



Cosmology with imaging surveys – from **precision** to *accuracy*–

Interpretable photometric redshifts with no spectroscopy



NYU

Boris Leistedt – @ixkael, www.ixkael.com

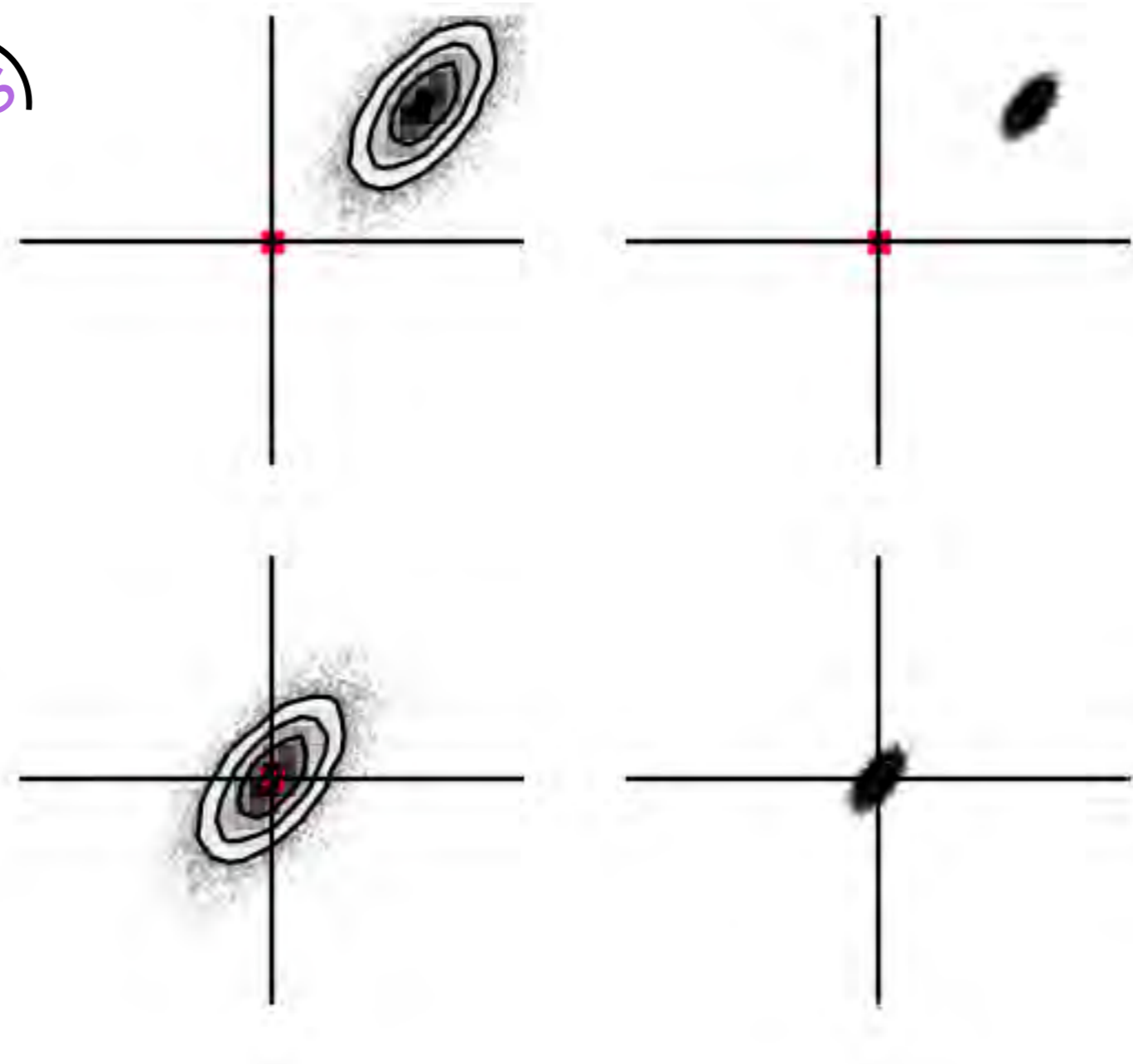
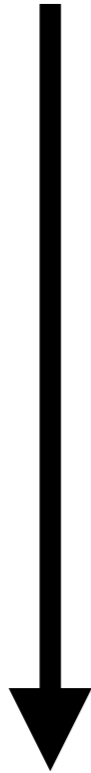
NASA Einstein Fellow @ CCPP, New York University



precision (*good data*)



accuracy
(*good methods*)



Frontiers: probabilistic models & observational systematics

ROAD MAP

Cosmology with galaxy surveys

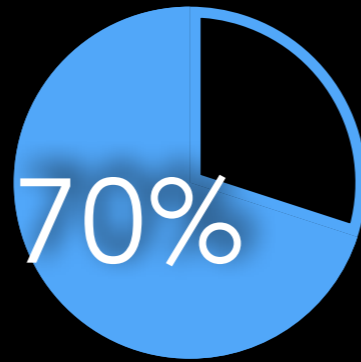
Photometric surveys and spatial systematics

Redshift distributions via hierarchical modeling

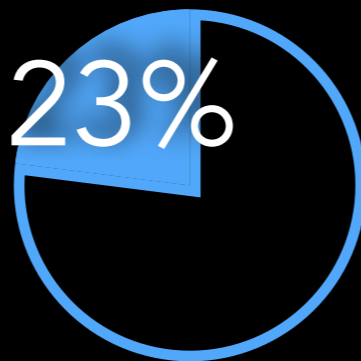
Photometric redshifts likelihood functions

Conclusions & the future (LSST)

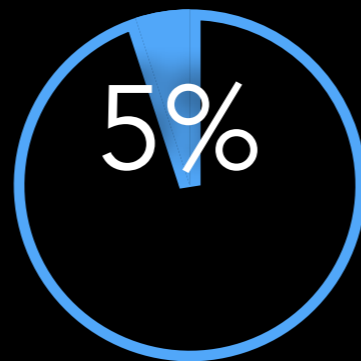
*Energy density
budget of the
universe (today)
and related
questions*



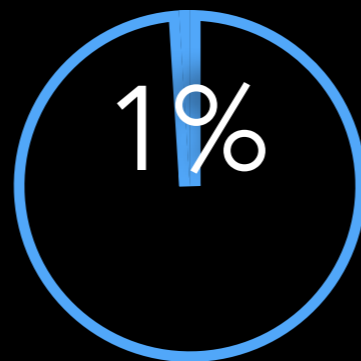
dark energy™
microscopic origin?
cosmological constant?



dark matter™
microscopic origin?
cosmological impact?

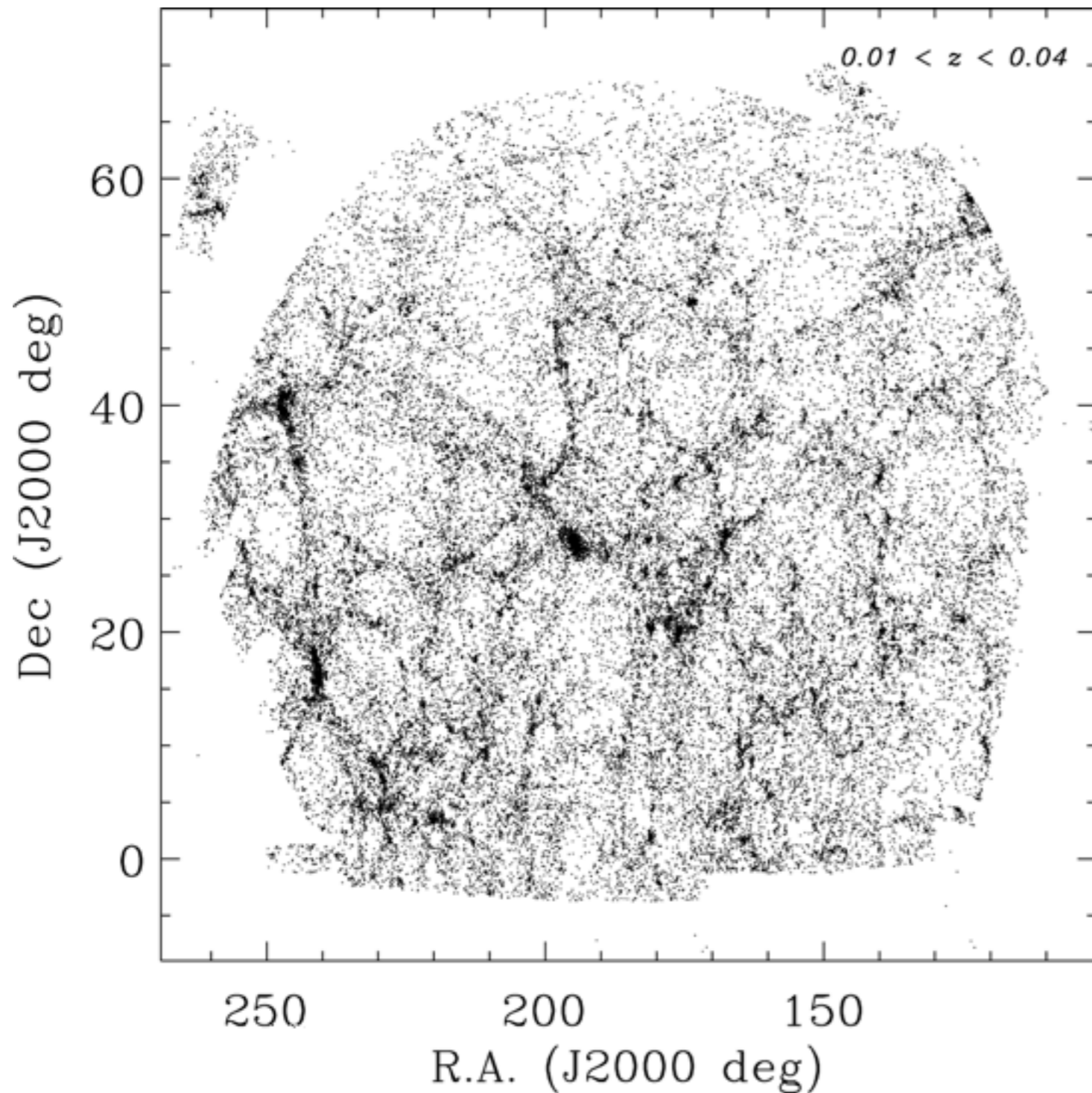


baryons
galaxy formation?
gastrophysics?



the rest
neutrinos, radiation, etc

Galaxy Surveys



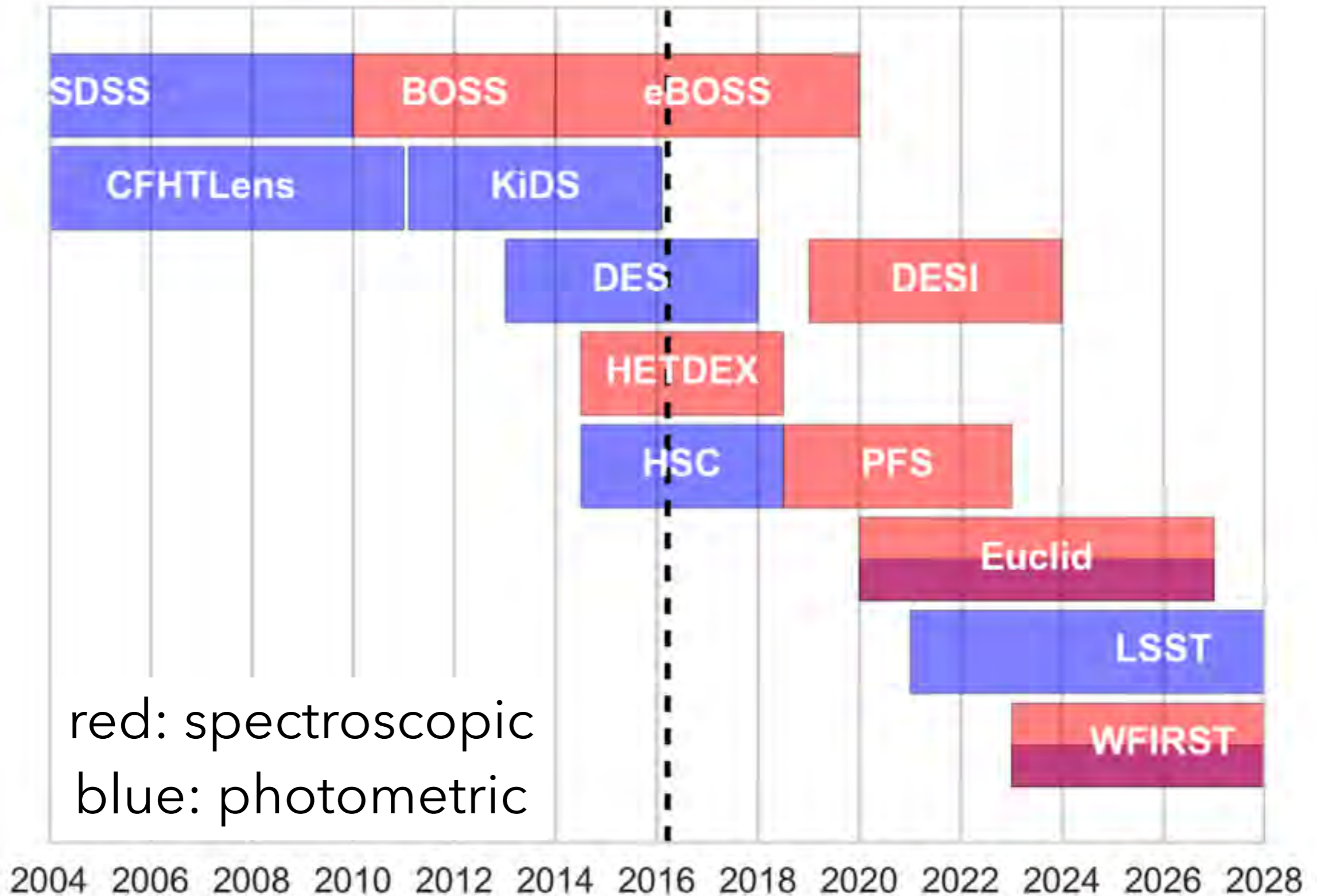
Rich space of models

early universe physics
astroparticles
gravity, etc

and observables

galaxy clustering
cosmic shear
galaxy-galaxy lensing
cross-correlations with
CMB, tSZ, etc

experimental landscape

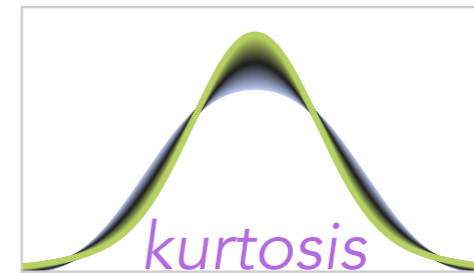
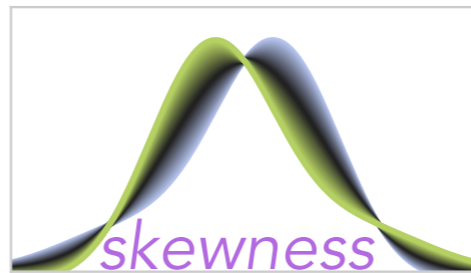
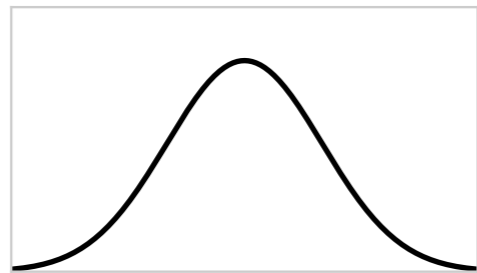


Galaxy surveys: goals

- ▶ Measure the expansion rate $H(z)$, growth rate f , and other cosmological parameters
- ▶ Constrain primordial non-Gaussianity (f_{NL} , g_{NL} , ...)
- ▶ Constrain primordial potential / power spectrum
- ▶ Constrain gravity from large-scale effects

Example: primordial non-Gaussianity

- ▶ Local type: $\Phi = \phi + f_{\text{NL}}[\phi^2 - \langle\phi^2\rangle] + g_{\text{NL}}[\phi^3 - 3\phi\langle\phi^2\rangle]$



- ▶ Imprinted in 3+ point correlations of the CMB and in large-scale galaxy power spectrum (2-point!)
- ▶ **Current limits:** $-3.1 < f_{\text{NL}} < 8.5$ from CMB (Planck collaboration 2015)
 $-16 < f_{\text{NL}} < 26$ from SDSS galaxies (Giannantonio+ 2014)
 $-39 < f_{\text{NL}} < 23$ from SDSS quasars (Leistedt+ 2014)
- ▶ Measurement ***plagued*** by spatial and redshift systematics
- ▶ LSST & SphereX could detect $f_{\text{NL}} \sim 1$ but will be difficult to reach!

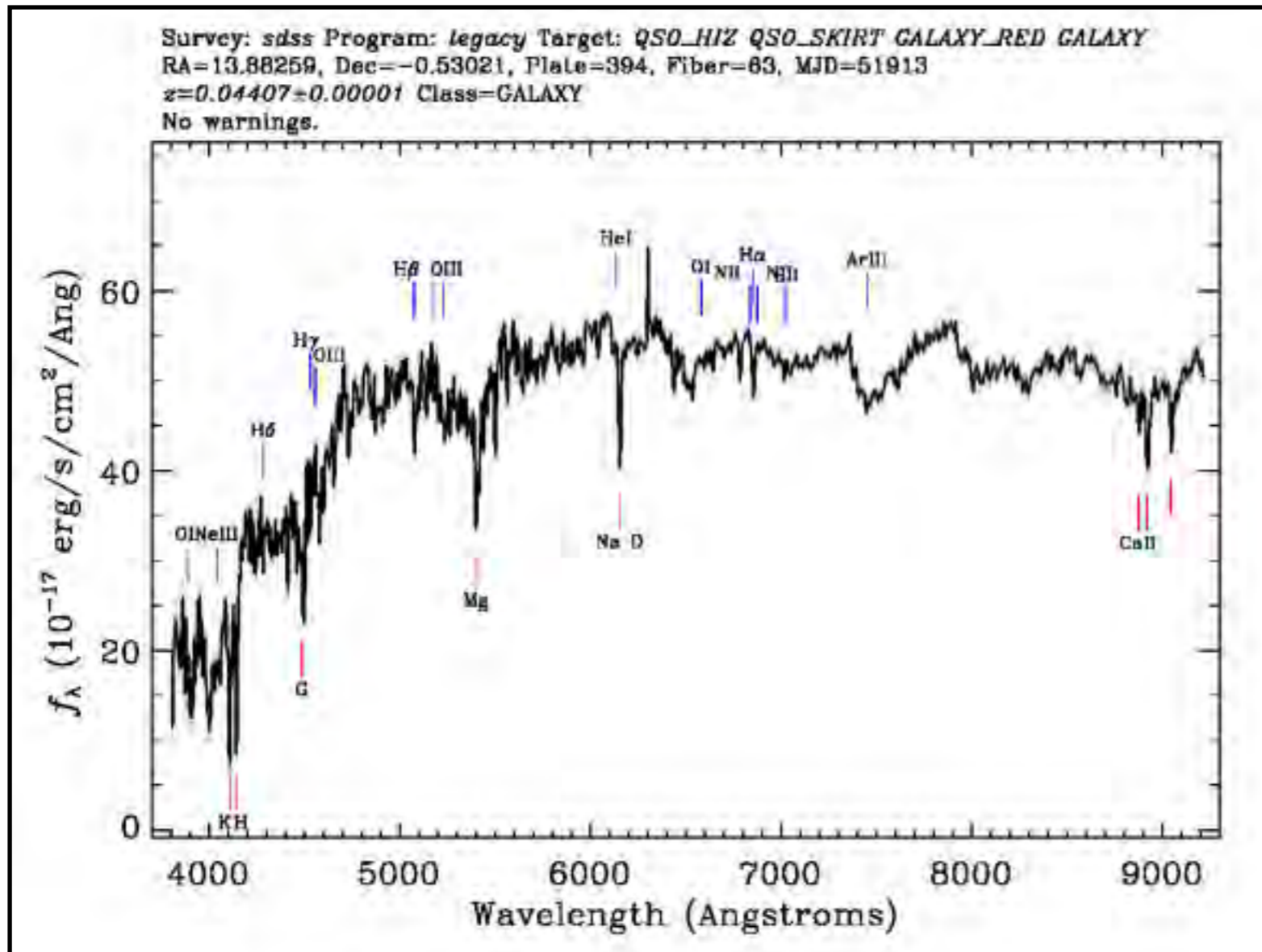
spectroscopic

VS

photometric

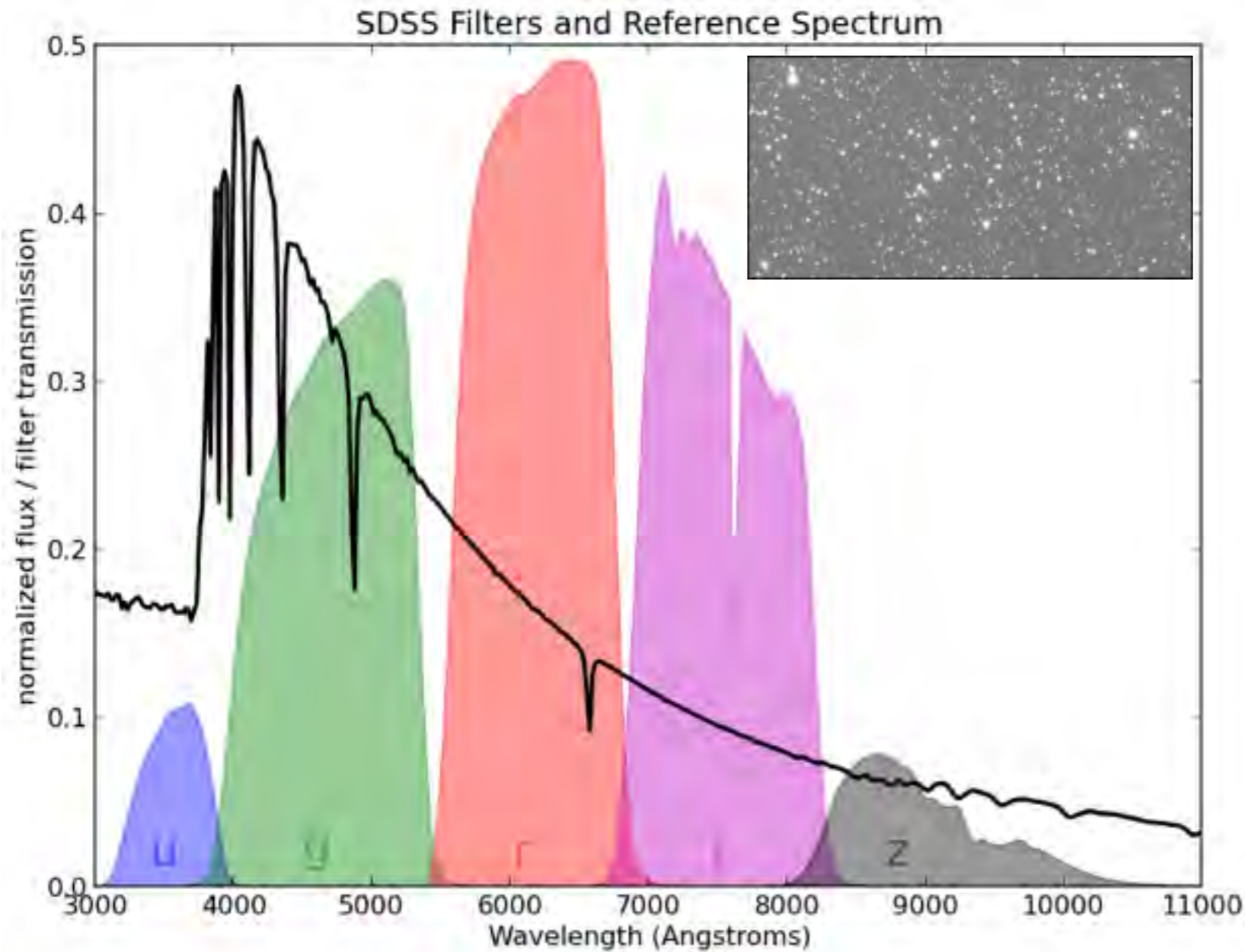
surveys

spectroscopic



- ✓ SEDs
- ✓ types
- ✓ redshifts
- ✗ shallow
- ✗ no shear

photometric



- ✓ CCD images
- ✓ deep
- ✓ shear
- ✗ no types
- ✗ no redshifts

spectroscopic

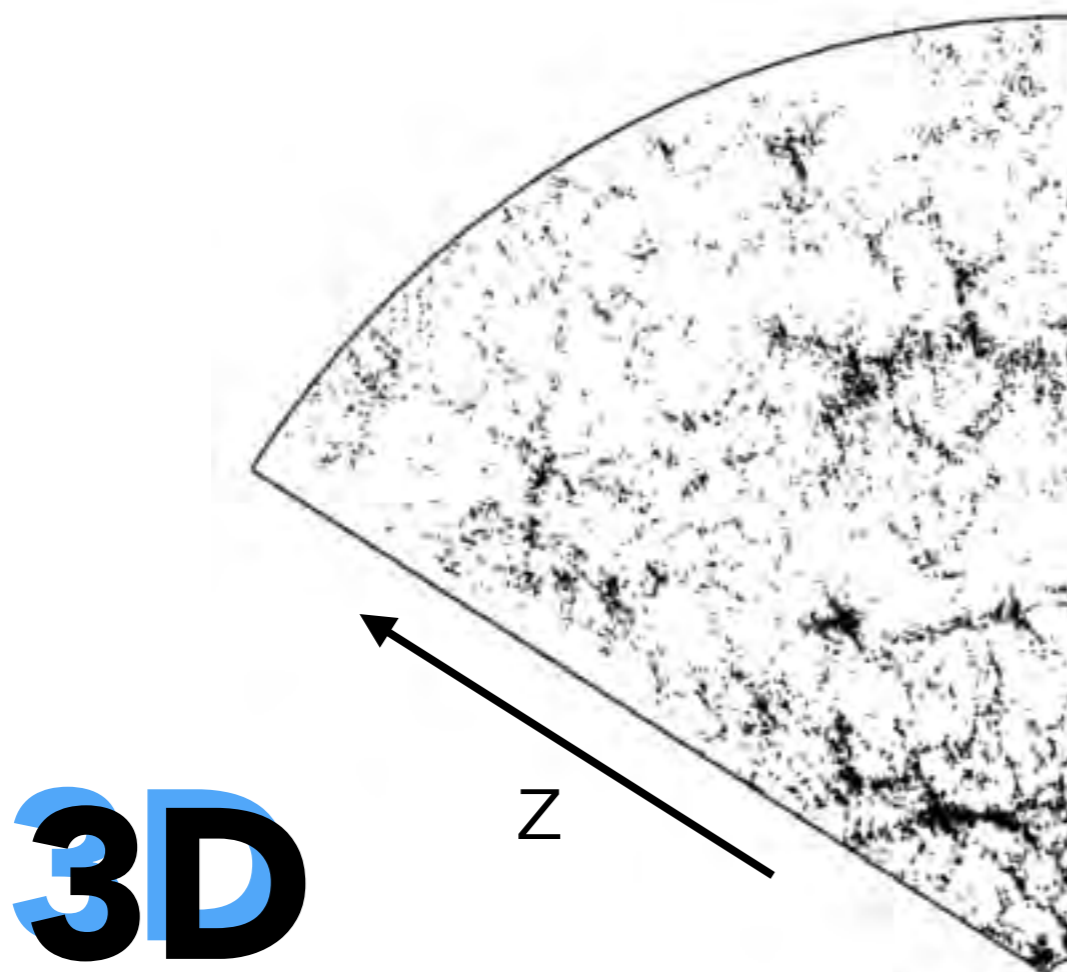
✓ types + redshifts

✗ shallow

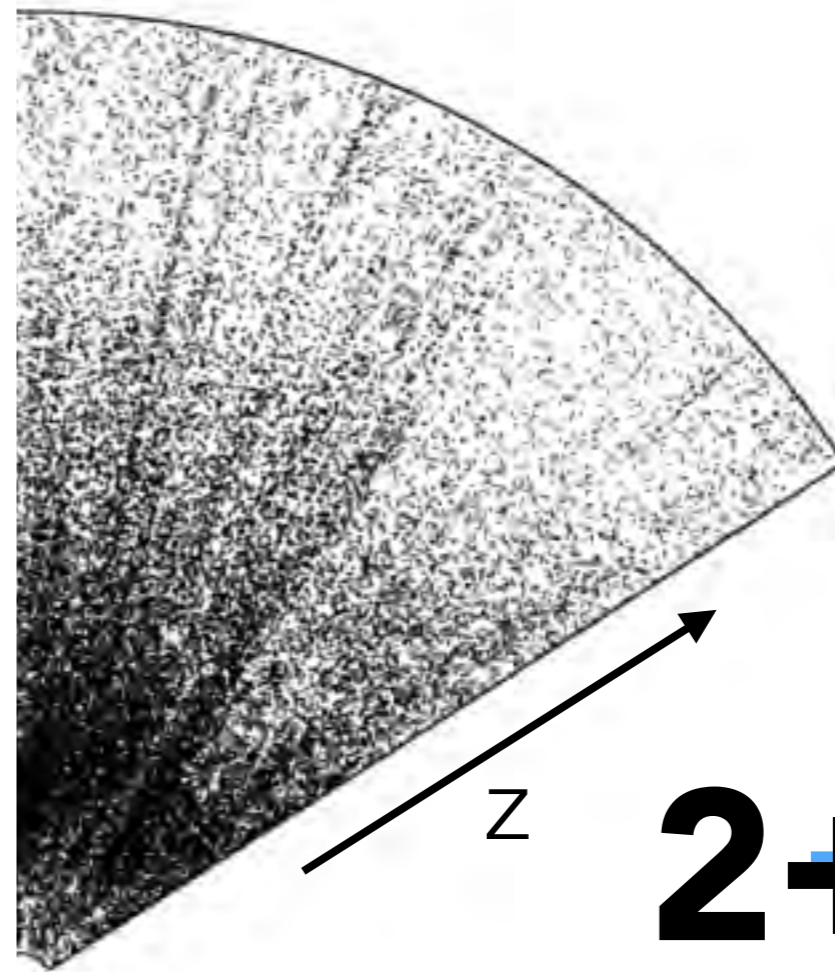
photometric

✗ no types / redshifts

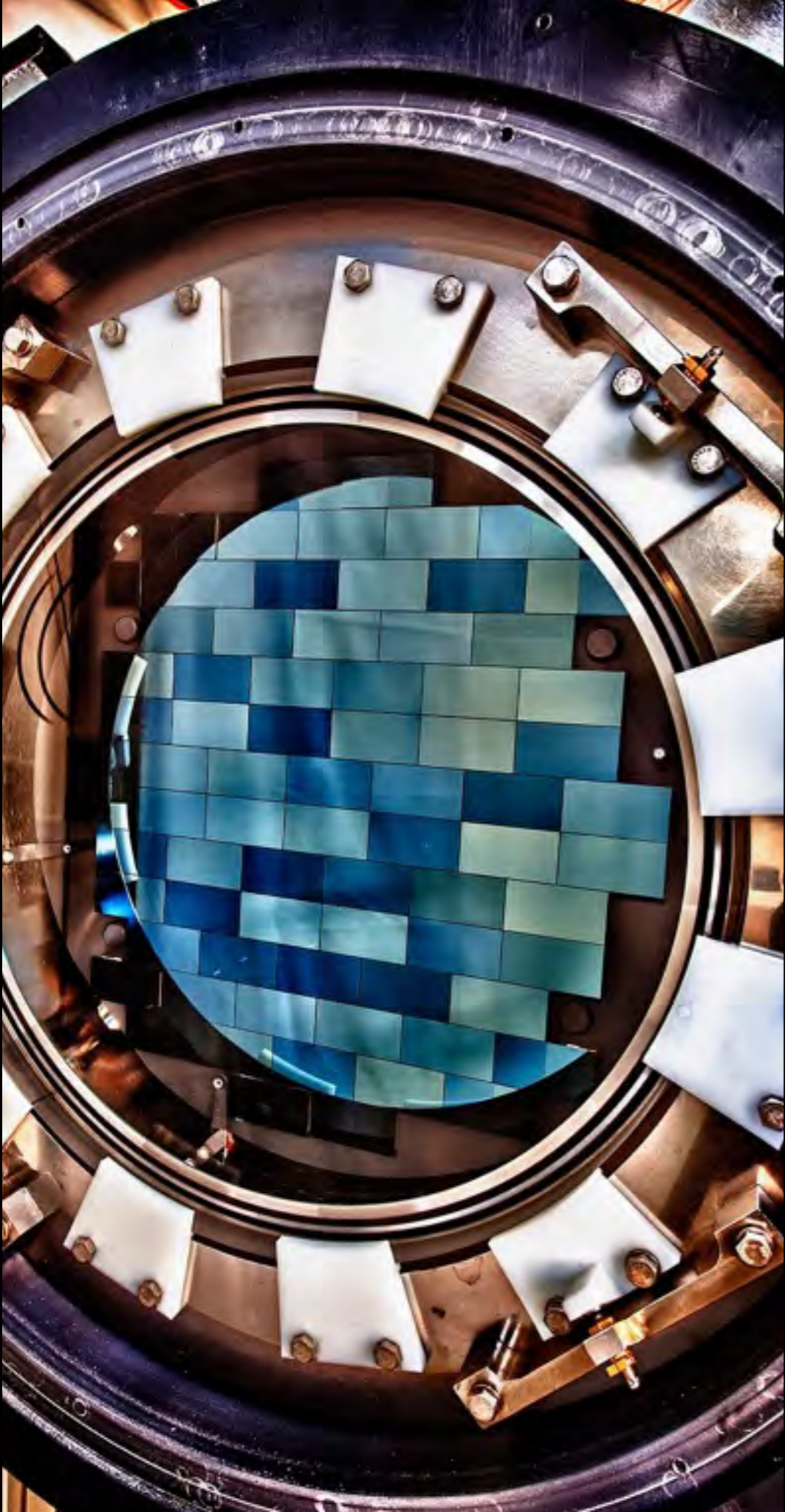
✓ deep



3D



2+1D



DES — The Dark Energy Survey

*3 deg² FOV, 570 Mpixel camera on
Blanco telescope (4m, CTIO)*

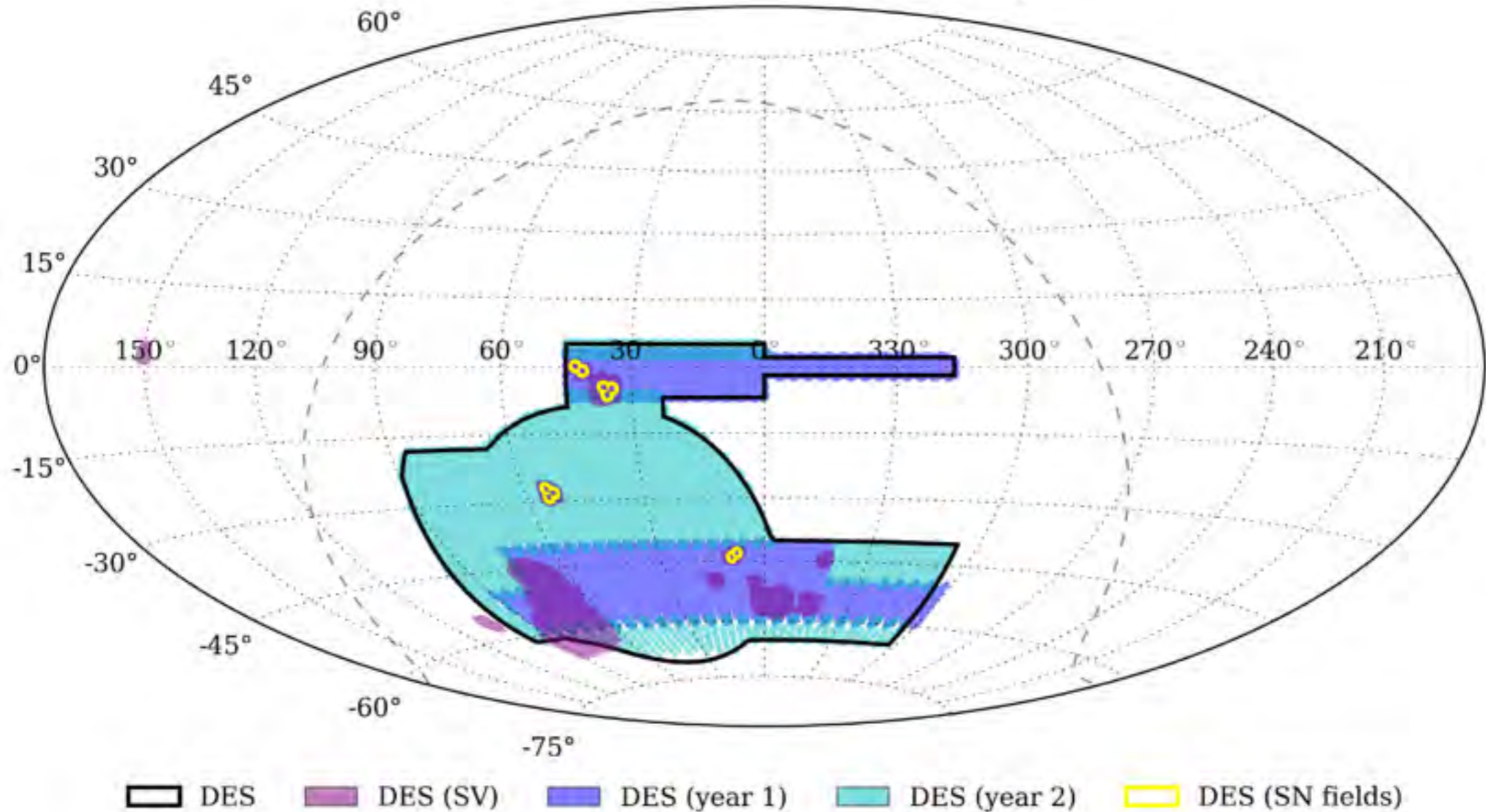
2013-2018, 300 collaborators

grizY bands

4 standard observational probes:

galaxy clustering, galaxy lensing,
supernovae, clusters

DES OBSERVING STRATEGY



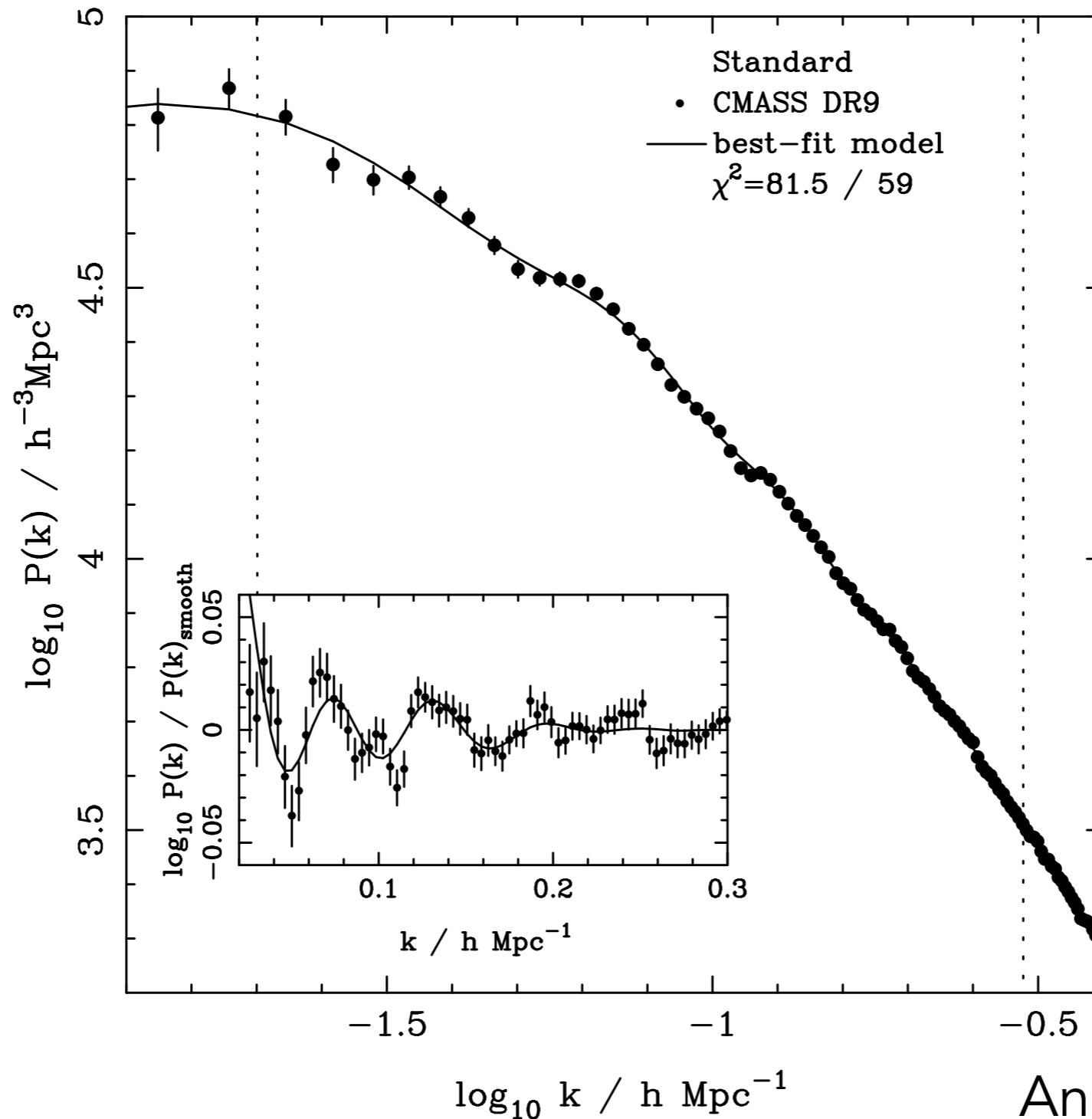
Full survey: 5000 deg² **grizY** to 24th mag

Science Verification: ~150 deg² to full depth

Year 1: ~1500 deg² but shallower

In the can (Y2-3): 5000 deg² deeper than Y1

3D matter power spectrum measured with spectroscopic survey



Anderson et al (2012)

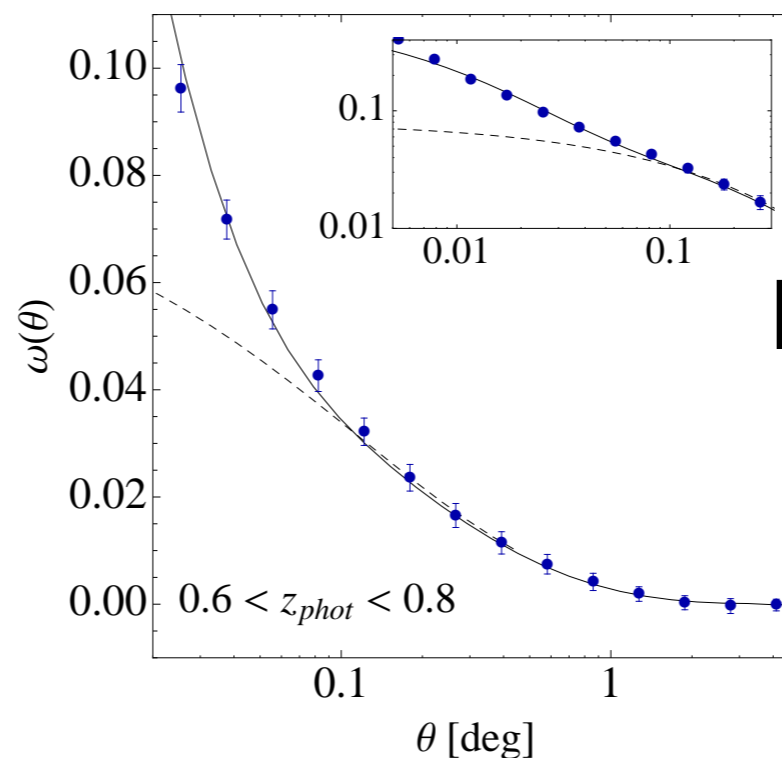
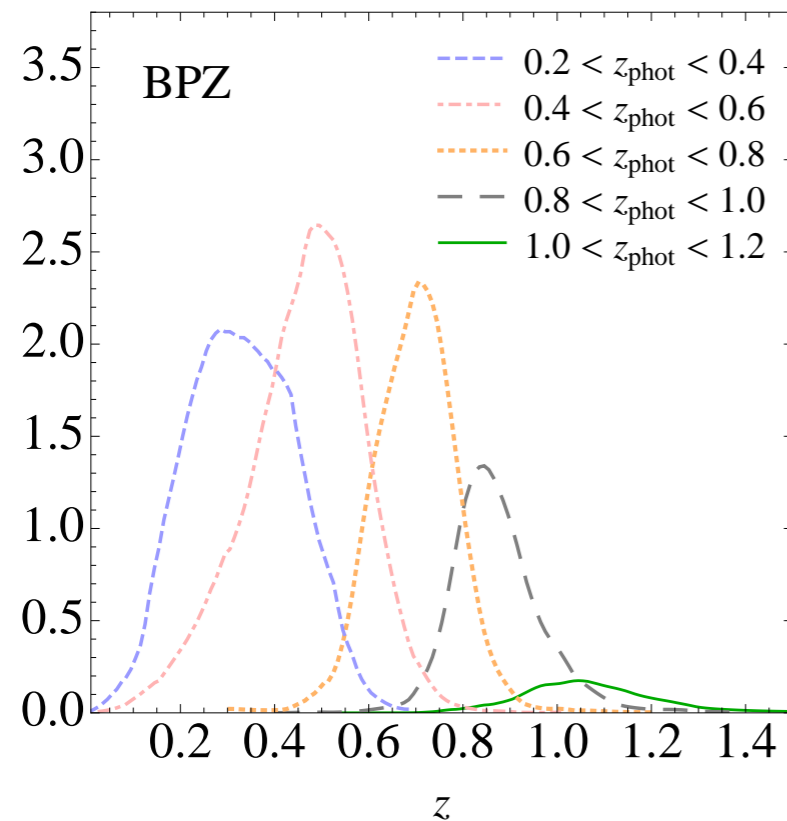
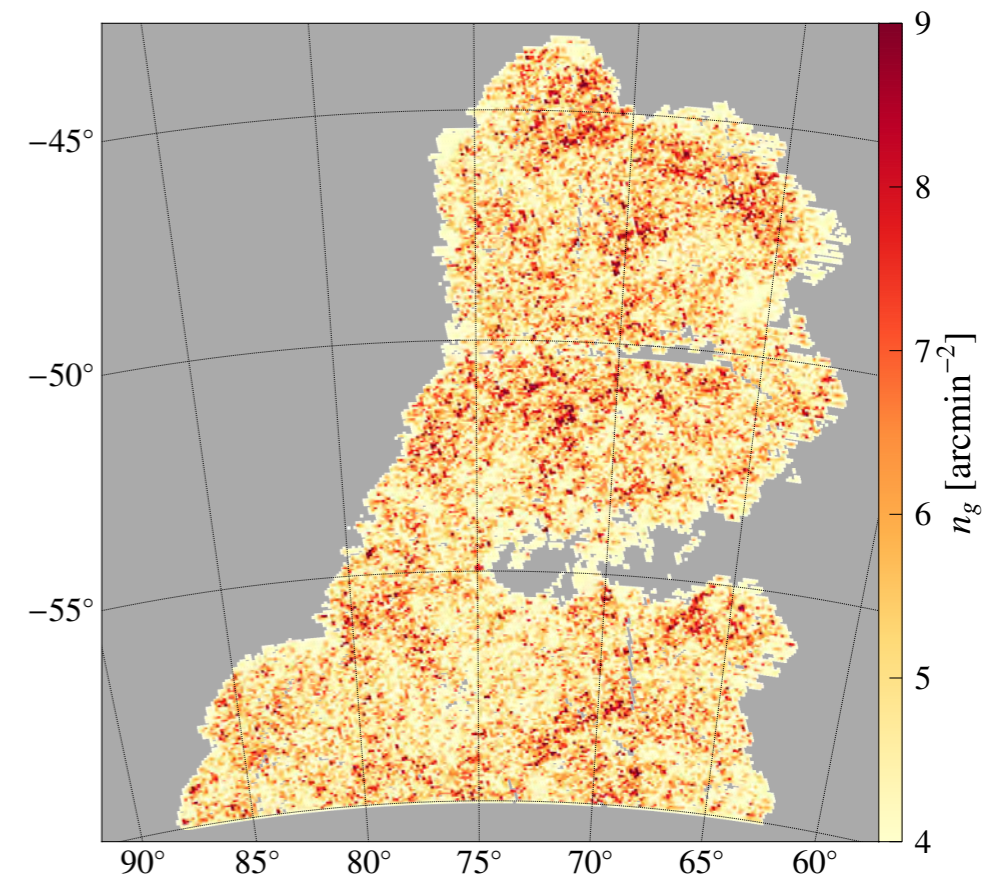
Photometric surveys

Redshifts are estimated from broad-band photometry and are **uncertain**

Typical approach: **tomographic analysis**
group galaxies into bins using a redshift estimate,
then do 2D angular 2-point correlation analysis

DES Science Verification (SV) data

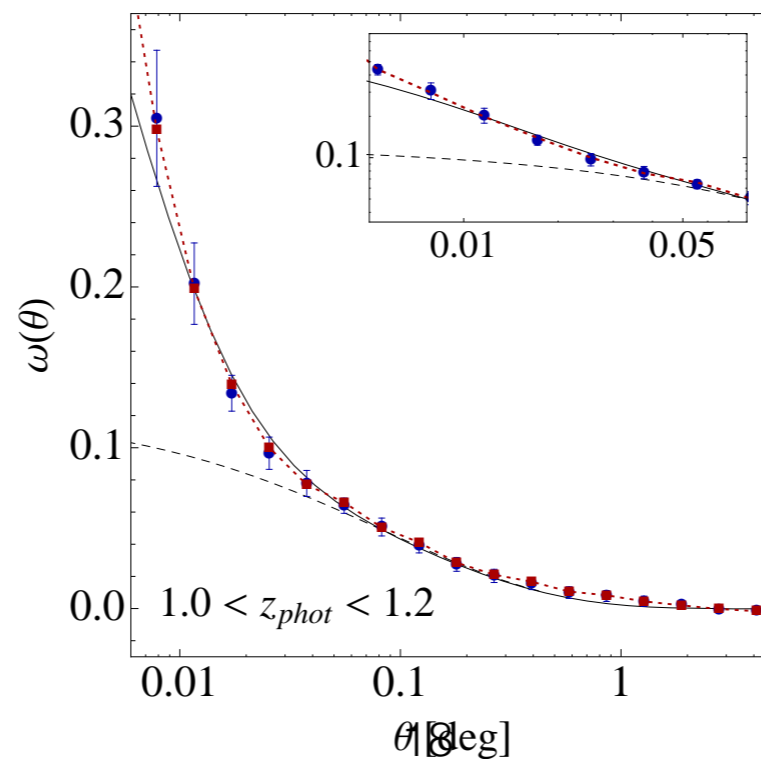
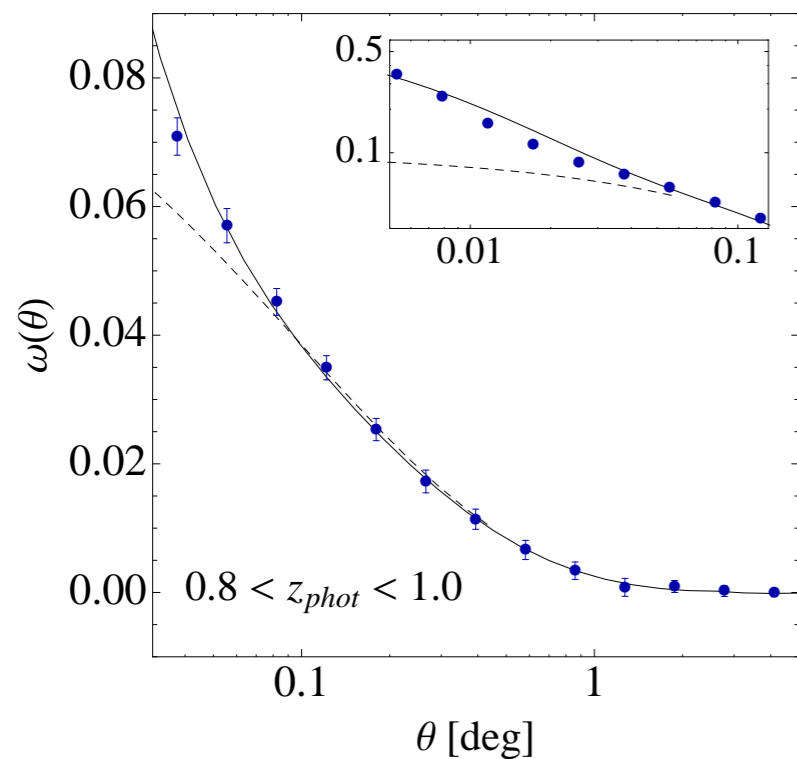
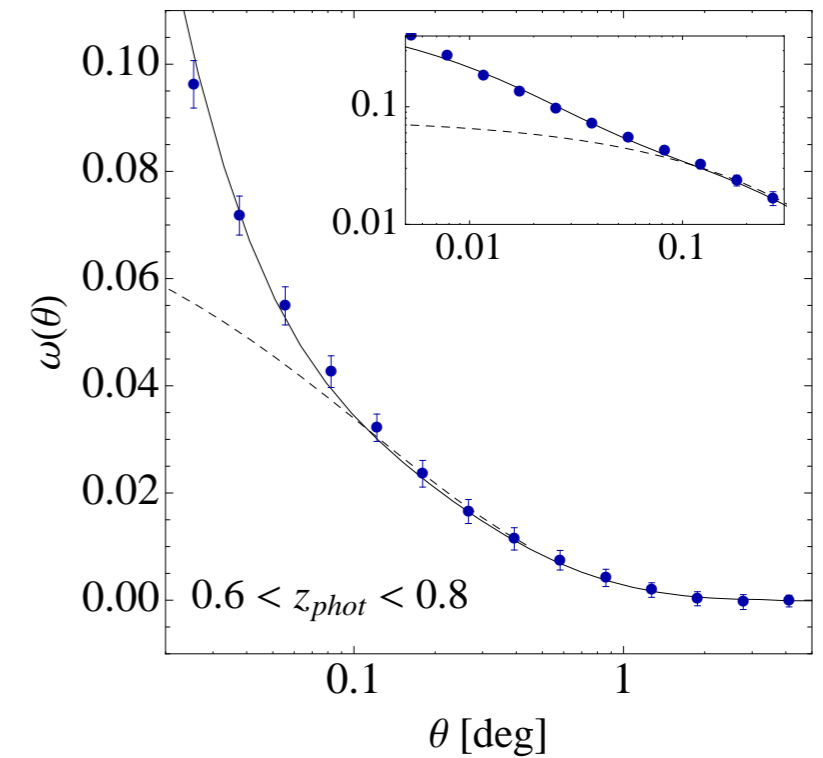
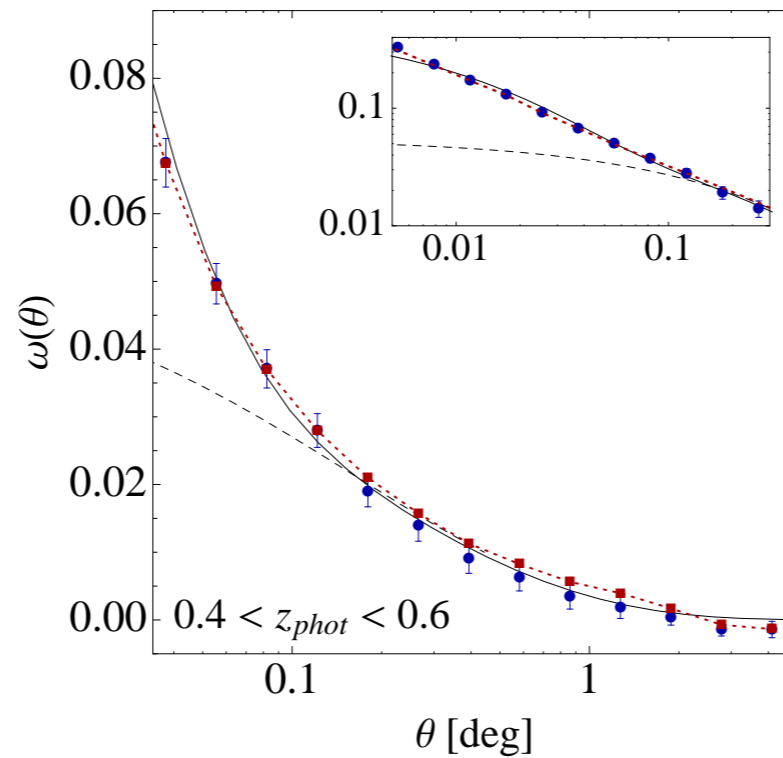
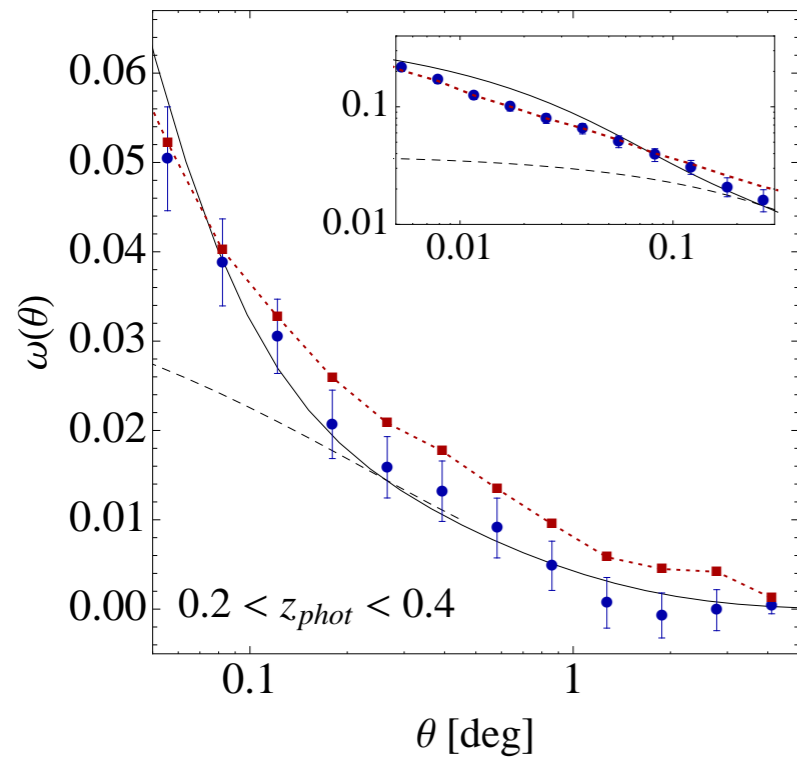
- ▶ Early data, excellent quality
- ▶ Benchmark galaxy sample:
 $\sim 120 \text{ deg}^2$, $i_{\text{mag}} < 22.5$,
cuts similar to CFHTLenS



BAO with Y1: 3-4%!

DES SV angular clustering

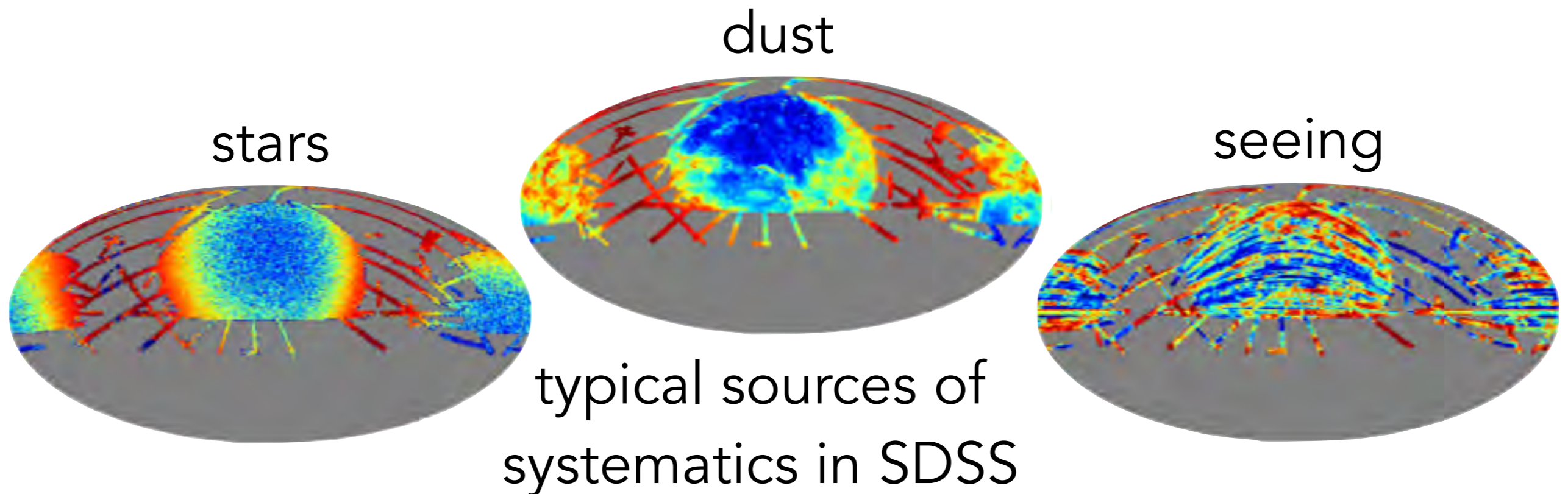
Crocce et al (MNRAS 2015, arXiv:1507.05360)



With corrections for spatial systematics

Spatial systematics?

- ▶ Anything that affects the measured galaxy properties
e.g. dust extinction, seeing, airmass, zero points, ...
- ▶ Create spatially varying depth & stellar contamination

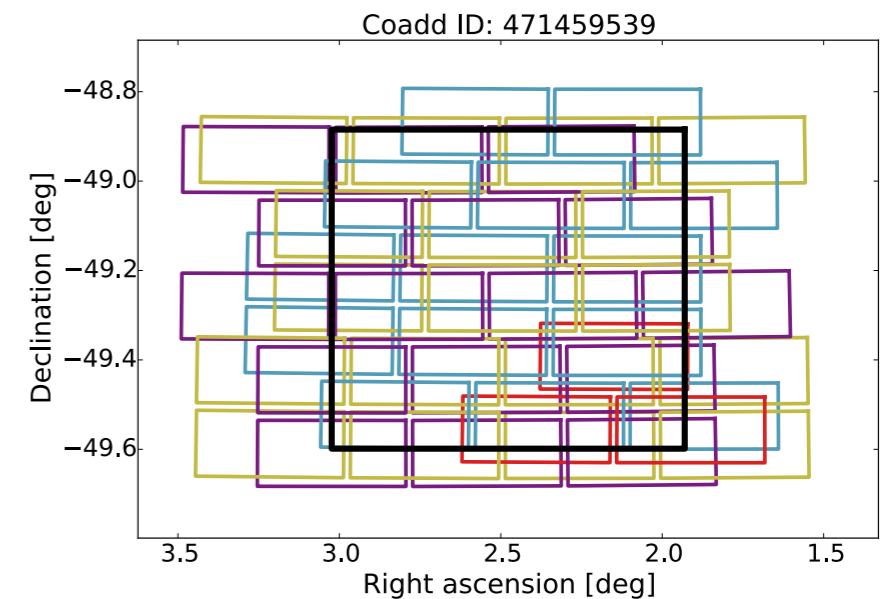


Observing conditions & systematics

Leistedt et al (ApJS 2015, arXiv:1507.05647)

Mapping & projecting image properties

Useful for null tests, systematics checks, ...

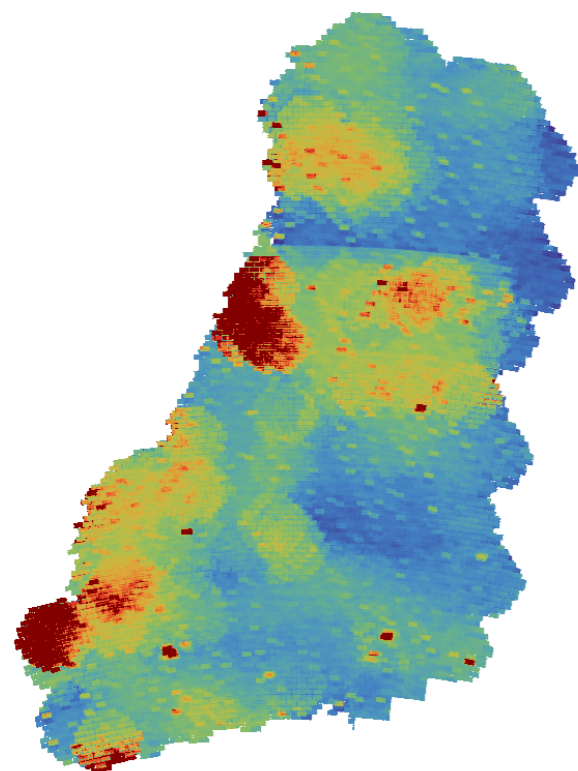


Mean sky sigma (i band)

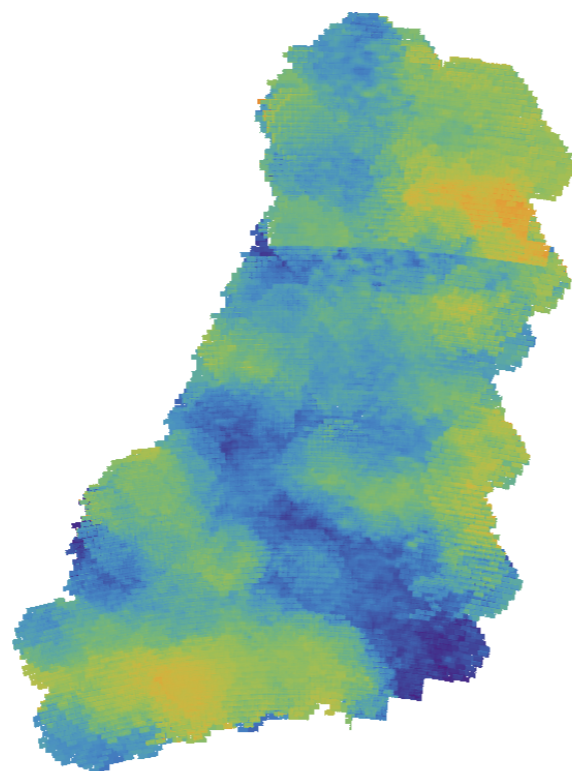
Mean seeing (i band)

Total exposure time (i band)

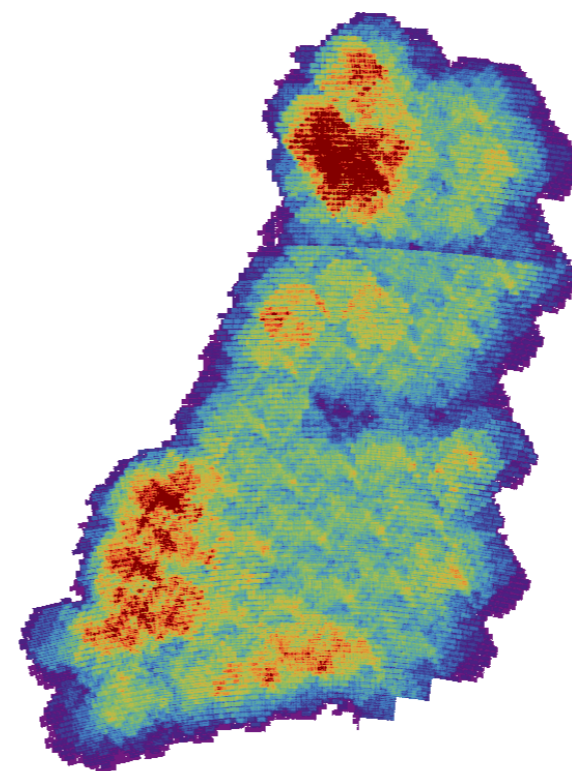
Total sky sigma (i band)



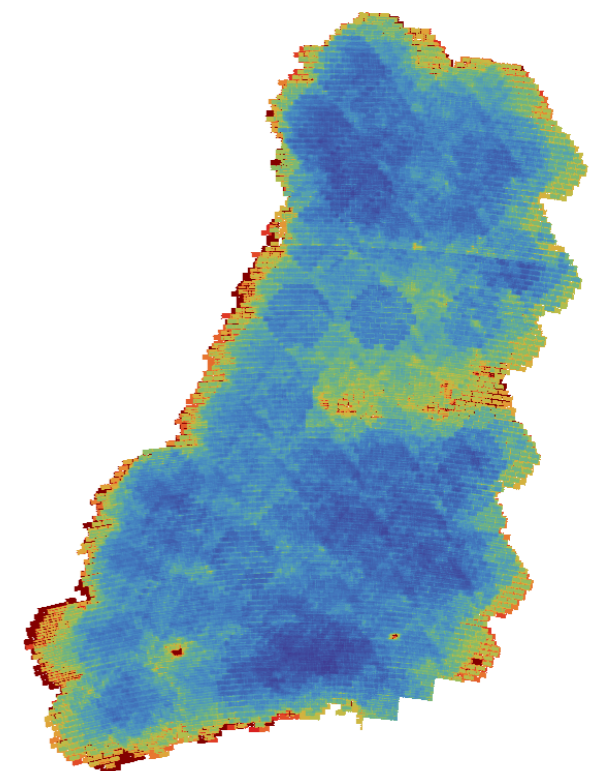
10.0 25.0 ADU
(75.5, -55.0) Equatorial



0.80 1.6 arcsec
(75.5, -55.0) Equatorial



90 1500 sec
(75.5, -55.0) Equatorial



1.0 15.0 ADU
(75.5, -55.0) Equatorial

Spatial systematics: state of the art

- ▶ Used to be main limiting systematic.
- ▶ **Confirmation bias** is dangerous (=fiddle with data and pipeline until results agree with expectations)
Blinding + meaningful statistical techniques essential.
- ▶ We now routinely **map** and **simulate** spatial systematics
Techniques to correct clustering measurements:
[Elsner, Leistedt & Peiris](#): arXiv:1609.03577, 1509.08933,
1507.05647, 1404.6530

photometric redshifts

(the elephant in the room)

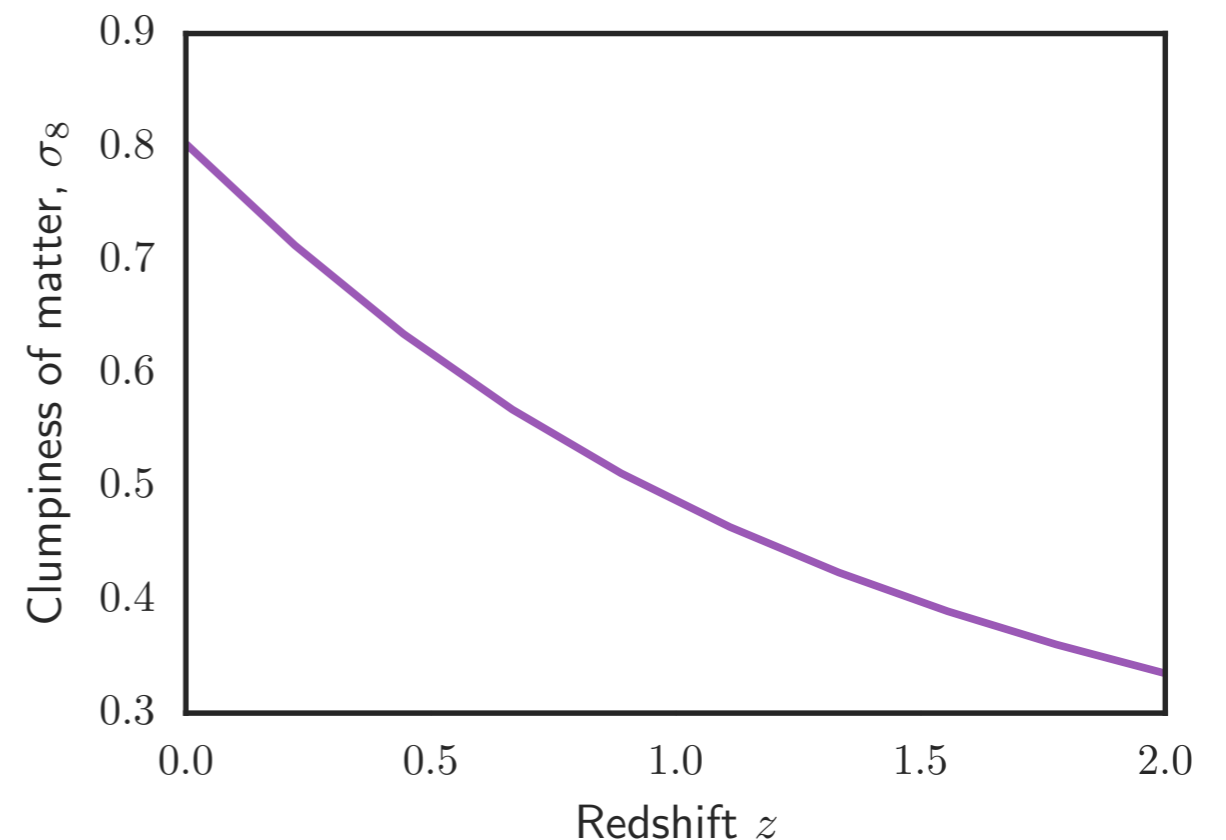
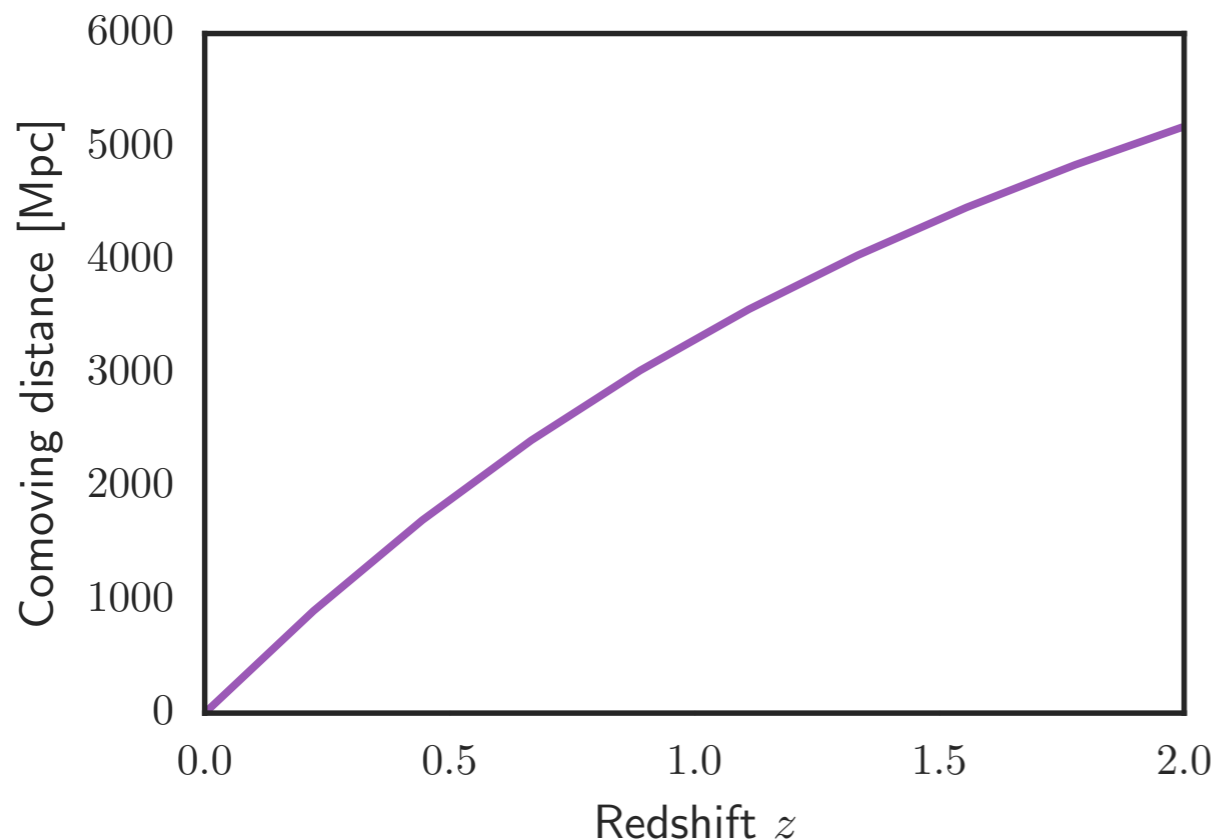
Redshift: doppler shift of electromagnetic radiation due to expansion of the universe = *indication of distance*

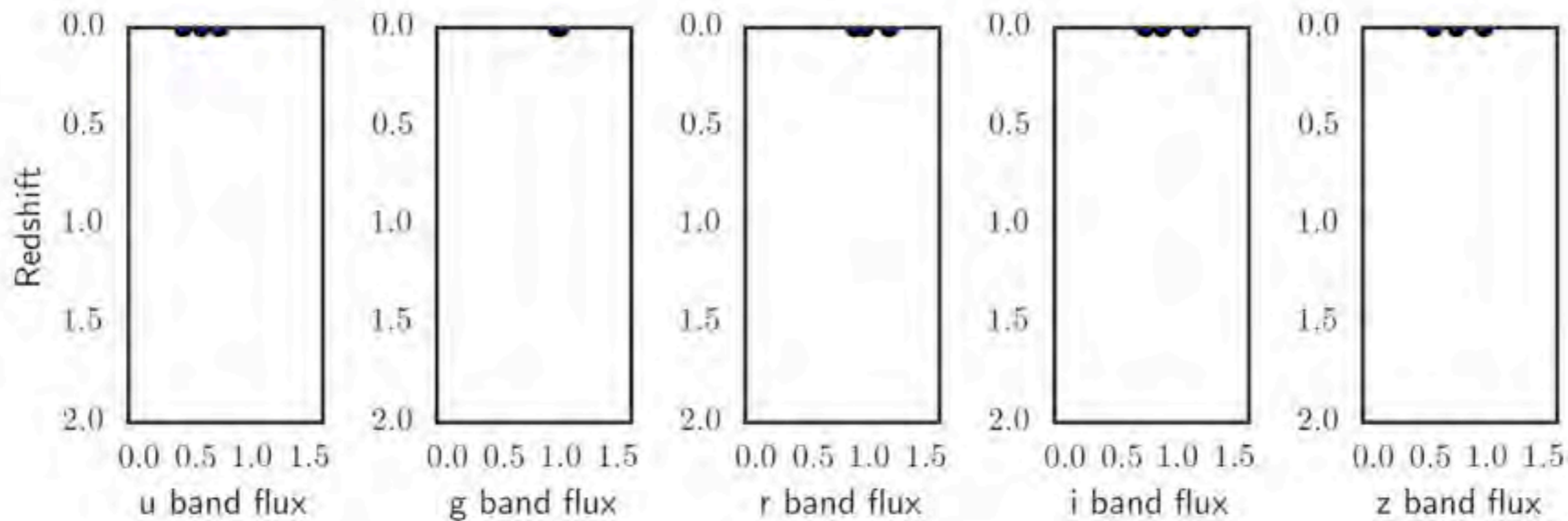
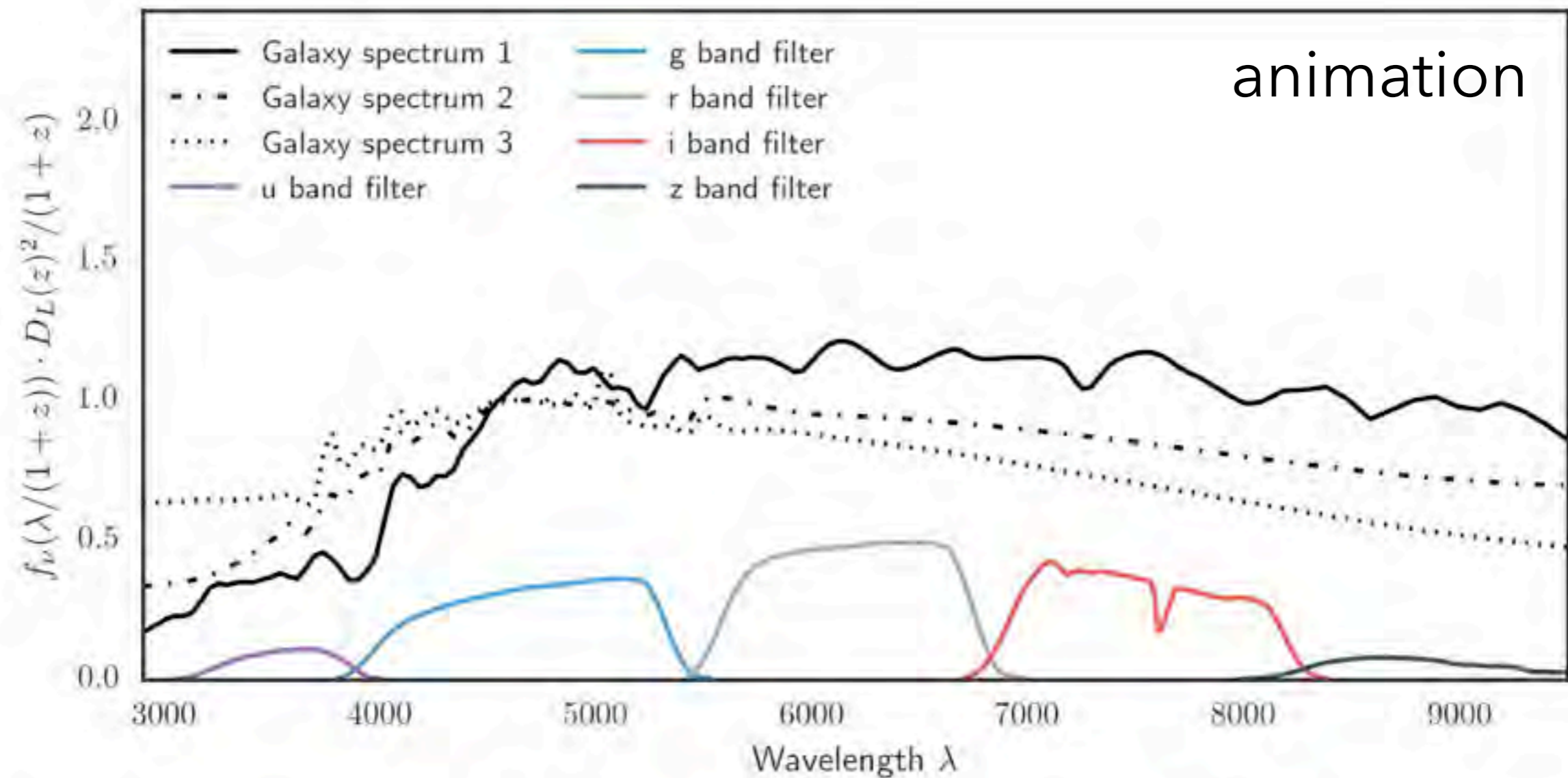
$$f_{\nu}(\lambda_{\text{obs}}, z) = \frac{(1+z)}{4\pi D_L^2(z)} L_{\nu} \left(\frac{\lambda_{\text{obs}}}{(1+z)} \right)$$

flux of redshifted object

intrinsic luminosity

Critical because redshift depends on cosmology

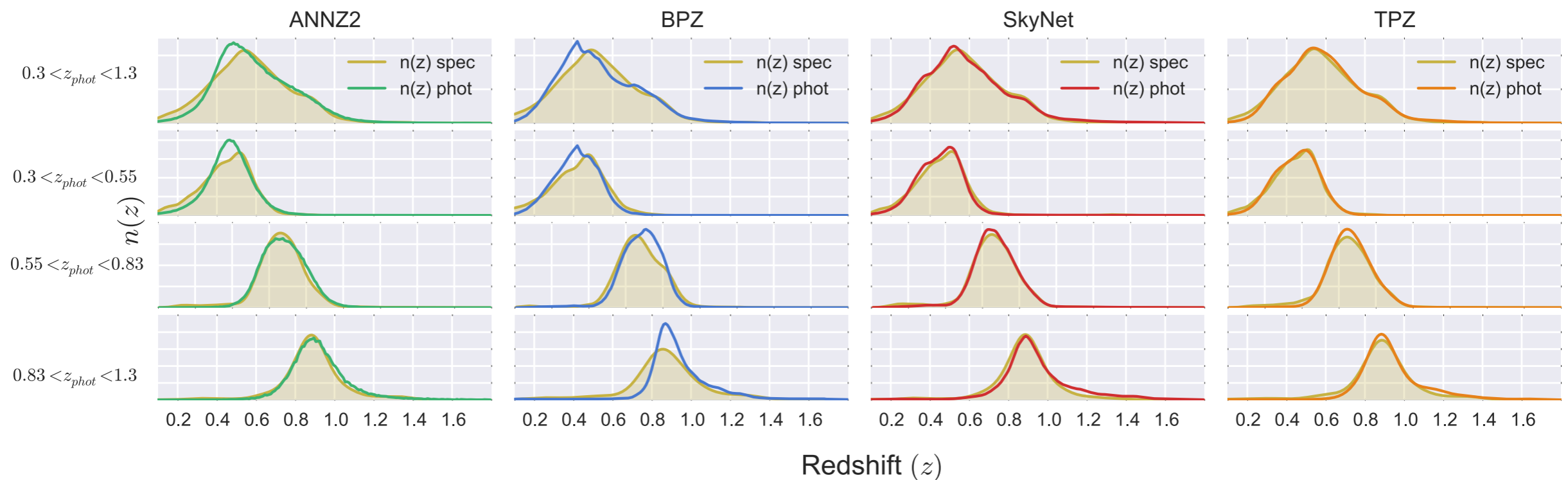
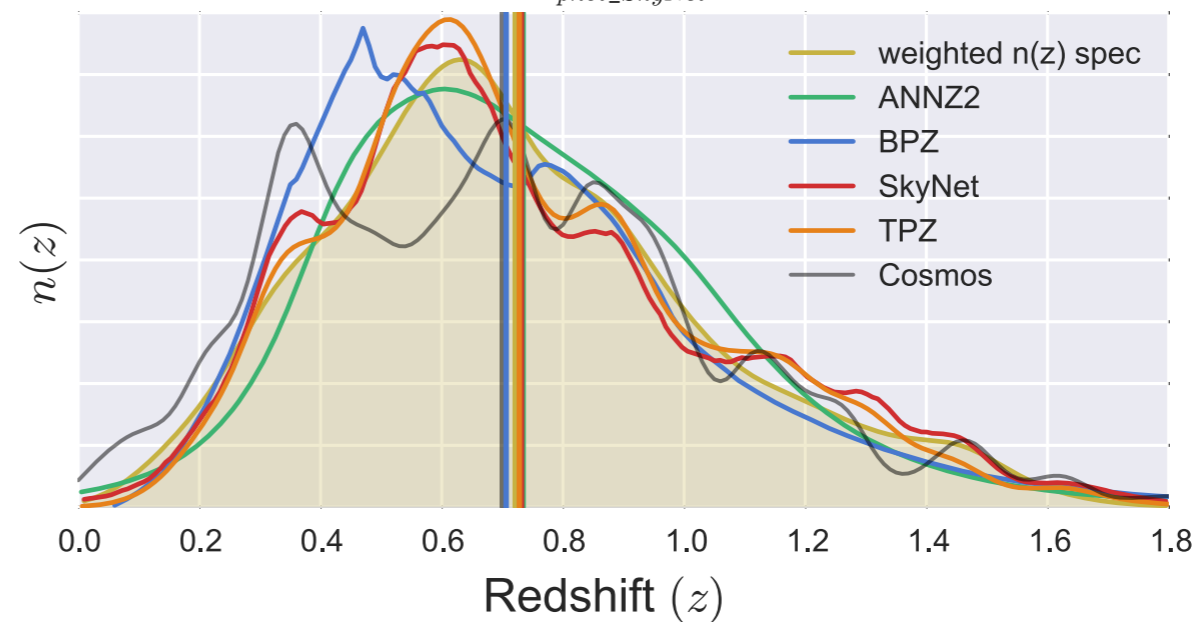




State of the art (DES)

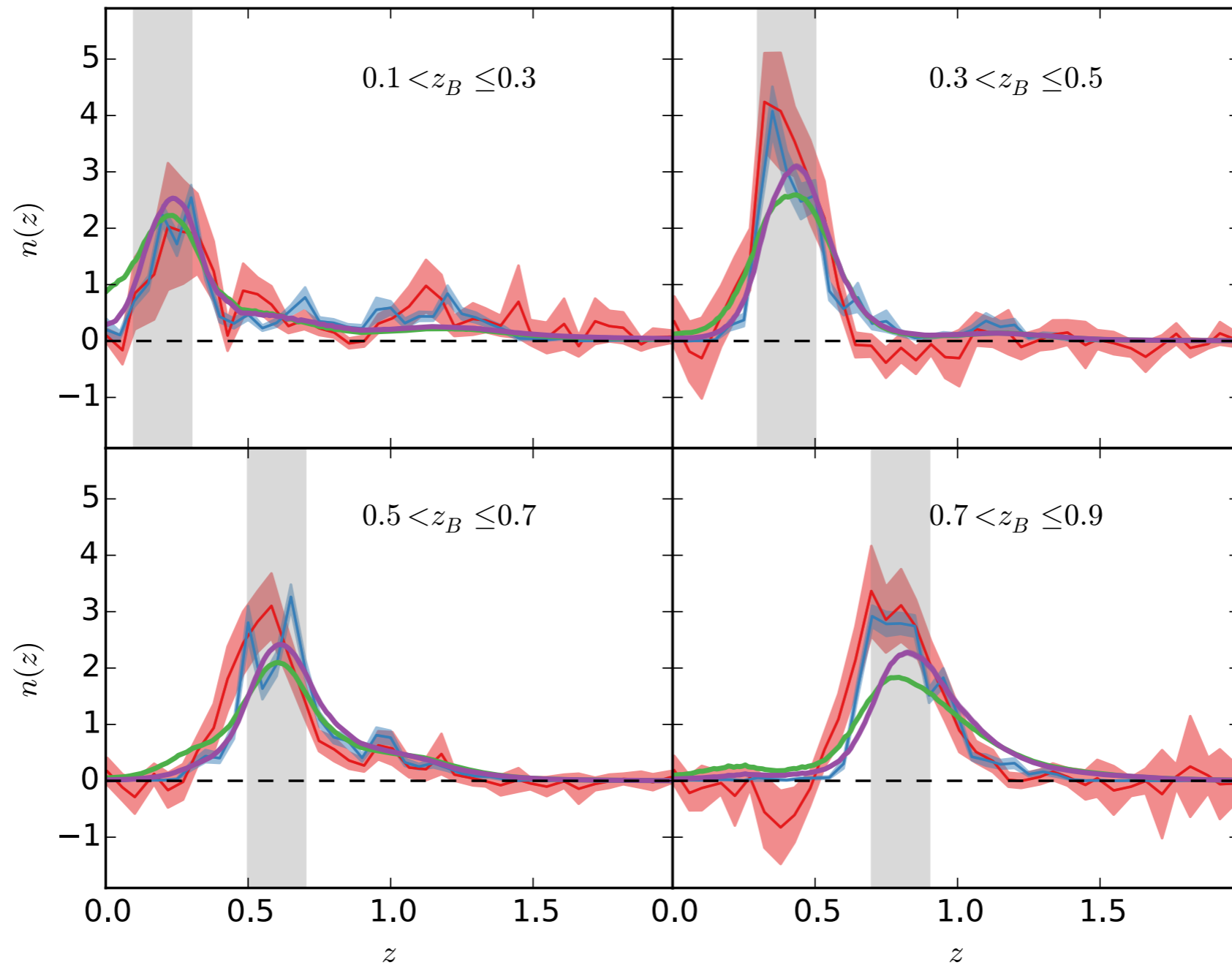
Weak Lensing Sample (NGMIX)

$$0.3 < z_{phot_SkyNet} < 1.3$$



Redshift distributions for DES SV galaxies (1507.05909)

State of the art (KIDS)



Redshift distributions for KIDS galaxies (1606.05338)

Ongoing surveys don't meet
photo-z requirements

LSST requires
insanely precise photo-z's

Why is it so hard?

Bayes Theorem

$$\underbrace{p(\mathbf{P}|\mathbf{D}, \mathbf{M})}_{\text{posterior}} = \underbrace{p(\mathbf{D}|\mathbf{P}, \mathbf{M})}_{\text{likelihood}} \times \underbrace{p(\mathbf{P}|\mathbf{M})}_{\text{prior}} / \underbrace{p(\mathbf{D}|\mathbf{M})}_{\text{evidence}}$$

Application to redshift distributions:

$$p(N(z), \{z_i\} | \{\text{Fluxes}_i\}) \propto p(N(z)) \prod_{i=1}^N p(z_i | N(z)) p(\text{Fluxes}_i | z_i)$$

full posterior prior population likelihood

Photo-z (likelihood) methods:

template fitting

vs

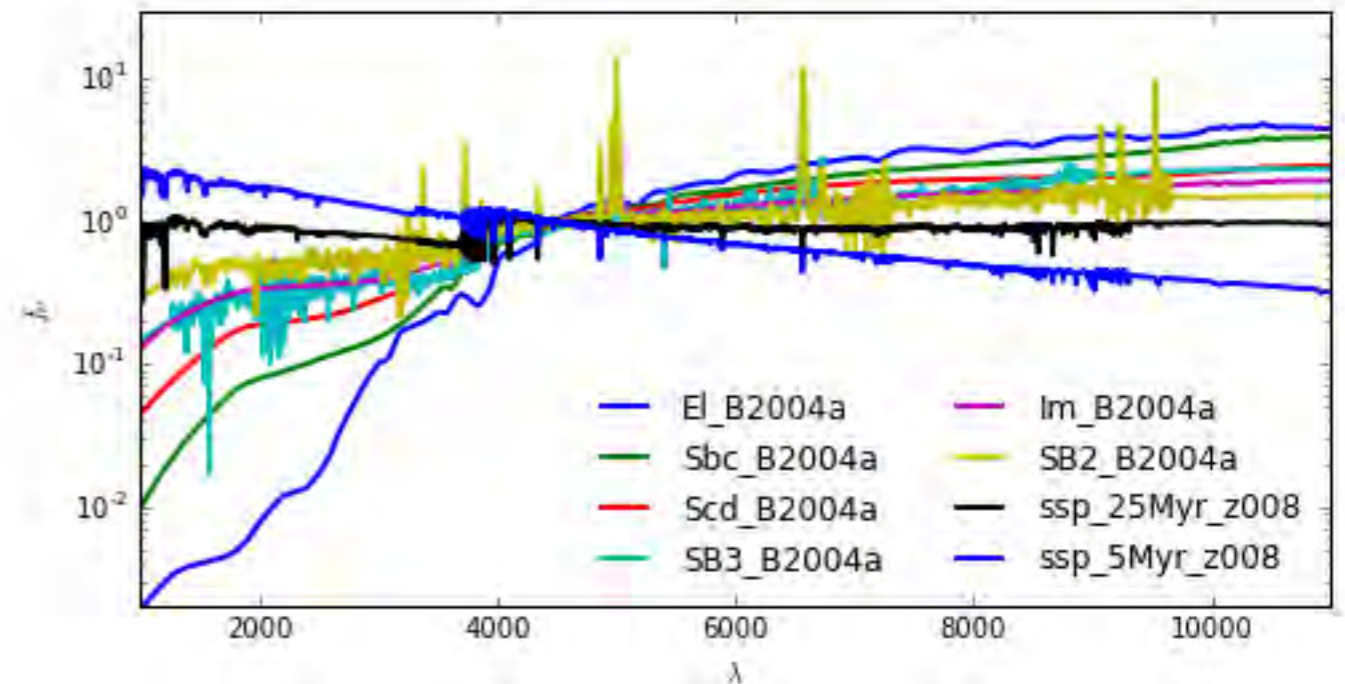
machine learning

(+new contestant: clustering redshifts)

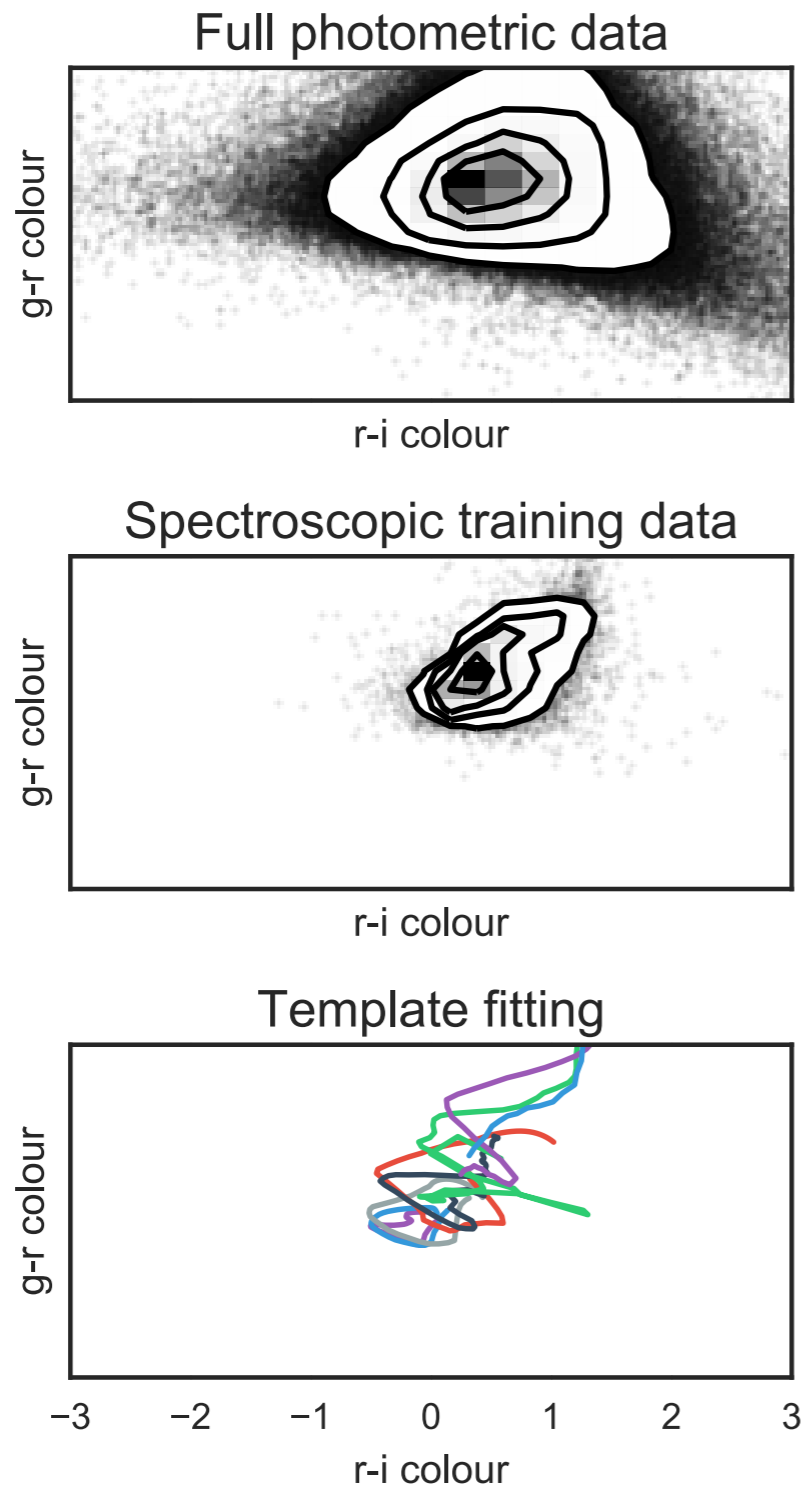
template fitting

- ✓ *physical model*
- ✓ *probabilistic*
- ✗ *need template set*
- ✗ *hard to capture data complexity*
- ✗ *sensitive to priors*

template set (CWW)



machine learning



✓ captures data complexity

✓ very flexible

✗ no physical model,
solves for flux \Rightarrow z,
cannot extrapolate

✗ not probabilistic

✗ requires representative
training data

Why is it so hard?

$$p(N(z), \{z_i\} | \{\text{Fluxes}_i\}) \propto p(N(z)) \prod_{i=1}^N p(z_i | N(z)) p(\text{Fluxes}_i | z_i)$$

full posterior prior population likelihood

- ▶ Galaxy SED models are inaccurate (high redshift, dust, star formation, variability, etc) ⇒ **likelihood is unreliable**
- ▶ Standard analyses stack redshift PDFs to obtain $N(z)$.
⇒ **$N(z)$ is biased and has no uncertainties.**
- ▶ My goals:
Create & calibrate SED models and likelihood function
Correctly infer $N(z)$ and propagate errors into cosmology

Hierarchical inference of redshift distributions

arXiv:1602.05960
with



Daniel
Mortlock



Hiranya
Peiris

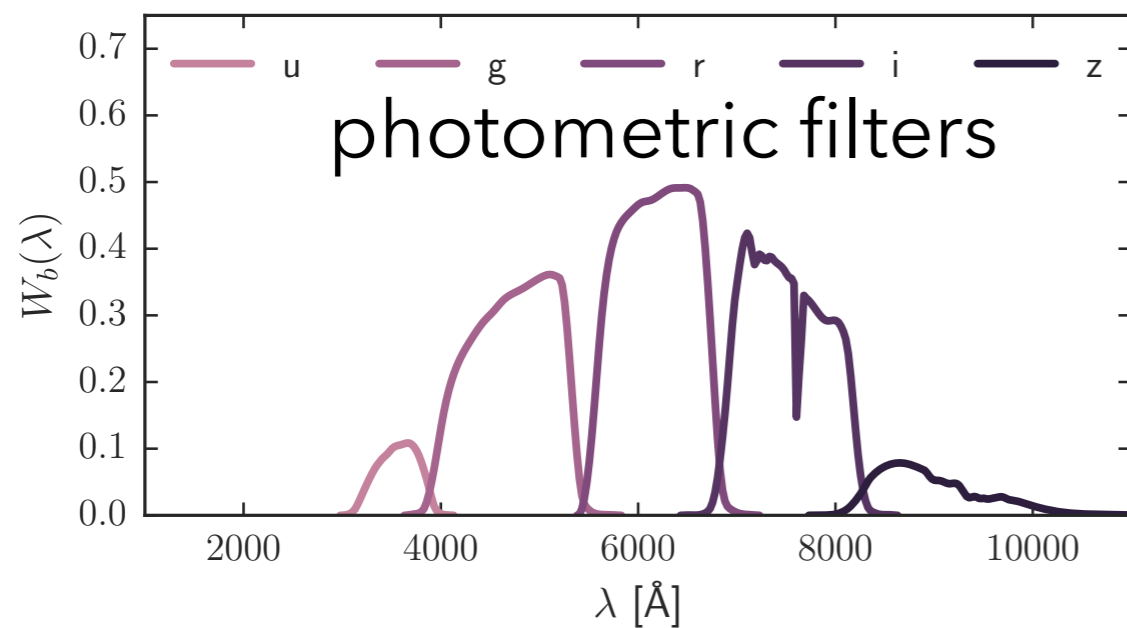
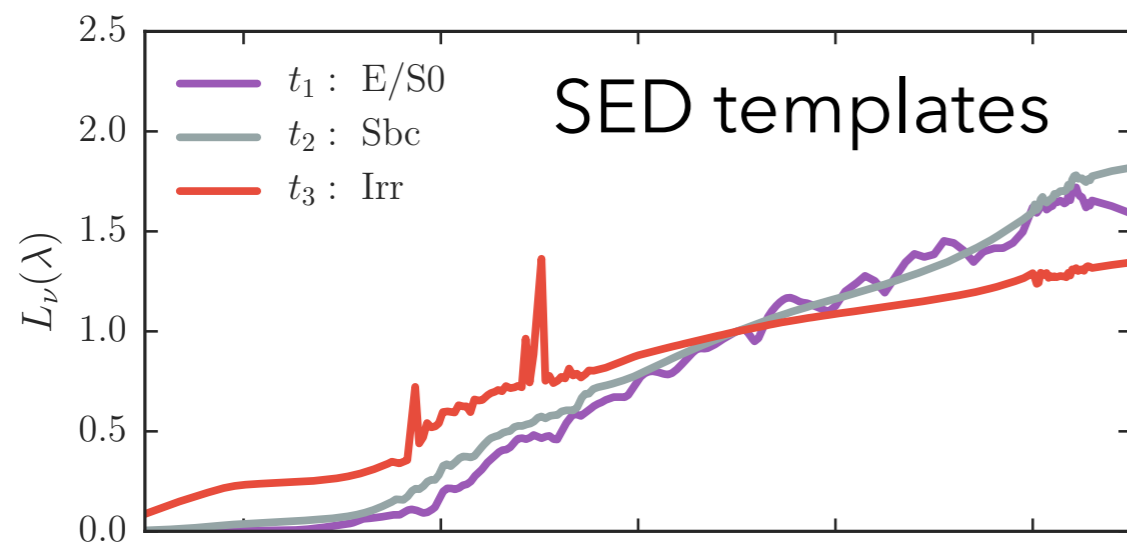
Hierarchical N(z) inference

$$p(N(z, t, m), \{z_i, t_i, m_i\} | \{\text{Fluxes}_i\})$$
$$\propto \underbrace{p(N(z, t, m))}_{\text{prior}} \prod_{i=1}^N \underbrace{p(z_i, t_i, m_i | N(z))}_{\text{population}} \underbrace{p(\text{Fluxes}_i | z_i, t_i, m_i)}_{\text{likelihood}}$$

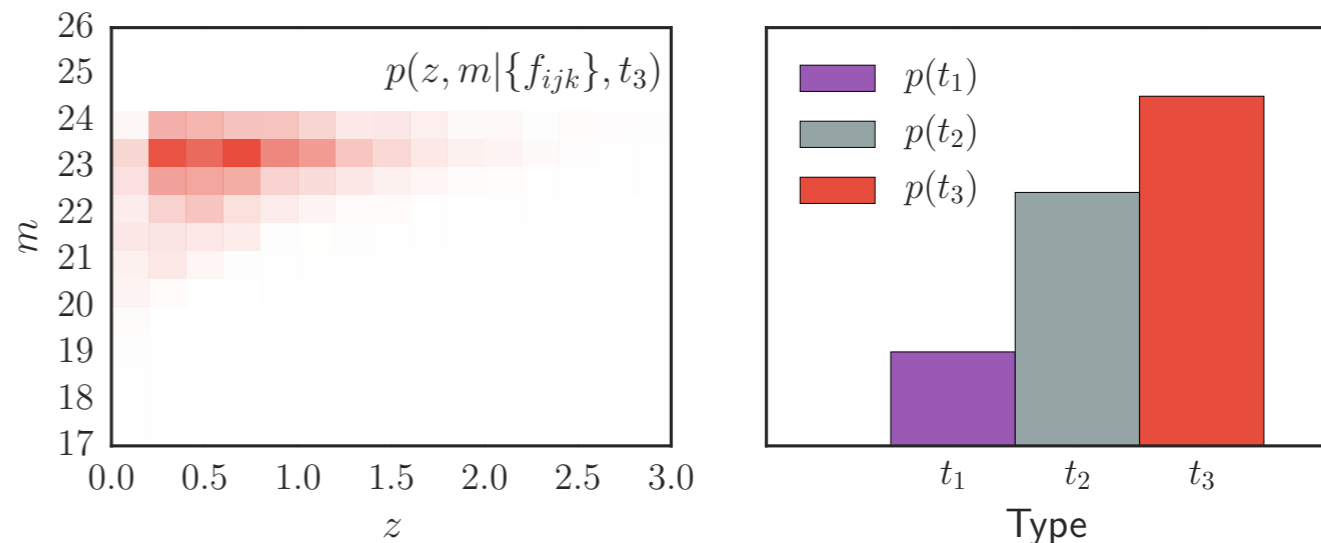
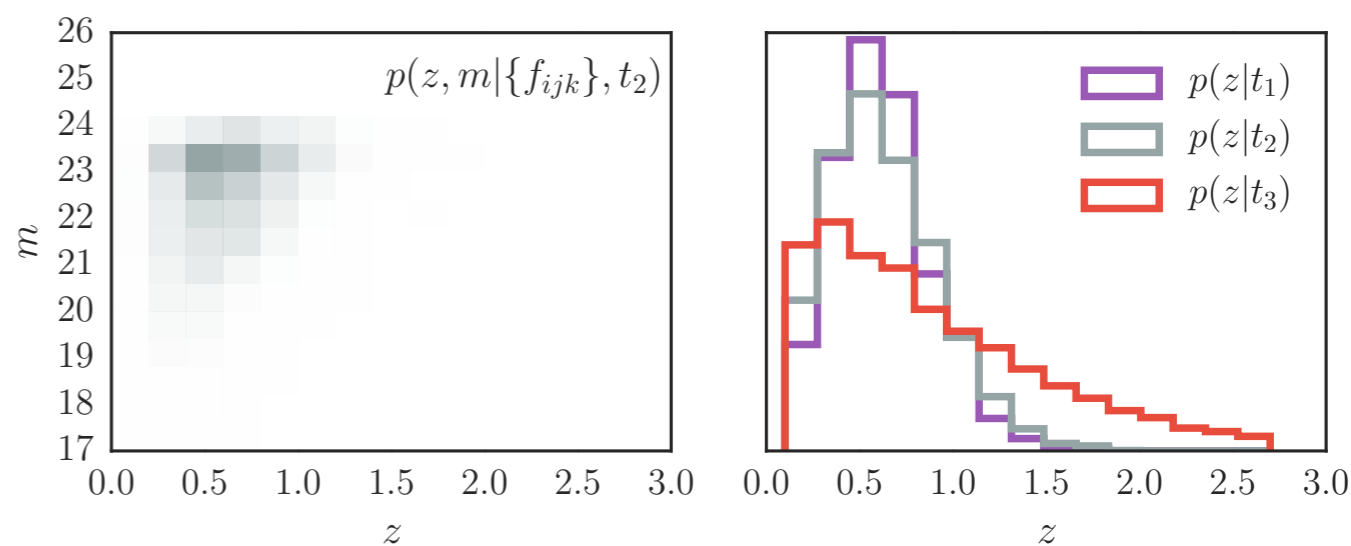
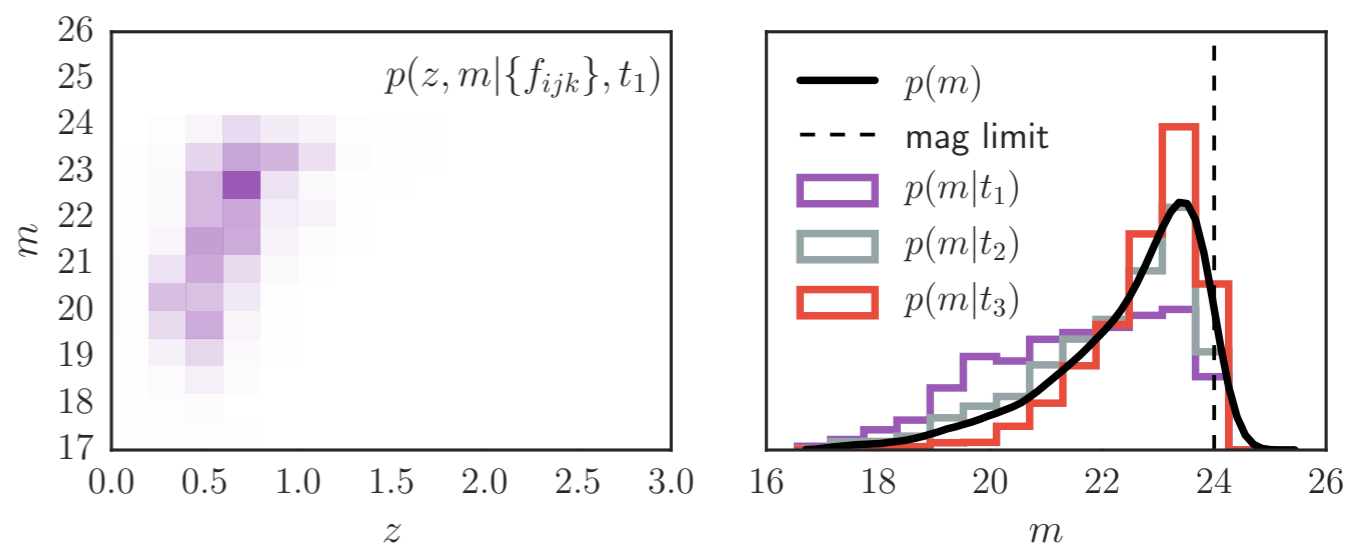
Likelihood based on SEDs, assumed to be **correct**.
Histogram model of $N(z, t, m)$ parameterized by $\{f_{ijk}\}$

Jointly infer $\{z, t, m\}_{\text{objects}}$ and $\{f_{ijk}\}$ using Gibbs sampler

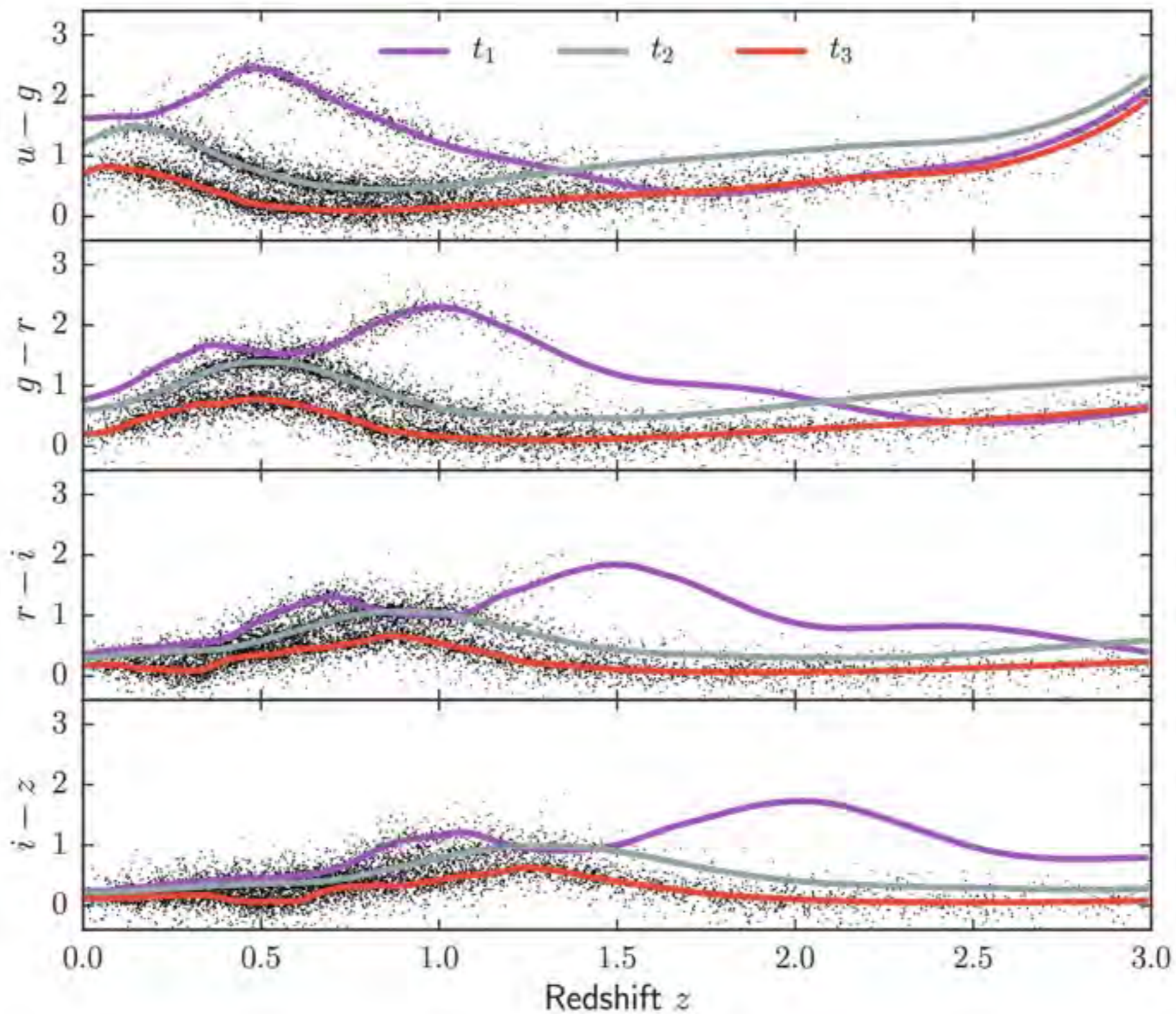
- ▶ 3 templates
- ▶ ugriz filters
- ▶ realistic distributions



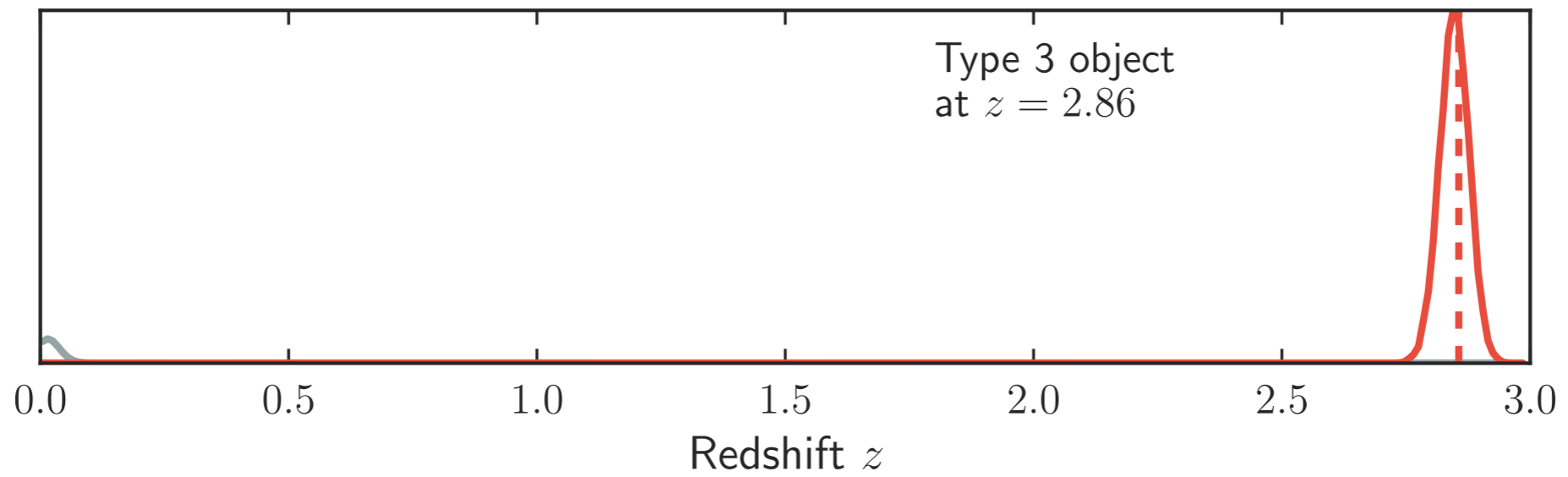
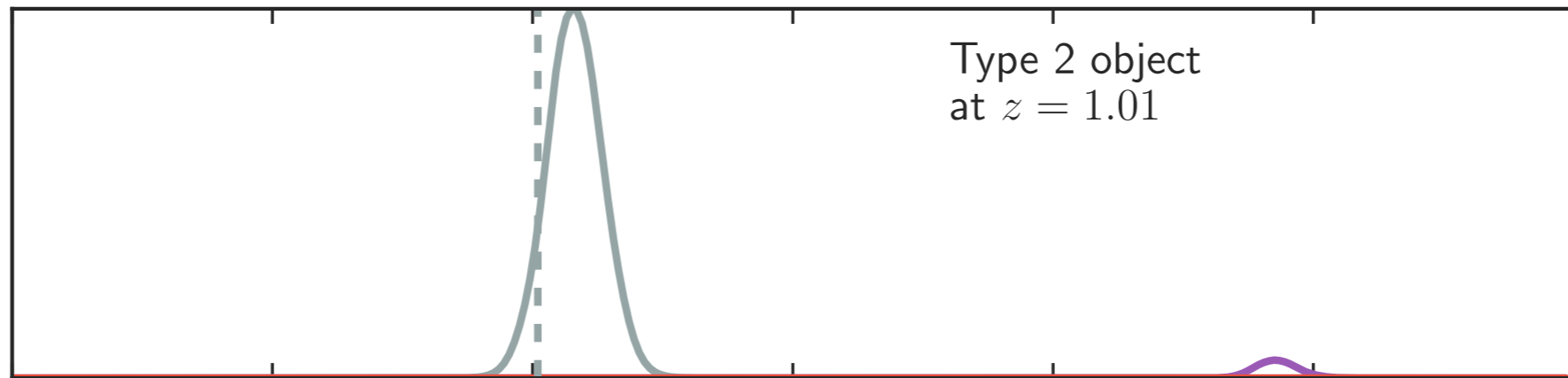
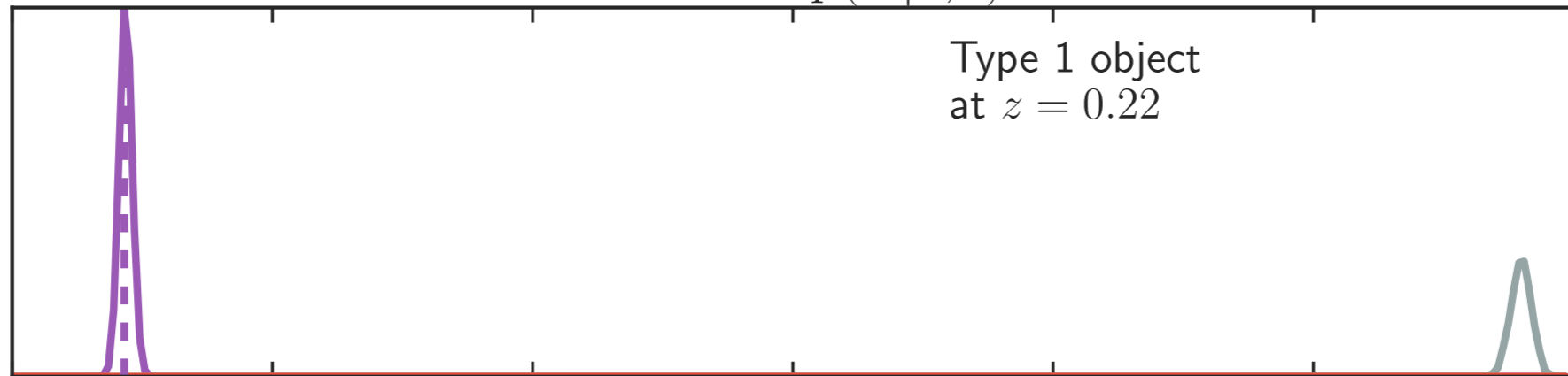
distribution $N(z, t, m)$ of the simulation



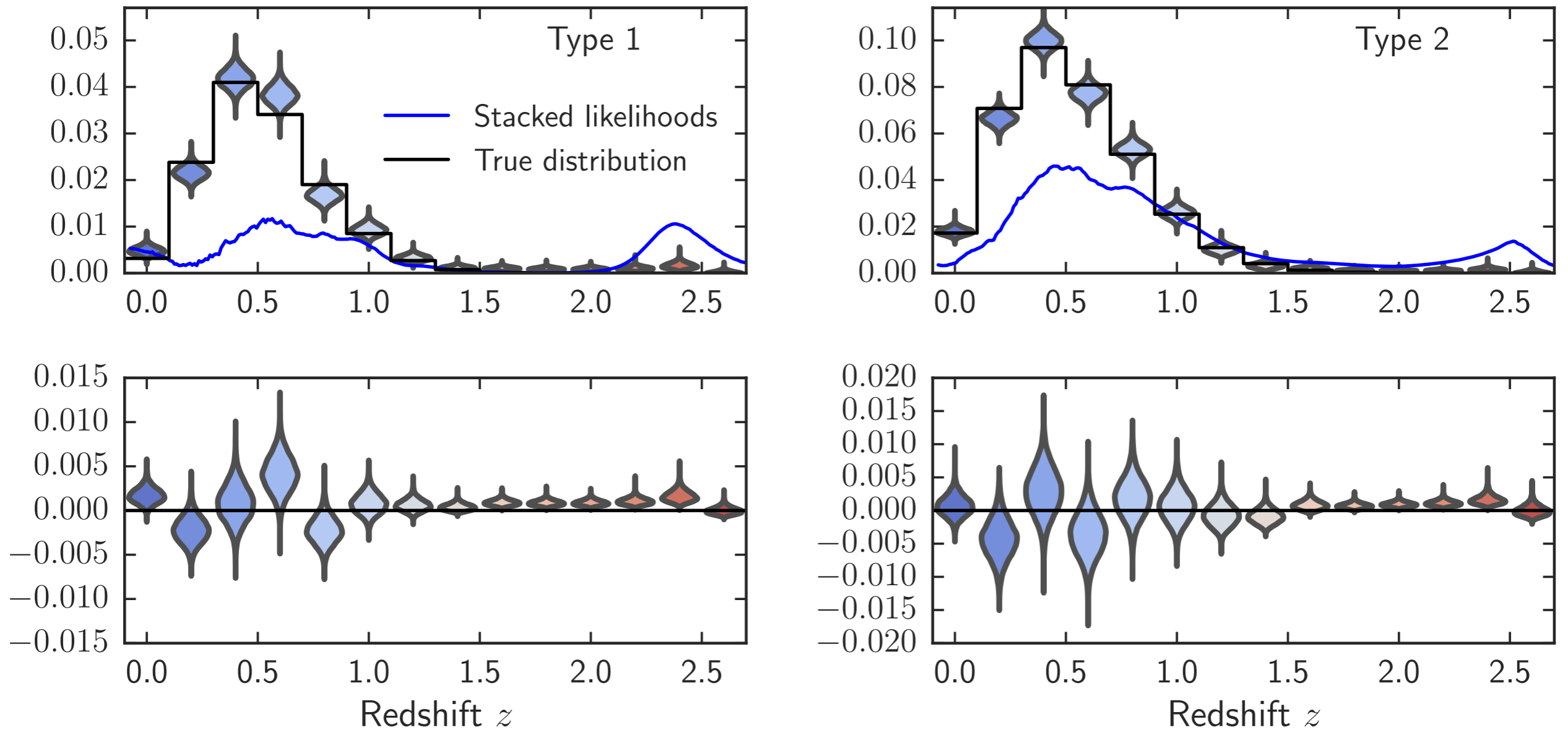
Simulated colors of 10^4 galaxies



Likelihood $p(F|z, t)$

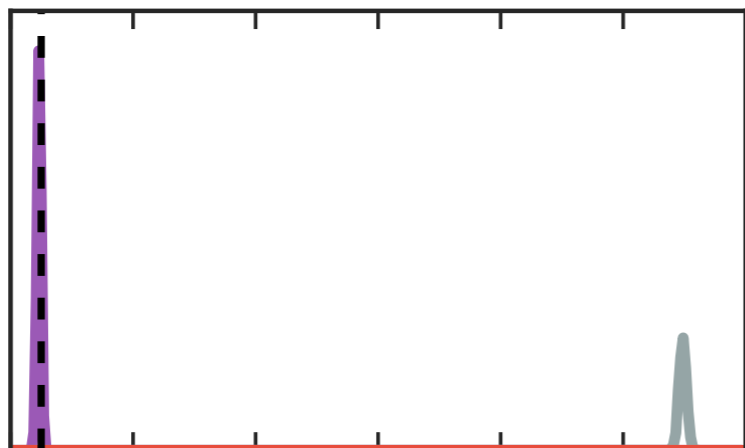


after inference: samples of the full posterior distributions

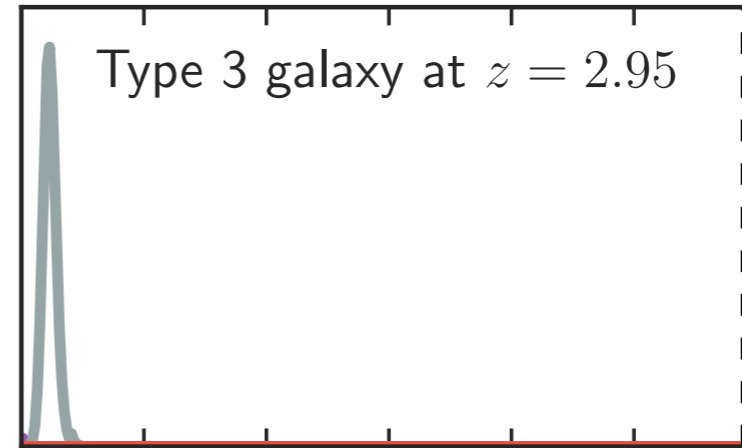
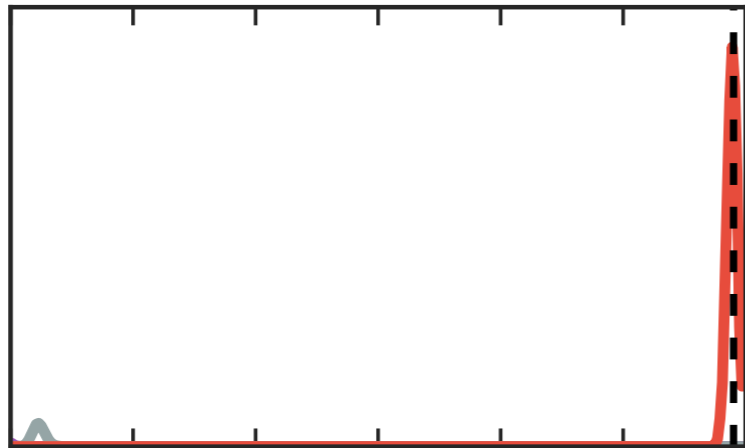
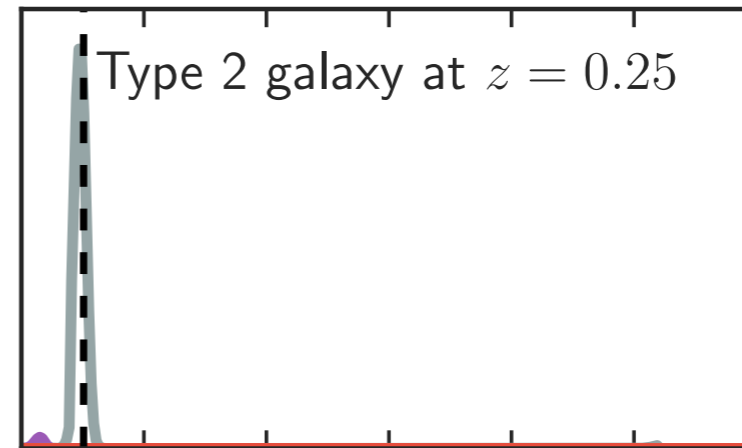
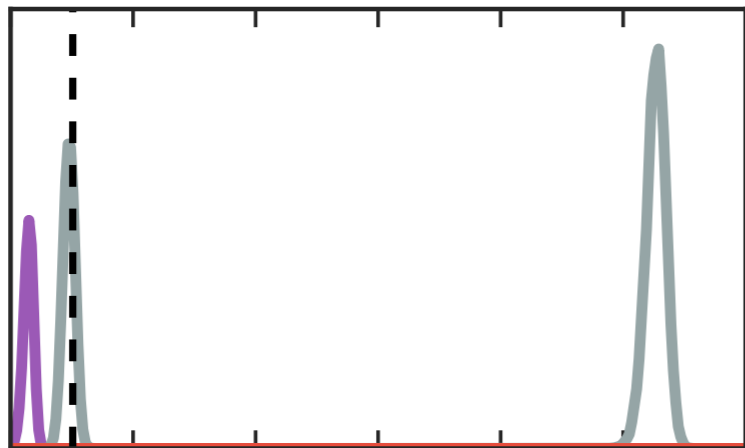
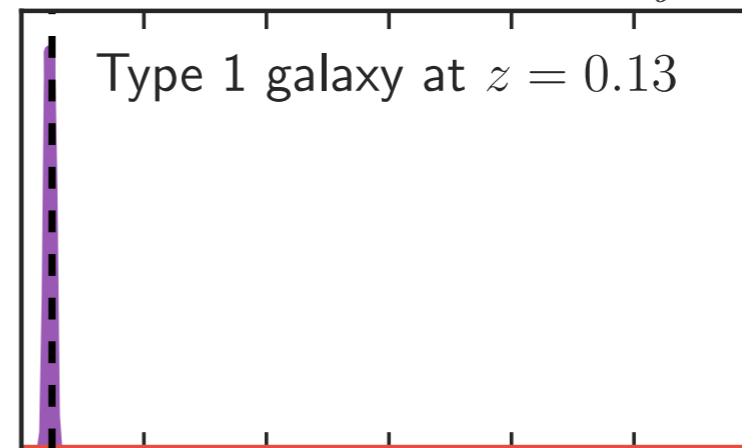


Redshift distributions correctly recovered
despite strong degeneracies

Likelihood $p(\{\hat{F}_b\}|z, t)$

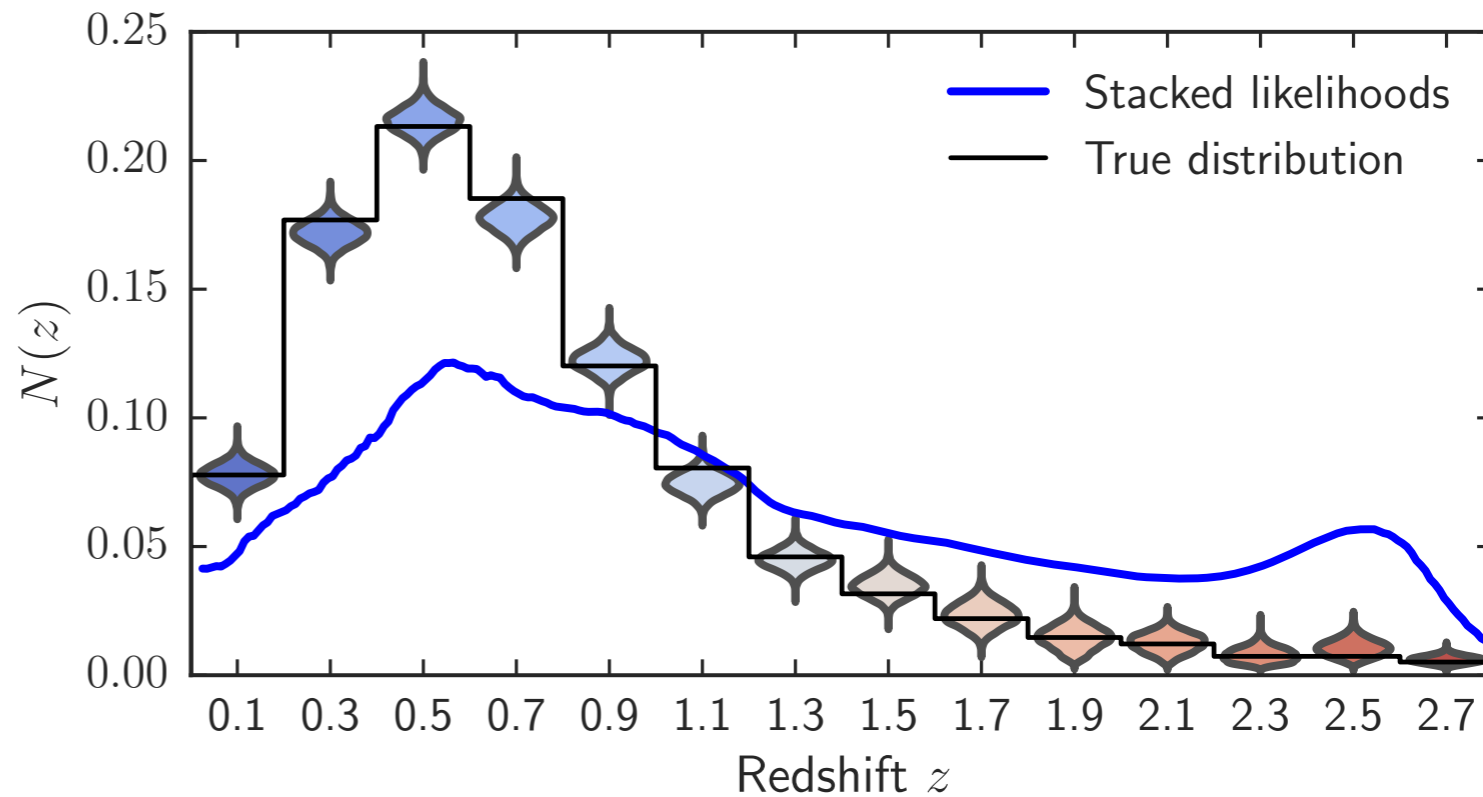
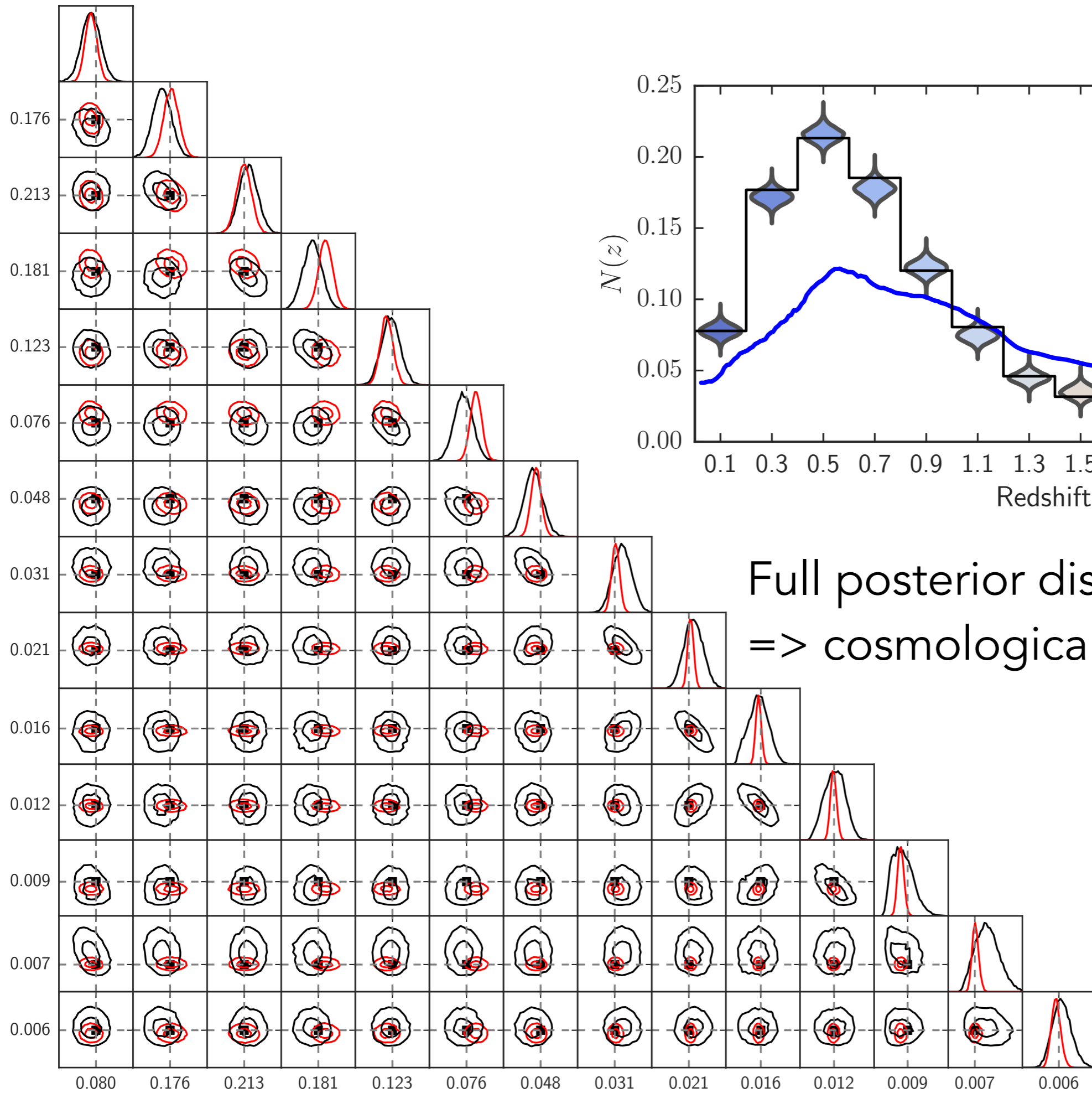


Posterior $p(z, t|\{\hat{F}_b\}, \{f_{ijk}\})$



0.0 0.5 1.0 1.5 2.0 2.5 3.0
Redshift z

0.0 0.5 1.0 1.5 2.0 2.5 3.0
Redshift z



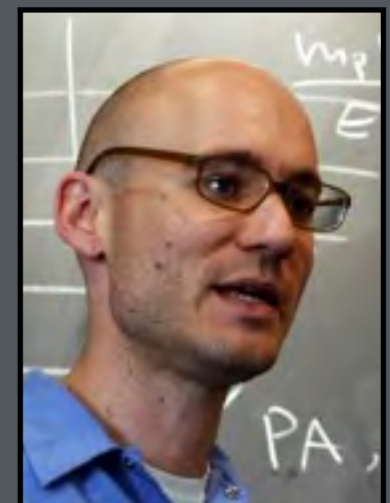
Full posterior distribution for $N(z)$
 => cosmological likelihoods

- Hierarchical probabilistic inference of $N(z)$
- $N(z)$ uncertainties captured and can be propagated into cosmological likelihoods

However, real likelihoods are incorrect/biased!

*Data-driven, interpretable
photometric redshifts
trained on heterogeneous and
unrepresentative data*

with David Hogg (NYU)



Will **never** have representative spectroscopic data

Galaxy SED models are not precise enough

Only deep spectroscopic & many-band surveys available

True PDFs needed with data **and** model uncertainties

Machine learning constrained by physics of the problem?

Idea

Target set: photometric survey

Training set: many-band or spectroscopic set
= deeper, heterogeneous version of target

No complete physical model for galaxy spectra
=> construct spectra compatible with training set

^

$$p(z|\text{Fluxes}) = \sum_j p(\text{Fluxes}|\text{Fluxes}(t_j, z)) p(z|t_j) p(t_j)$$

Classical template fitting *(e.g., BPZ, EAZY, ZEBRA, etc)*

- ▶ Use a small set of fixed templates based on low-redshift bright spectra or physical models

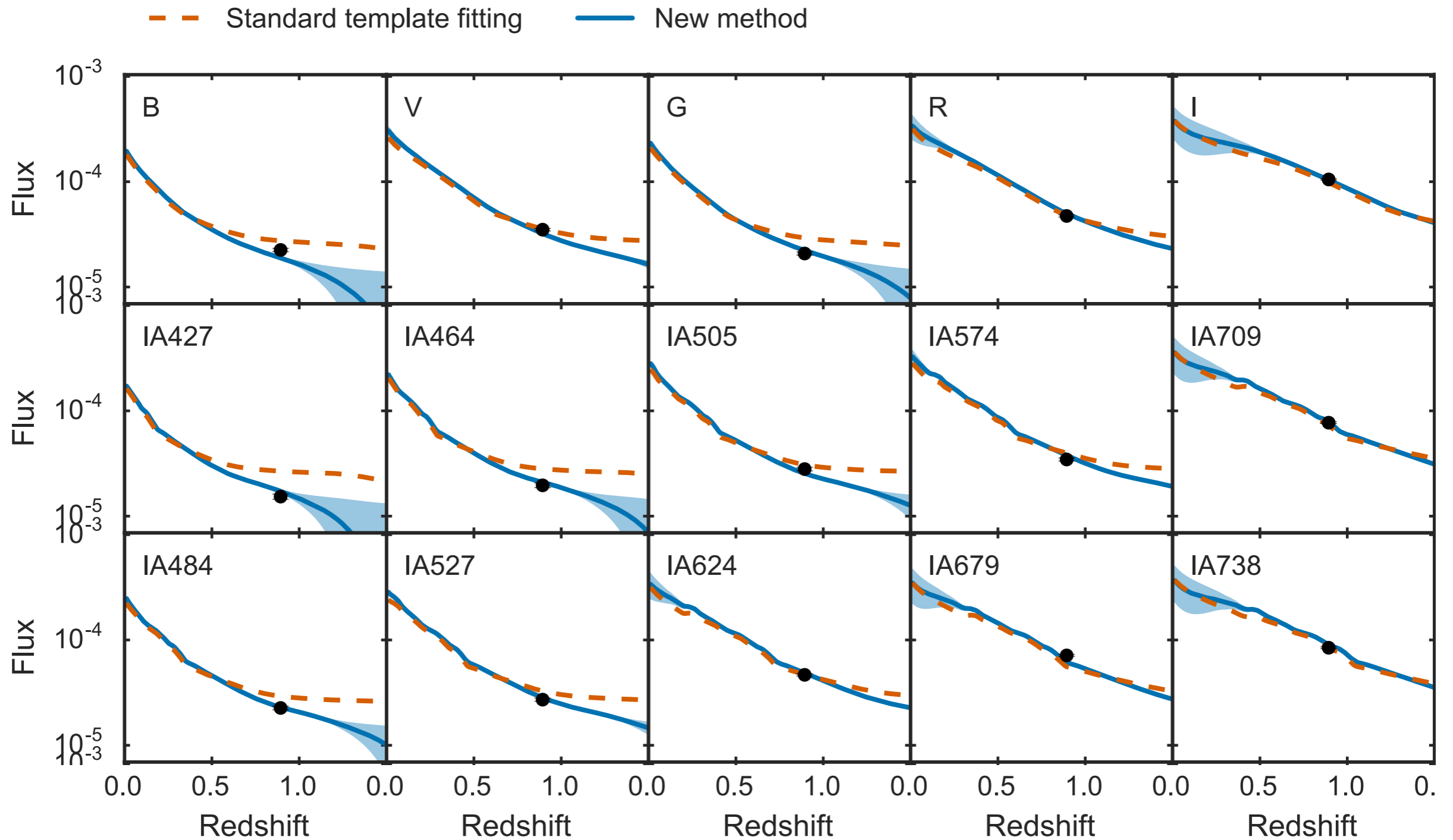
New data driven approach 1 *(work in progress...)*

- ▶ Forward model a probabilistic system of templates and priors, to be constrained from the training data.

New data driven approach 2 *(Leistedt & Hogg, arXiv:1612.00847)*

- ▶ Construct one probabilistic template per training galaxy. Pairwise comparison of target galaxies (redshifts unknown) with training galaxies (redshift known or constrained)

Example of a model per training galaxy



New method: DELIGHT™

Leistedt & Hogg (arXiv:1612.00847) – github.com/ixkael/Delight

Concept: implicitly fitting and redshifting SEDs to each training galaxy for pairwise comparison with target galaxies
= *machine learning + template fitting*

Probabilistic, physical, and data driven

Interpretable model & PDFs. Flexibility via parameters.

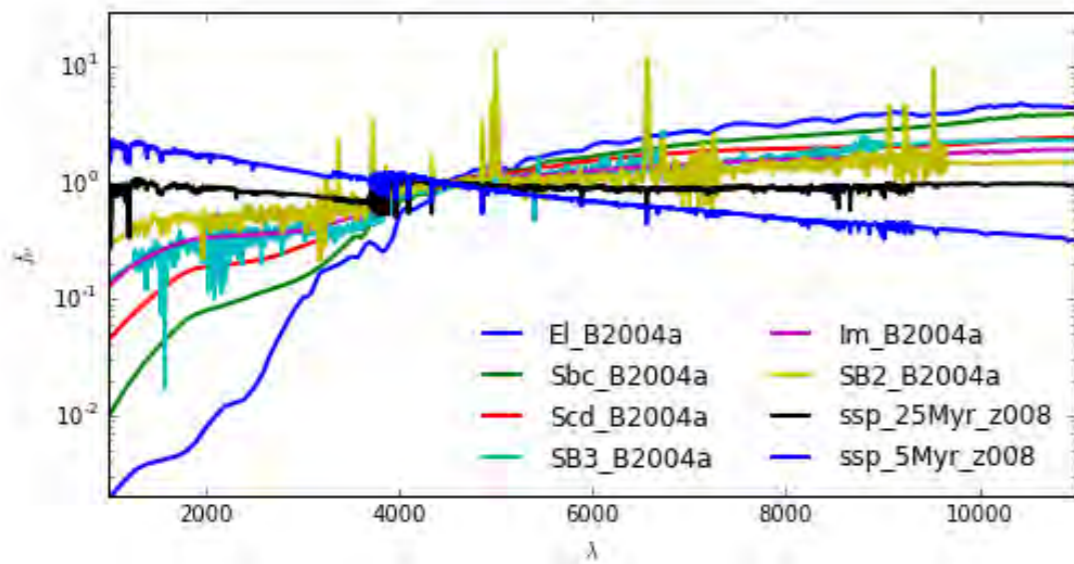
Use much more data than existing methods: heterogeneous combination of spectroscopic or deeper photometric data

Fast to (re-)train/apply. No need to store tabulated PDFs.

SED model

How to quickly construct SED model and make predictions?

$$p(\underbrace{\{\hat{F}'_b\}}_{\text{target}} | z', t_i) = p(\underbrace{\{\hat{F}_b\}}_{\text{training}} | z', z, \{\hat{F}_b\})$$

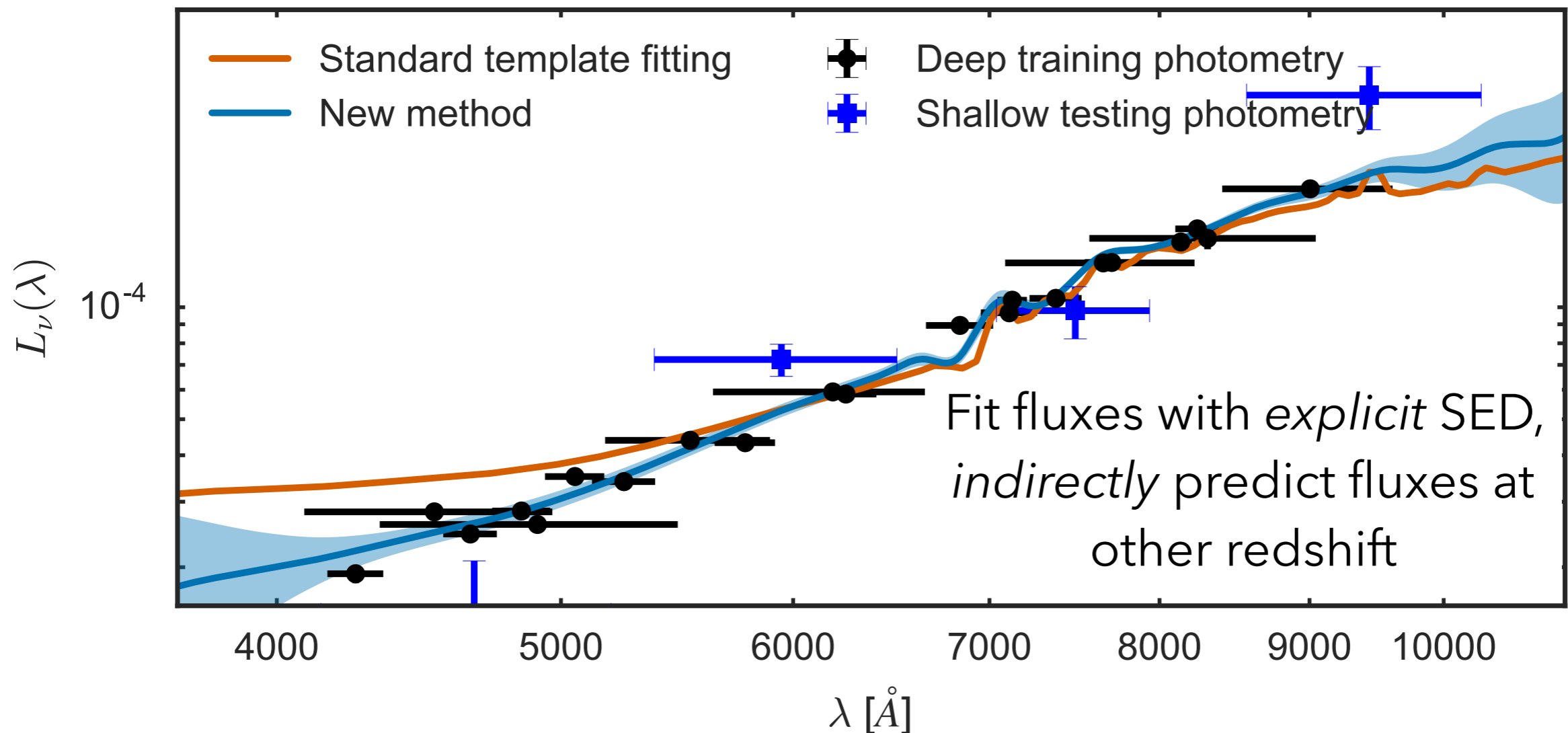


SED model:

$$L_\nu(\lambda) = \underbrace{\sum_k \alpha_k T_\nu^k(\lambda)}_{\text{templates}} + \underbrace{R_\nu(\lambda)}_{\text{residuals}}$$

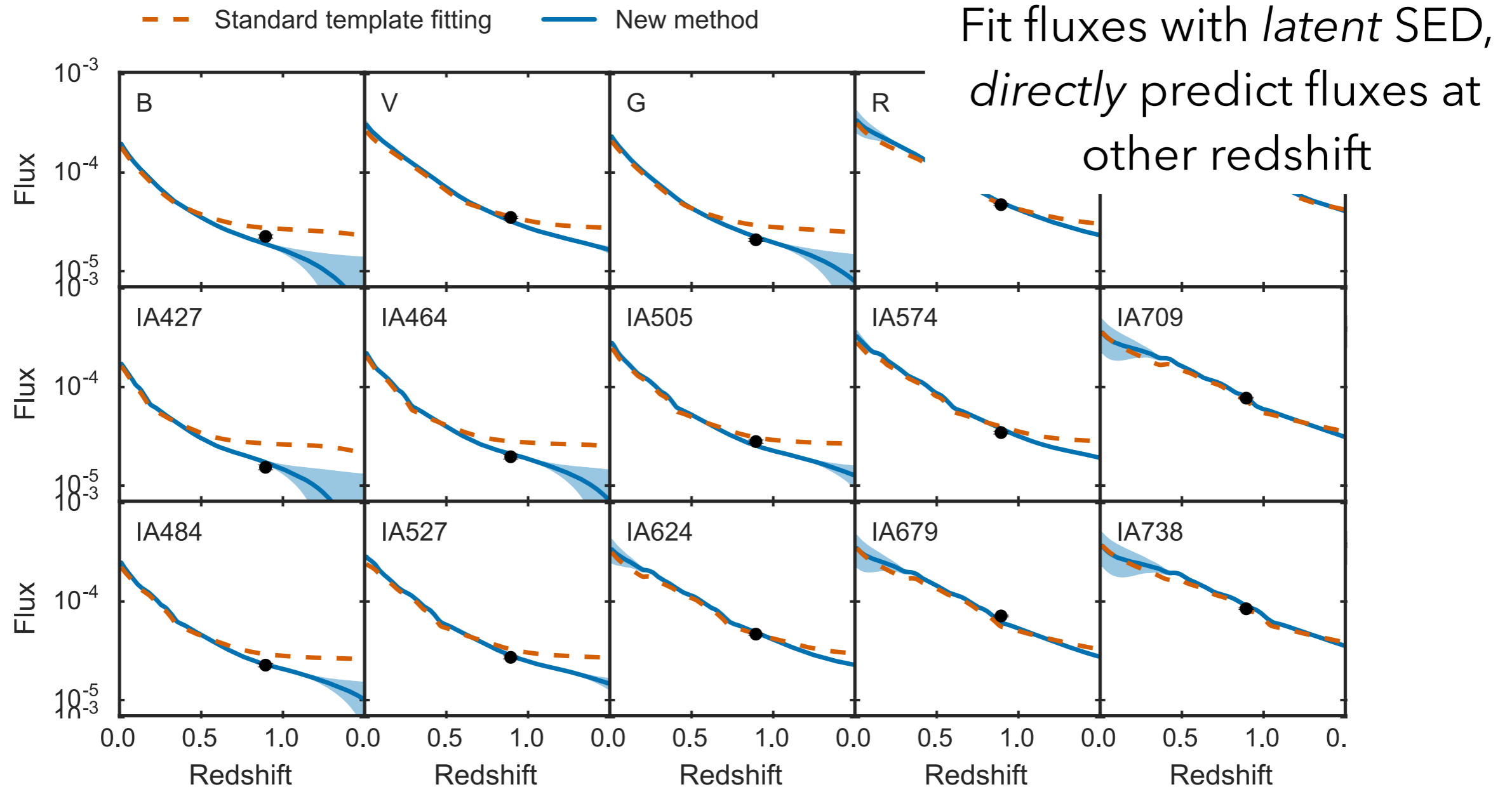
The ~~crazy~~ intractable way

Explore all SEDs compatible with training galaxy
(noisy fluxes + spec-z) via MCMC



The ~~elegant~~ efficient way

Directly fit for training galaxy in flux-redshift space
+ *force the fit to correspond to underlying SEDs*



*<introduction
to Gaussian Processes>*

Gaussian processes

$$f \sim \mathcal{GP} \iff p(f(\vec{x}), f(\vec{x}')) \text{ is Gaussian } \forall \vec{x}, \vec{x}'$$

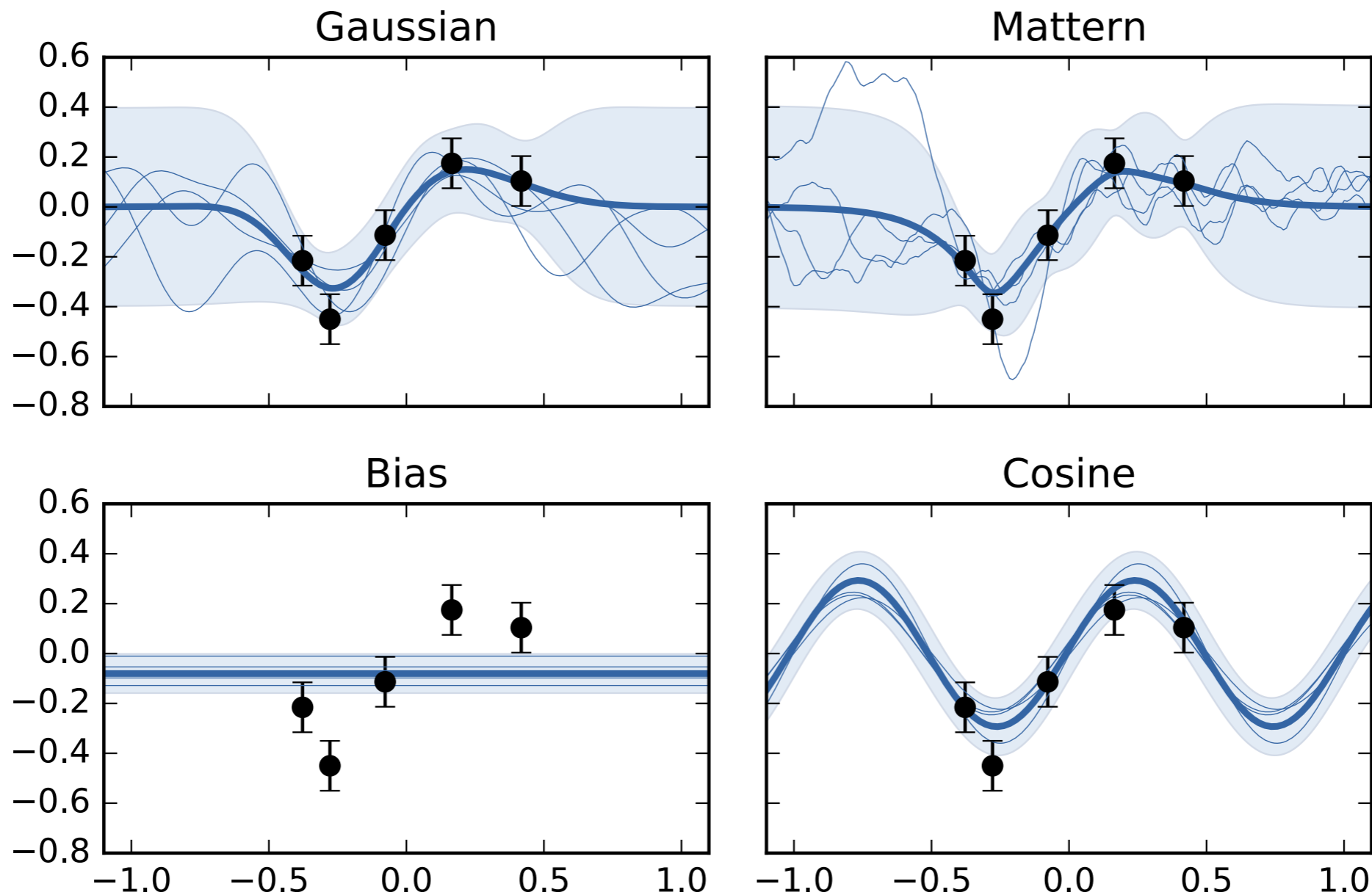
characterized by mean and kernel

$$m(\vec{x}) = \mathbb{E}[f(\vec{x})]$$

$$k(\vec{x}, \vec{x}') = \mathbb{E}[(f(\vec{x}) - m(\vec{x}))(f(\vec{x}') - m(\vec{x}'))]$$

for Gaussian likelihood, posterior/predictions **tractable**
see Rasmussen & Williams (2006)

Fitting with GPs = using priors over functions
Modelling correlated signal and/or noise
Choice of kernel is key (captures correlations)



*</introduction
to Gaussian Processes>*

Photo-z gaussian process

if SED model is: $L_{\nu}(\lambda) \sim \mathcal{GP}\left(\sum_k \alpha_k T_{\nu}^k(\lambda), k(\lambda, \lambda')\right)$
 k templates residuals

then the fluxes: $F(b, z) \sim \mathcal{GP}\left(\mu^F(b, z), k^F(b, b', z, z')\right)$
mean flux and covariance

GP with physical mean function **and residuals**

Fitting and predicting photometric fluxes while capturing the physics of redshifts

Analytically tractable under simple assumptions

Photo-z gaussian process (proof)

Redshifted
galaxy SED

$$f_\nu(\lambda_{\text{obs}}, z) = \frac{(1+z)}{4\pi D_L^2(z)} L_\nu\left(\frac{\lambda_{\text{obs}}}{(1+z)}\right)$$

Photometric
fluxes

$$F(b, z) = \frac{\int_0^\infty f_\nu(\lambda, z) W_b(\lambda) d\lambda/\lambda}{\int_0^\infty g^{\text{AB}} W_b(\lambda) d\lambda/\lambda}$$

SED model

$$L_\nu(\lambda) = \underbrace{\sum_k \alpha_k T_\nu^k(\lambda)}_{\text{templates}} + \underbrace{R_\nu(\lambda)}_{\text{residuals}}$$

Photo-z gaussian process (proof)

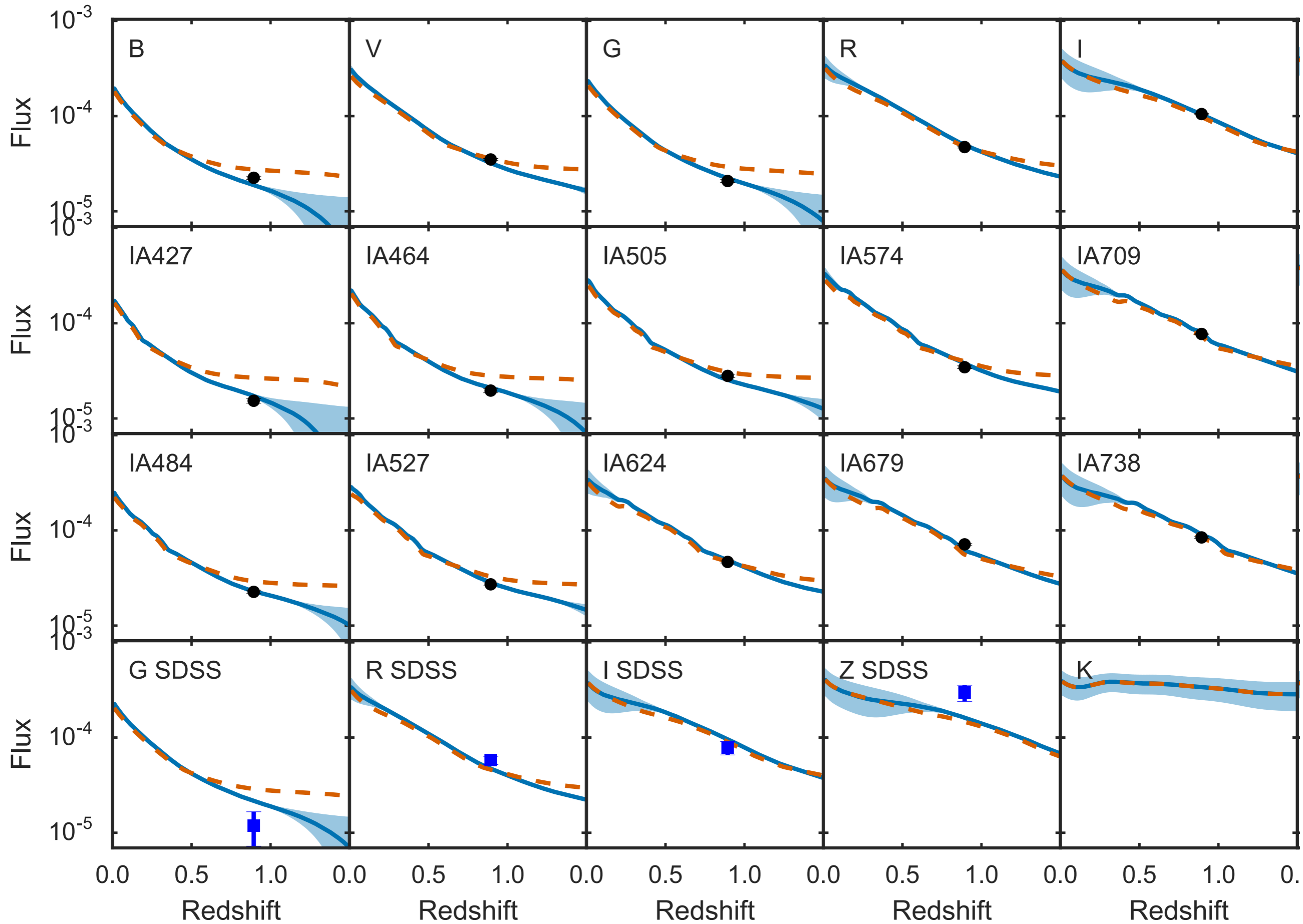
$$R_\nu(\lambda) \sim \mathcal{GP}\left(0, k(\lambda, \lambda')\right)$$

$$\implies L_\nu(\lambda) \sim \mathcal{GP}\left(\sum_k \alpha_k T_\nu^k(\lambda), k(\lambda, \lambda')\right)$$

$$F(b, z) = \frac{(1+z)}{4\pi D_L^2(z) C_b} \int_0^\infty L_\nu\left(\frac{\lambda}{(1+z)}\right) W_b(\lambda) d\lambda/\lambda$$

$$\implies F(b, z) \sim \mathcal{GP}\left(\mu^F(b, z), k^F(b, b', z, z')\right)$$

Standard template fitting New method



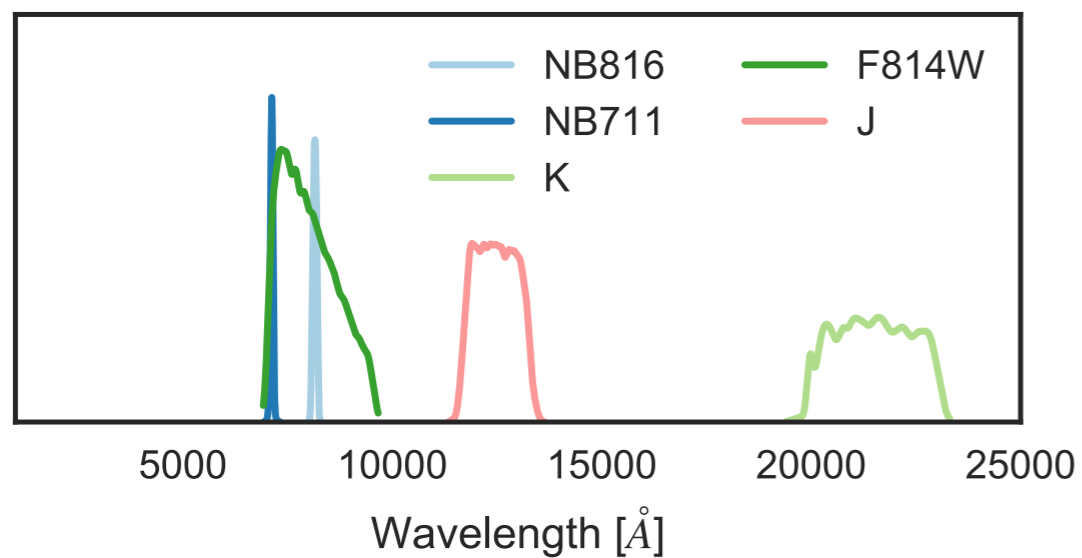
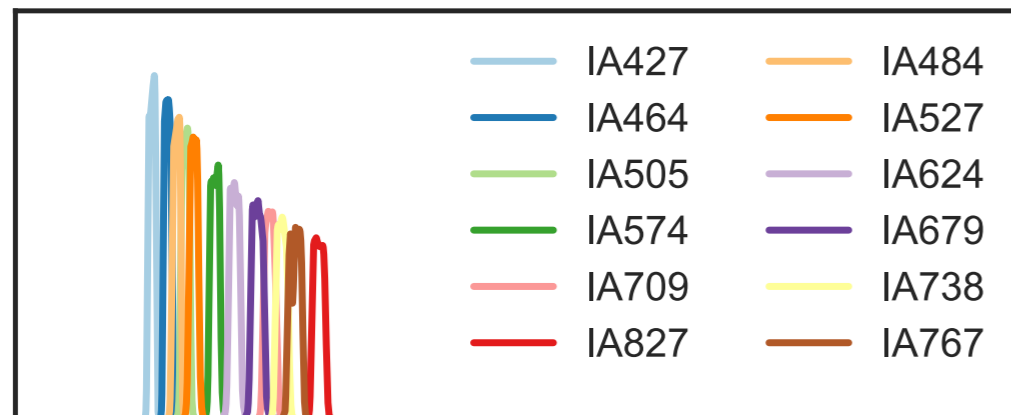
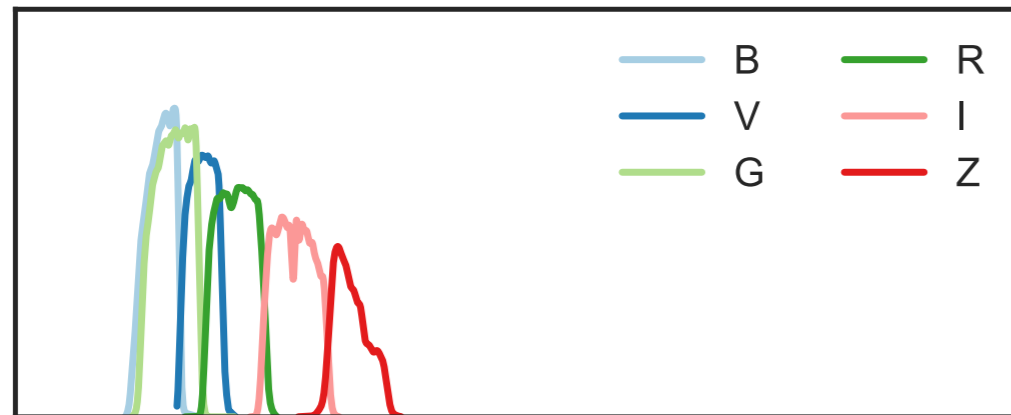
G10 / COSMOS data

training: deep SUBARU/HST bands
with spectroscopic redshifts

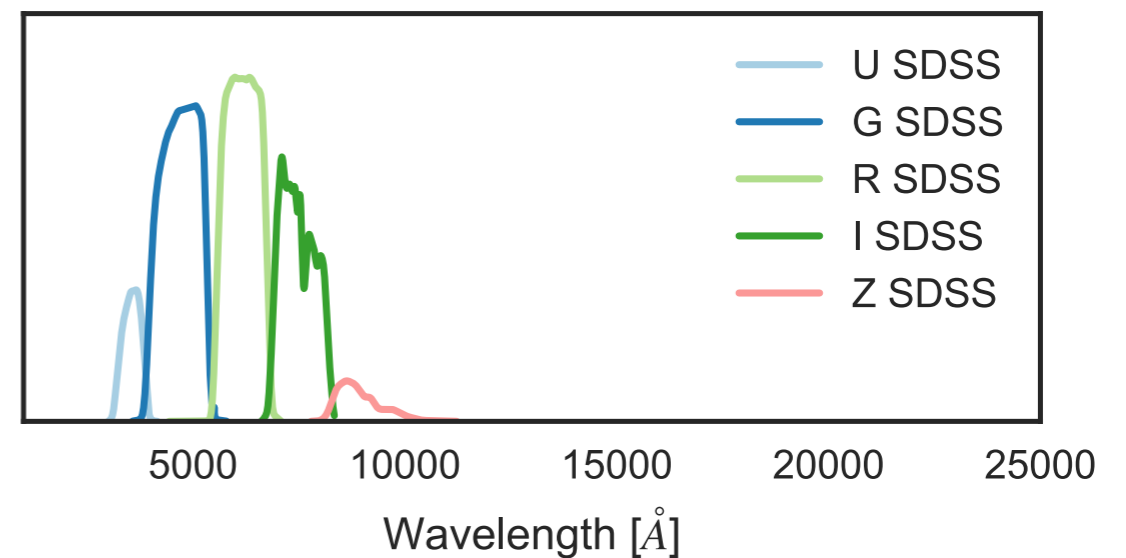
target: *ugriz* SDSS bands

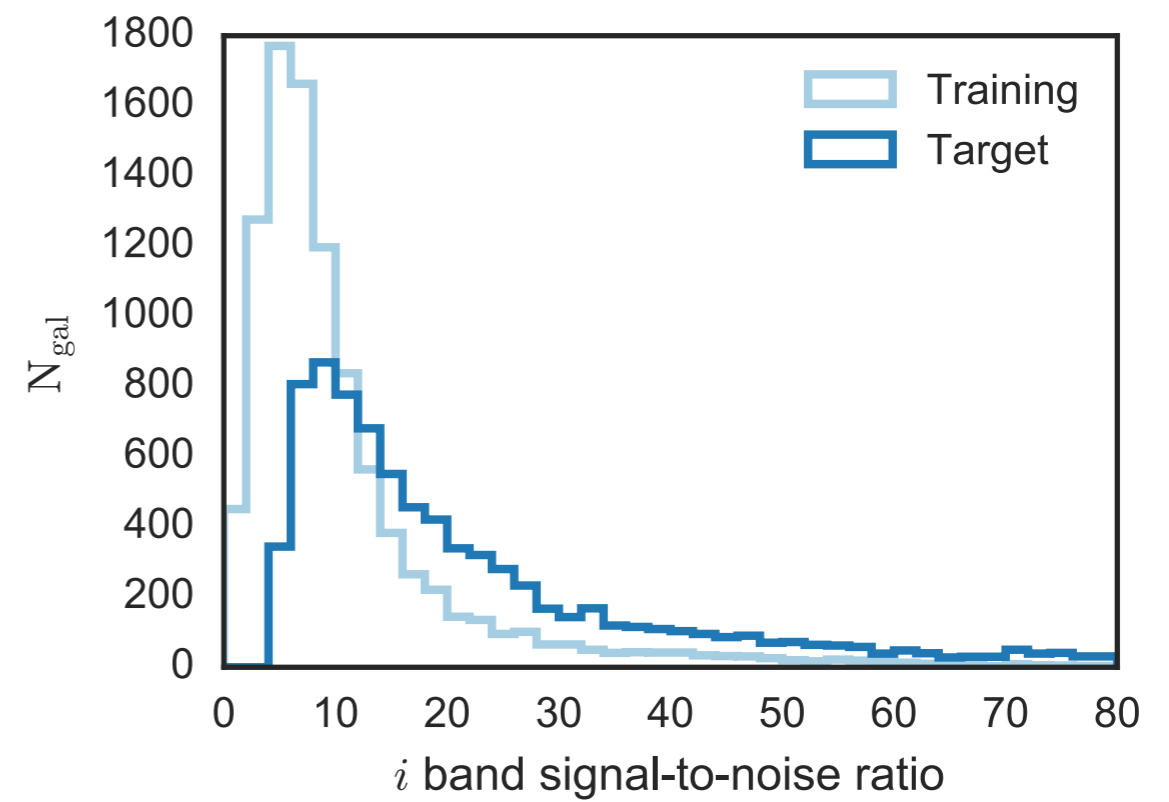
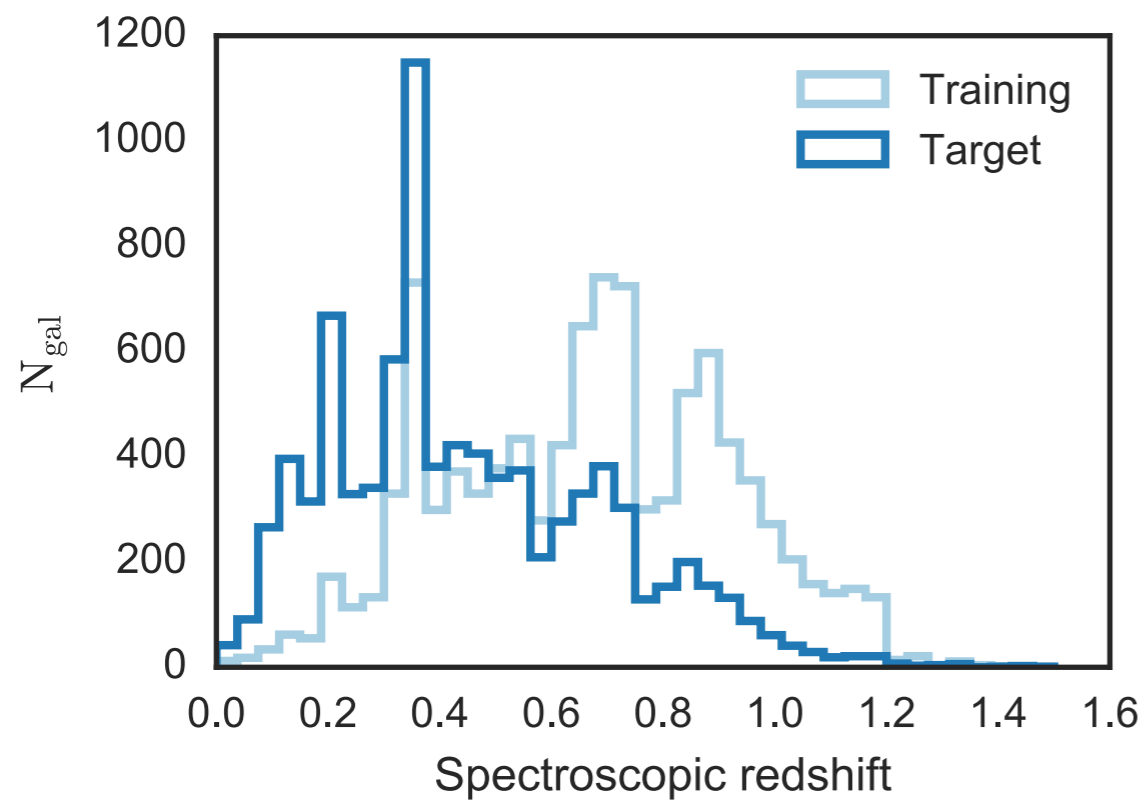
training/target: 10k/10k objects

Photometric filters (training)



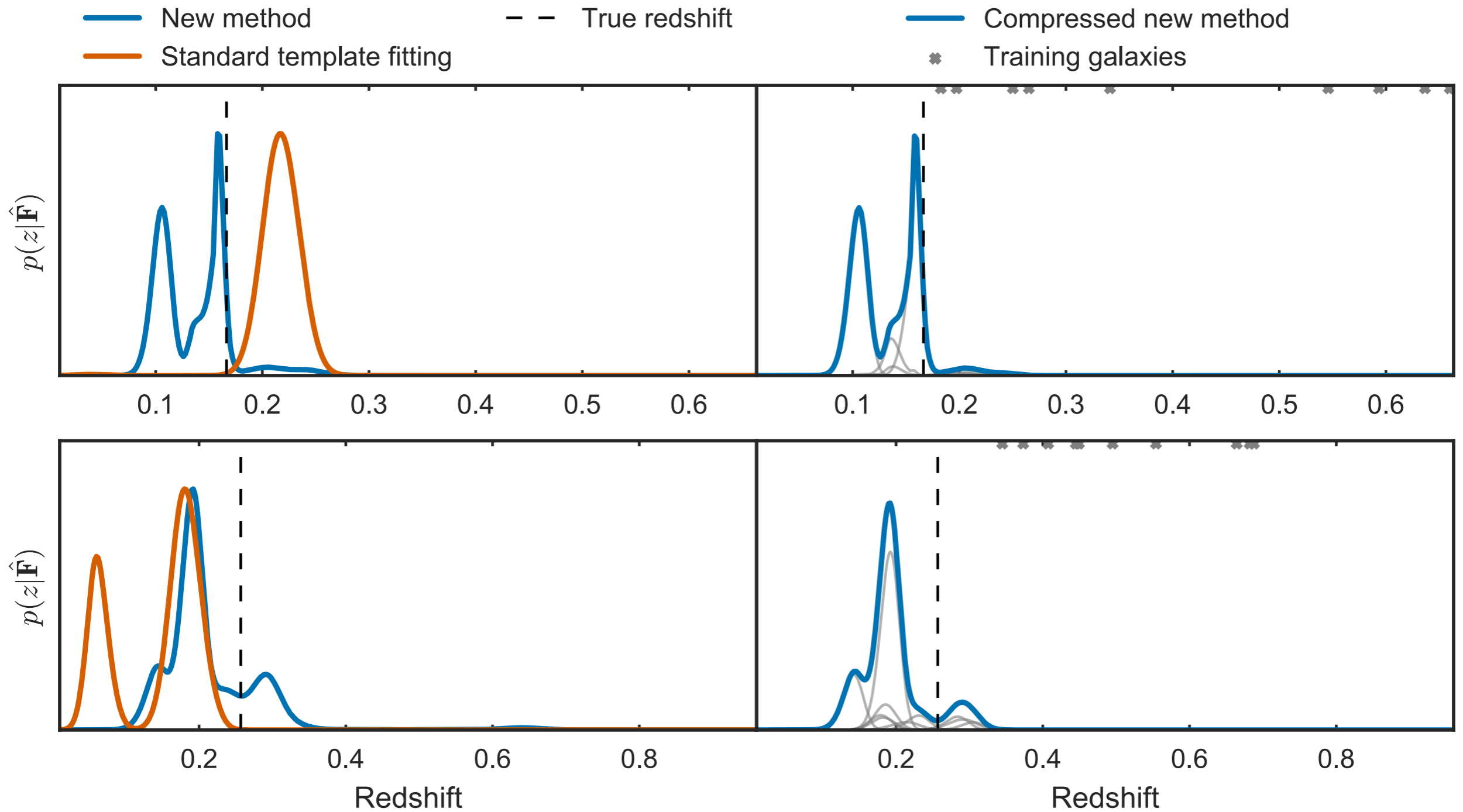
Photometric filters (target)

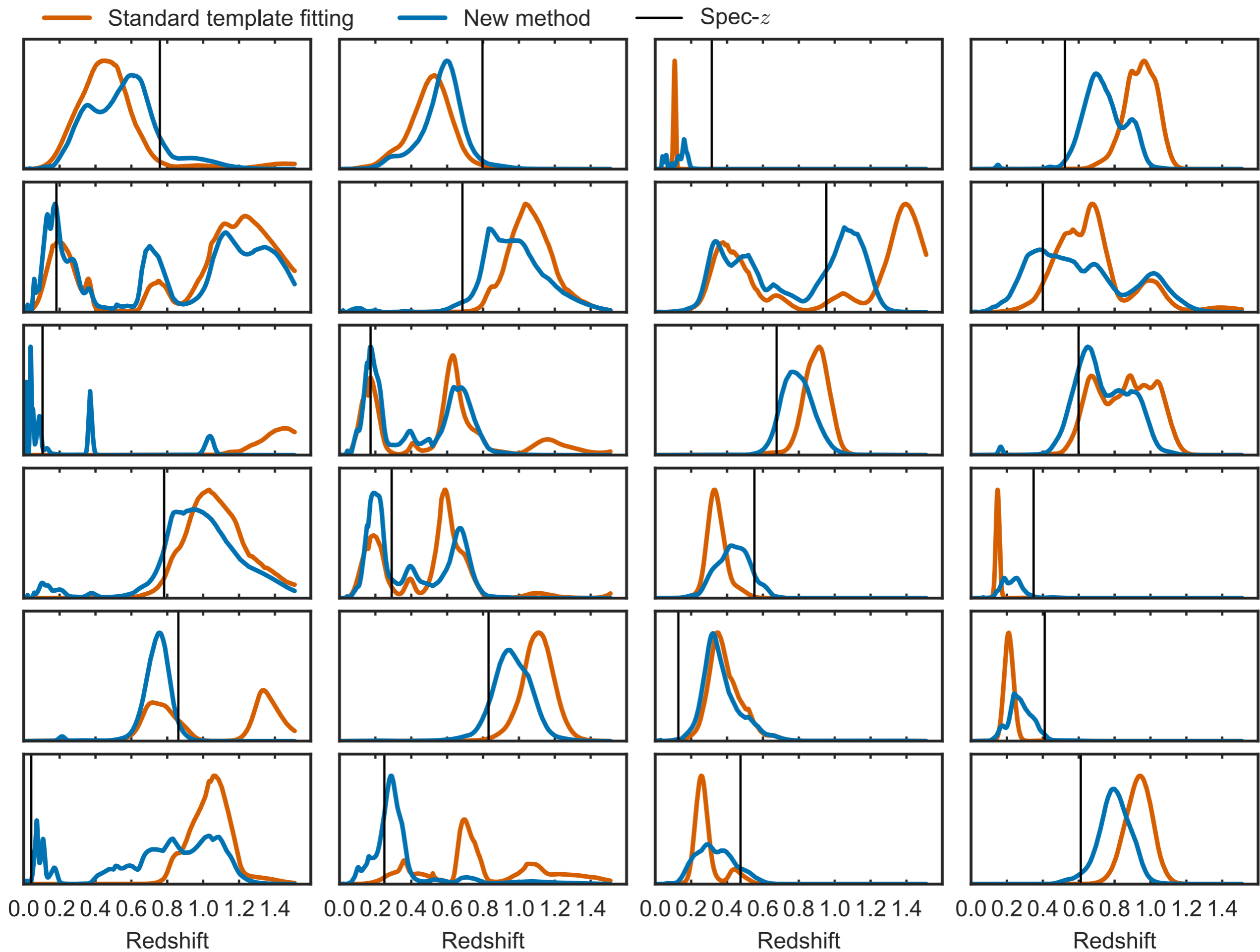


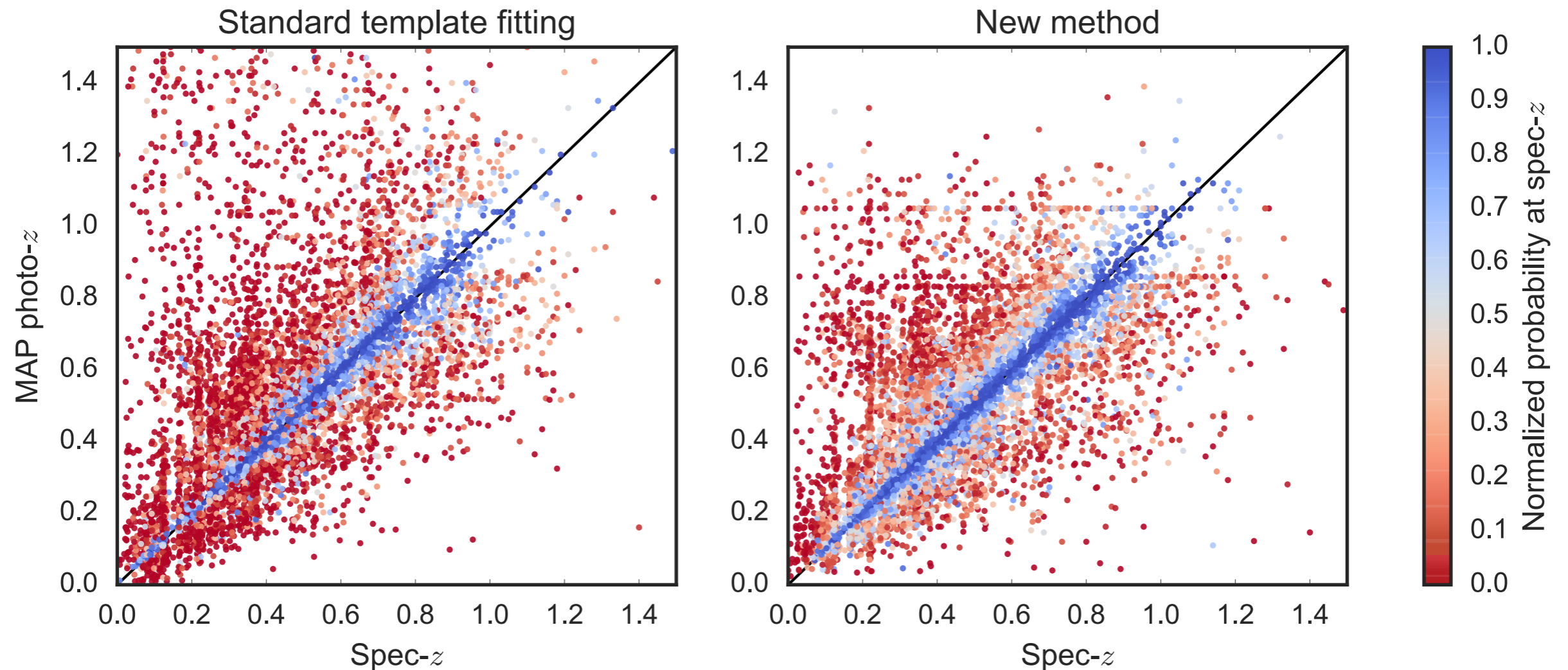


unrepresentative training set with different bands & noise

a closer look at two PDFs...







7 fixed templates \Rightarrow 10,000 probabilistic templates
 (a system of types) (one per training galaxy)

*Improvement, but more data/flexibility required.
 Not exploiting low-dimensionality of galaxy types.*

Conclusions

Imaging surveys

*diverse science: fundamental physics, astrophysics
systematics limited — require exquisite photo-z's*

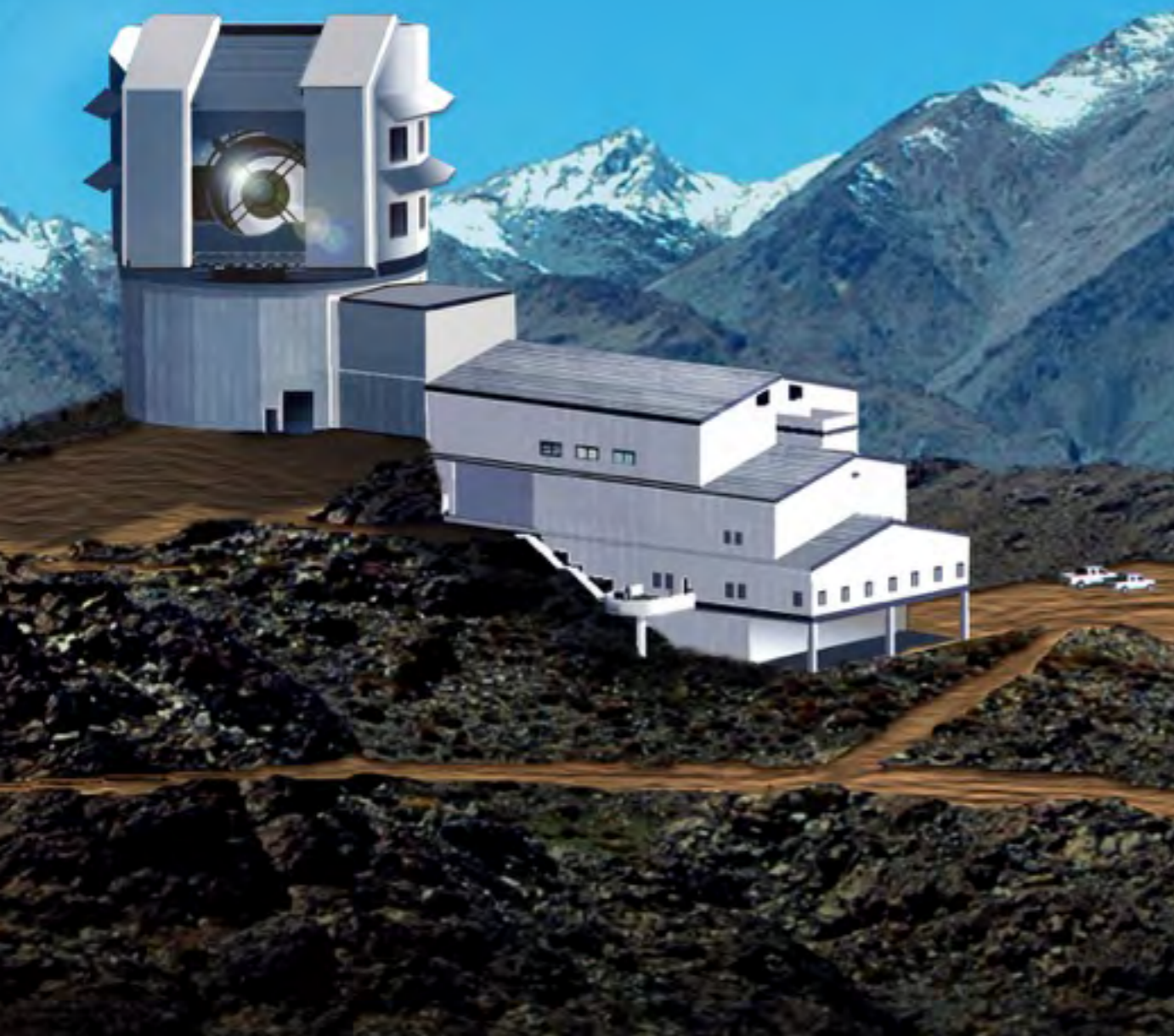
DELIGHT — [GITHUB.COM/IXKAEL/DELIGHT](https://github.com/ixkael/delight)

*data-driven method with physics & machine learning
delivers accurate, interpretable redshifts probabilities*

What's next?

*fit SED templates and luminosity functions, calibrate
photo-z likelihood without spectroscopic redshifts*

LSST
Large Synoptic Survey Telescope



20 billion galaxies

17 billion stars

*7 trillion sources detected
in single epochs*

30 trillion forced photometry

10 million alerts per night