



# Photo-z's PDF I/O, DB, DES (Lessons learned ... so far)

Matías Carrasco Kind

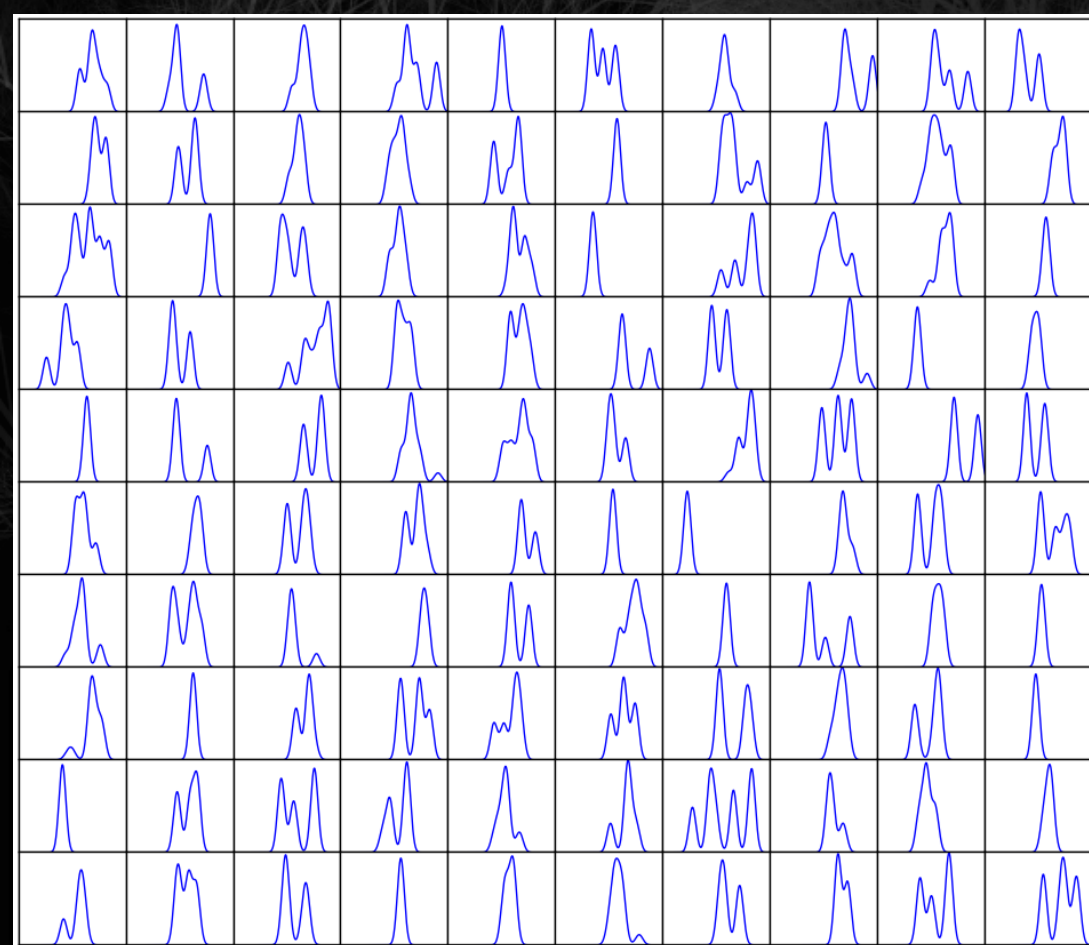
NCSA/Department of Astronomy  
University of Illinois at Urbana-Champaign

LSST\_DESC @ Argonne  
October 29<sup>th</sup>, 2015

- Photo- $z$  PDF becoming important
- 10 DES papers using photo- $z$  (WL, 2pt WL, 2pt LSS, DES X CMB, Systematics, redmagic, etc)
- Asorey et al 2015 in prep. includes PDF in clustering analysis.
- Storage and distribution still an issue
- Compression techniques



# Photo- $z$ PDF representation and storage in DES DB



- Single Gaussian fit
- Multi-Gaussian fit
- Monte Carlo sampling
- Sparse representation techniques
- Reduce number of points while increasing accuracy

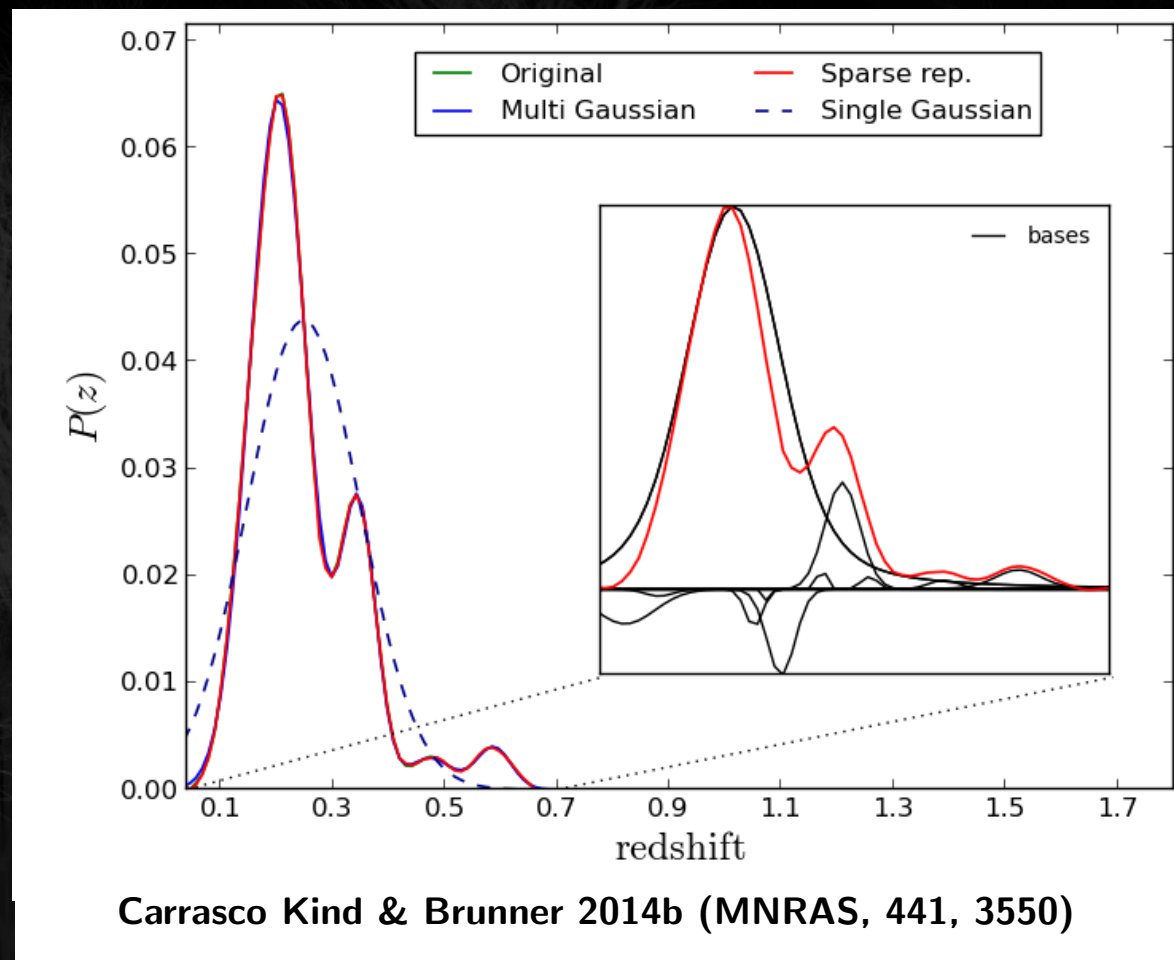
Single Gaussian fit

Multi-Gaussian fit

Monte Carlo sampling

Sparse representation  
techniques

Reduce number of points  
while increasing accuracy



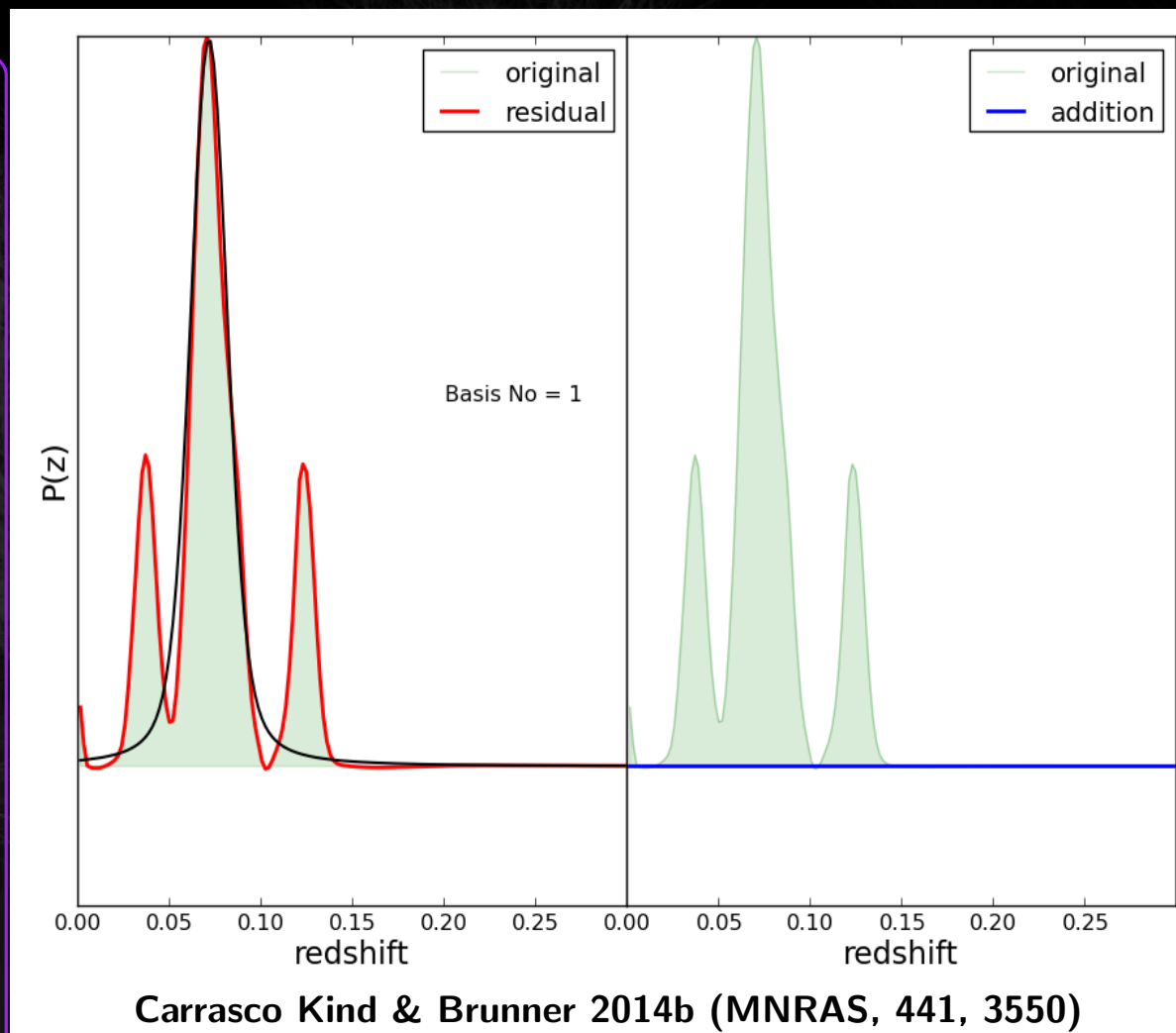


Use Gaussian and Voigt profiles as bases, need  $N_{\text{original}}^2$  bases

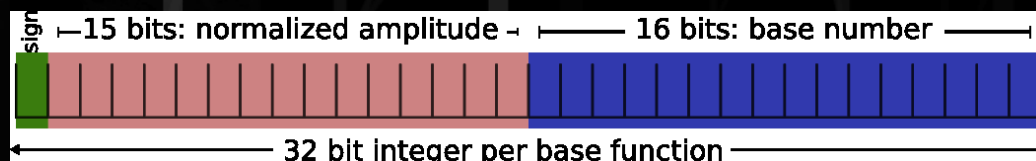
With only 10-20 bases achieve 99.9 % accuracy

Use 32-bits integer per basis, compression

Store Multiple PDFs



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

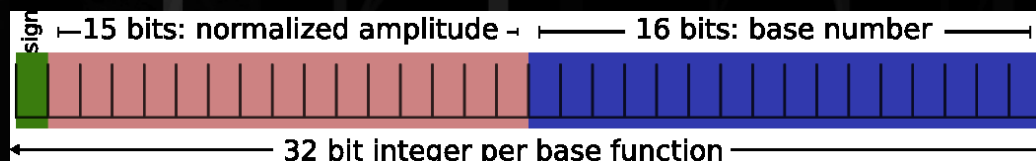
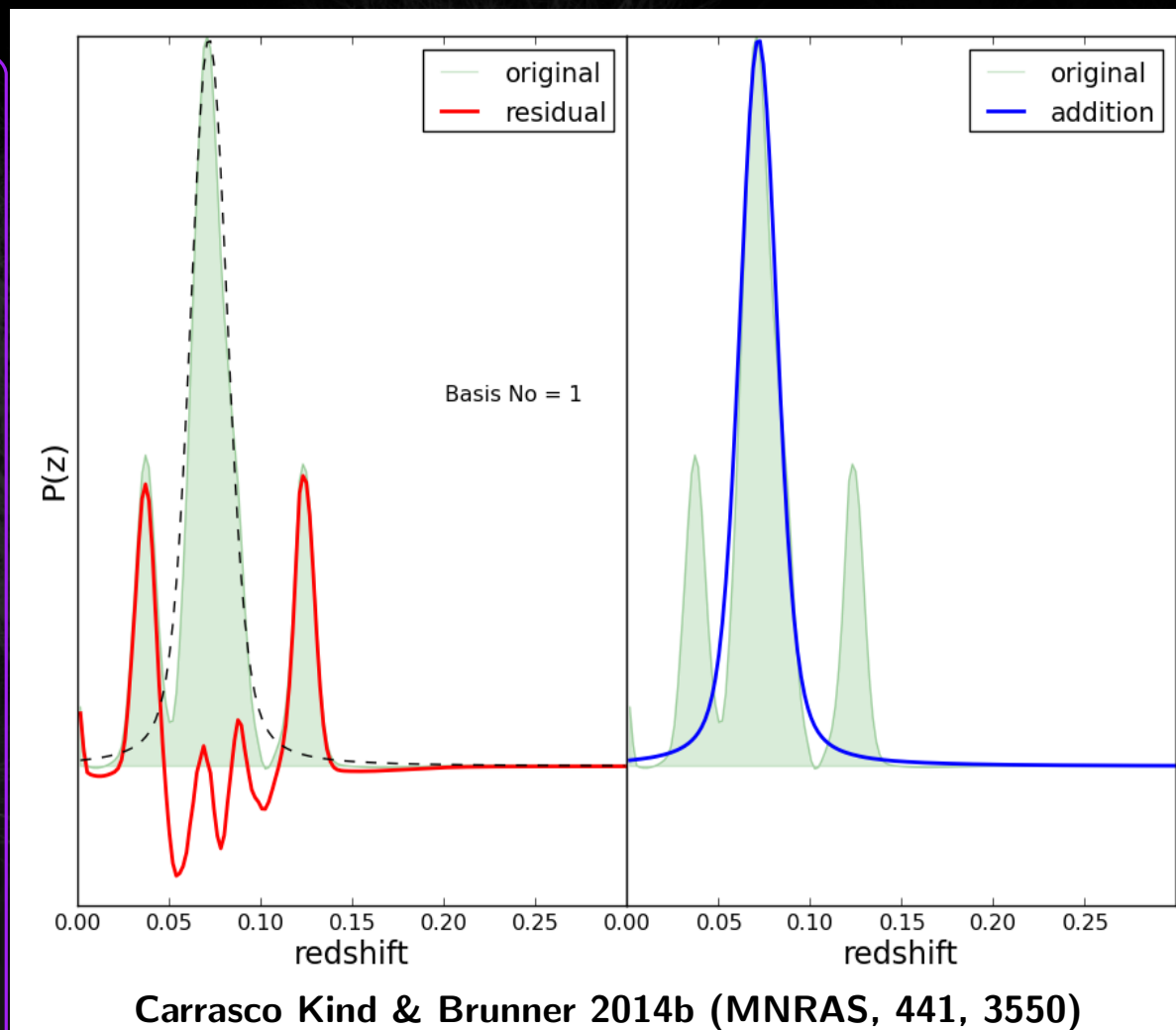


Use Gaussian and Voigt profiles as bases, need  $N_{\text{original}}^2$  bases

With only 10-20 bases achieve 99.9 % accuracy

Use 32-bits integer per basis, compression

Store Multiple PDFs

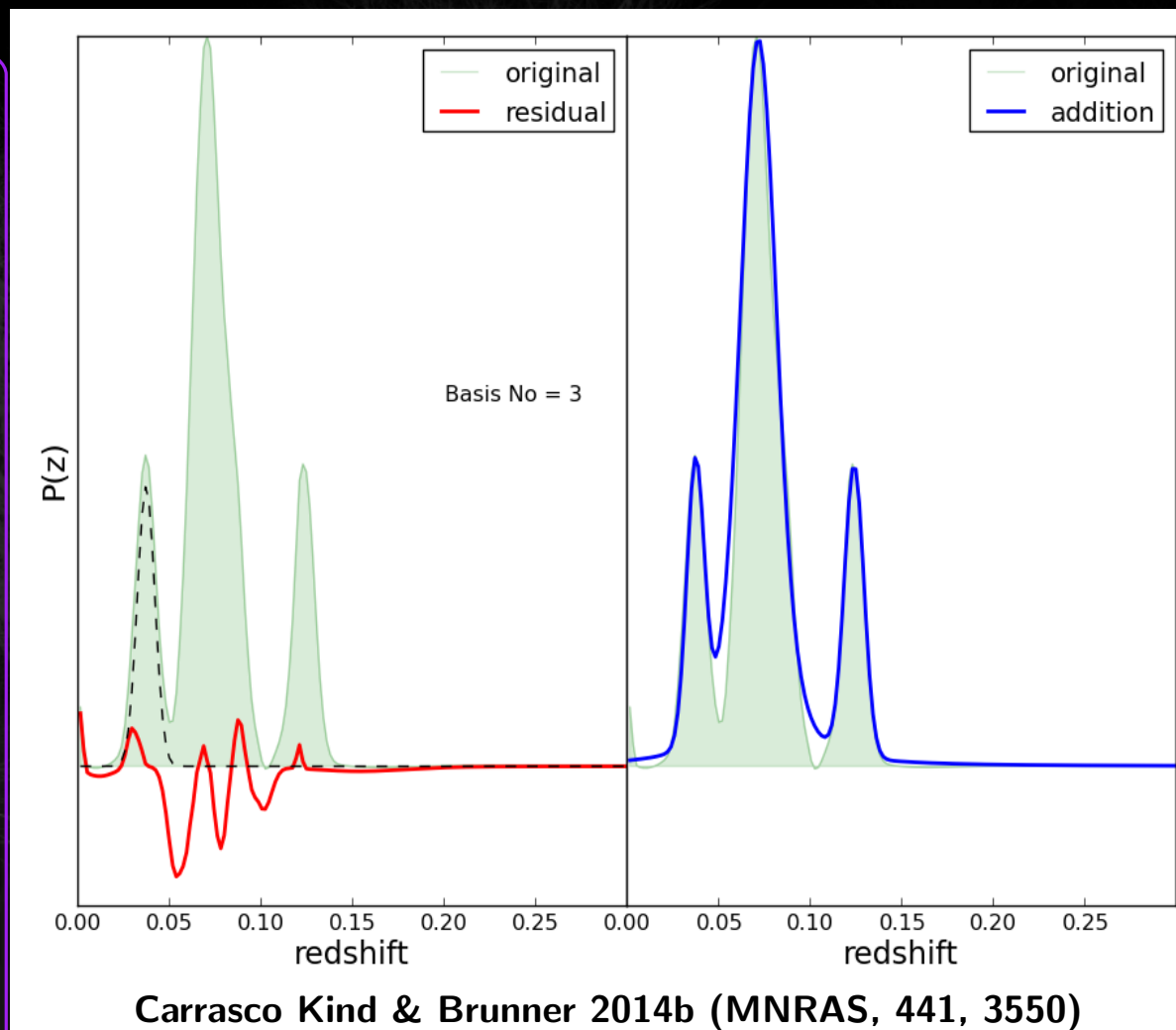


Use Gaussian and Voigt profiles as bases, need  $N_{\text{original}}^2$  bases

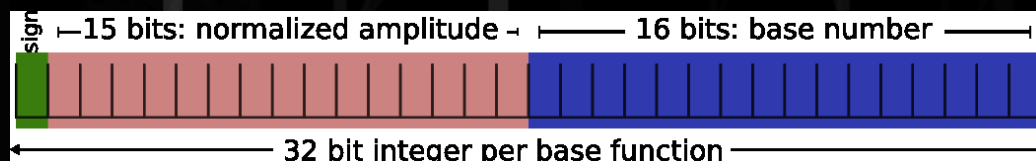
With only 10-20 bases achieve 99.9 % accuracy

Use 32-bits integer per basis, compression

Store Multiple PDFs



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



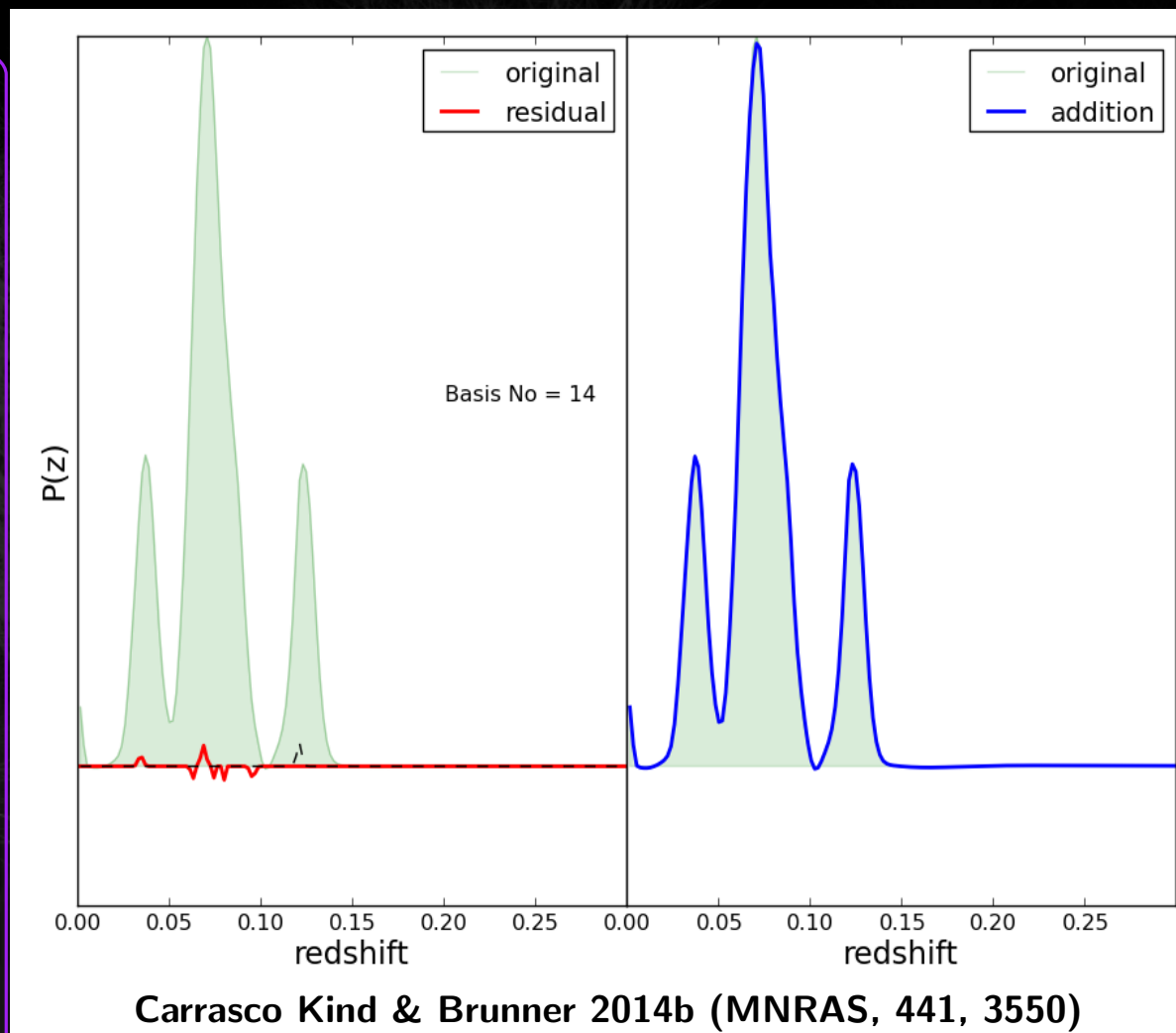


Use Gaussian and Voigt profiles as bases, need  $N_{\text{original}}^2$  bases

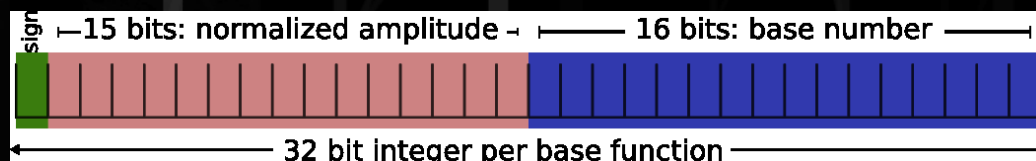
With only 10-20 bases achieve 99.9 % accuracy

Use 32-bits integer per basis, compression

Store Multiple PDFs



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



# How we can do it?

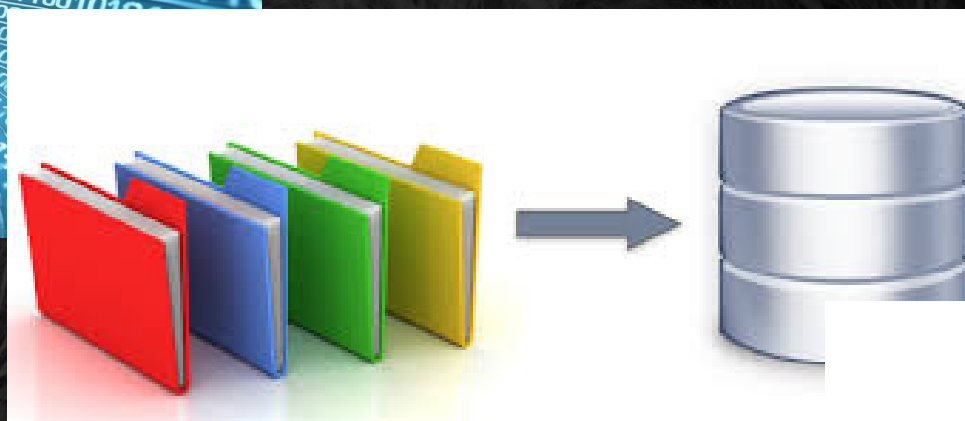
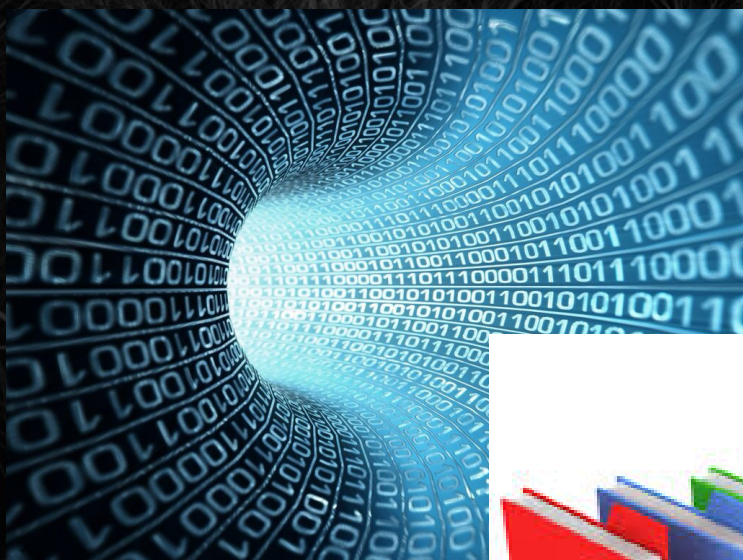




Photo-z for ALL objects!

Y1A1  $\sim$  130 M objects

Y2A1 expected to have more than  
100 M objects

PDF methods require lot of  
space, sparse rep as alternative.



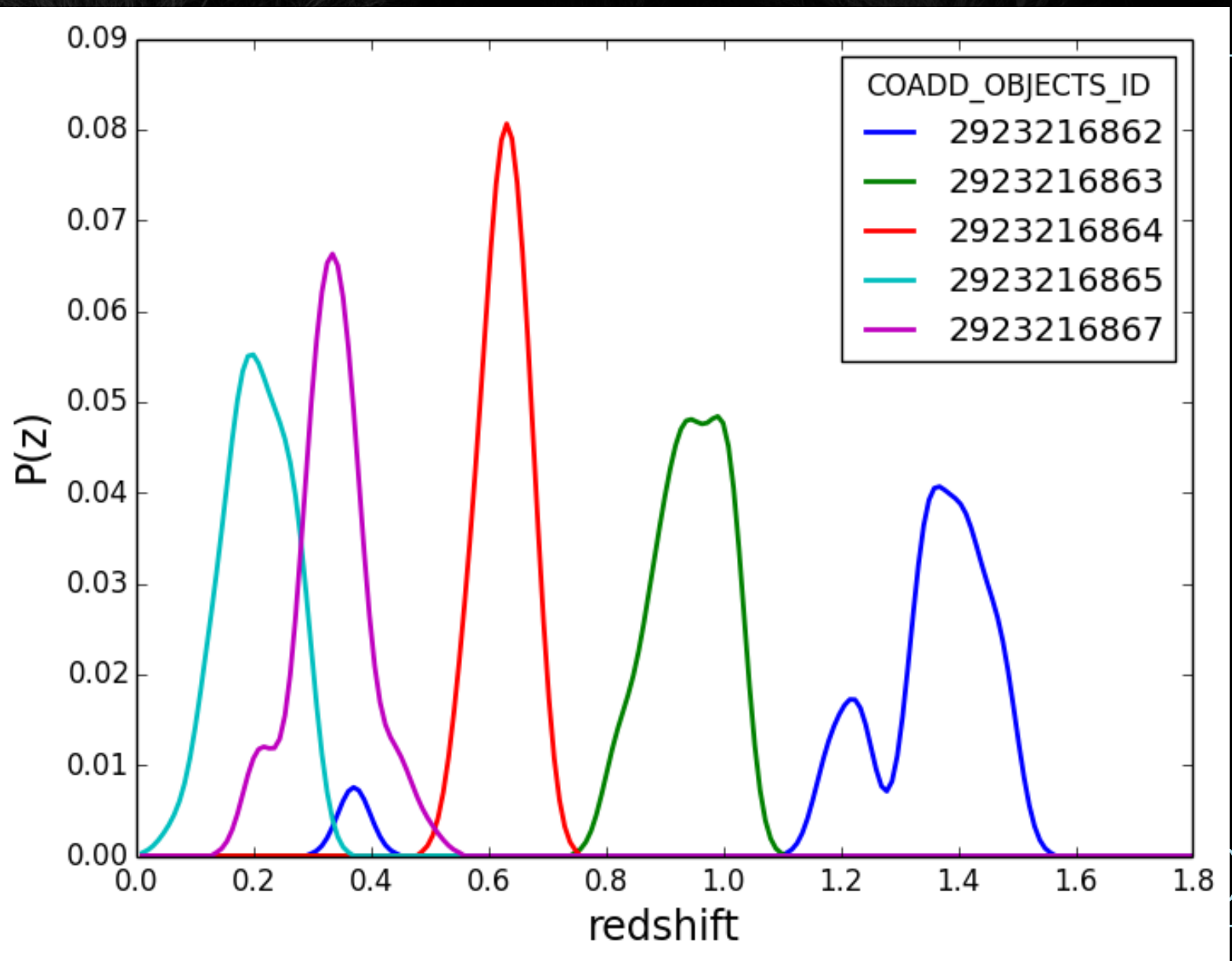
<https://opensource.ncsa.illinois.edu/confluence/display/DESDM/Access+to+photo-z+from+DB>

- New data types
  - \* PFULL  $\rightarrow$  200-vector
  - \* PSPARSE  $\rightarrow$  20-vector
  - \* PFULL\_TB (ancillary)
- New functions
  - \* GET\_PDF (PSPARSE TYPE)
  - \* MAX (PFULL TYPE)
  - \* MEAN (PFULL TYPE)
  - \* PEAK (PFULL TYPE)
  - \* MEDIAN (PFULL TYPE)
  - \* SUM (PFULL TYPE)
  - \* NZ aggregate function (select NZ() from ...)

```
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')
```

```
query
select
PHOT
cc=cc
#Hand
df=ea
for

plt.
plt.
plt.
```

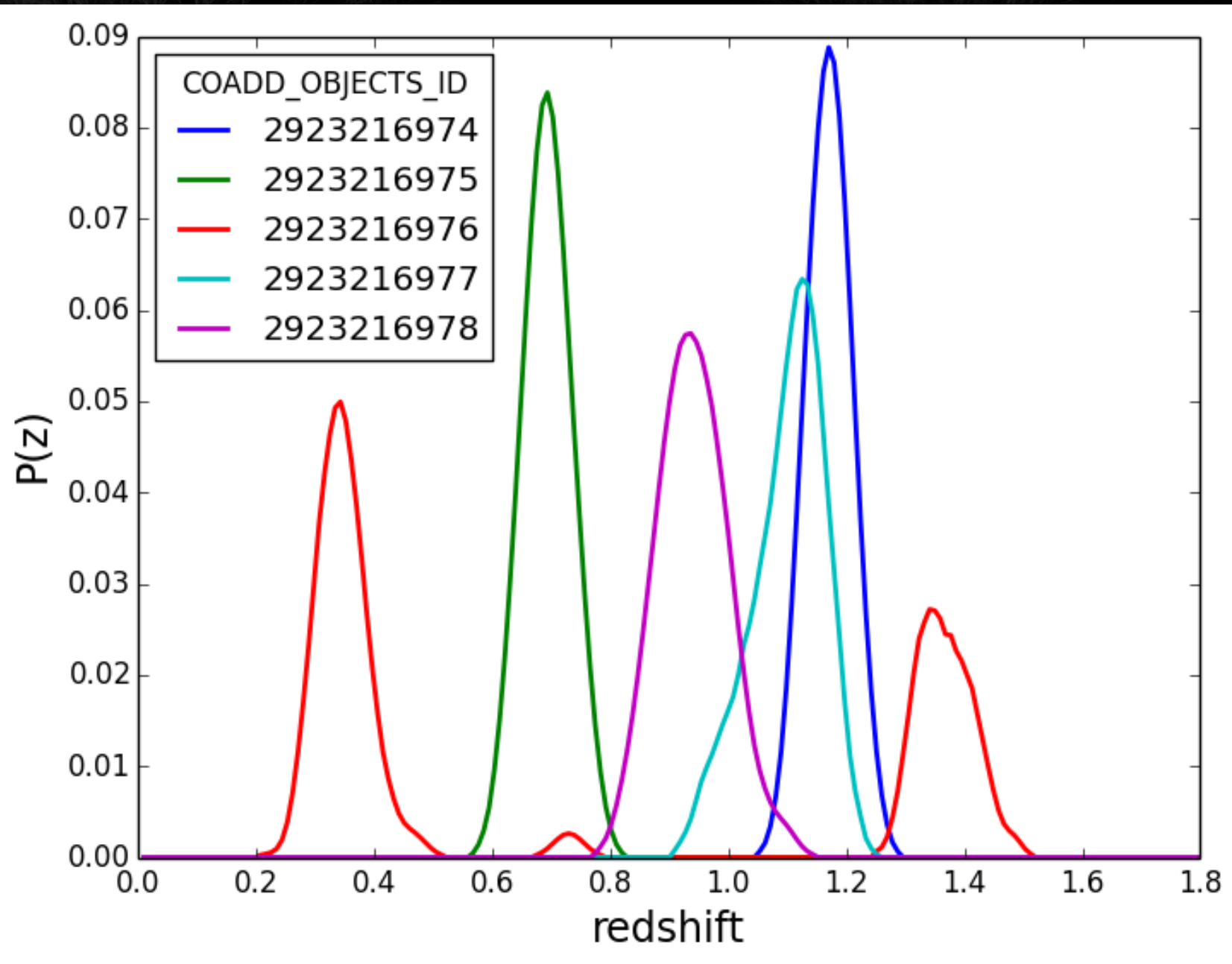




```
query="""
select COADD_OBJECTS_ID,PHZ.GET_PDF(TPZ) as
TPZ from PHOTOZ_SPARSE_SVA1_GOLD
where rownum < 6"""
cc=cursor.execute(query)
#Handling and plot
df=ea.to_pandas(cc)
for i in xrange(5):
    cid=df.COADD_OBJECTS_ID.values[i]
    plt.plot(zbins,df.TPZ.values[i],
             lw=2,label=cid)
plt.xlabel('redshift',fontsize=17)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

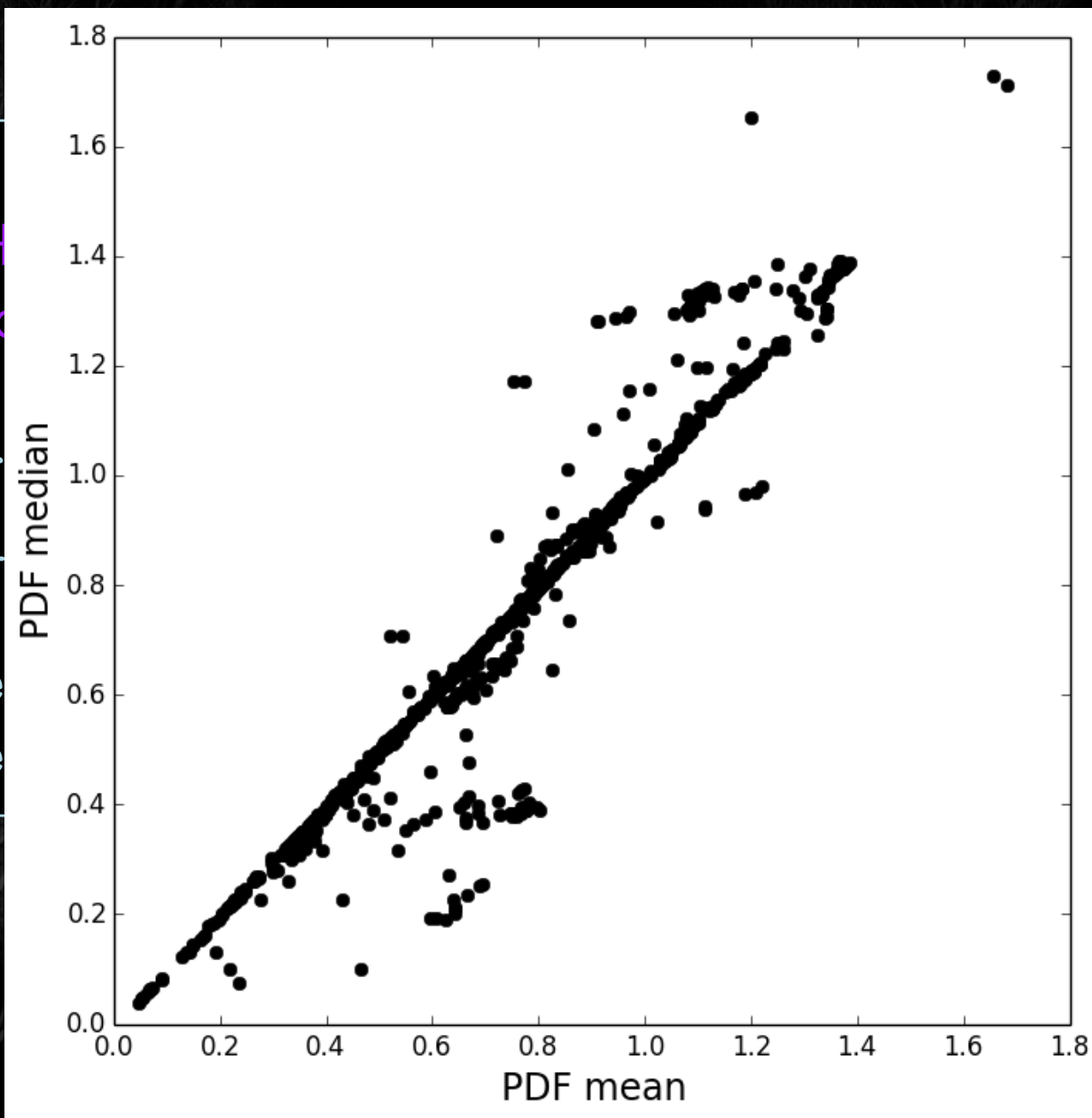
que  
 sele  
 TPZ  
 whe  
 cc=  
 #Ha  
 df=  
 for  
  
 plt  
 plt  
 plt





```
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)
```

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



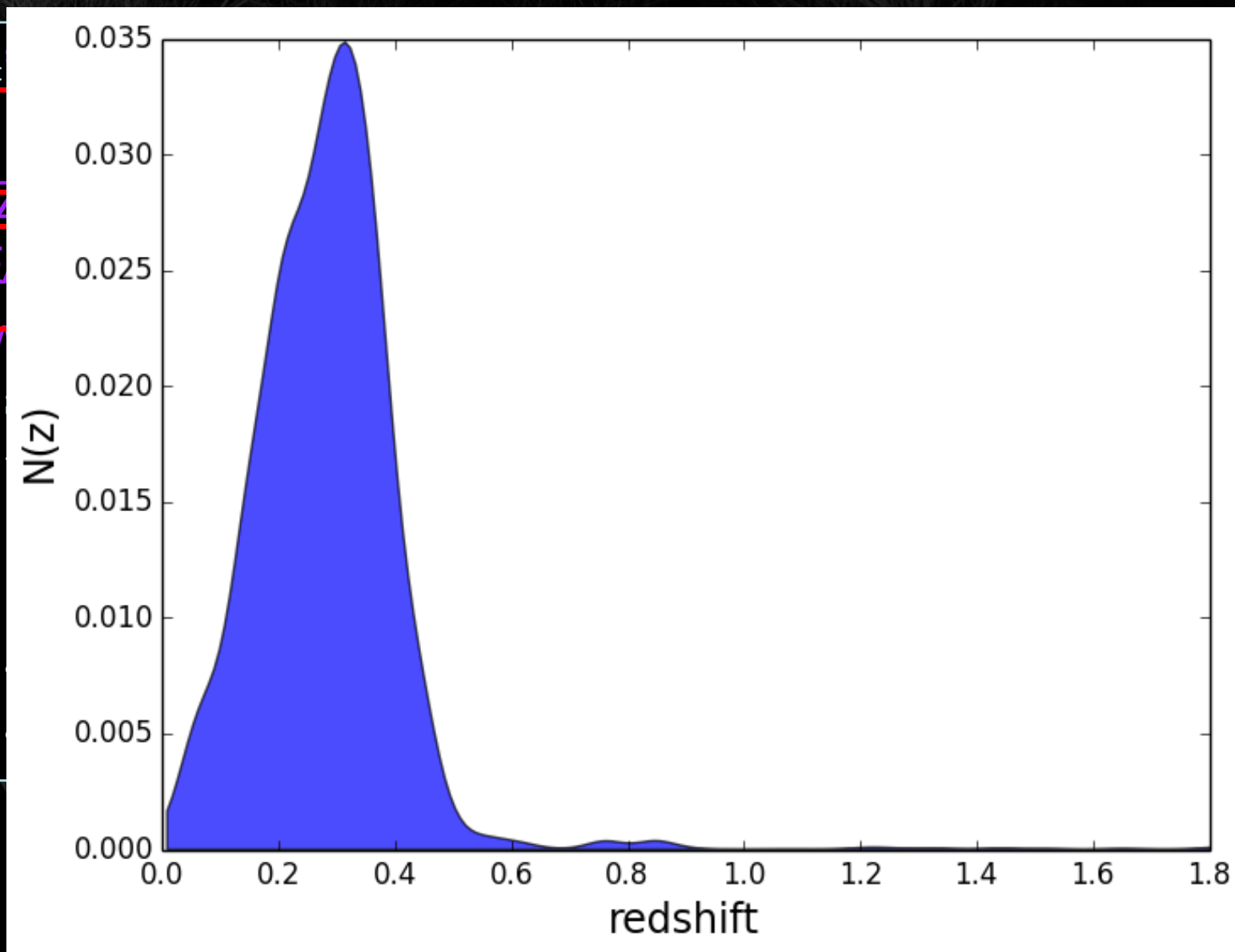
pz )

```
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)
```

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



- Photo-z tables in DB!
- Access to photo-z is easier and coordinated
- Use sparse representation for PDFs
- Bring analysis (software) to DB!
- [github.com/mgckind/SparsePz](https://github.com/mgckind/SparsePz)



# Questions?

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

