



# Tales of photo- $z$ , DBs and DES

Matías Carrasco Kind

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

CEFCA Seminar

October 23<sup>th</sup>, 2016

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

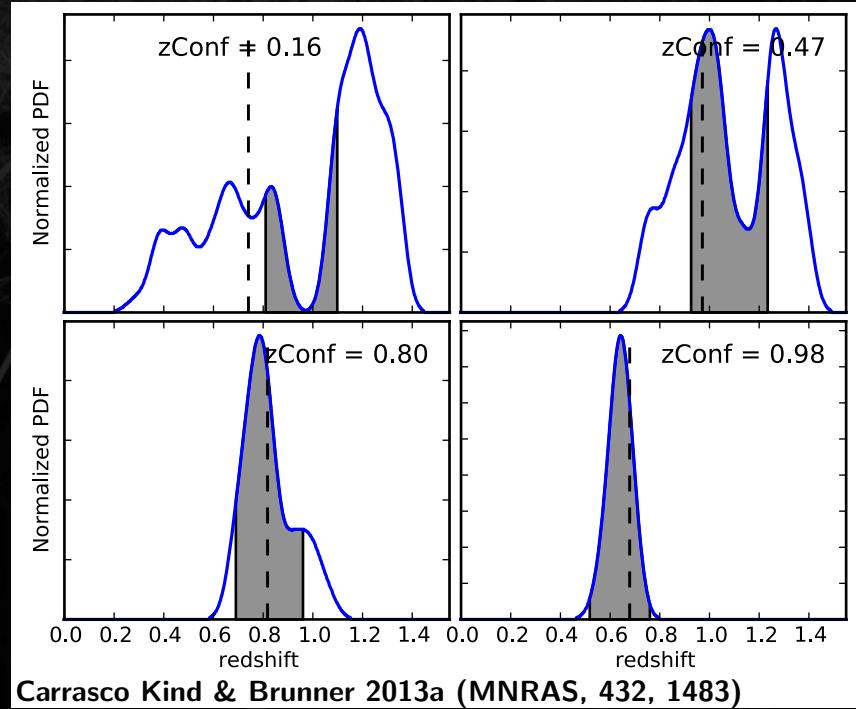
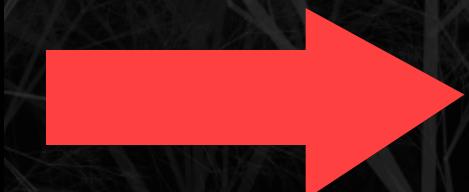
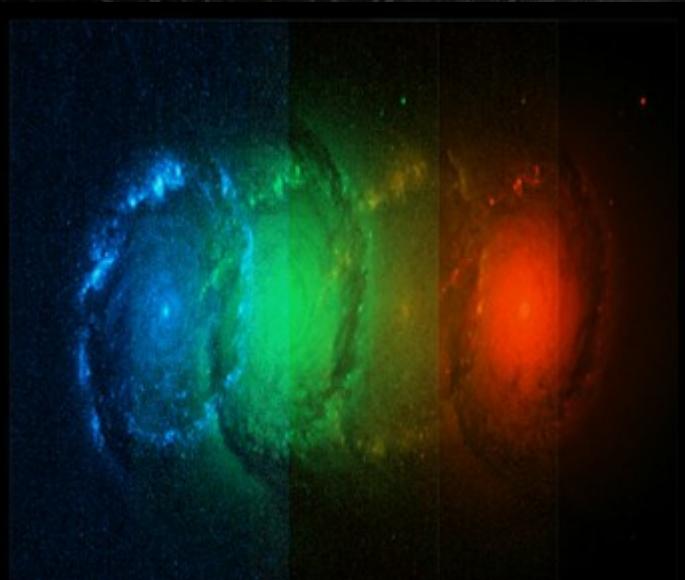
- Photo- $z$  Probability Density Functions needed
- Several methods/codes to compute photo- $z$
- Need for a meta-algorithm that combines multiple techniques
- PDFs are good but for large datasets, storage and I/O will be an issue
- Machine Learning and statistical tools

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

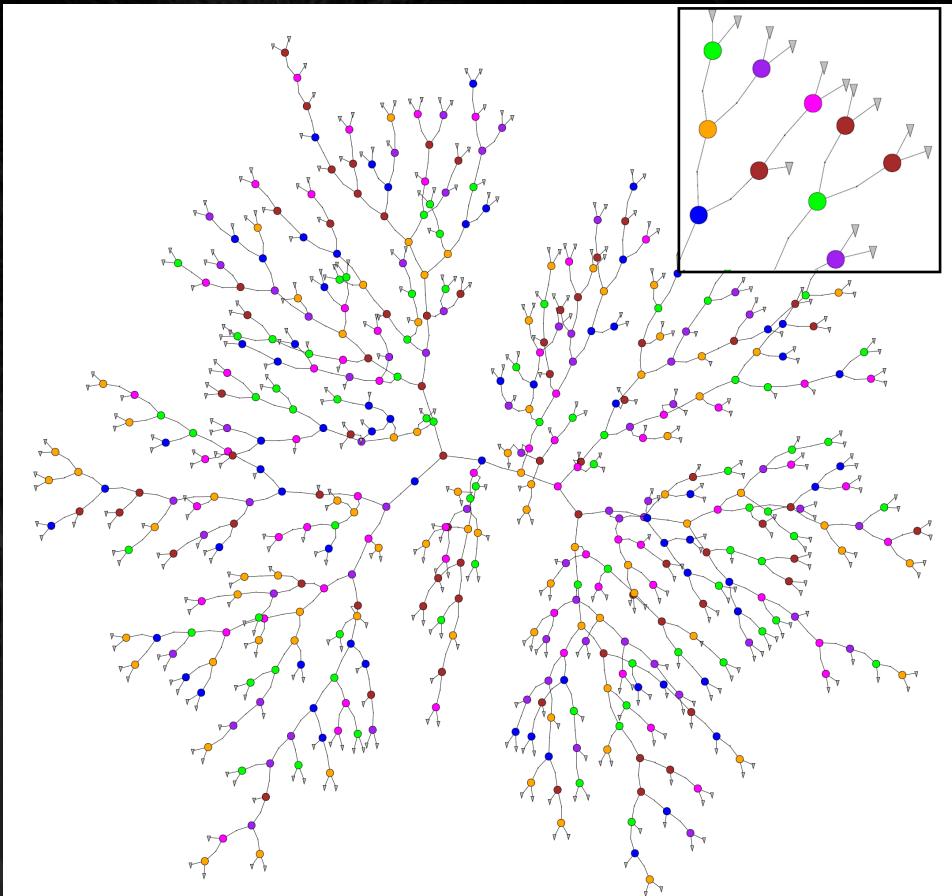
- Photo- $z$  Probability Density Functions needed
- Several methods/codes to compute photo- $z$
- Need for a meta-algorithm that combines multiple techniques
- PDF are good but for large datasets, storage and I/O will be an issue
- Machine Learning and statistical tools

- 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.)

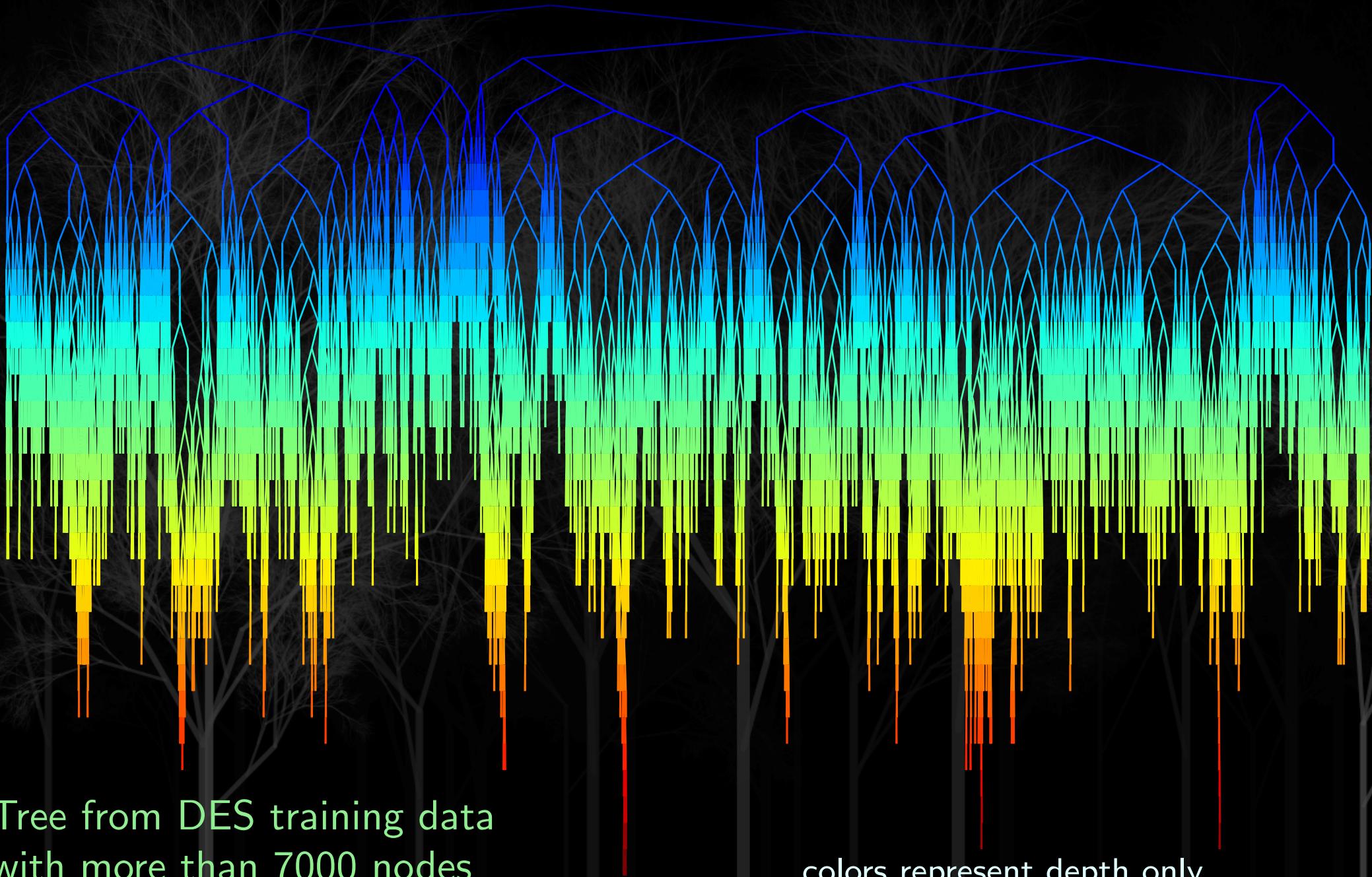


- TPZ (Trees for Photo-Z) is a supervised machine learning code
- Prediction trees and random forest
- Incorporate measurements errors and deals with missing values
- Ancillary information: expected errors, attribute ranking and others
- Application to the S/G

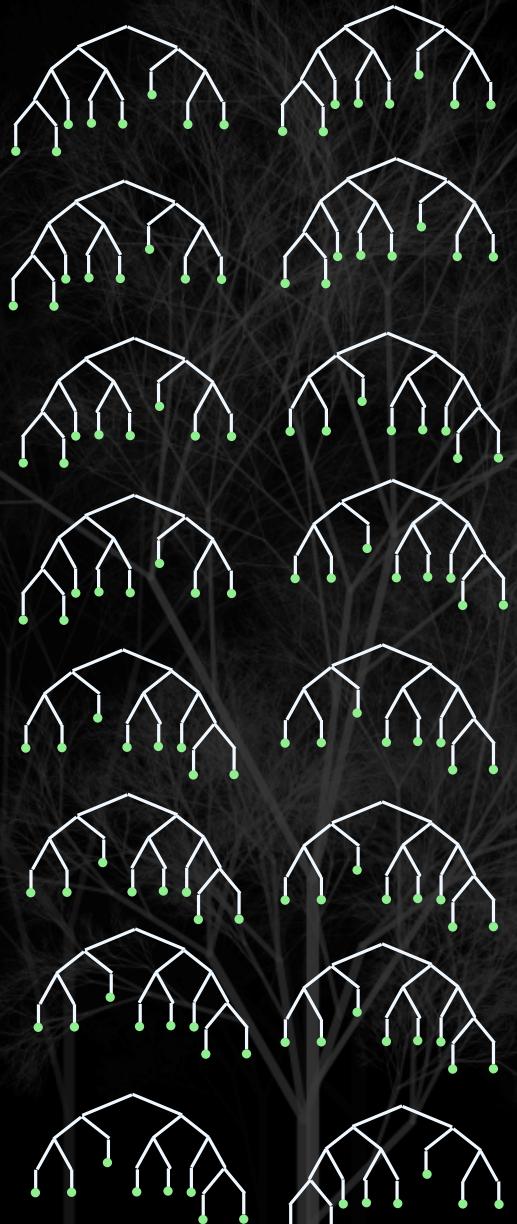


Carrasco Kind & Brunner 2013a (MNRAS, 432, 1483)

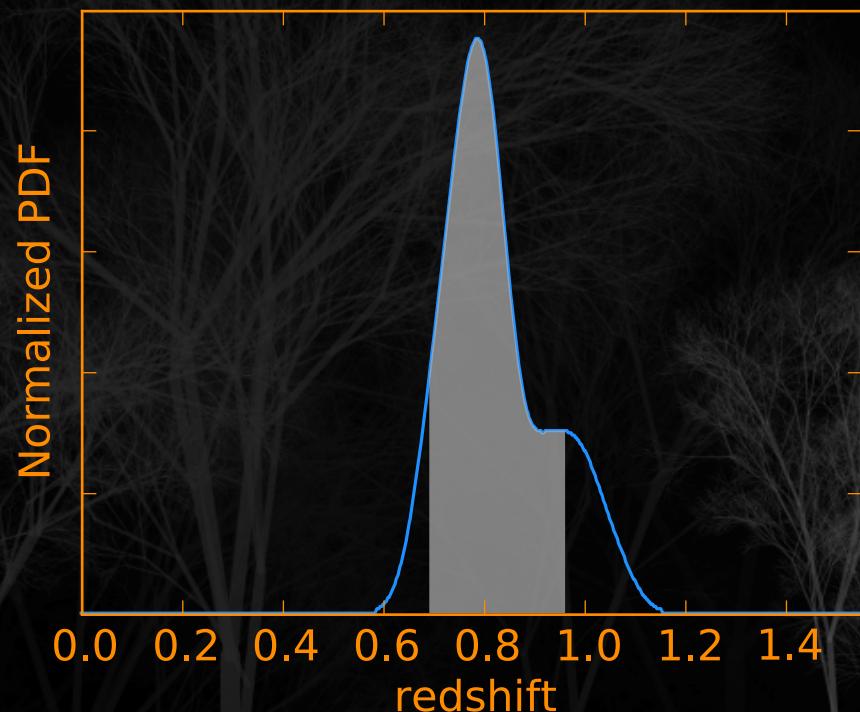
<http://lcdm.astro.illinois.edu/code/mlz.html>



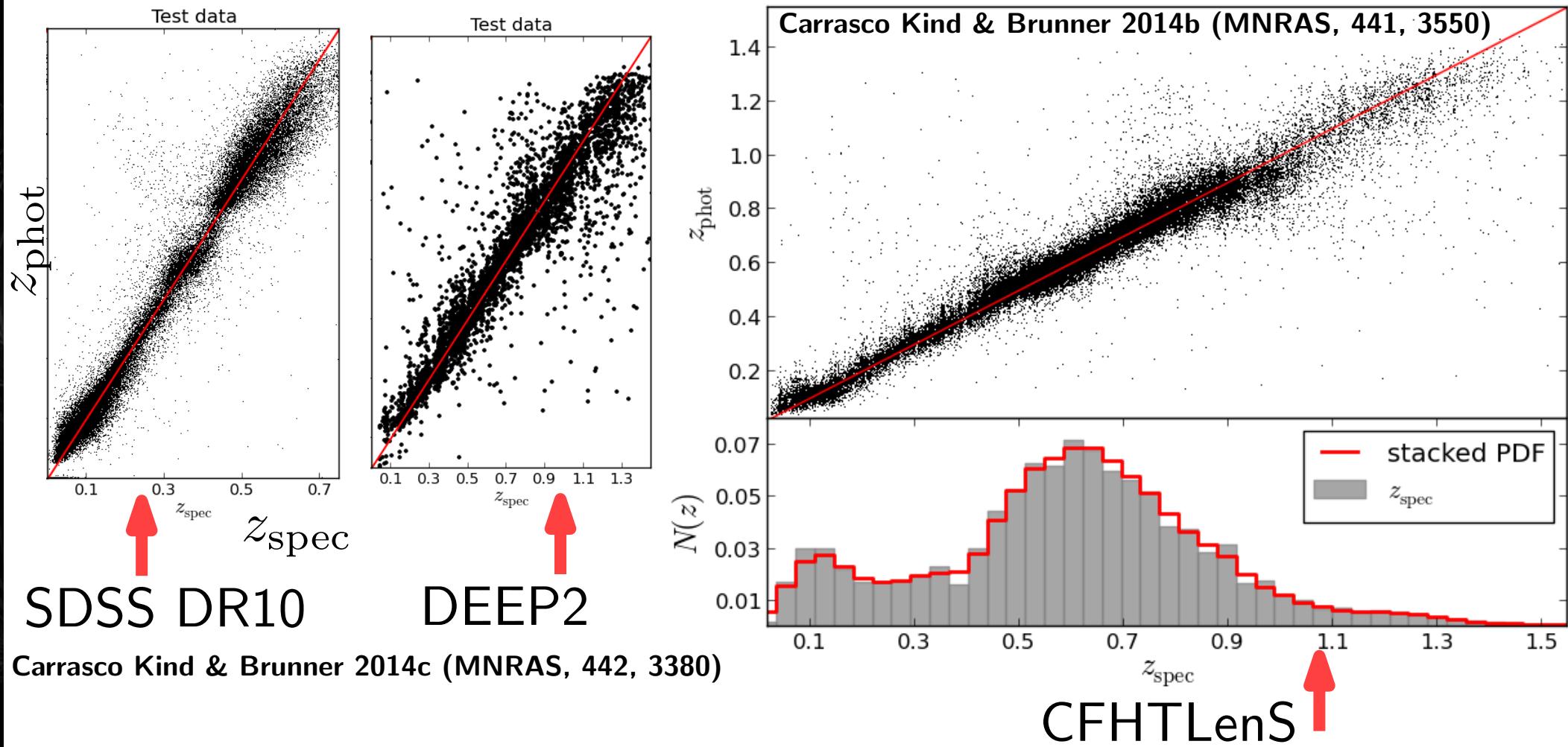
# Photo- $z$ PDF estimation: Random forest



Combine predictions  
from trees

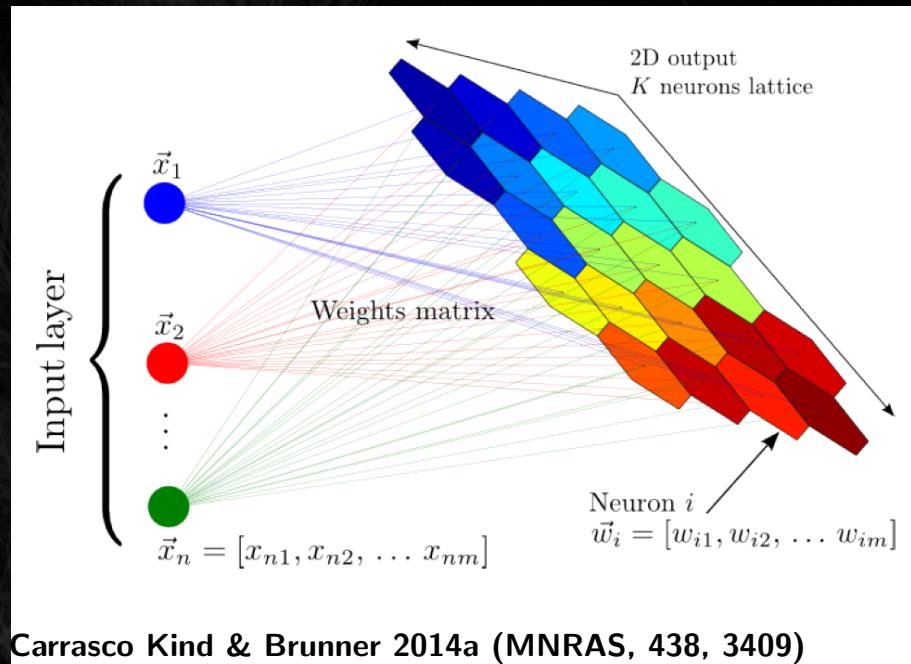


Trees are ideally uncorrelated and strong  
Bootstrapping and error sampling  
Random features at each node

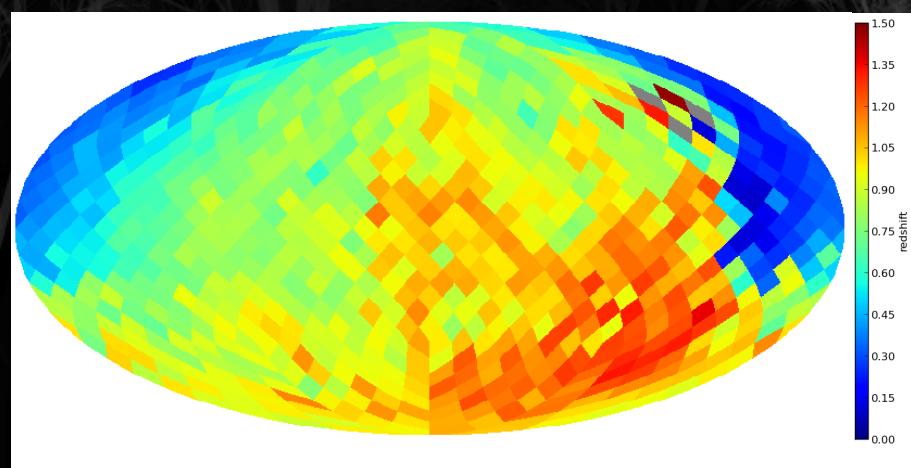


TPZ has been tested in several databases with remarkable results

- 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

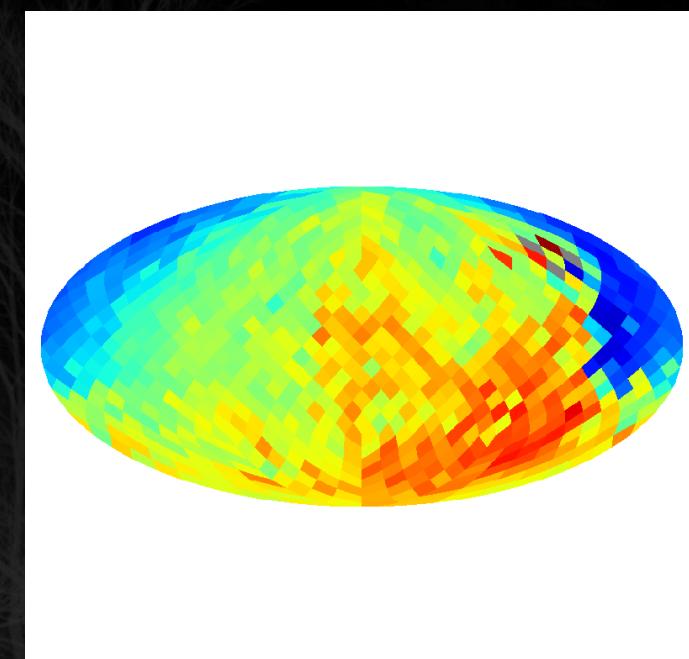
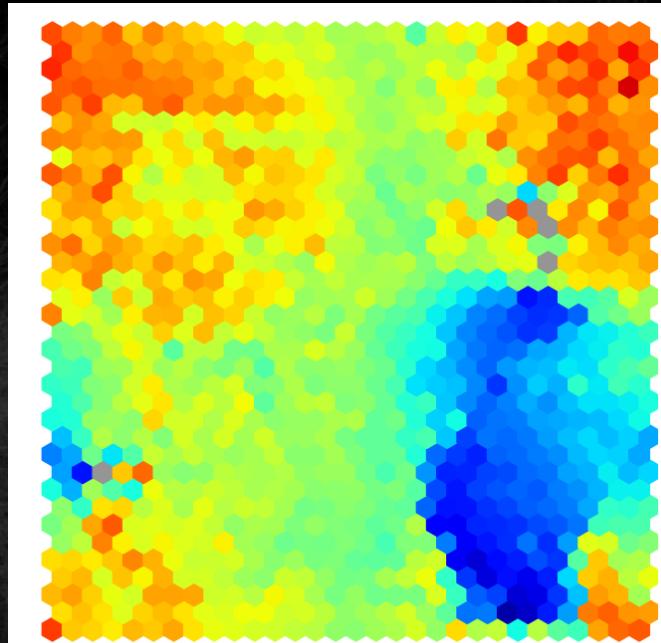
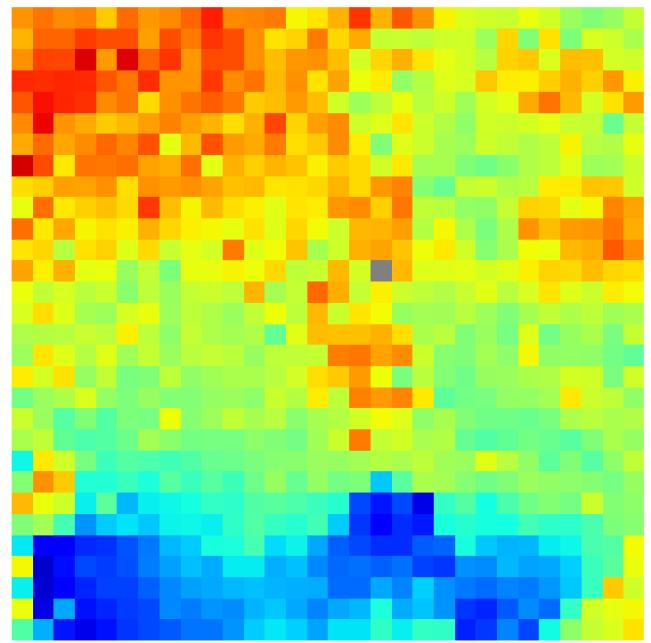


Carrasco Kind & Brunner 2014a (MNRAS, 438, 3409)



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

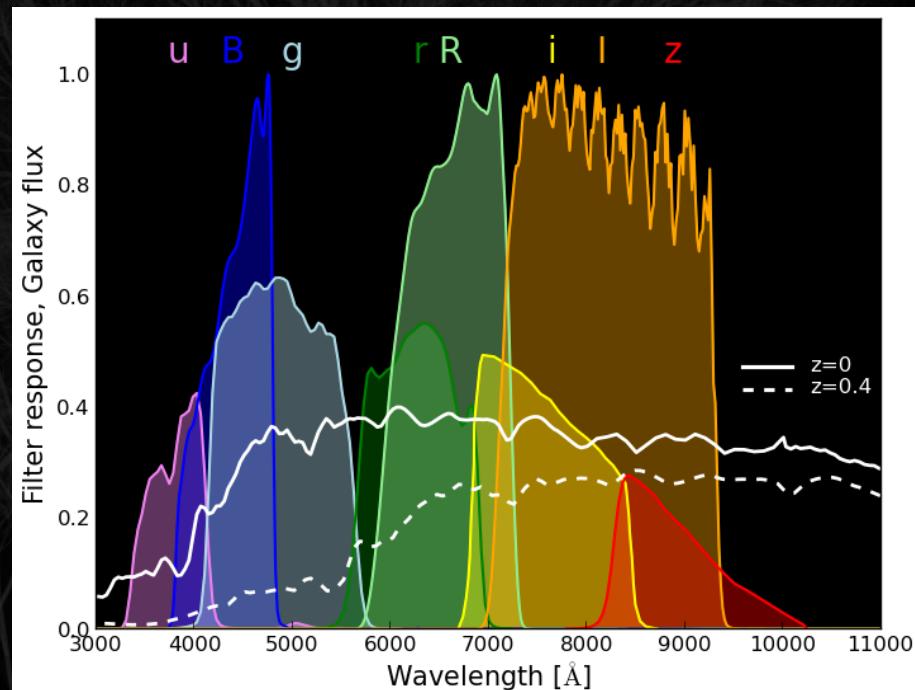
# SOM topologies



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

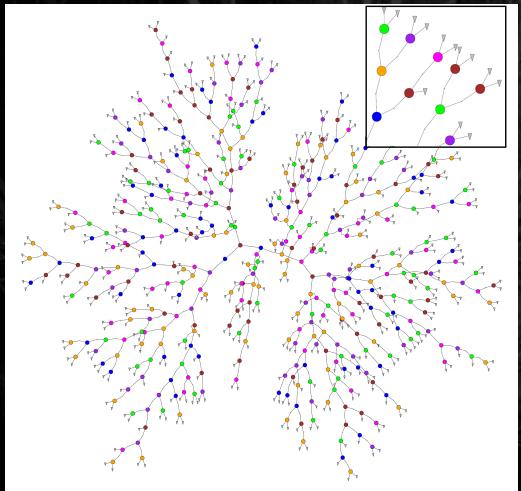
Different topologies can be used with or without periodic boundary conditions

- 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

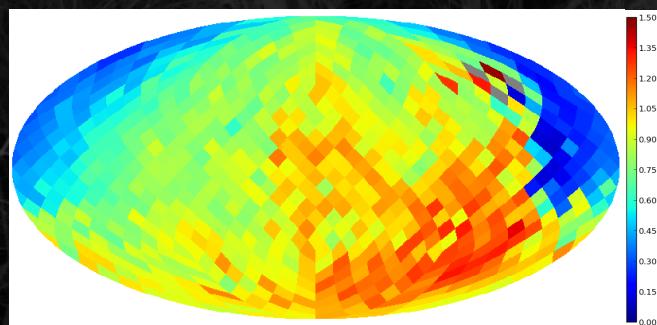


Carrasco Kind & Brunner 2014c (MNRAS, 442, 3380)

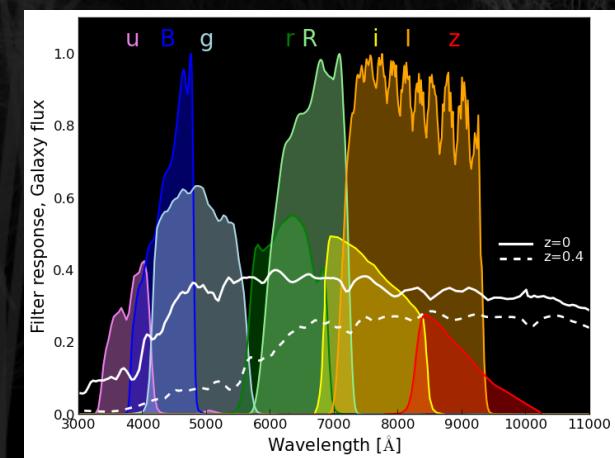
# Photo- $z$ PDF combination



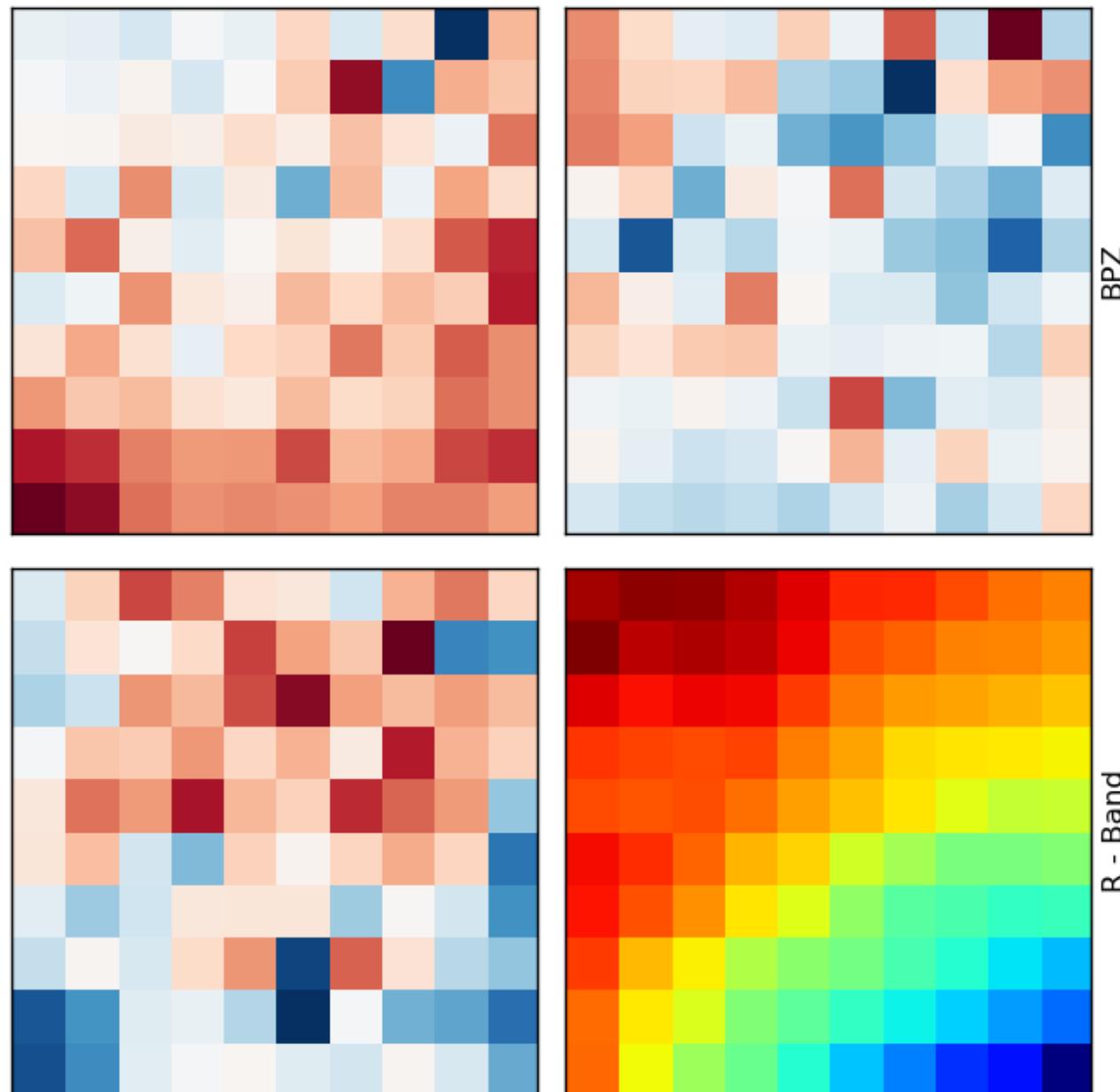
+



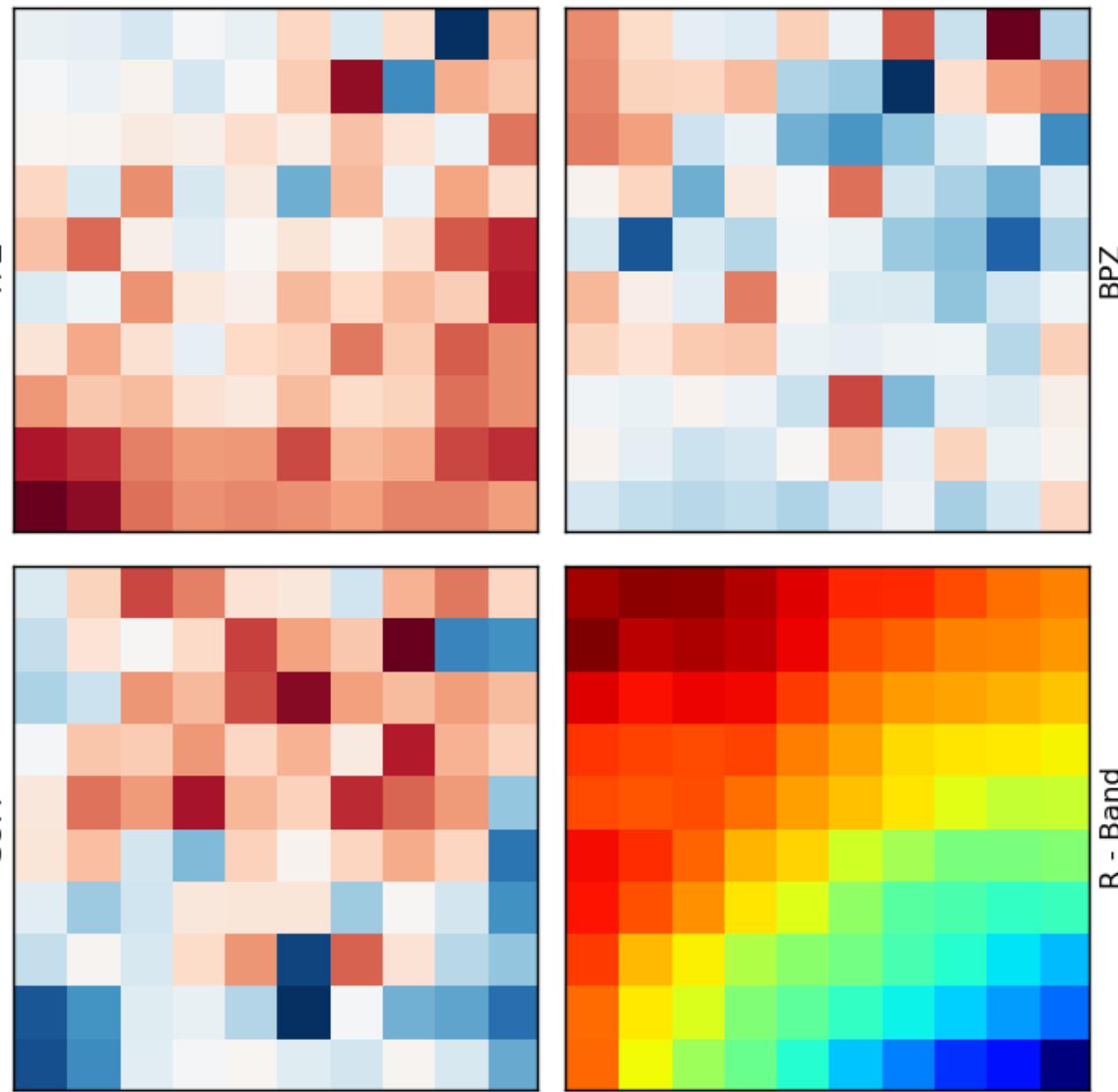
+



# Bayesian framework



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



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

Our approach

Supervised method

+

Unsupervised method

+

Template fitting

+

Weighting scheme

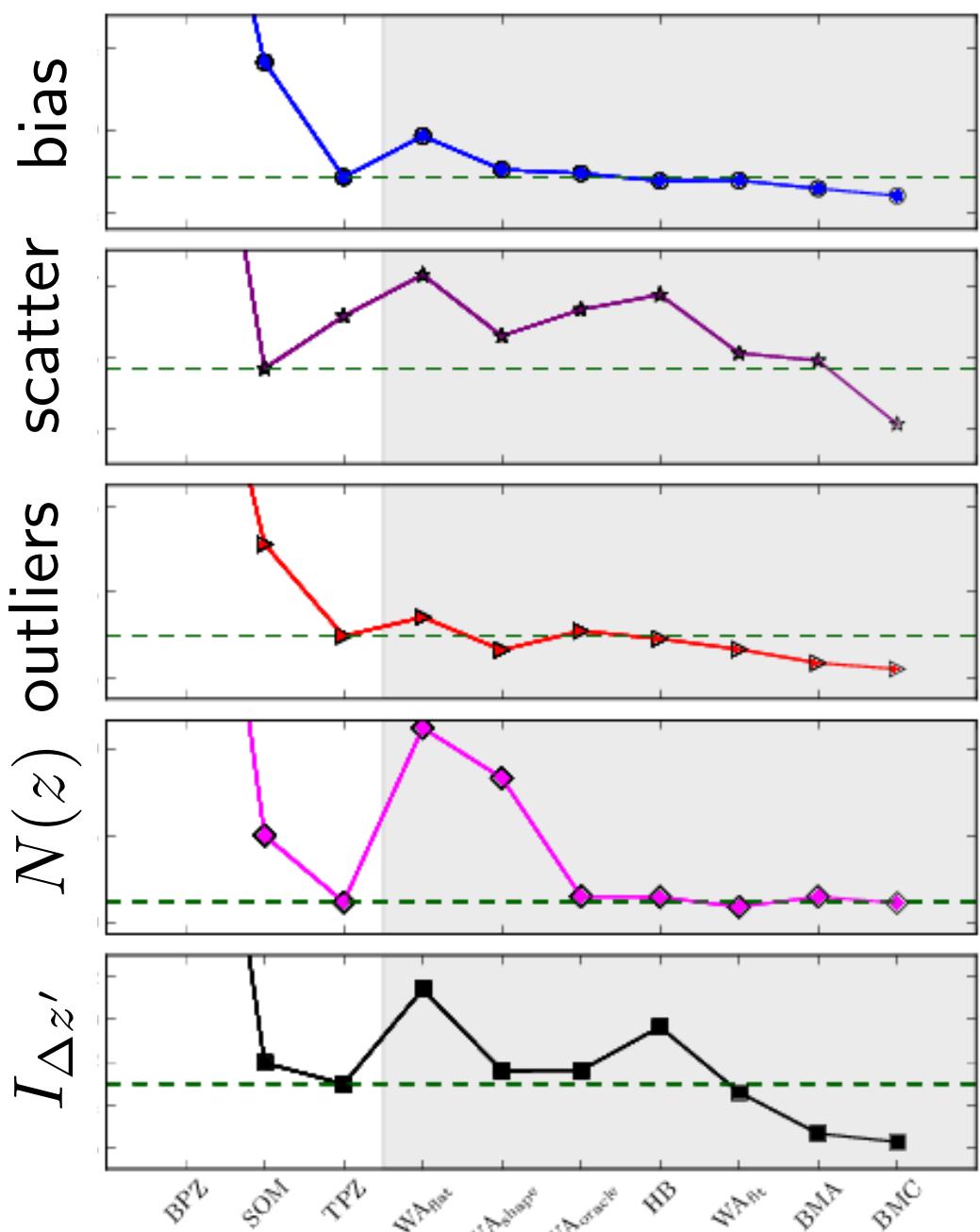
↓

photo- $z$  PDF

+

Outliers

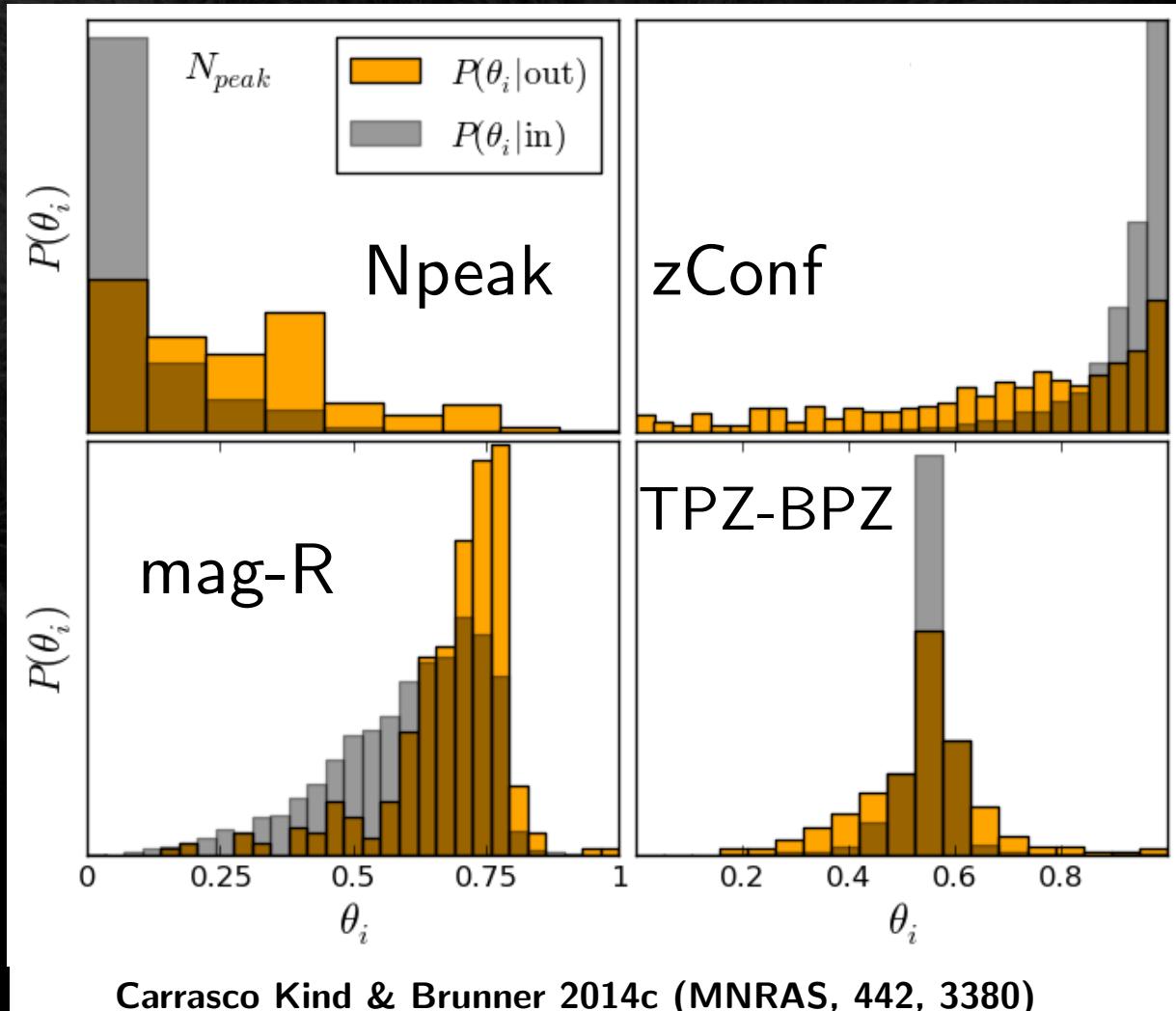
# Photo- $z$ PDF combination: Results



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

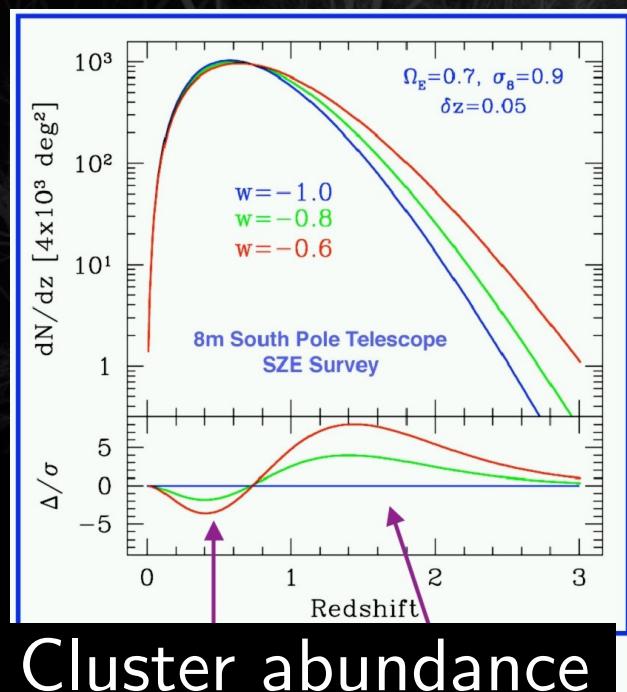
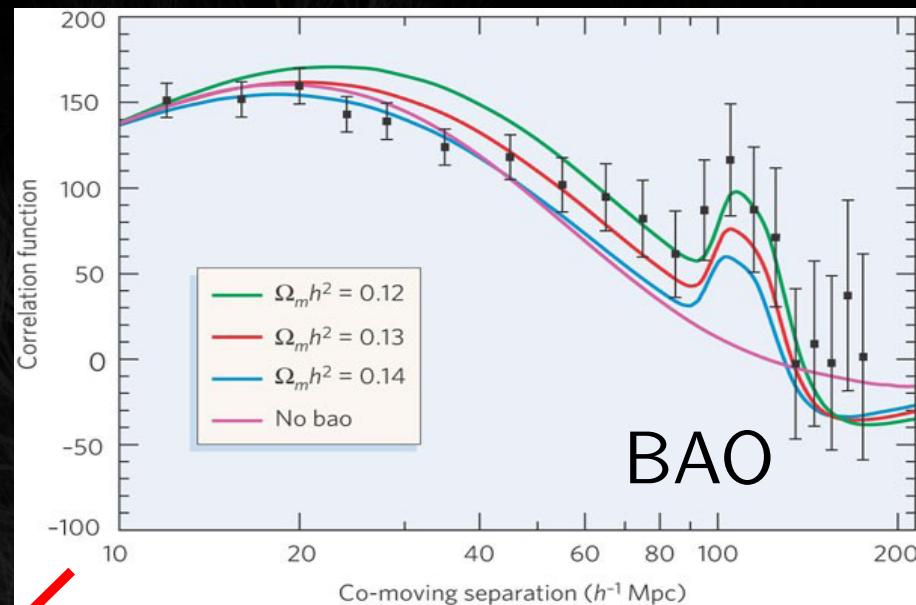
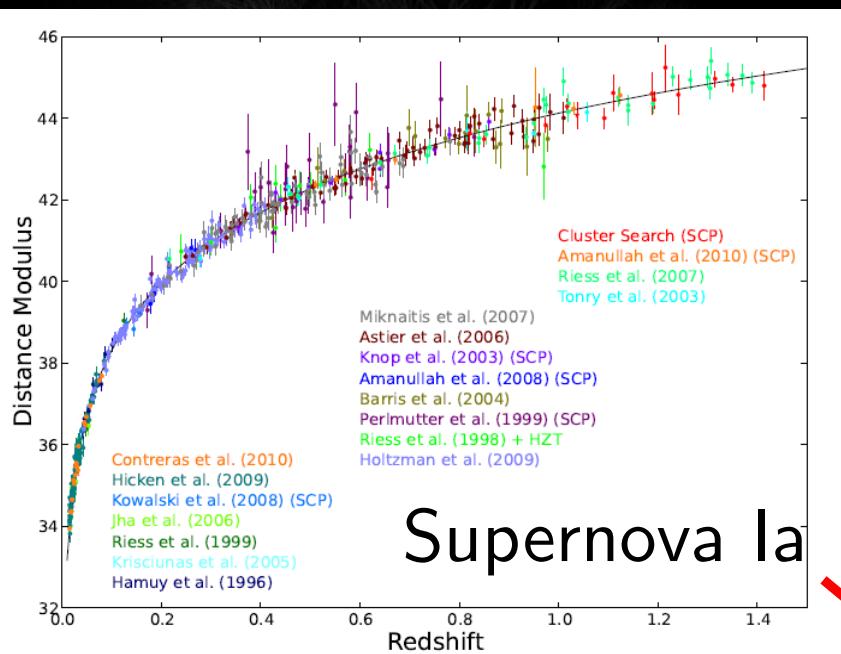
Naïve Bayes Classifier (same used for spam emails) to identify "spam" galaxies using information from multiple techniques

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



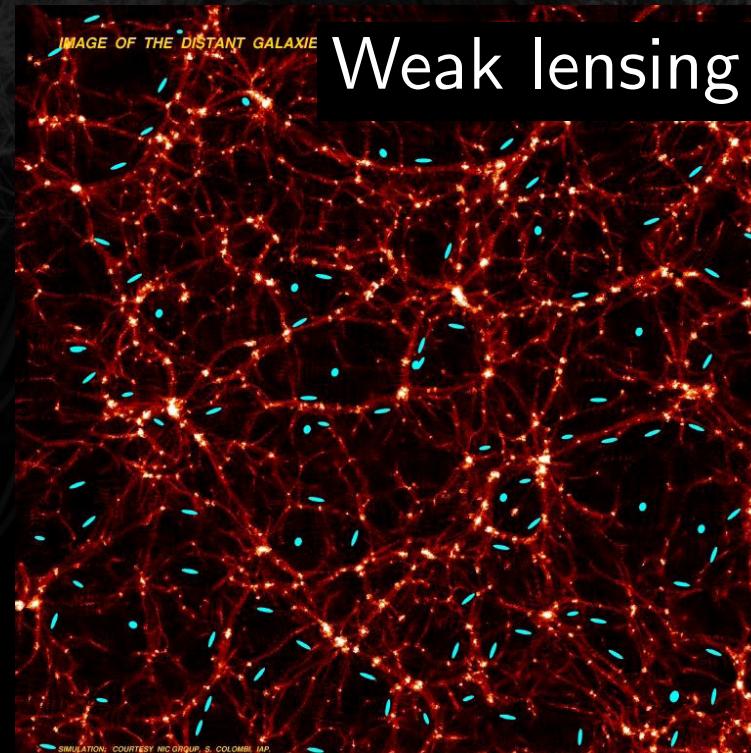
Carrasco Kind & Brunner 2014c (MNRAS, 442, 3380)

# Observational Probes of Dark Energy

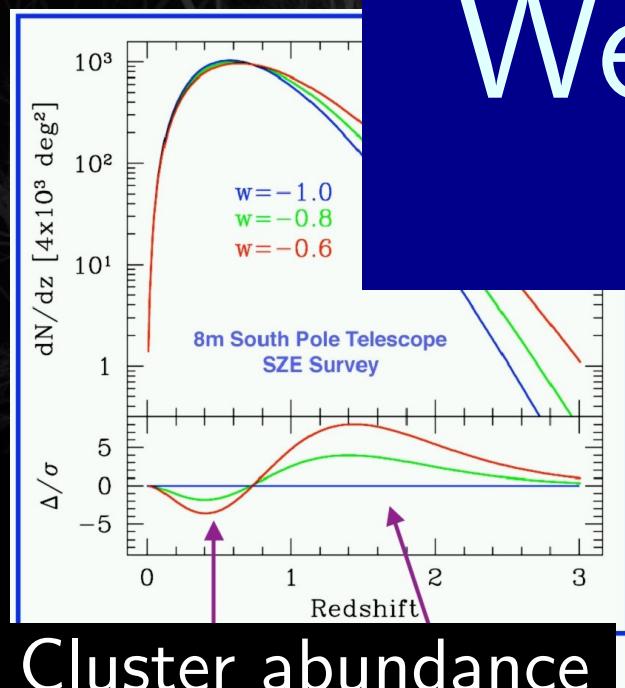
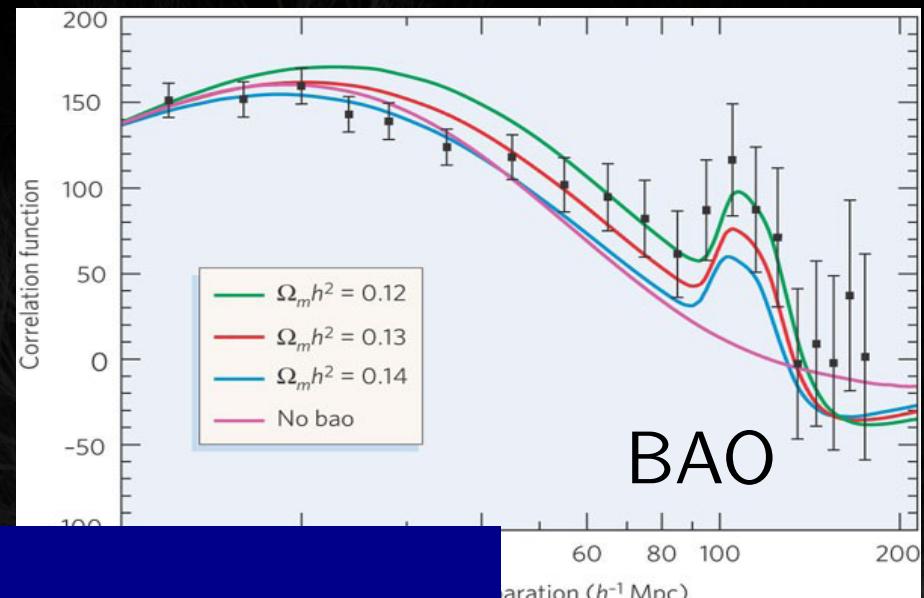
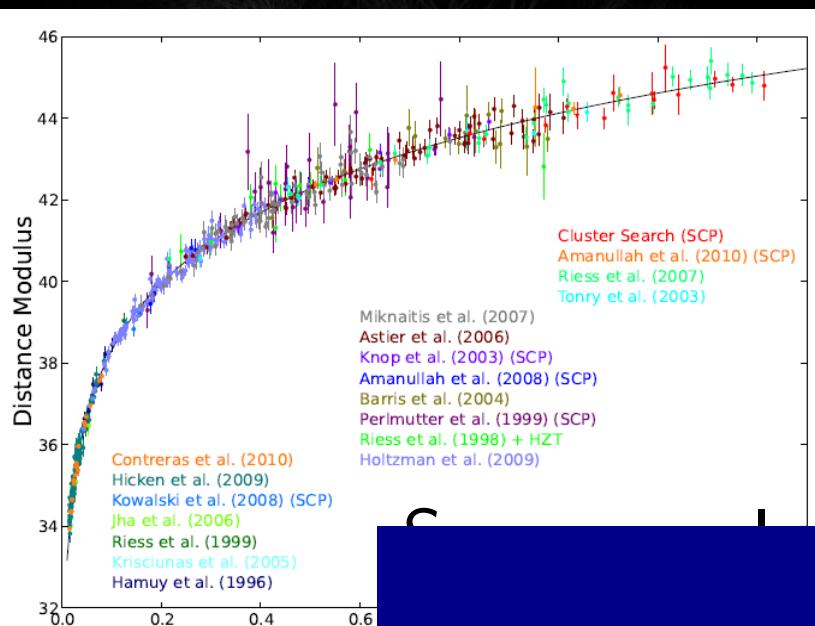


Geometry

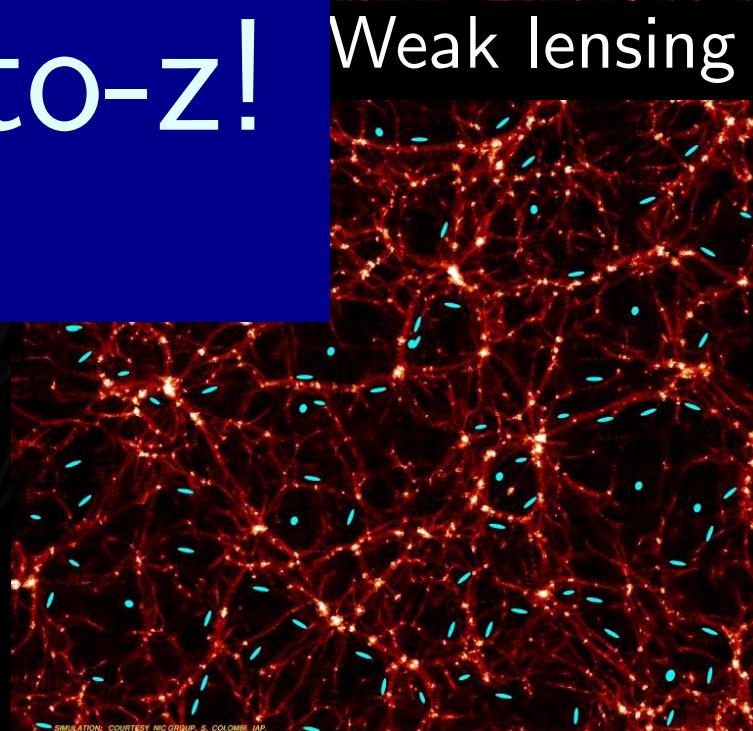
Growth of structure



# Observational Probes of Dark Energy



Growth of structure



# Photo-z for DES SV data



Monthly Notices  
of the  
ROYAL ASTRONOMICAL SOCIETY

MNRAS 445, 1482–1506 (2014)



doi:10.1093/mnras/stu1836

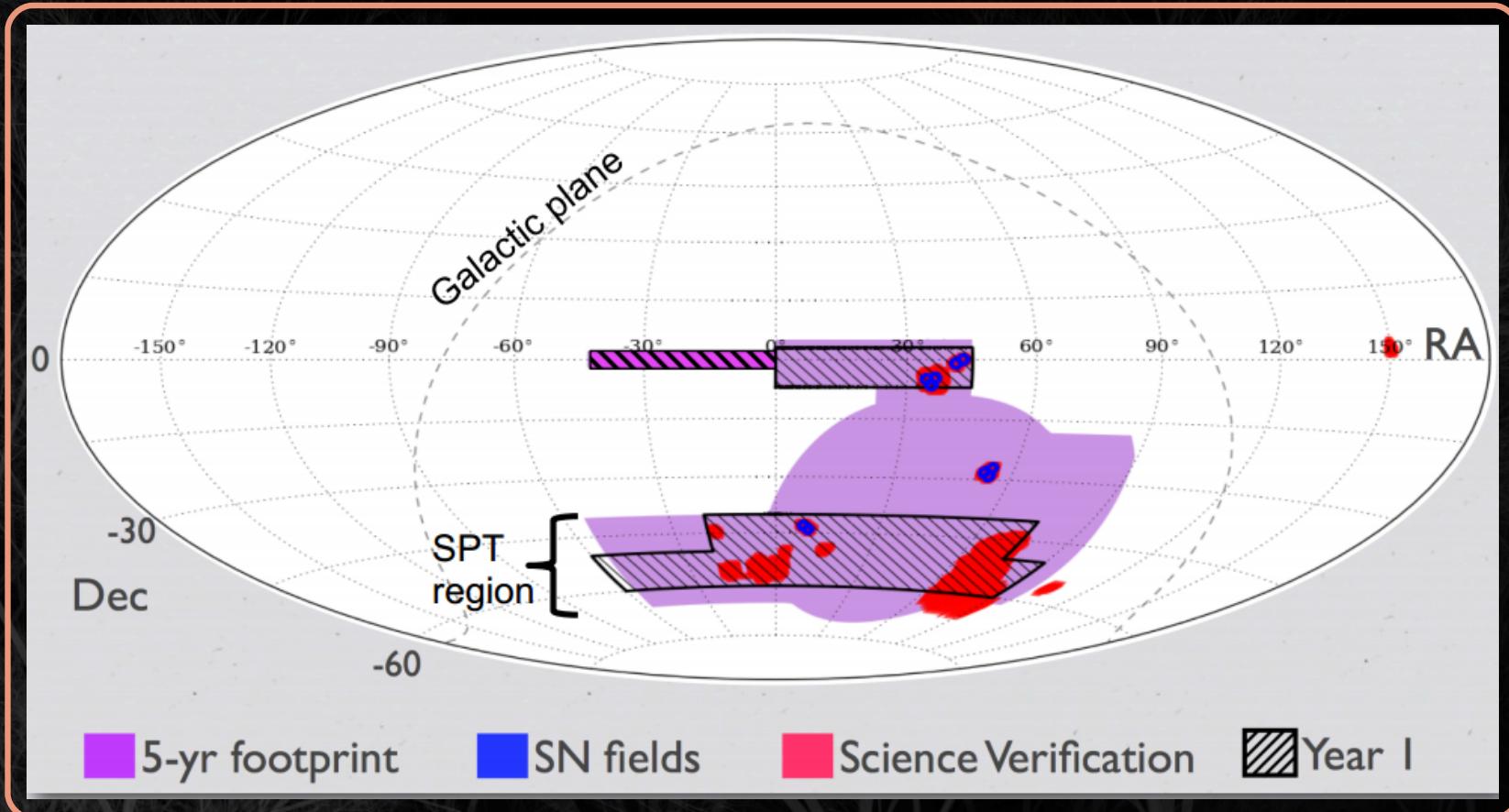
## Photometric redshift analysis in the Dark Energy Survey Science Verification data

C. Sánchez,<sup>1\*</sup> M. Carrasco Kind,<sup>2,3</sup> H. Lin,<sup>4</sup> R. Miquel,<sup>1,5</sup> F. B. Abdalla,<sup>6</sup>  
 A. Amara,<sup>7</sup> M. Banerji,<sup>6</sup> C. Bennett,<sup>1</sup> R. Brunner,<sup>2</sup> D. Capozzi,<sup>8</sup> A. Carnero,<sup>9,10</sup>  
 F. J. Castander,<sup>11</sup> L. A. N. da Costa,<sup>9,10</sup> C. Cunha,<sup>12</sup> A. Fausti,<sup>10</sup> D. Gerdes,<sup>13</sup>  
 N. Greisel,<sup>14,15</sup> J. Gschwend,<sup>9,10</sup> W. Hartley,<sup>7,16</sup> S. Jouvel,<sup>6</sup> O. Lahav,<sup>6</sup> M. Lima,<sup>10,17</sup>  
 M. A. G. Maia,<sup>9,10</sup> P. Martí,<sup>1</sup> R. L. C. Ogando,<sup>9,10</sup> F. Ostrovski,<sup>9,10</sup> P. Pellegrini,<sup>9</sup>  
 M. M. Rau,<sup>14,15</sup> I. Sadeh,<sup>6</sup> S. Seitz,<sup>14,15</sup> I. Sevilla-Noarbe,<sup>18</sup> A. Sypniewski,<sup>13</sup>  
 J. de Vicente,<sup>18</sup> T. Abbot,<sup>19</sup> S. S. Allam,<sup>4,20</sup> D. Atlee,<sup>21</sup> G. Bernstein,<sup>22</sup>  
 J. P. Bernstein,<sup>23</sup> E. Buckley-Geer,<sup>4</sup> D. Burke,<sup>12,24</sup> M. J. Childress,<sup>25,26</sup> T. Davis,<sup>26,27</sup>  
 D. L. DePoy,<sup>28,29</sup> A. Dey,<sup>21,30</sup> S. Desai,<sup>31,32</sup> H. T. Diehl,<sup>4</sup> P. Doel,<sup>6</sup> J. Estrada,<sup>4</sup>  
 A. Evrard,<sup>13,33,34</sup> E. Fernández,<sup>1</sup> D. Finley,<sup>4</sup> B. Flaugher,<sup>4</sup> J. Frieman,<sup>4</sup>  
 E. Gaztanaga,<sup>11</sup> K. Glazebrook,<sup>35</sup> K. Honscheid,<sup>36</sup> A. Kim,<sup>37</sup> K. Kuehn,<sup>38</sup>  
 N. Kuropatkin,<sup>4</sup> C. Lidman,<sup>38</sup> M. Makler,<sup>39</sup> J. L. Marshall,<sup>28,29</sup> R. C. Nichol,<sup>8</sup>  
 A. Roodman,<sup>12,24</sup> E. Sánchez,<sup>18</sup> B. X. Santiago,<sup>10,40</sup> M. Sako,<sup>22</sup> R. Scalzo,<sup>25</sup>  
 R. C. Smith,<sup>19</sup> M. E. C. Swanson,<sup>3</sup> G. Tarle,<sup>13</sup> D. Thomas,<sup>8,41</sup> D. L. Tucker,<sup>4</sup>  
 S. A. Uddin,<sup>26,35</sup> F. Valdés,<sup>21</sup> A. Walker,<sup>19</sup> F. Yuan<sup>25,26</sup> and J. Zuntz<sup>42</sup>

*Affiliations are listed at the end of the paper*

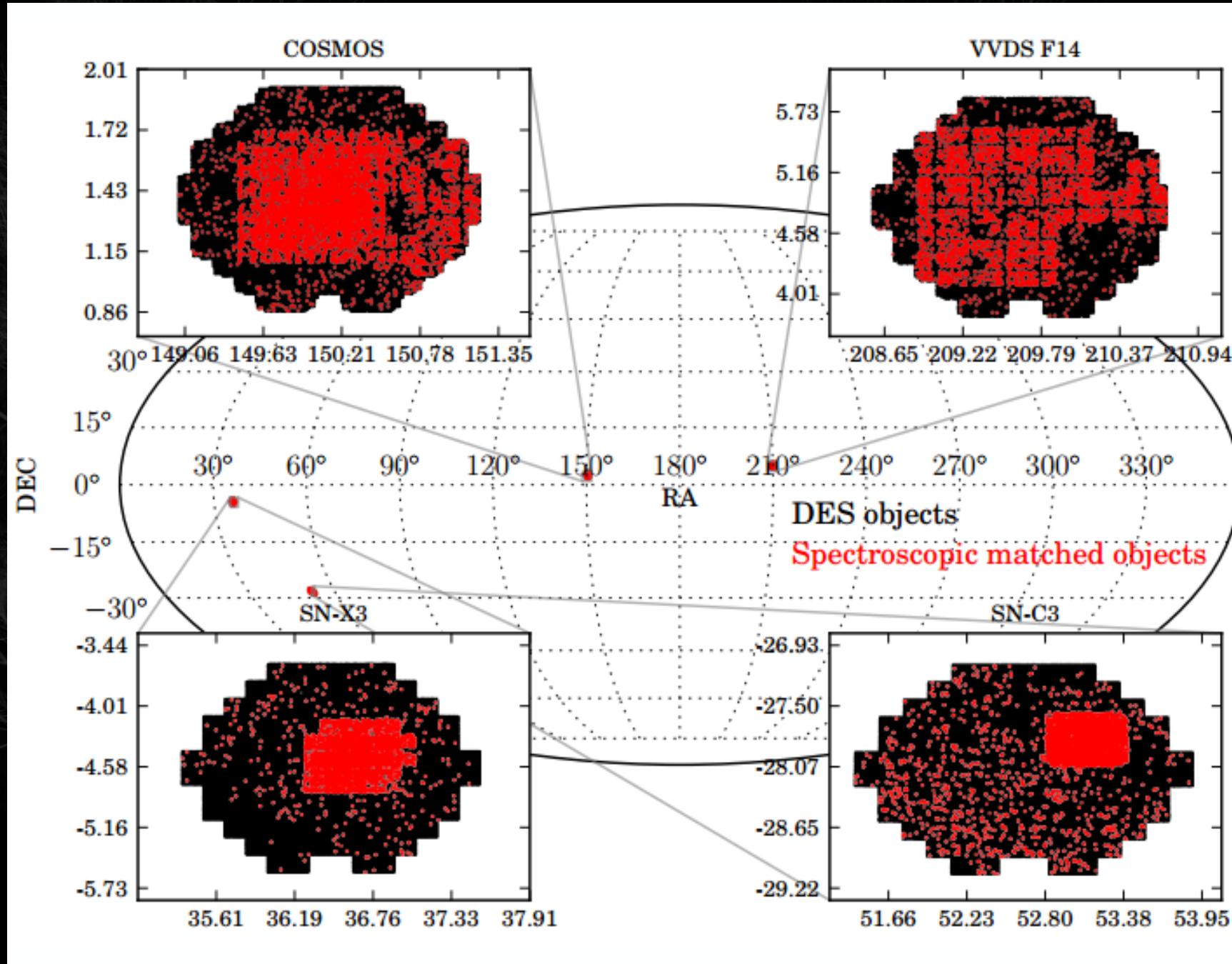
Sanchez et al. (2015). First published paper using DES data!

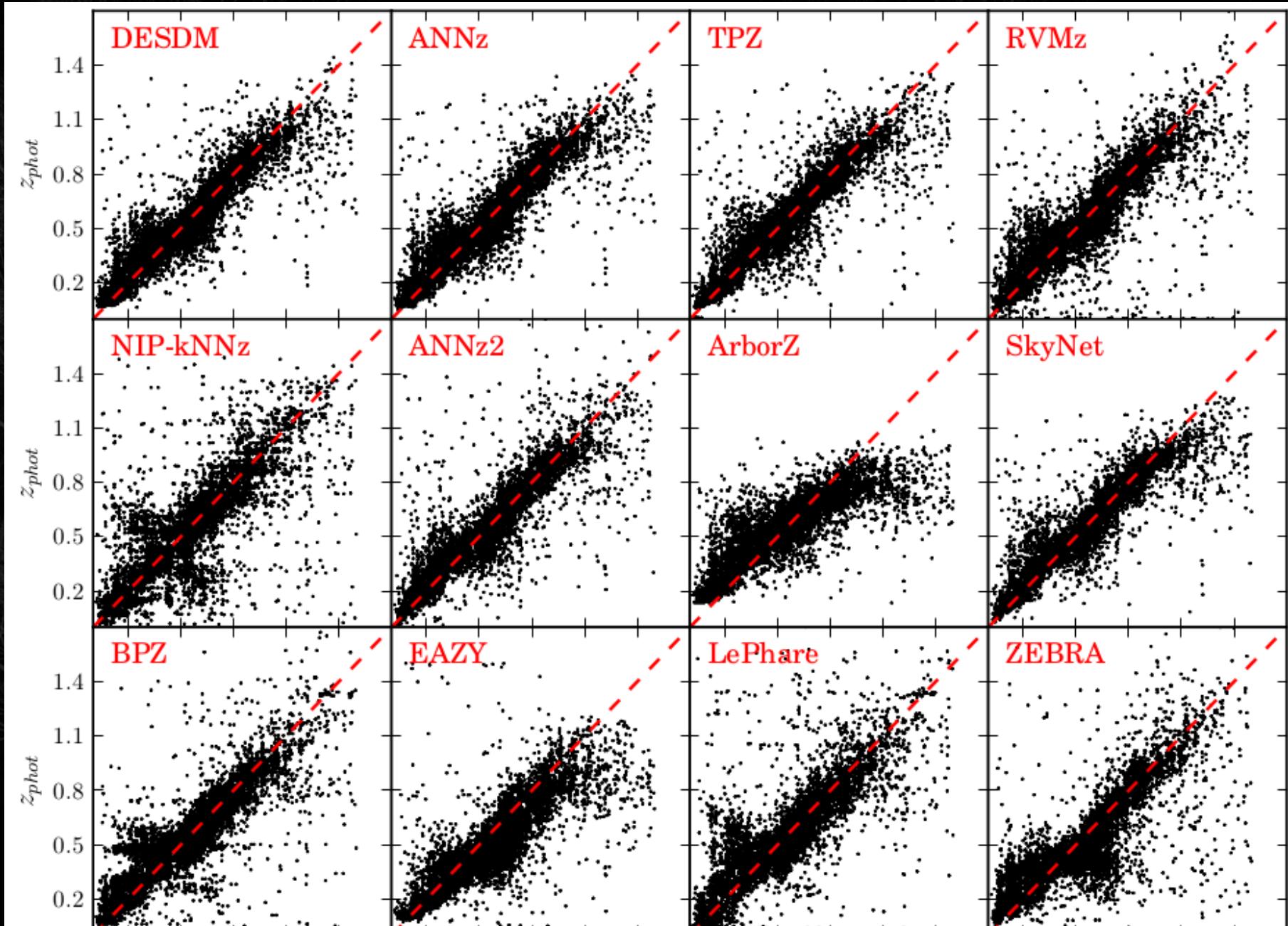
# Photo-z for DES SV data



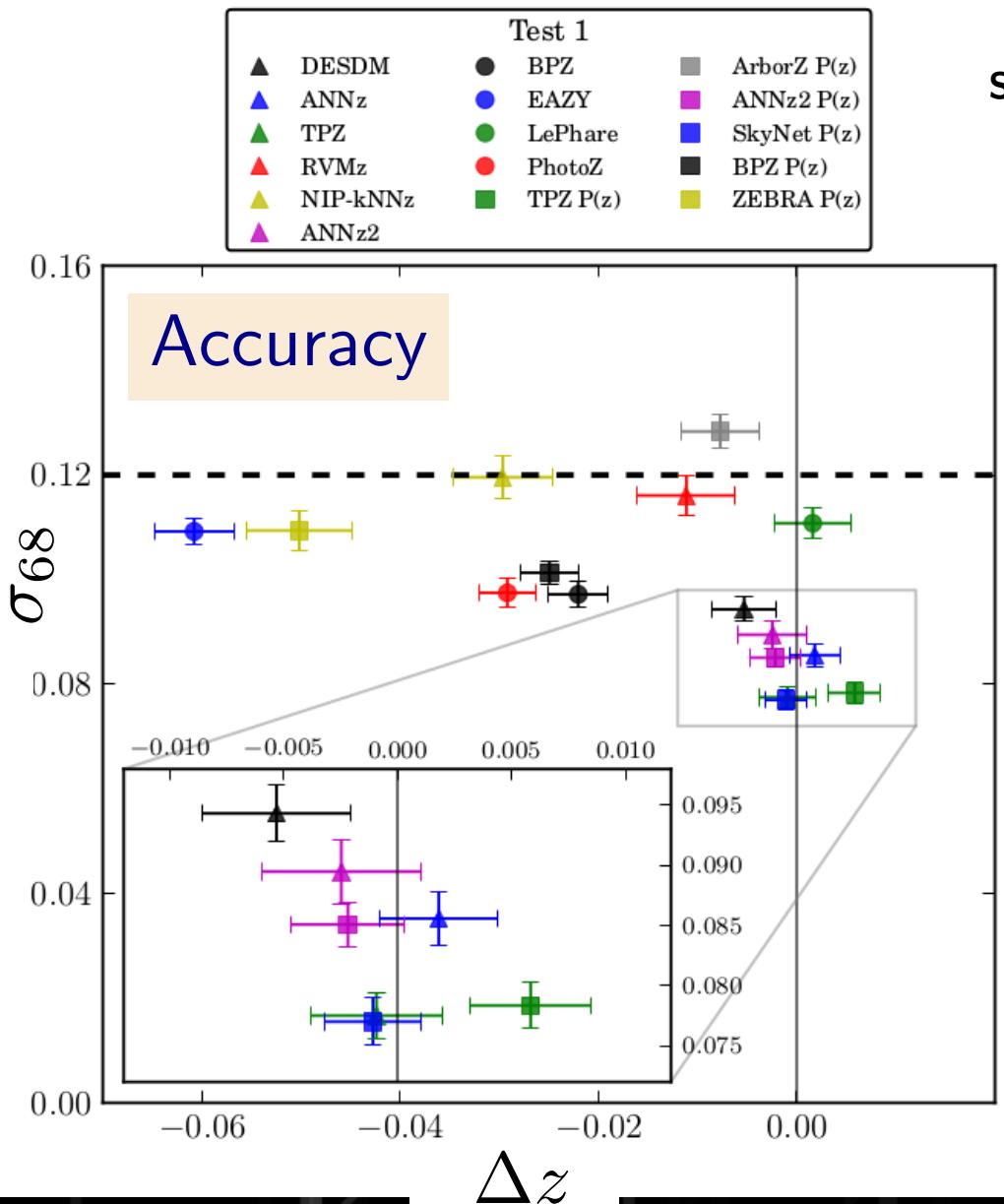
- DES Science Verification data
- Photo-z code comparison and analysis
- Good benchmark results for future releases

# Spectroscopic data

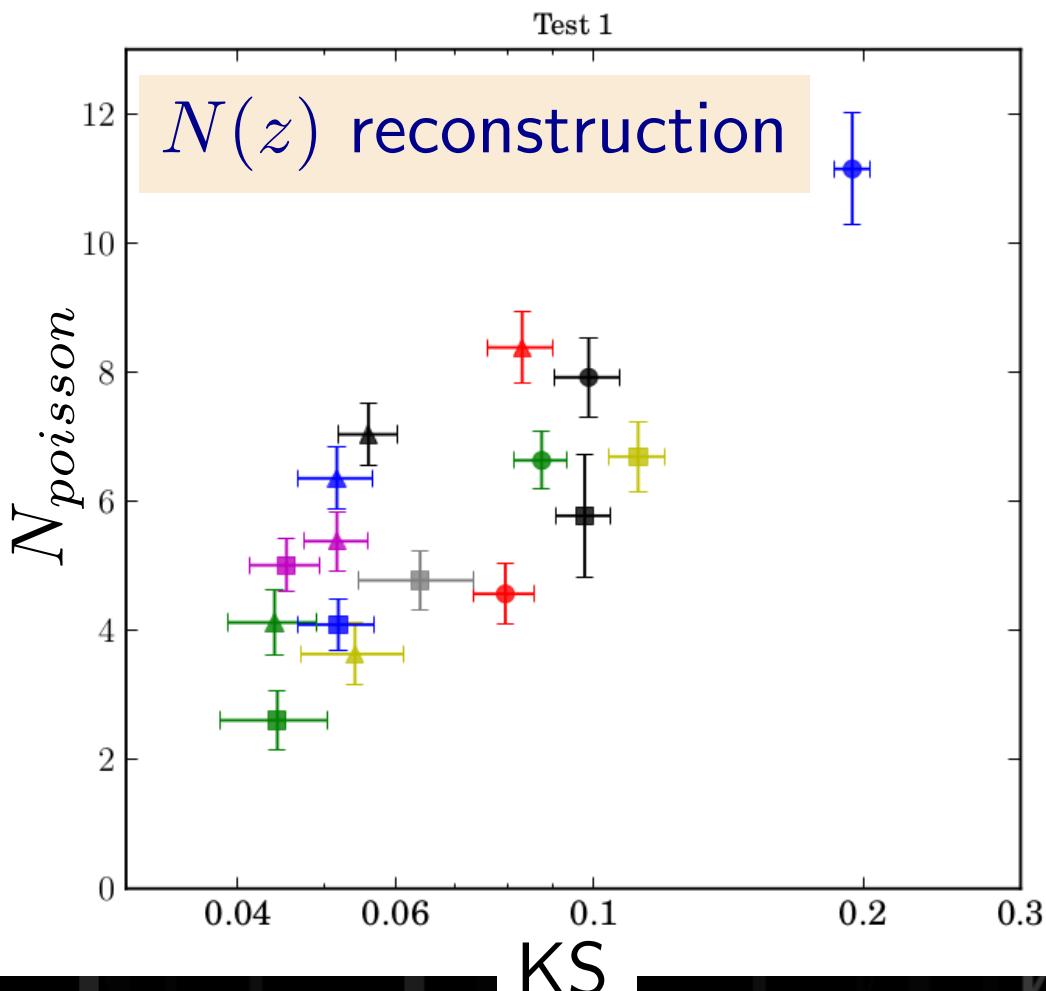




Sánchez, Carrasco Kind, et al. 2014 (MNRAS, 445, 1482)

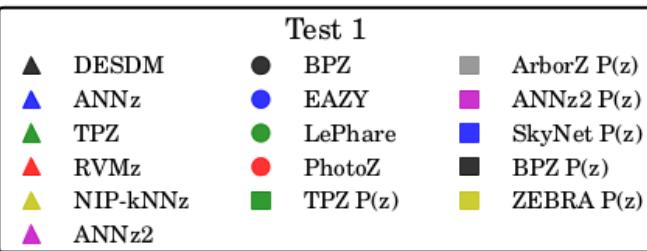


Sánchez, Carrasco Kind, et al. 2014 (MNRAS, 445, 1482)

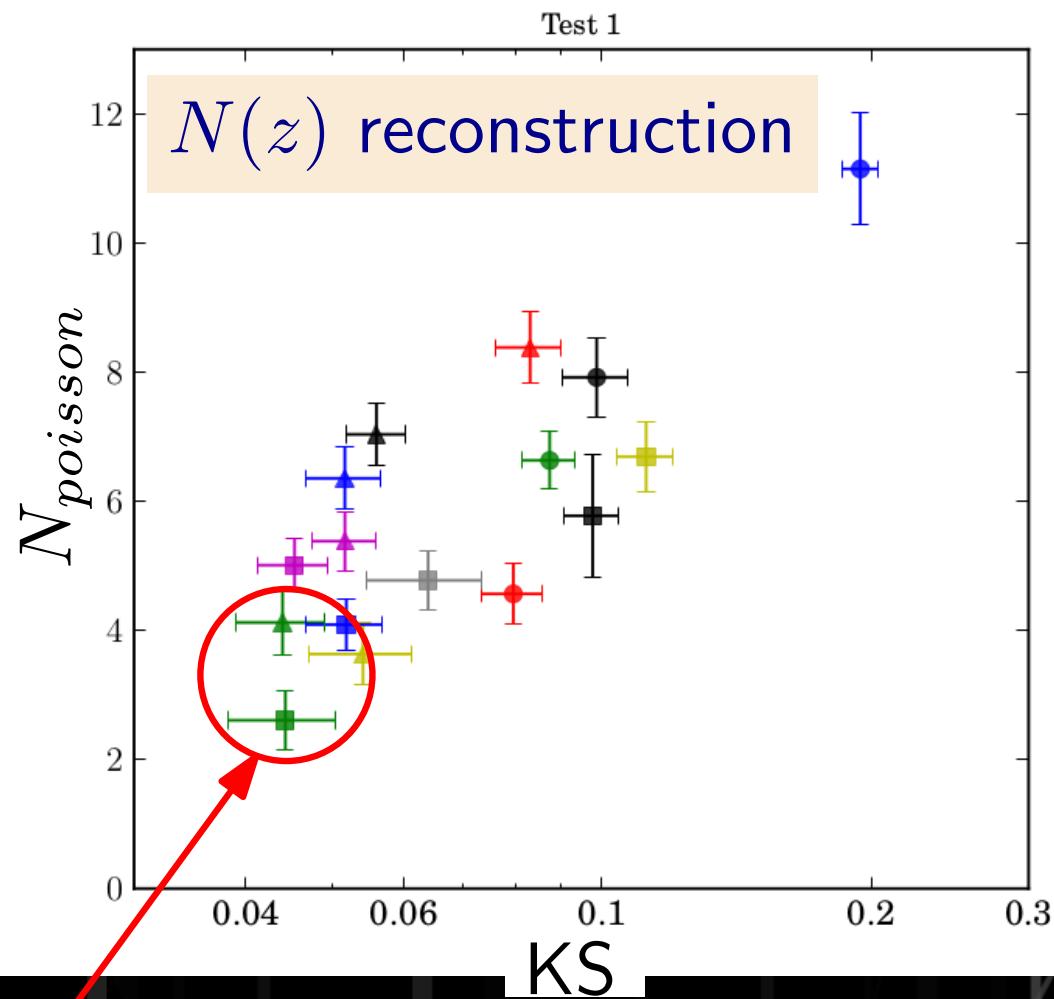
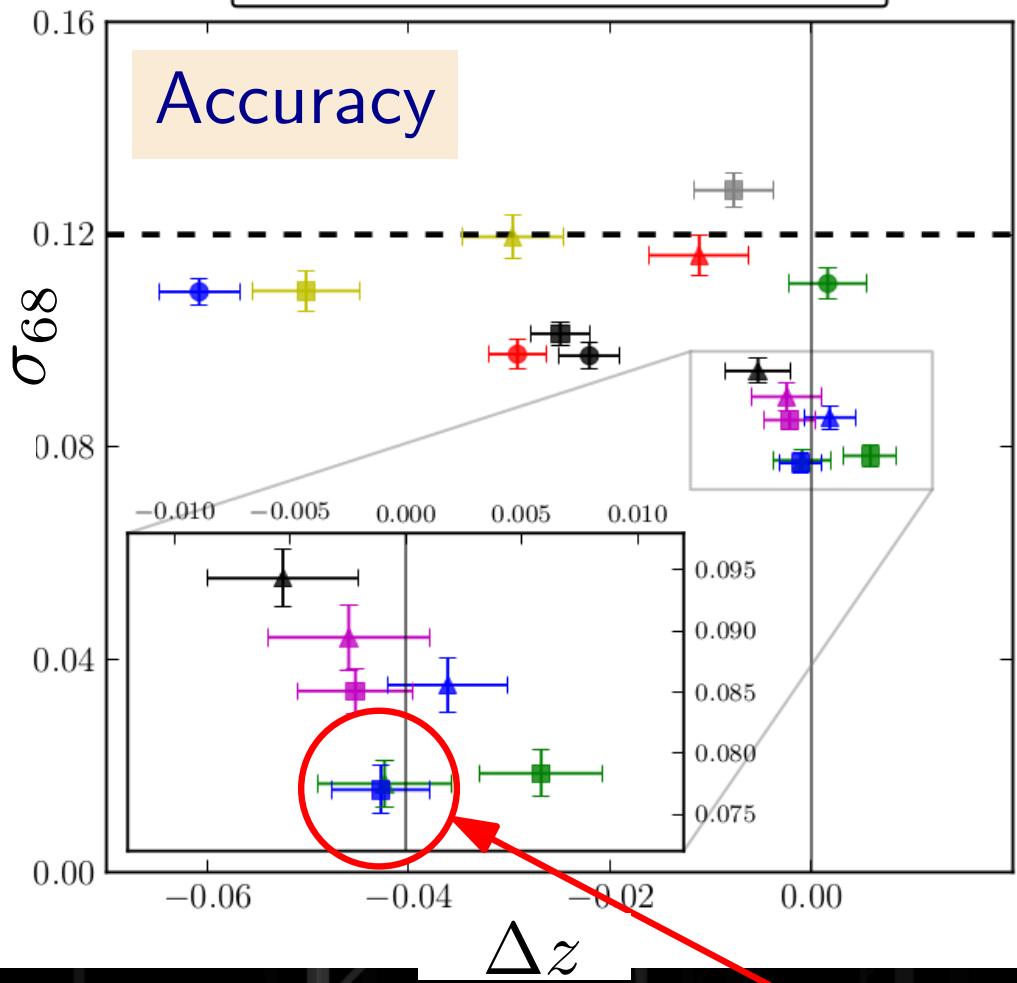


# 13 photo-z codes comparison

# Photo-z for DES SV data



Sánchez, Carrasco Kind, et al. 2014 (MNRAS, 445, 1482)



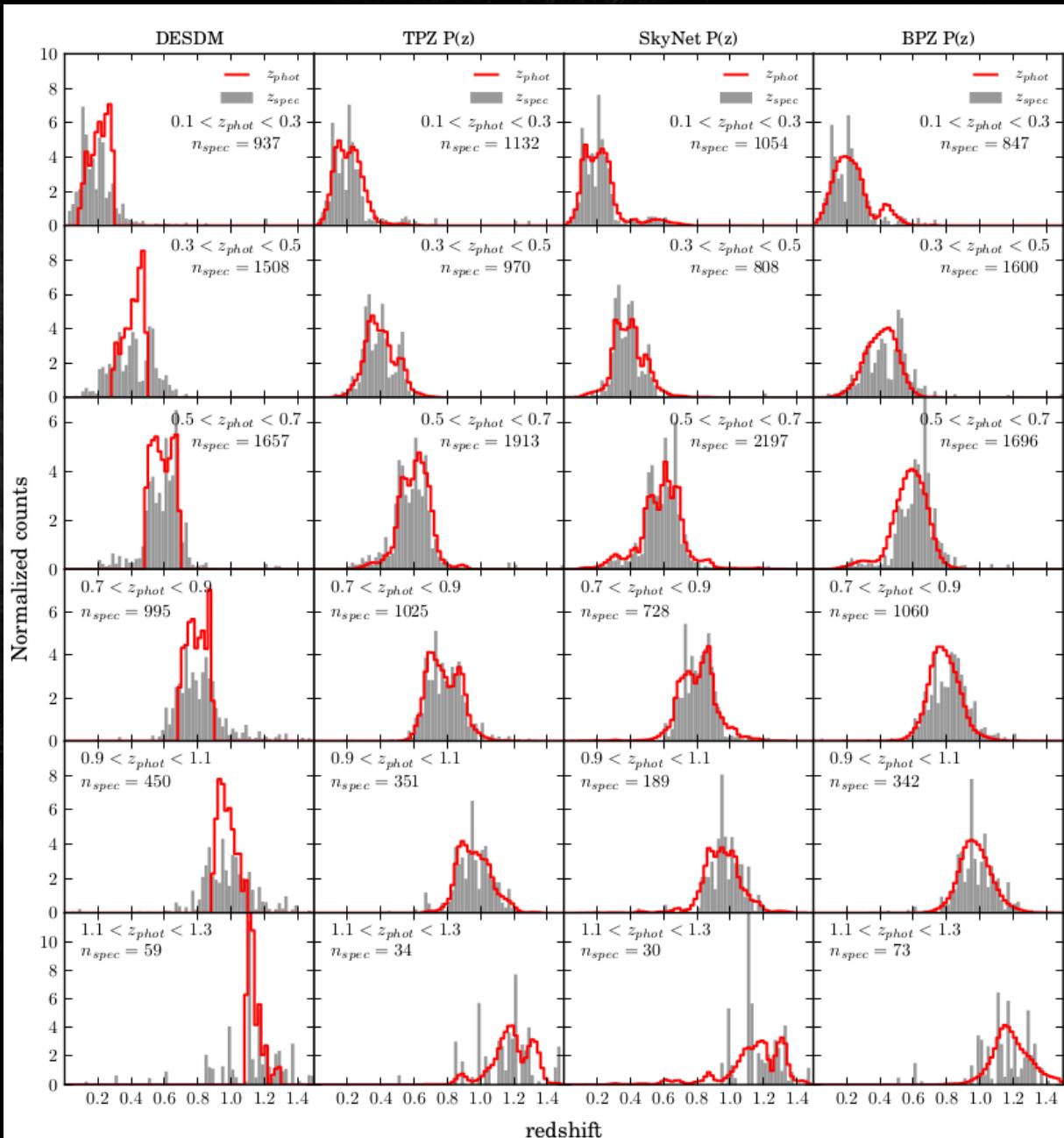
# Photo-z for DES SV data

4 codes  
recommendation

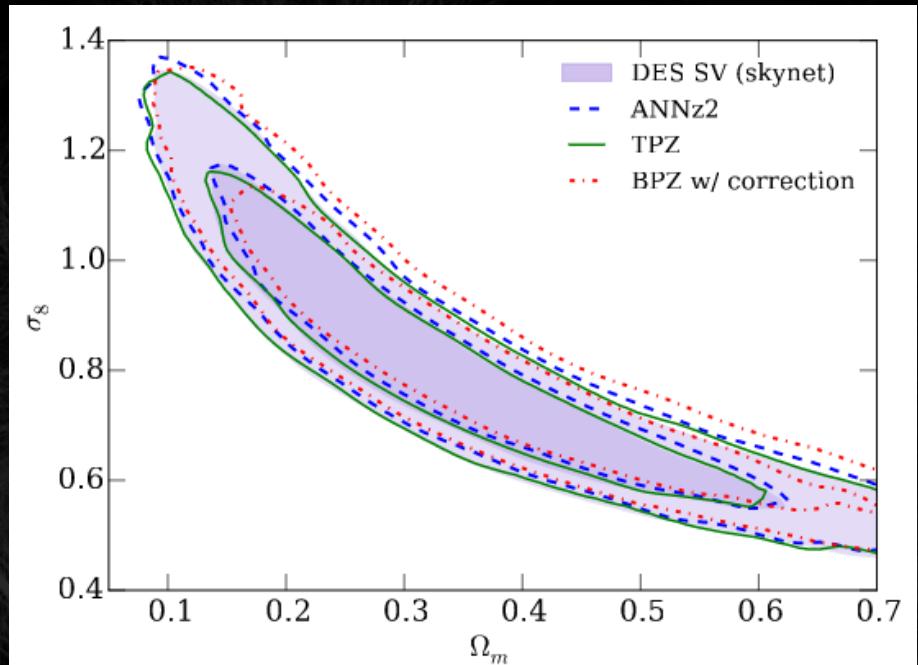
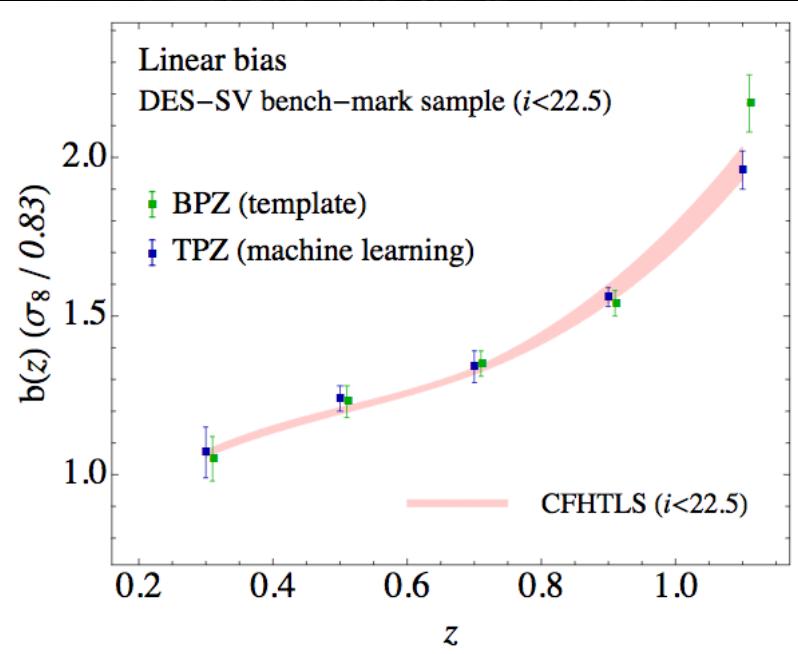
Default, 2 training and  
1 template

PDF methods are  
better for  $N(z)$

Combination methods!  
(not used here)



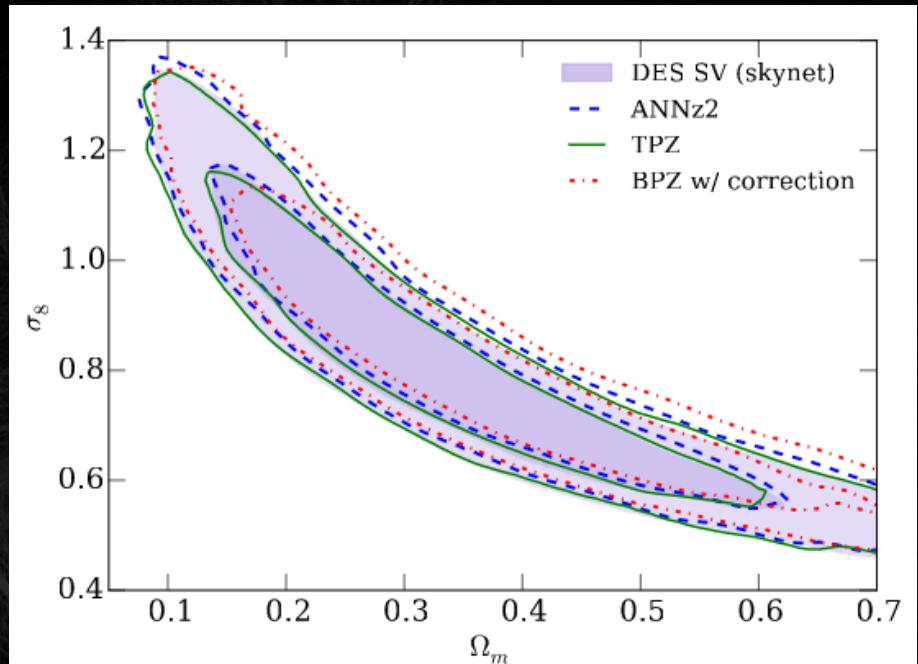
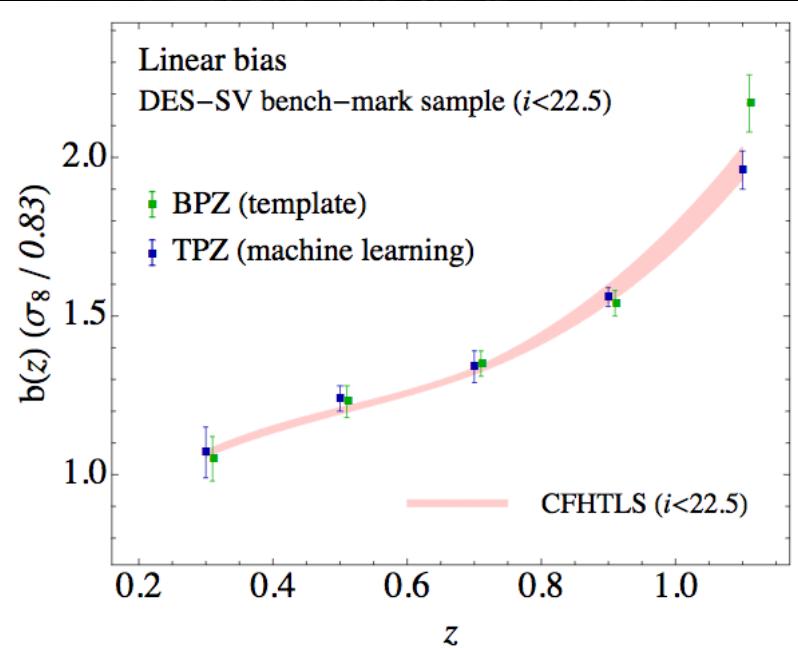
Sánchez, Carrasco Kind, et al. 2014 (MNRAS, submitted)



Galaxy clustering, photometric redshifts and diagnosis of systematics in the DES SV data (arXiv:1507.05360)

Cosmology from Cosmic Shear with DES SV Data (arXiv:1507.05552)

# Photo-z for DES SV data



Galaxy clustering, photometric redshifts and diagnosis of systematics in the DES SV data (arXiv:1507.05360)

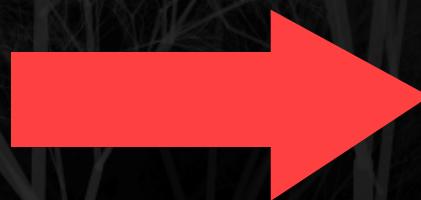
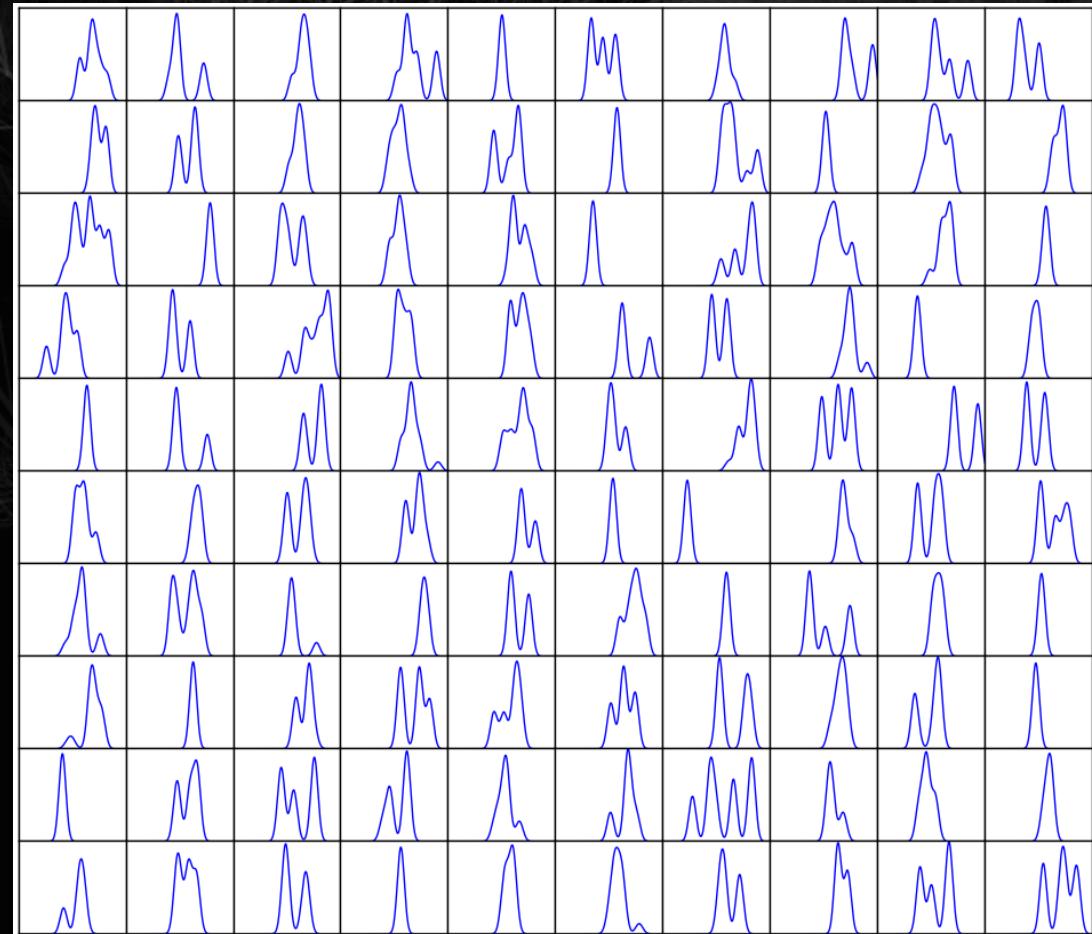
... and 14 other DES papers!

Cosmology from Cosmic Shear with DES SV Data (arXiv:1507.05552)

For the first time using Machine Learning photo-z in cosmology!

- Volume of data (SV: 25M, Y1: 130M, Y3: 250M)
- Training not only garbage in → garbage out BUT features in → features out, overfitting, manifolds, etc...
- Cosmic variance, biases, magnitude errors, uniformity of spec samples, cross matching, etc
- Software sustainability, distribution and license.
- PDF generation and storage. Multiple methods, multiple versions!

# 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

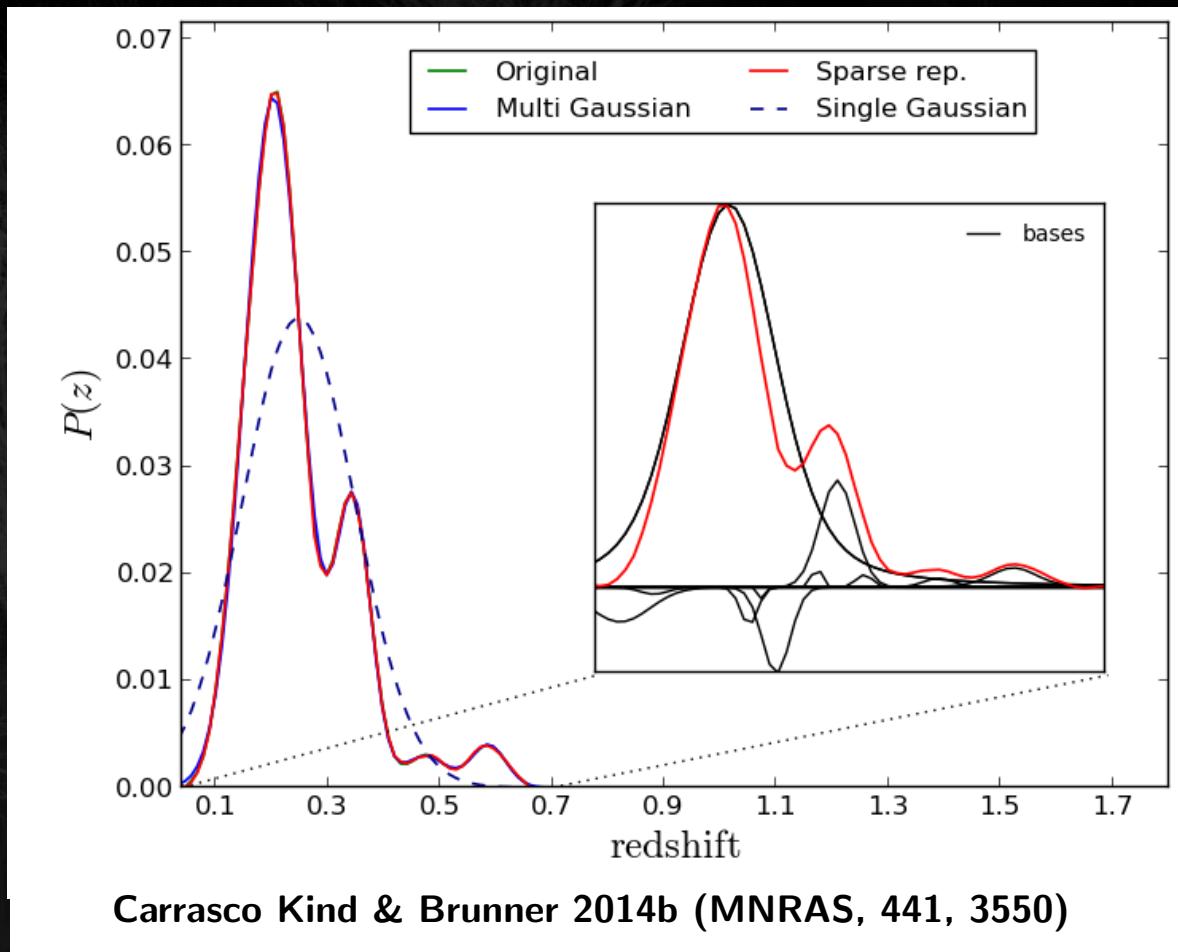
Single Gaussian fit

Multi-Gaussian fit

Monte Carlo sampling

Sparse representation  
techniques

Reduce number of points  
while increasing accuracy



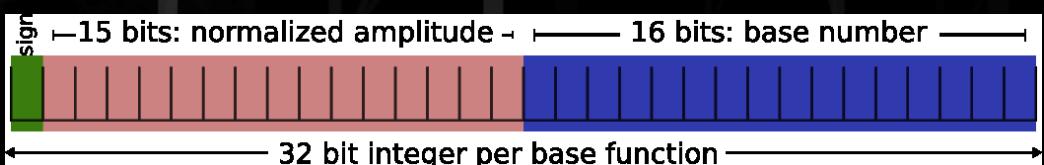
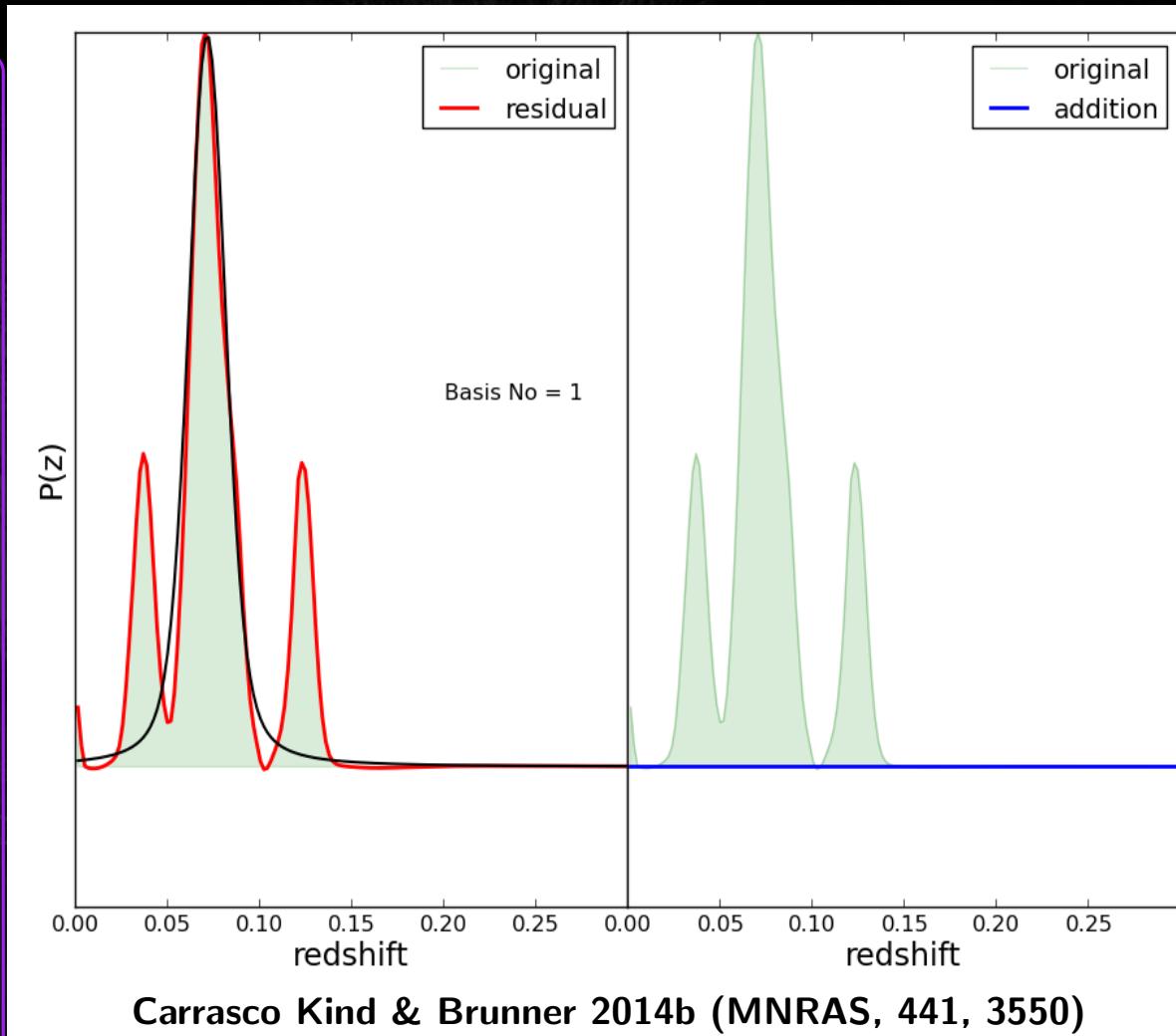
# Photo- $z$ PDF storage: Sparse representation

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

With only 10-20 bases achieve 99.9 % accuracy

Use 32-bits integer per basis, compression

Store Multiple PDFs



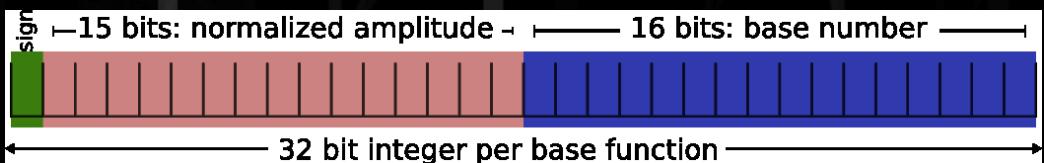
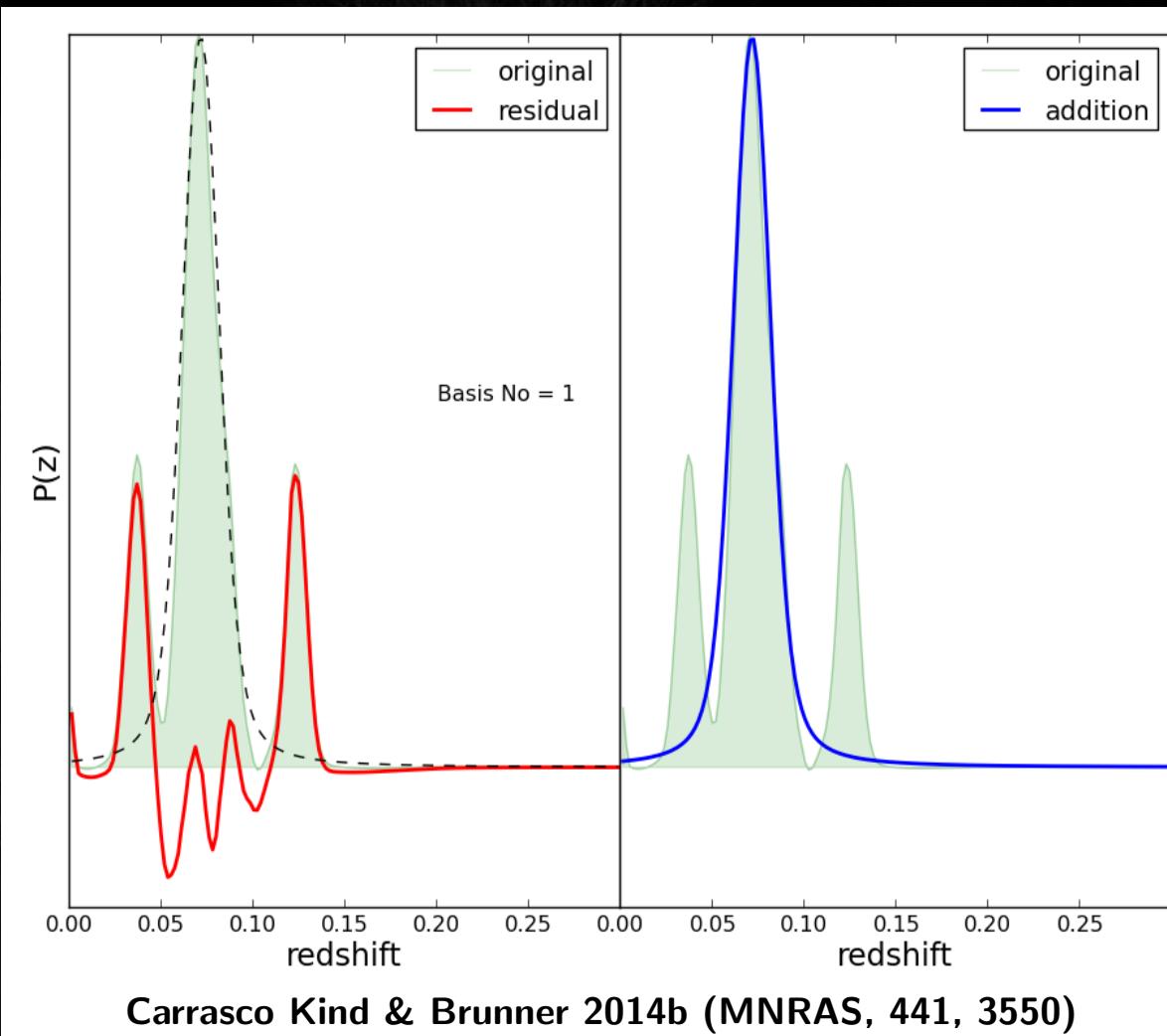
# Photo- $z$ PDF storage: Sparse representation

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

With only 10-20 bases achieve 99.9 % accuracy

Use 32-bits integer per basis, compression

Store Multiple PDFs



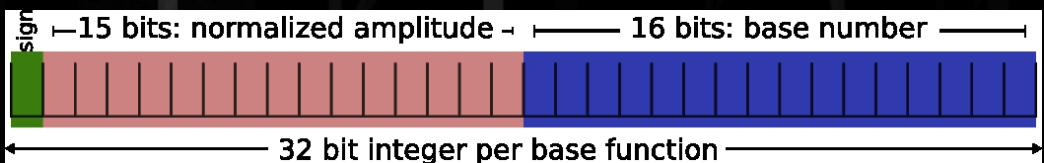
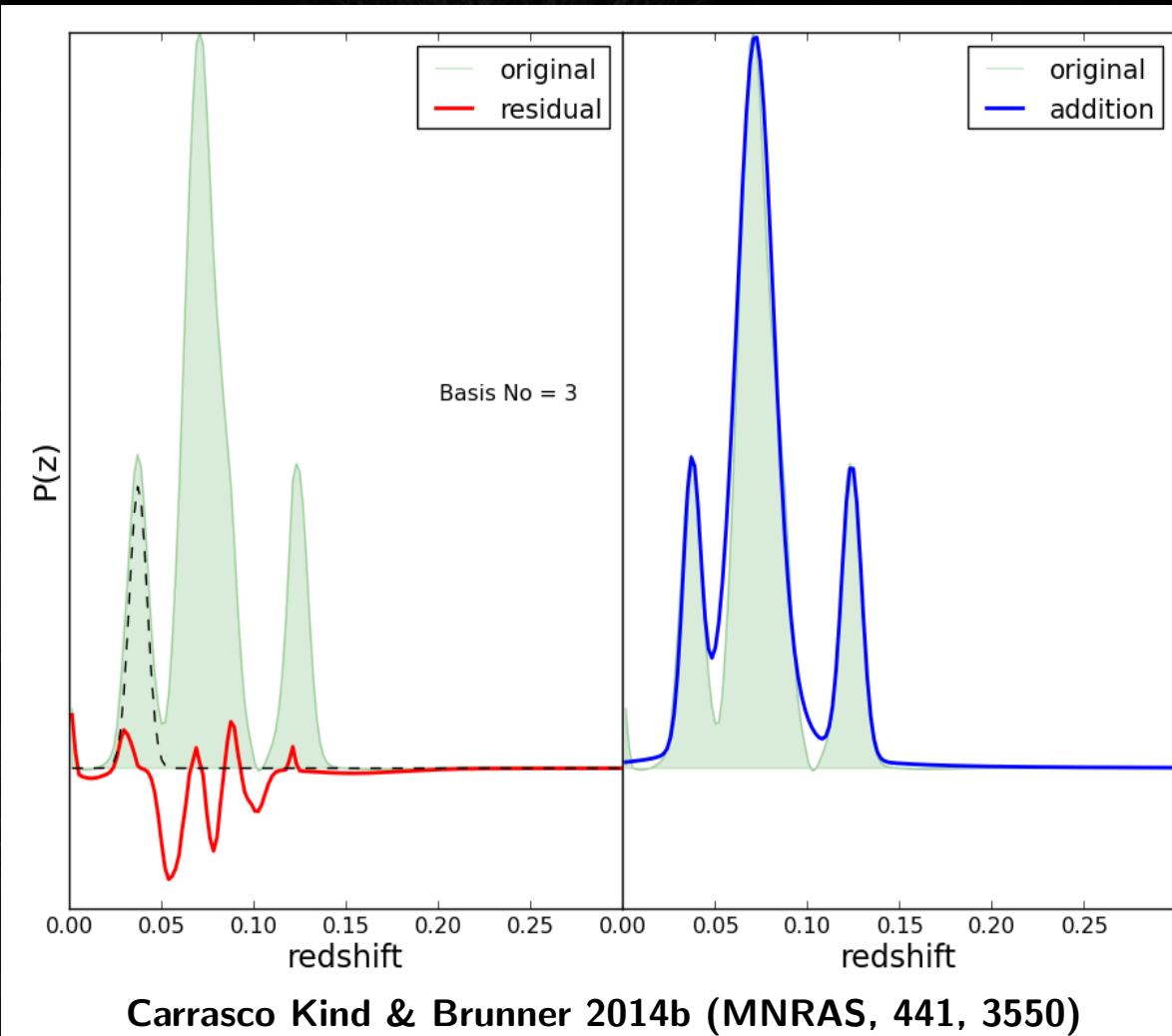
# Photo- $z$ PDF storage: Sparse representation

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

With only 10-20 bases achieve 99.9 % accuracy

Use 32-bits integer per basis, compression

Store Multiple PDFs



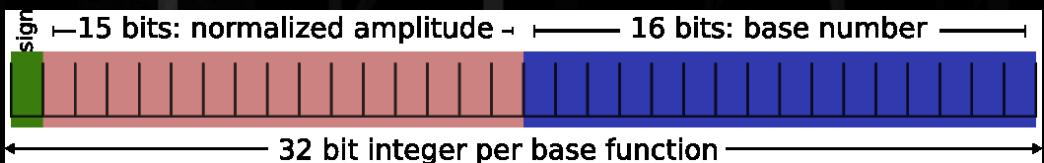
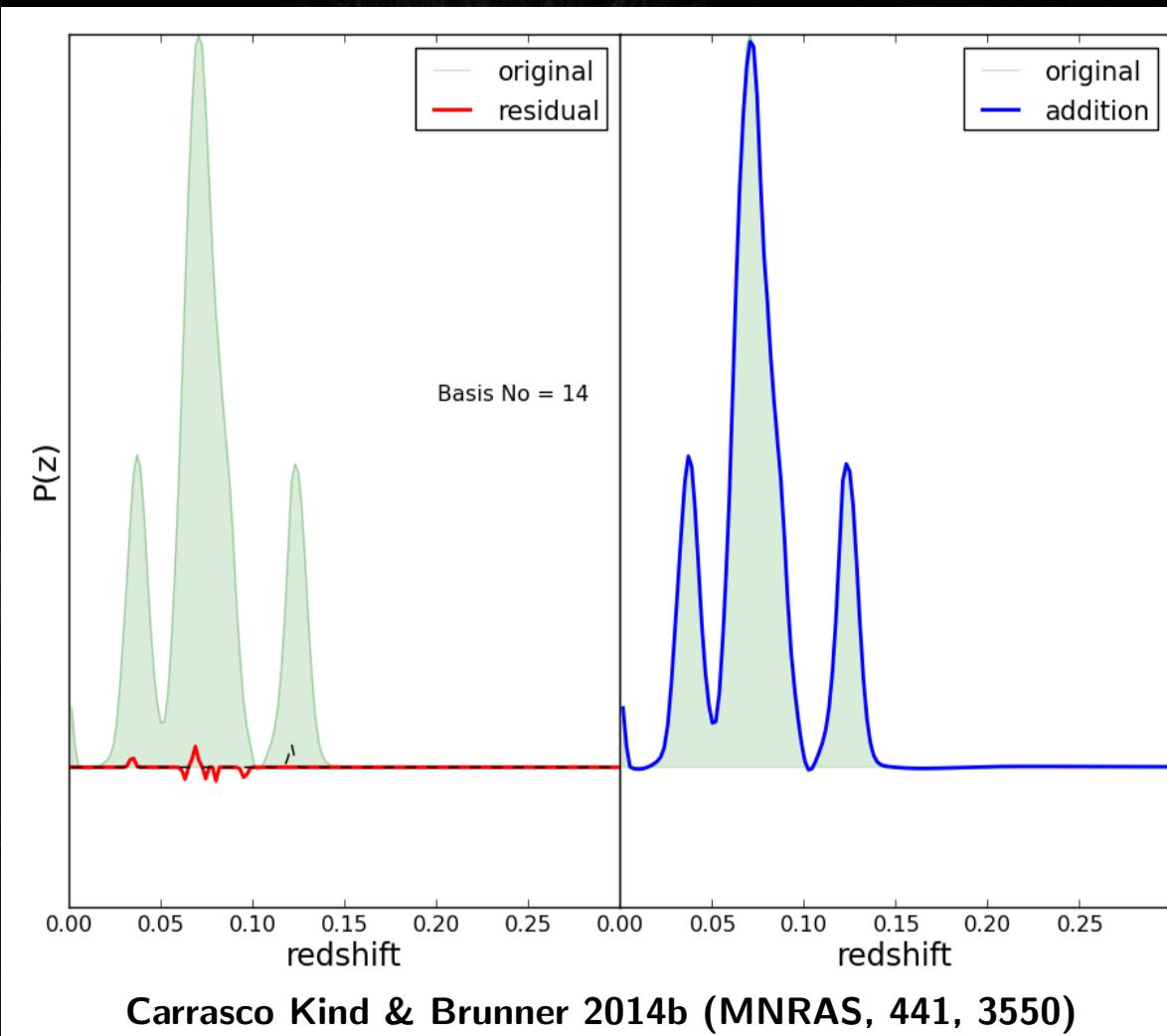
# Photo- $z$ PDF storage: Sparse representation

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

With only 10-20 bases achieve 99.9 % accuracy

Use 32-bits integer per basis, compression

Store Multiple PDFs



By definition:

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

# $N(z)$ and sparse representation

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  $pz_k$  as:

$\mathbf{pz}_k \approx \mathbf{D} \cdot \boldsymbol{\delta}_k$   $\mathbf{D}$  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:

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

# $N(z)$ and sparse representation

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:

$\mathbf{p}_{\mathbf{z}_k} \approx \mathbf{D} \cdot \boldsymbol{\delta}_k$   $\mathbf{D}$  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

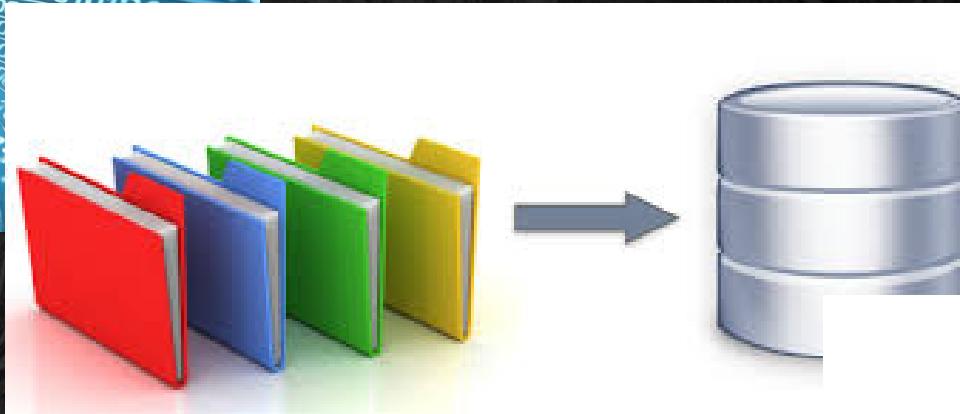
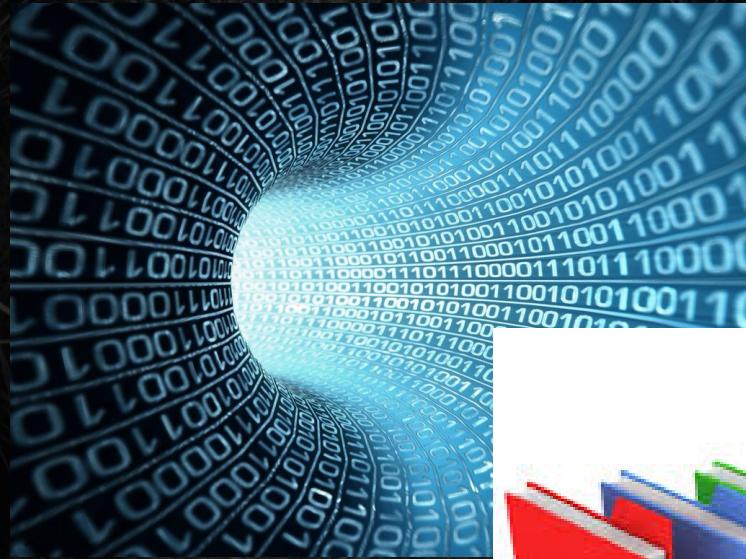
by precomputing:

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

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

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

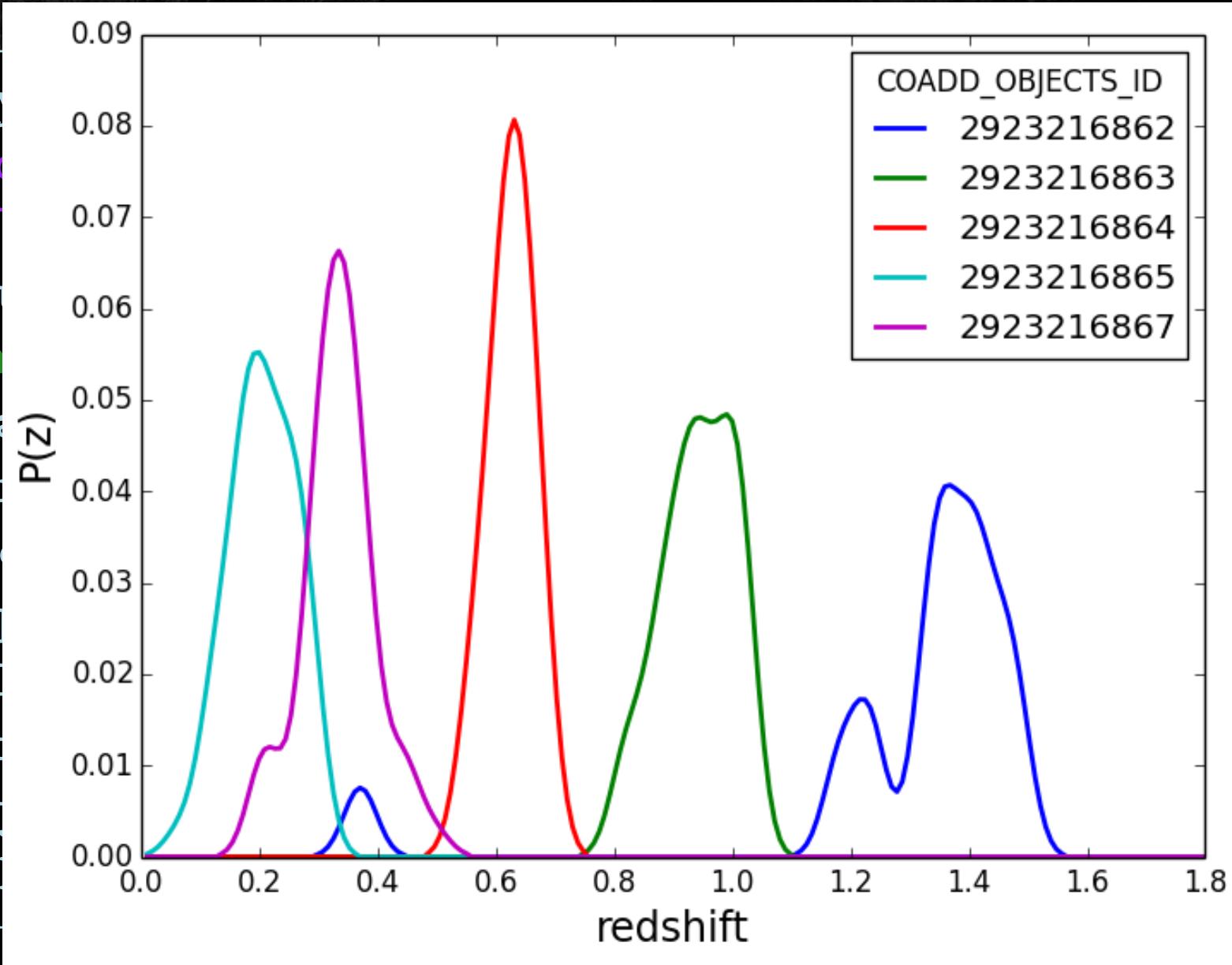
# How we can do it?



```
query="""
select COADD_OBJECTS_ID ,TPZ from
PHOTOZ_PDF_SVA1_GOLD where rownum < 6"""
cc=cursor.execute(query)
#Handling and plot
df=ea.to_pandas(cc)
for i in xrange(5):
    cid=df.COADD_OBJECTS_ID.values[i]
    plt.plot(zbins,df.TPZ.values[i],
              lw=2,label=cid)
plt.xlabel('redshift',fontsize=17)
plt.ylabel('P(z)',fontsize=17)
plt.legend(loc=0, title='COADD_OBJECTS_ID')
```

# Getting some PDFs from DB

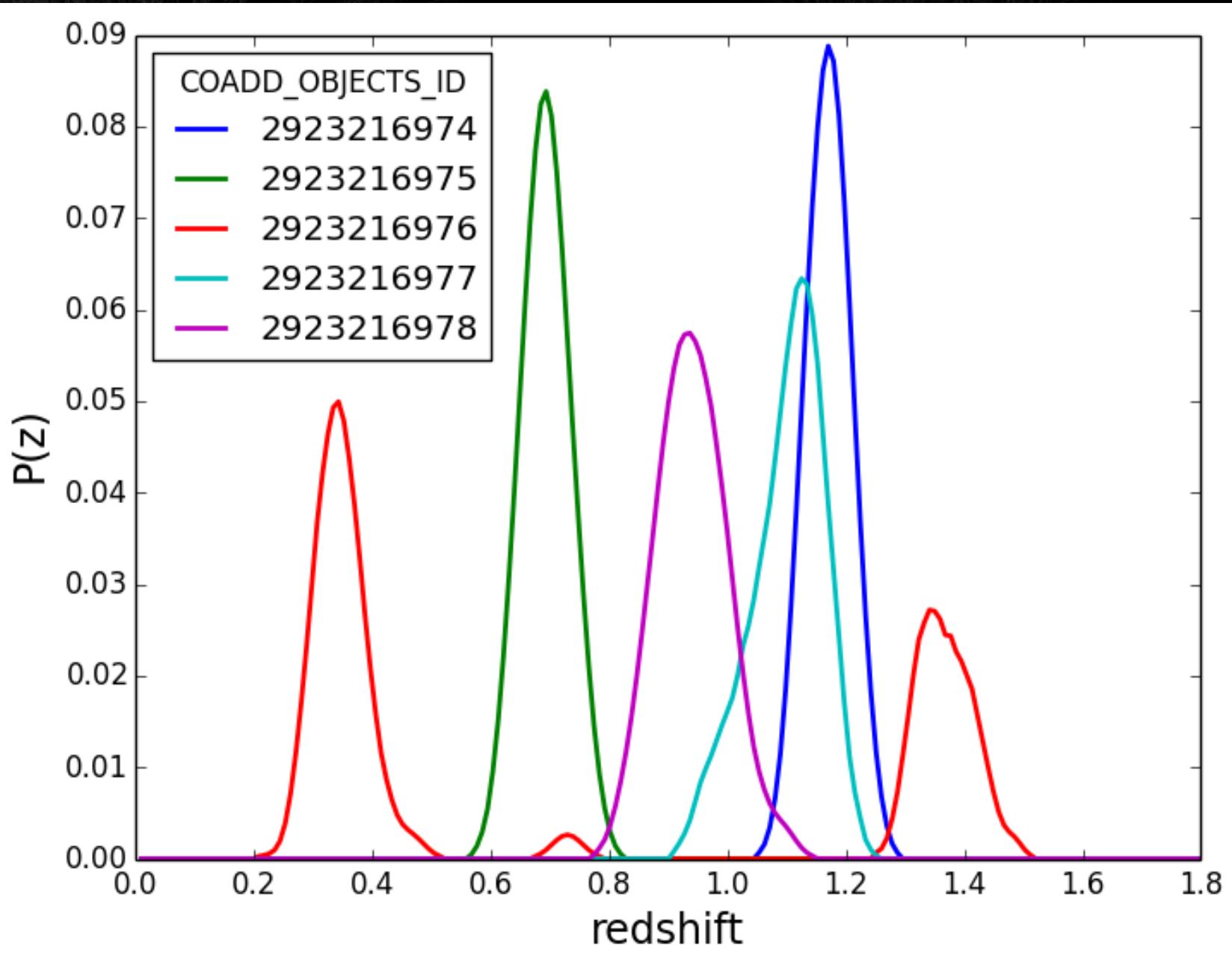
```
query  
selected  
PHOT  
cc=cc  
#Hand  
df=e  
for  
plt.  
plt.  
plt.
```



```
query="""
select COADD_OBJECTS_ID ,PHZ.GET_PDF(TPZ) as
TPZ from PHOTOZ_SPARSE_SVA1_GOLD
where rownum < 6"""
cc=cursor.execute(query)
#Handling and plot
df=ea.to_pandas(cc)
for i in xrange(5):
    cid=df.COADD_OBJECTS_ID.values[i]
    plt.plot(zbins,df.TPZ.values[i],
              lw=2,label=cid)
plt.xlabel('redshift',fontsize=17)
plt.ylabel('P(z)',fontsize=17)
plt.legend(loc=0, title='COADD_OBJECTS_ID')
```

```
query="""
select COADD_OBJECTS_ID, PHZ.GET_PDF(TPZ) as
TPZ from PHOTOZ_SPARSE_SVA1_GOLD
where rownum < 6"""
cc=cursor.execute(query)
#Handling and plot
df=ea.to_pandas(cc)
for i in xrange(5):
    cid=df.COADD_OBJECTS_ID.values[i]
    plt.plot(zbins,df.TPZ.values[i],
              lw=2,label=cid)
plt.xlabel('redshift', fontsize=17)
plt.ylabel('P(z)', fontsize=17)
plt.legend(loc=0, title='COADD_OBJECTS_ID')
```

# Now using Sparse rep.



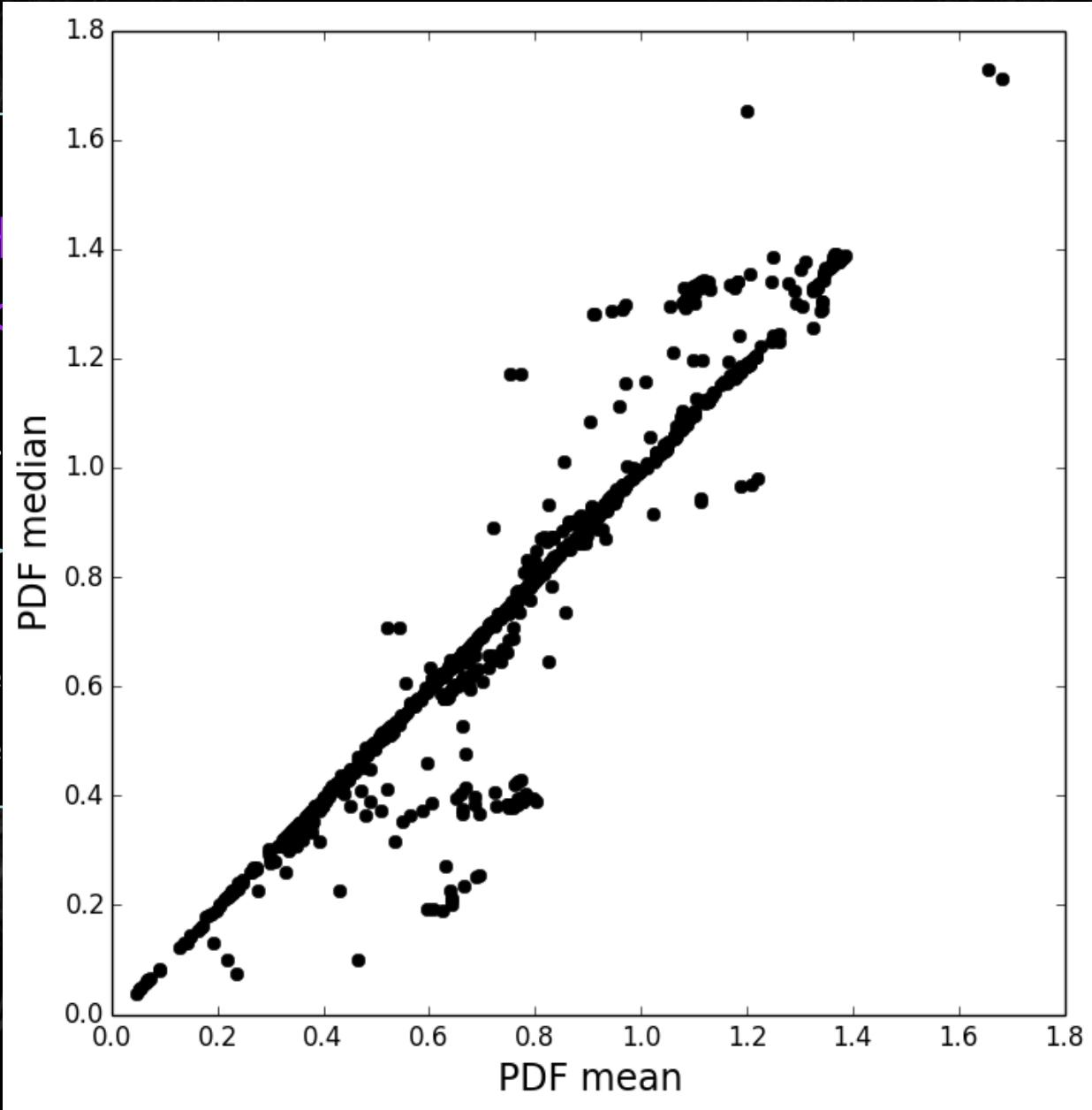
```
query="""
Select PHZ.MEAN(tpz) mean, PHZ.MEDIAN(tpz)
median from PHOTOZ_PDF_SVA1_GOLD
where rownum < 1000"""

cc=cursor.execute(query)
df=ea.to_pandas(cc)
plt.plot(df.MEAN, df.MEDIAN, 'ko')
plt.xlabel('PDF mean', fontsize=17)
plt.ylabel('PDF median', fontsize=17)
```

# Getting metrics on the fly

```
query="""
Select PH
median from
where row
cc=cursor
df=ea.to_
plt.plot(
plt.xlabel(
plt.ylabel
```

pz )



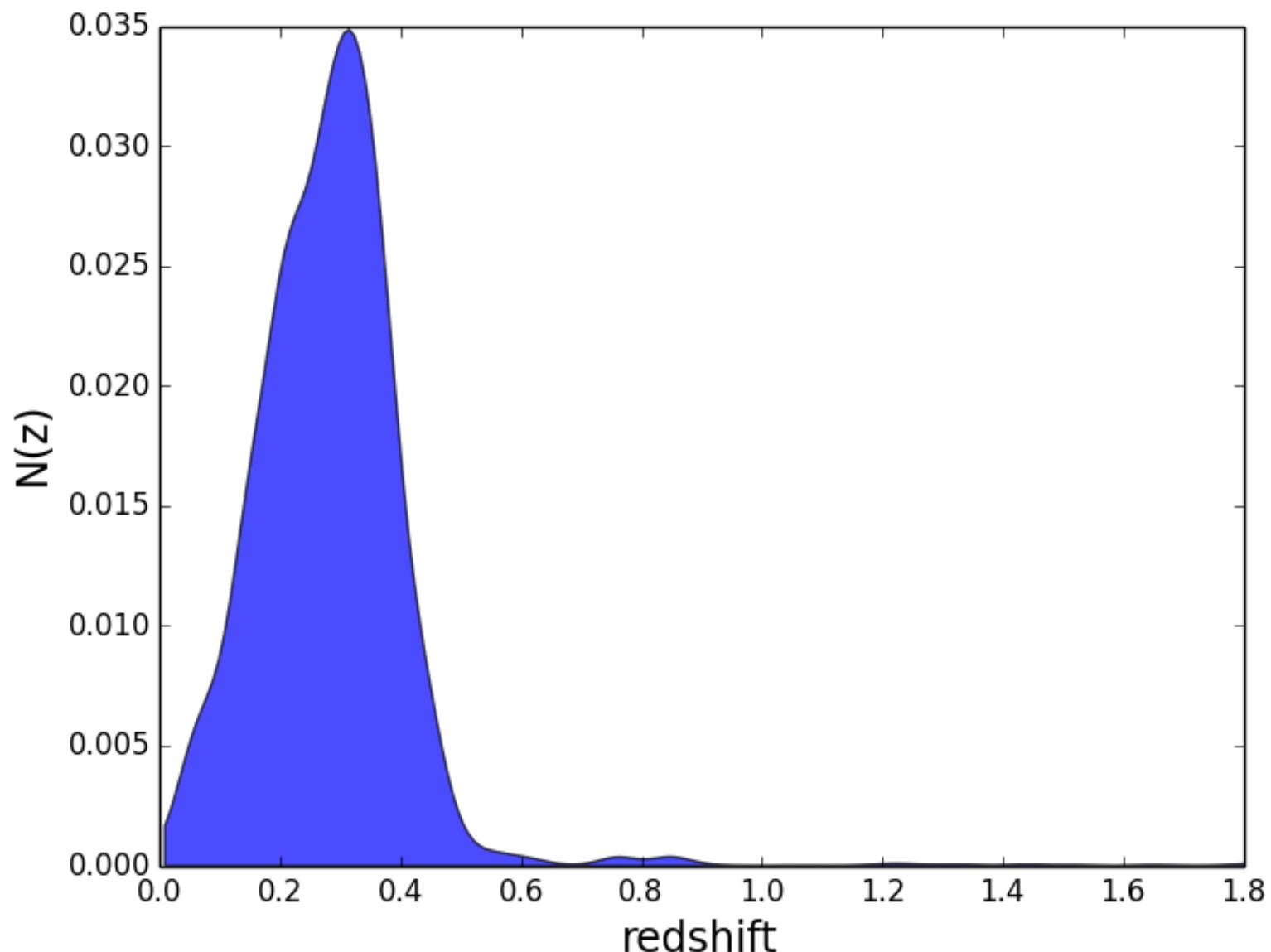
```
query="""
Select NZ(PHZ.TOTABLE(tpz)) as NZ from
PHOTOZ_PDF_SVA1_GOLD where
PHZ.MEAN(tpz) BETWEEN 0.1 and 0.4
and rownum < 100000"""
cc=cursor.execute(query)
df=ea.to_pandas(cc)
plt.fill_between(zbins,df.NZ.values[0],
                 facecolor='blue',alpha=0.7)
plt.xlabel('redshift',fontsize=17)
plt.ylabel('N(z)',fontsize=17)
```

```
query="""
Select NZ(PHZ.TOTABLE(tpz)) as NZ from
PHOTOZ_PDF_SVA1_GOLD where
PHZ.MEAN(tpz) BETWEEN 0.1 and 0.4
and rownum < 100000"""

cc=cursor.execute(query)
df=ea.to_pandas(cc)
plt.fill_between(zbins,df.NZ.values[0],
                 facecolor='blue',alpha=0.7)
plt.xlabel('redshift',fontsize=17)
plt.ylabel('N(z)',fontsize=17)
```

# Stacking PDFs in DB cluster!

```
query=  
Select  
PHOTOZ  
PHZ.ME  
and row  
cc=cur  
df=ea.  
plt.fi  
f  
plt.xl  
plt.yl
```



# ✓ 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!

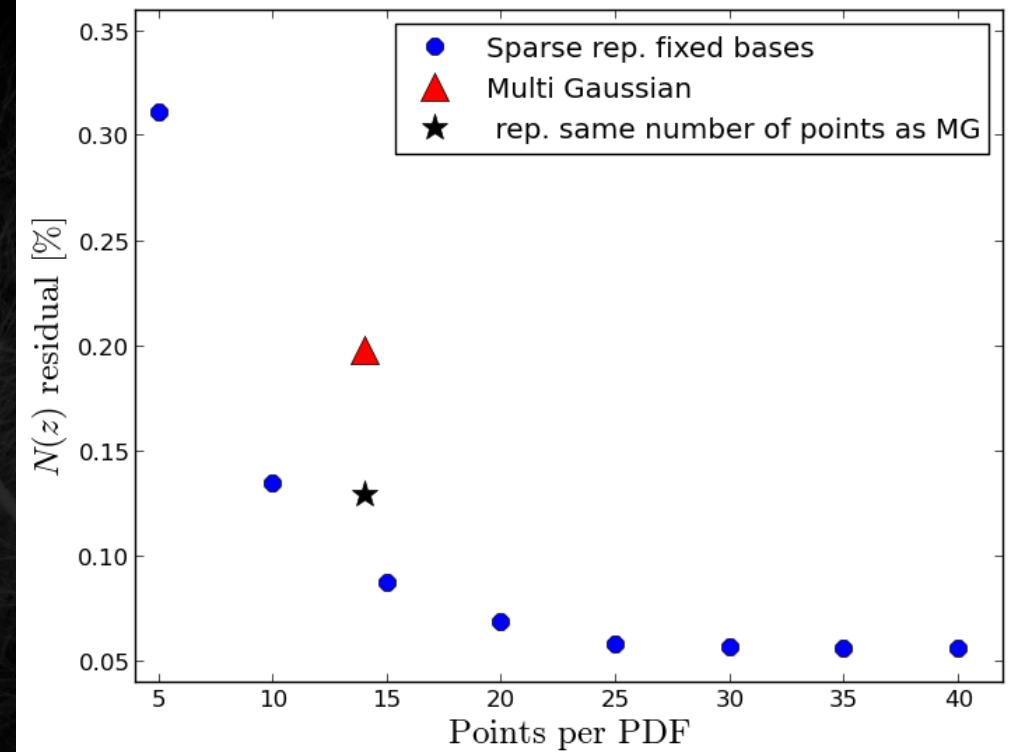
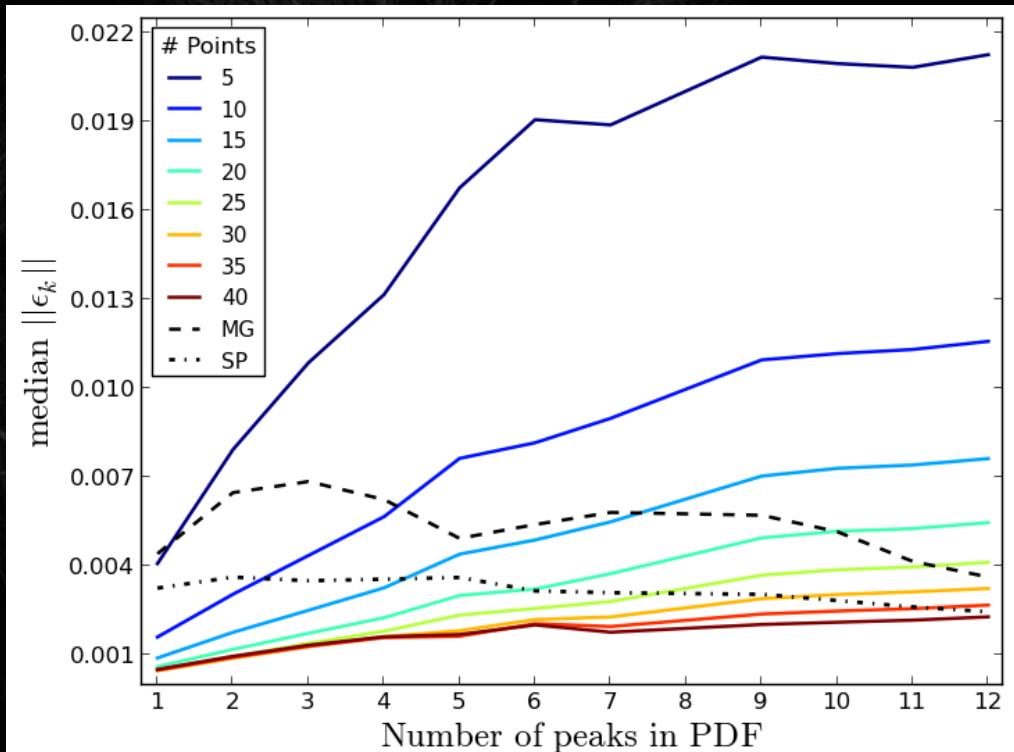
Clustering, weak lensing, DES, DESDM, etc...

# Questions?

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

# EXTRA SLIDES

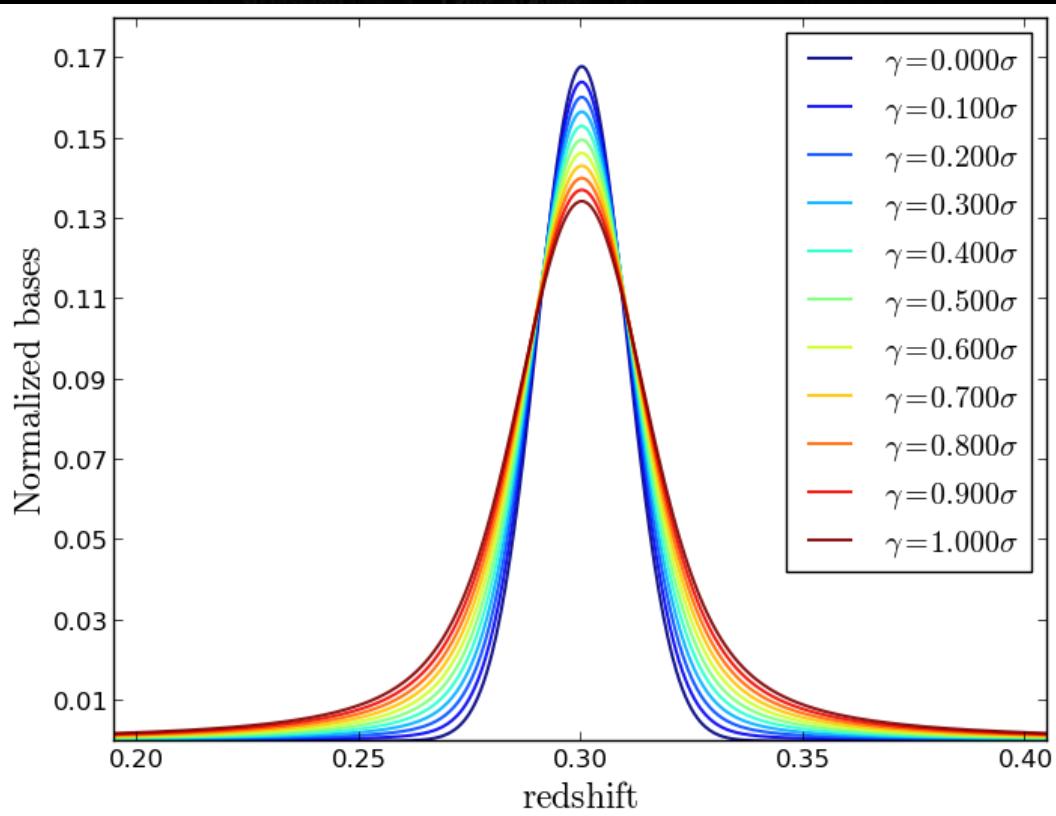
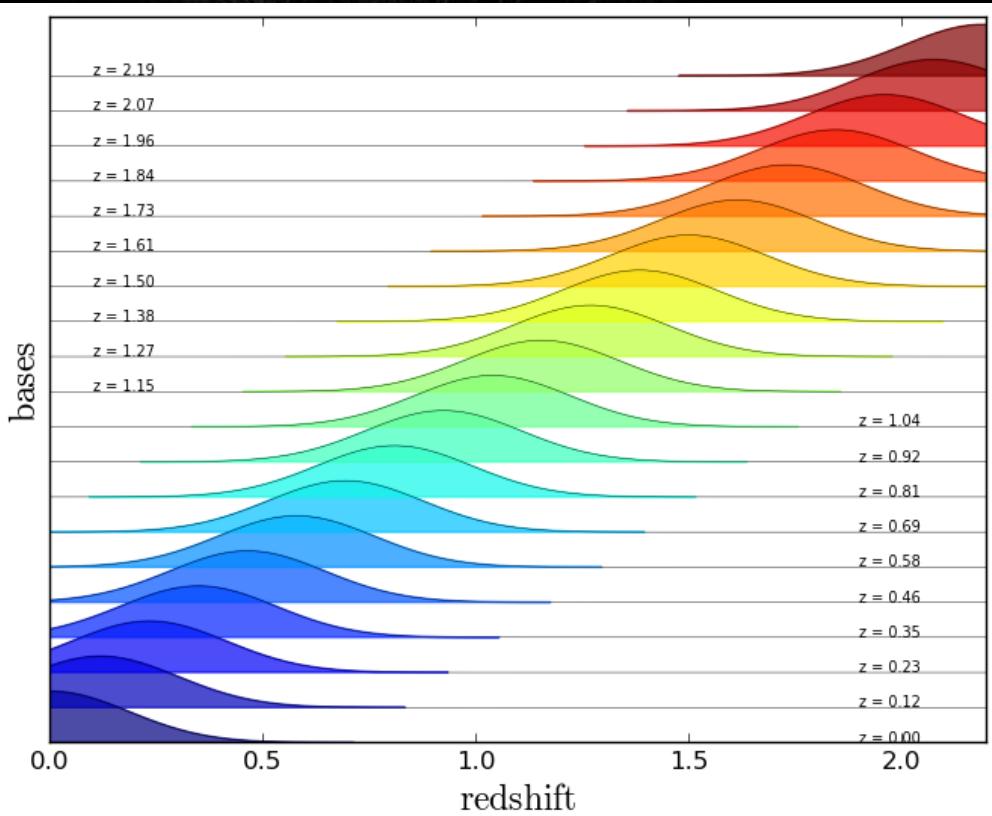
# Photo- $z$ PDF storage: Results



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

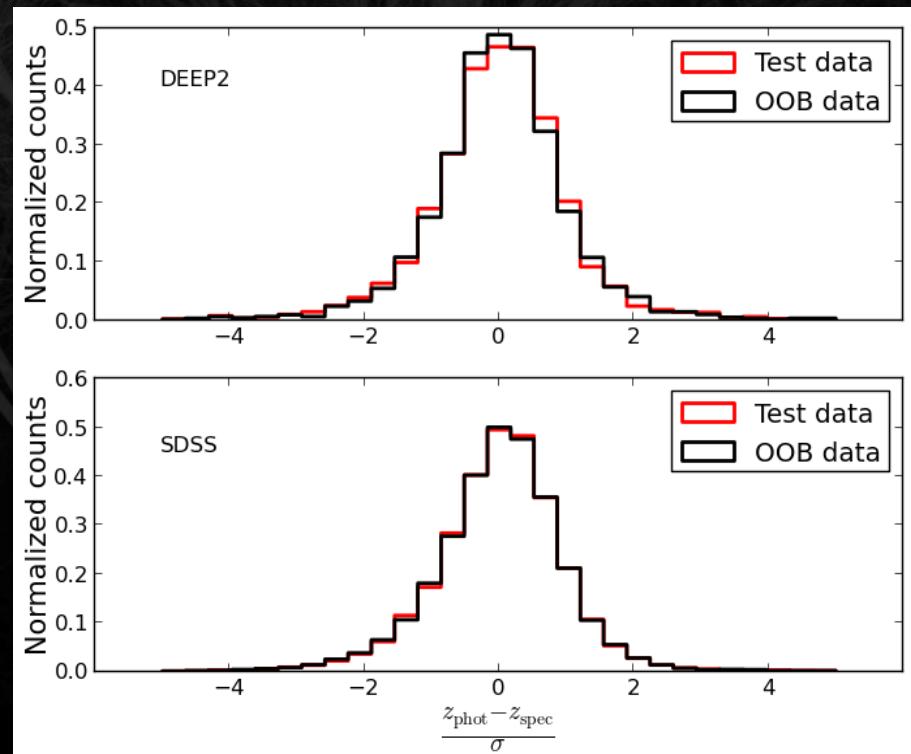
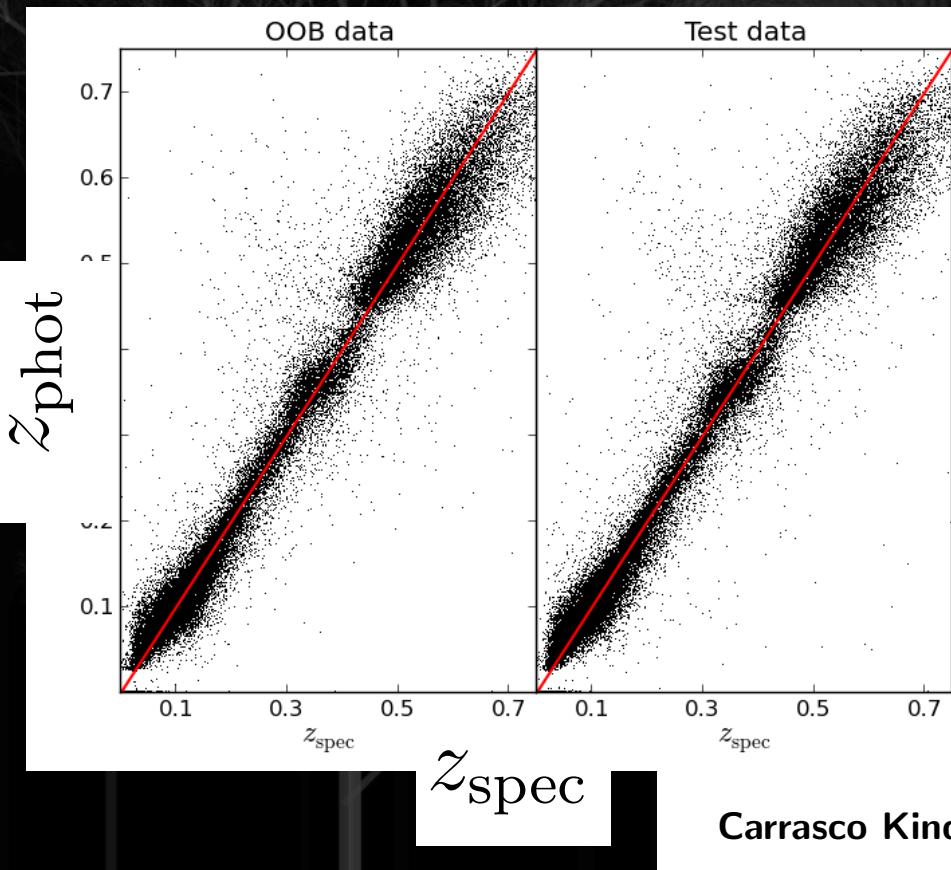
We can learn from the cross-validation data!

# Photo- $z$ PDF estimation: Error and validation

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

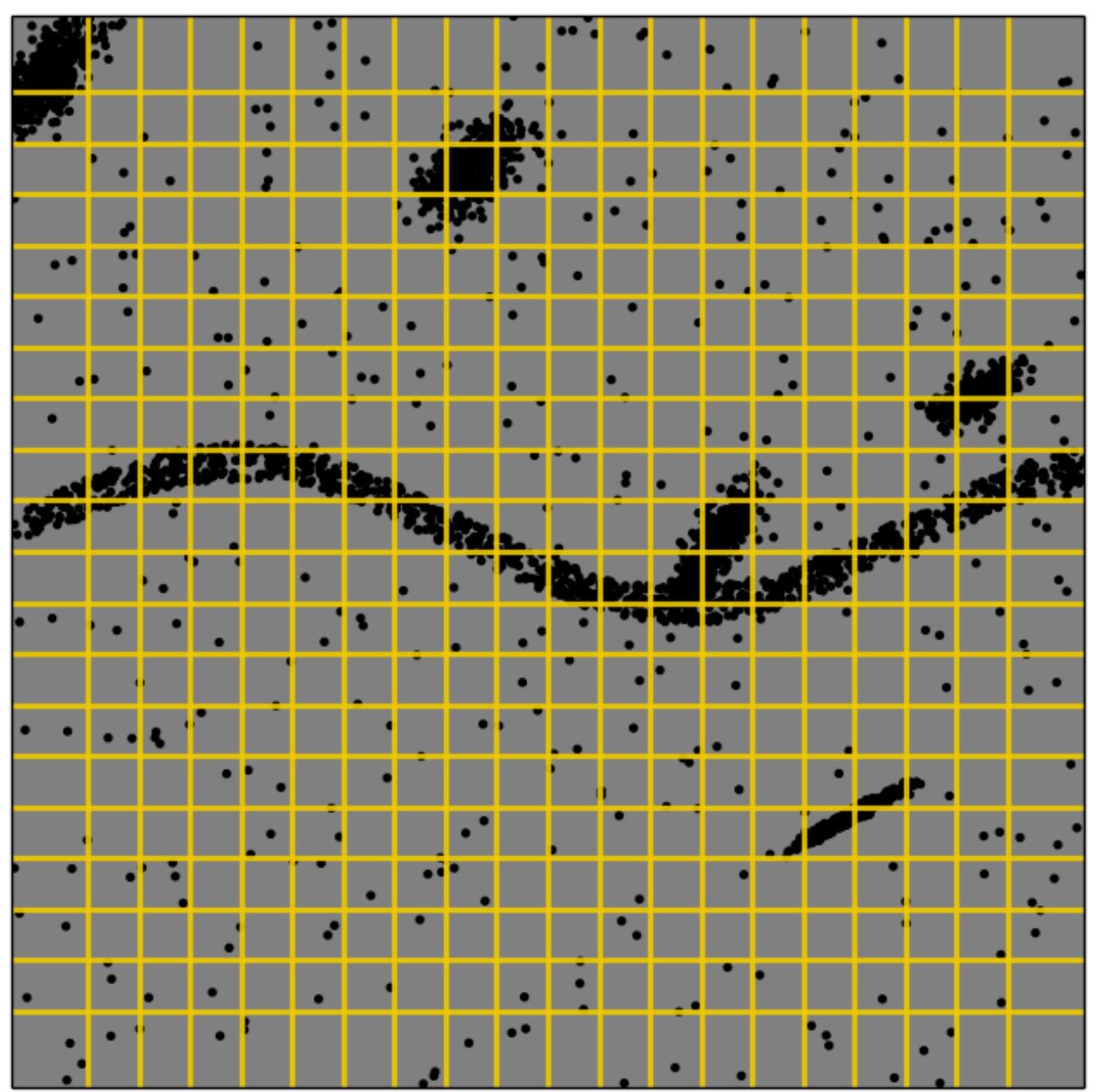
Changes for every tree/map and is not used during training

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

