

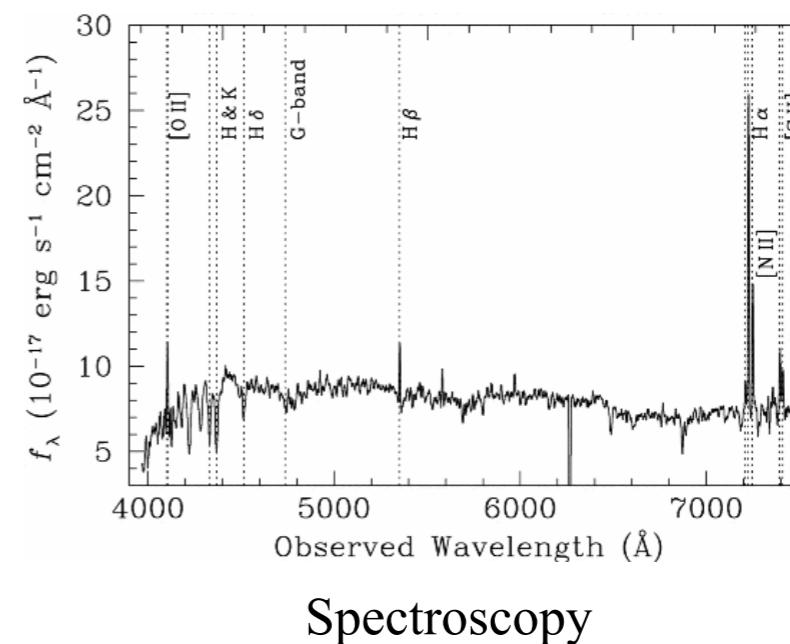
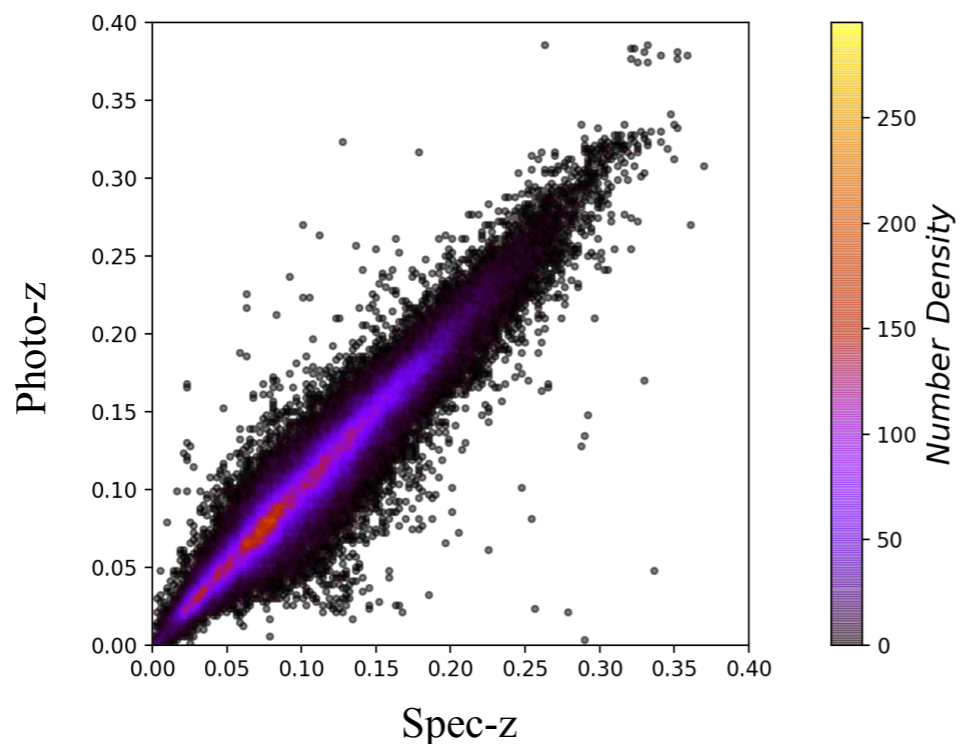
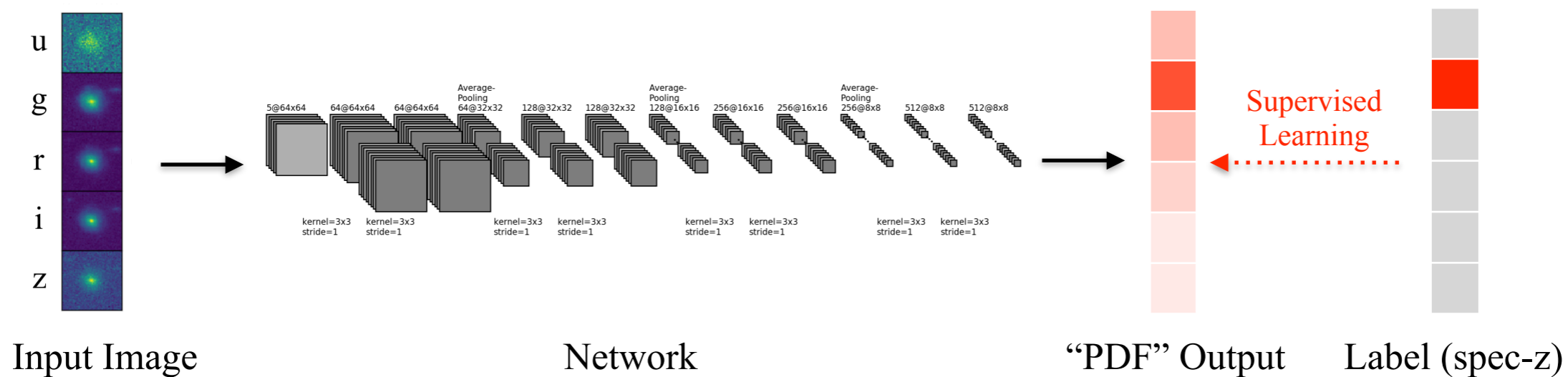
Estimation of photo-z probability density functions via deep learning with statistical basis

Presenter: Qiufan Lin (林秋帆)



Photo-z estimation as a computer vision problem supervised by spec-z

- ✓ State-of-the-art, best accuracy
- ✗ Direct PDF prediction lacks statistical basis, and may suffer from biases
- ✗ The network lacks interpretability



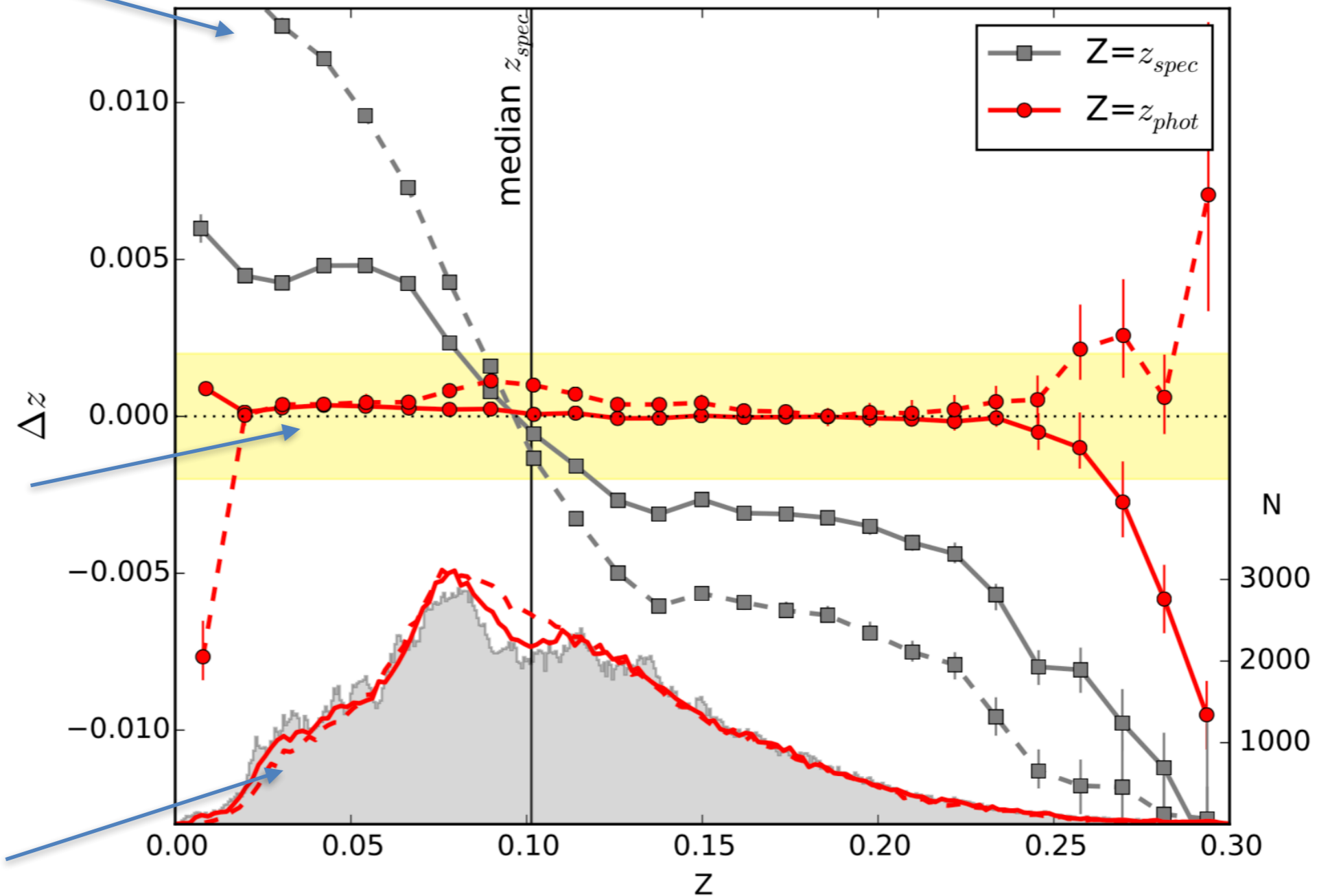
One form of bias: mean redshift residuals as a function of spec-z or photo-z

$$\Delta z = \frac{z_{photo} - z_{spec}}{1 + z_{spec}}$$

w.r.t spec-z
(biased)

w.r.t photo-z
(not guaranteed
to be unbiased)

Unbalanced
Distribution



SDSS $z < 0.4$
(Pasquet et al. 2019)

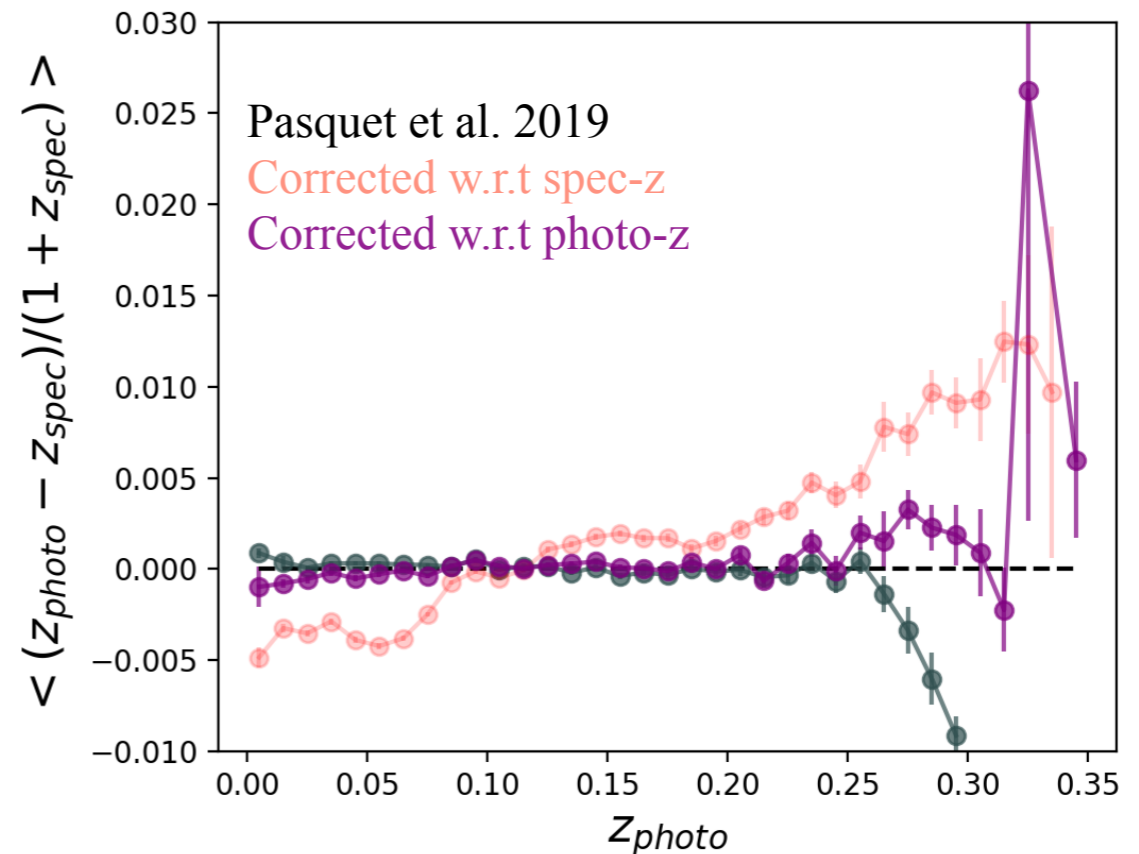
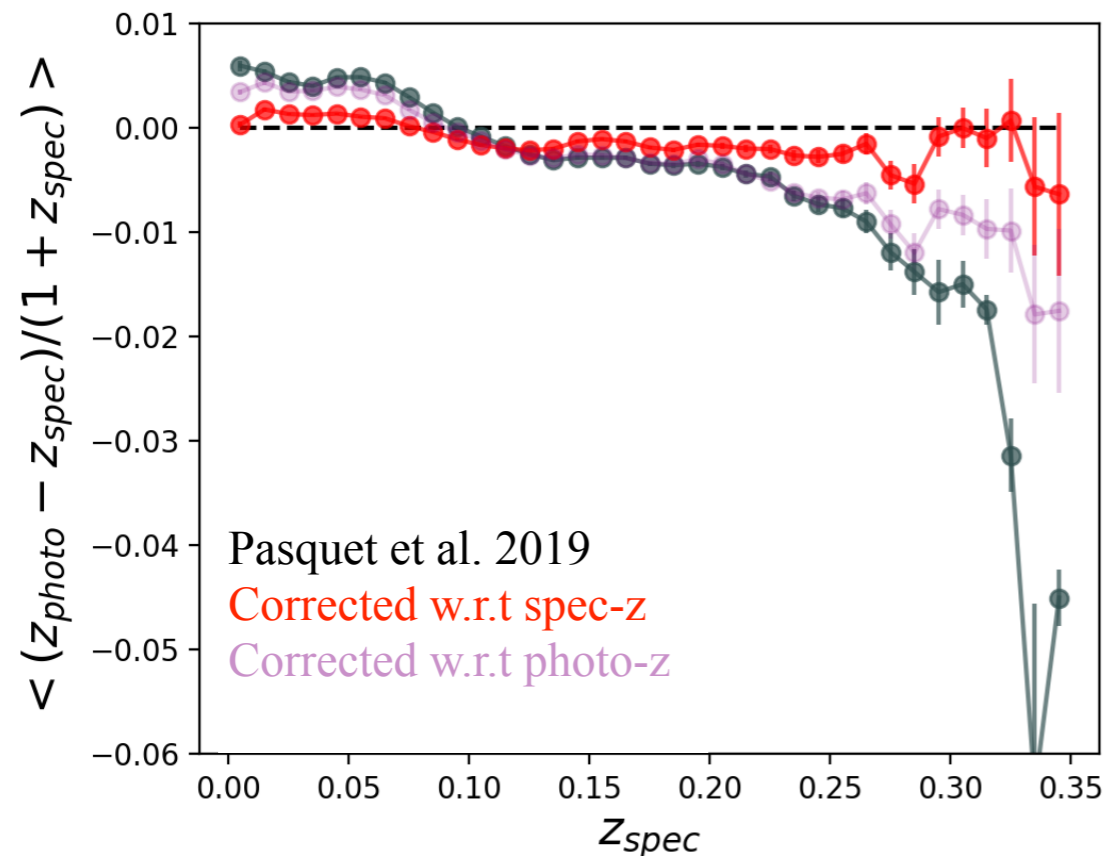
Bias correction via splitting representation and estimation (Lin et al. 2022)

Representation Learning (all data)

Estimation (a near-balanced subset)

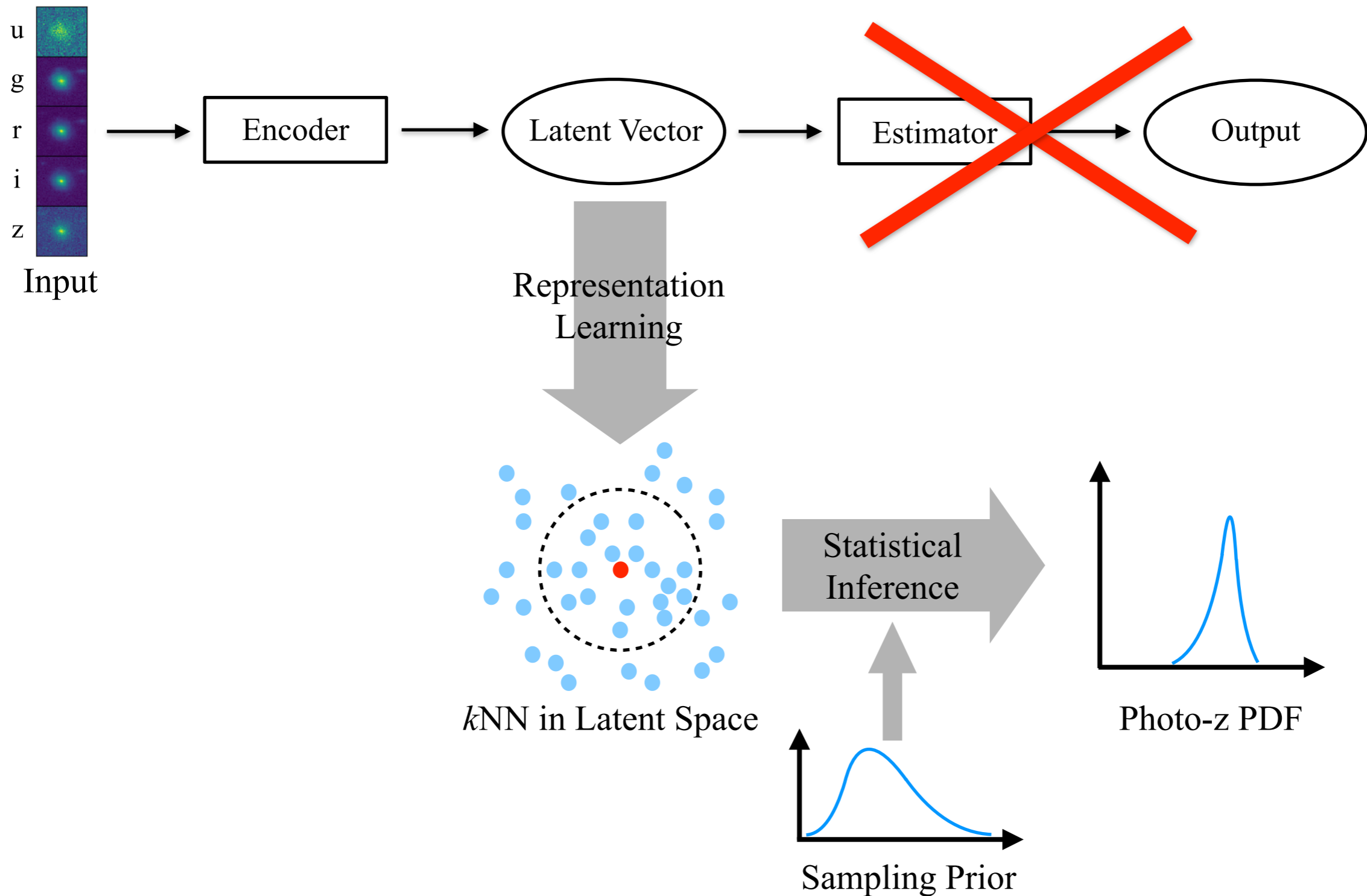


- Treat spectroscopic & photometric spaces separately:



Current work: empower deep learning with statistical basis

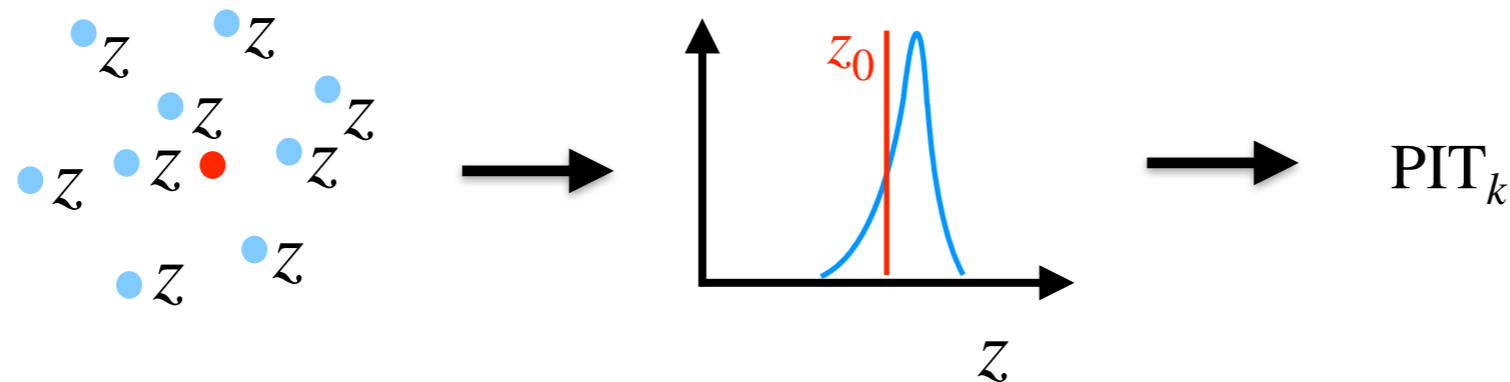
- Representation learning + statistical inference (or k NN)



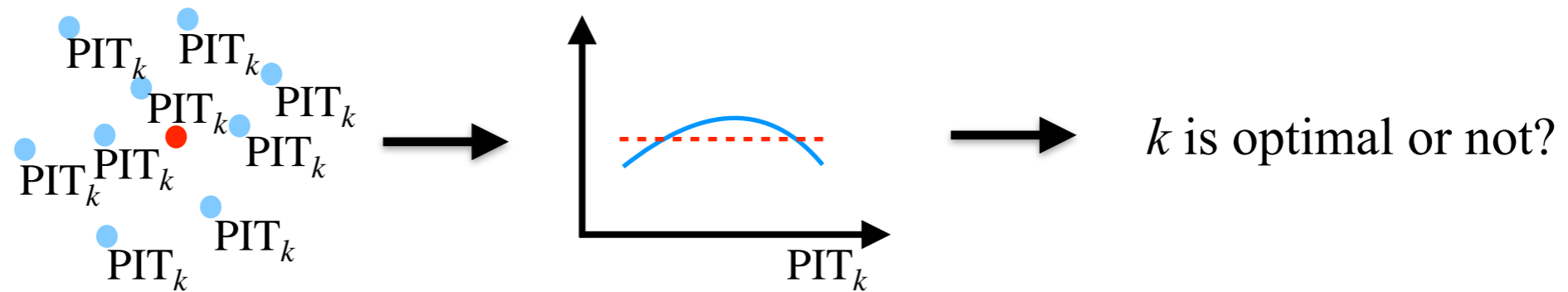
Determine the optimal k via local Probability Integral Transform (PIT) diagnostics

- Select k from [5, 10, 15, ..., 2000]
- For each labeled galaxy:

$$\text{PIT}(z_{\text{spec}}) = \int_0^{z_{\text{spec}}} p(z) dz$$



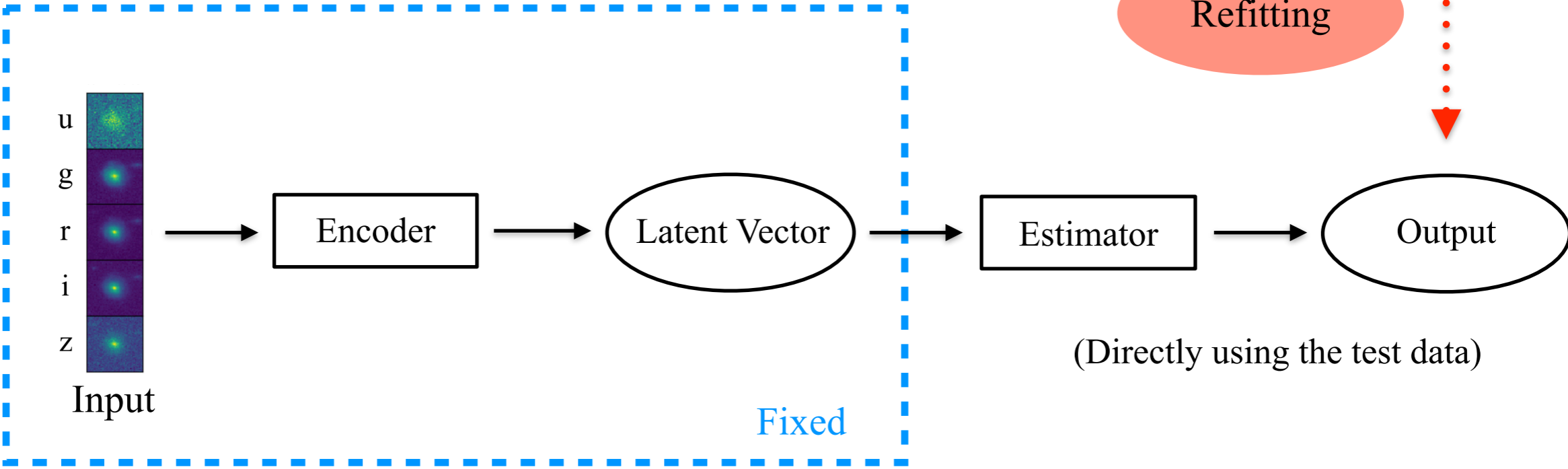
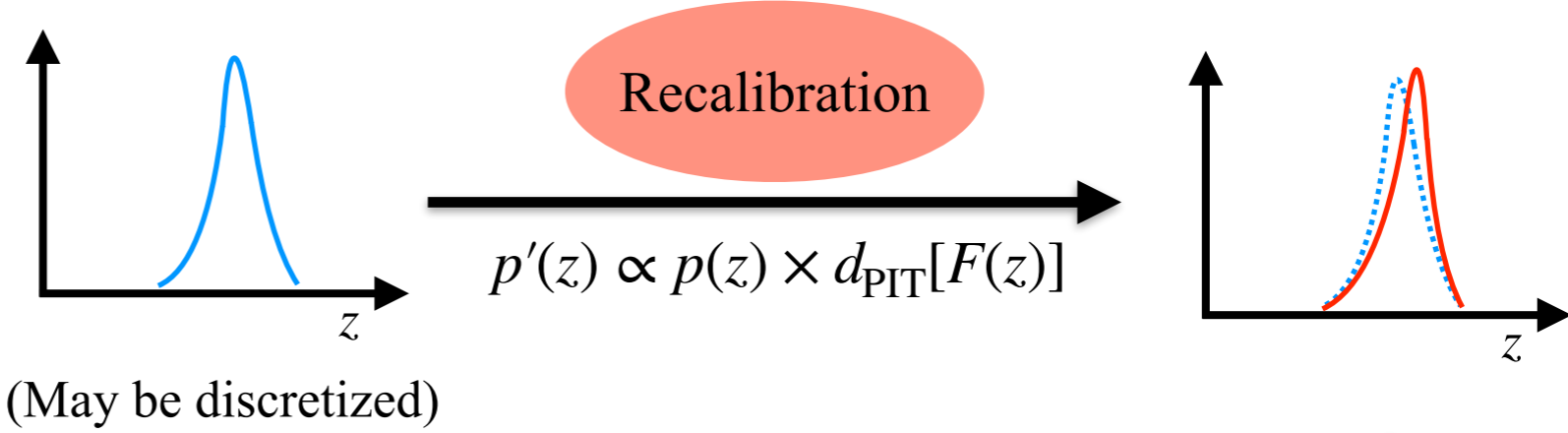
- For each query (unlabeled) galaxy:



- k is optimal when the PIT_k distribution is closest to a **uniform** distribution

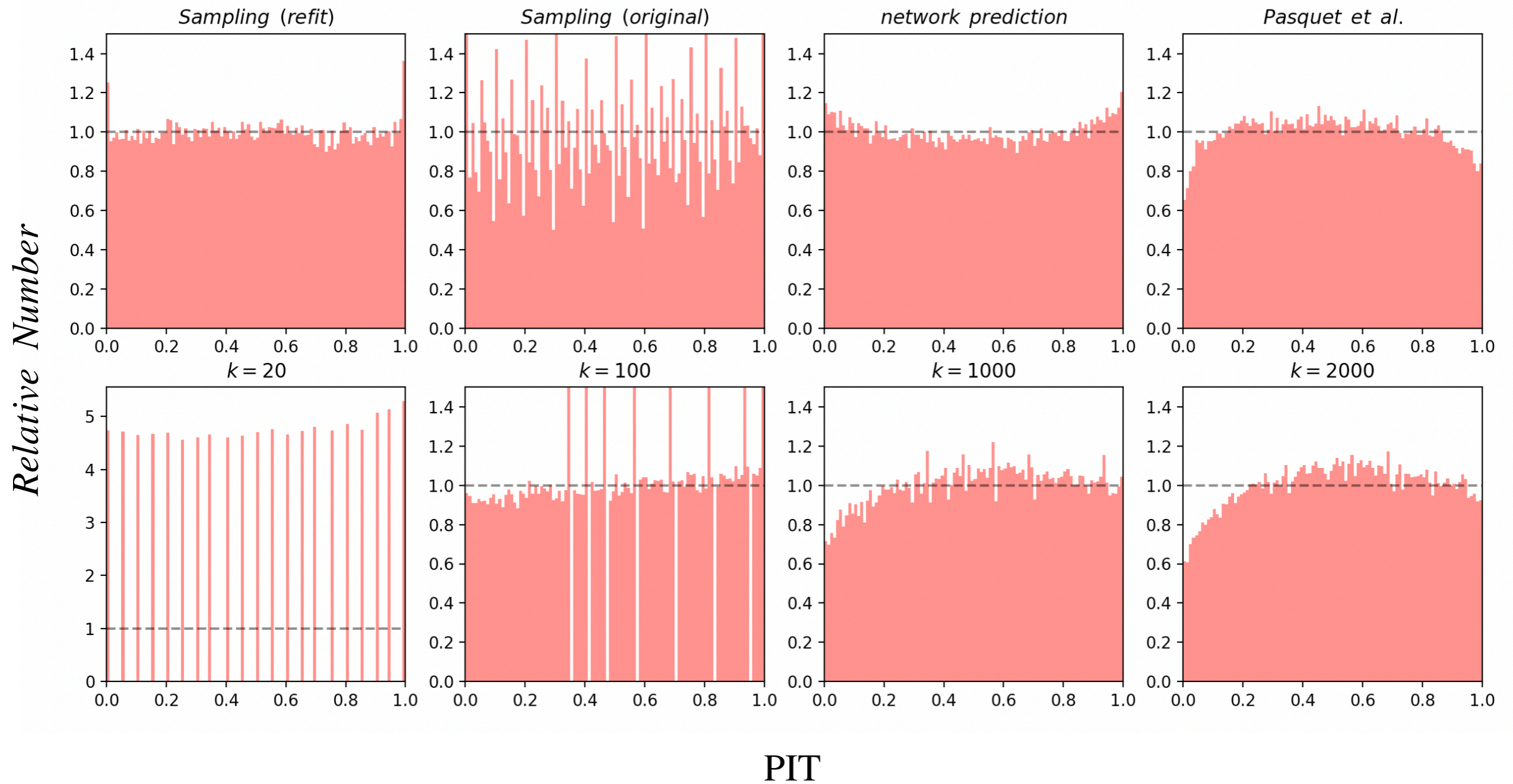
Recalibration + Refitting

- Discretized \longrightarrow Recalibration
- Non-uniform \longrightarrow Refitting



Results: PITs

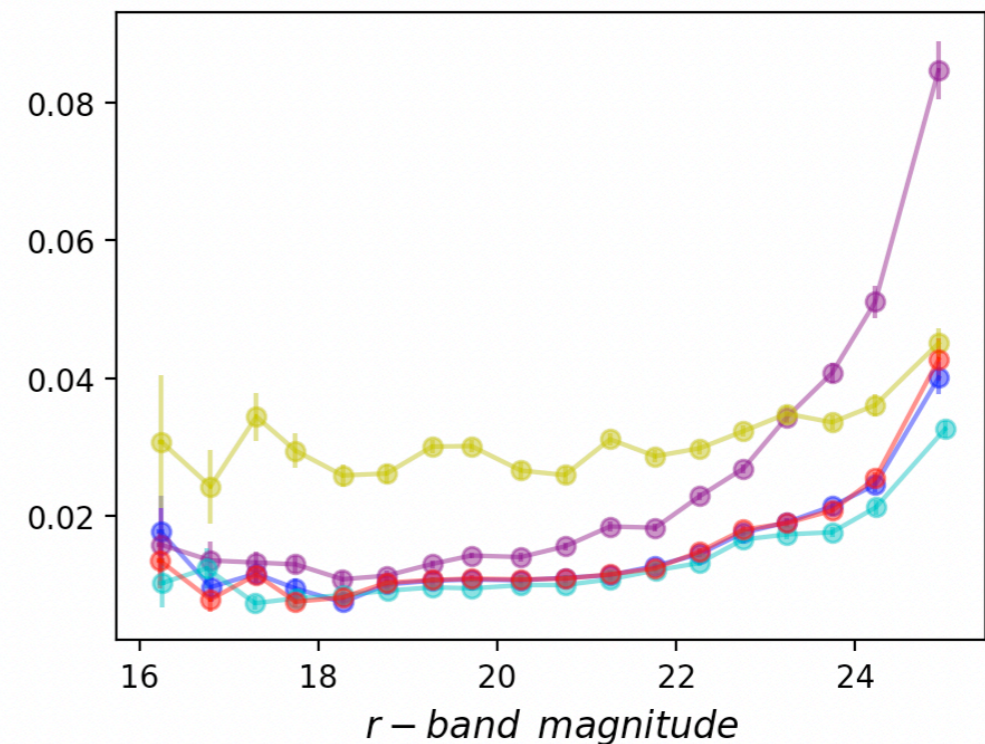
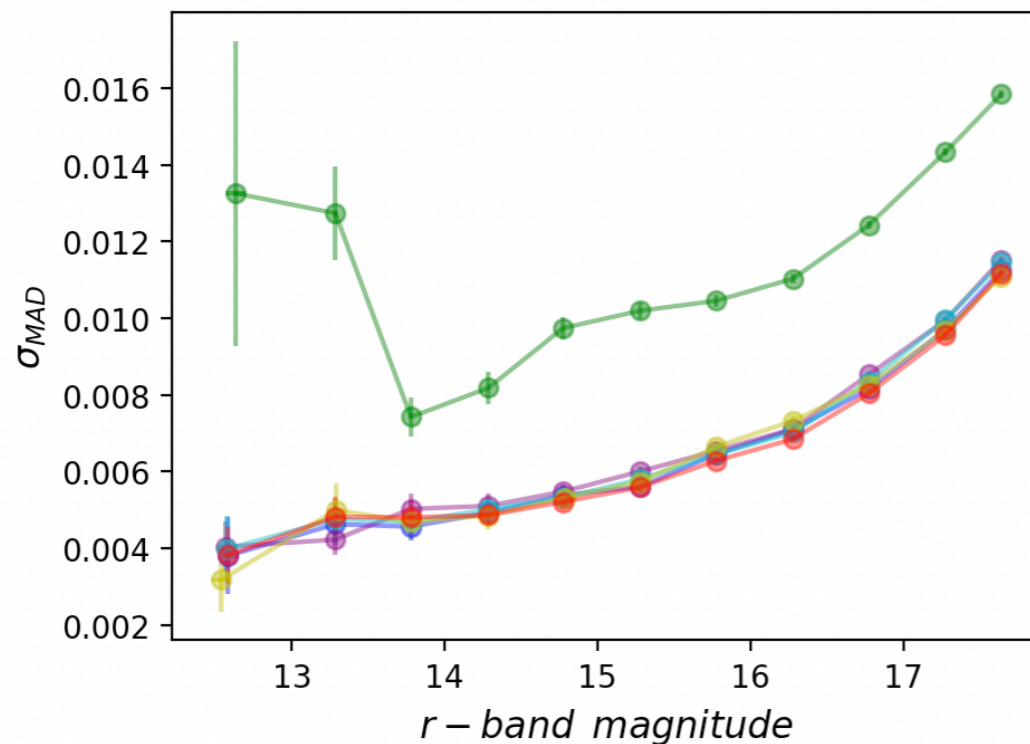
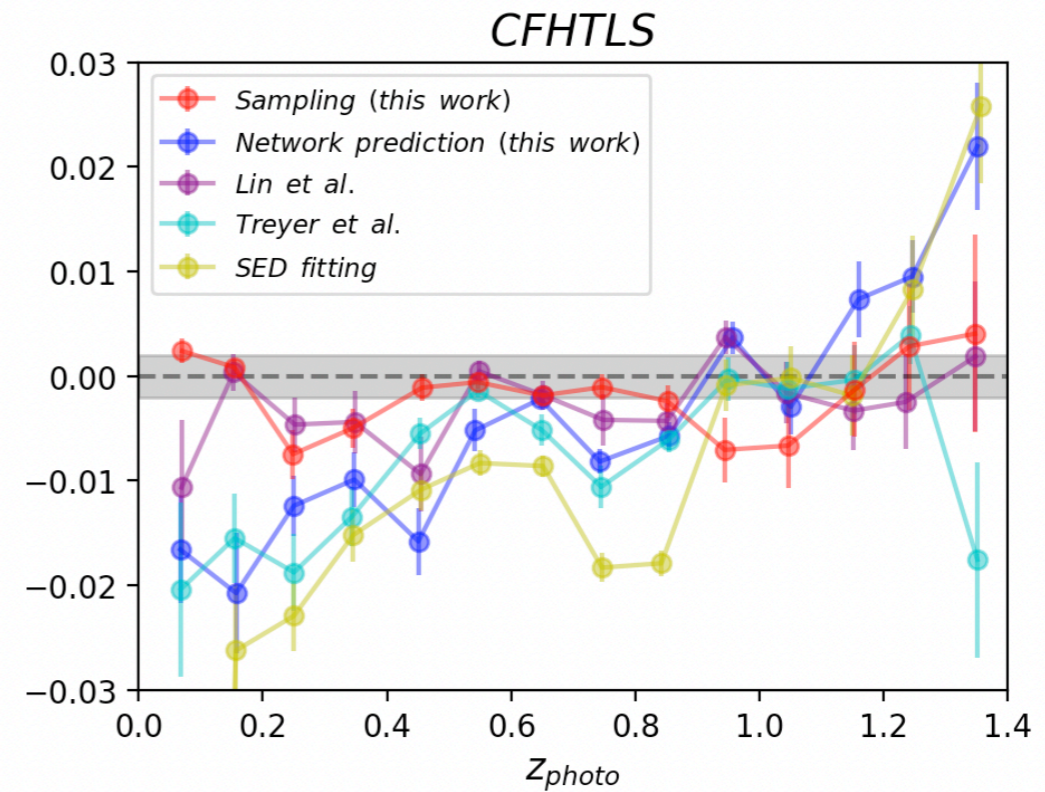
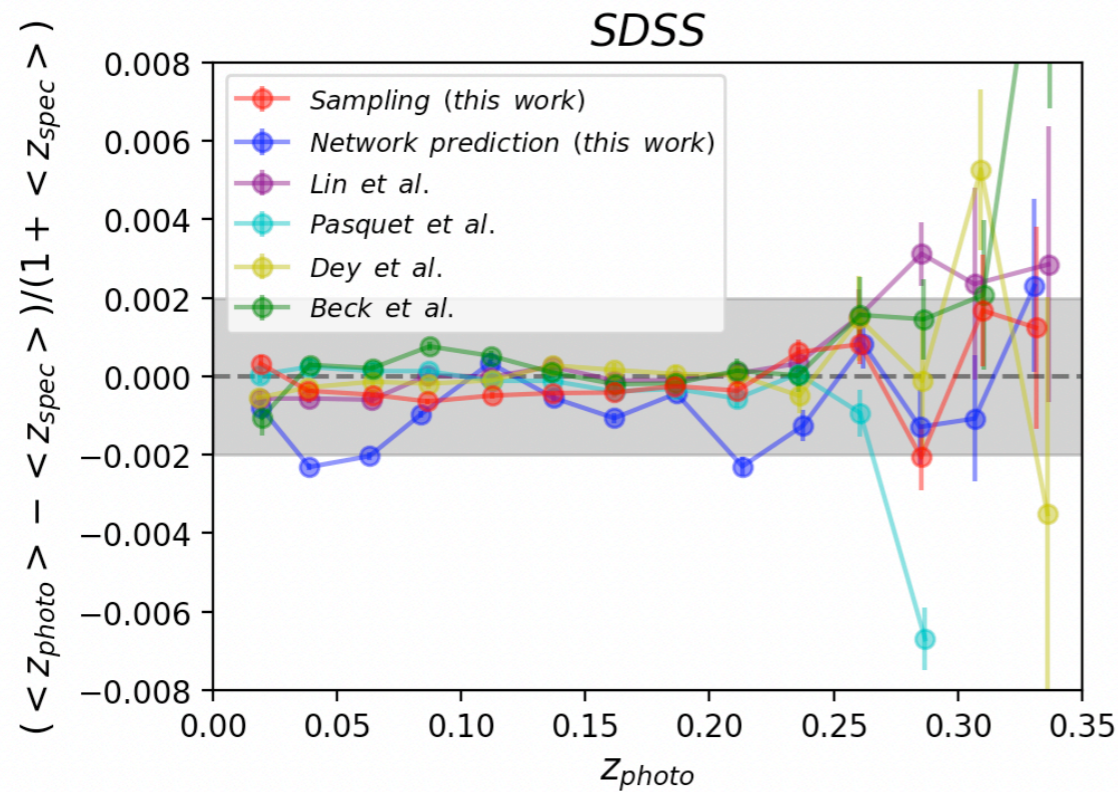
- Good PDF calibration achieved by sampling/inference (shown for the SDSS data)



Results: point estimates

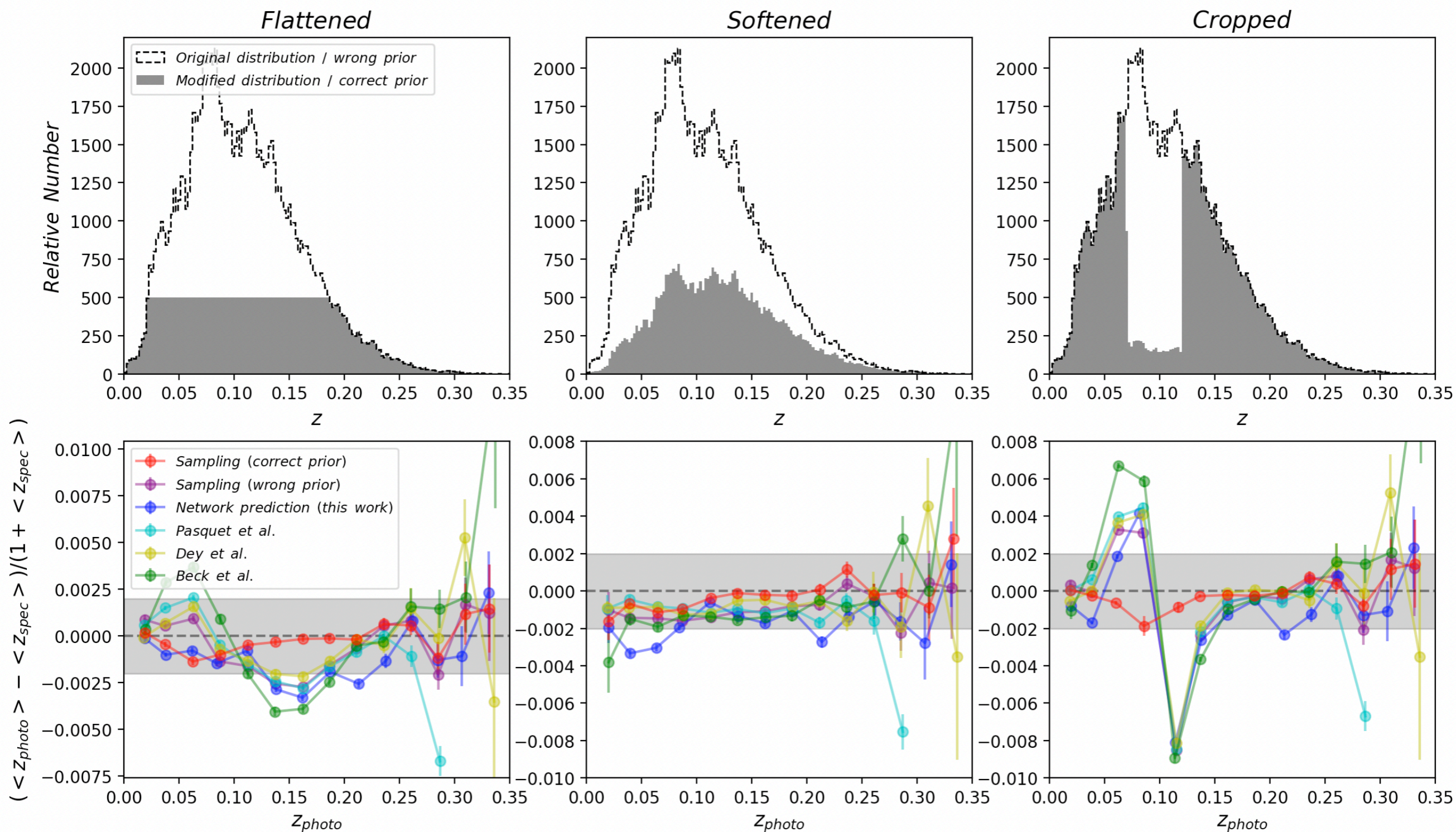
- Mean redshift bias correction achieved by sampling/inference (similar to Lin et al.)
- No loss in accuracy (contrary to Lin et al.)

$$z_{photo} = \int_0^{z_{max}} z \times p(z) dz$$



Results: the impact of distribution mismatch

- Robustness under distribution mismatch with correct sampling prior

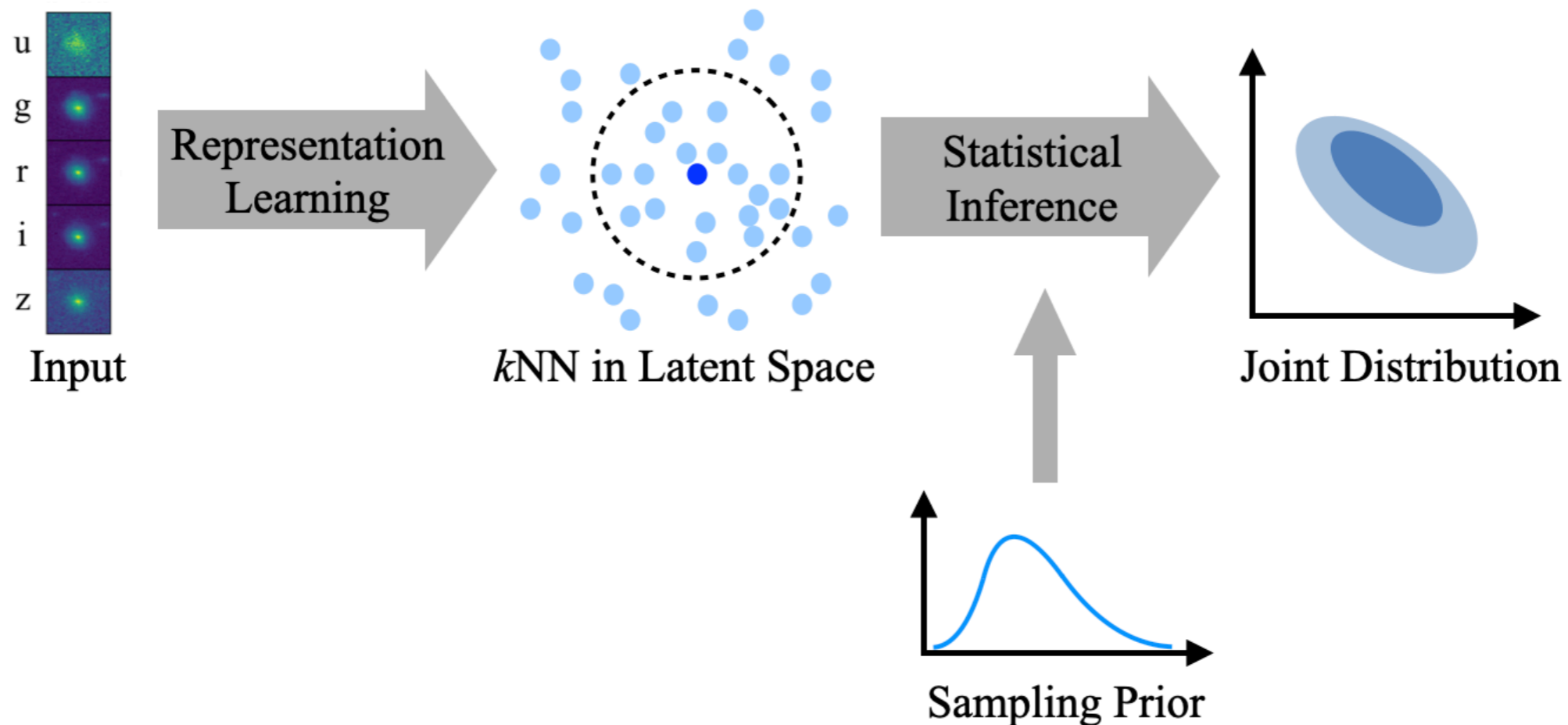


Summary

- Key idea: combine deep learning and statistical basis
- Representation learning, statistical inference, recalibration & refitting
- Better results over benchmark methods:
 - Well-calibrated PDFs
 - Good control of photo-z-dependent residuals without compromising accuracy
 - Robustness under distribution mismatch

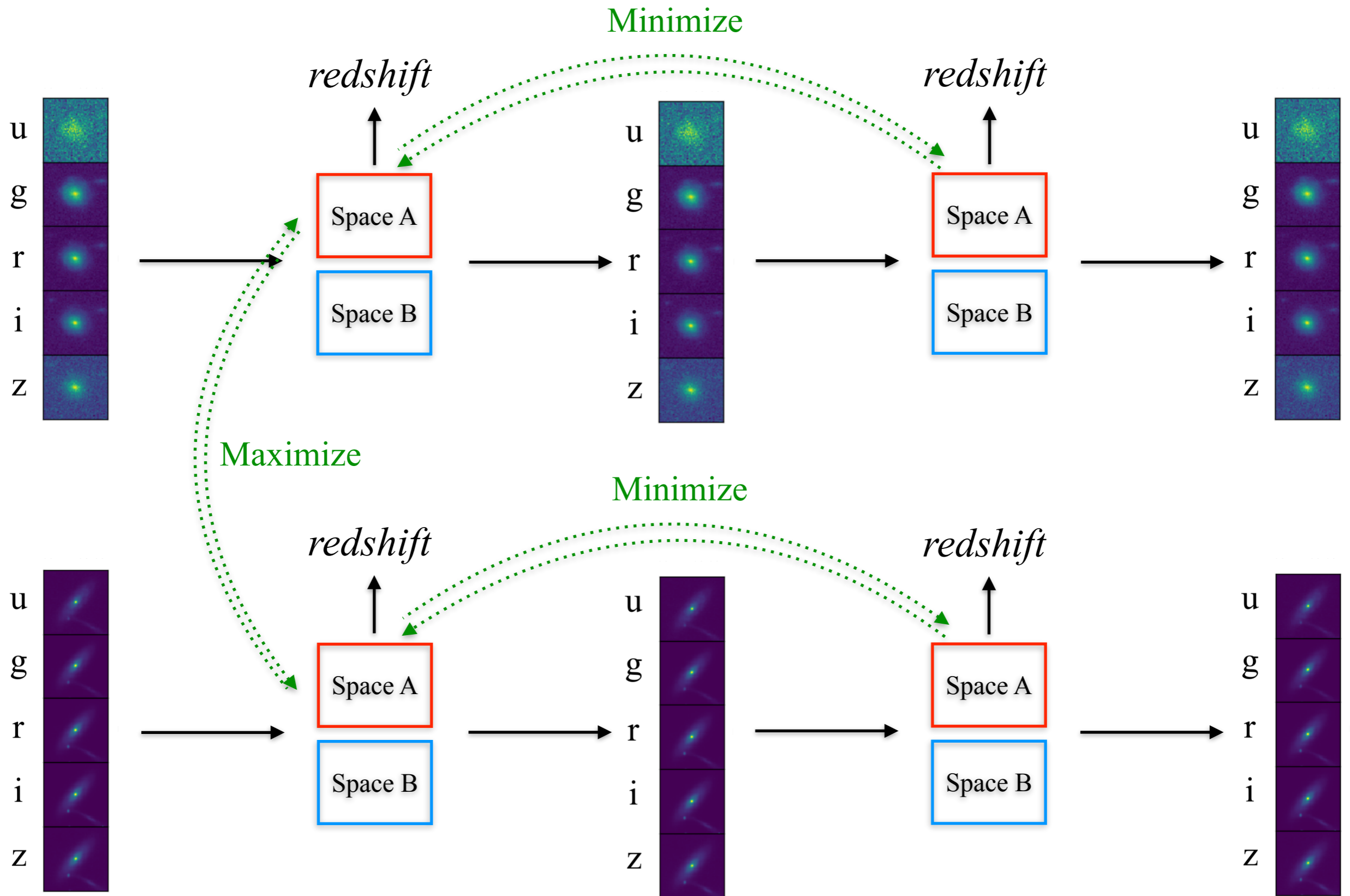
For interpretability: analyze *redshift-variable* correlations

- Information/variables to be exploited for reducing redshift residuals (e.g., galaxy structures, environmental properties, etc.)
- Relations between variables

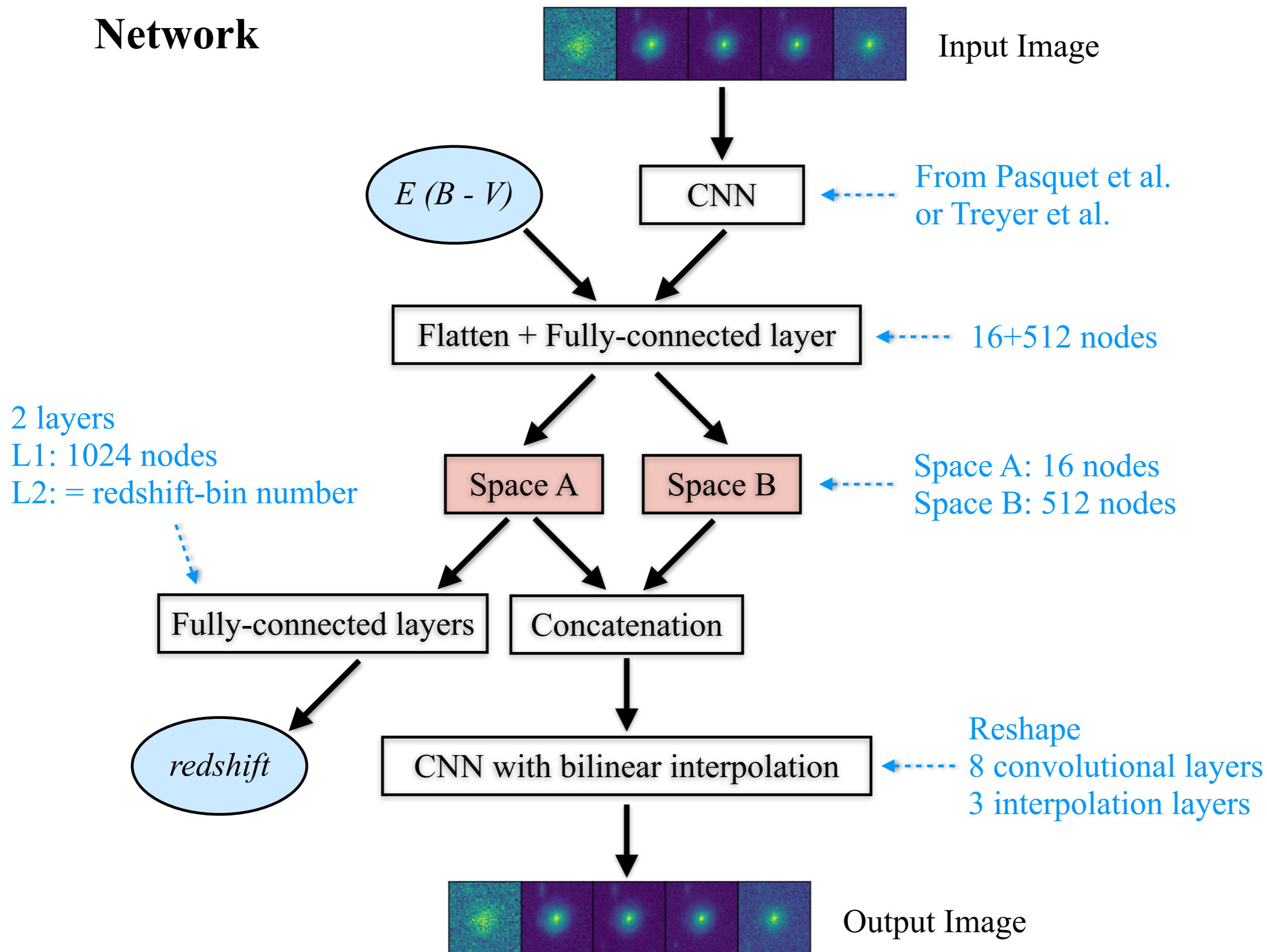


Back-up slides

Representation learning



Network



Optimal metric for PIT diagnostics: Wasserstein distance

$$D_{Wasserstein}[P_{PIT_k}, P_{uniform}] = \int_0^1 |F_{PIT_k}(p) - F_{uniform}(p)| dp$$

$$D_{CrossEntropy}[P_{PIT_k}, P_{uniform}] = \int_0^1 [P_{uniform}(p) \log P_{PIT_k}(p) + (1 - P_{uniform}(p)) \log(1 - P_{PIT_k}(p))] dp$$

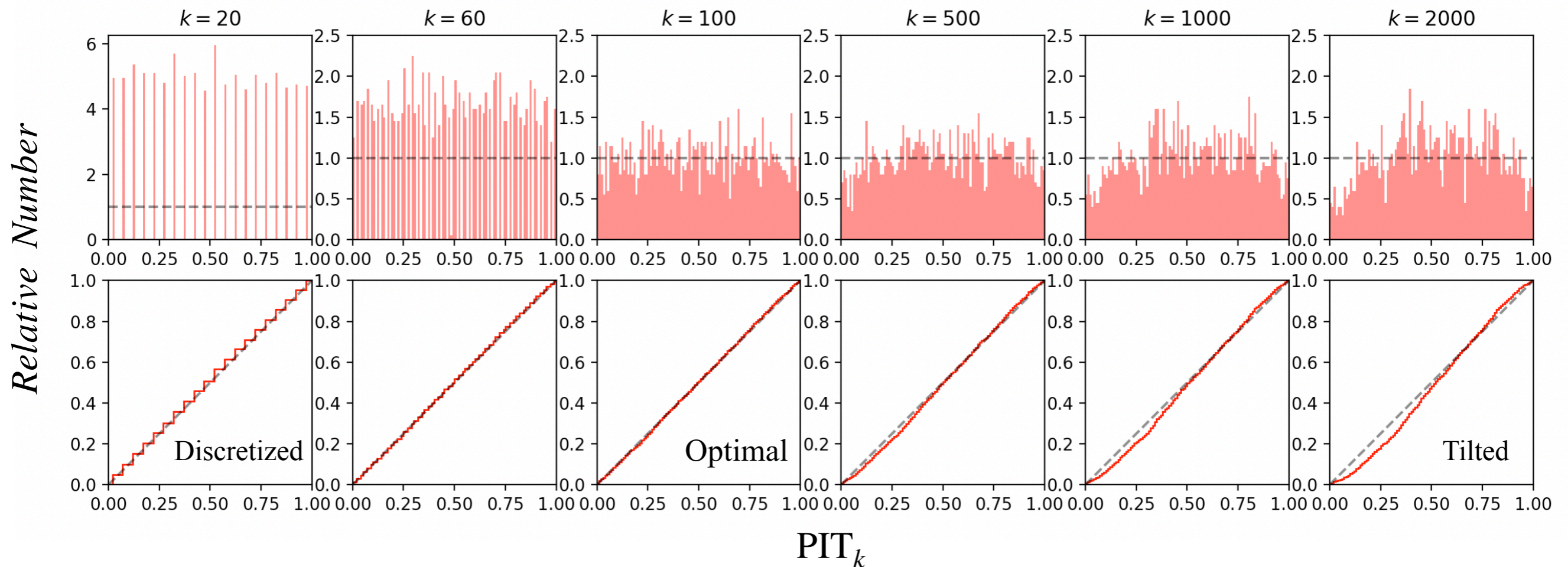
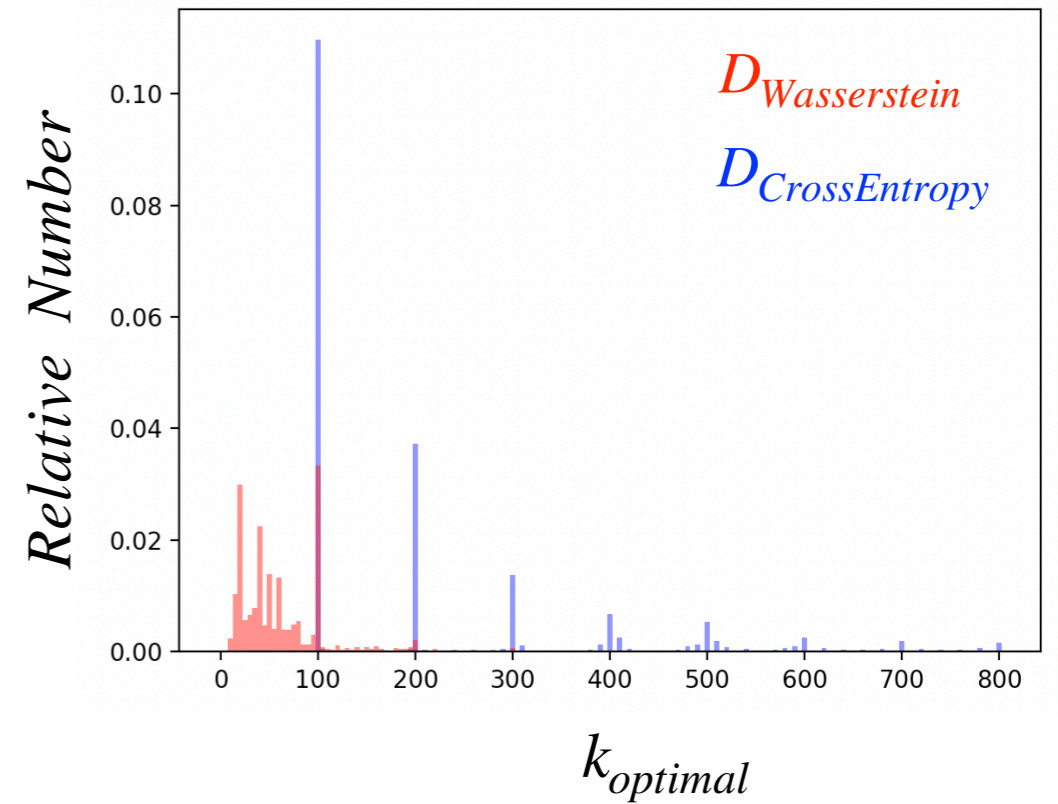
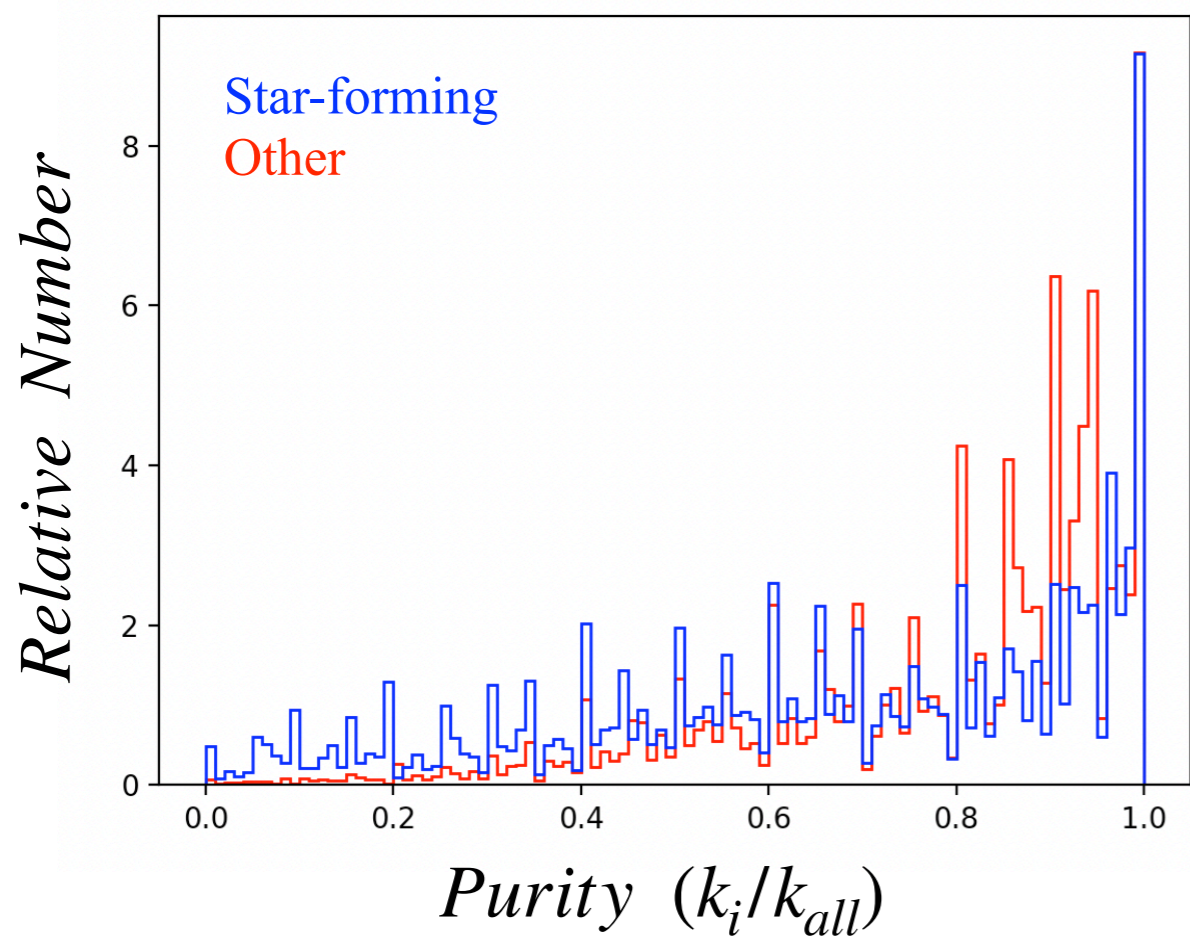
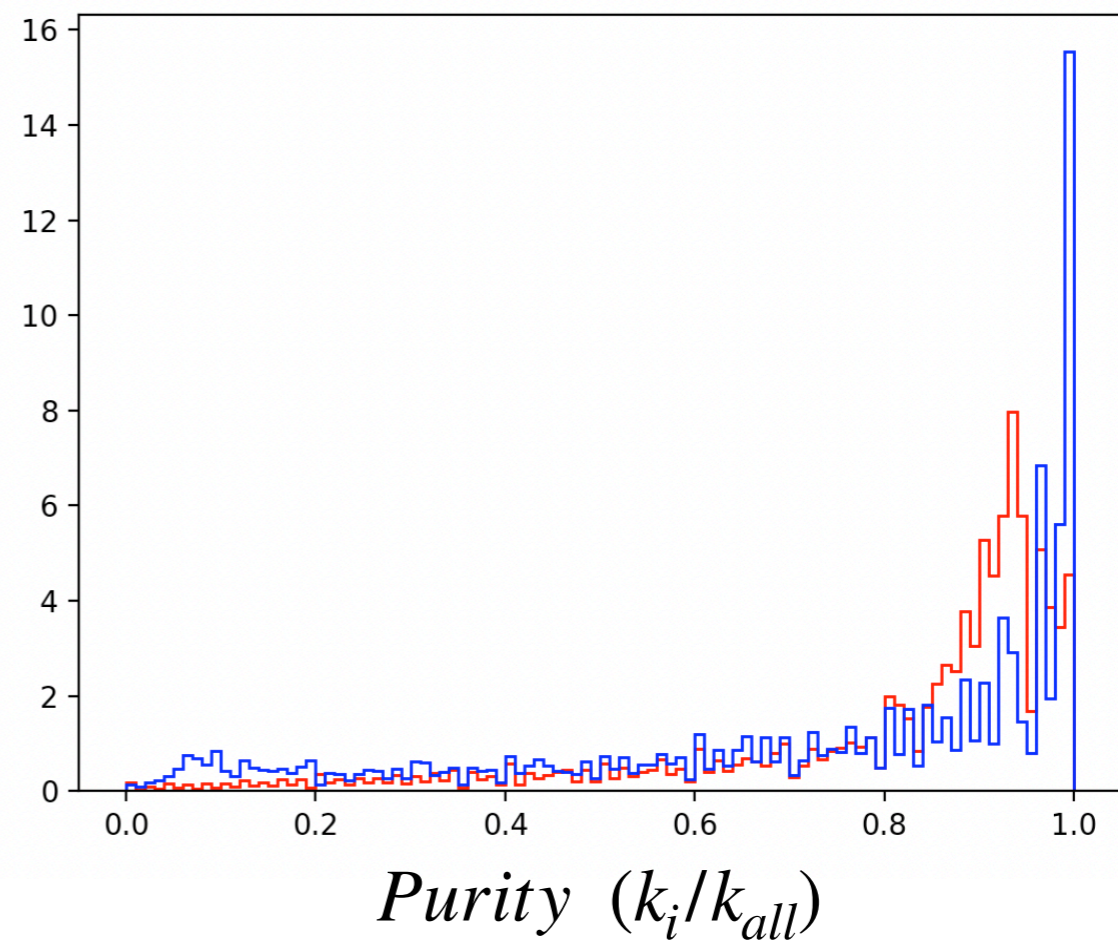


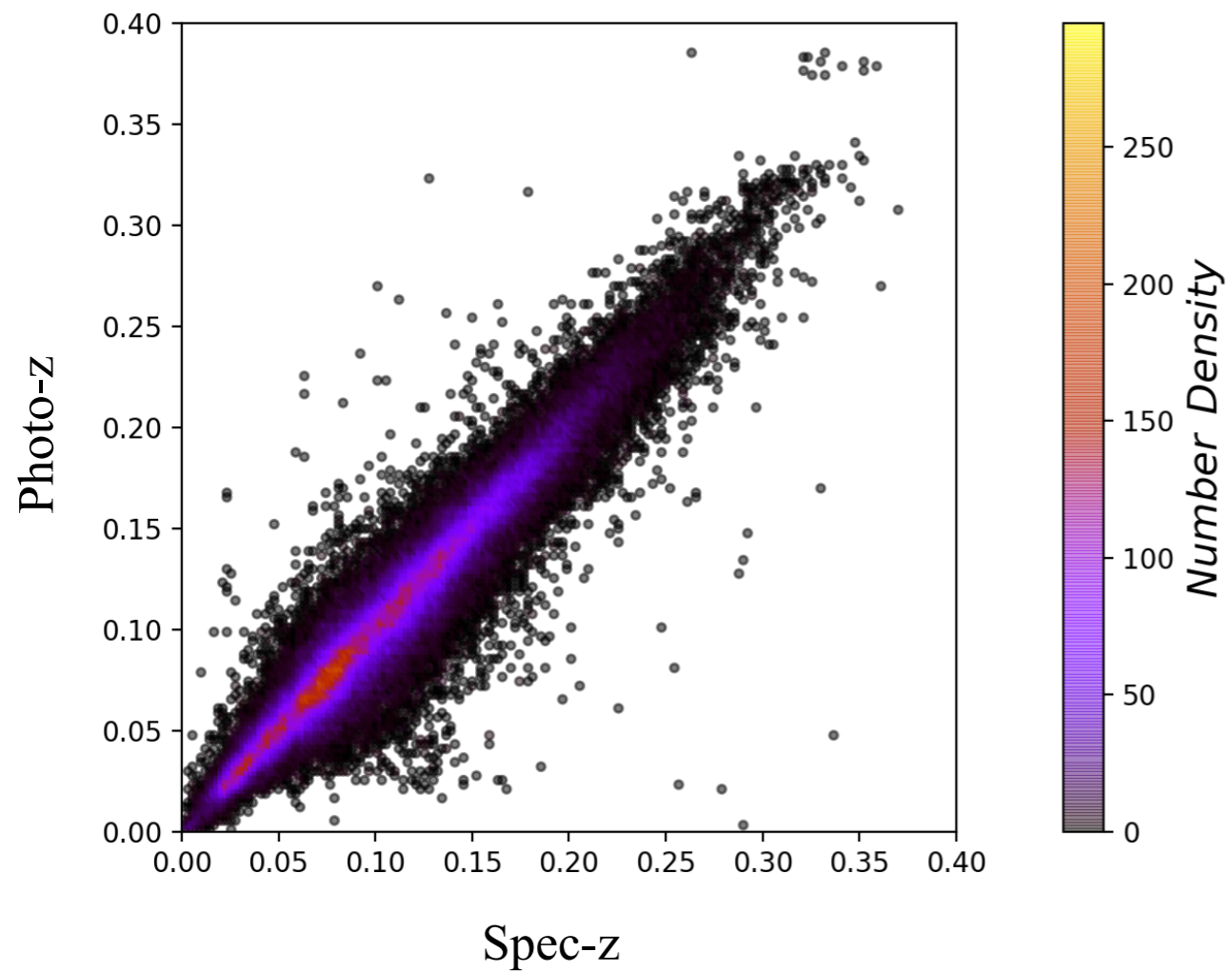
Image-based



Photometry-based



SDSS



CFHTLS

