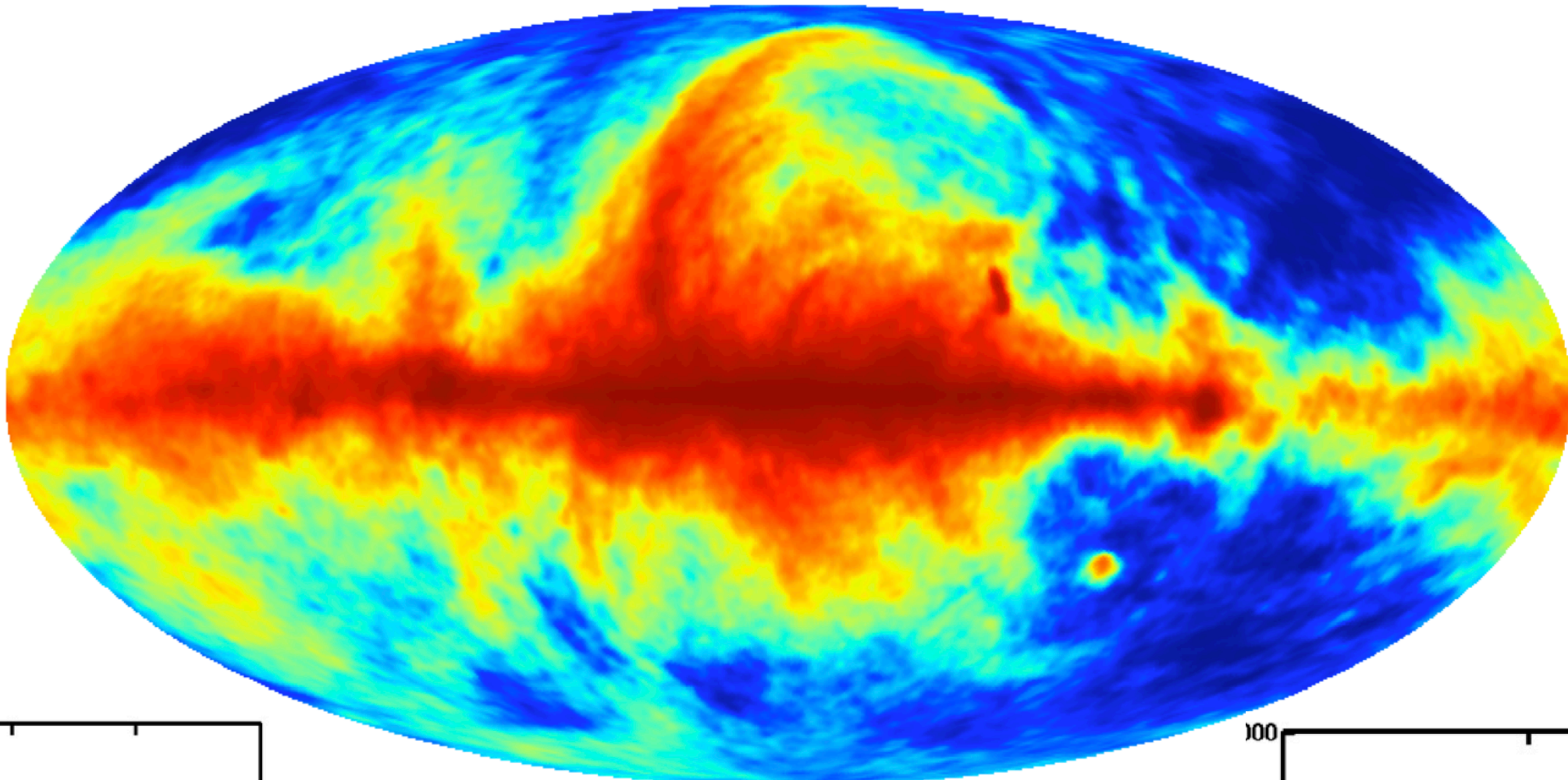


Enforcing sparsity: in which domain?

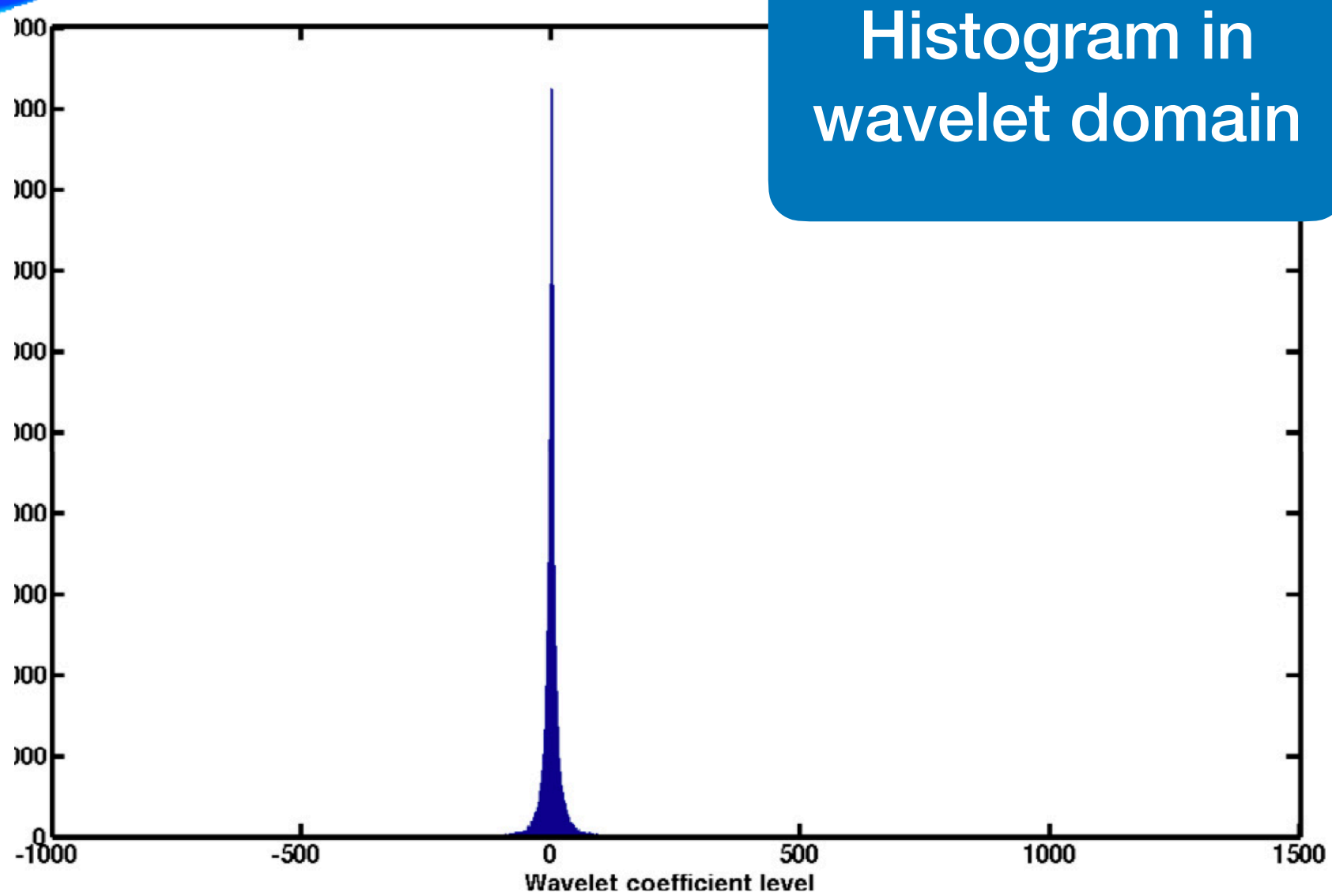


Histogram in pixel domain

Enforcing sparsity: in which domain?



Histogram in pixel domain



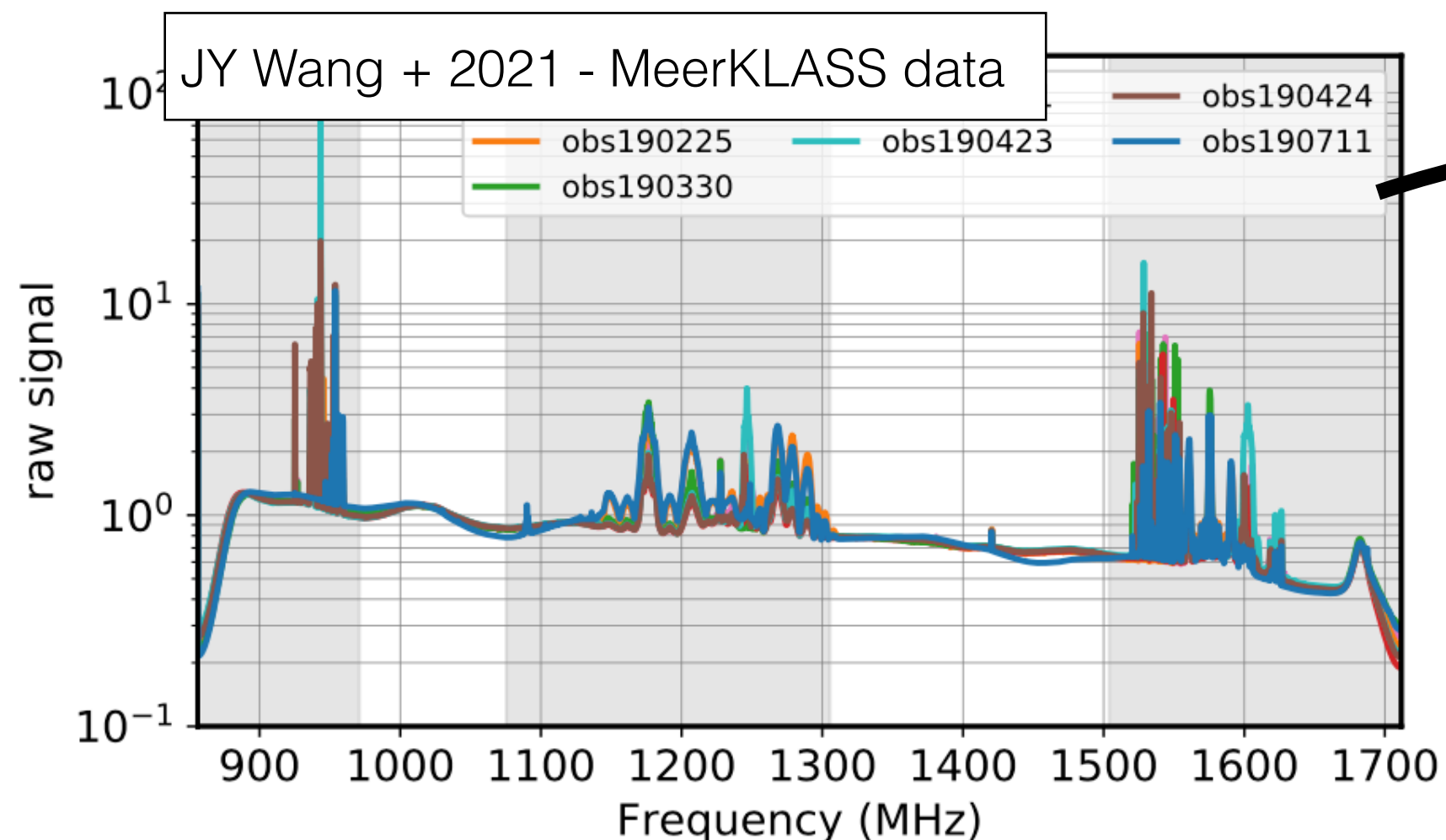
Histogram in wavelet domain

Sparsity-based component-separation for 21-cm IM

GMCA: Generalised Morphological Component Analysis

Bobin+ 2007, 2008, 2012,.. Applied on data in different astro-context: CMB (e.g. Bobin+2016), EoR (e.g. Hothi+2020), X-ray (Picquenot+2019), ...

- wavelet decomposition → **multi-scale** approach
- **No priors on signal**

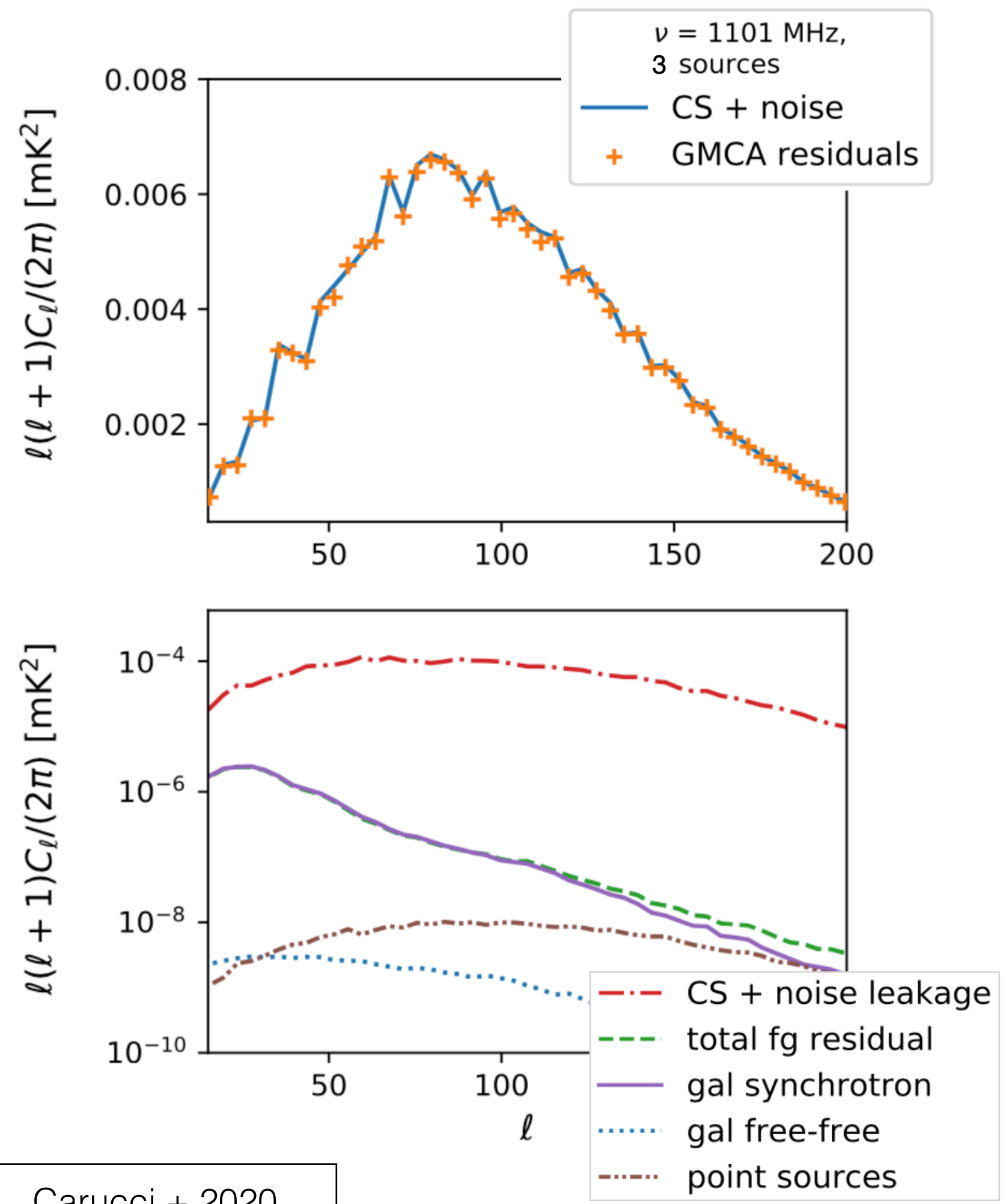


in Carucci+ 2020,
for the first time in the literature:

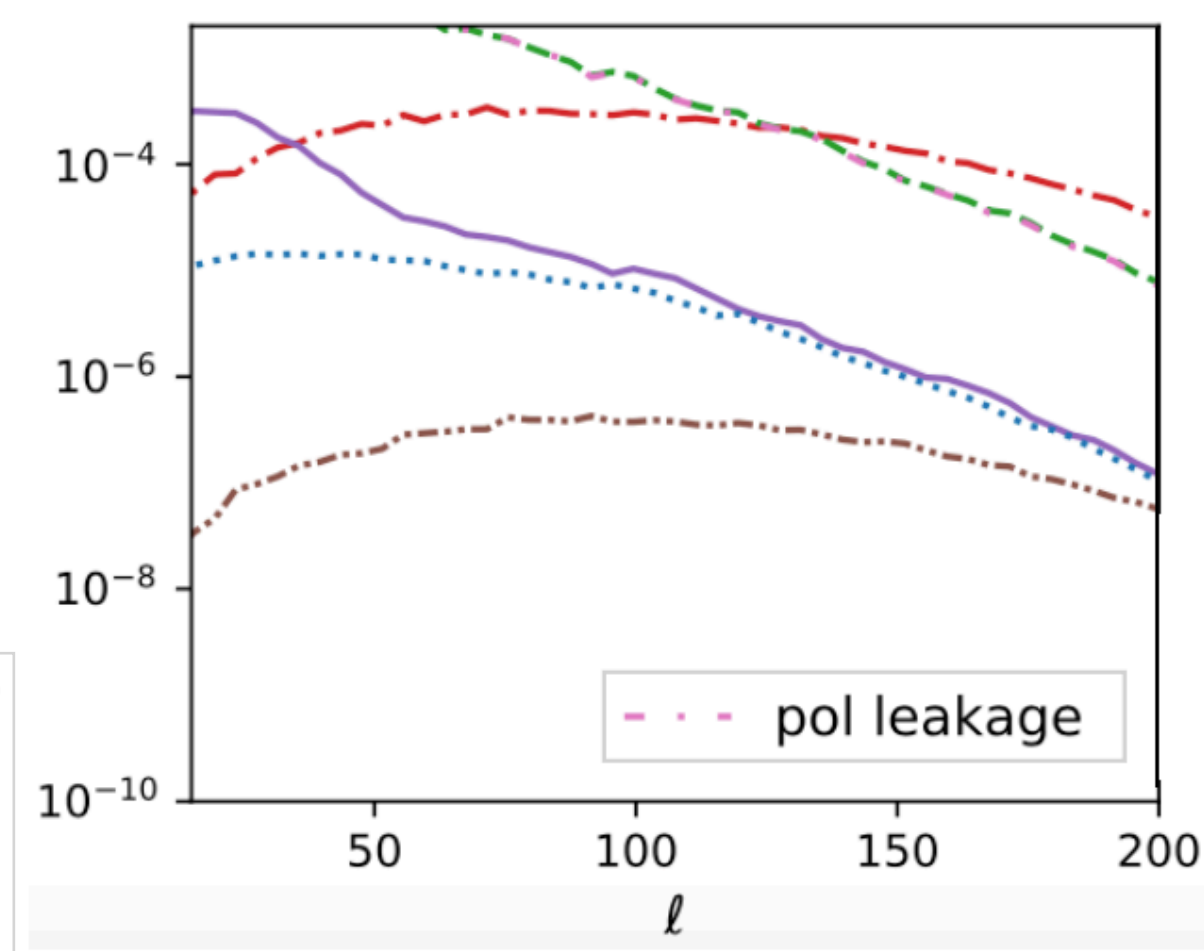
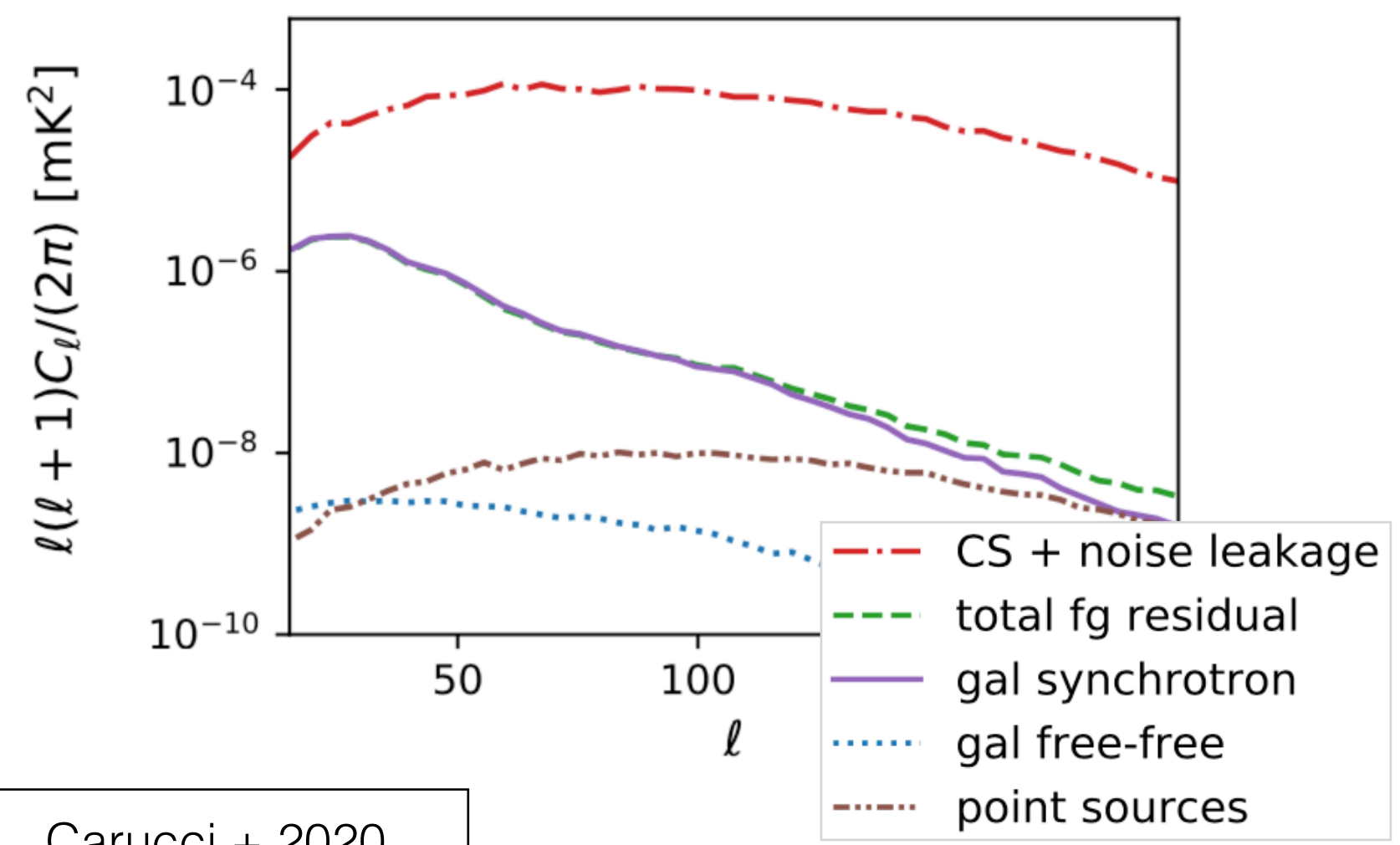
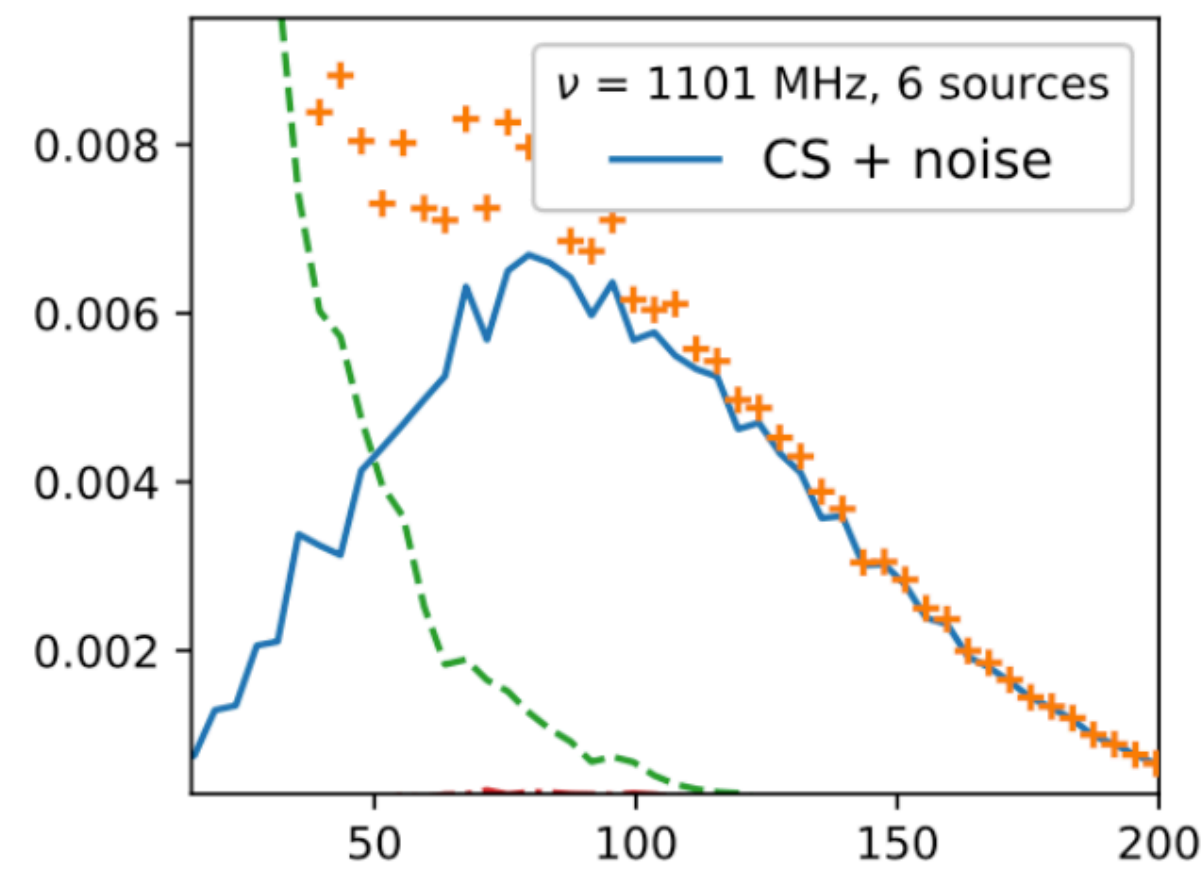
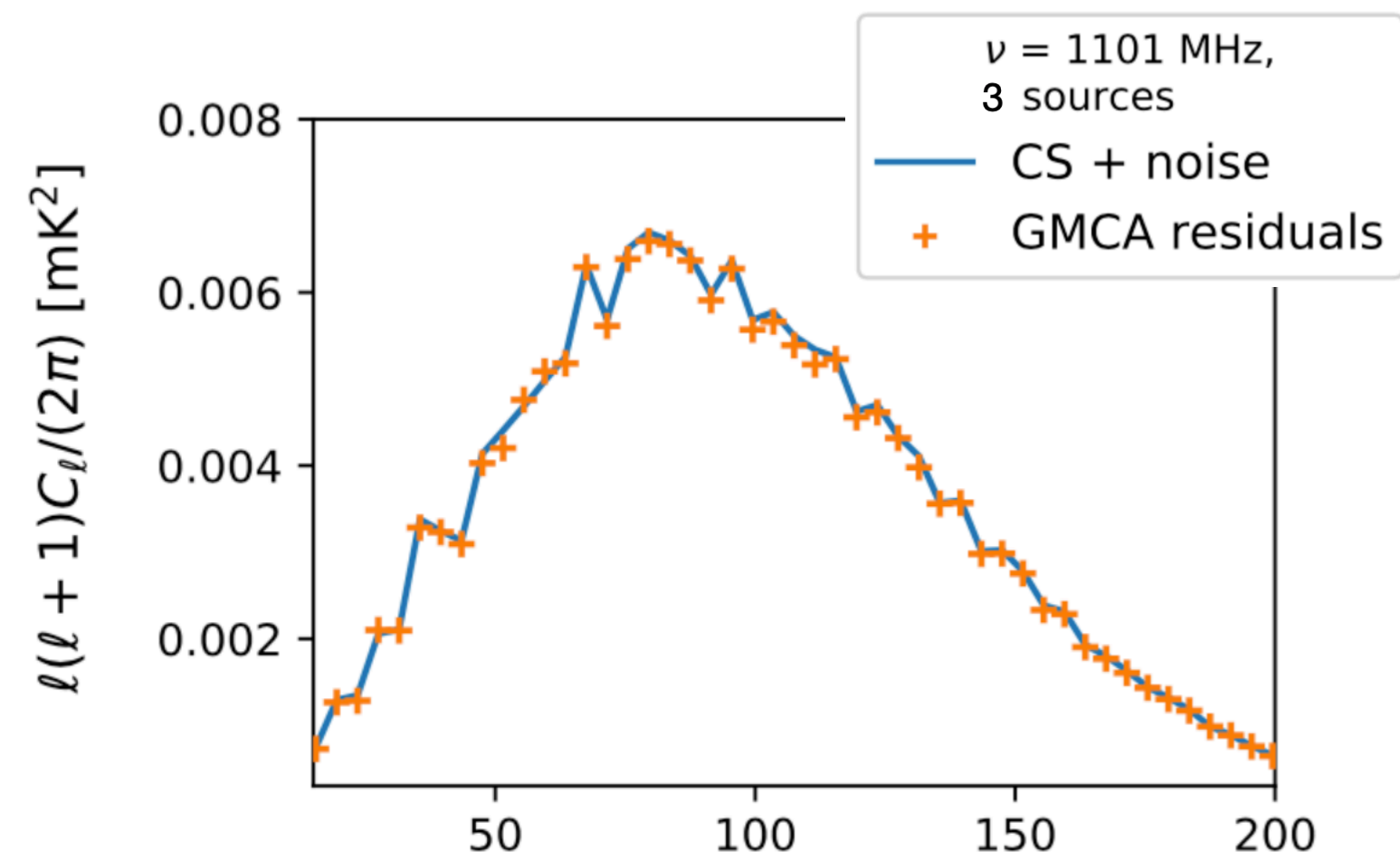
1. Good performance also with **RFI-flagged** data cubes!
(TV stations, telecommunication, satellites,..)
2. **Pol leakage:** greater complexity of data
(higher number of sources needed, convergence not assured, mode-mixing assured)

To reproduce these results:
codes and sims available online

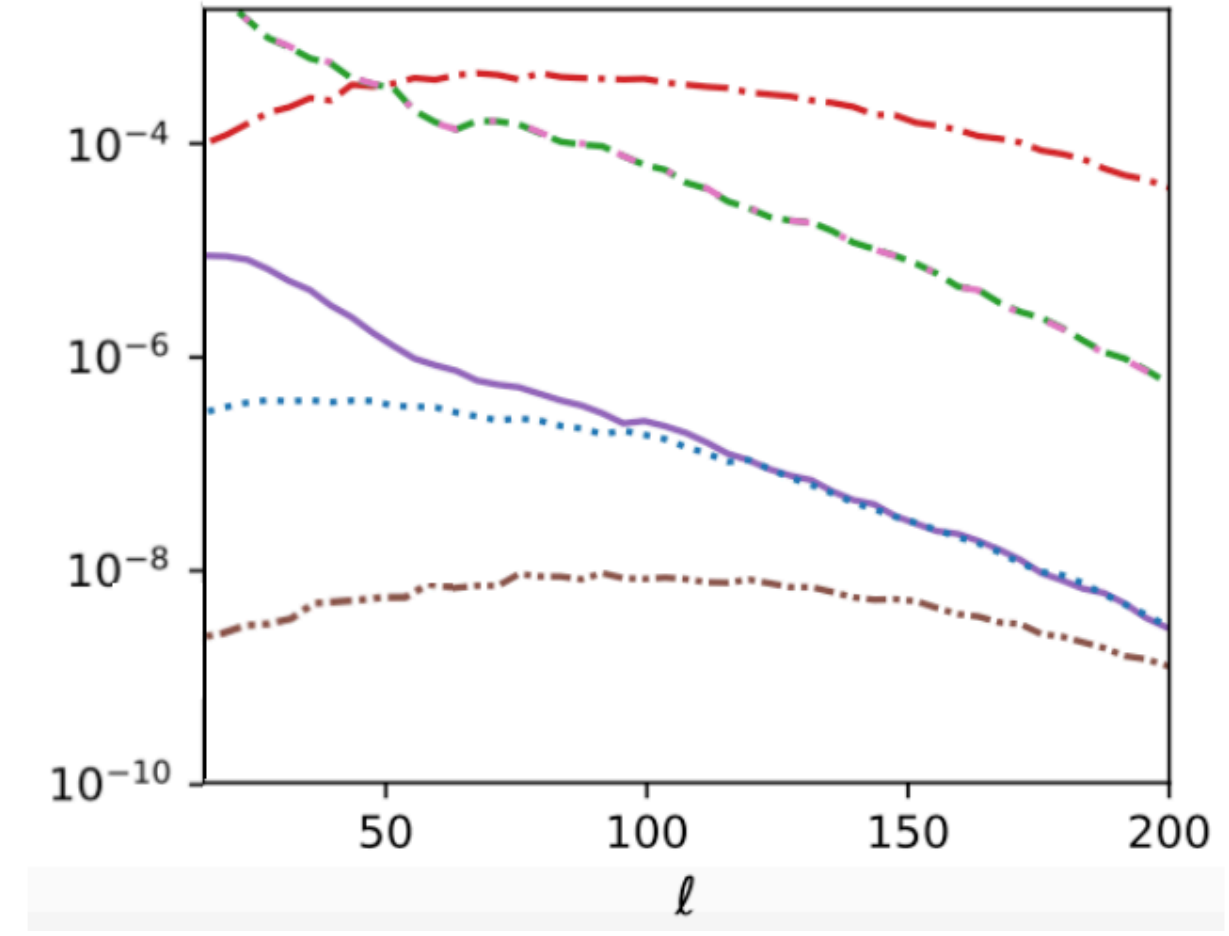
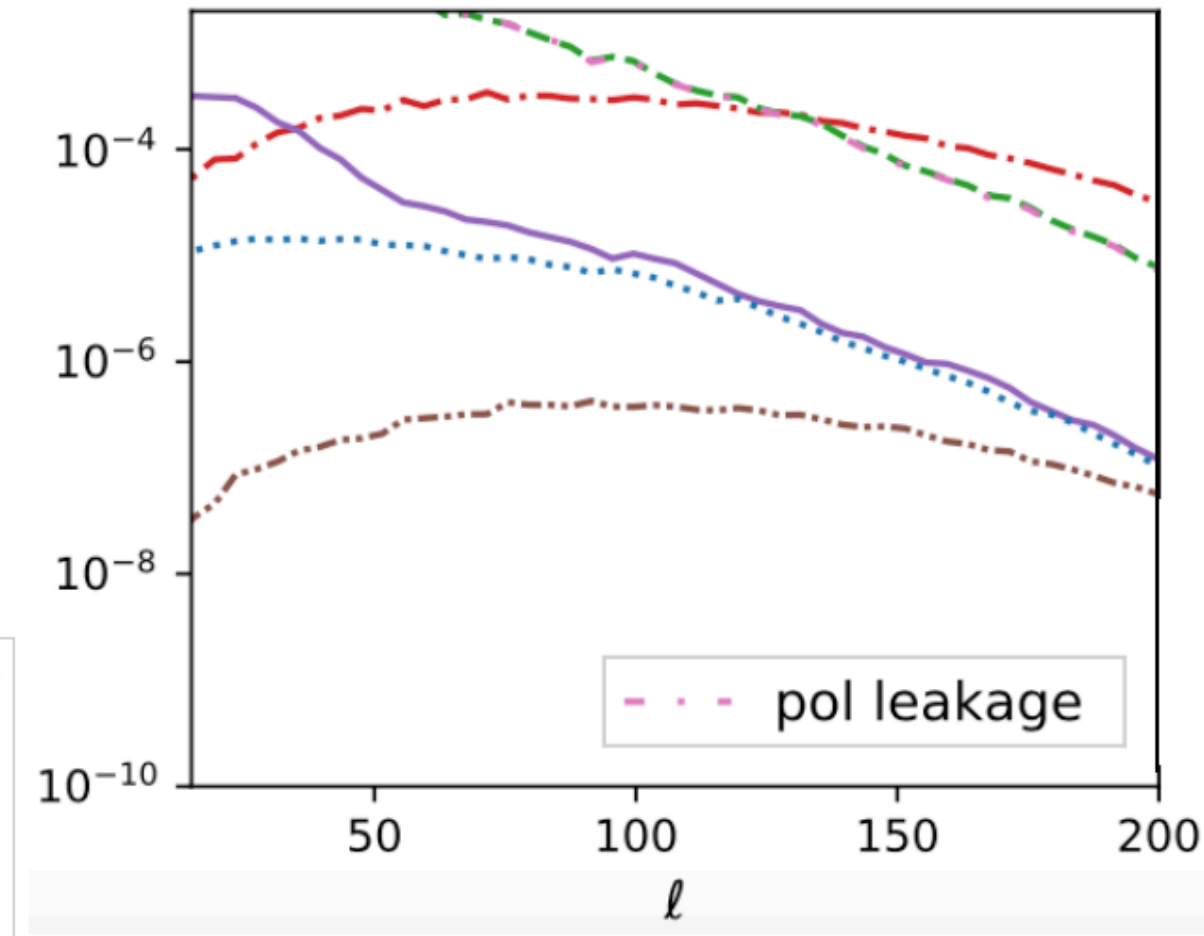
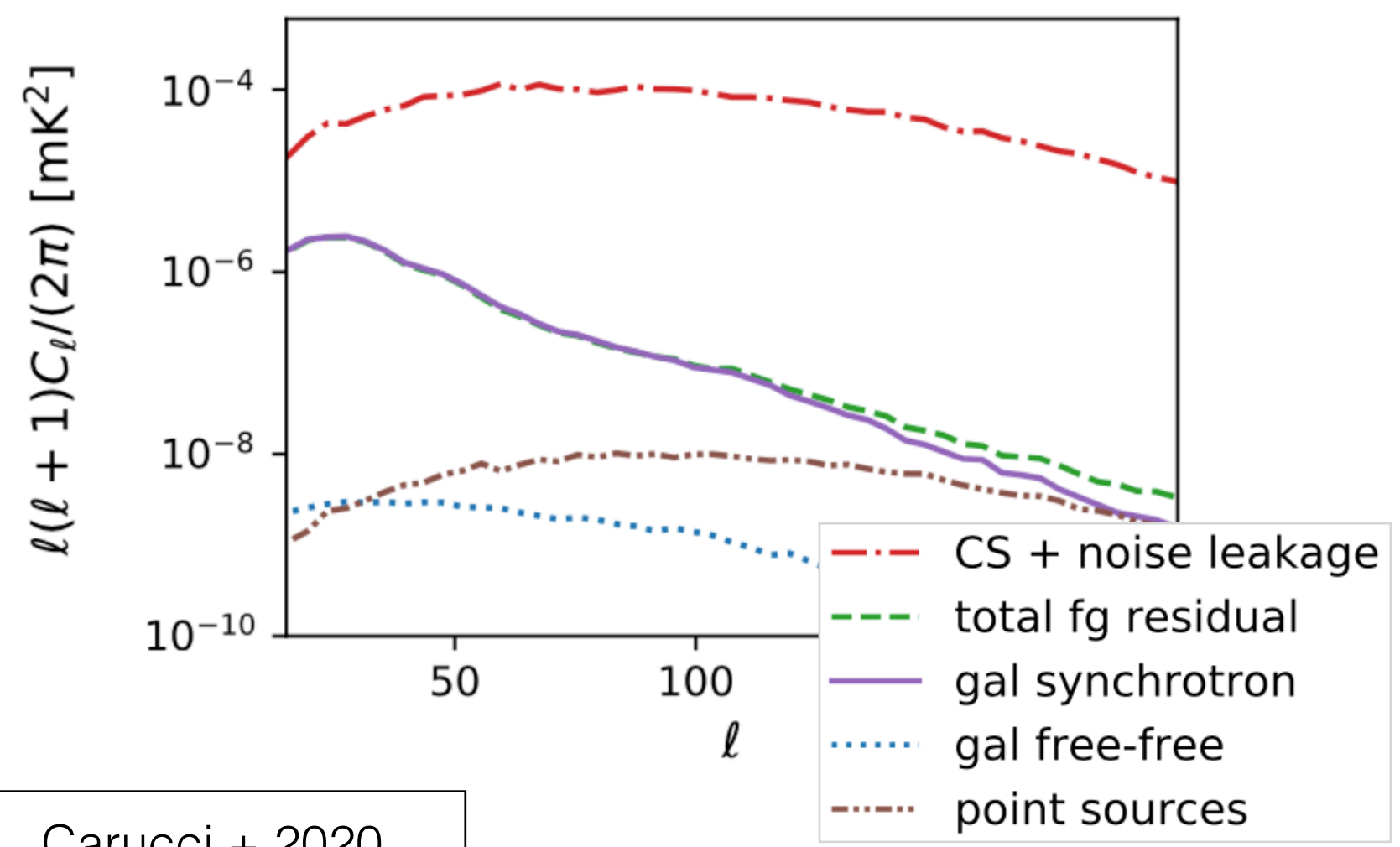
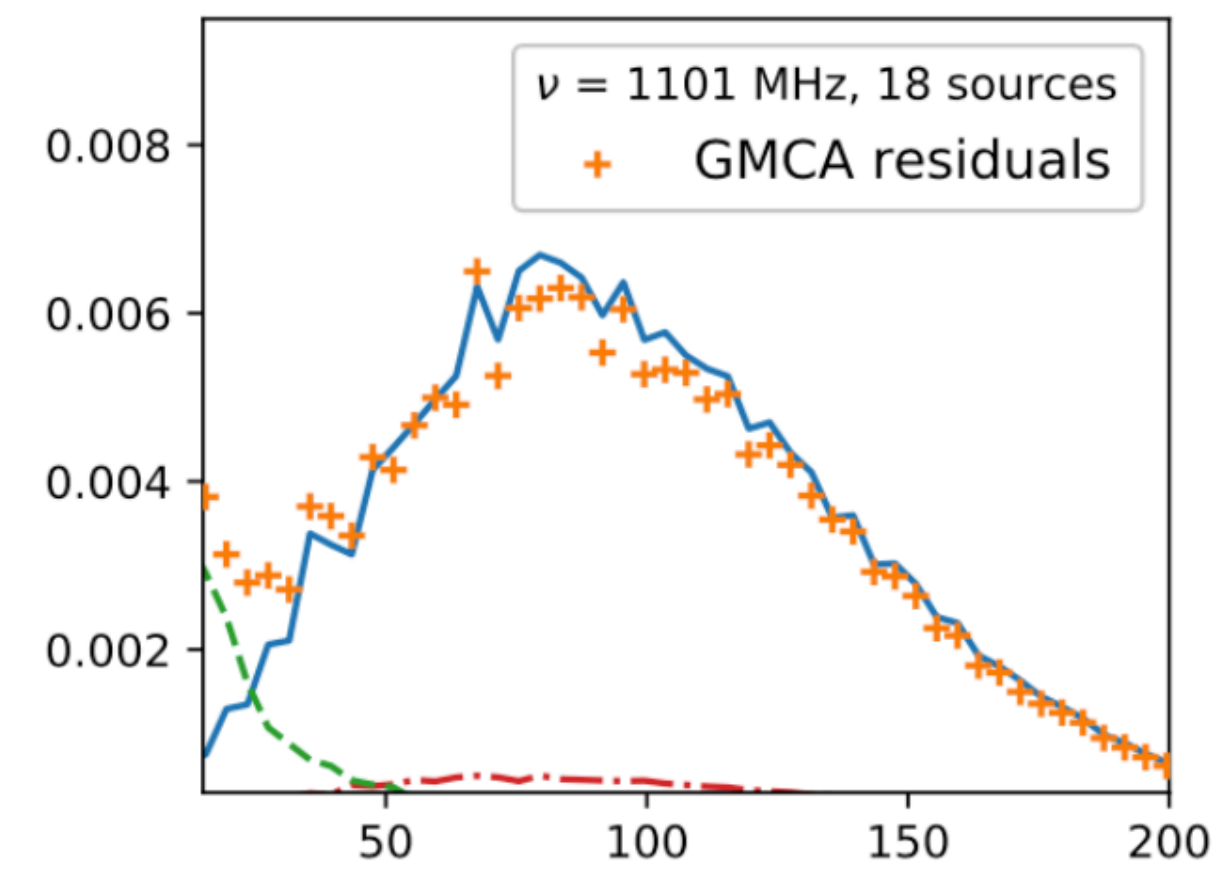
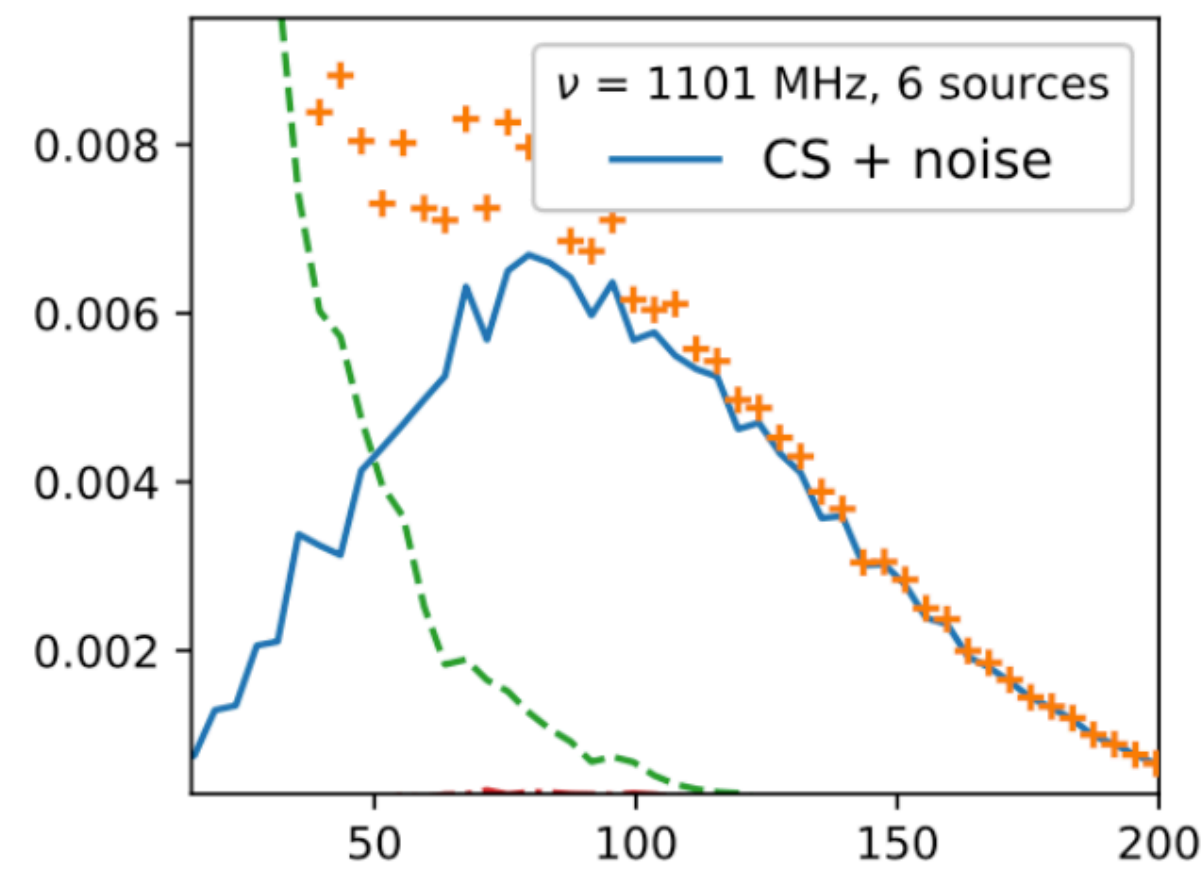
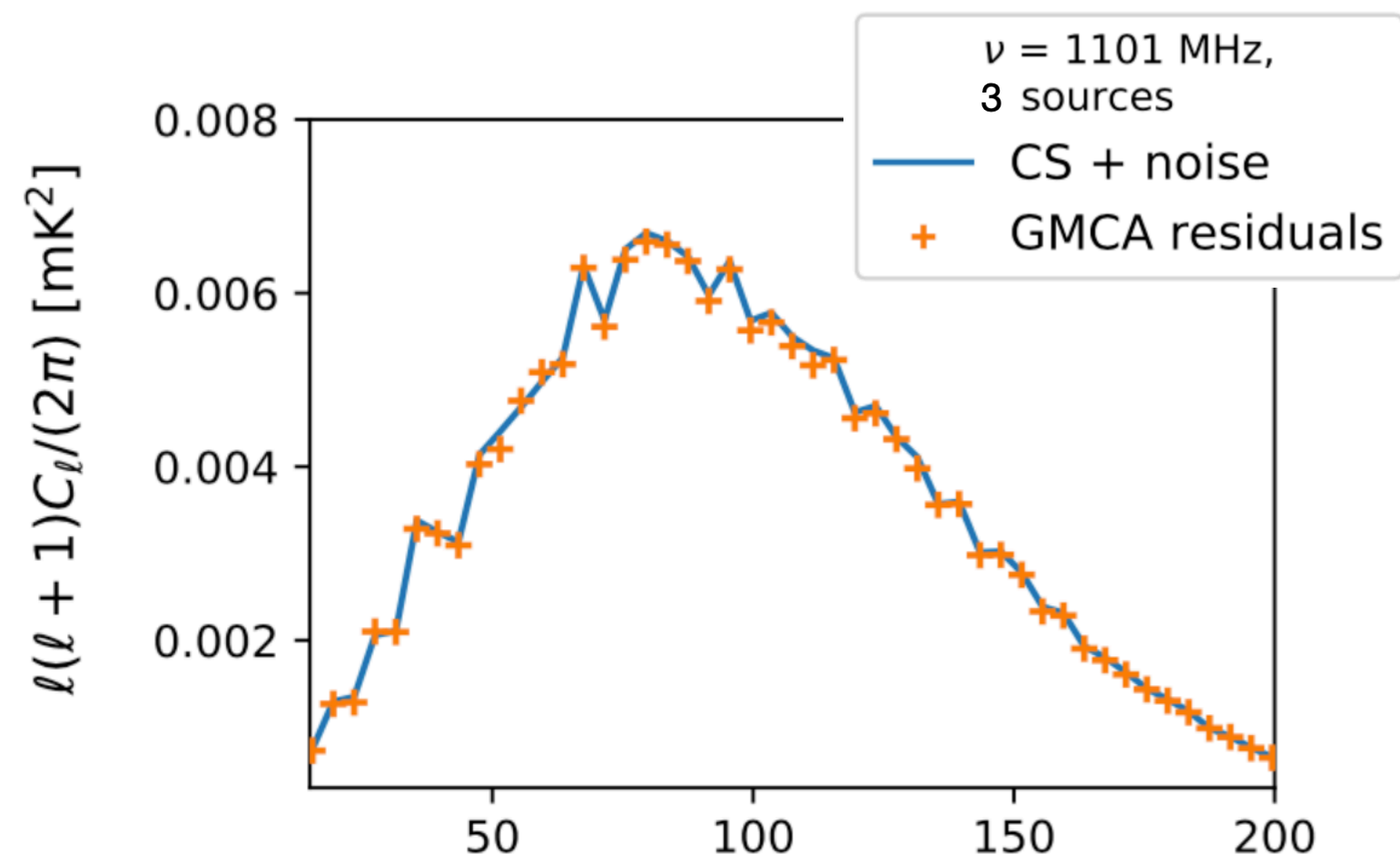
a quick interlude on mixGMCA



Carucci + 2020

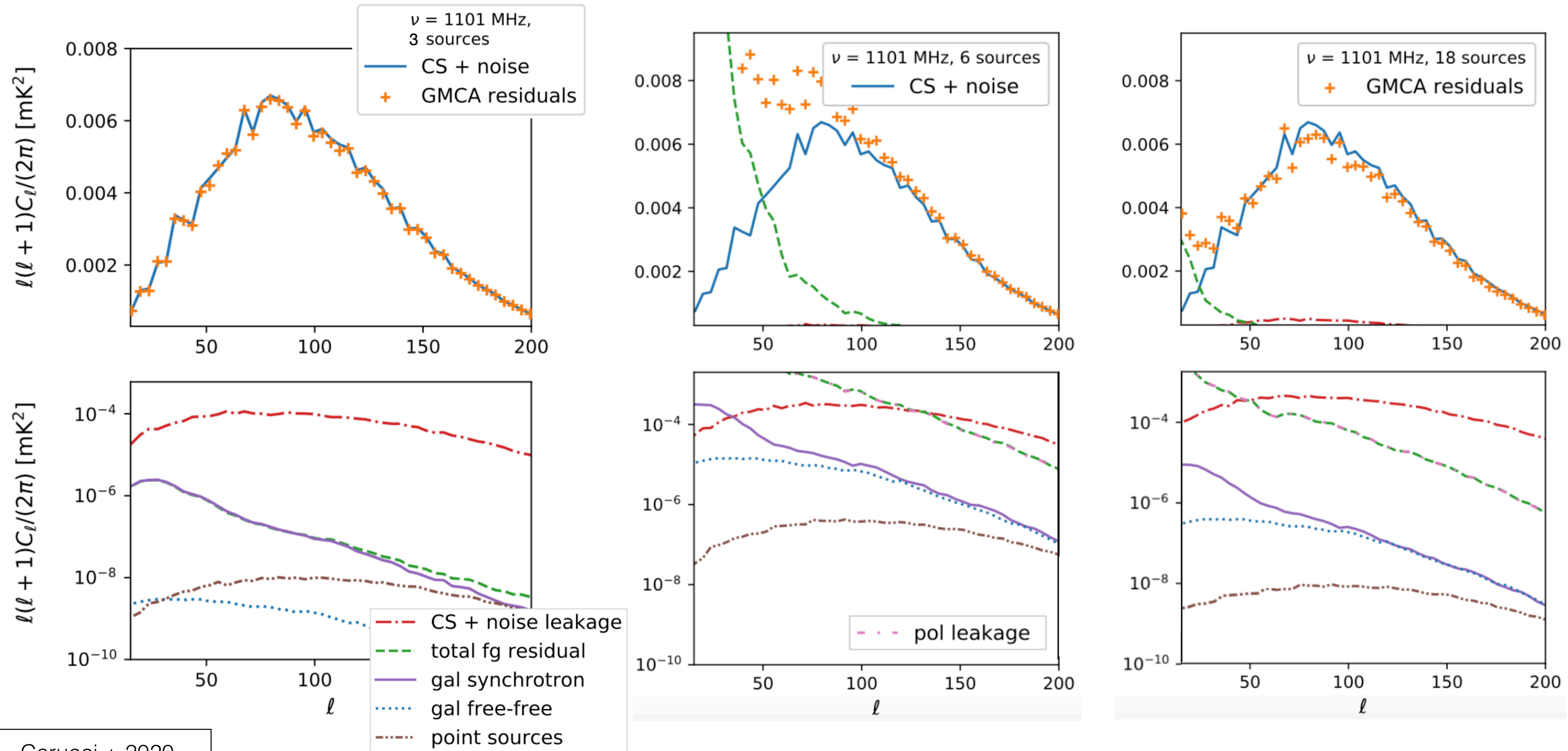


Carucci + 2020



Carucci + 2020

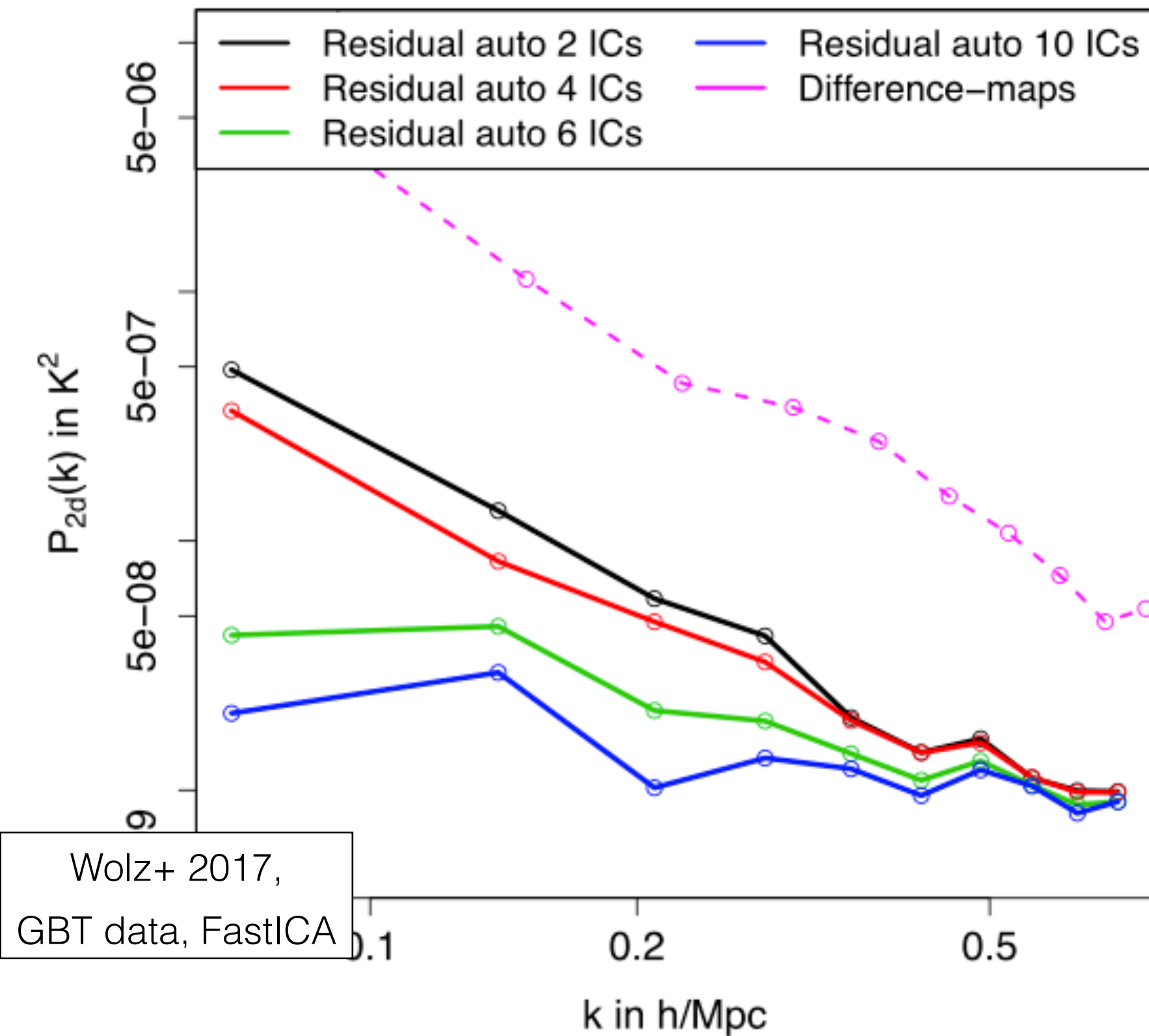
Different scales need different care



Carucci + 2020

Different scales need different care

The wavelet domain is a multi-scale framework!



- GMCA performs very well on small scales
- PCA / ICA \rightarrow overfit the large scales

PCA on the large scale
+
GMCA on the small scales

mixGMCA

See also Hothi+2020 with LOFAR data

MeerKAT



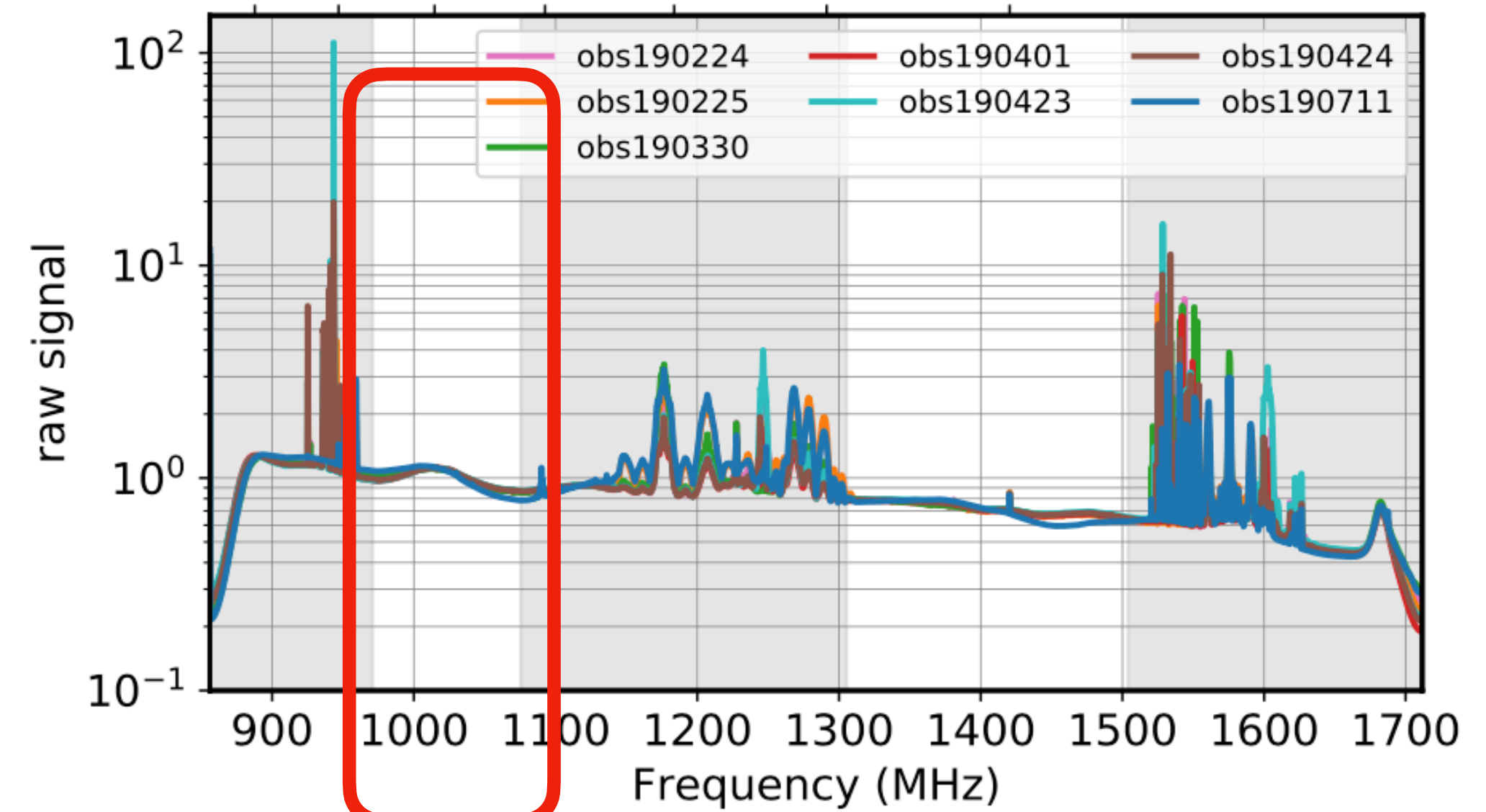
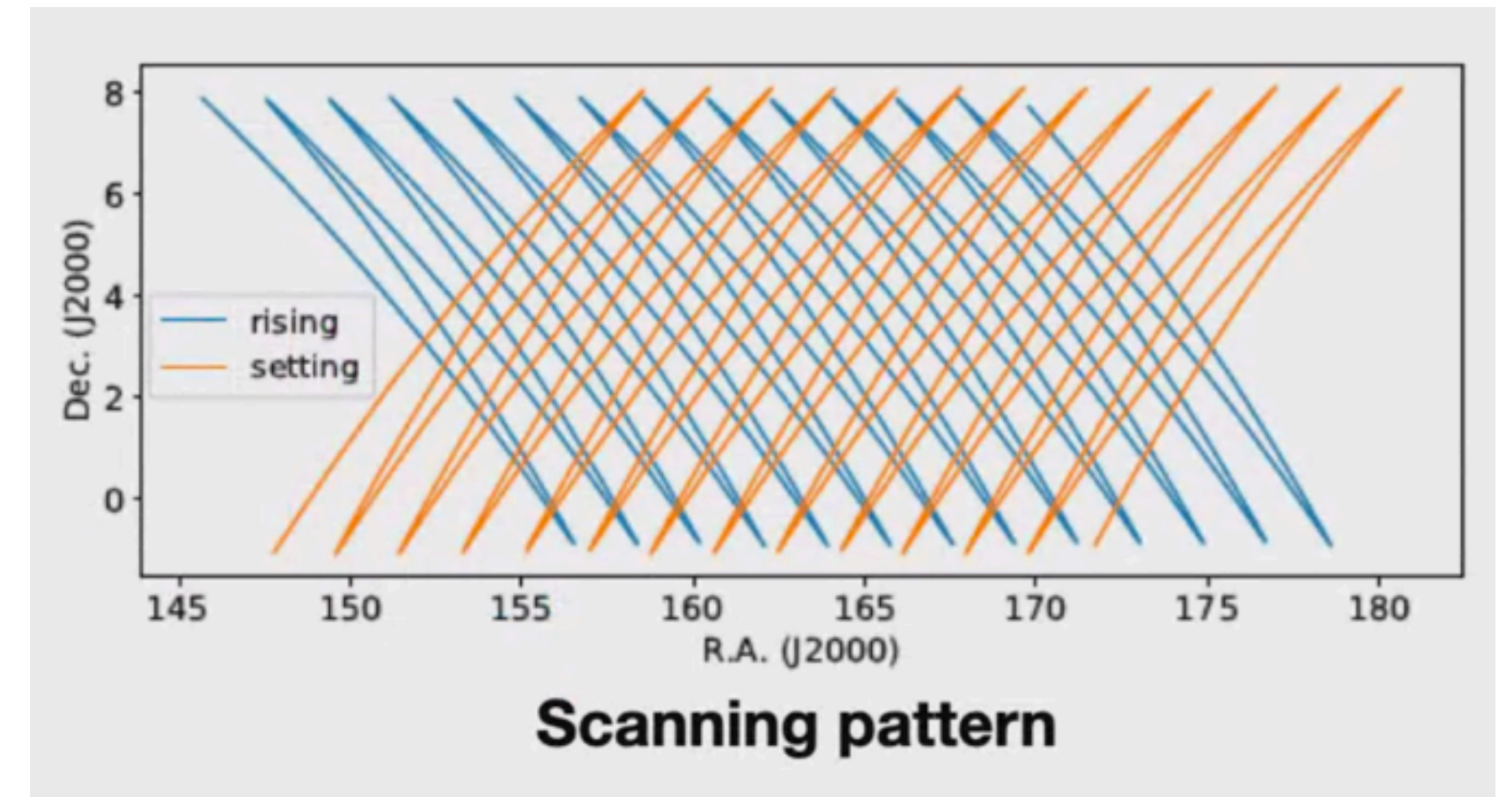
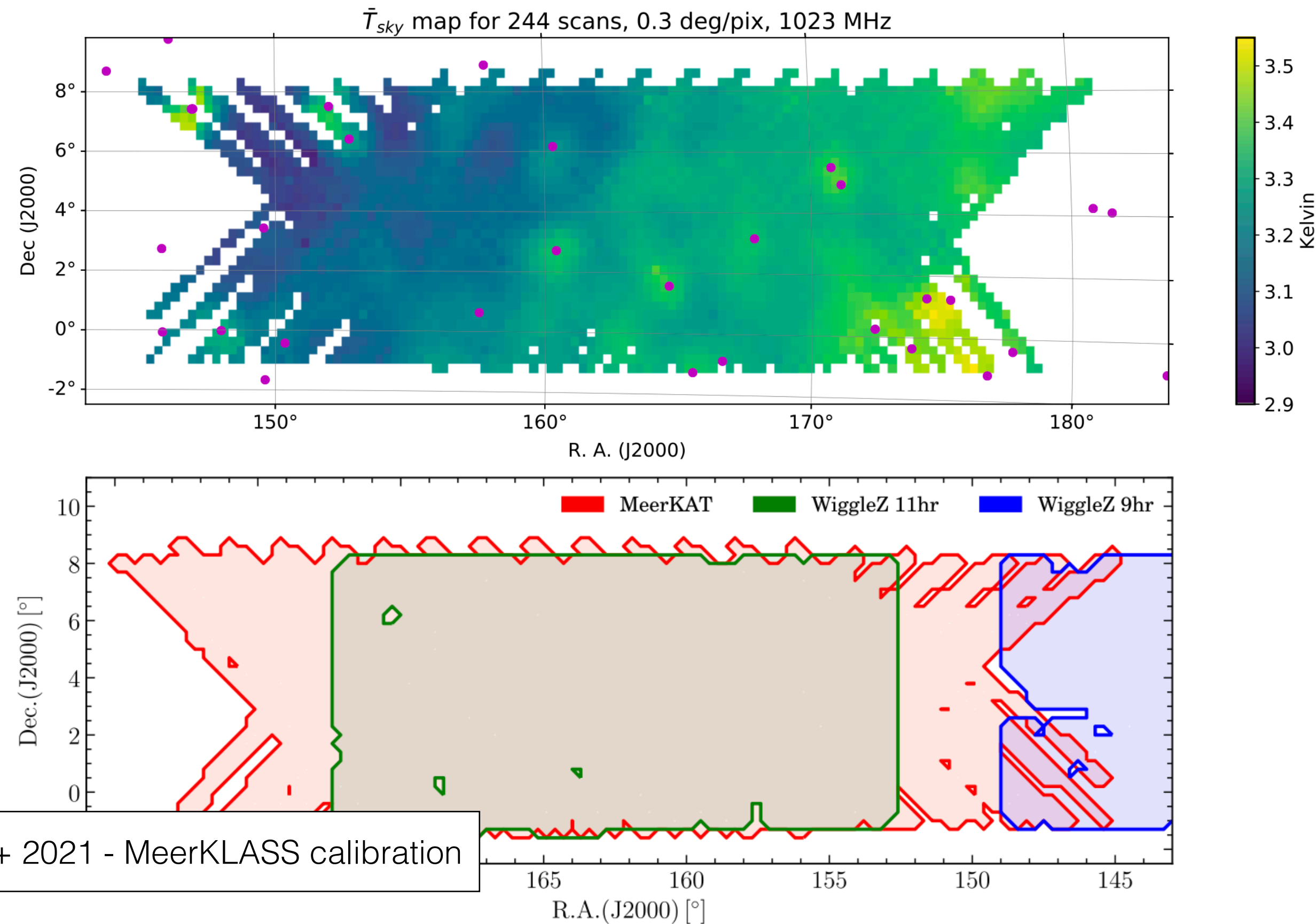
MeerKLASS: MeerK_{AT} Large Area Synoptic Survey

ArXiv: 1709:06099

**Alkistis Pourtsidou, Amadeus Wild, Brandon Engelbrecht, Isabella Carucci,
Jingying Wang, Keith Grainge, Laura Wolz, Mario Santos, Marta Spinelli, Mel Irfan,
Phil Bull, Stefano Camera, Steve Cunnington, Tamirat Gebeyehu, Zé Fonseca, ...**

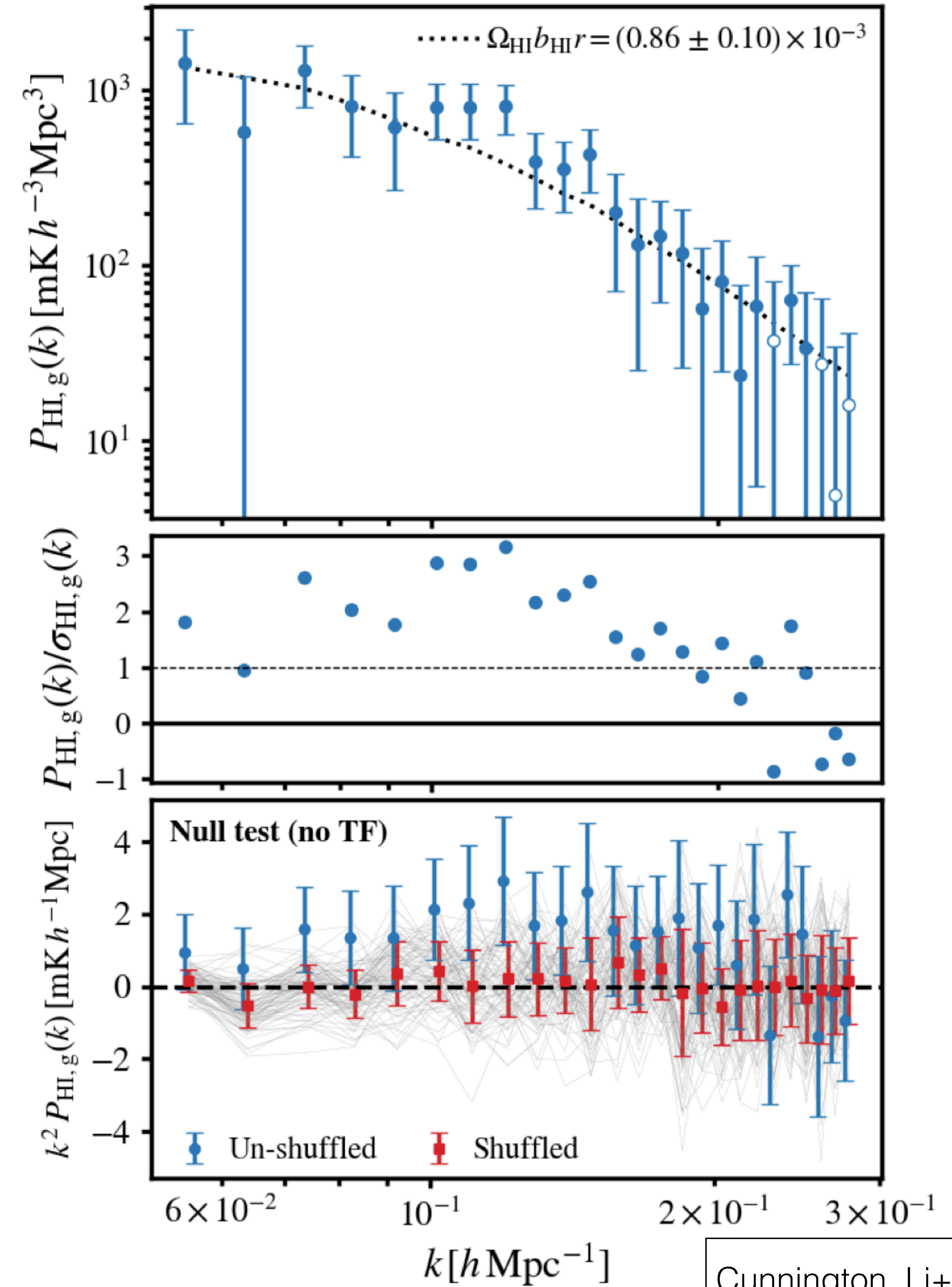
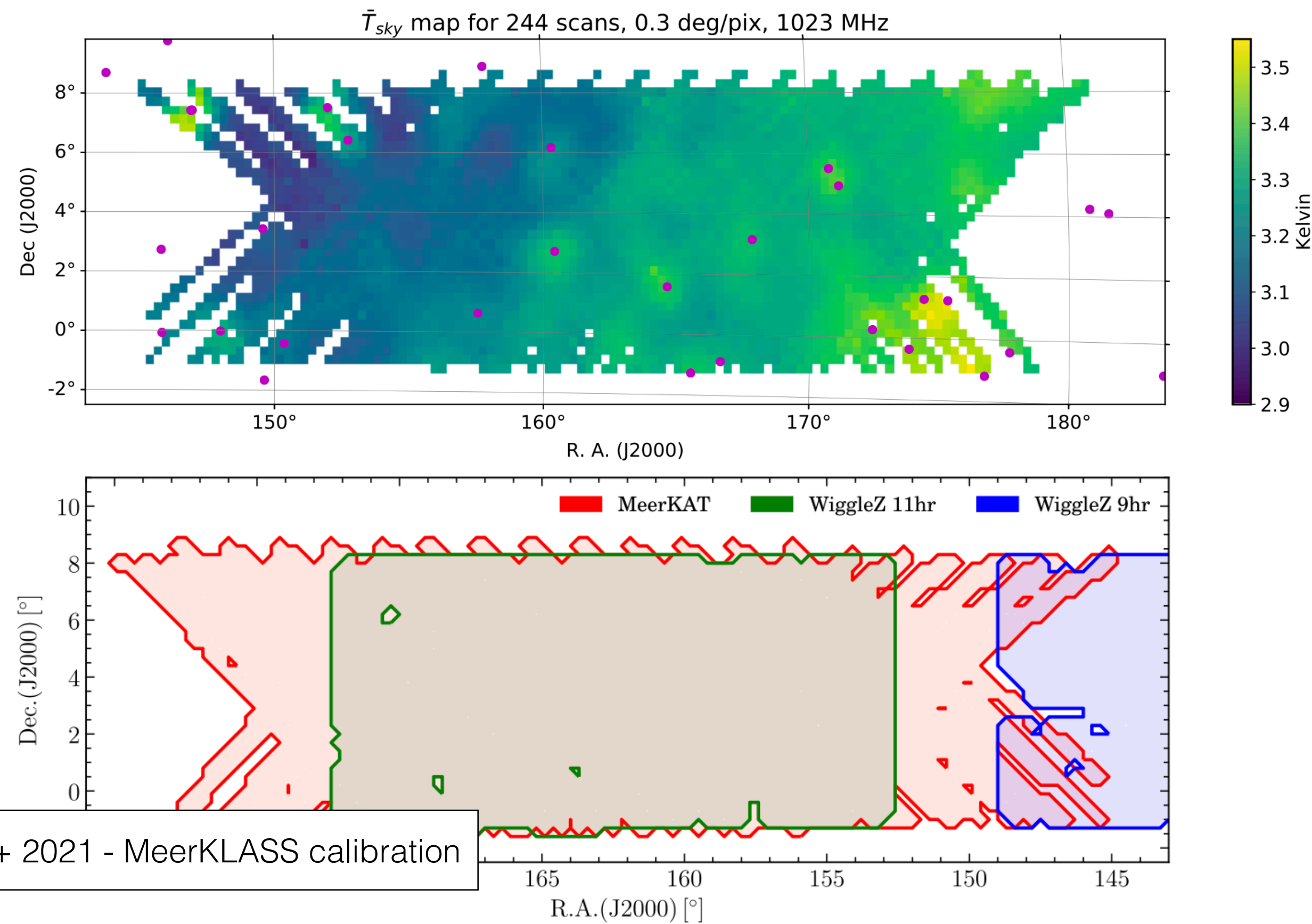
Pilot survey data (2019):

- 10.5 hours of data from six nights of observations
- Overlapping with the WiggleZ 11hr field (~200 deg²)
- We use data in range 973-1015 MHz ($0.40 < z < 0.46$)

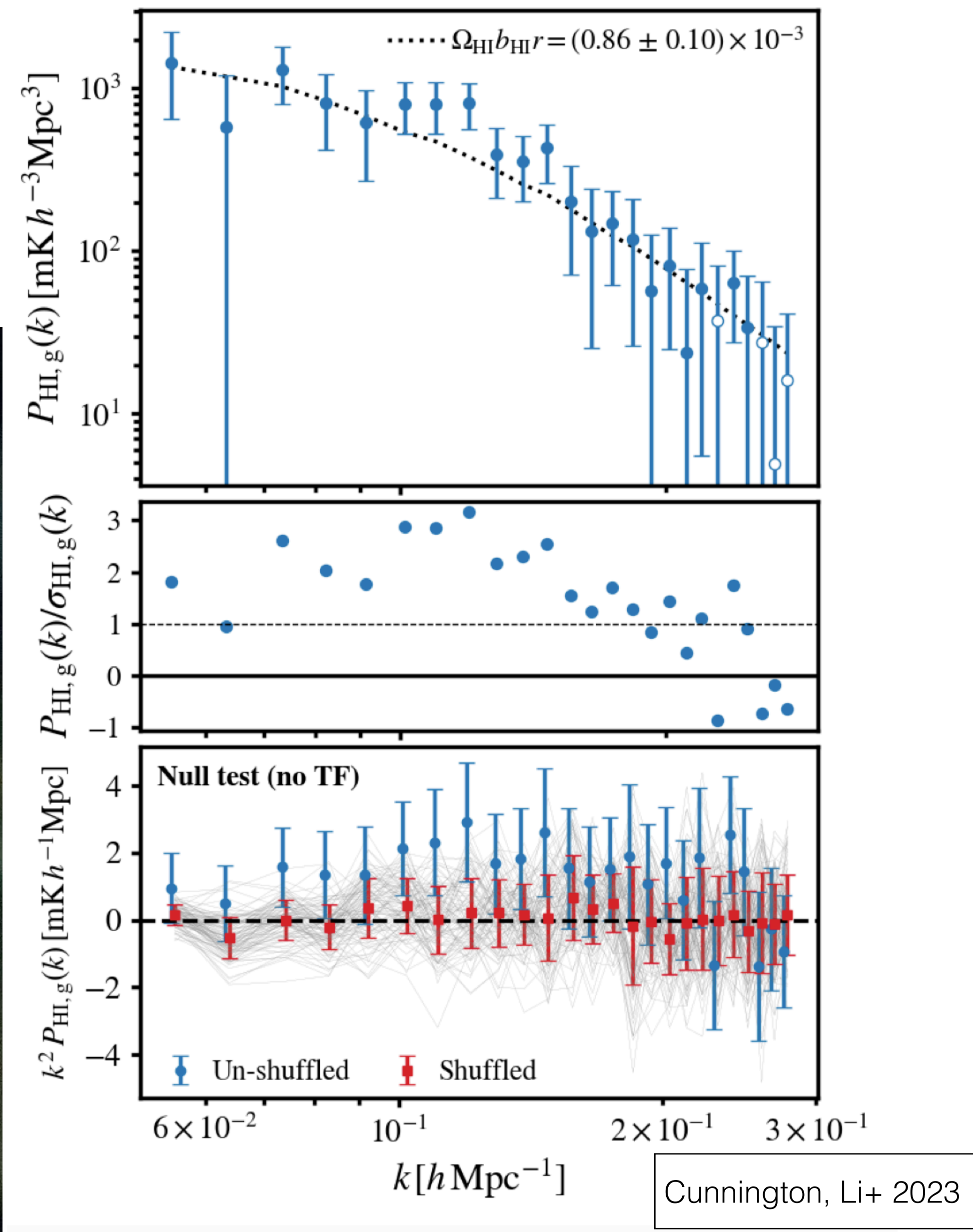


Pilot survey data (2019):

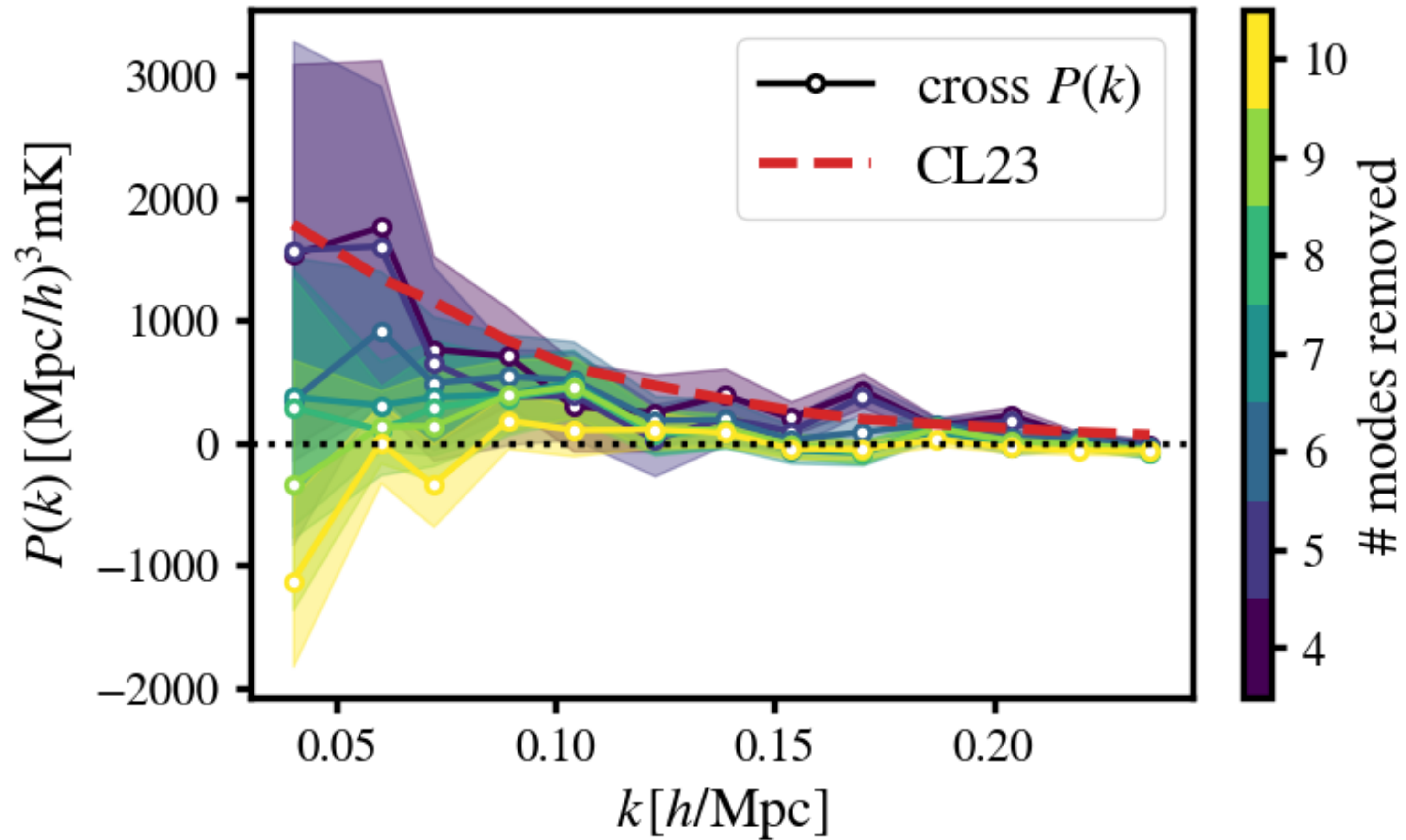
- 10.5 hours of data from six nights of observations
- Overlapping with the WiggleZ 11hr field (~200 deg²)
- We use data in range 973-1015 MHz ($0.40 < z < 0.46$)



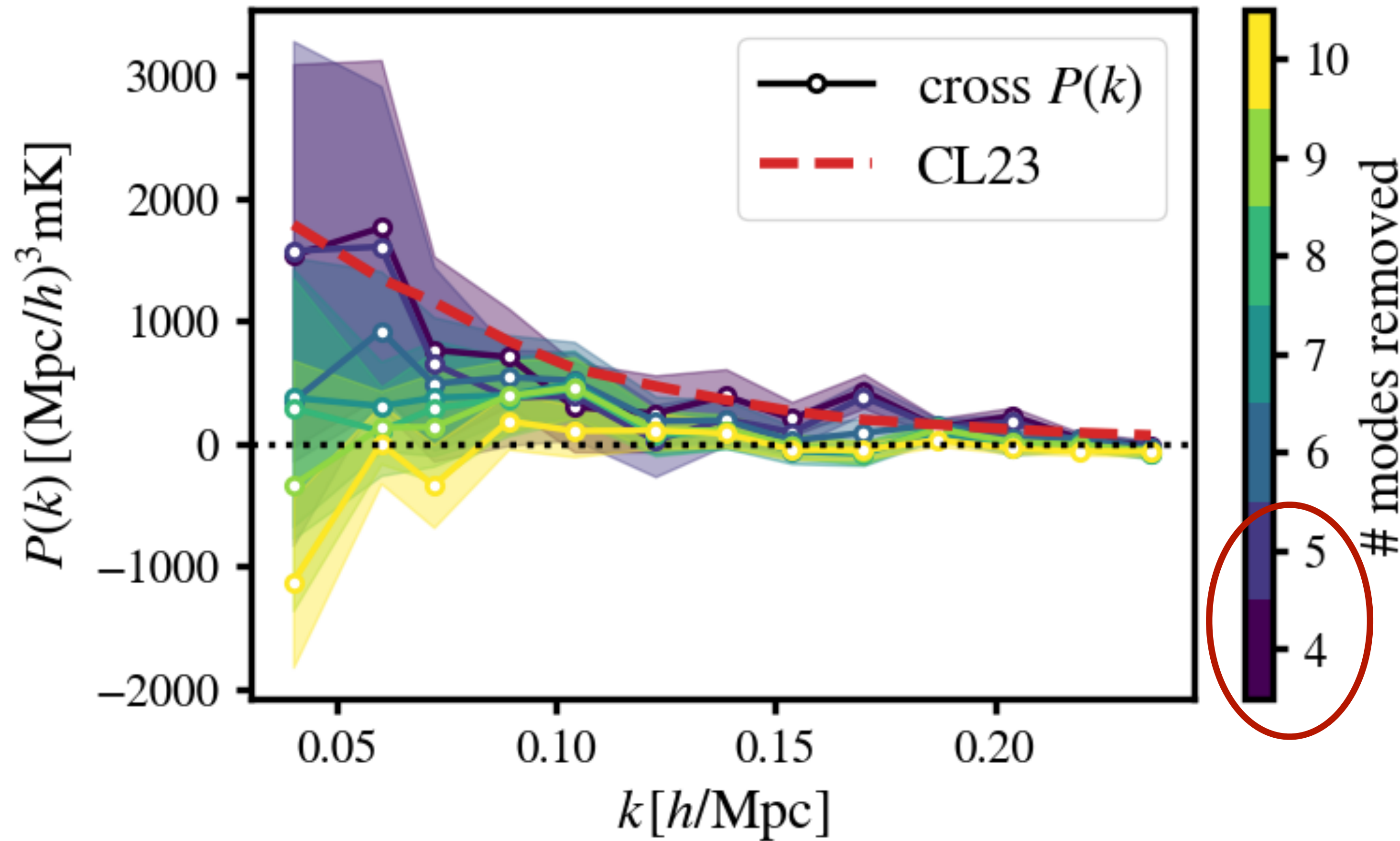
Our precious!!



Use the cross-correlation as a benchmark



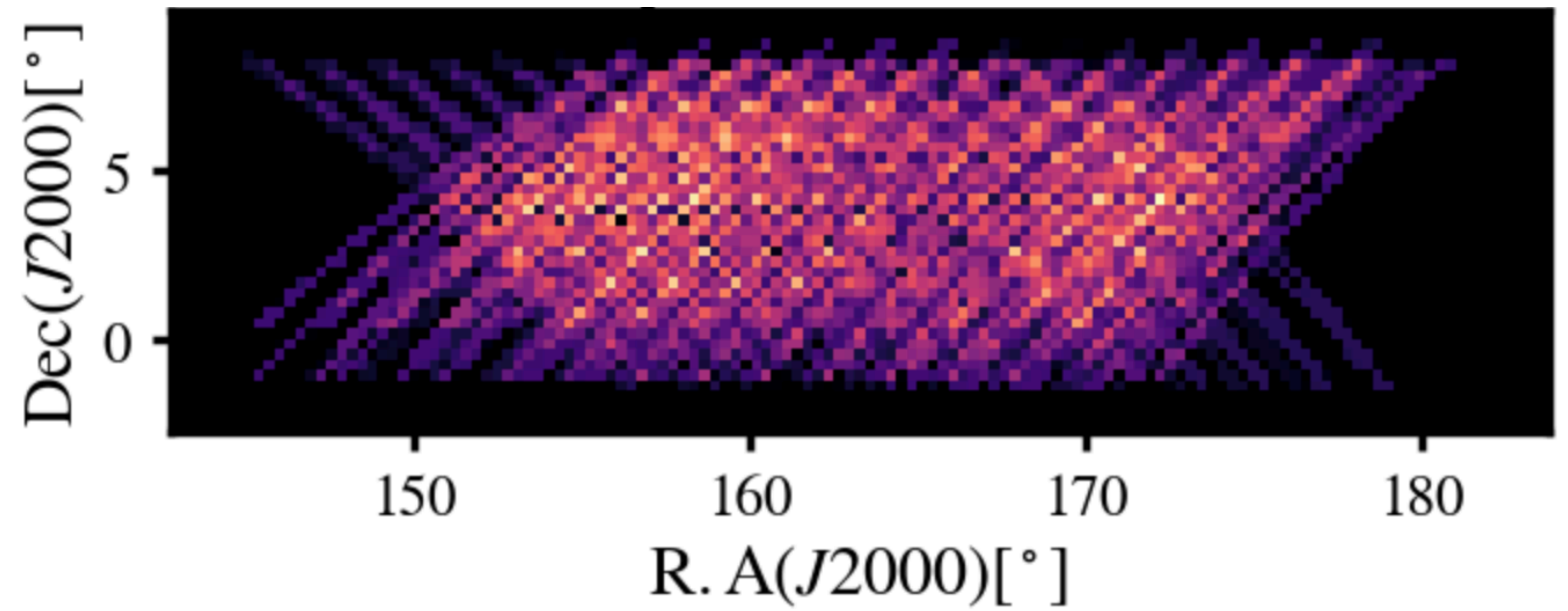
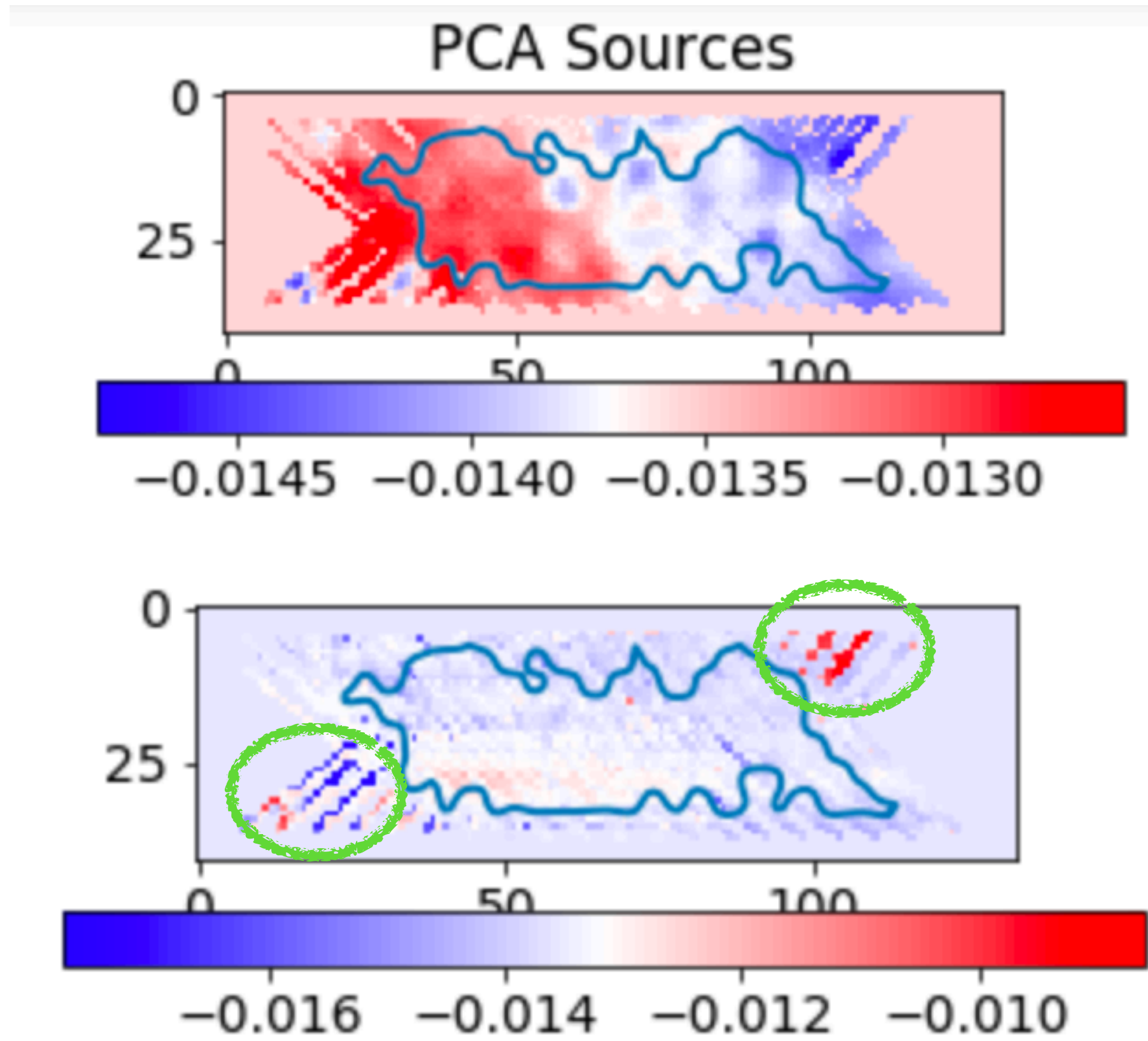
Use the cross-correlation as a benchmark



To be compared with 30 !

Re-analysis of 2019 data

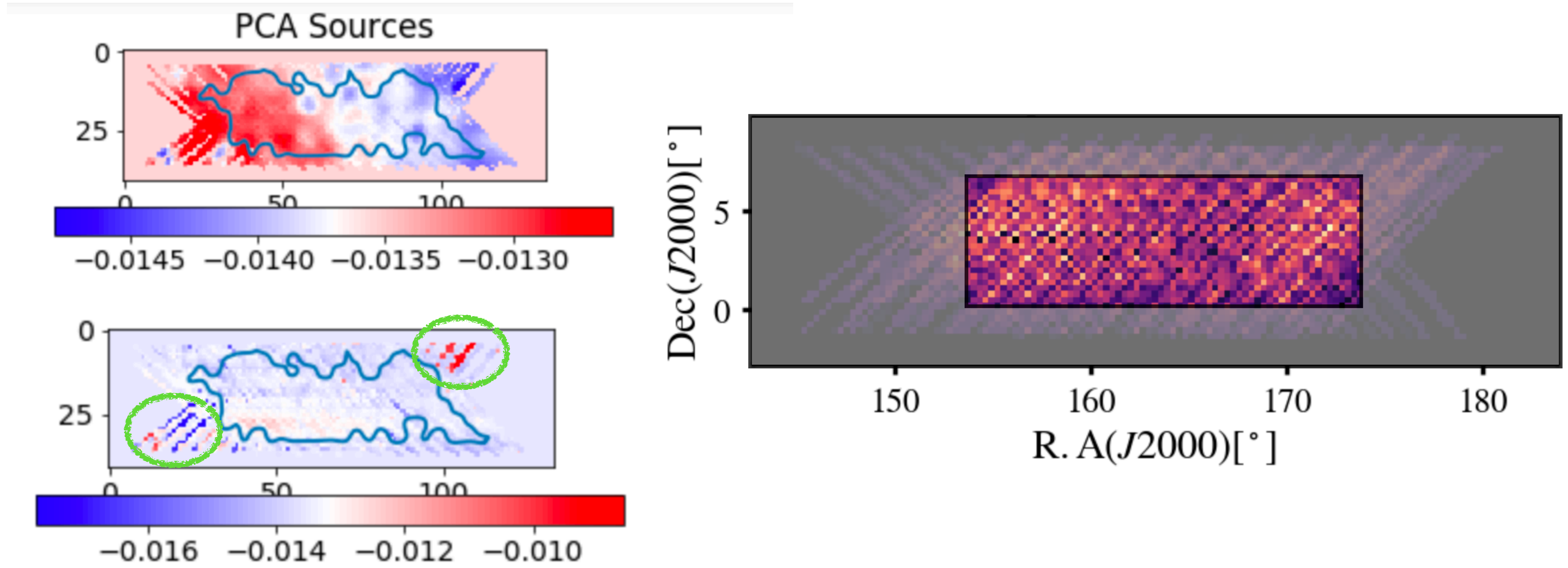
$$X = A S$$



Re-analysis of 2019 data

- 1. PCA-informed pixel flagging

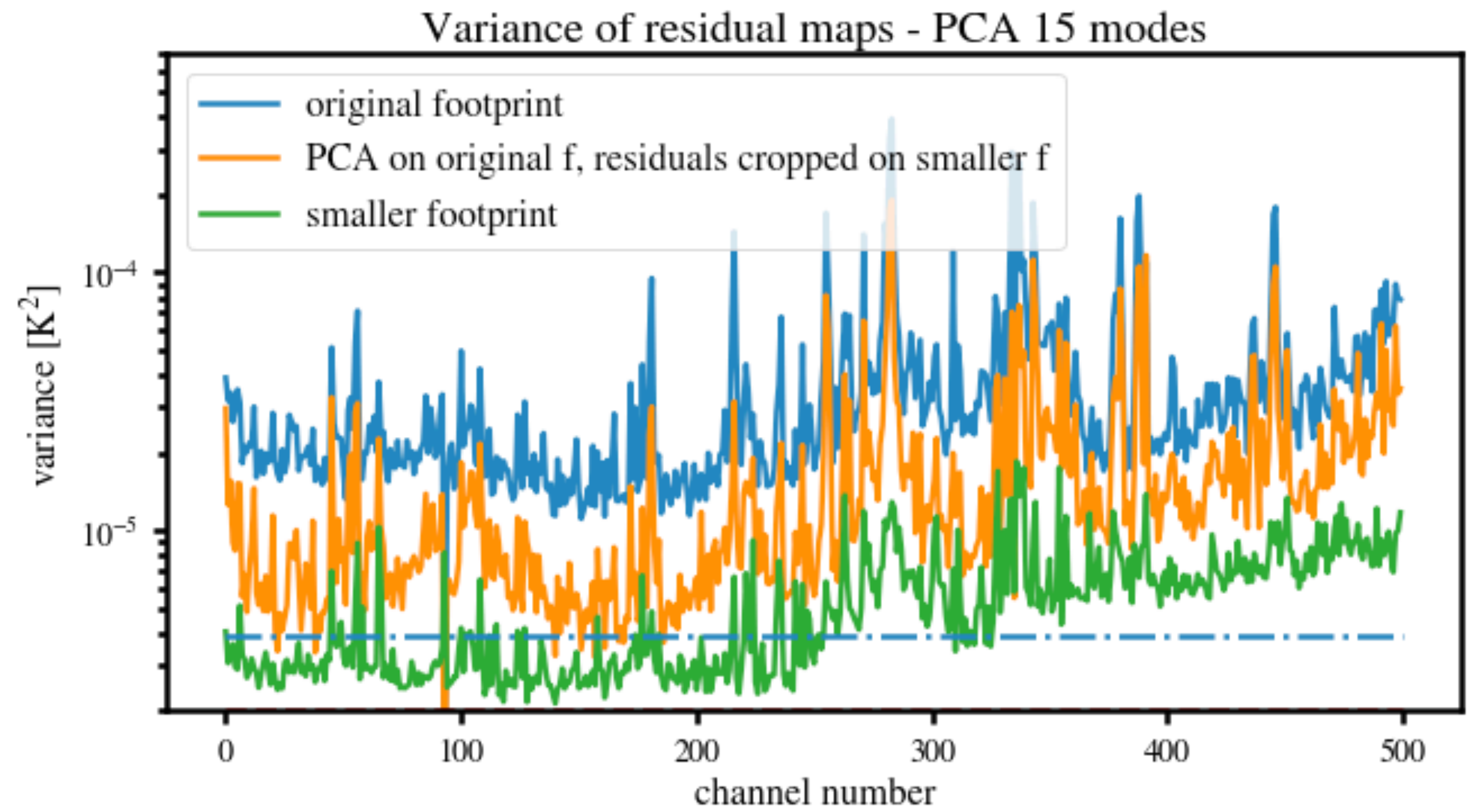
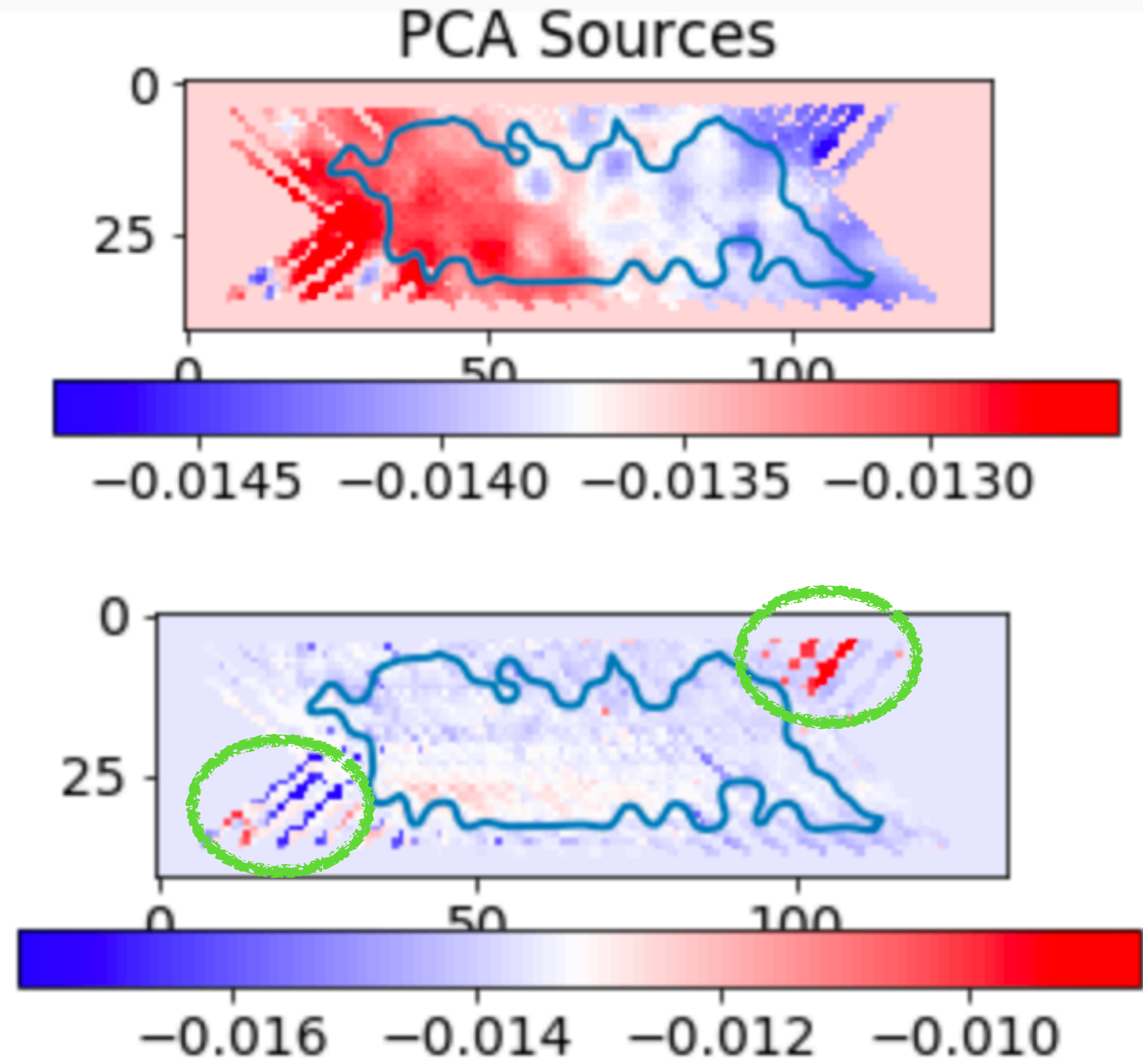
$$X = A S$$



Re-analysis of 2019 data

1. PCA-informed pixel flagging

$$X = A S$$

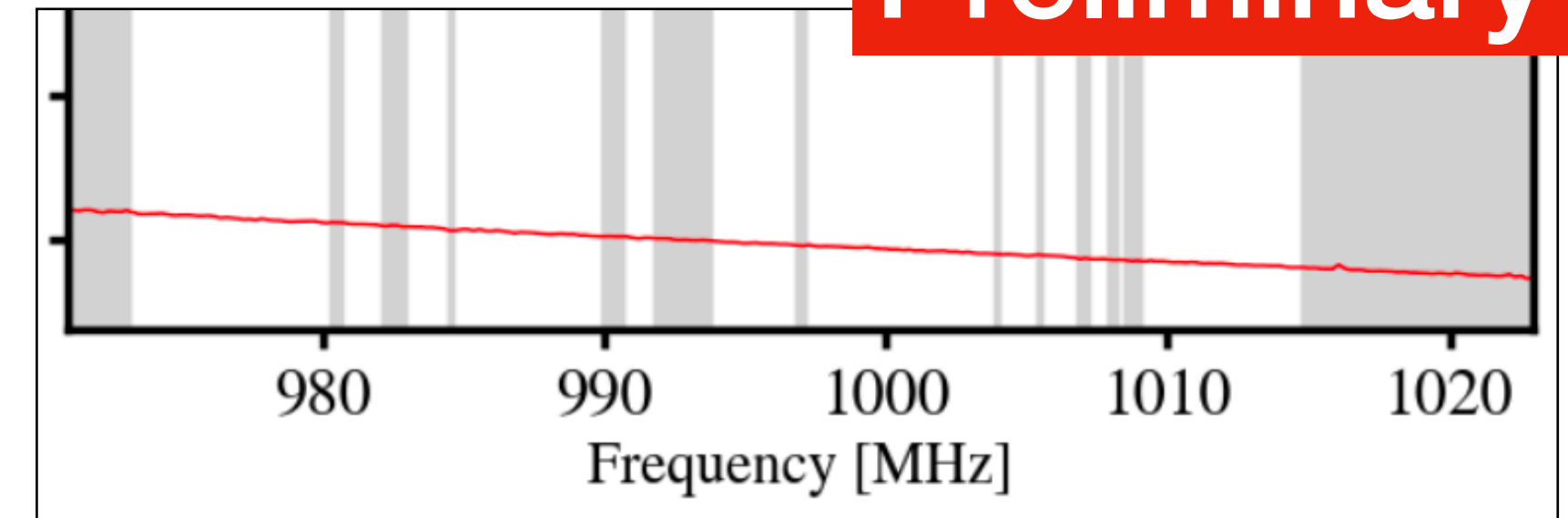


Re-analysis of 2019 data

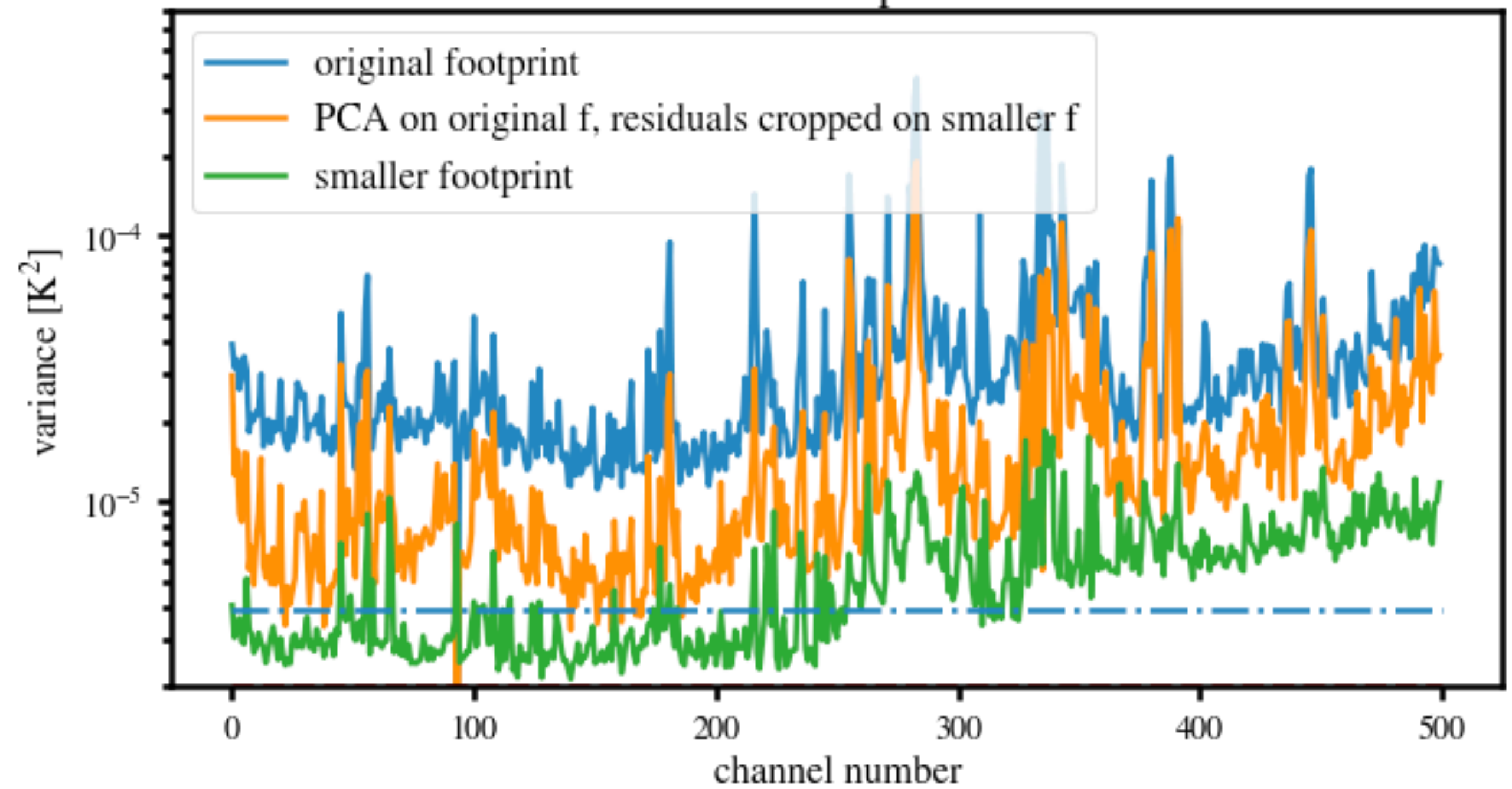
- 1. PCA-informed pixel flagging
- 2. Keep *bad* channels

see also discussion in Carucci+ 2020

$$X = A S$$



Variance of residual maps - PCA 15 modes



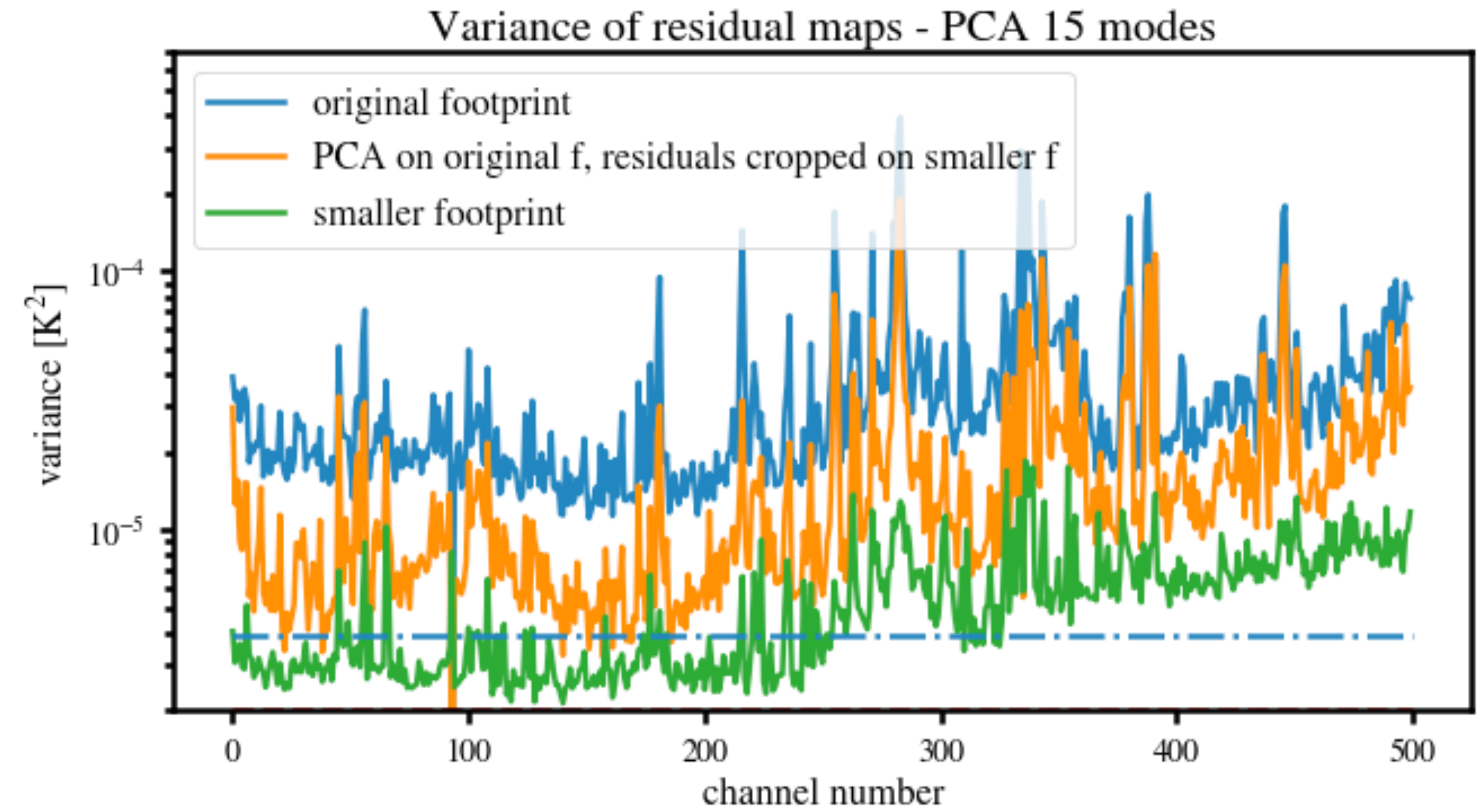
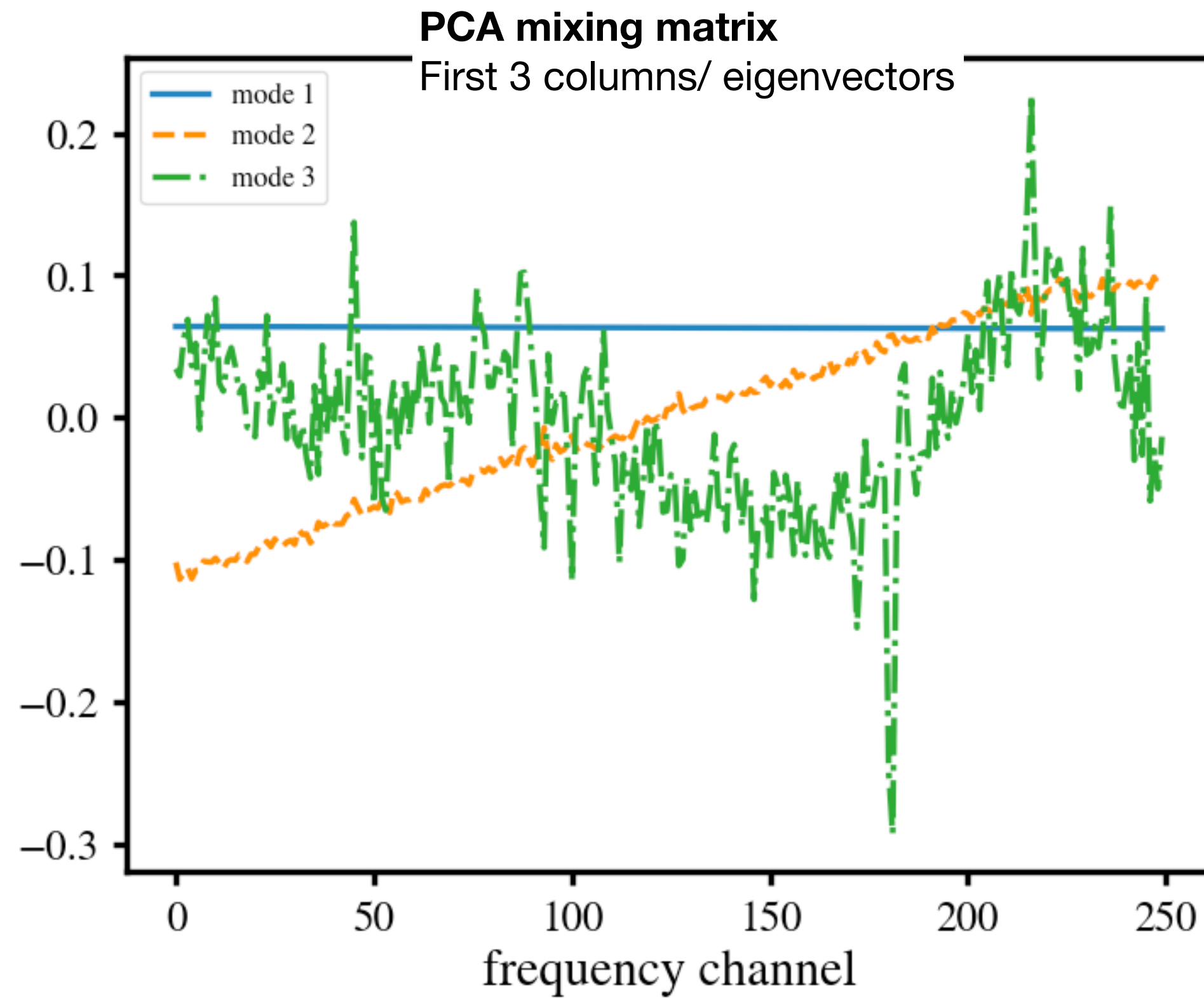
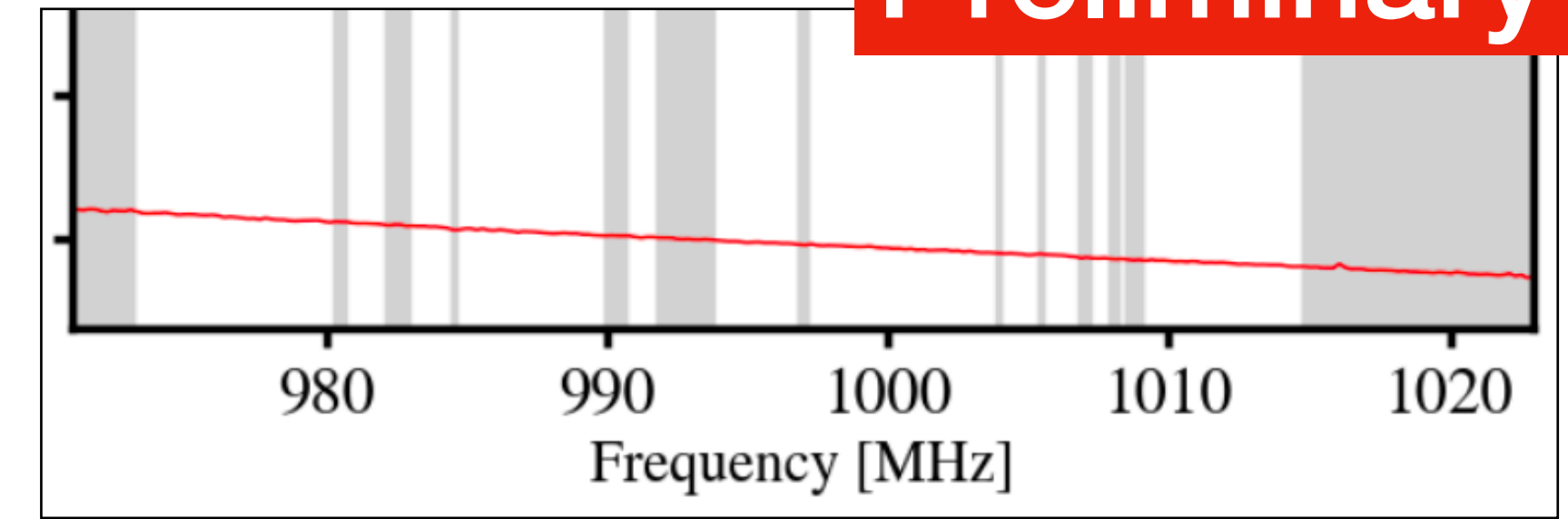
Re-analysis of 2019 data

1. PCA-informed pixel flagging
2. Keep *bad* channels

see also discussion in Carucci+ 2020

$$\mathbf{X} = \mathbf{A} \mathbf{S}$$

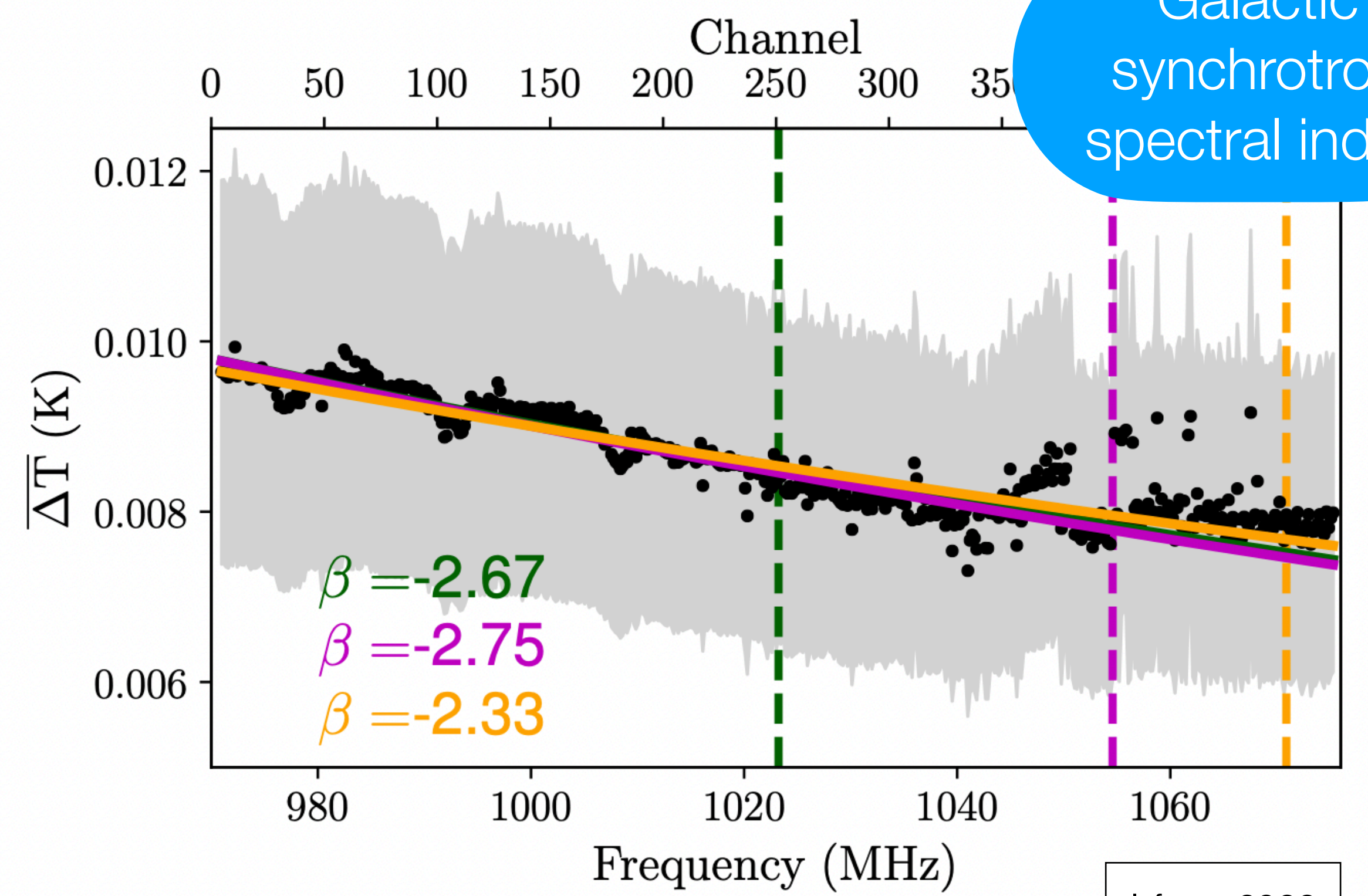
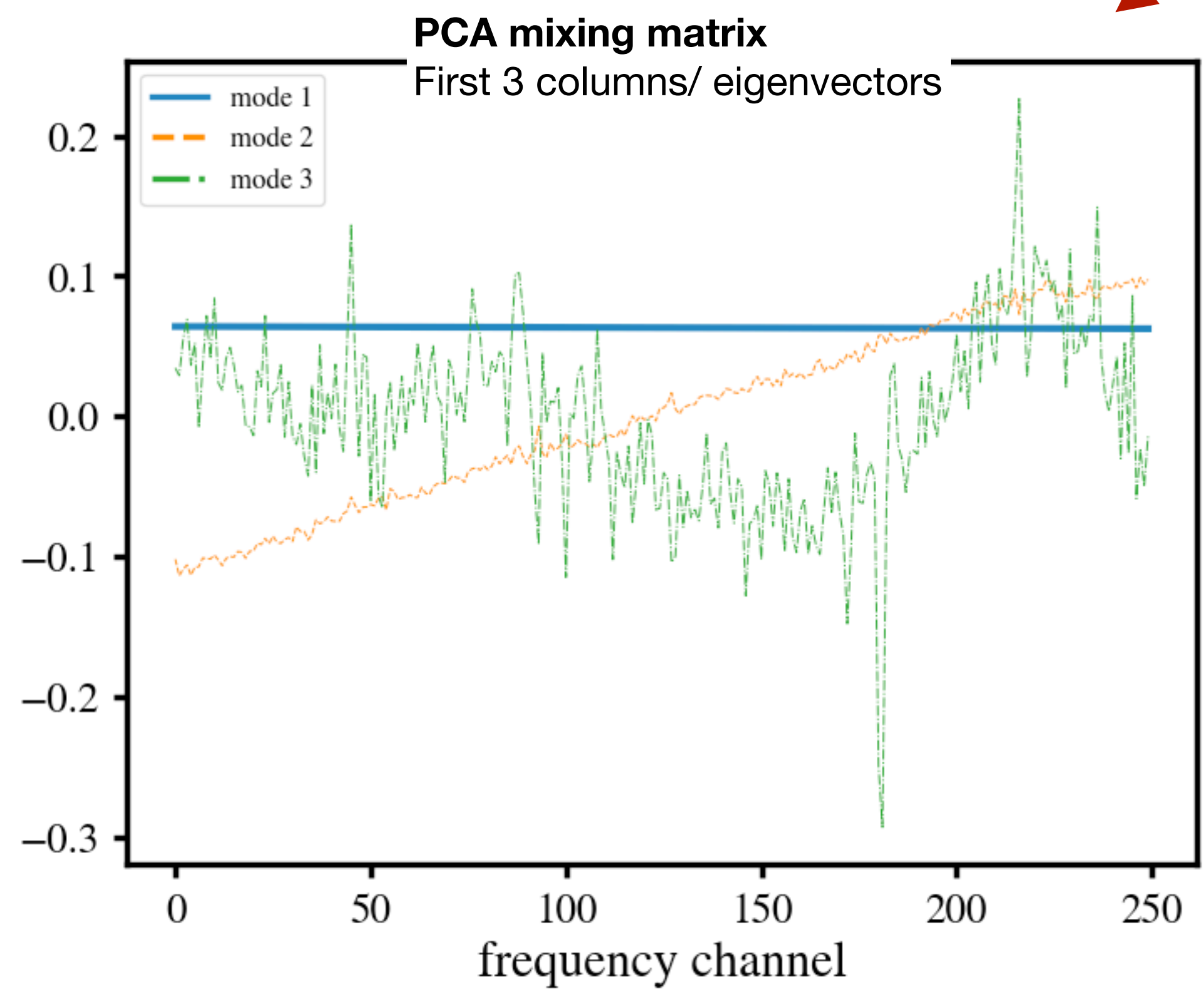
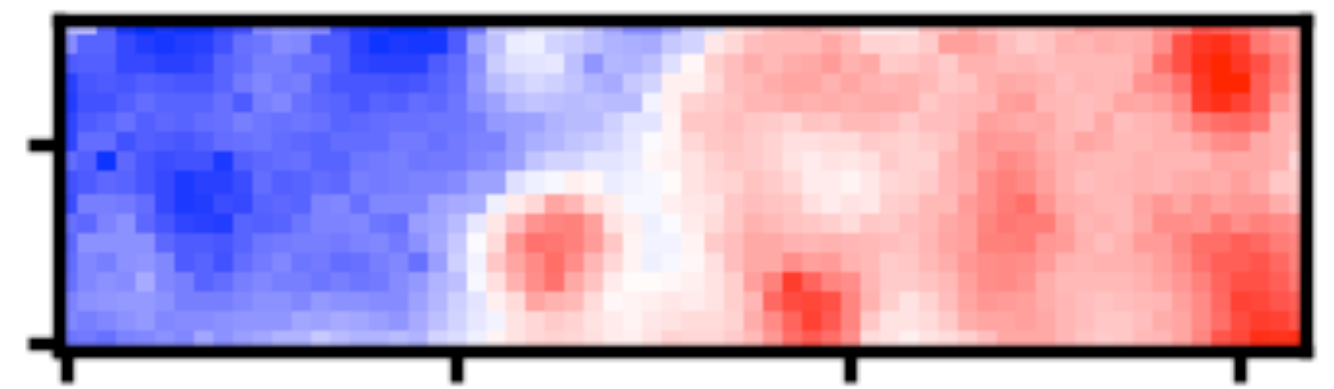
Preliminary



Re-analysis of 2019 data

1. PCA-informed pixel flagging
2. Keep *bad* channels

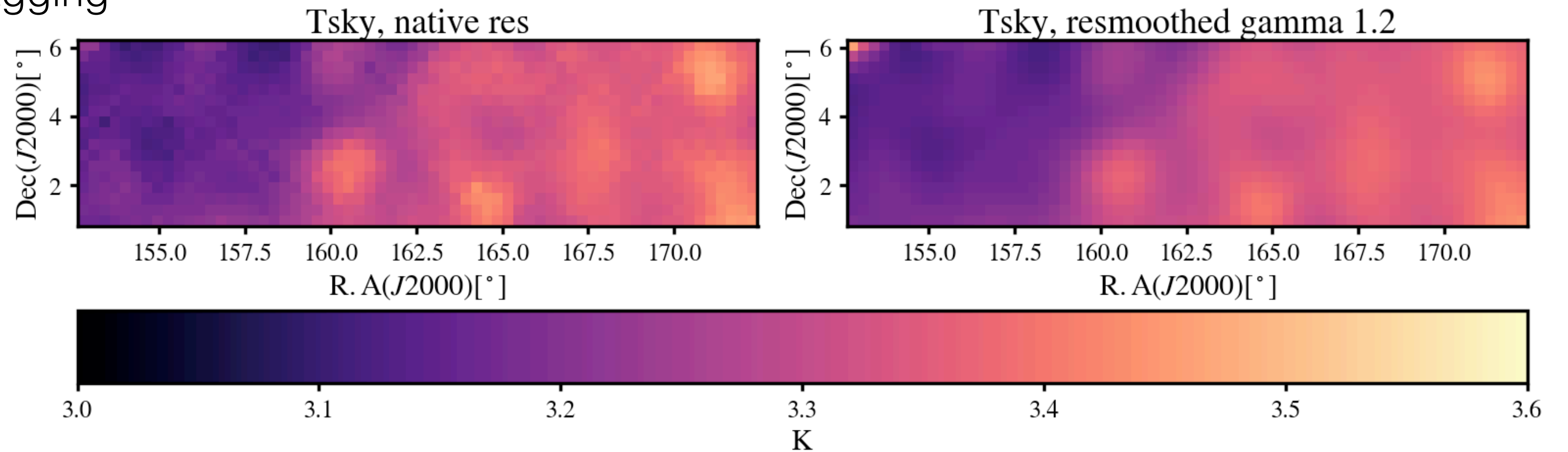
$$X = A S$$



Irfan+ 2022

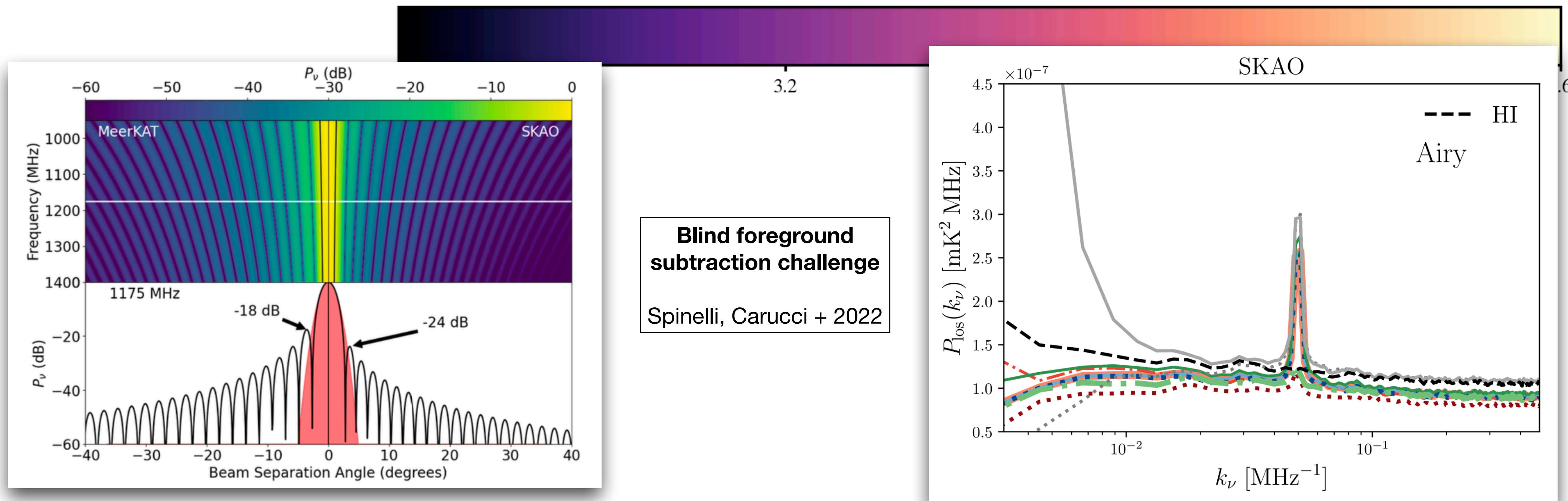
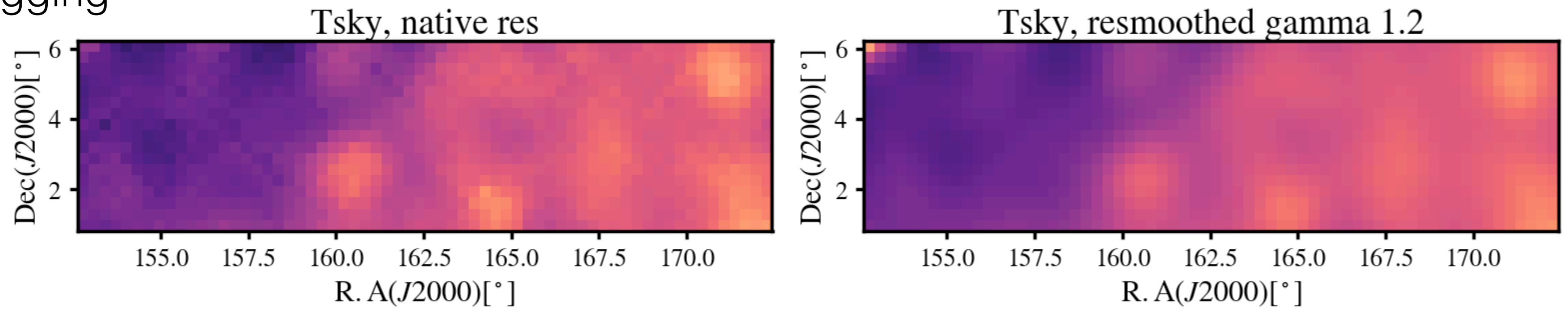
Re-analysis of 2019 data

1. PCA-informed pixel flagging
2. Keep *bad* channels
3. No re-smoothing



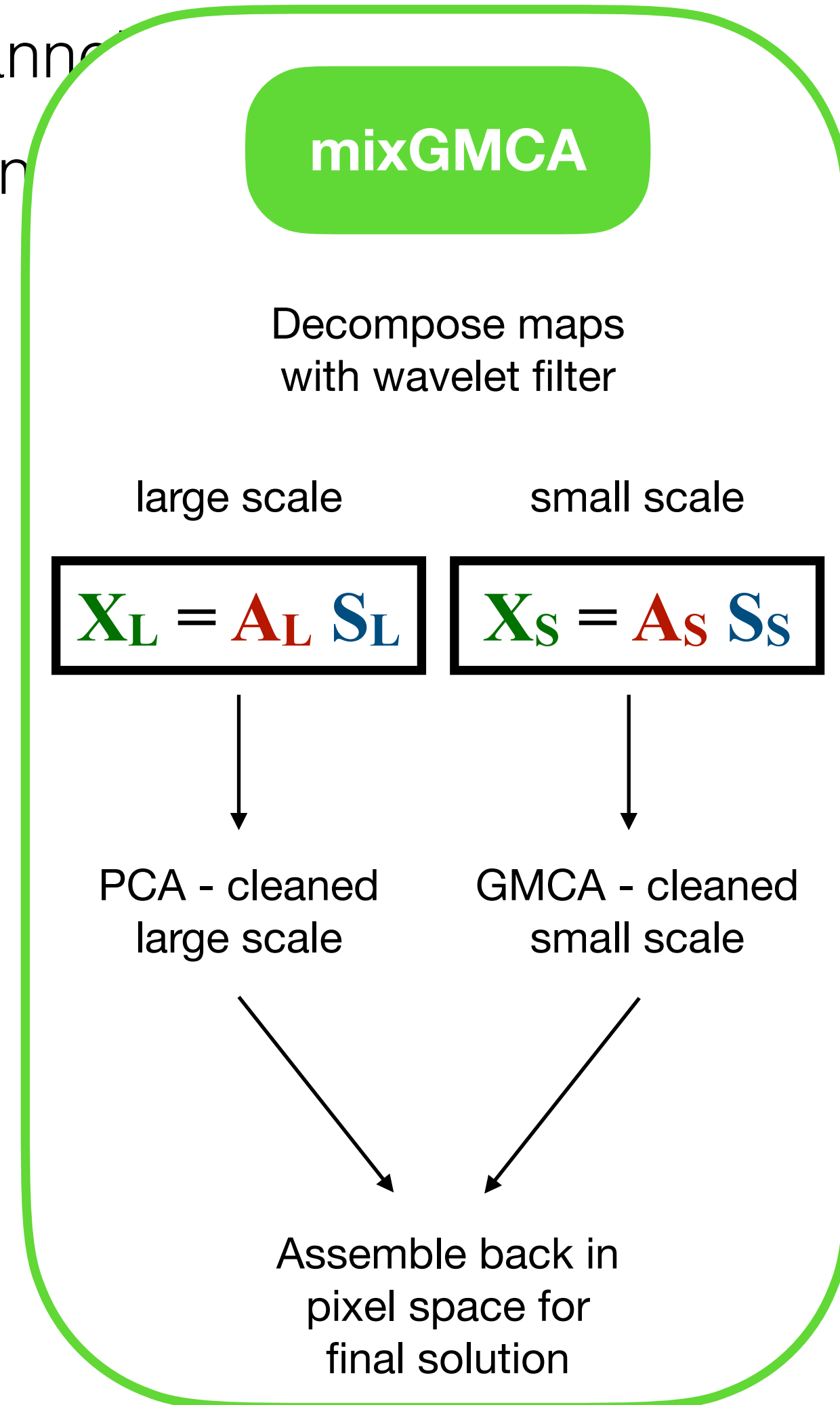
Re-analysis of 2019 data

1. PCA-informed pixel flagging
2. Keep *bad* channels
3. No re-smoothing



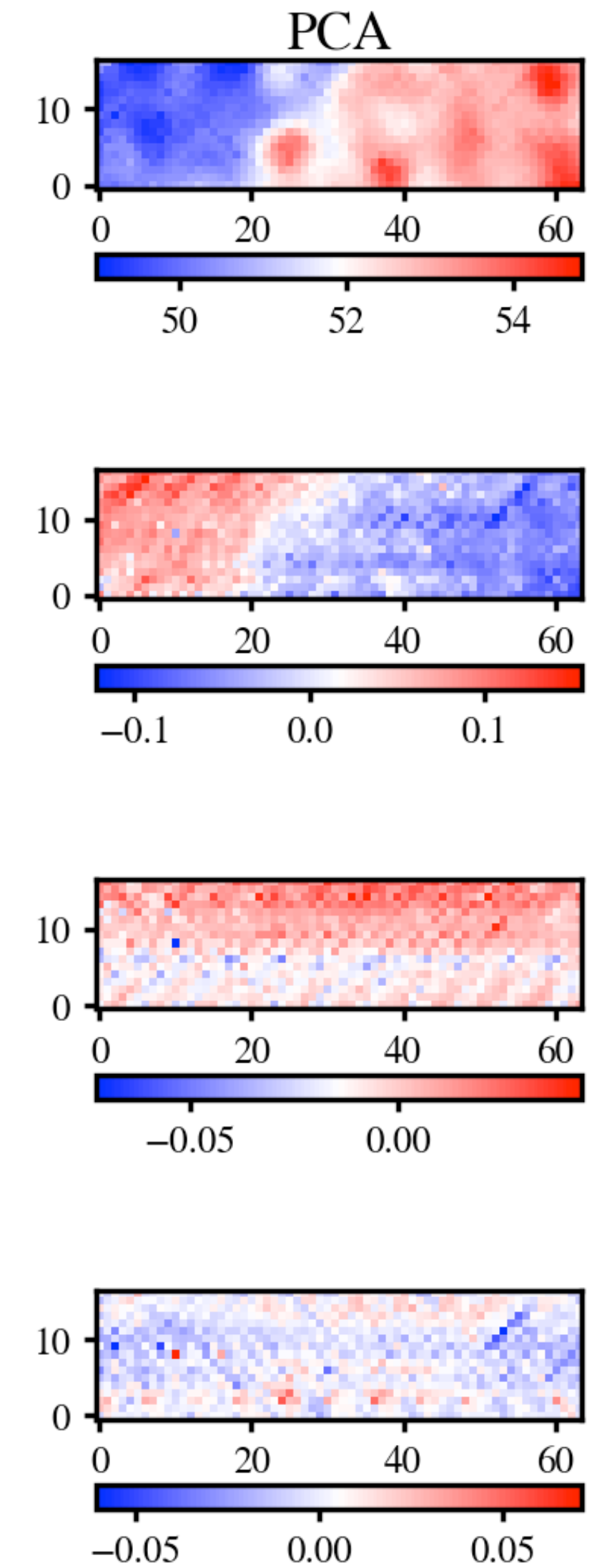
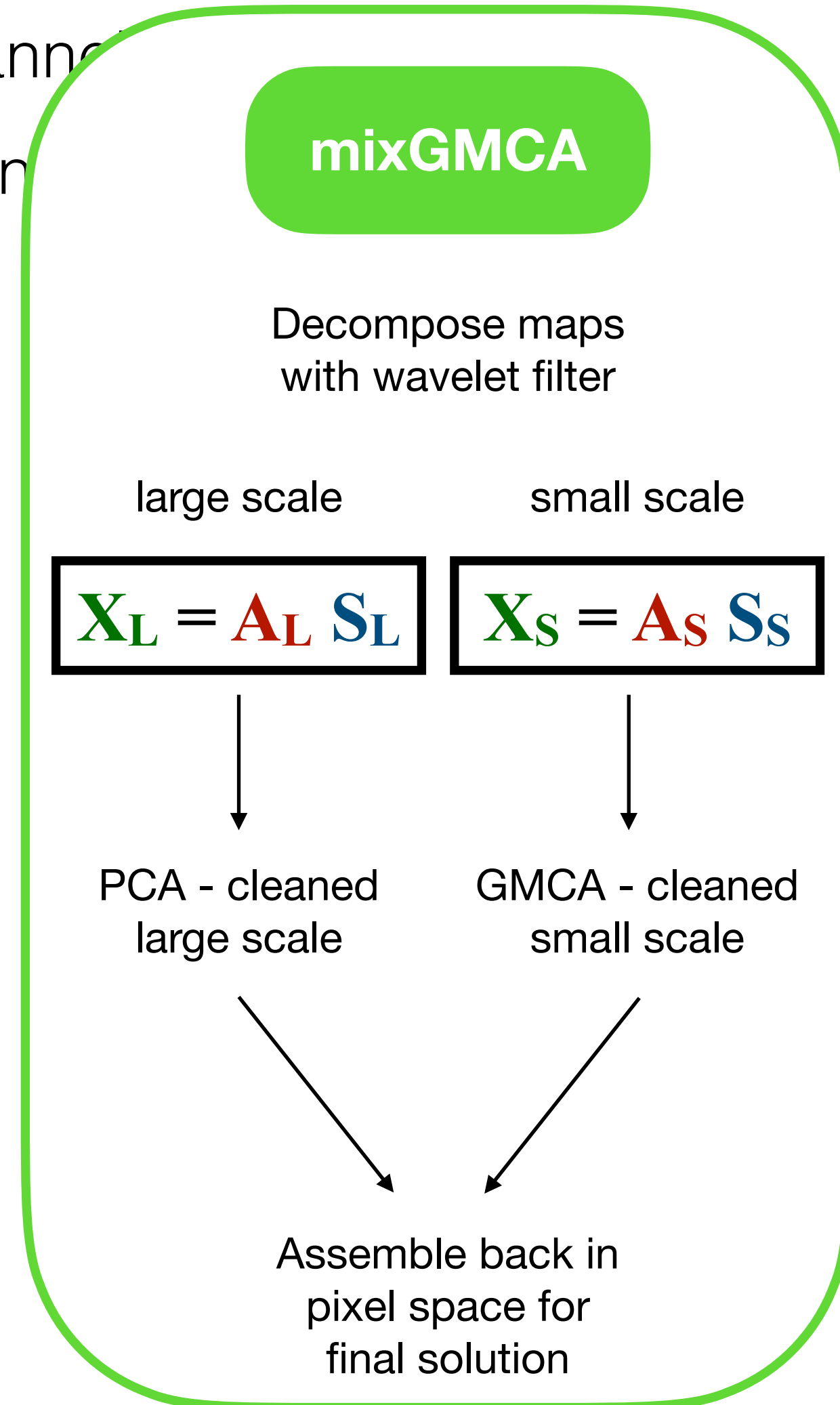
Re-analysis of 2019 data

1. PCA enforced pixel-flagging
2. Keep *bad* channels
3. No re-smoothing
4. Beyond PCA?



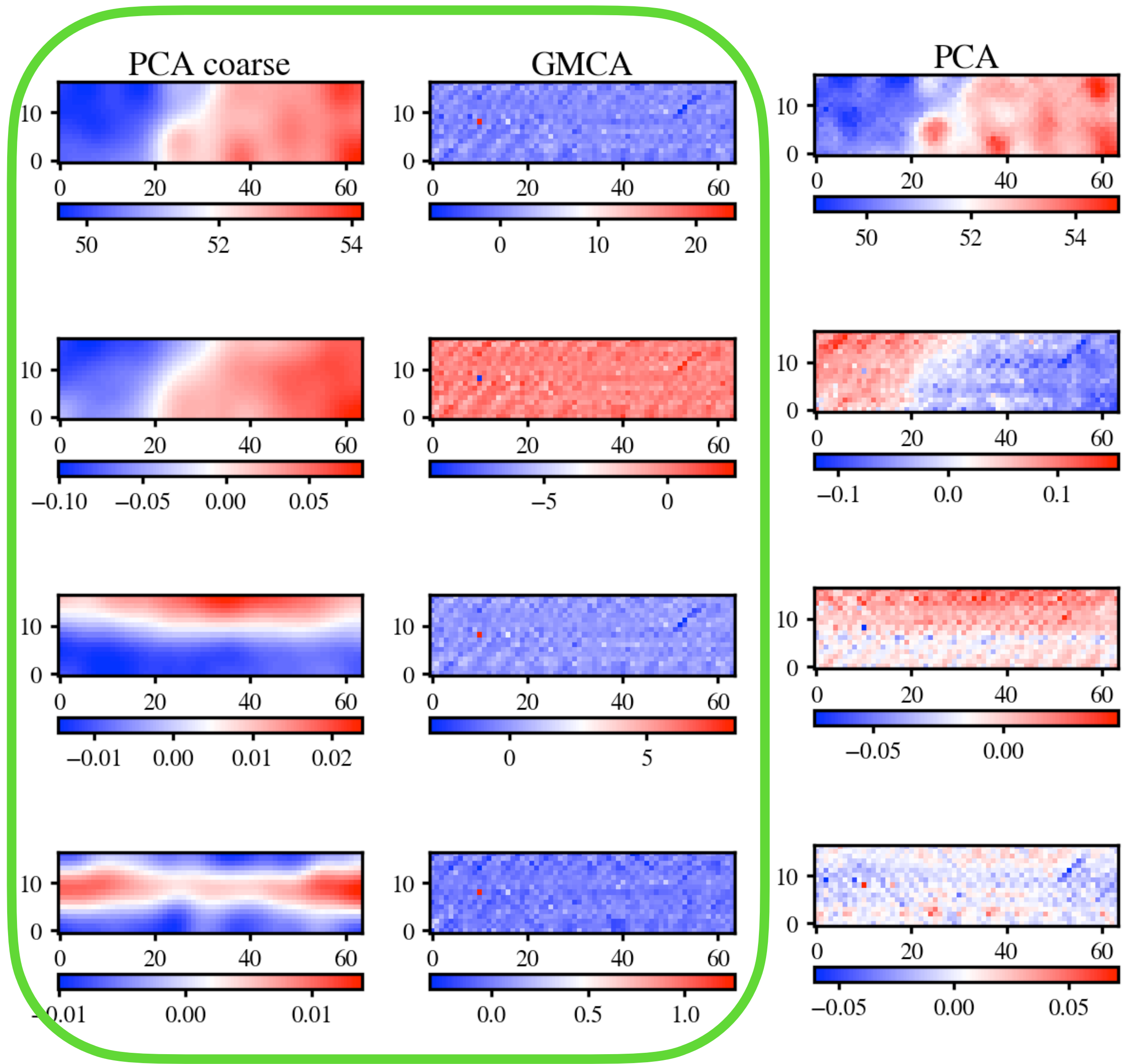
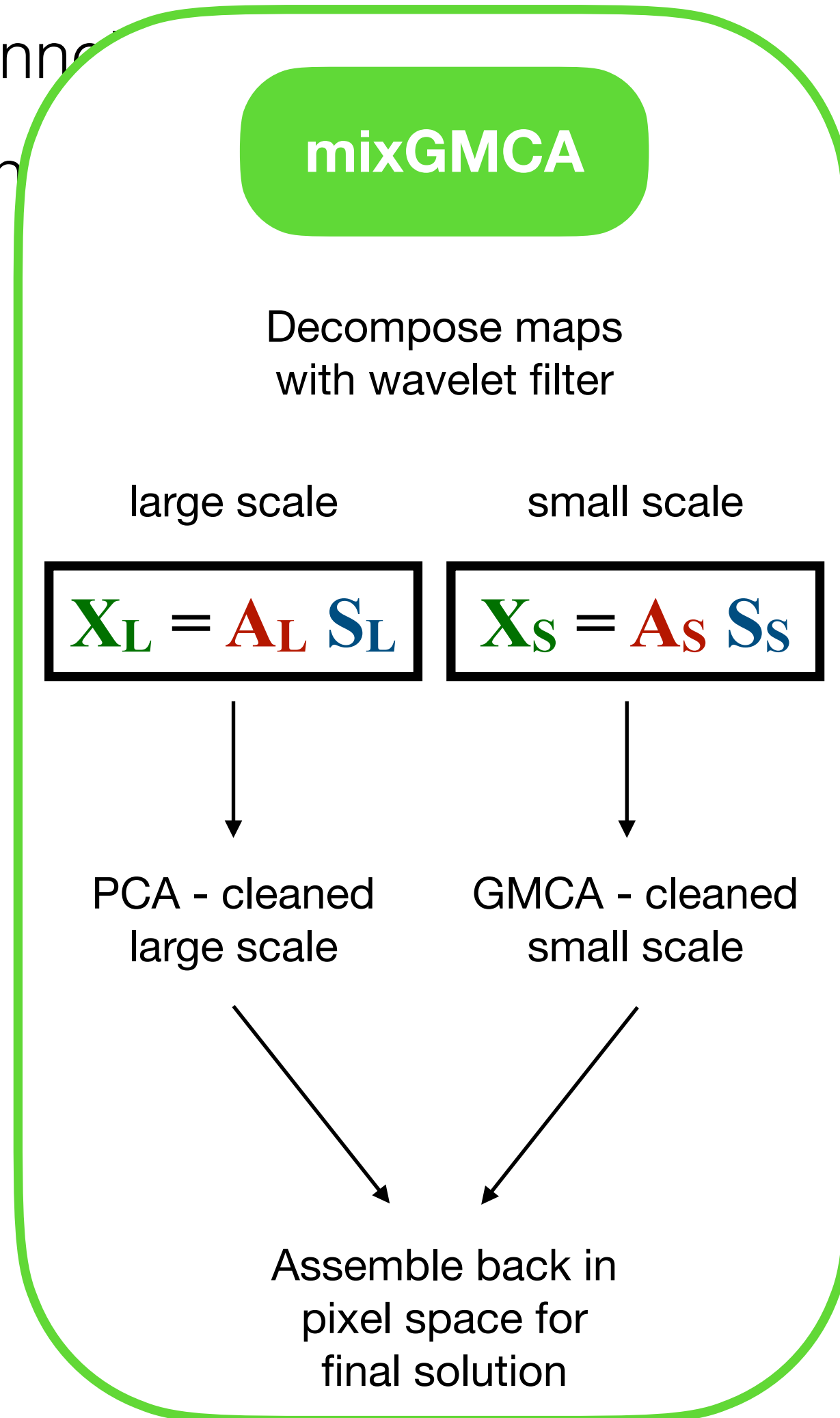
Re-analysis of 2019 data

1. PCA enforced pixel-flagging
2. Keep *bad* channels
3. No re-smoothing
4. Beyond PCA?



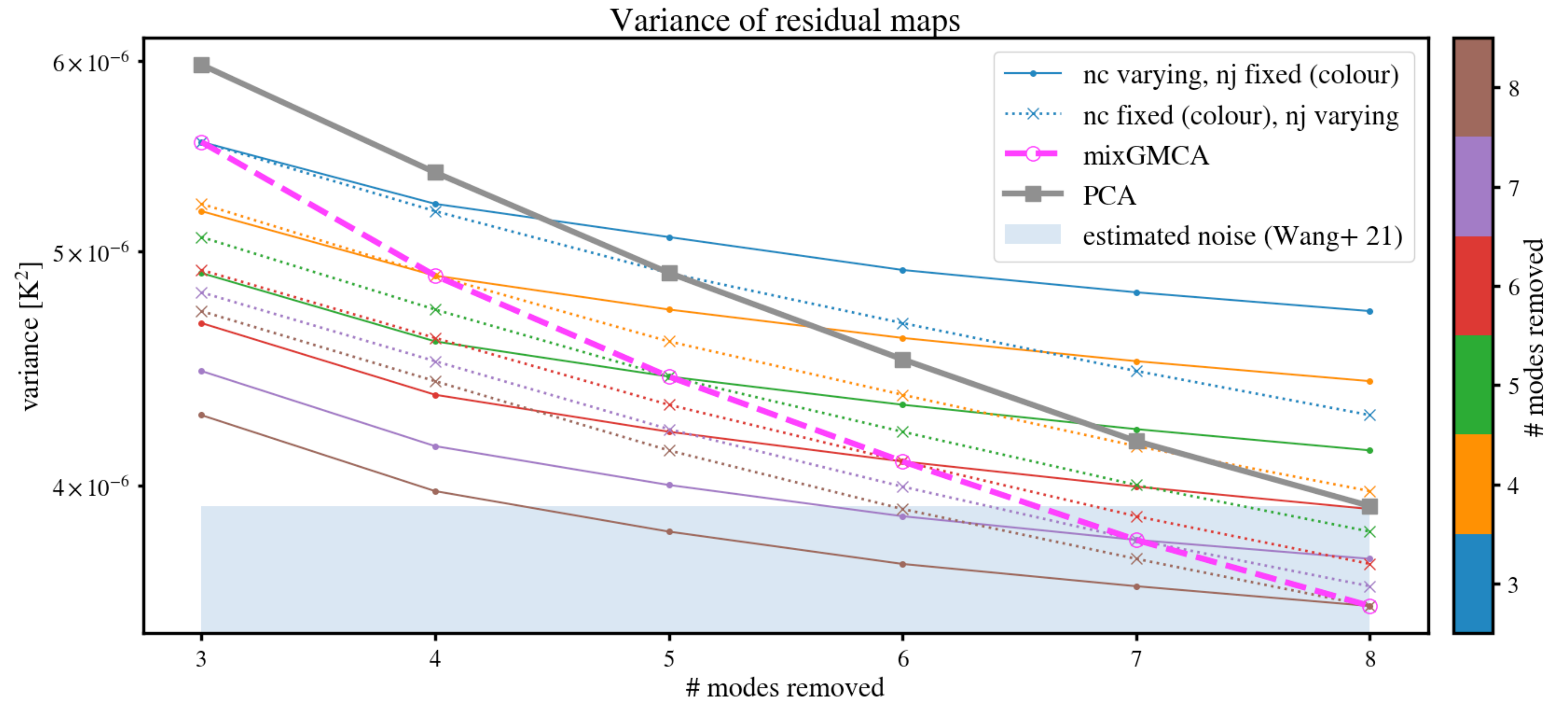
Re-analysis of 2019 data

1. PCA enforced pixel-flagging
2. Keep *bad* channels
3. No re-smoothing
4. Beyond PCA?

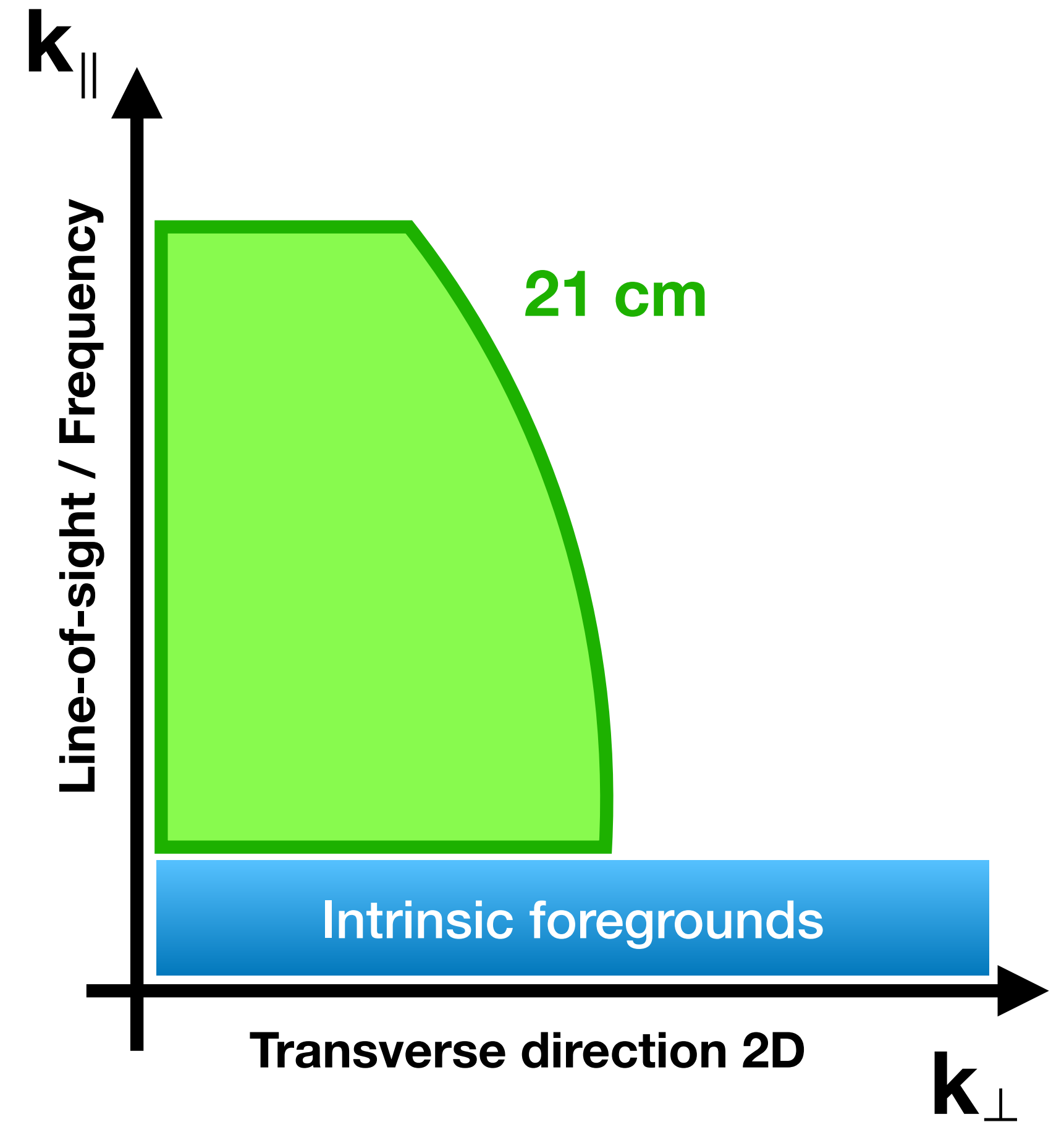


Re-analysis of 2019 data

1. PCA enforced pixel-flagging
2. Keep *bad* channels
3. No re-smoothing
4. Beyond PCA?



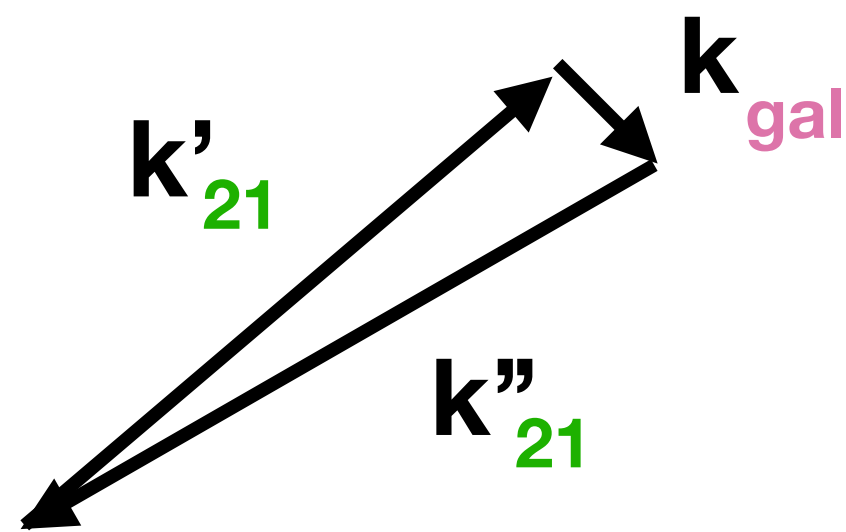
Beat systematics: cross correlations



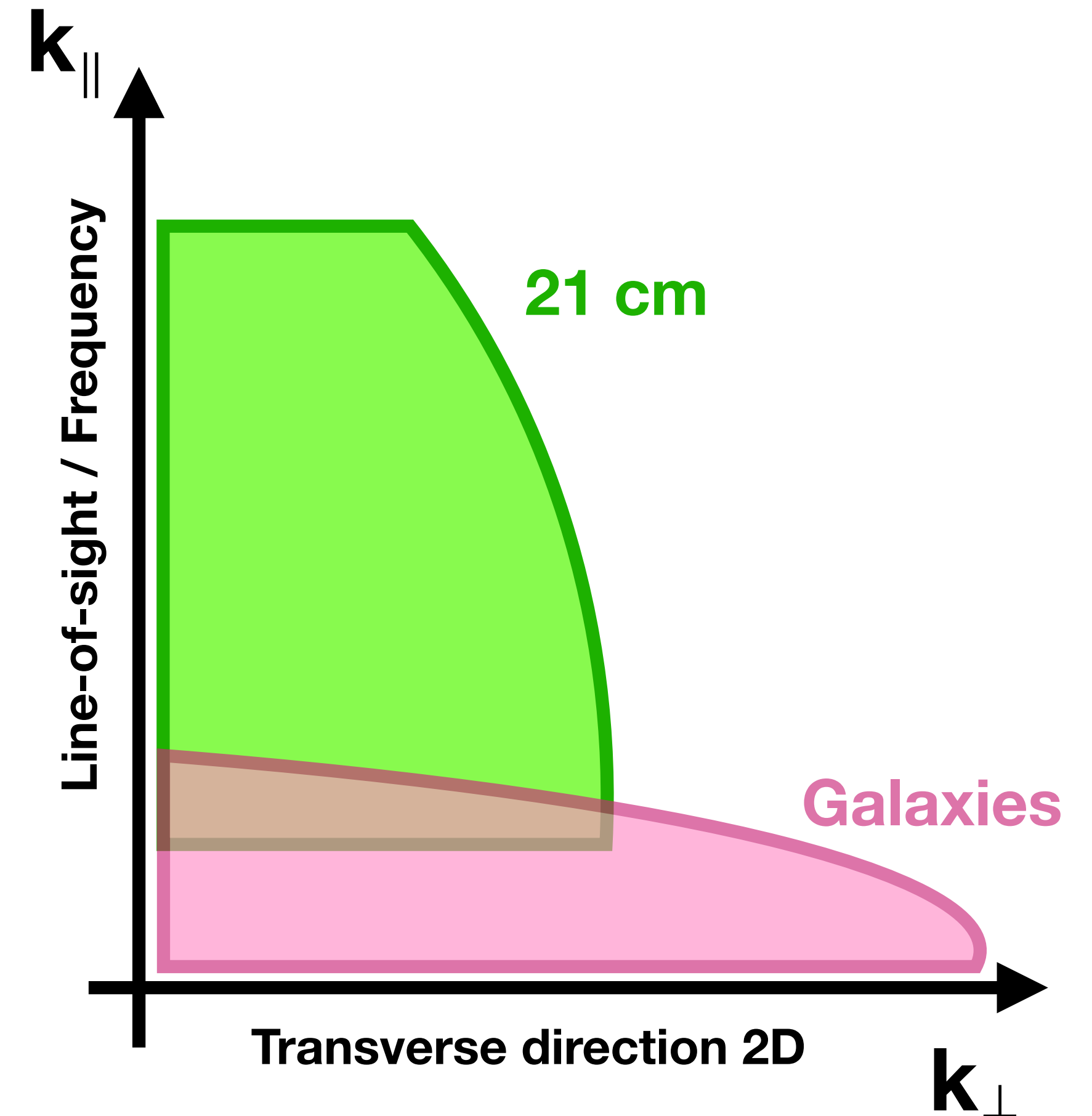
Beat systematics: cross correlations

Direct 21 cm x galaxies signal vanishes due to foregrounds in long wavelength line-of-sight modes. Need to use higher order correlations.

- e.g., a *squeezed* bispectrum estimator:
1 low-k mode from galaxy survey X 2 high-k 21 cm modes.

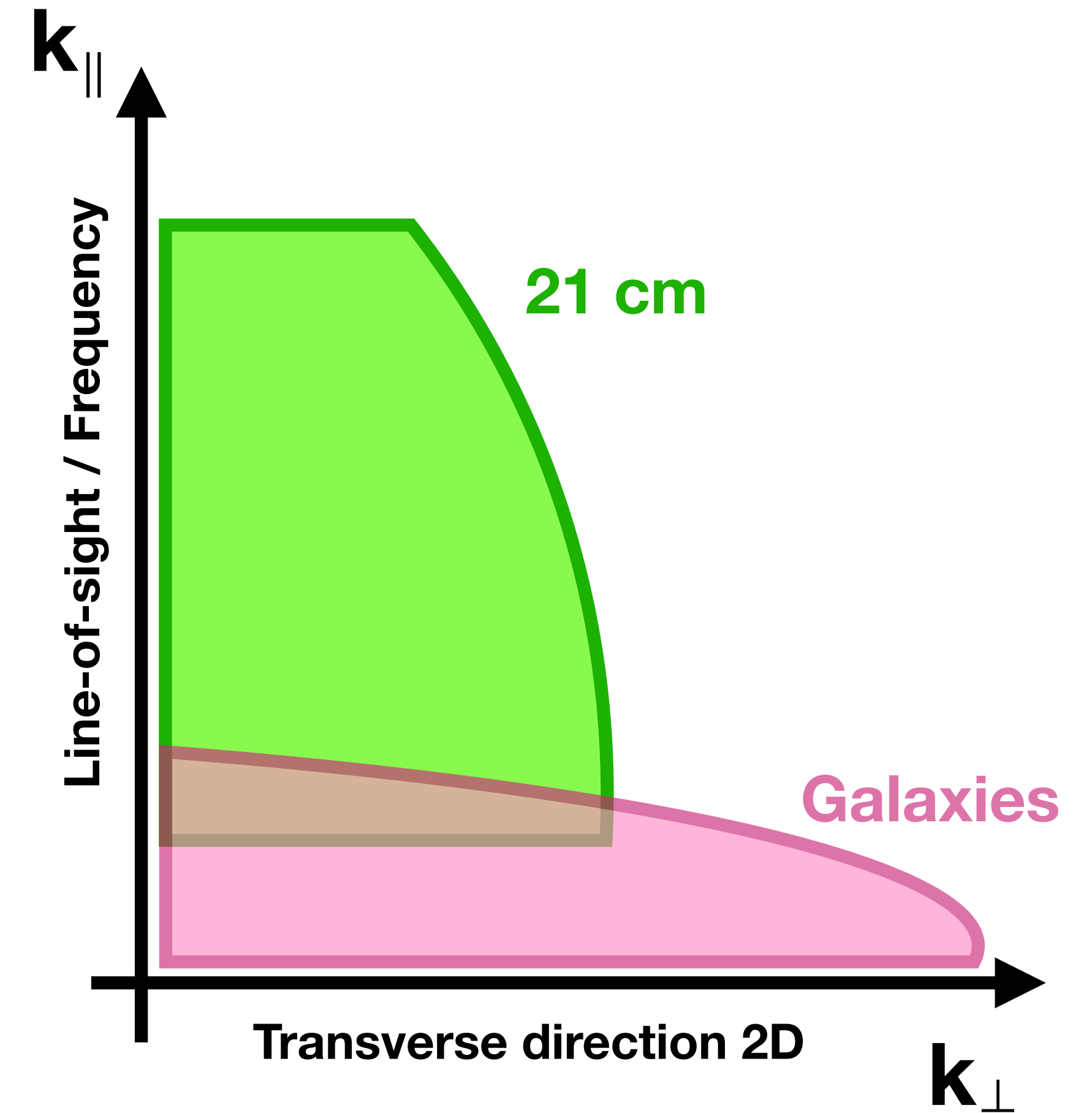
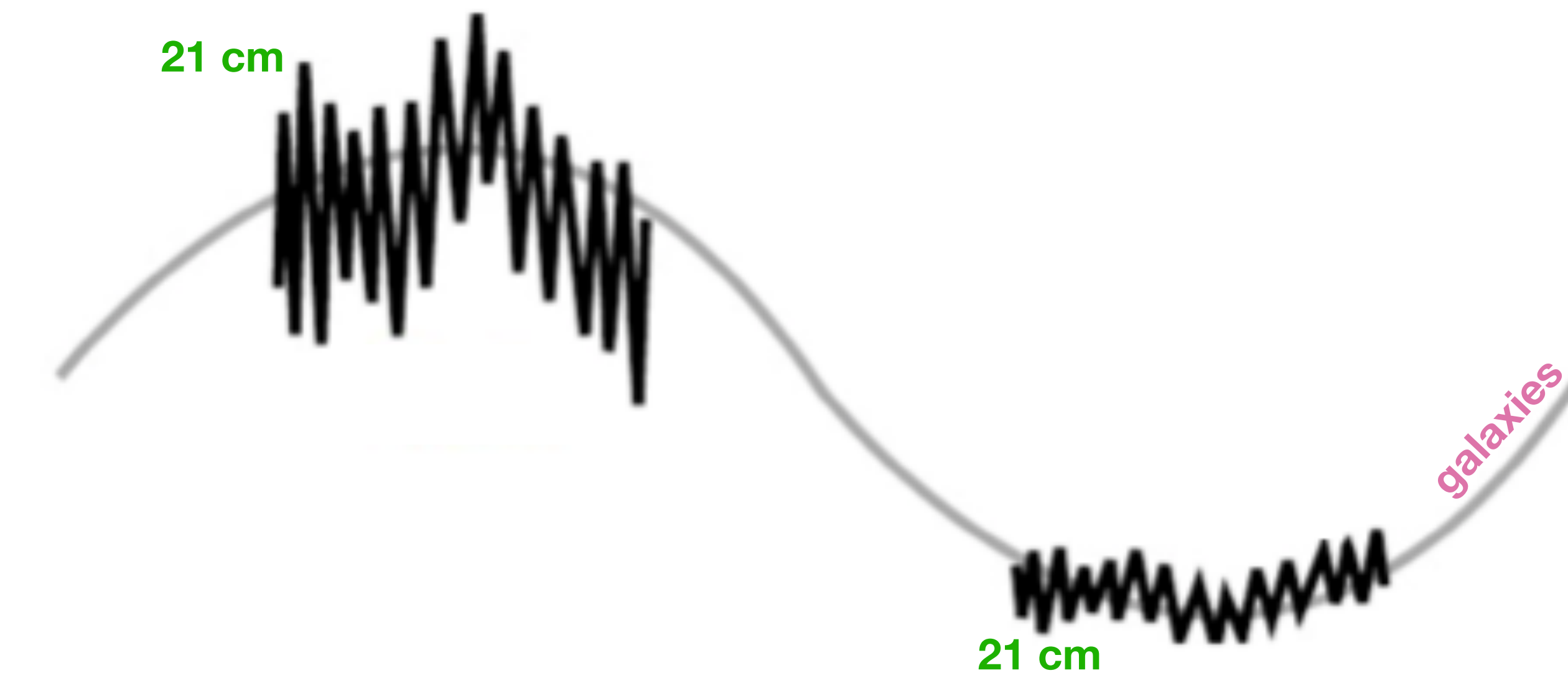


$$\langle \delta(\mathbf{k}) \delta(\mathbf{k}') \delta(\mathbf{k}'') \rangle = \delta_D(\mathbf{k} + \mathbf{k}' + \mathbf{k}'') B(\mathbf{k} + \mathbf{k}')$$



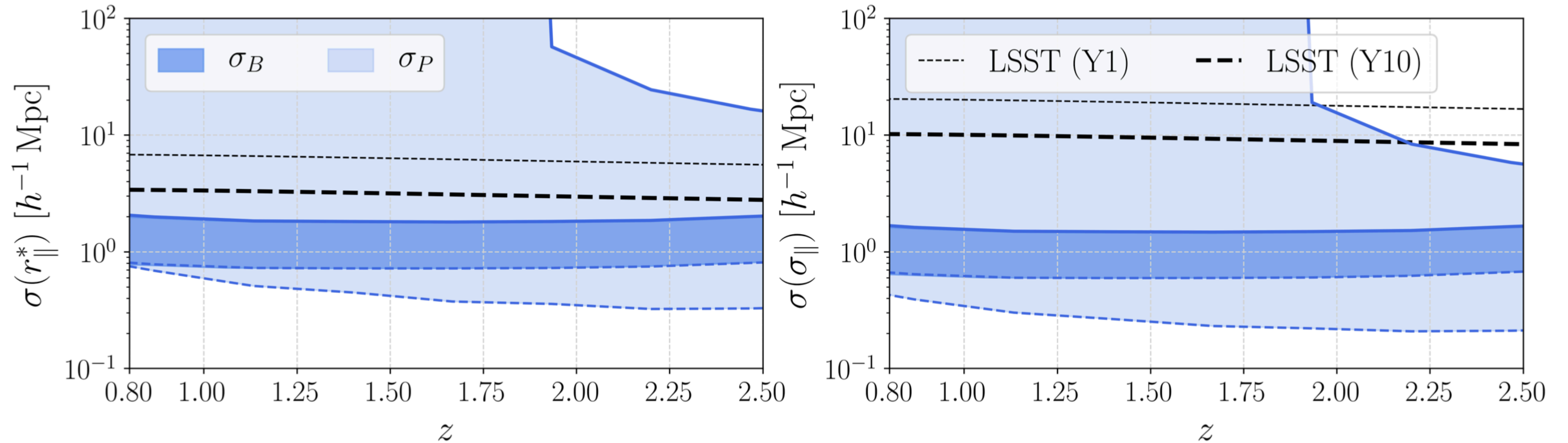
Beat systematics: cross correlations

Direct 21 cm x galaxies signal vanishes due to foregrounds in long wavelength line-of-sight modes. Need to use higher order correlations.



Beat systematics: cross correlations

Direct 21cm x galaxies signal vanishes due to foregrounds in long wavelength line-of-sight modes.
Need to use higher order correlations.



Guandalin+ 2022

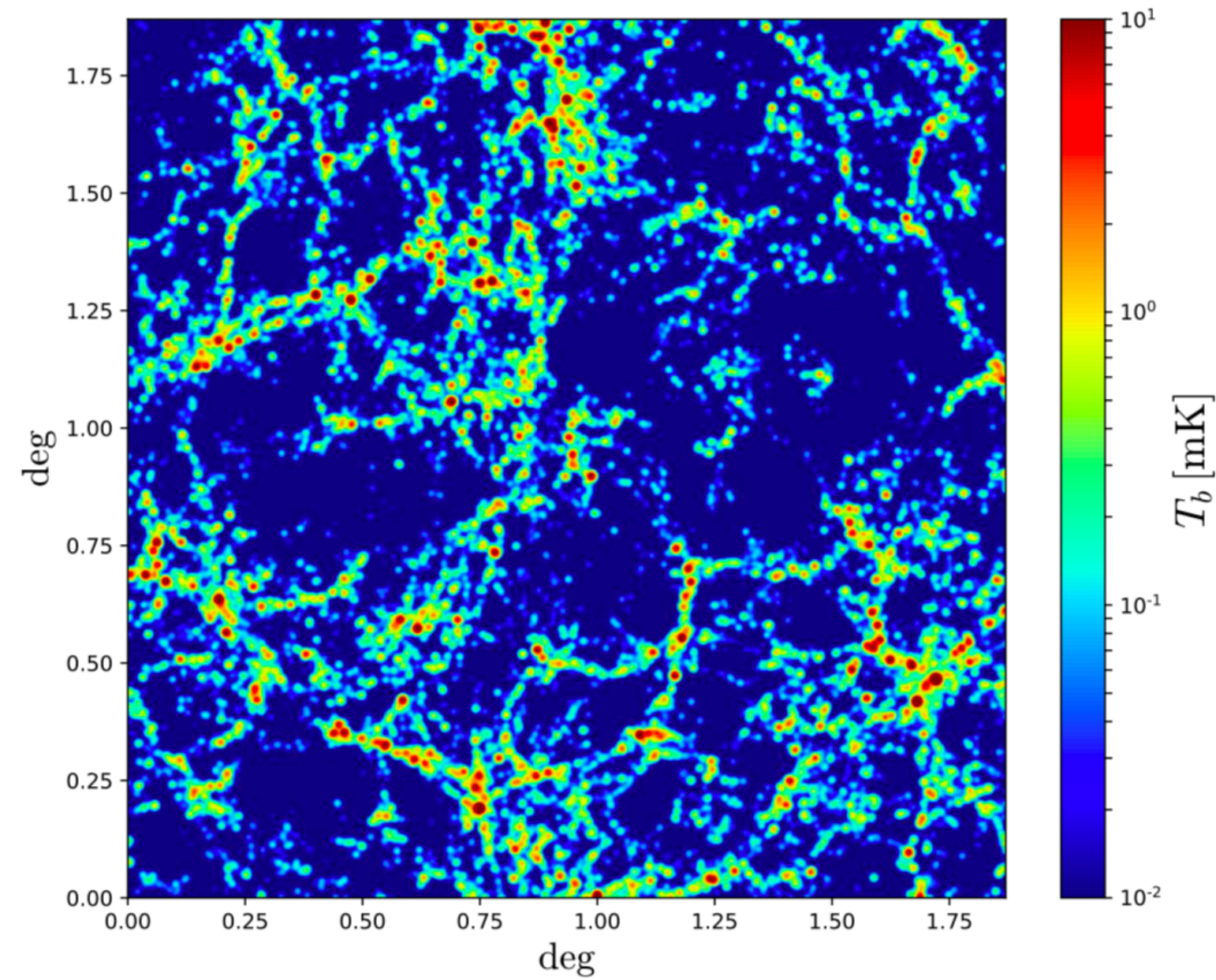
Summary

- **HI IM** will be game changer in cosmology
- Contaminants-removal is the biggest problem, lots of efforts devoted to this
- **MeerKLASS** ongoing!
- We are detecting (again!) the cross signal with WZ galaxies to test different pre-processing steps and cleaning algorithms
- Cross-Pk as sweet-spot as we play with # of components removed
- Cross-Pk in agreement among different cleaning methods
- Separating scales for the cleaning is more efficient at reducing the cube variance (mixGMCA)
- Smart estimators + clever cross possibilities: we can get plenty of new/complementary info out of the maps

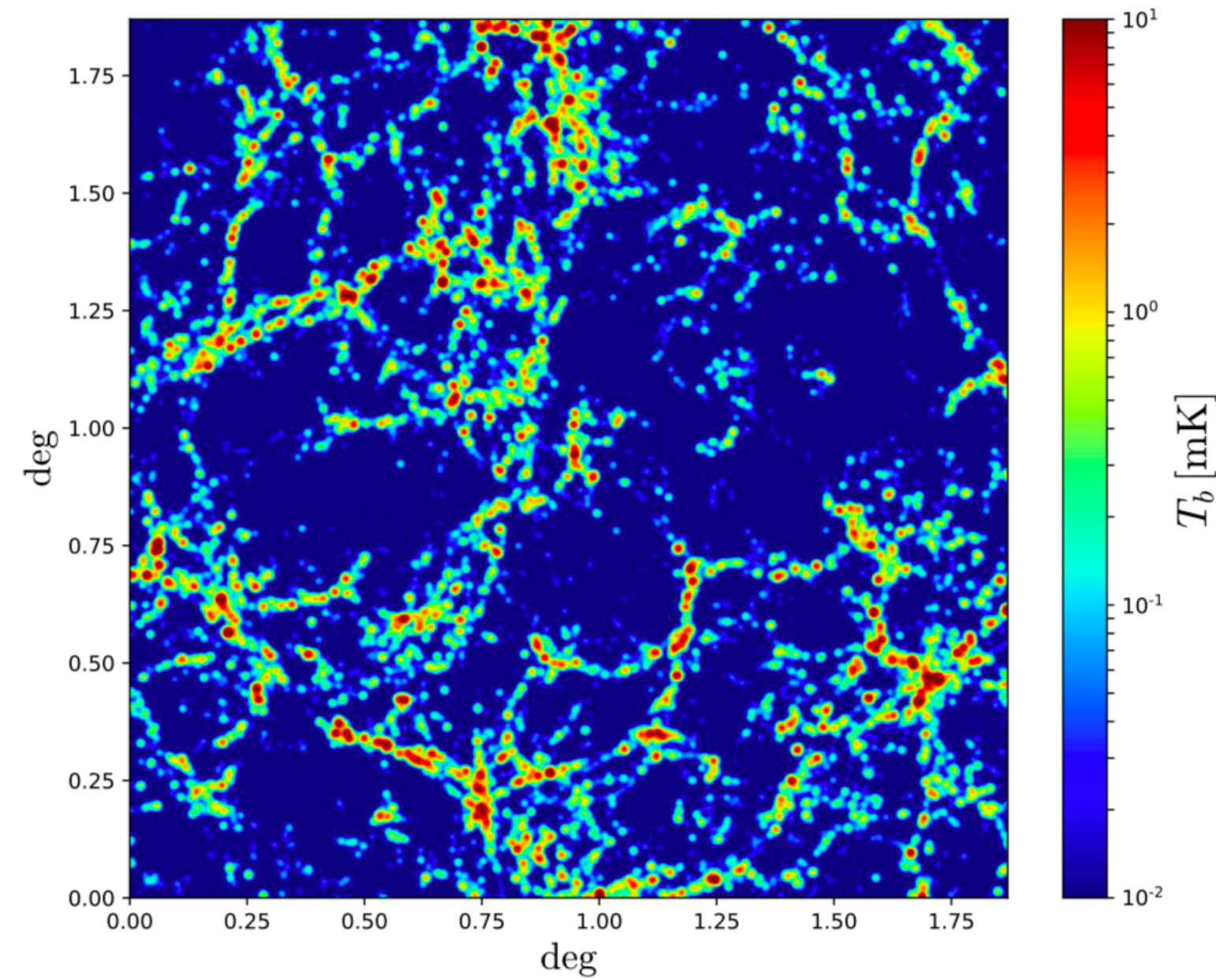
Getting ready for the SKAO HI IM science

Distribution of HI in the Universe

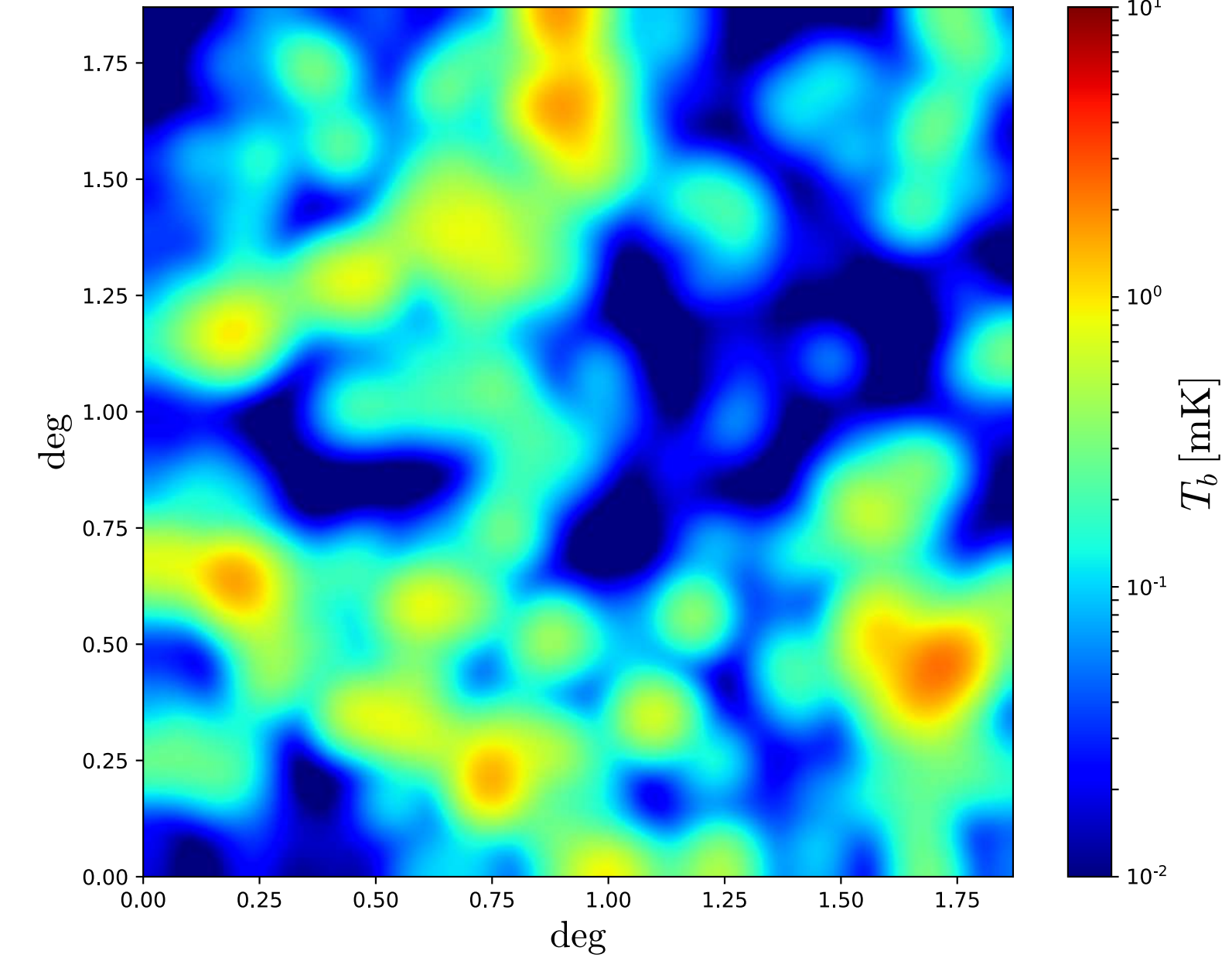
cheap Nbody



IllustrisTNG



forecasts



Villaescusa-Navarro, .. , IPC + 2018

HI intensity mapping

