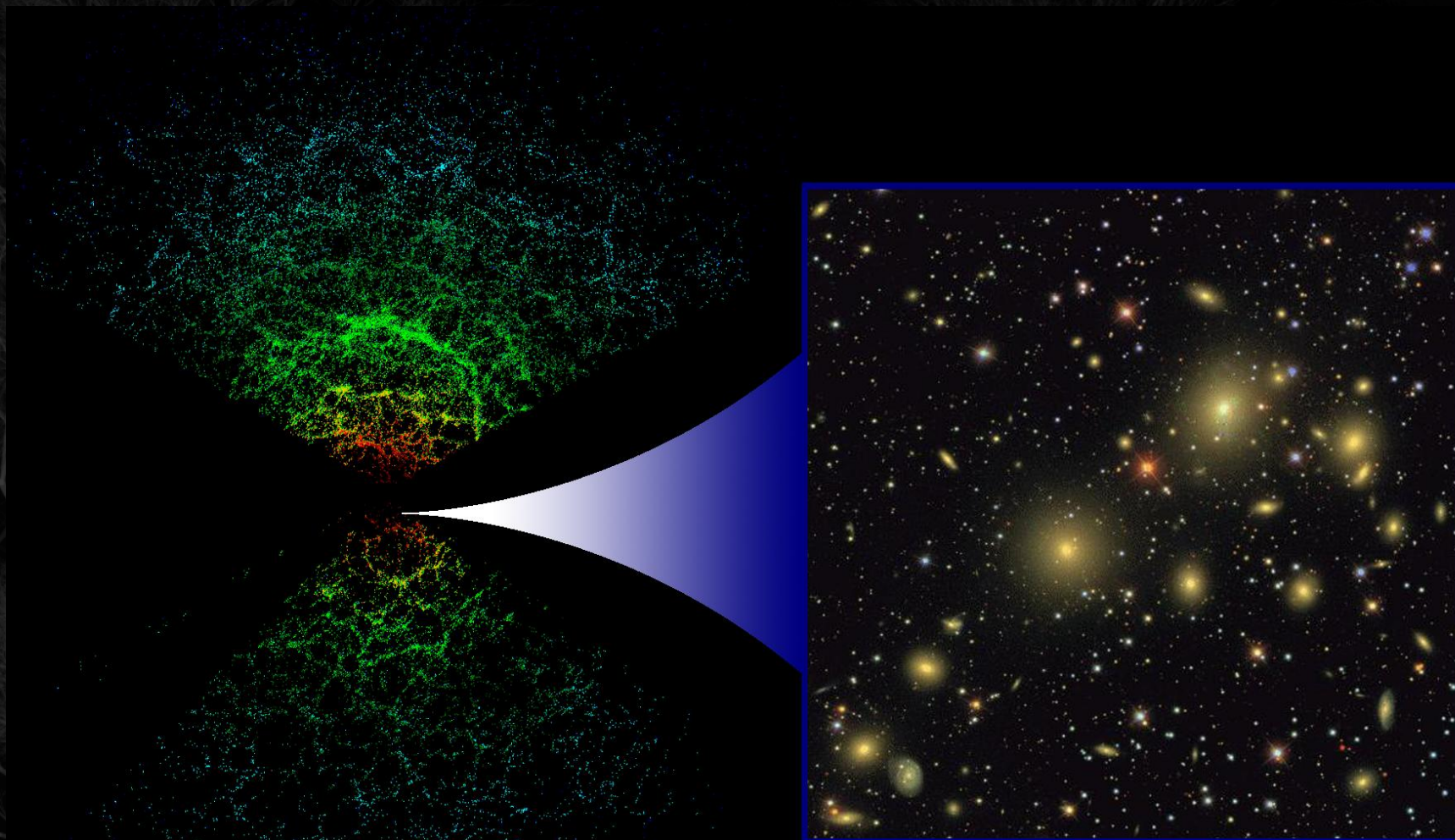


Probabilistic photometric redshifts in the era of Petascale Astronomy

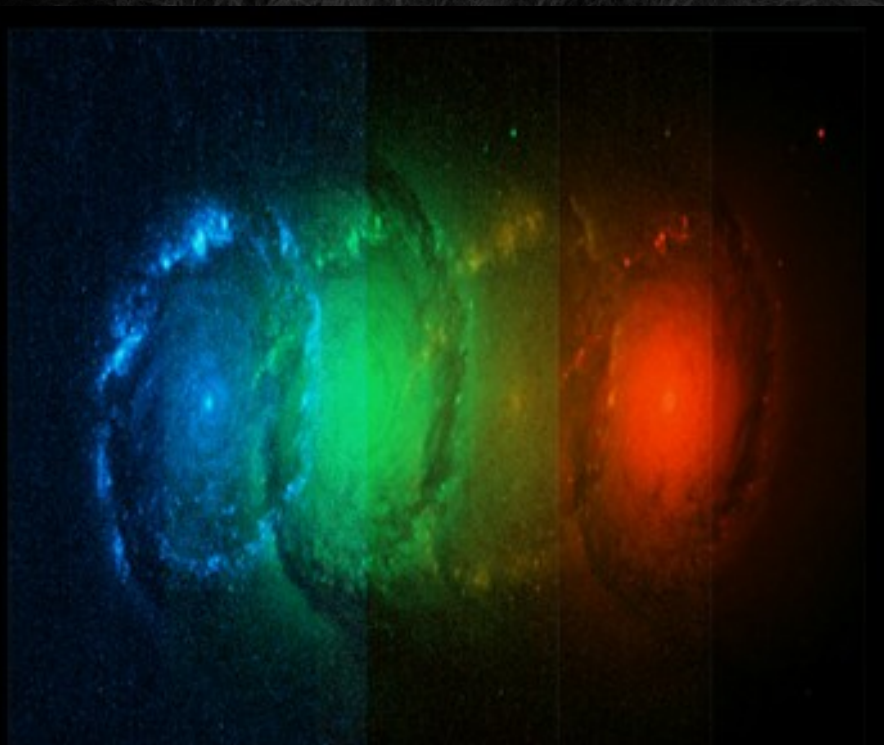
Matías Carrasco Kind

NCSA/Department of Astronomy
University of Illinois at Urbana-Champaign

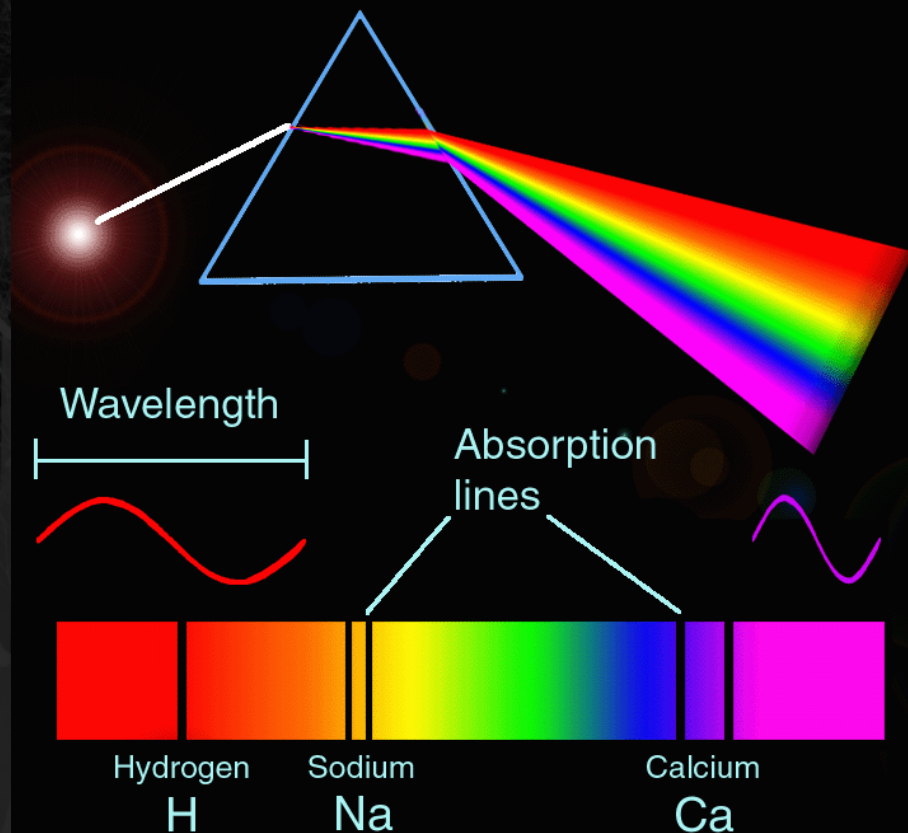


Credit: SDSS Collaboration

3D Clustering of galaxies as a probe in cosmology, e.g., 2 point correlation function, power spectrum of the galaxy distribution, etc.



VS.





It's happening! 😊

~ 300 millions galaxies up to $z = 1.5$

5,000 squares degrees (1/8 sky)

Data management at NCSA

DES specially designed to probe the origin of dark energy

S/G class and photo- z needed

1 TB of data per day

2 years completed, 3 more to go



It's happening! 😊

~ 300 millions galaxies up to $z = 1.5$

5,000 squares degrees (1/8 sky)

Data management at NCSA

DES specially designed to probe the origin of dark energy

S/G class and photo- z needed

1 TB of data per day

2 years completed, 3 more to go

Large Synoptic Survey Telescope

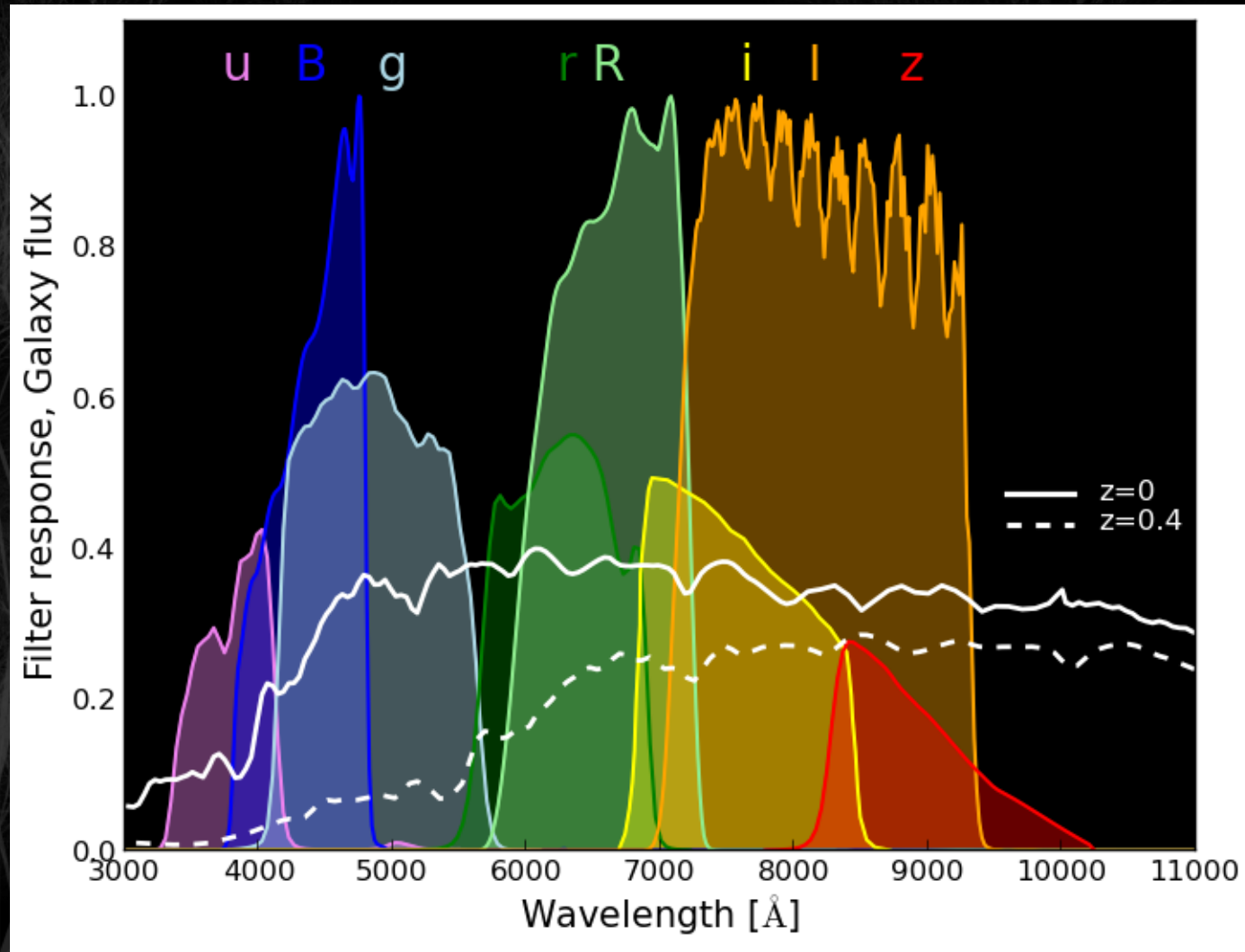
~ 2020 first light

Half of sky

30 TB of data nightly for 10 years!!

NCSA involved in data analysis

Big Data Challenge, ~1B galaxies



Determine redshift using limited information.
8 points instead of thousands!

- Photo- z Probability Density Functions needed
- Several methods/codes to compute photo- z
- Need for a meta-algorithm that combines multiple techniques
- PDF are good but for large datasets, storage and I/O will be an issue
- Machine Learning and statistical tools

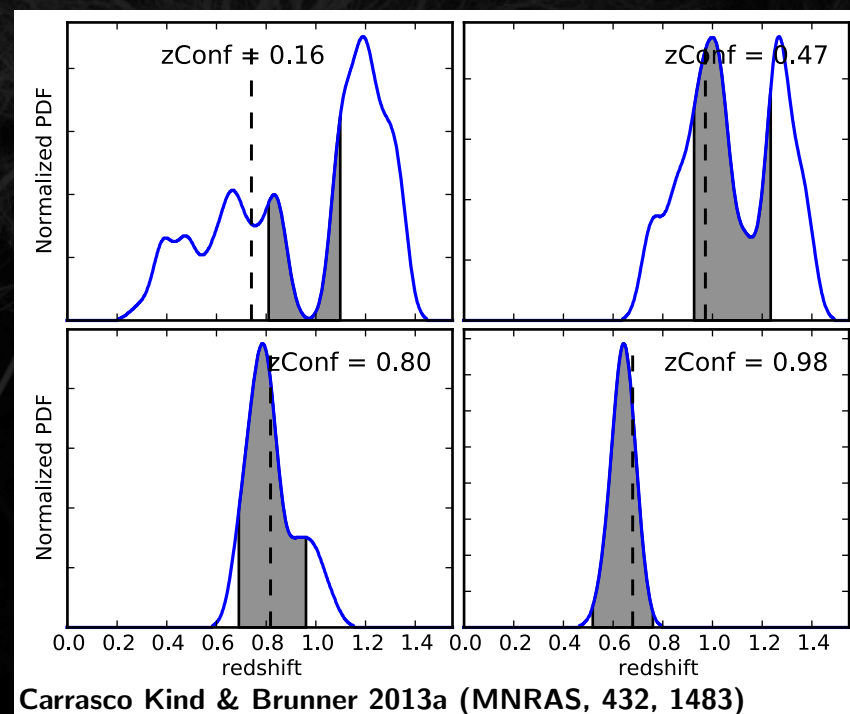
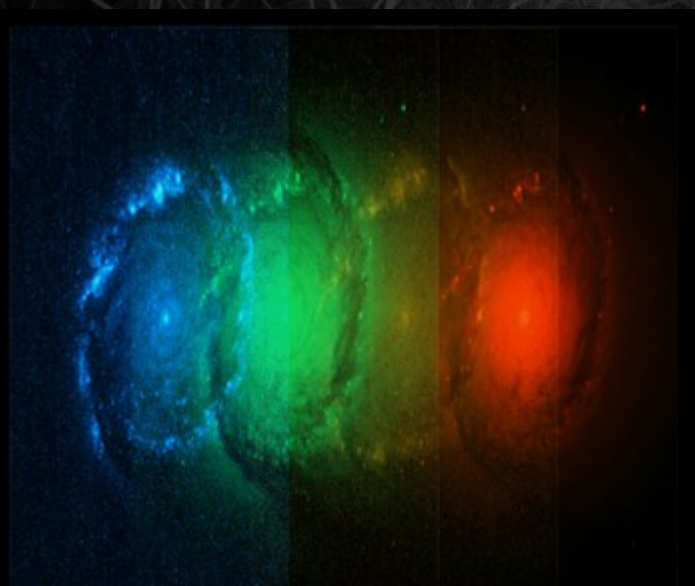
- Photo- z Probability Density Functions needed
- Several methods/codes to compute photo- z
- Need for a meta-algorithm that combines multiple techniques
- PDF are good but for large datasets, storage and I/O will be an issue
- Machine Learning and statistical tools

- Photo- z Probability Density Functions needed
- Several methods/codes to compute photo- z
- Need for a meta-algorithm that combines multiple techniques
- PDF are good but for large datasets, storage and I/O will be an issue
- Machine Learning and statistical tools

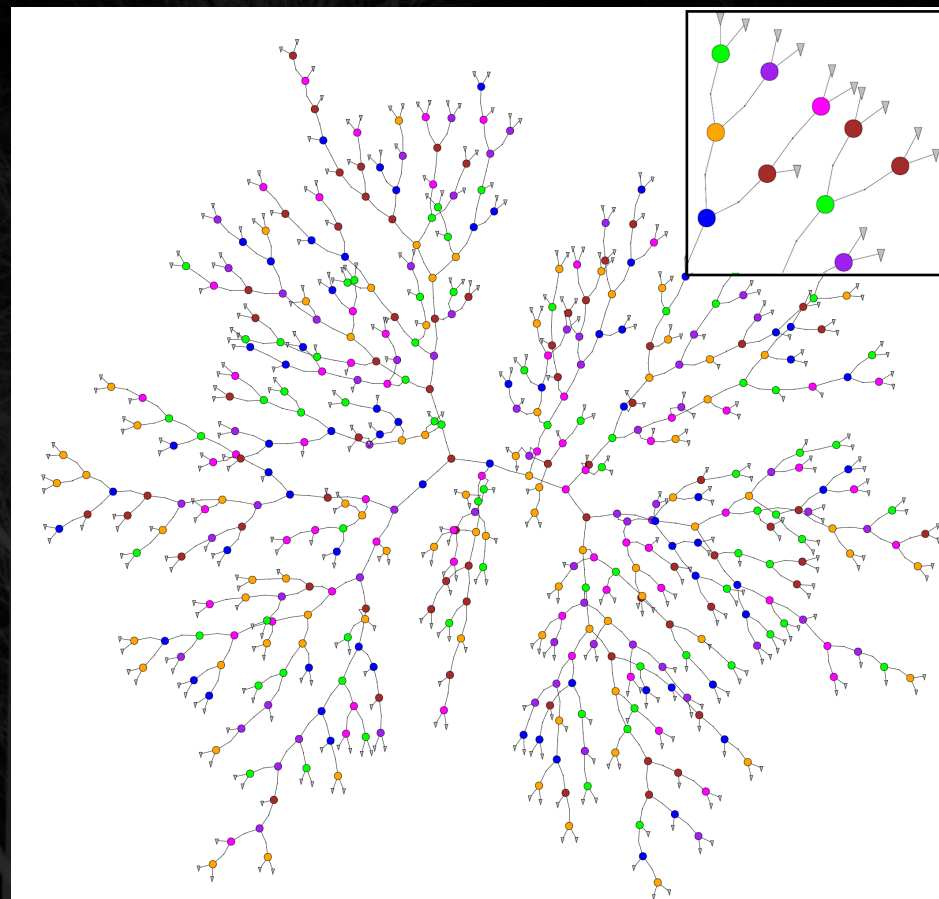
- Photo- z Probability Density Functions needed
- Several methods/codes to compute photo- z
- Need for a meta-algorithm that combines multiple techniques
- PDF are good but for large datasets, storage and I/O will be an issue
- Machine Learning and statistical tools

- Photo- z Probability Density Functions needed
- Several methods/codes to compute photo- z
- Need for a meta-algorithm that combines multiple techniques
- PDF are good but for large datasets, storage and I/O will be an issue
- Machine Learning and statistical tools

Photo- z PDF estimation (in 5 min.)



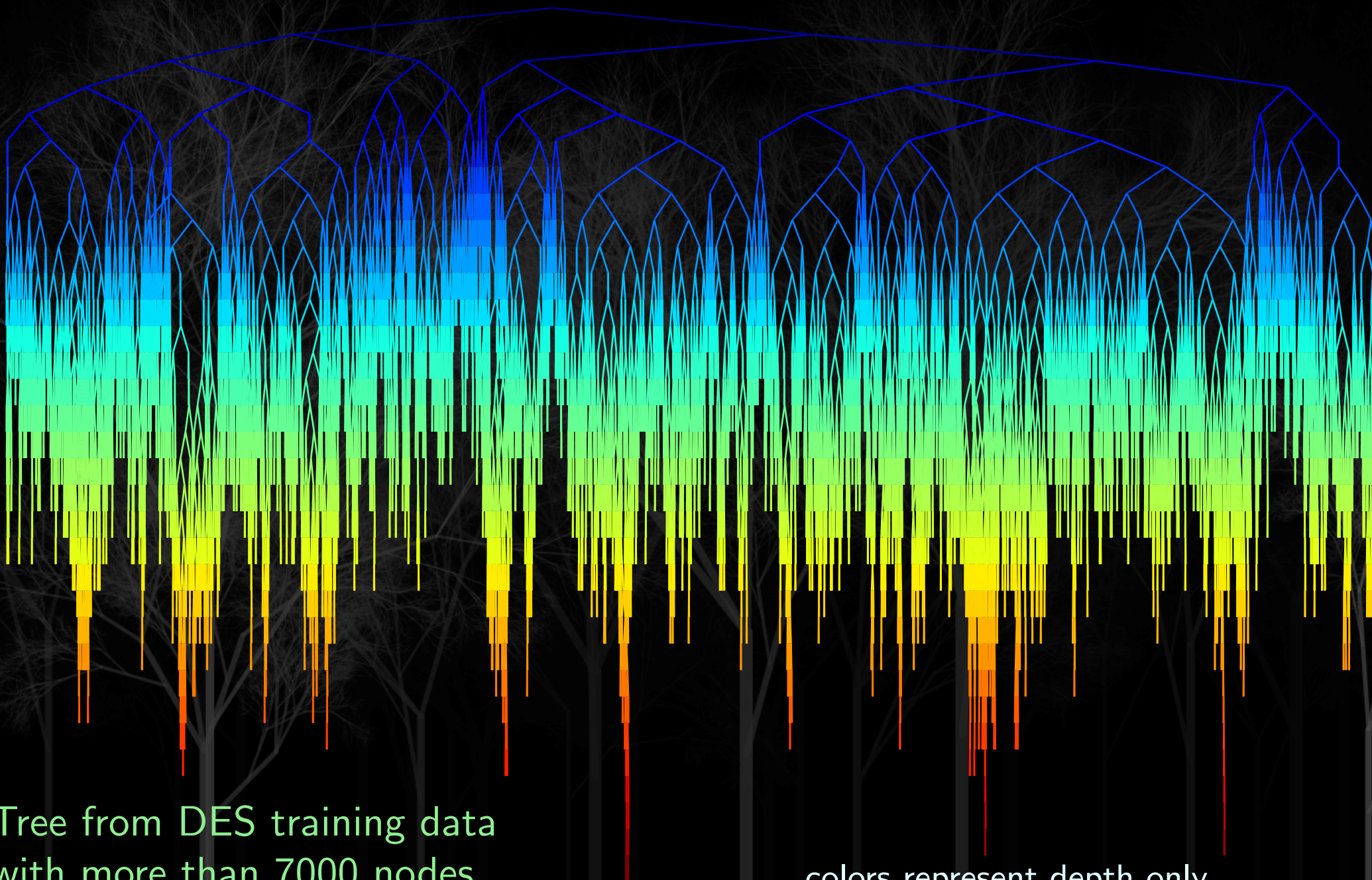
- TPZ (Trees for Photo-Z) is a supervised machine learning code
- Prediction trees and random forest
- Incorporate measurements errors and deals with missing values
- Ancillary information: expected errors, attribute ranking and others
- Application to the S/G

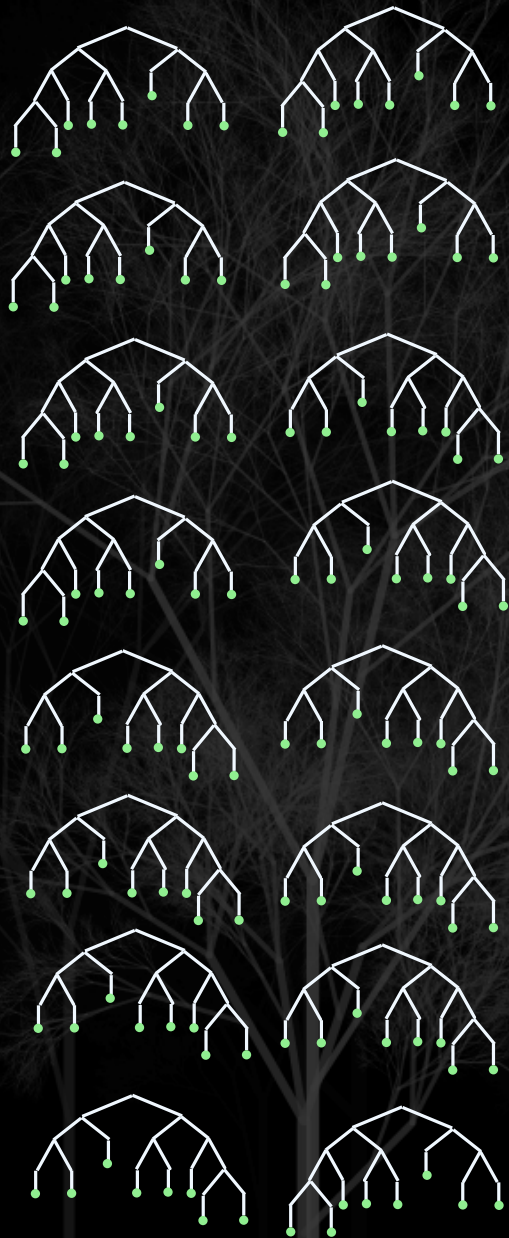


Carrasco Kind & Brunner 2013a (MNRAS, 432, 1483)

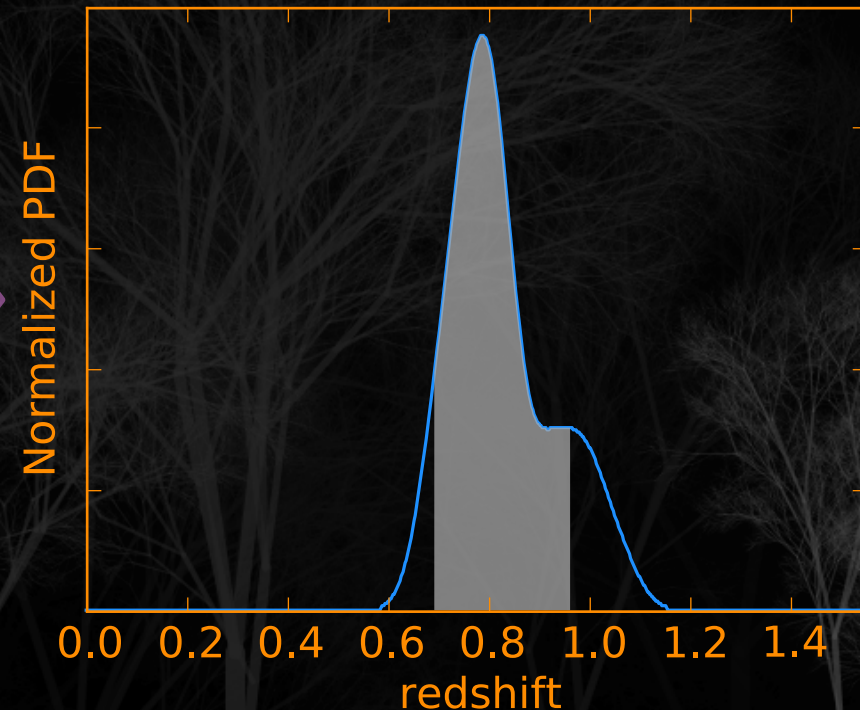
<http://lcdm.astro.illinois.edu/code/mlz.html>

Photo- z PDF estimation: TPZ example

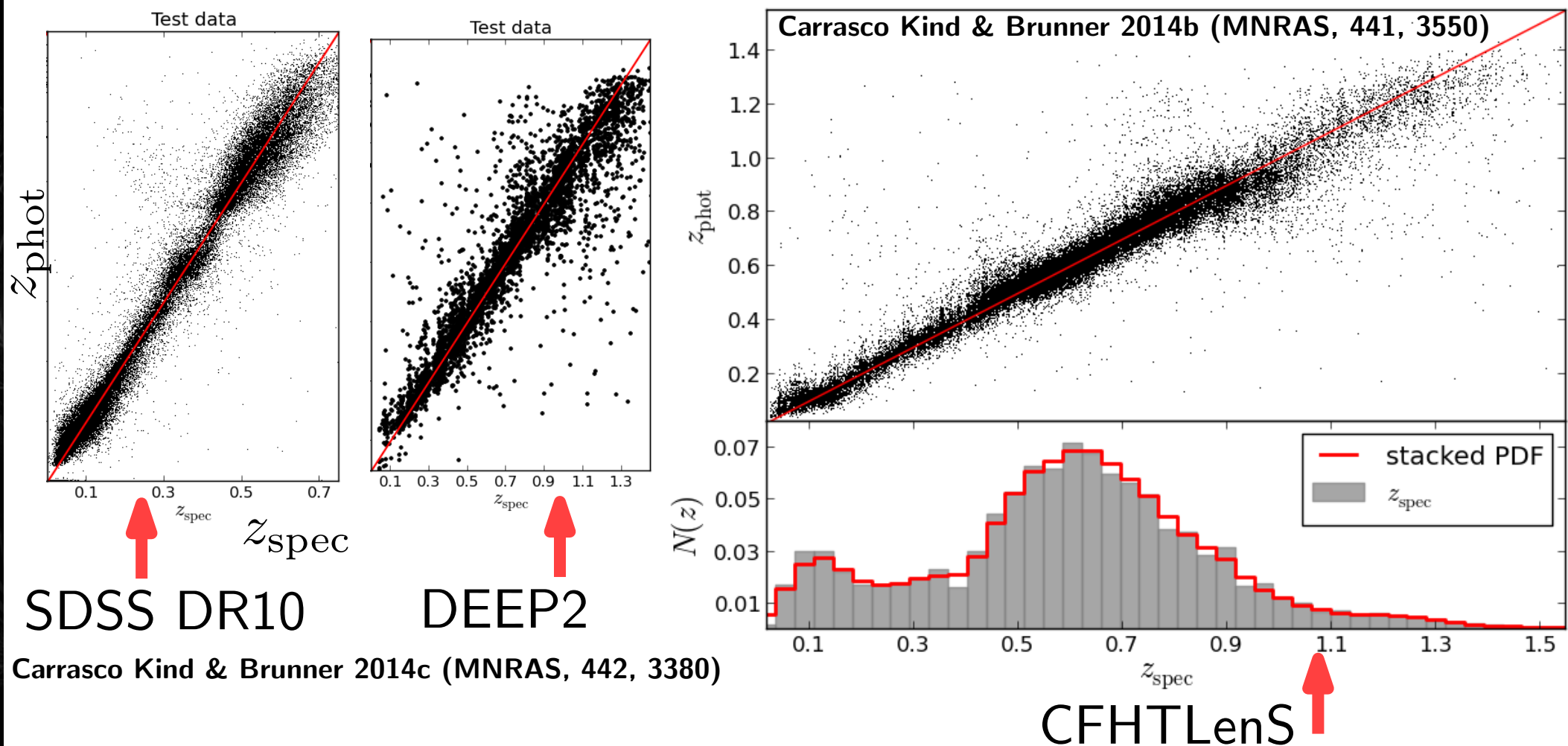




Combine predictions
from trees

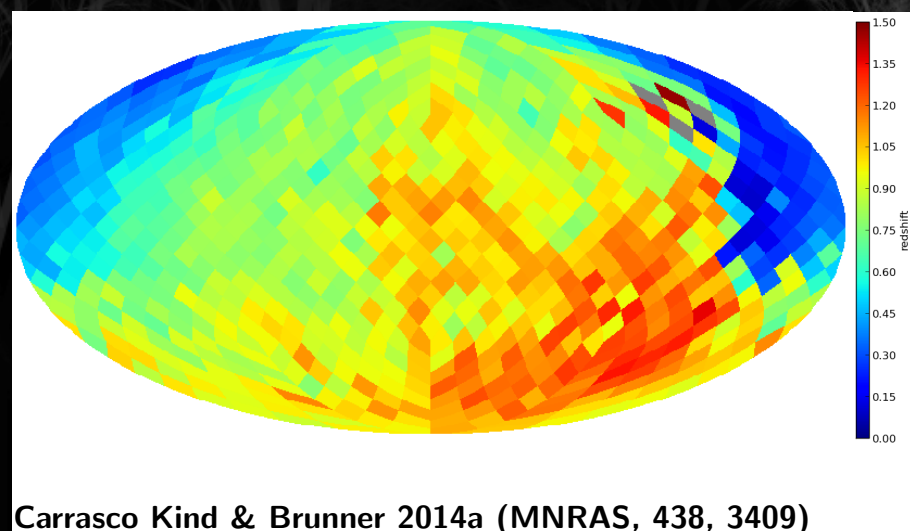
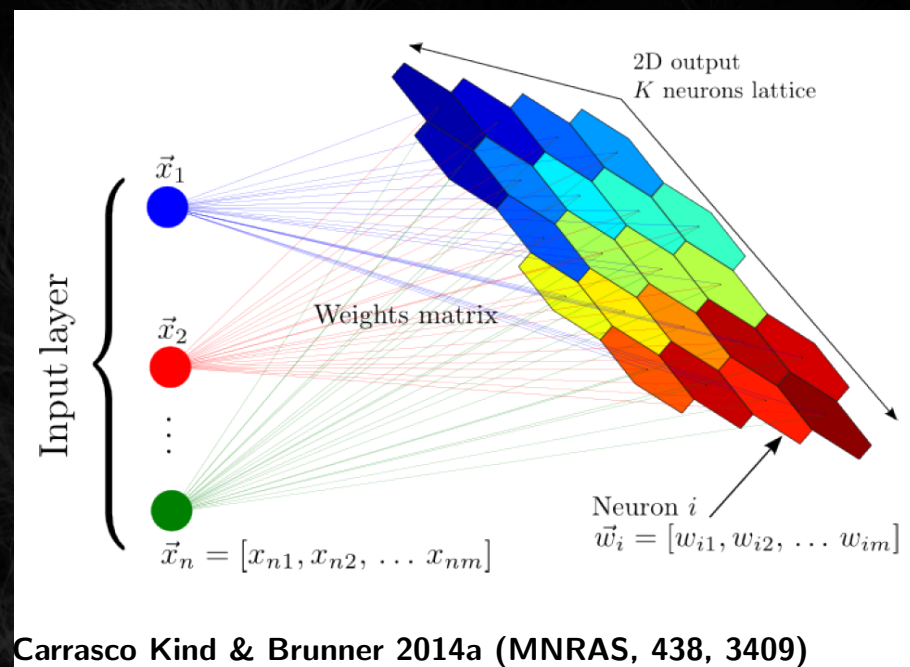


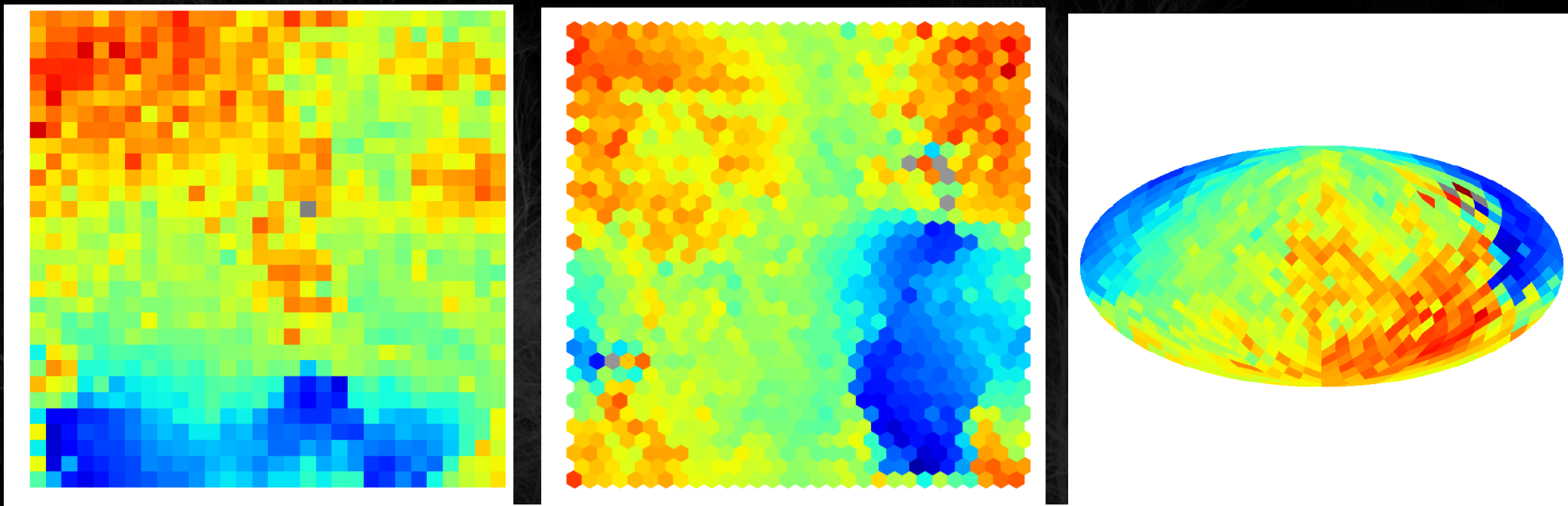
Trees are ideally uncorrelated and strong
Bootstrapping and error sampling
Random features at each node



TPZ has been tested in several databases with remarkable results

- SOM(Self Organized Map) is a unsupervised machine learning algorithm
- Competitive learning to represent data conserving topology
- 2D maps and *Random Atlas*
- Framework inherited from TPZ
- Application to the S/G

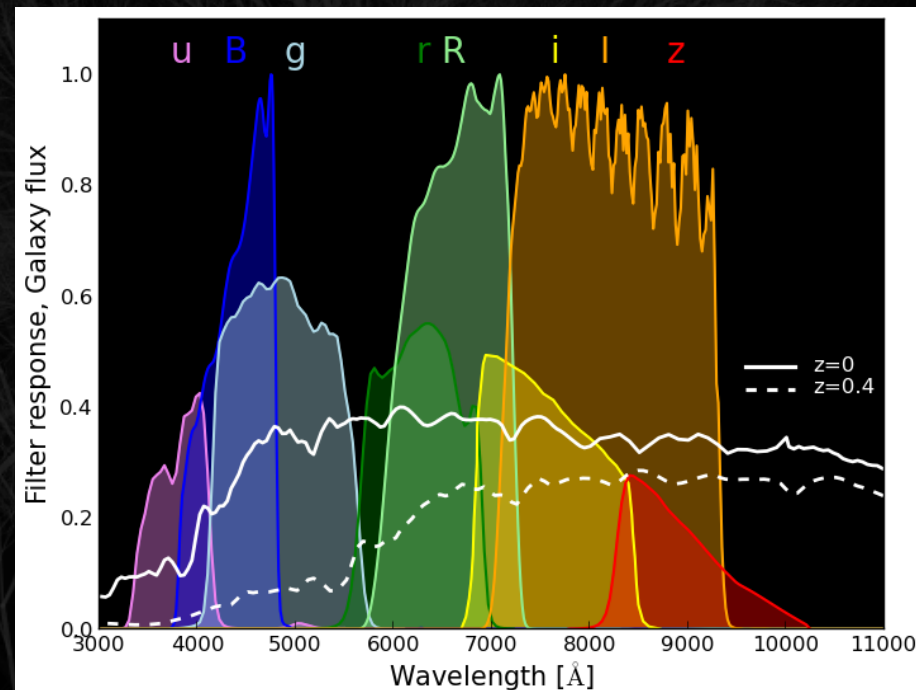




Carrasco Kind & Brunner 2014a (MNRAS, 438, 3409)

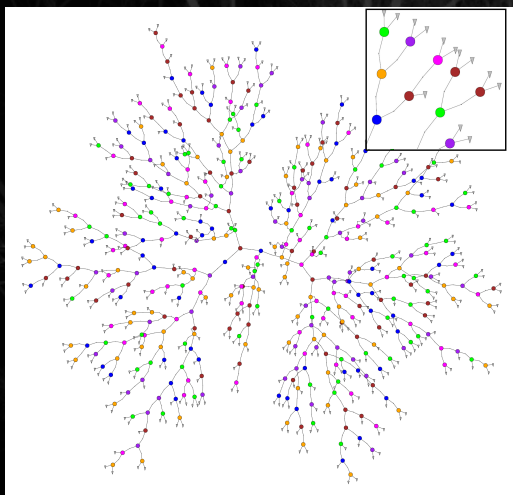
Different topologies can be used with or without periodic boundary conditions

- BPZ (Benitez, 2000) is a Bayesian template fitting method to obtain PDFs
- Set of calibrated SED and filters
- Doesn't need training data
- Priors can be included

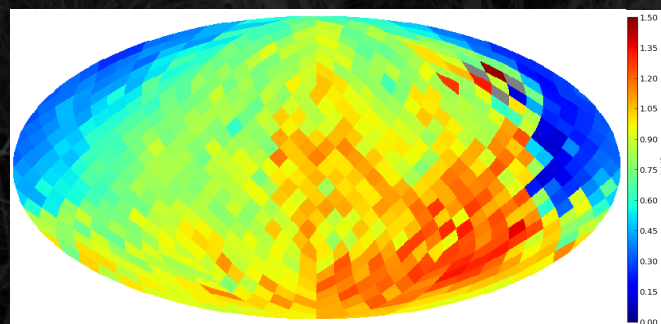


Carrasco Kind & Brunner 2014c (MNRAS, 442, 3380)

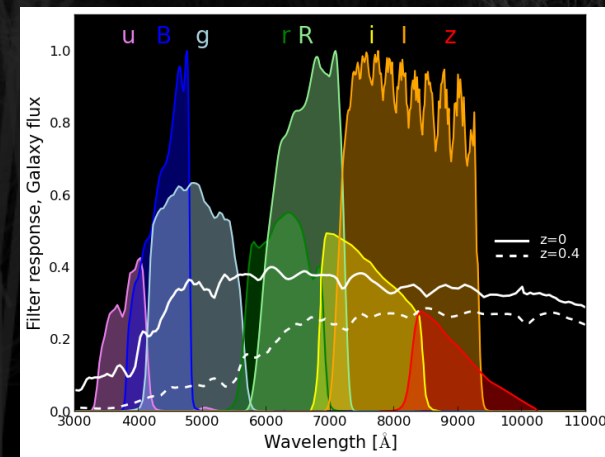
Photo- z PDF combination

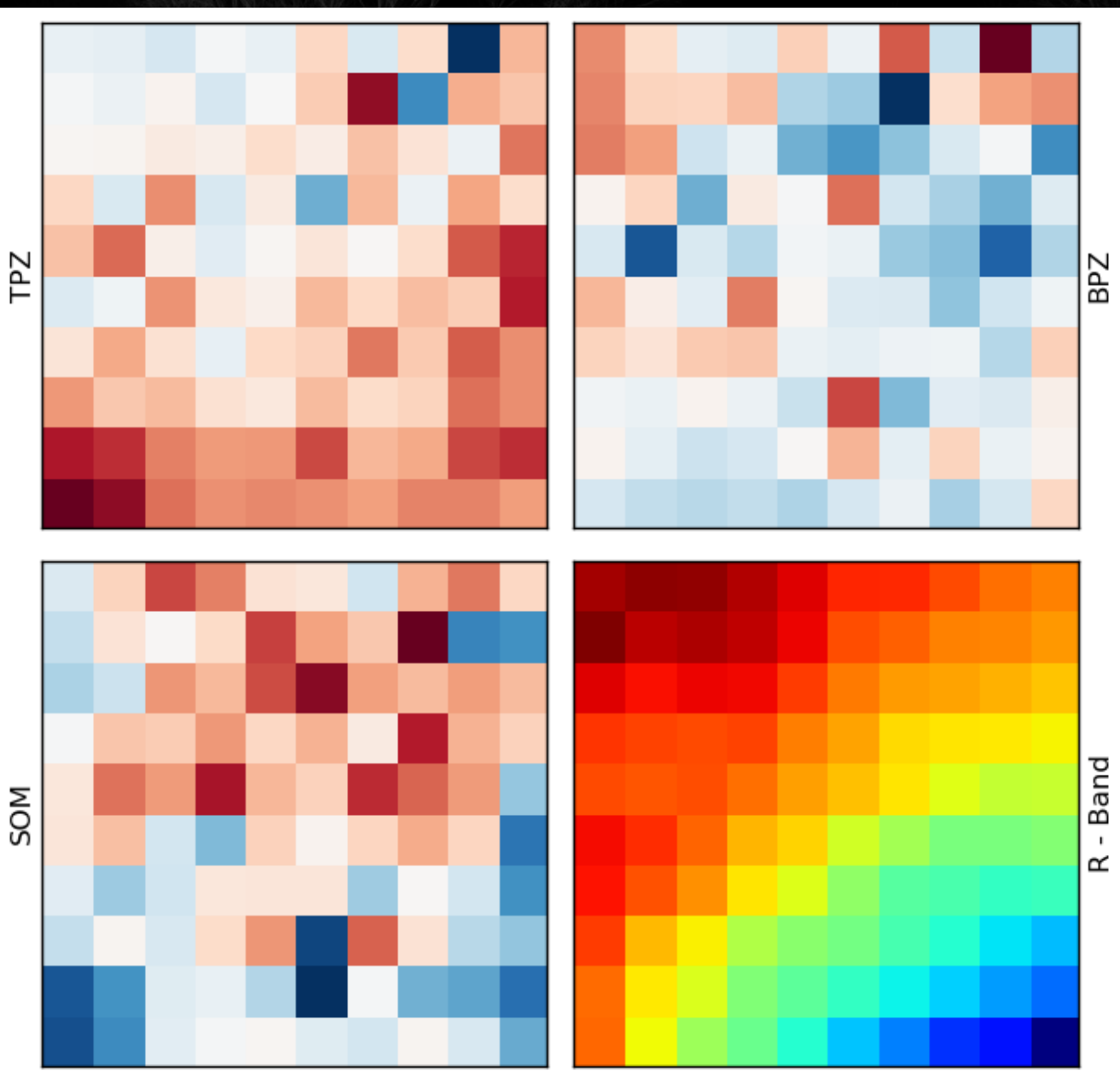


+

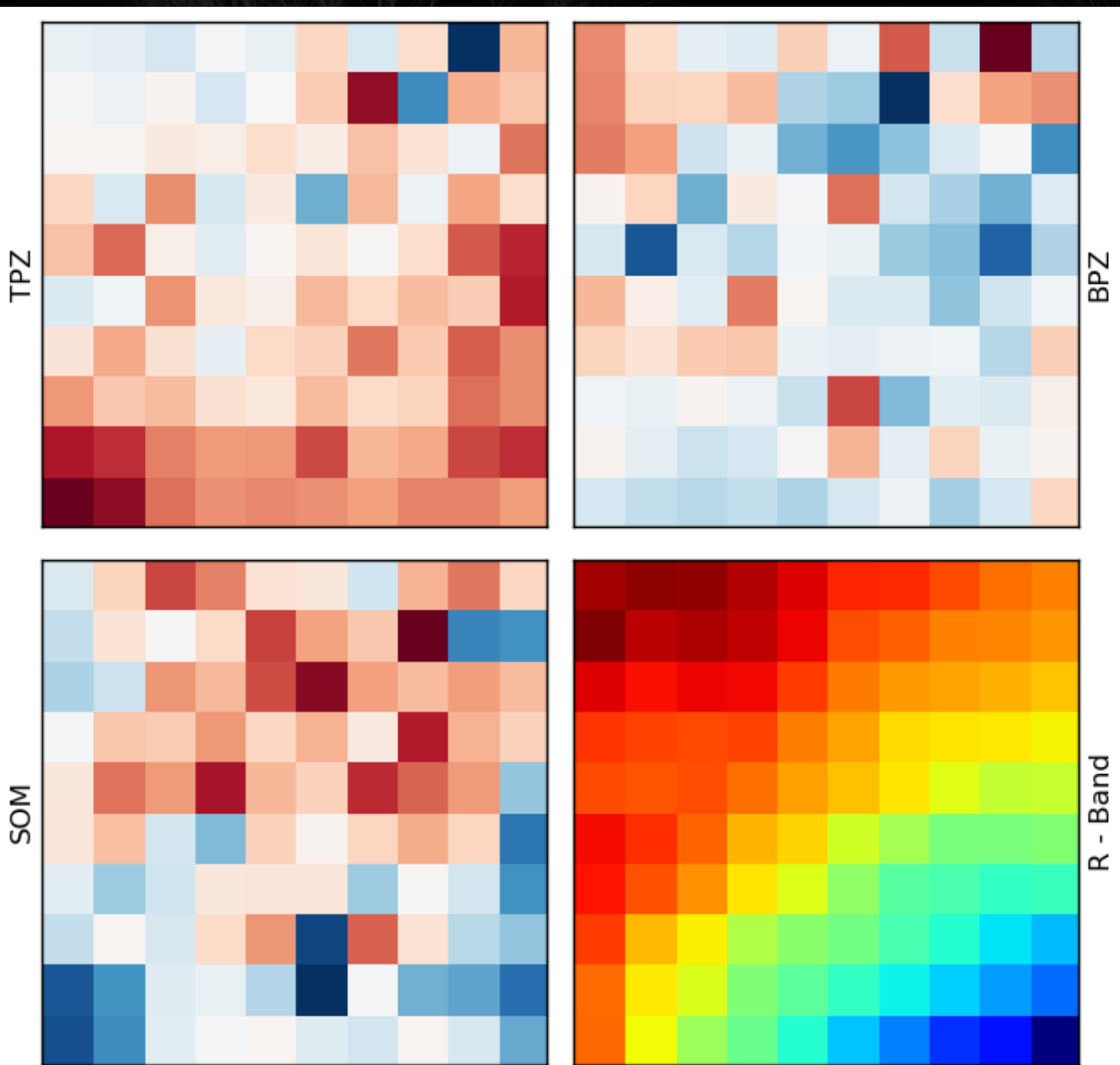


+

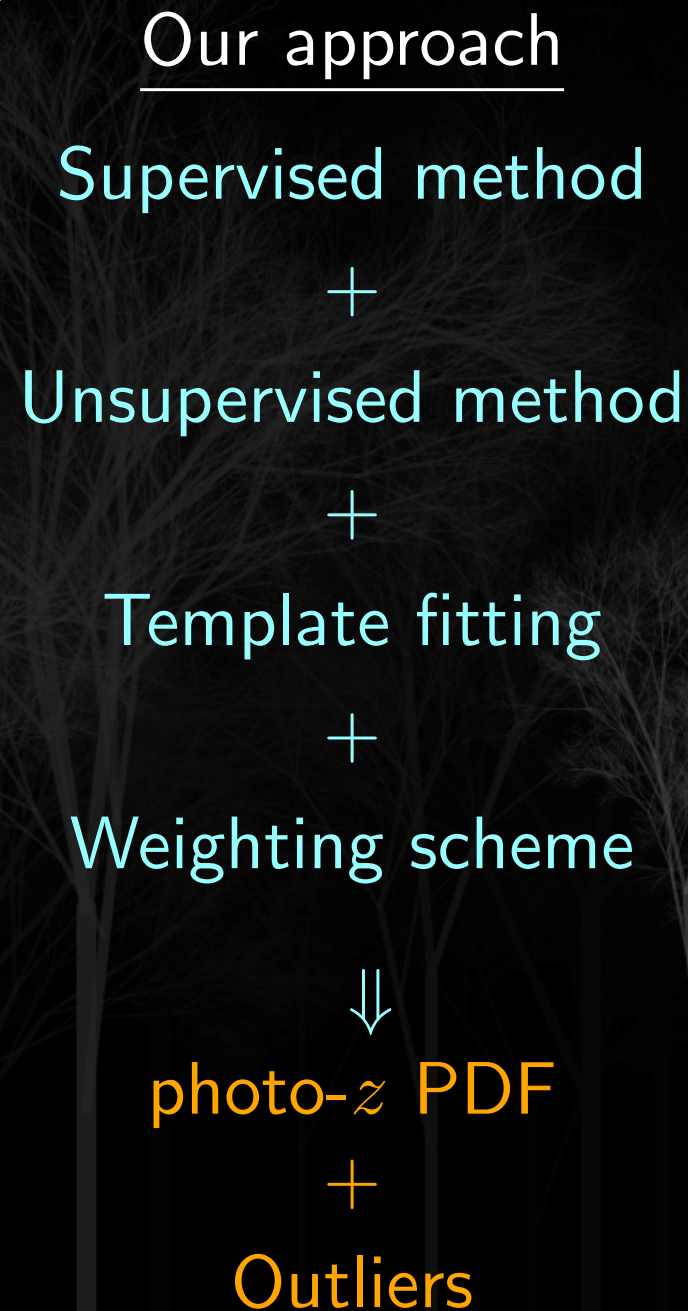


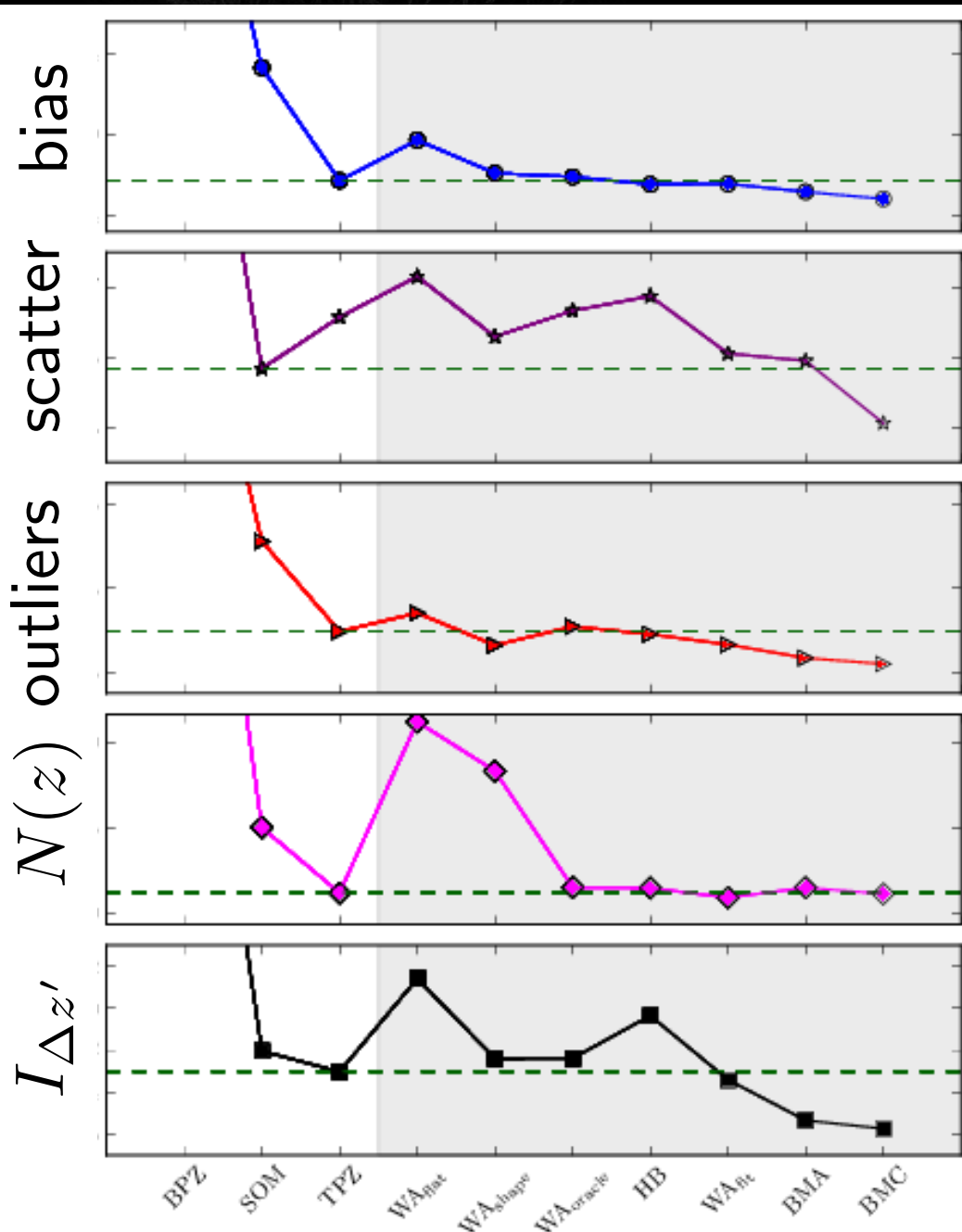


Carrasco Kind & Brunner 2014c (MNRAS, 442, 3380)



Carrasco Kind & Brunner 2014c (MNRAS, 442, 3380)





- Several combination methods
- Bayesian model averaging (BMA) and combination (BMC) are the best
- Same applies to S/G (Kim, Brunner & CK in prep.)

Naïve Bayes Classifier (same used for spam emails) to identify "spam" galaxies using information from multiple techniques

Each feature provides information about these two classes, and can be combined to make a stronger classifier

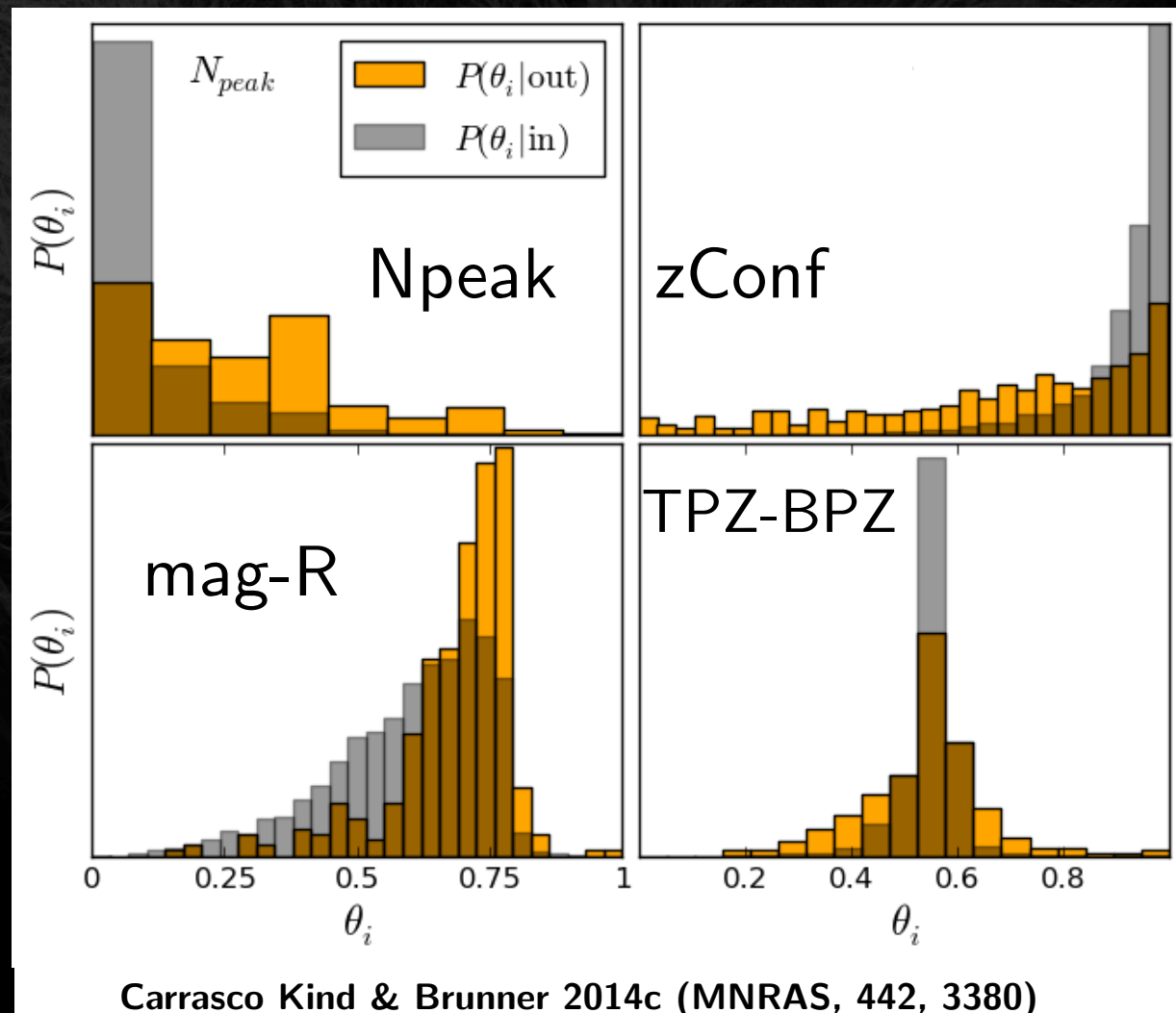
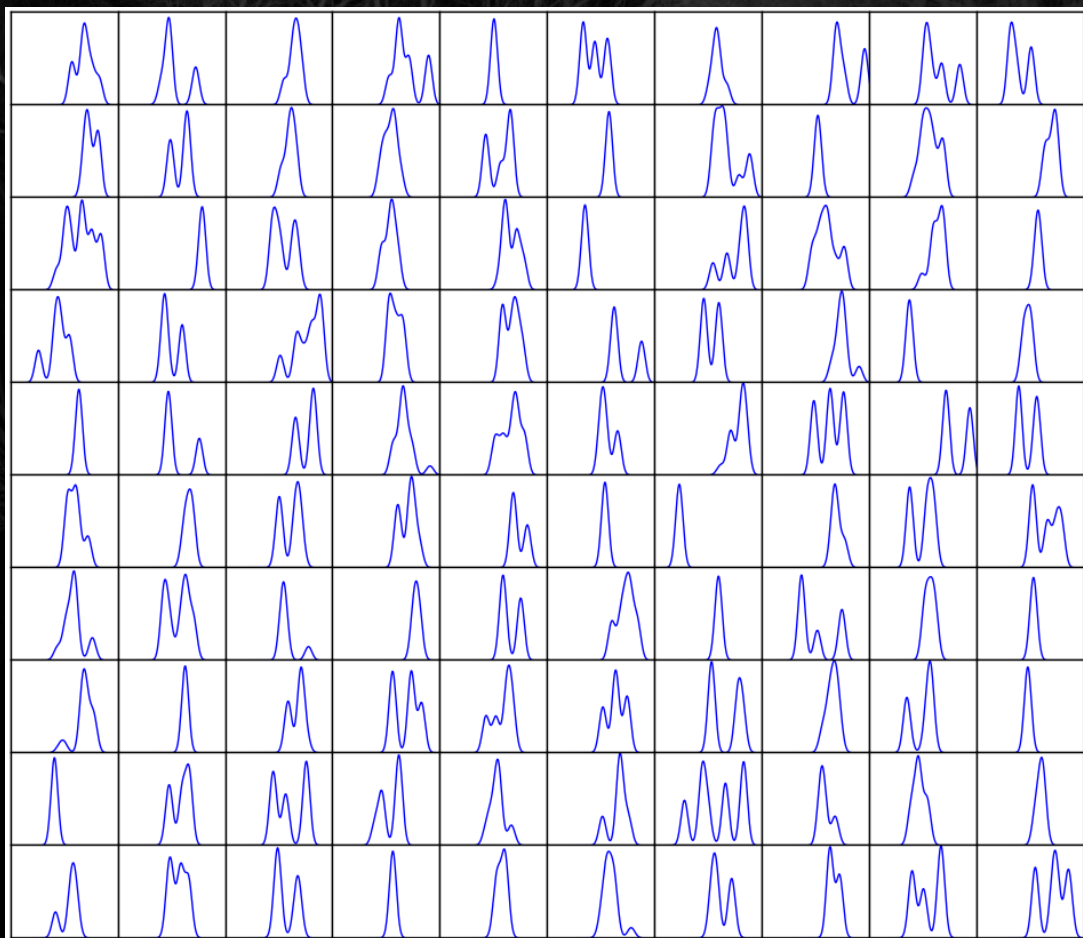


Photo- z PDF representation and storage



Single Gaussian fit

Multi-Gaussian fit

Monte Carlo sampling

Sparse representation
techniques

Reduce number of points
while increasing accuracy

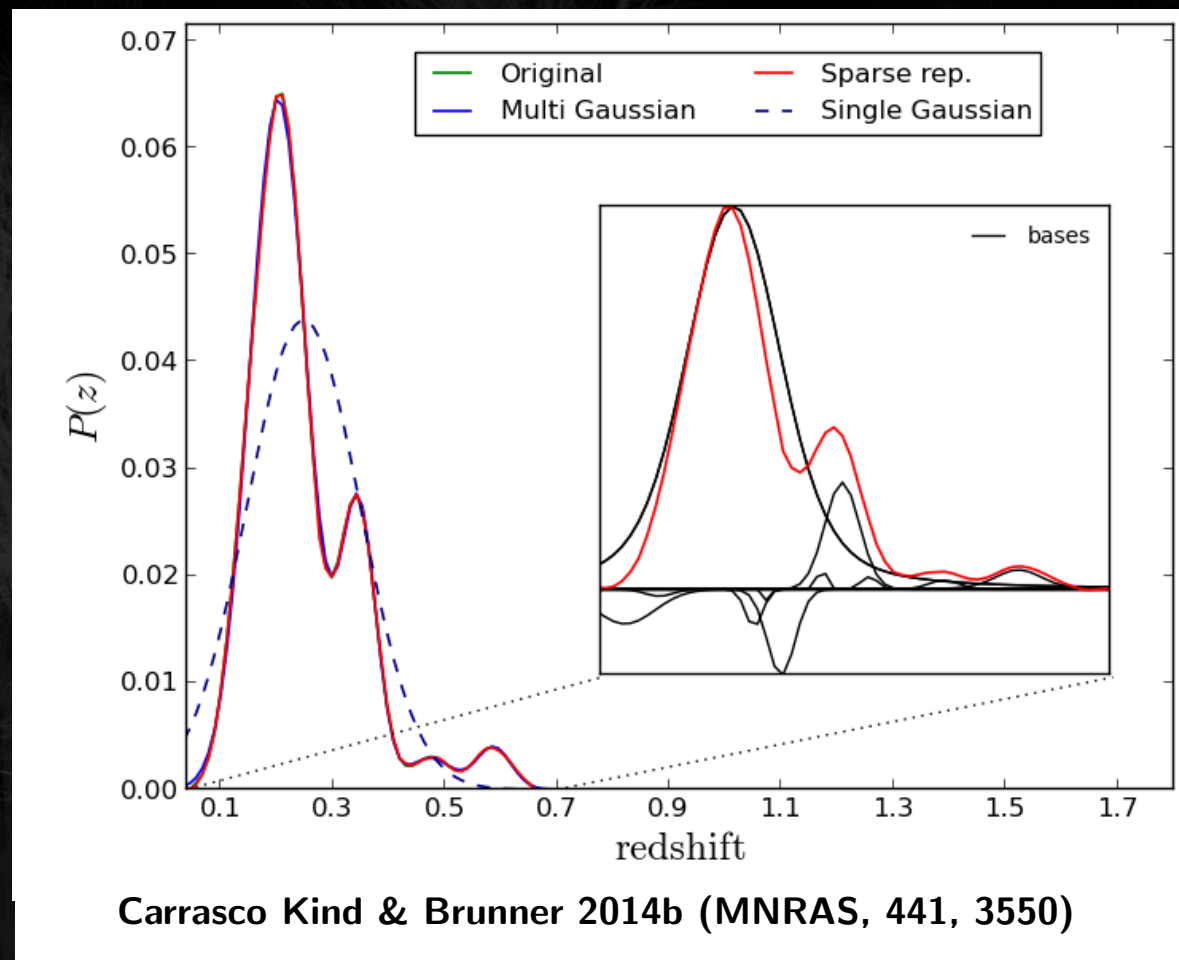
Single Gaussian fit

Multi-Gaussian fit

Monte Carlo sampling

Sparse representation techniques

Reduce number of points while increasing accuracy

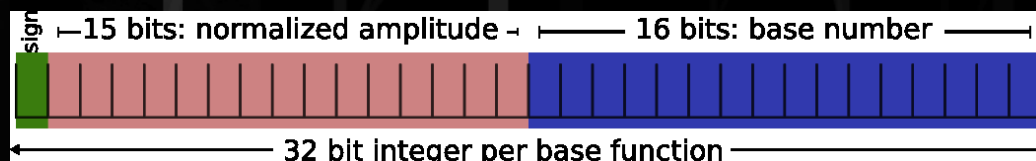
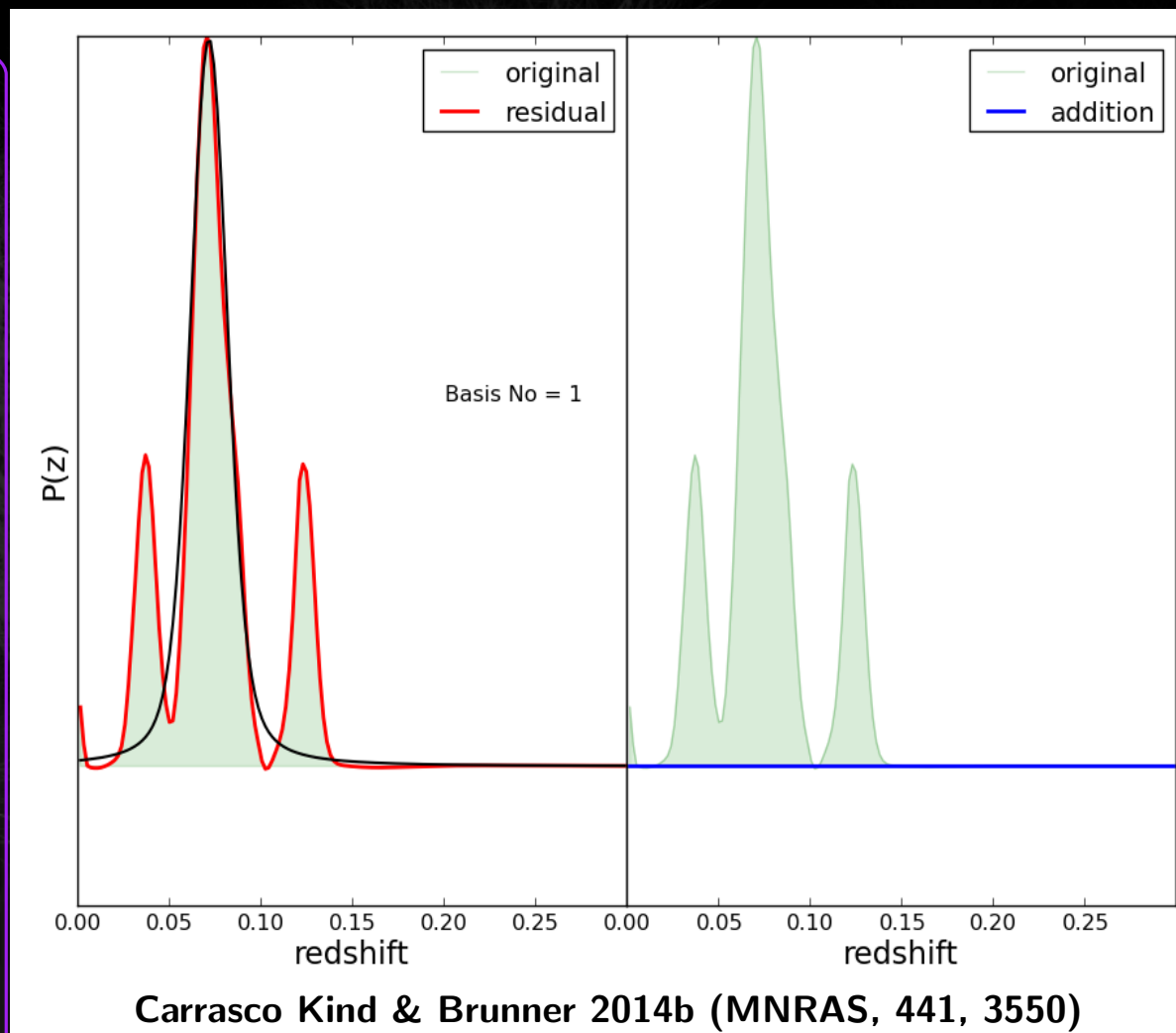


Use Gaussian and Voigt profiles as bases, need N_{original}^2 bases

With only 10-20 bases achieve 99.9 % accuracy

Use 32-bits integer per basis, compression

Store Multiple PDFs

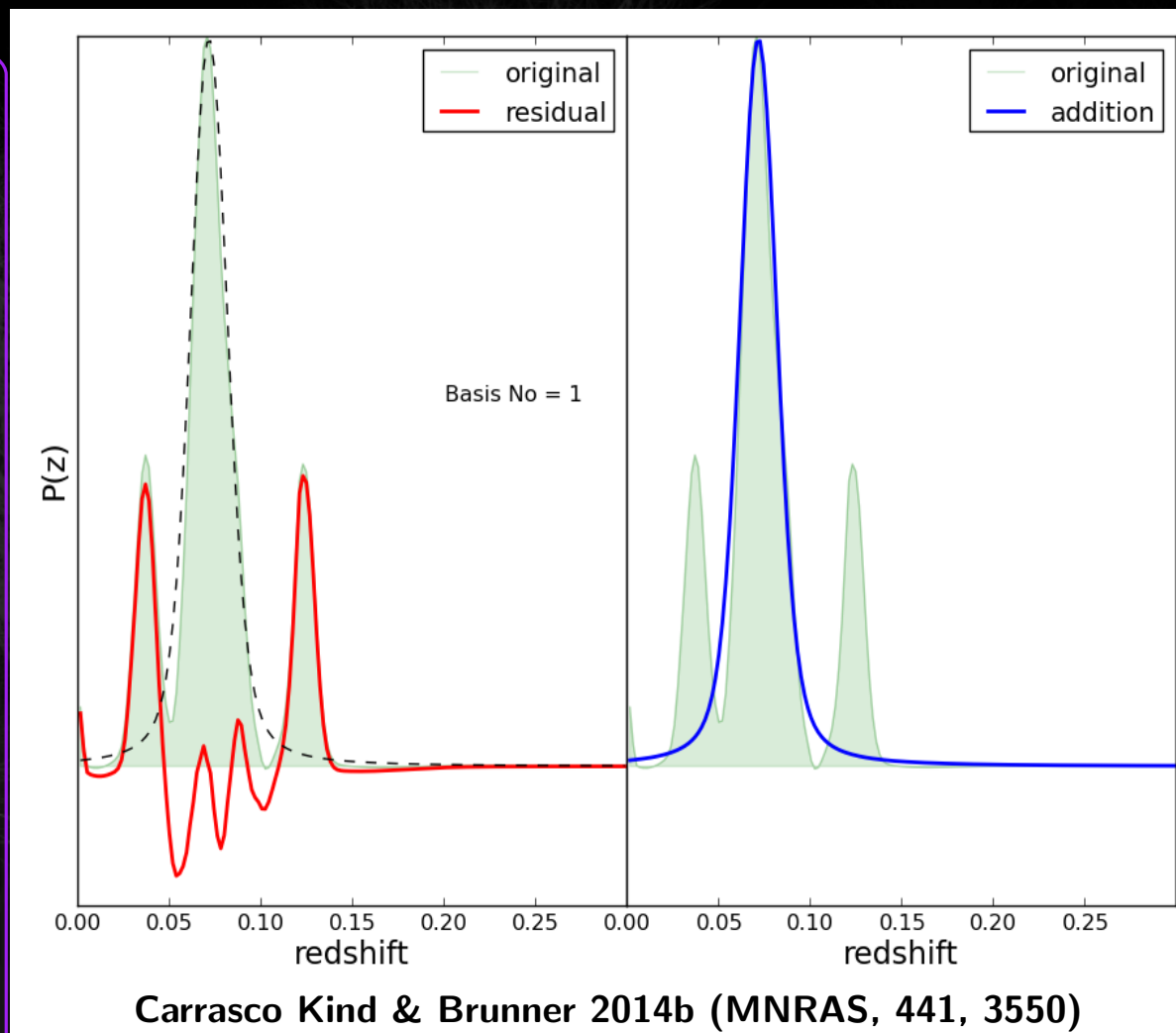


Use Gaussian and Voigt profiles as bases, need N_{original}^2 bases

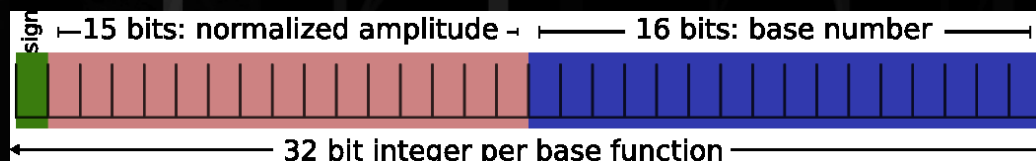
With only 10-20 bases achieve 99.9 % accuracy

Use 32-bits integer per basis, compression

Store Multiple PDFs



Carrasco Kind & Brunner 2014b (MNRAS, 441, 3550)

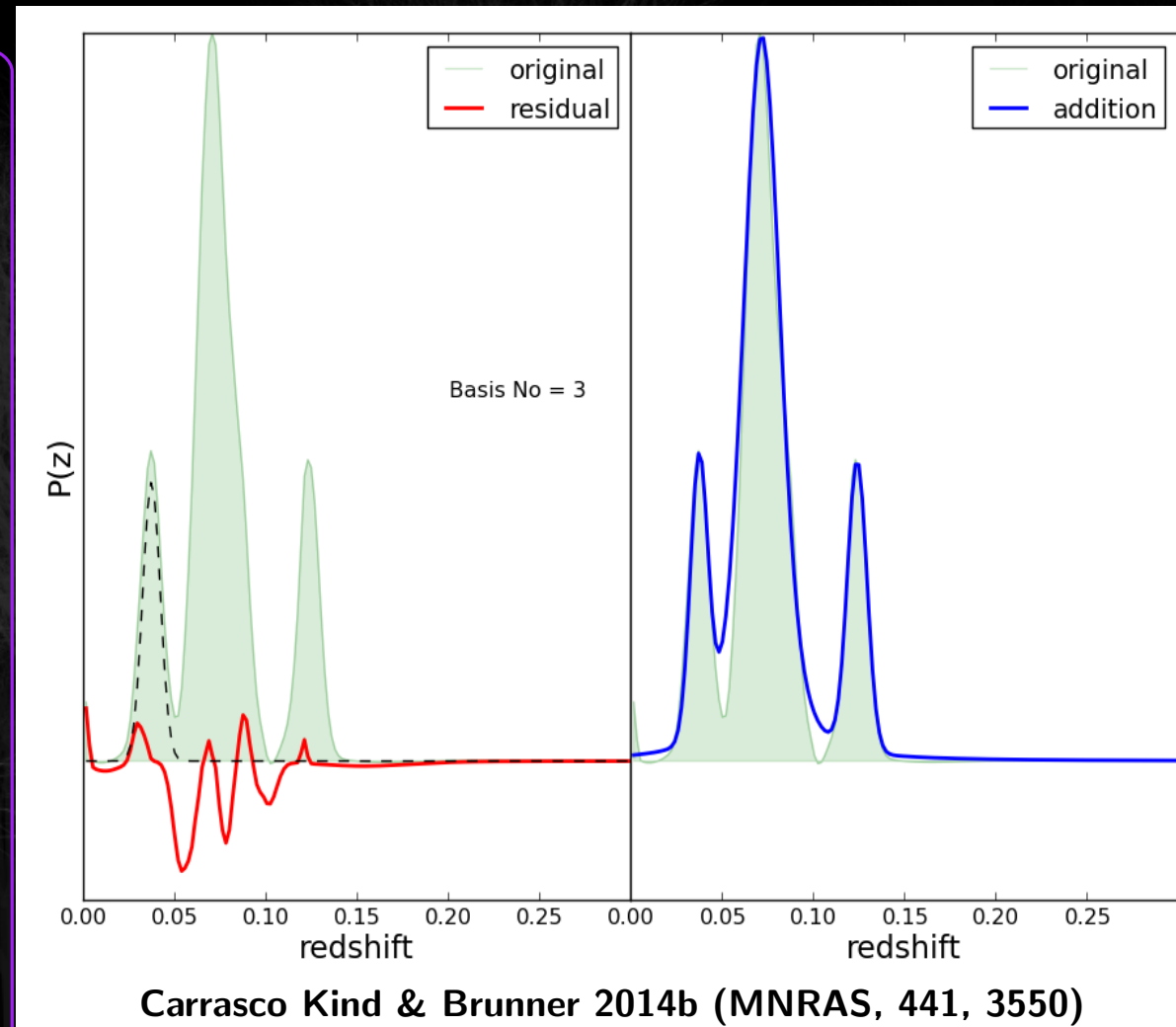


Use Gaussian and Voigt profiles as bases, need N_{original}^2 bases

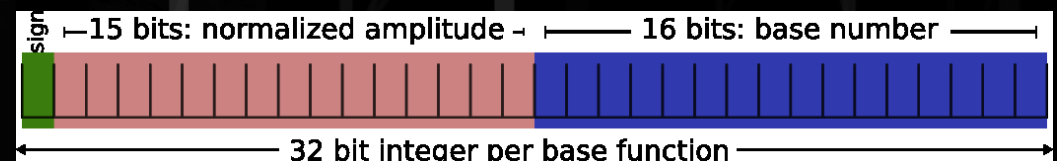
With only 10-20 bases achieve 99.9 % accuracy

Use 32-bits integer per basis, compression

Store Multiple PDFs



Carrasco Kind & Brunner 2014b (MNRAS, 441, 3550)

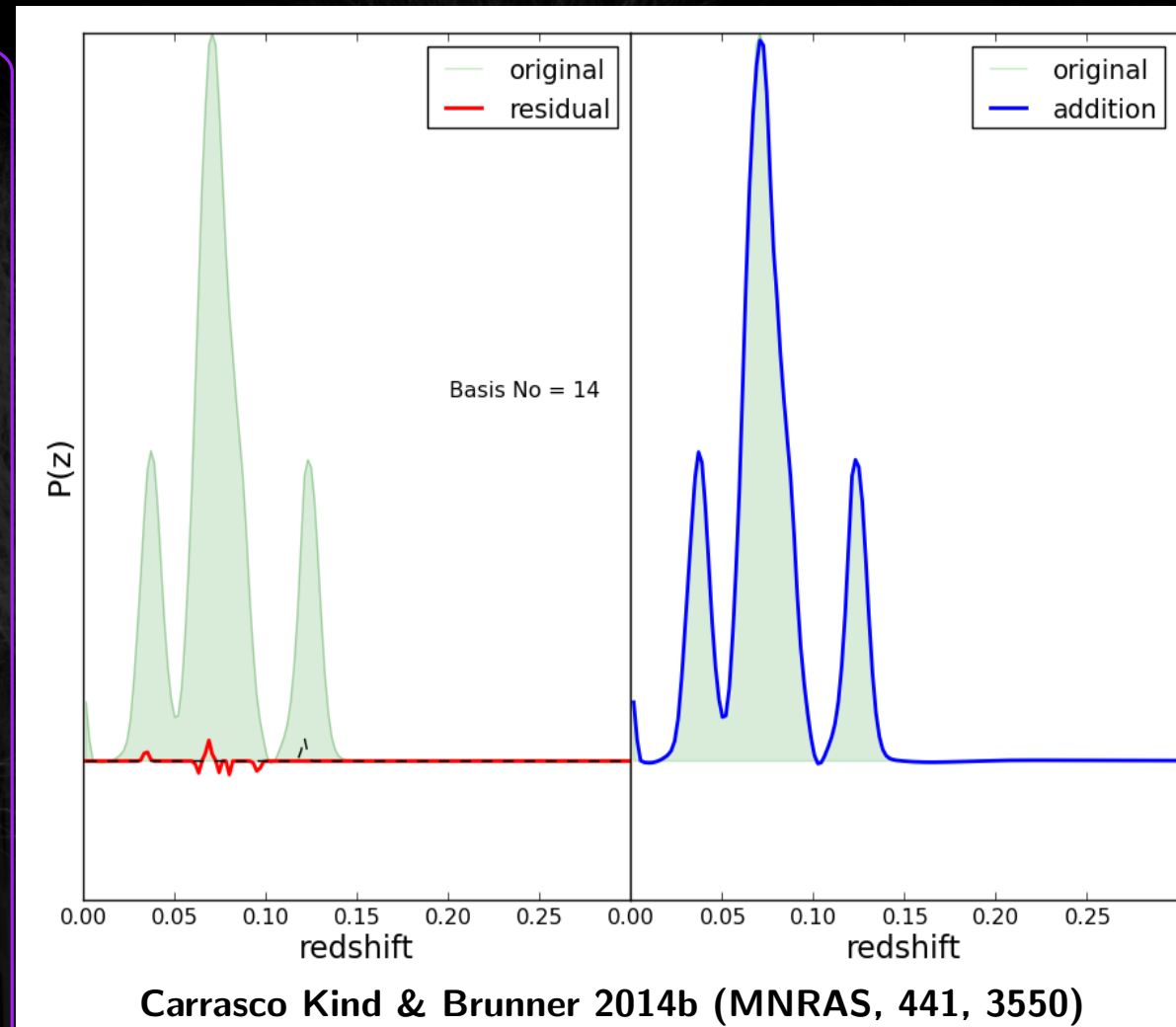


Use Gaussian and Voigt profiles as bases, need N_{original}^2 bases

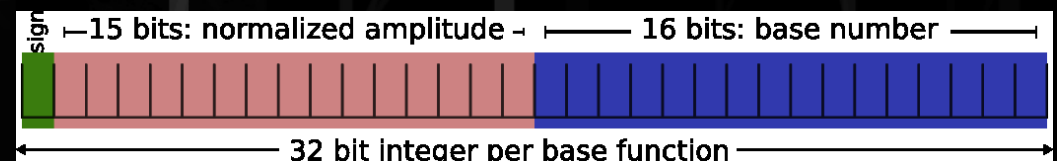
With only 10-20 bases achieve 99.9 % accuracy

Use 32-bits integer per basis, compression

Store Multiple PDFs



Carrasco Kind & Brunner 2014b (MNRAS, 441, 3550)



By definition:

$$N(z) = \sum_{k=1}^N \int_{z-\Delta z/2}^{z+\Delta z/2} P_k(z) dz$$

By definition:

$$N(z) = \sum_{k=1}^N \int_{z-\Delta z/2}^{z+\Delta z/2} P_k(z) dz$$

Using sparse representation, we represent each PDF p_{z_k} as:

$\mathbf{p}_{z_k} \approx \mathbf{D} \cdot \boldsymbol{\delta}_k$ \mathbf{D} is the dictionary, $\boldsymbol{\delta}_k$ is the sparse vector, then

$$N(z) = \sum_{k=1}^N \boldsymbol{\delta}_k \cdot \int_{z-\Delta z/2}^{z+\Delta z/2} \mathbf{D} dz$$

Only bases are integrated

by precomputing:

$$\boldsymbol{\delta}_N = \sum_{k=1}^N \boldsymbol{\delta}_k \quad \mathbf{I}_D(z) = \int_{z-\Delta z/2}^{z+\Delta z/2} \mathbf{d}_j dz \quad j = 1, 2, \dots, m$$

By definition:

$$N(z) = \sum_{k=1}^N \int_{z-\Delta z/2}^{z+\Delta z/2} P_k(z) dz$$

Using sparse representation, we represent each PDF p_{z_k} as:

$\mathbf{p}_{z_k} \approx \mathbf{D} \cdot \boldsymbol{\delta}_k$ \mathbf{D} is the dictionary, $\boldsymbol{\delta}_k$ is the sparse vector, then

$$N(z) = \sum_{k=1}^N \boldsymbol{\delta}_k \cdot \int_{z-\Delta z/2}^{z+\Delta z/2} \mathbf{D} dz$$

Only bases are integrated

by precomputing:

$$\boldsymbol{\delta}_N = \sum_{k=1}^N \boldsymbol{\delta}_k \quad \mathbf{I}_D(z) = \int_{z-\Delta z/2}^{z+\Delta z/2} \mathbf{d}_j dz \quad j = 1, 2, \dots, m$$

$N(z)$ is reduce to a simple dot product

$$N(z) = \mathbf{I}_D(z) \cdot \boldsymbol{\delta}_N$$

- Machine Learning in DES
- Photo- z in DES early data
- Photo- z PDF in DESDM
- New tools to access these from DB

✓ Compute photo-z PDF

Individual techniques (MLZ; arXiv:1303.7269, arXiv:1312.5753)

✓ Combine PDFs efficiently

Better than individual, outliers identification (arXiv:1403.0044)

✓ PDF Sparse Representation

99.9% accuracy in $P(z)$ and $N(z)$ with 15 points (arXiv:1404.6442)

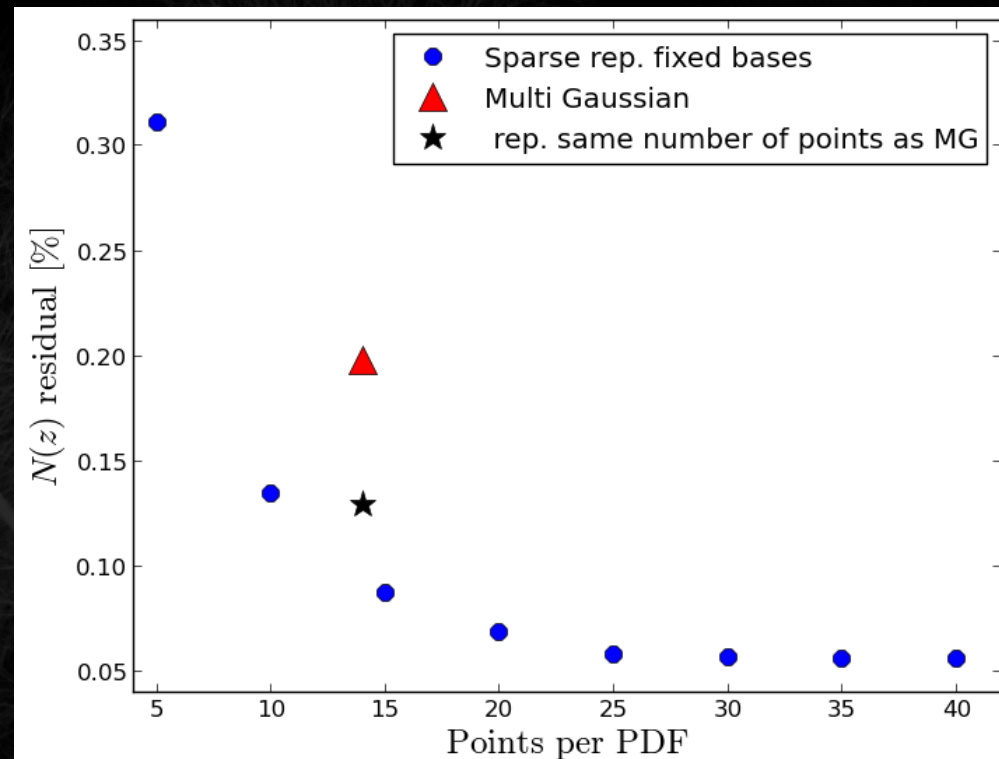
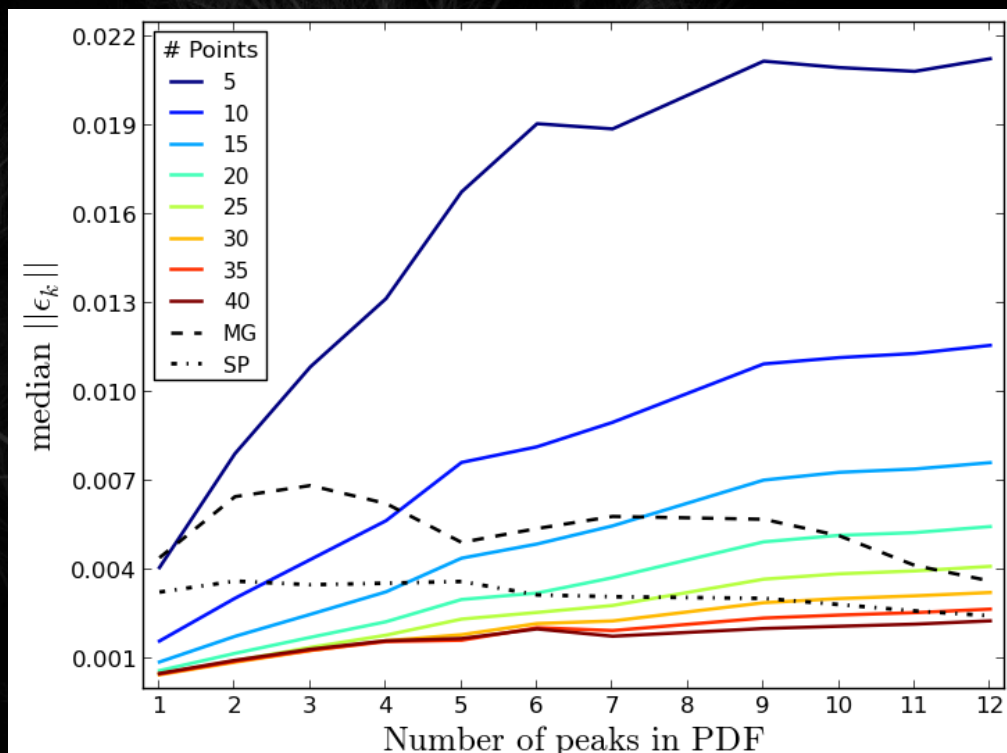
✓ Uses of these tools!

Clustering, weak lensing, DES, DESDM, etc...

Questions?

Matias Carrasco Kind
NCSA/UIUC
mcarras2@ncsa.illinois.edu
<http://matias-ck.com/>
<https://github.com/mgckind>

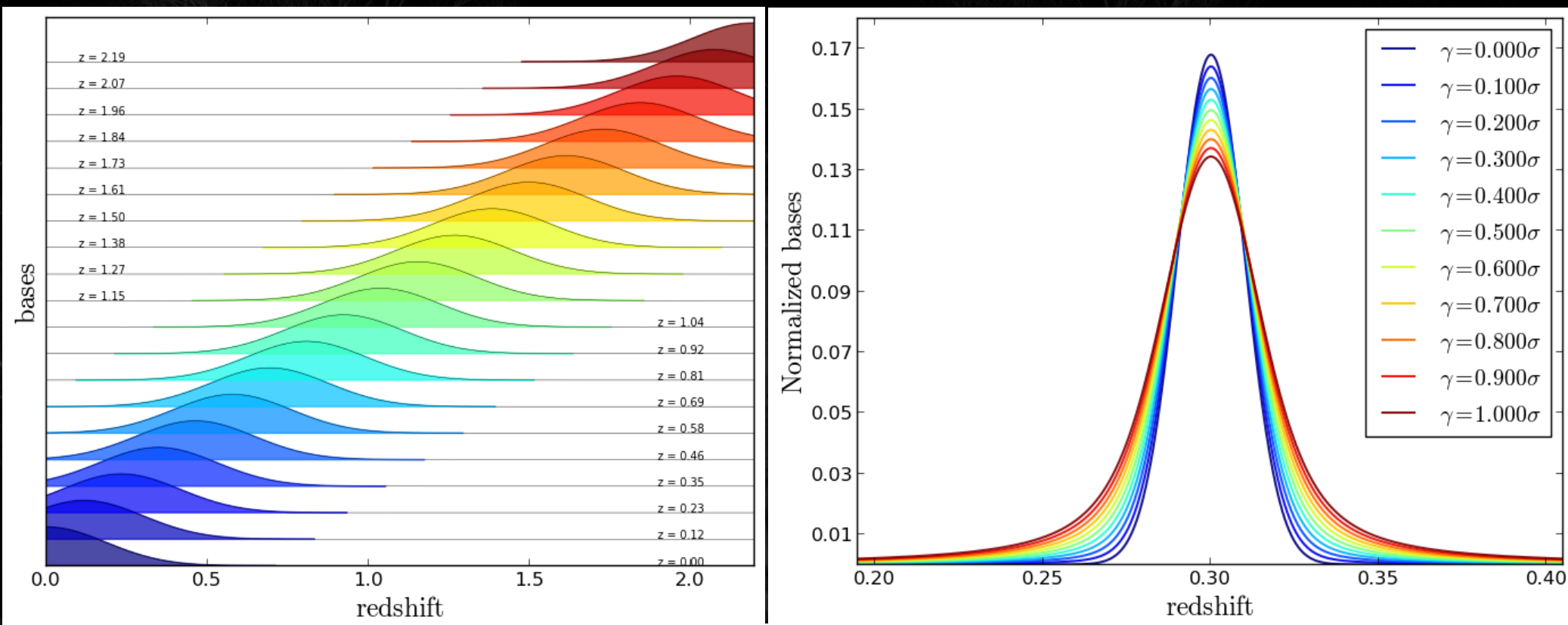




Carrasco Kind & Brunner 2014b (MNRAS, 441, 3550)

For PDFs with less than 4 peaks 5-10 points should be sufficient

Sparse representation gives more accurate and more compressed representation for $N(z)$, 99.9% accuracy with 15 points (200 points originally)



Carrasco Kind & Brunner 2014b (MNRAS, 441, 3550)

Combination of Gaussian and Voigt profiles

Covering the whole redshift space, at each location we have several bases

Out of Bag (cross-validation) data used to validate trees/maps

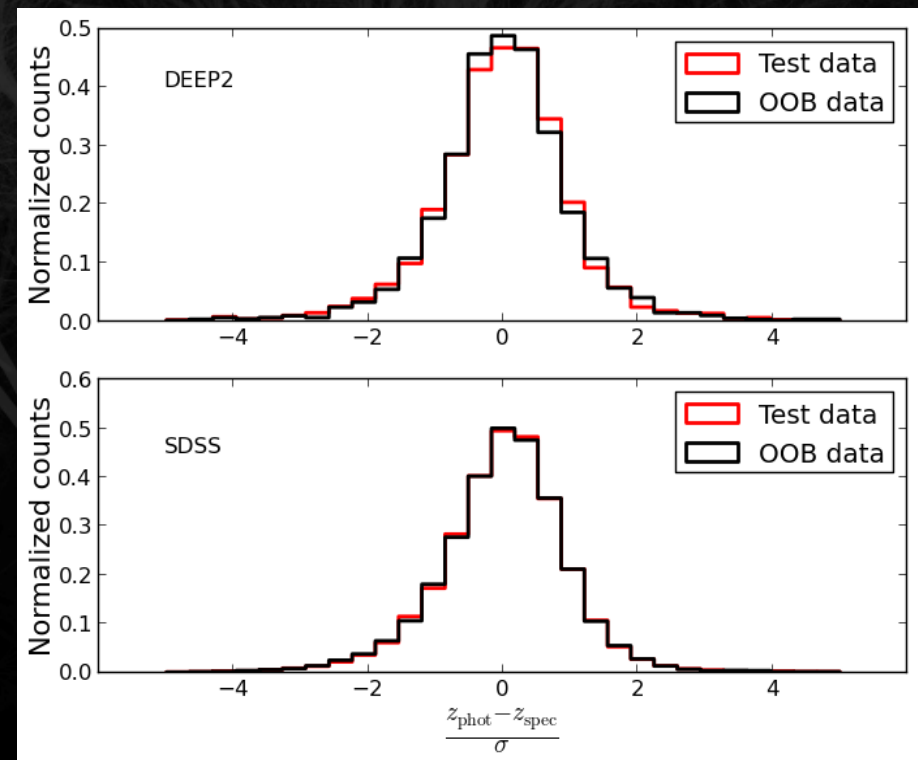
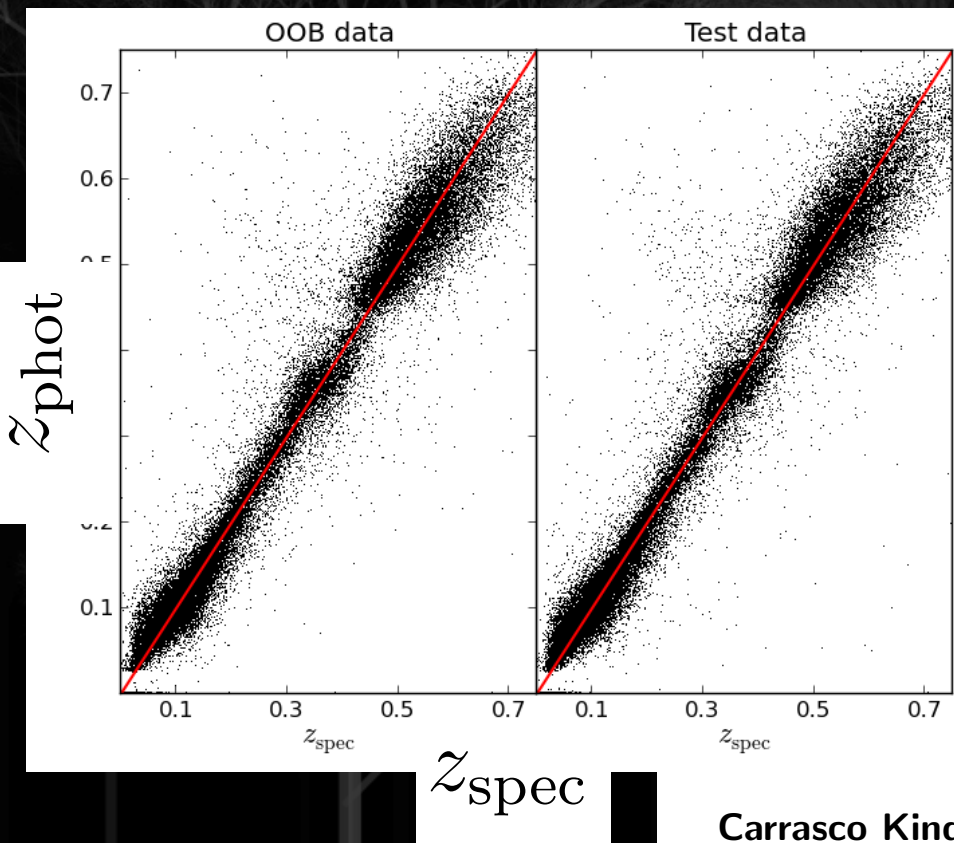
Changes for every tree/map and is not used during training

We can learn from the cross-validation data!

Out of Bag (cross-validation) data used to validate trees/maps

Changes for every tree/map and is not used during training

We can learn from the cross-validation data!

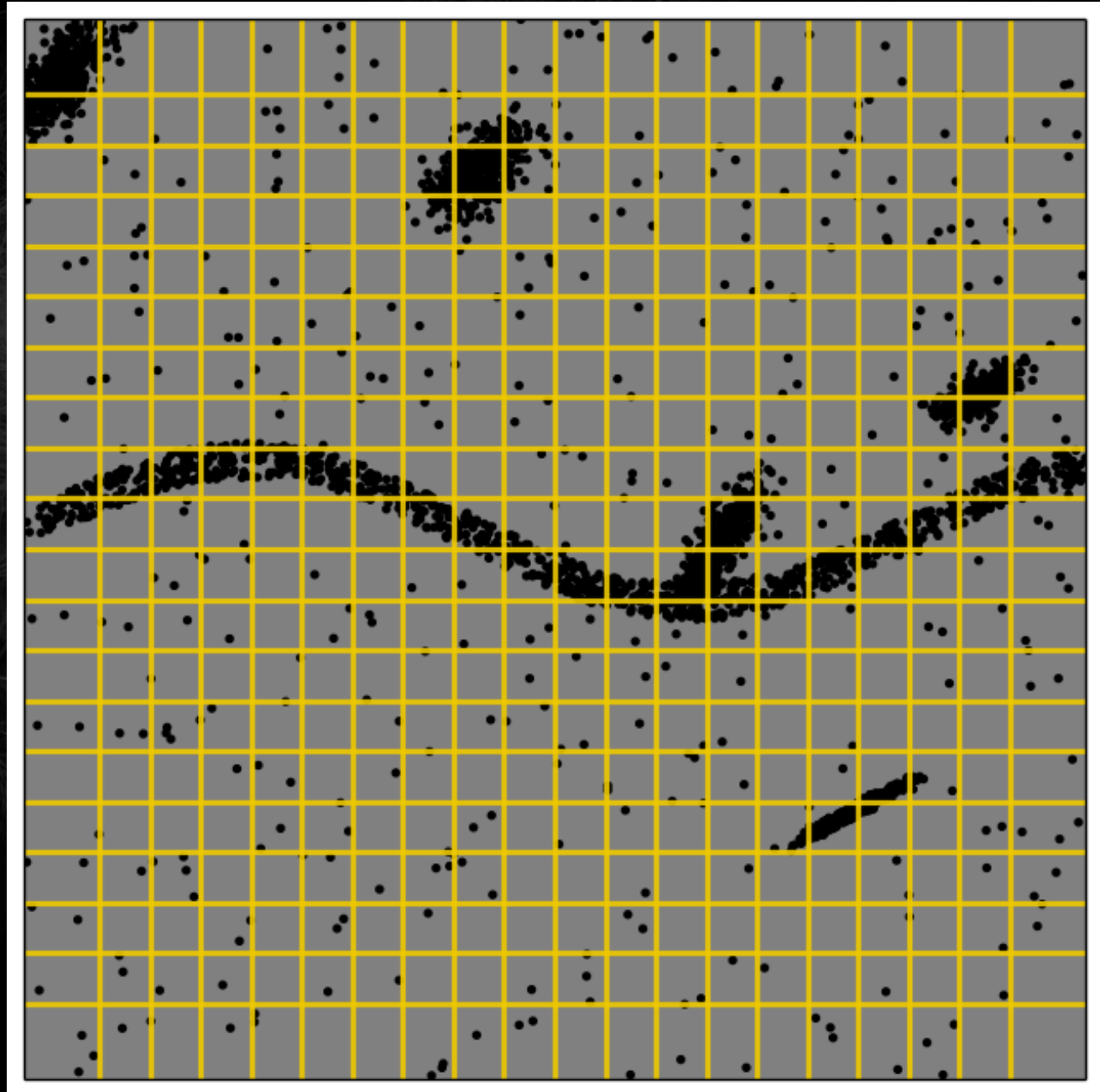


Carrasco Kind & Brunner 2014c (MNRAS, 442, 3380)

Suppose 2D data
distributed in a given
space

De-project the data
in a 2D map

Each cell will contain
objects with similar
properties



Suppose 2D data
distributed in a given
space

De-project the data
in a 2D map

Each cell will contain
objects with similar
properties

