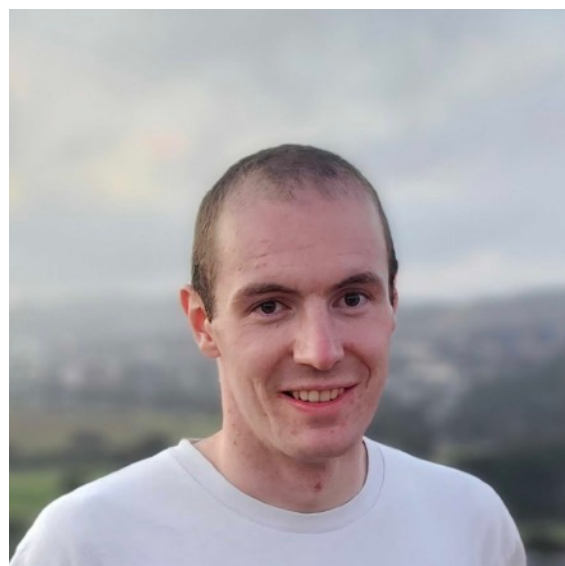


pop-cosmos:

Scaleable Bayesian inference of galaxy properties under a diffusion model prior



Stephen Thorp



Sinan Deger



Anik Halder



Gurjeet Jagwani



Hiranya Peiris



Madalina Tudorache



Justin Alsing



Boris Leistedt



Joel Leja



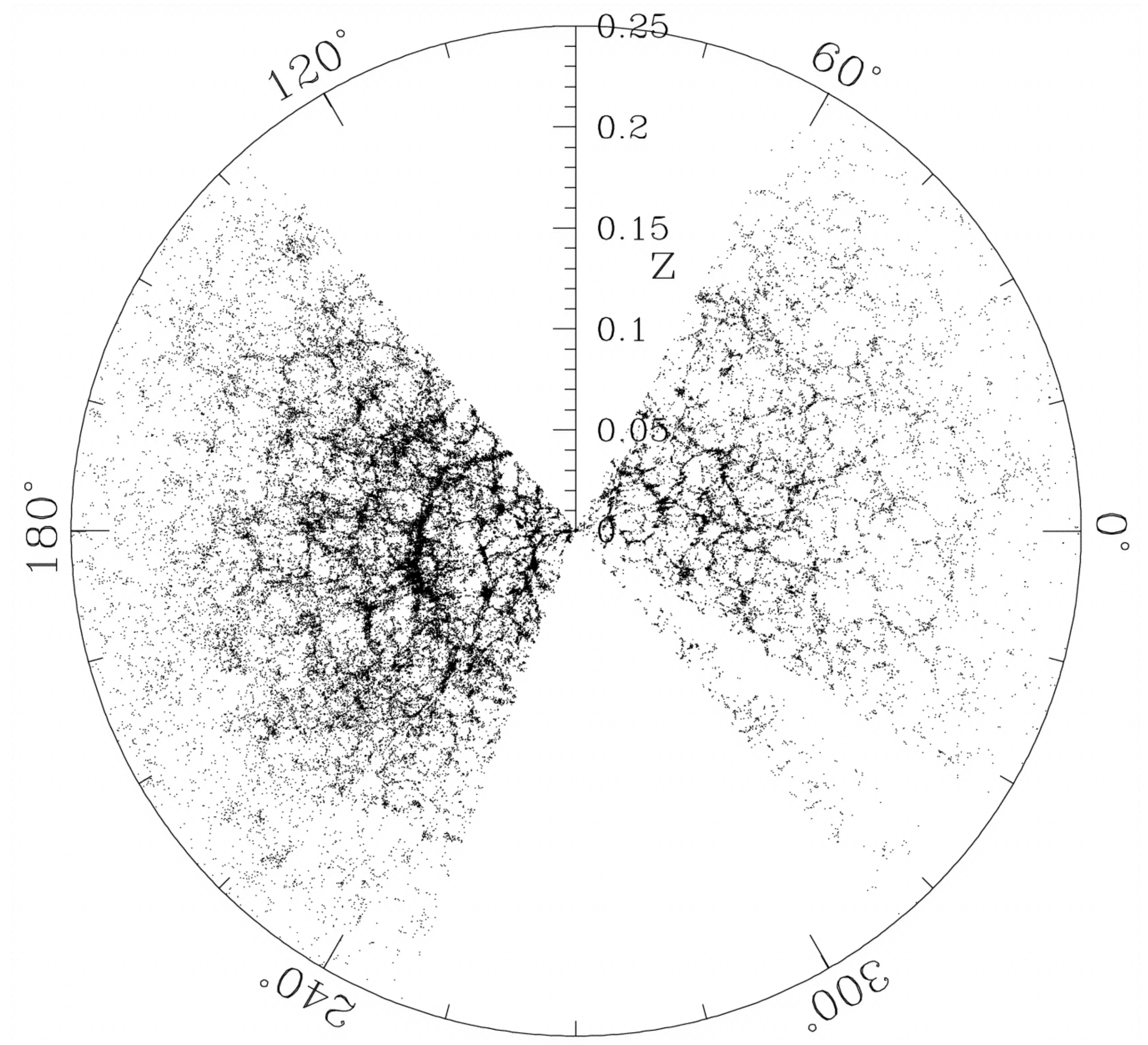
Arthur Loureiro



Daniel Mortlock

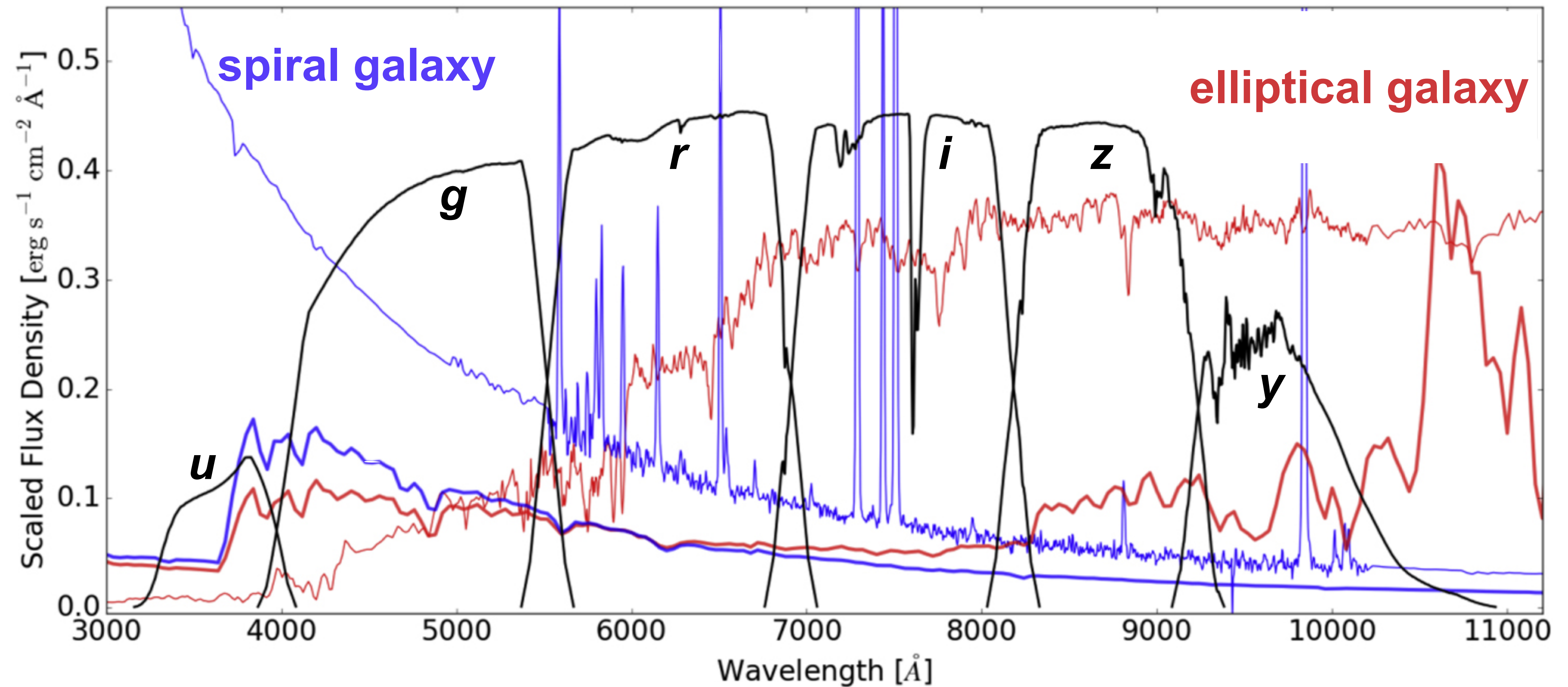
Large Scale Structure Cosmology

- Map the matter distribution in the universe using spectra or photometry of galaxies
- Galaxy clustering (positions), or weak lensing (shapes) sensitive to cosmological parameters
- Photometric surveys depend more on our ability to model galaxies well (either individually or as a population)



SDSS: *Blanton et al. (2003)*

Photometry vs. Spectroscopy



Coming (very) soon...

- First-look events June 23rd:
<https://rubinobservatory.org>
- First data June 30th!
- 20 billion galaxies; 18,000 deg²
- Deep imaging in *ugrizy*
- Single epoch: $r \lesssim 24$ mag;
10 year co-add: $r < 26.9$ mag
- Photometric data only; need to
model galaxies with very limited
information (just 6 numbers)



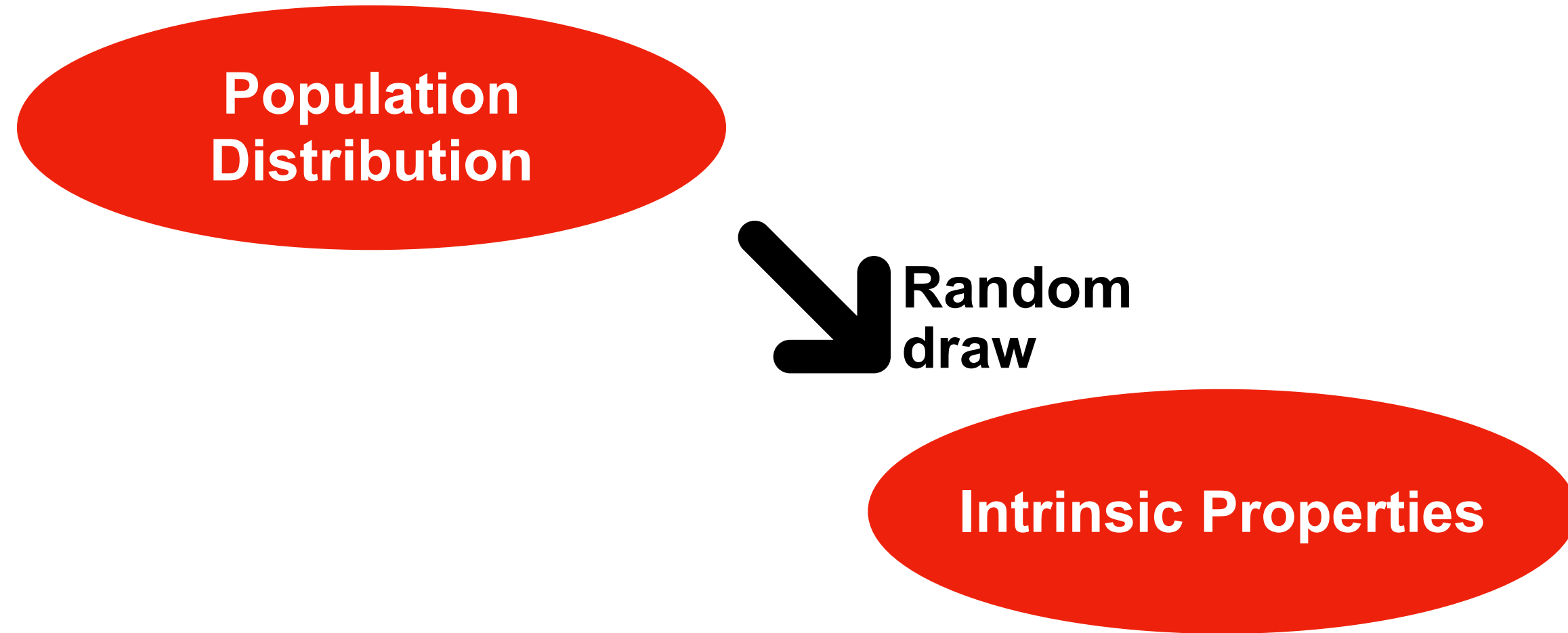
Credit: Rubin Observatory/NOIRLab/SLAC/DOE/NSF/AURA/B. Quint

Setting up a generative model

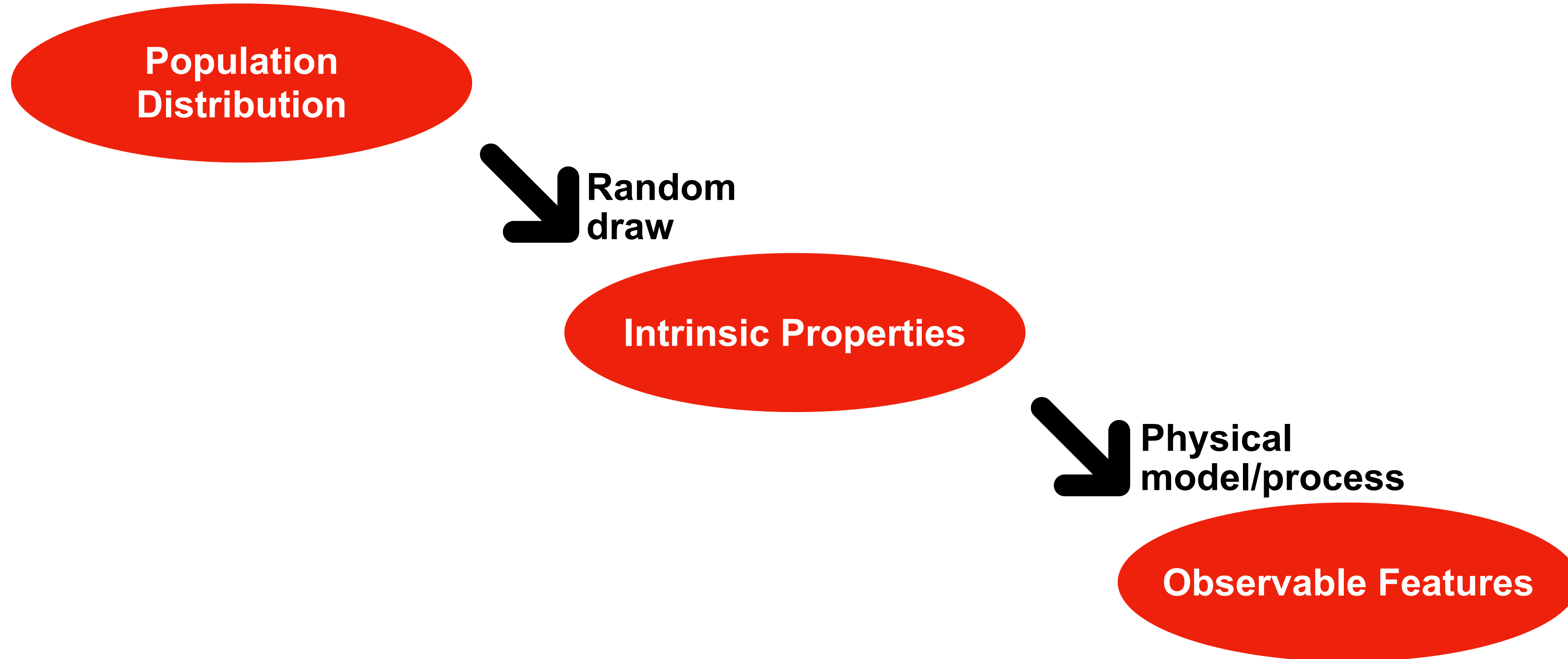


**Population
Distribution**

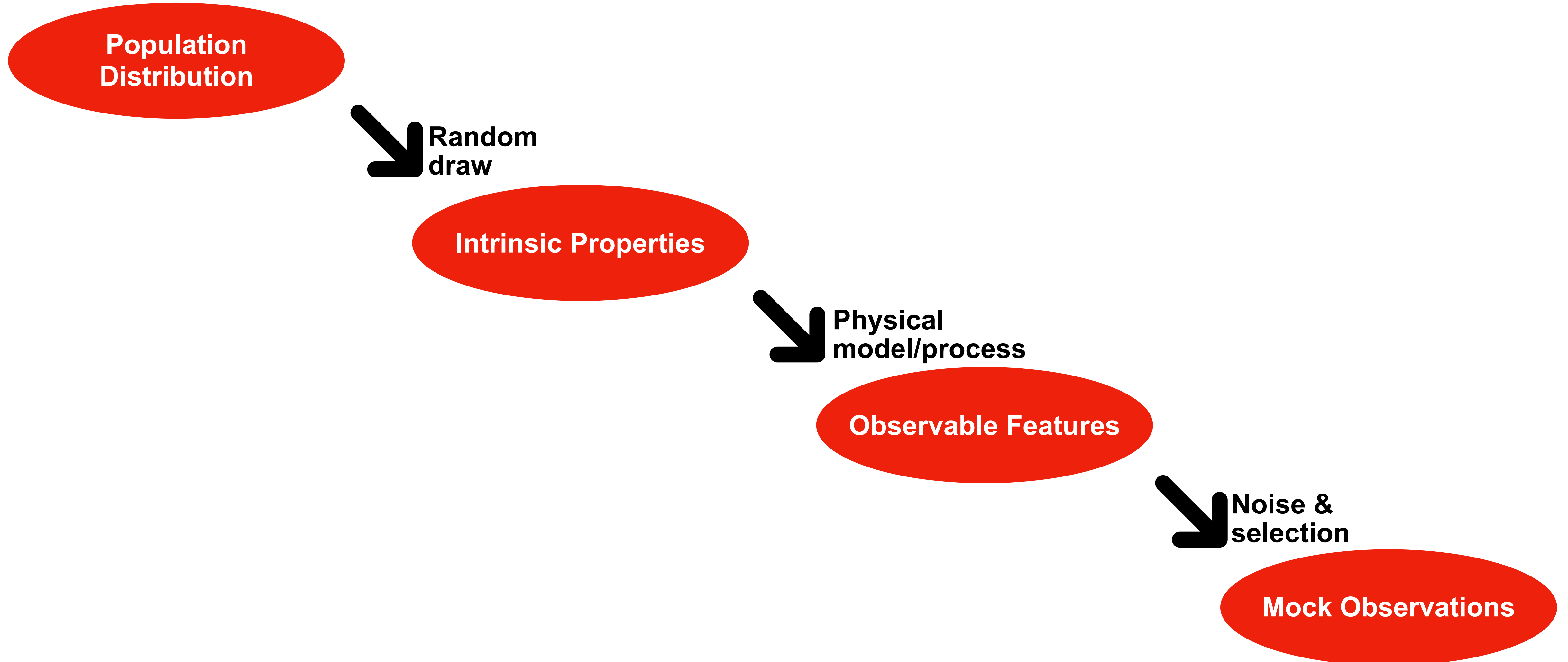
Setting up a generative model



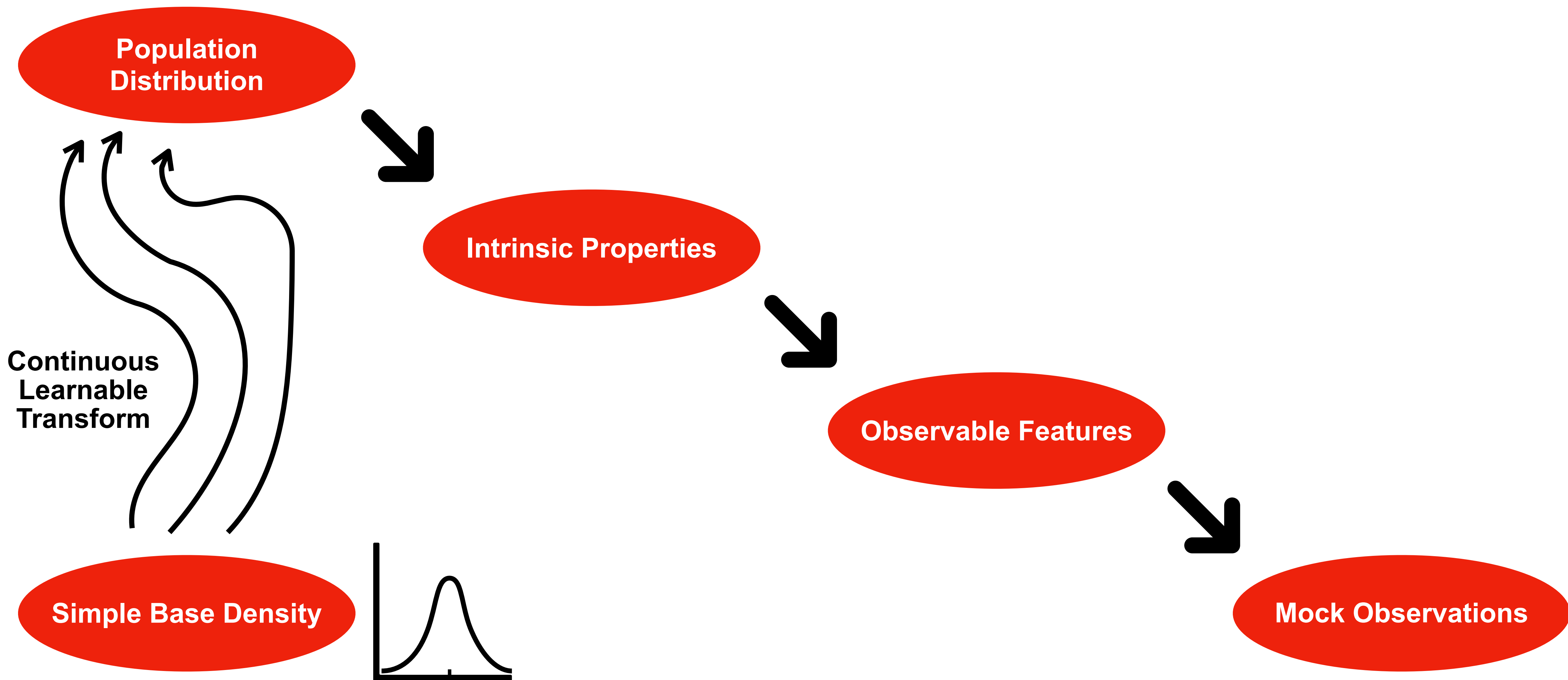
Setting up a generative model



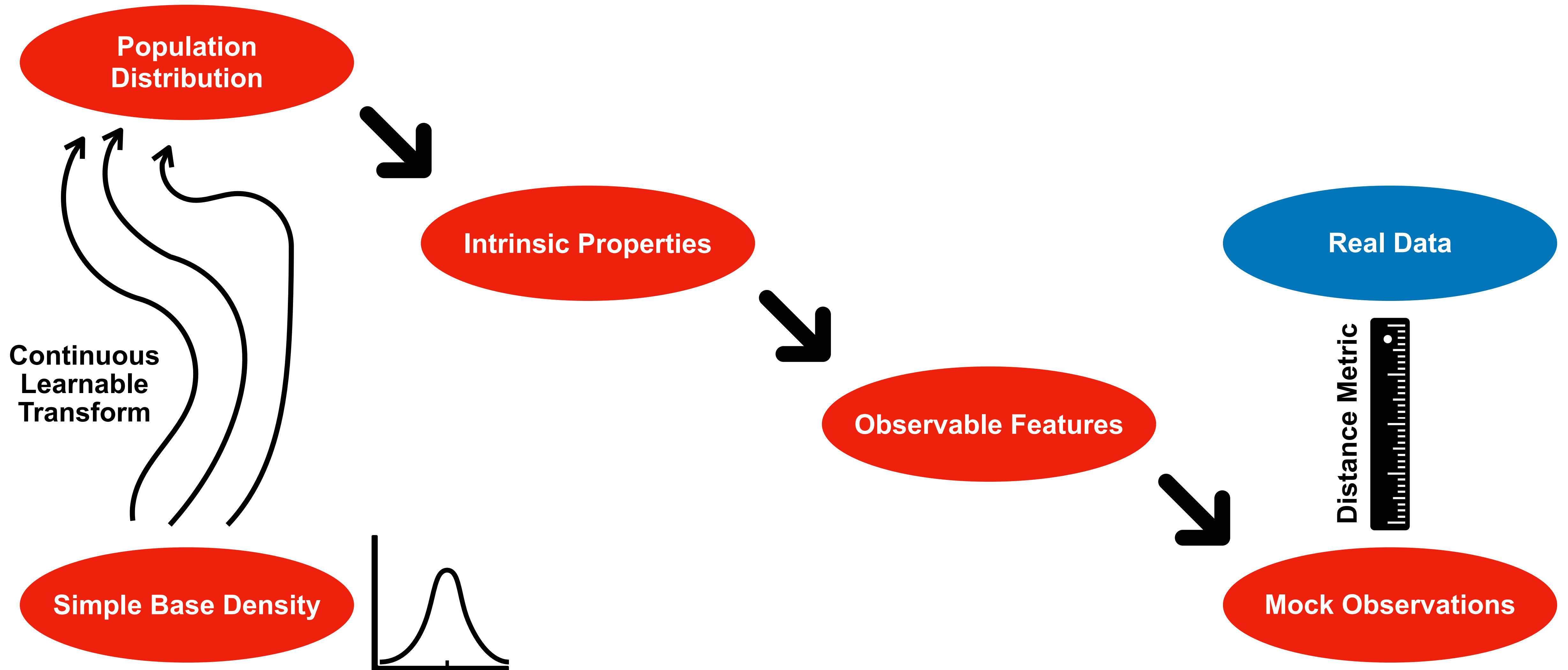
Setting up a generative model



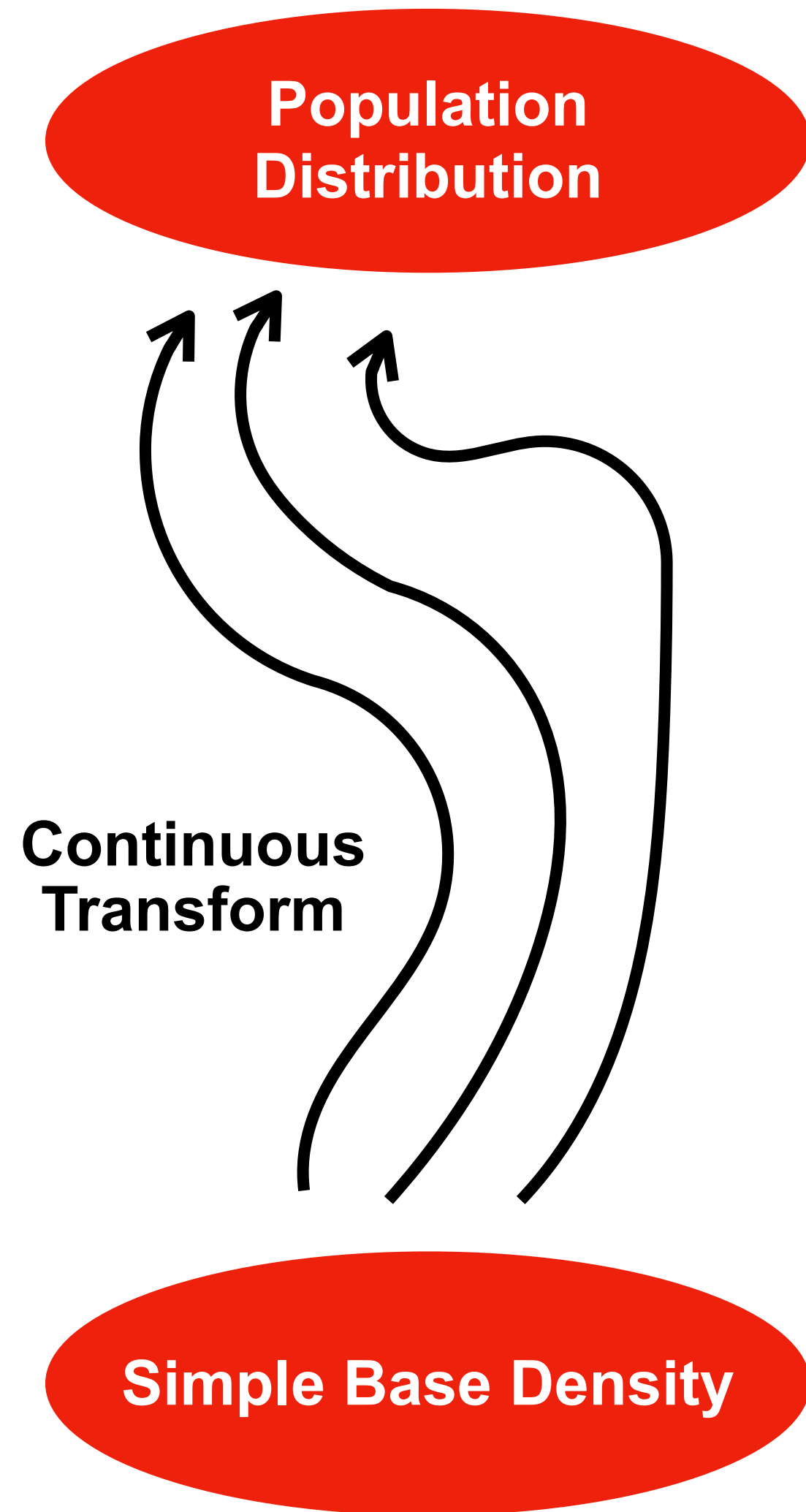
Diffusion models



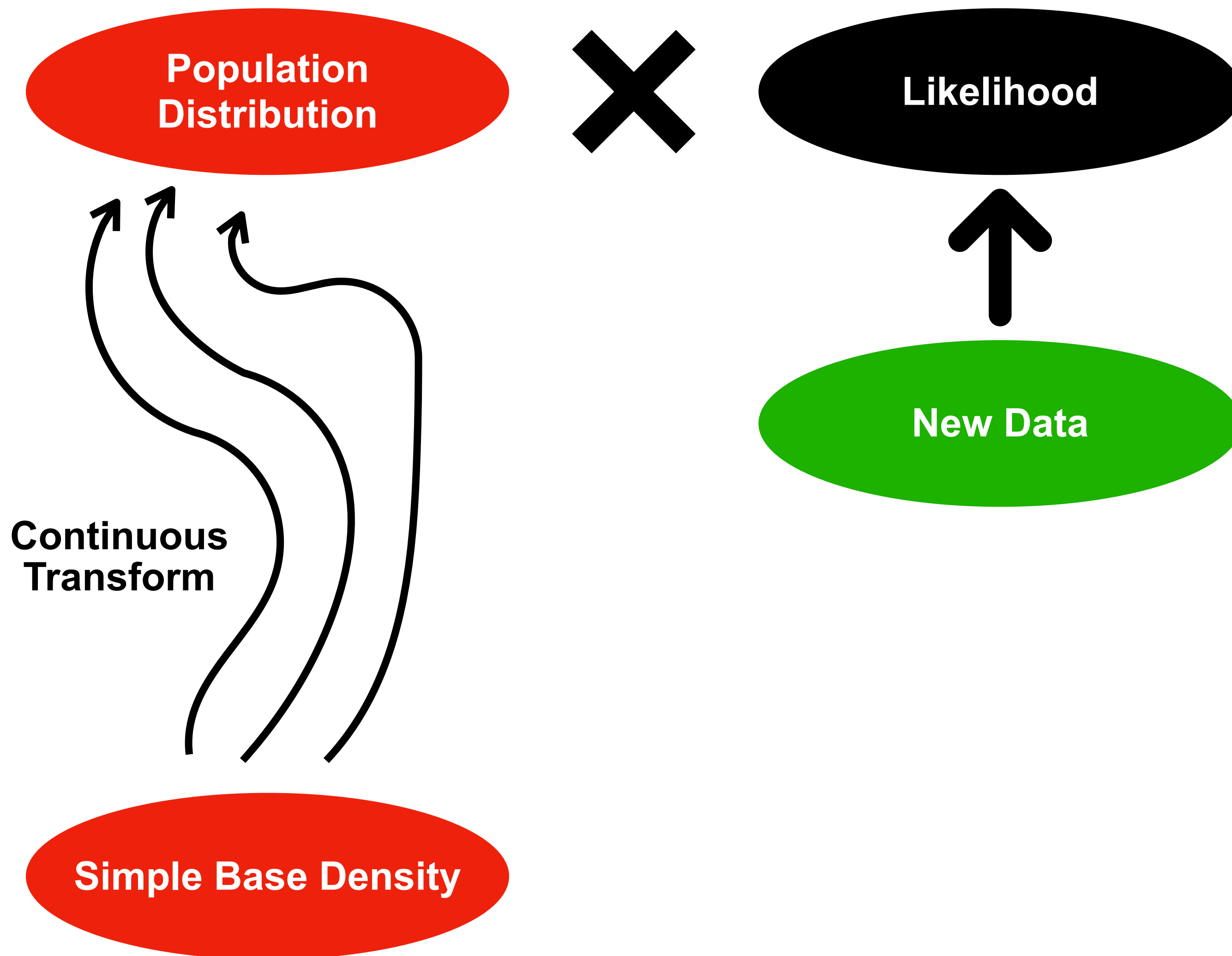
Fitting a generative model to data



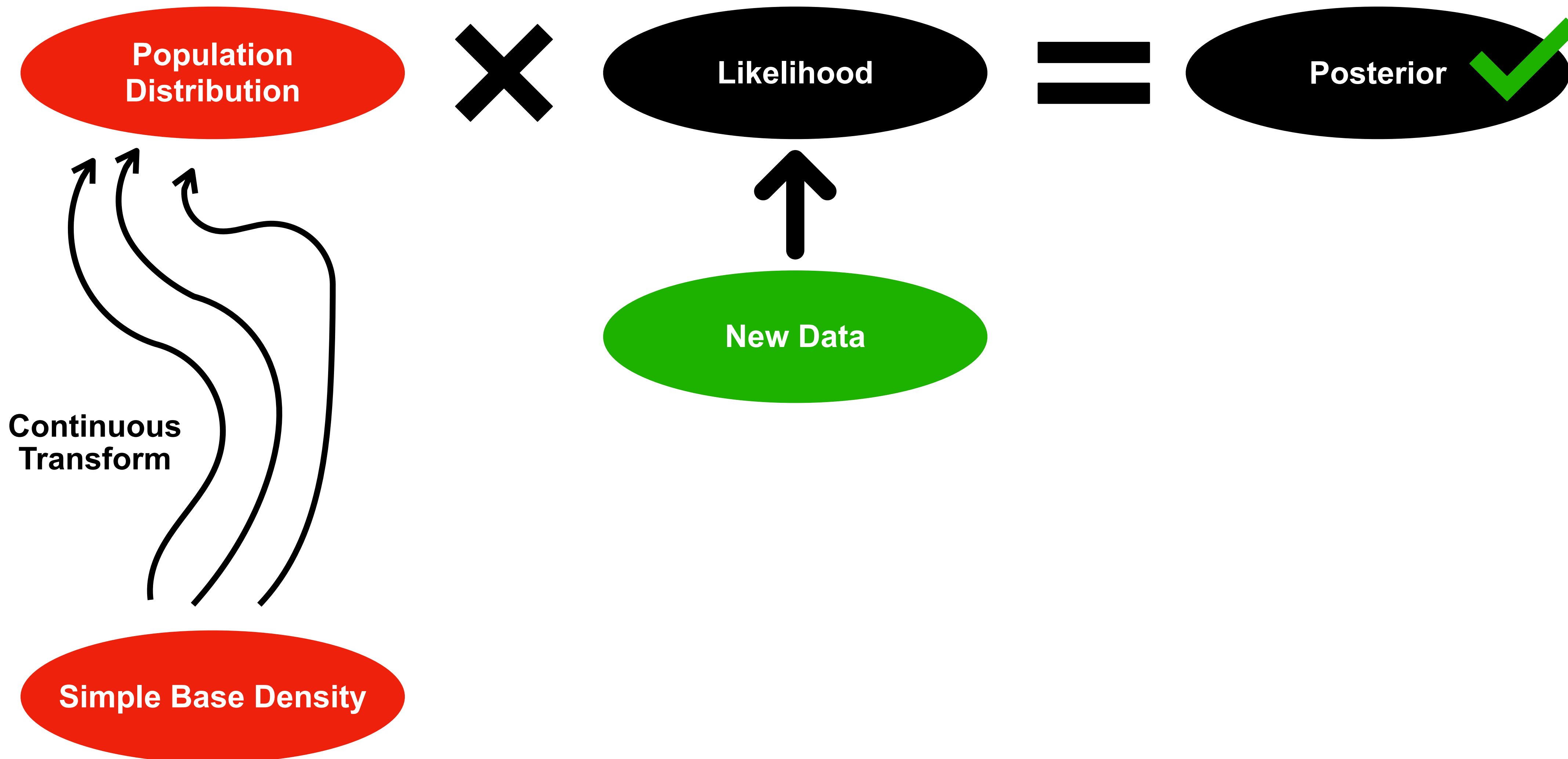
Inference under a diffusion model prior



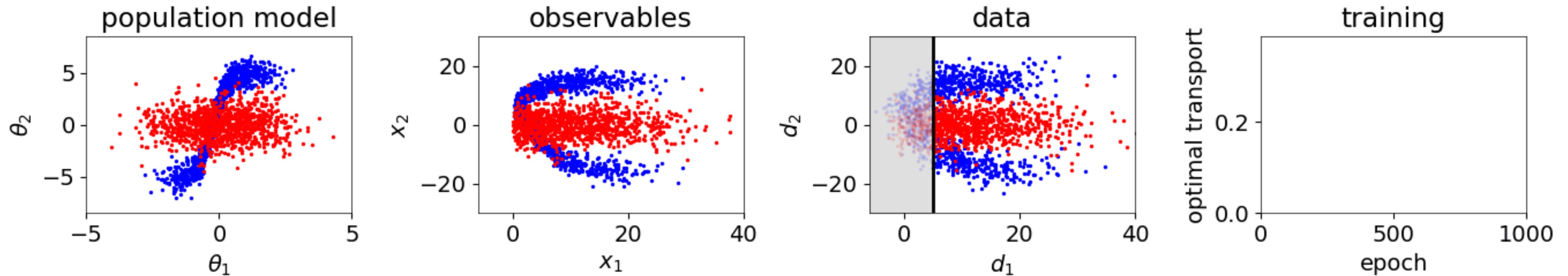
Inference under a diffusion model prior



Inference under a diffusion model prior

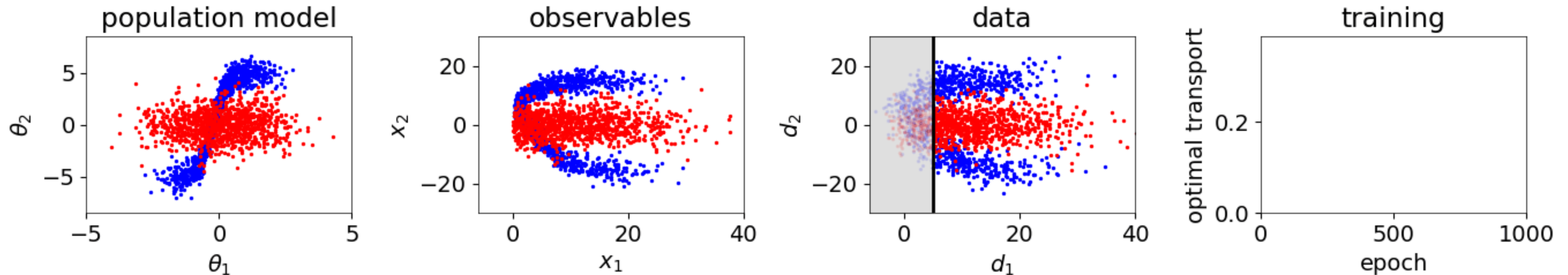


How do we forward model a galaxy catalog?



How do we forward model a galaxy catalog?

stellar population
synthesis (SPS)
parameters

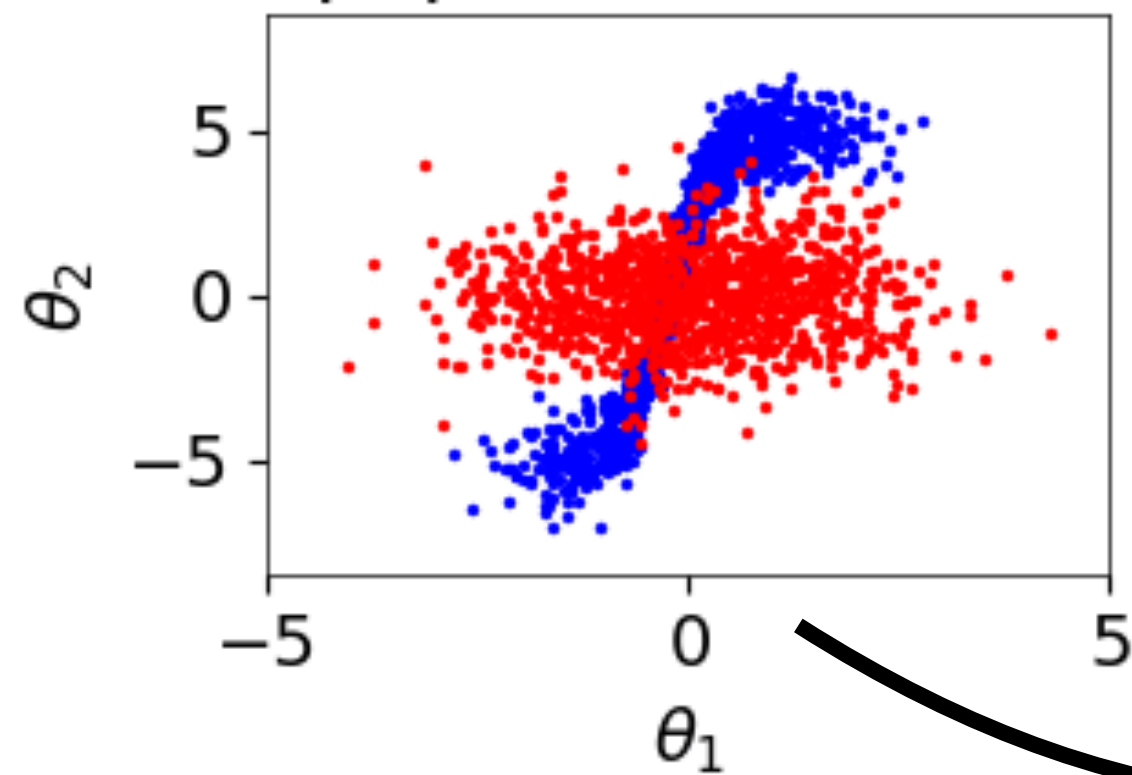


How do we forward model a galaxy catalog?

stellar population
synthesis (SPS)
parameters



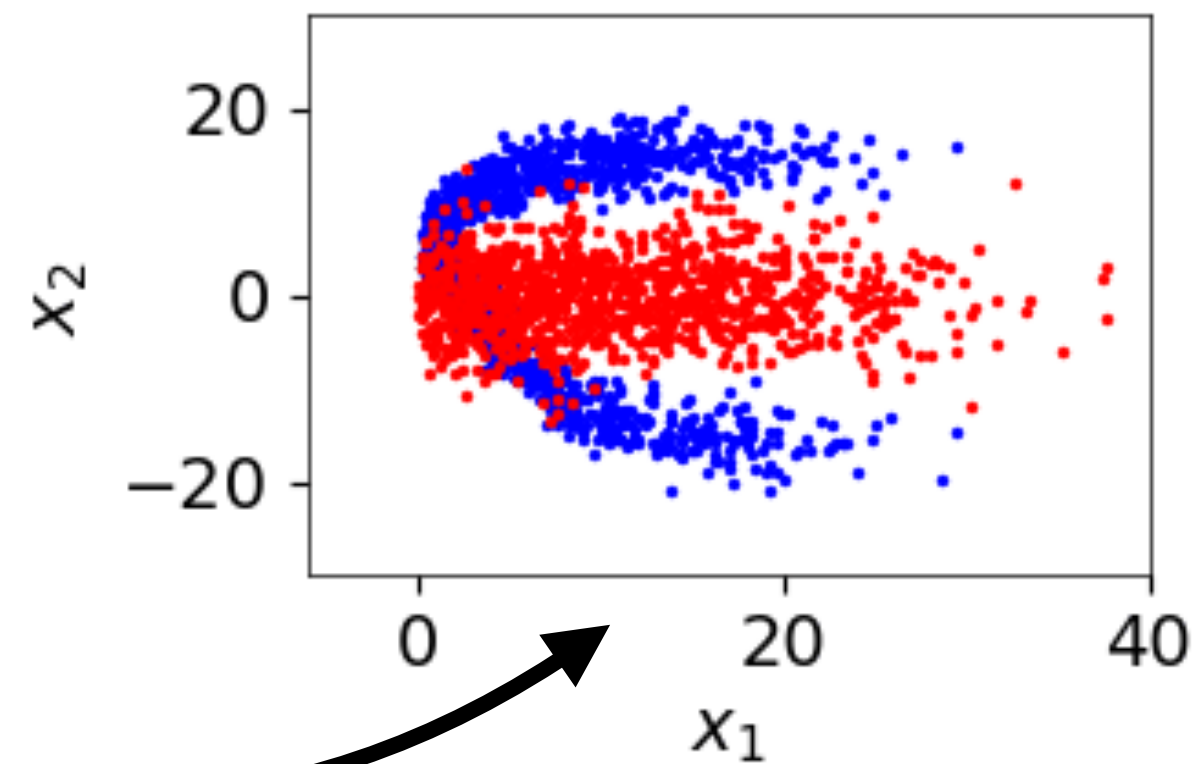
population model



noiseless fluxes



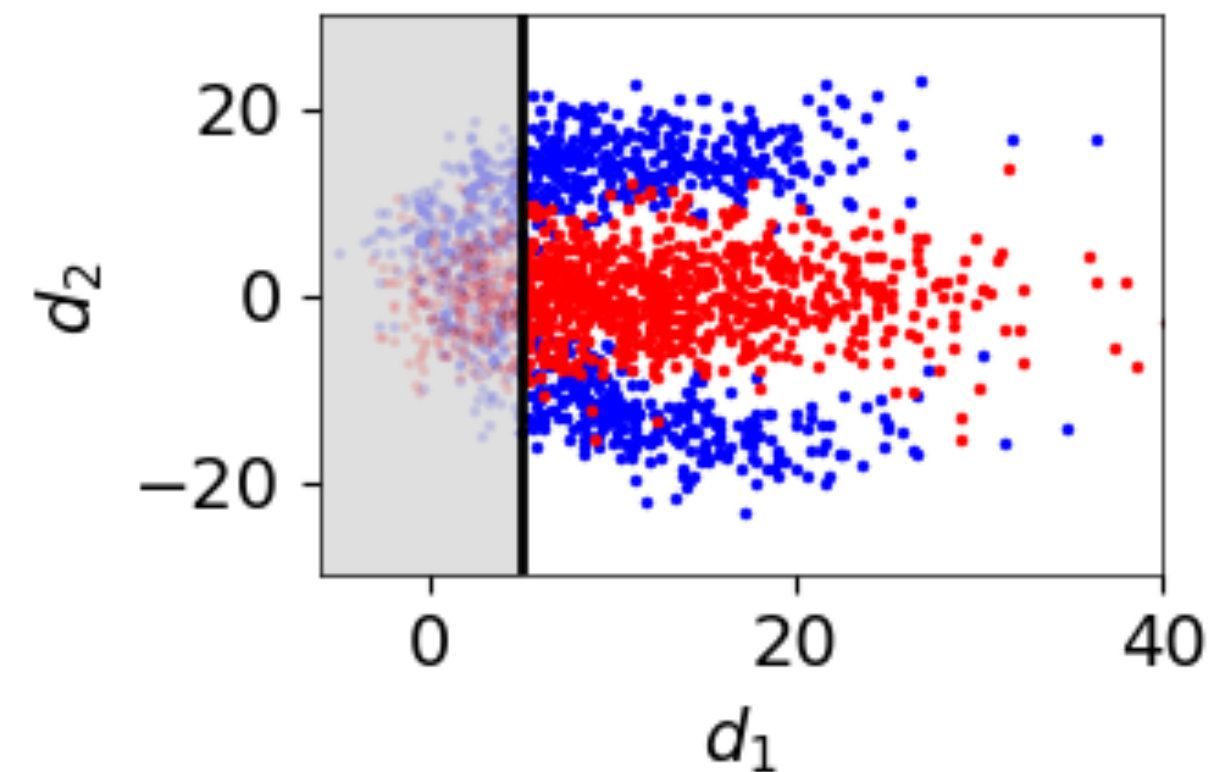
observables



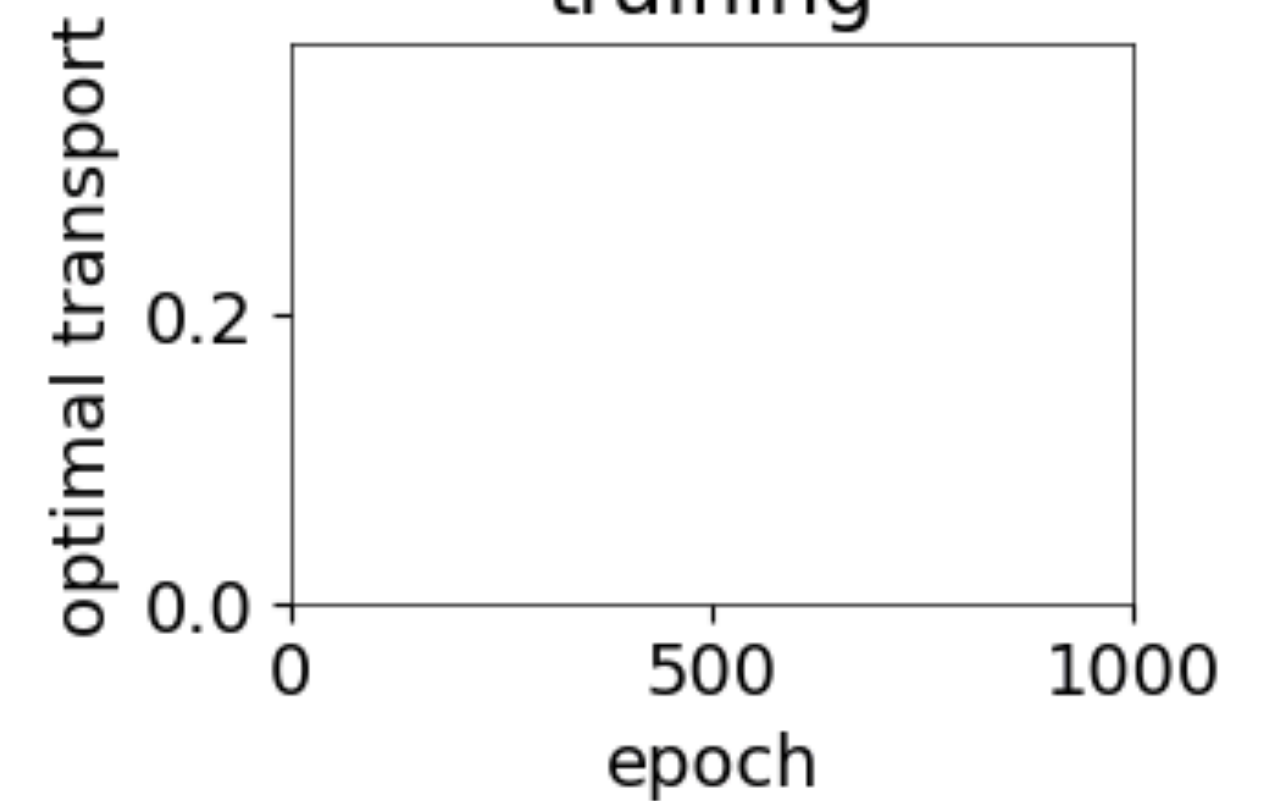
emulated SPS



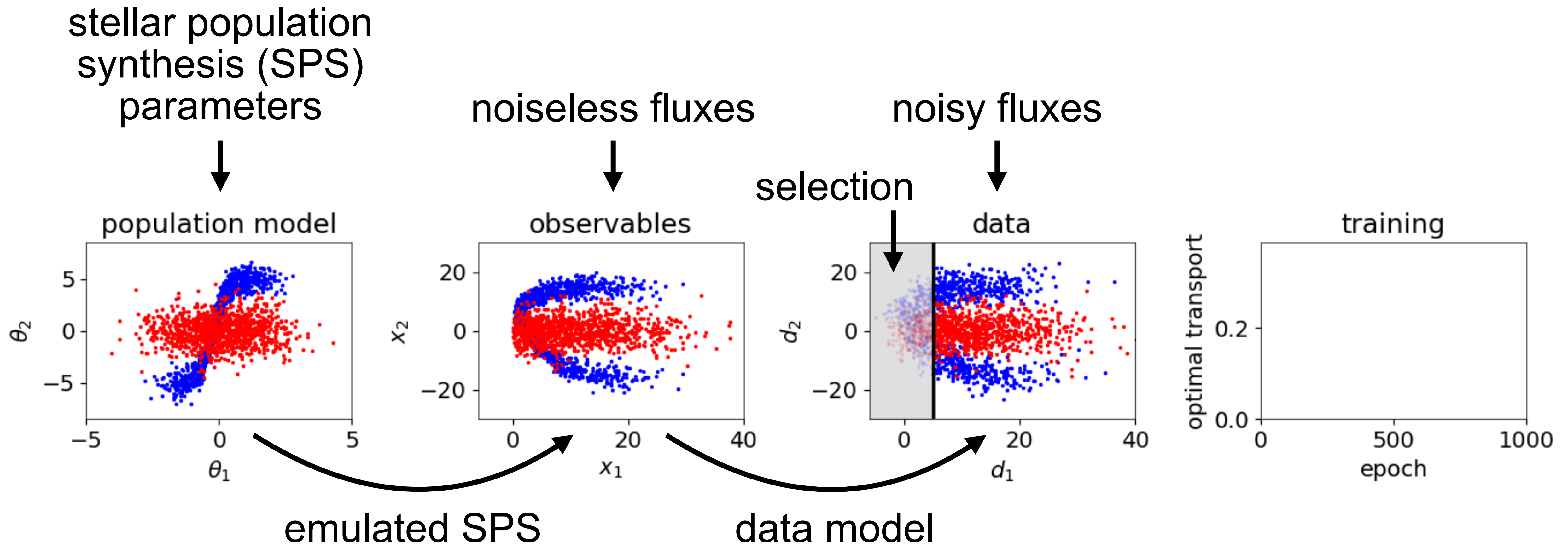
data



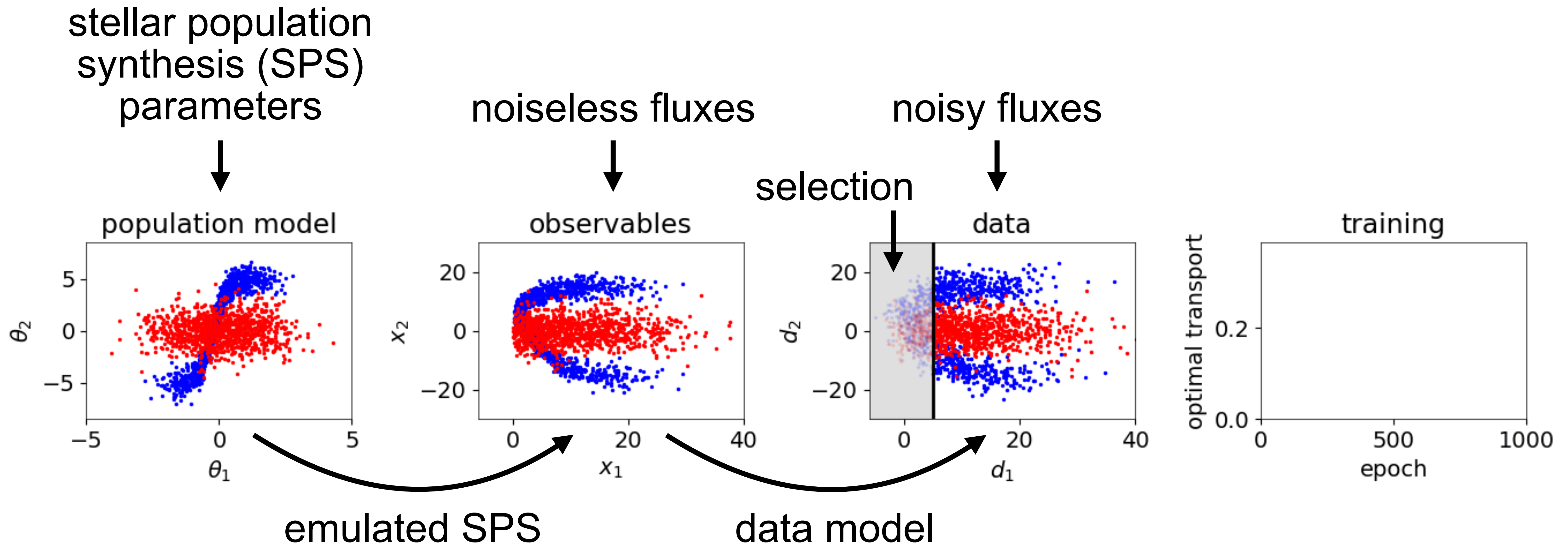
training



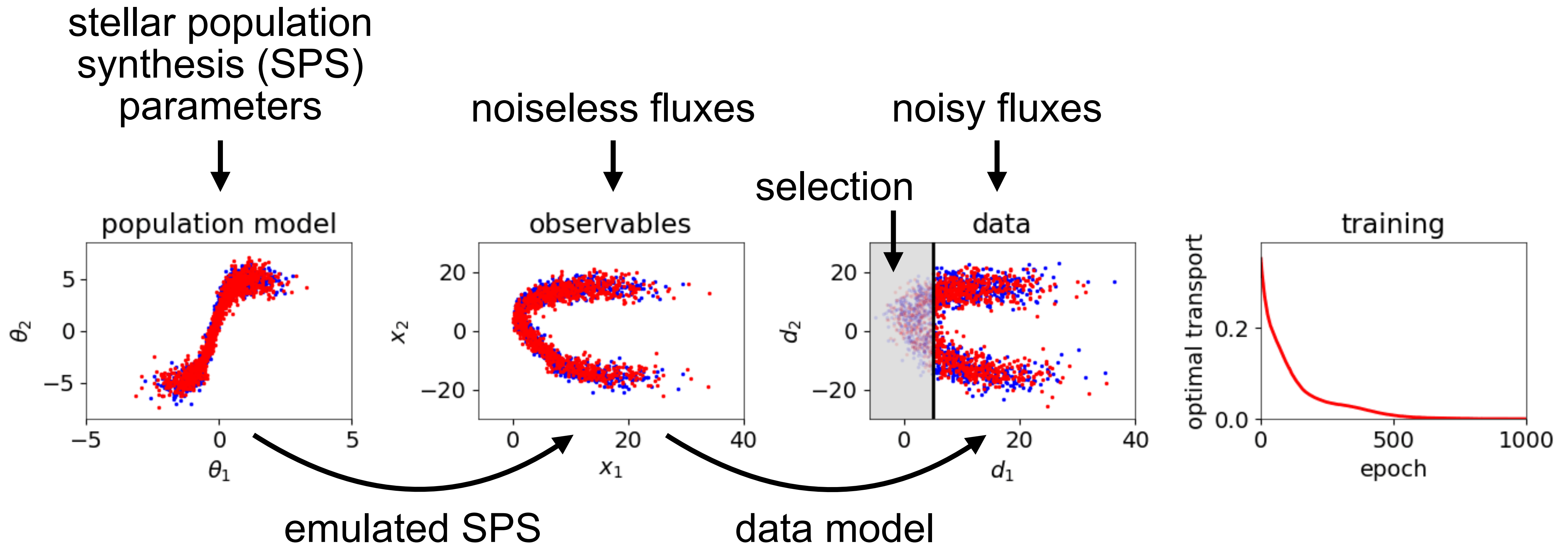
How do we forward model a galaxy catalog?



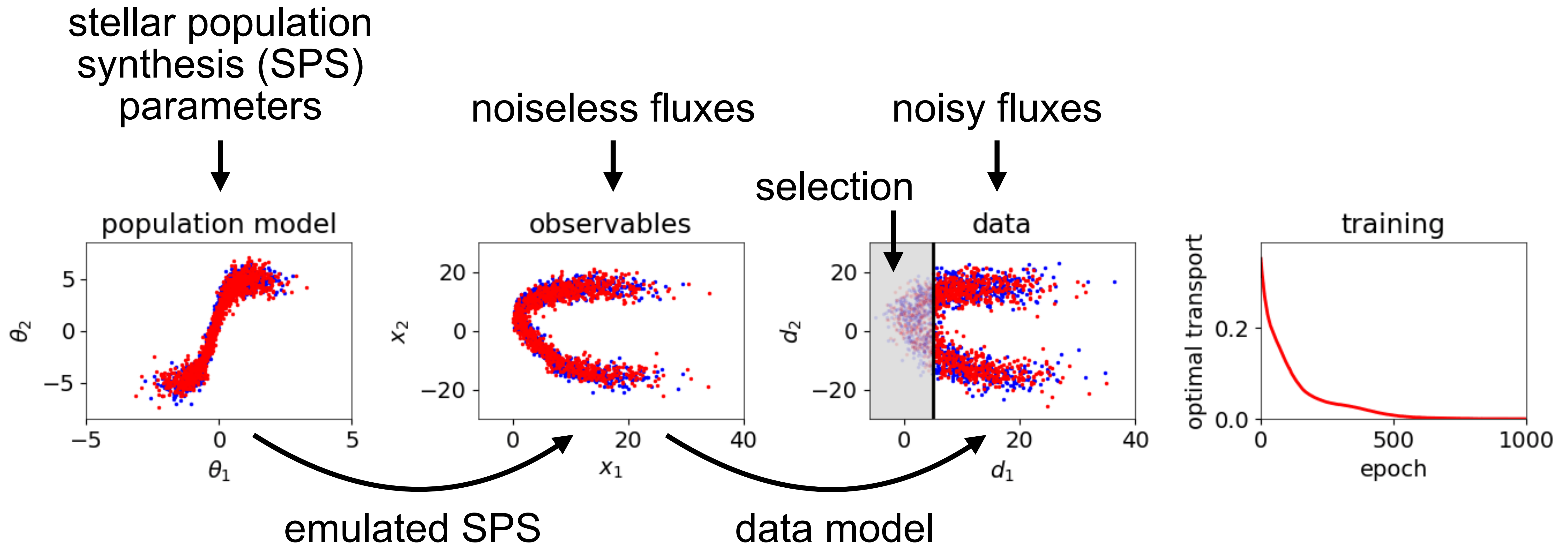
How do we forward model a galaxy catalog?



How do we forward model a galaxy catalog?



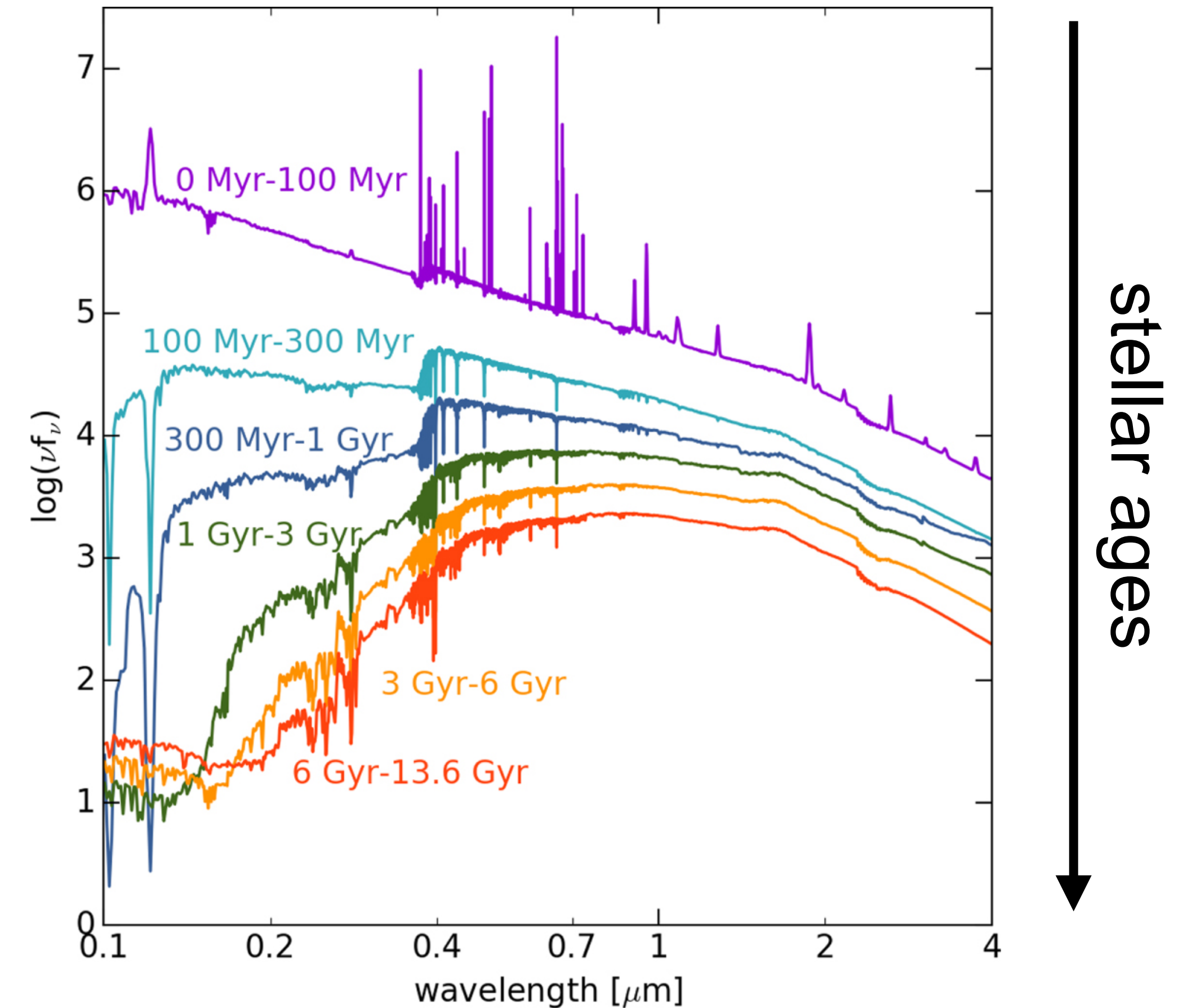
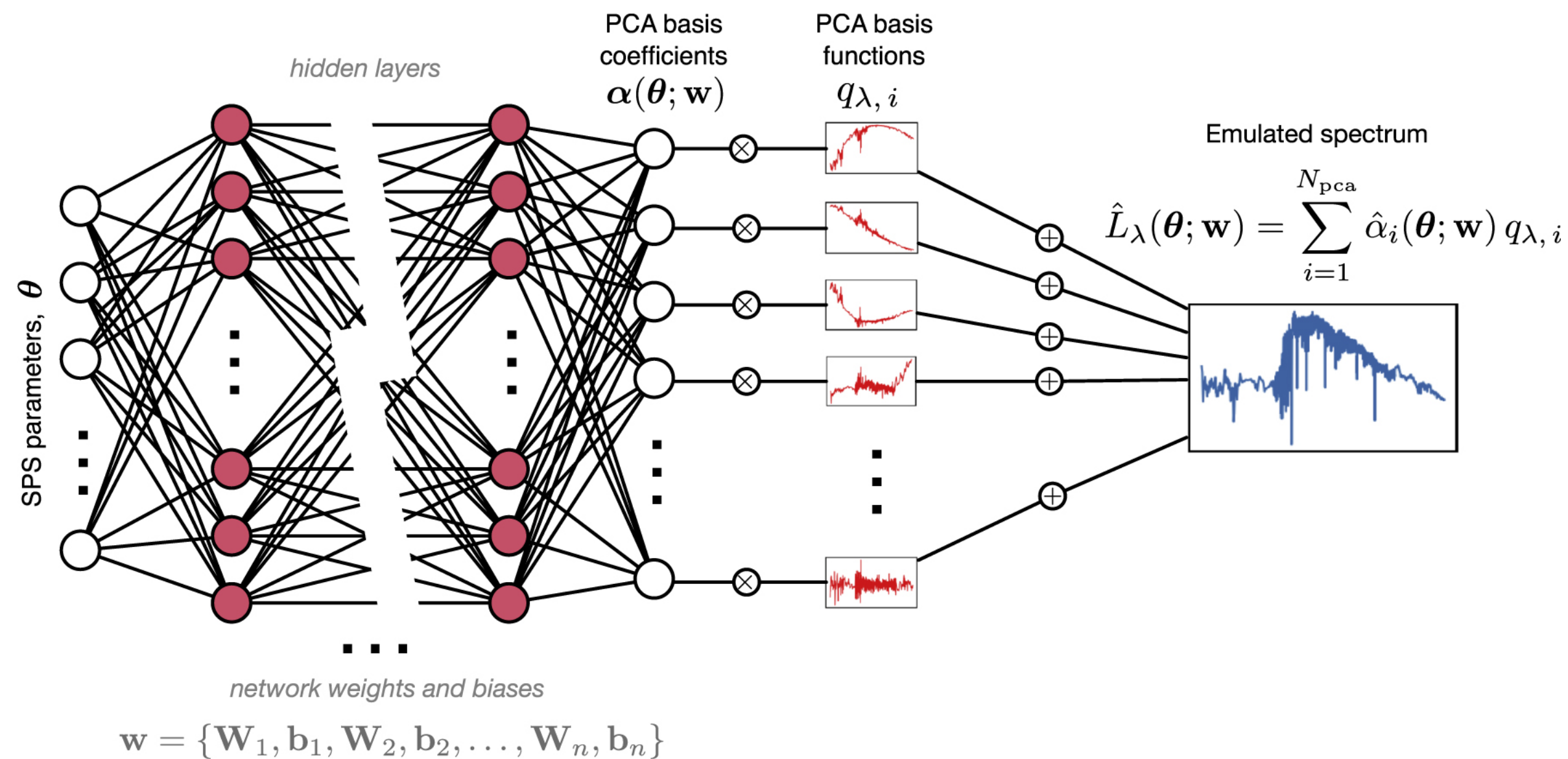
How do we forward model a galaxy catalog?



How do we represent a galaxy?

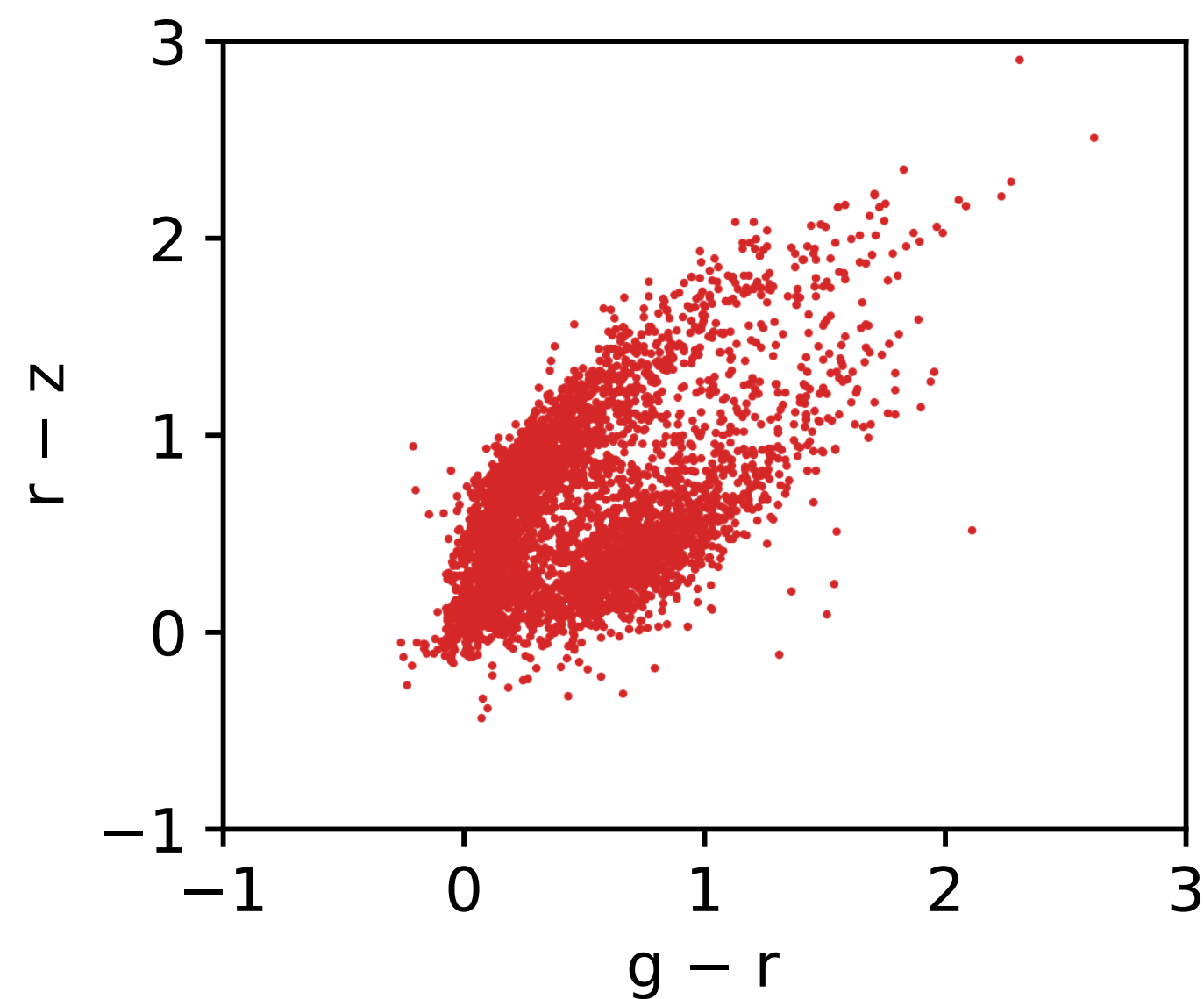
Stellar Population Synthesis (SPS)

- 16 parameters (\approx Prospector- α)
- Emulated with Speculator



What will our population model be?

A score-based diffusion model



target distribution
over physical
parameters ($t = 0$)

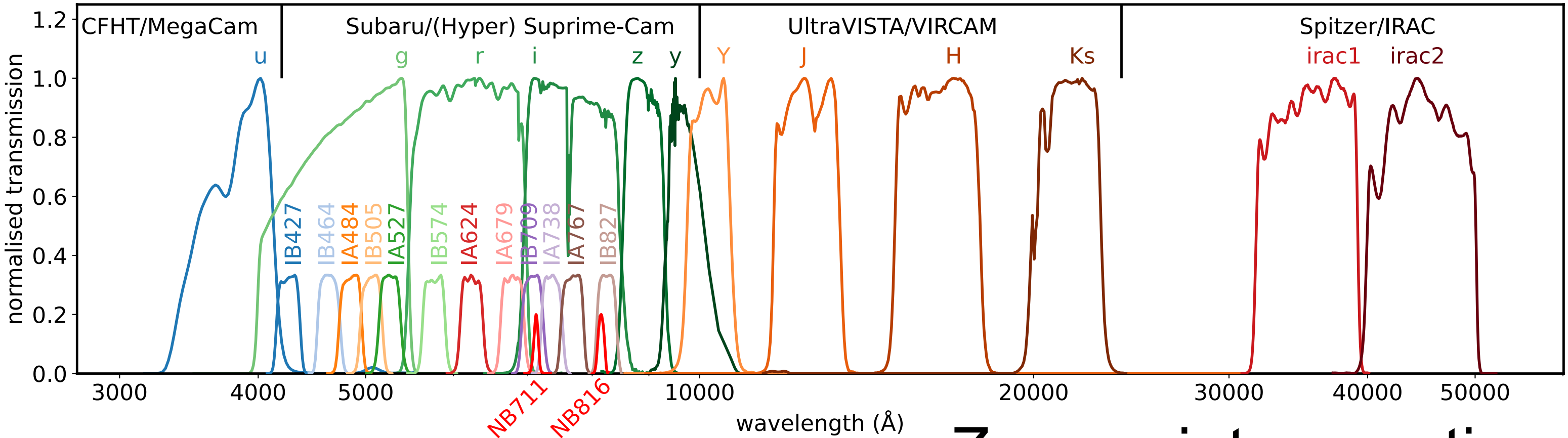
→
differential equation
←



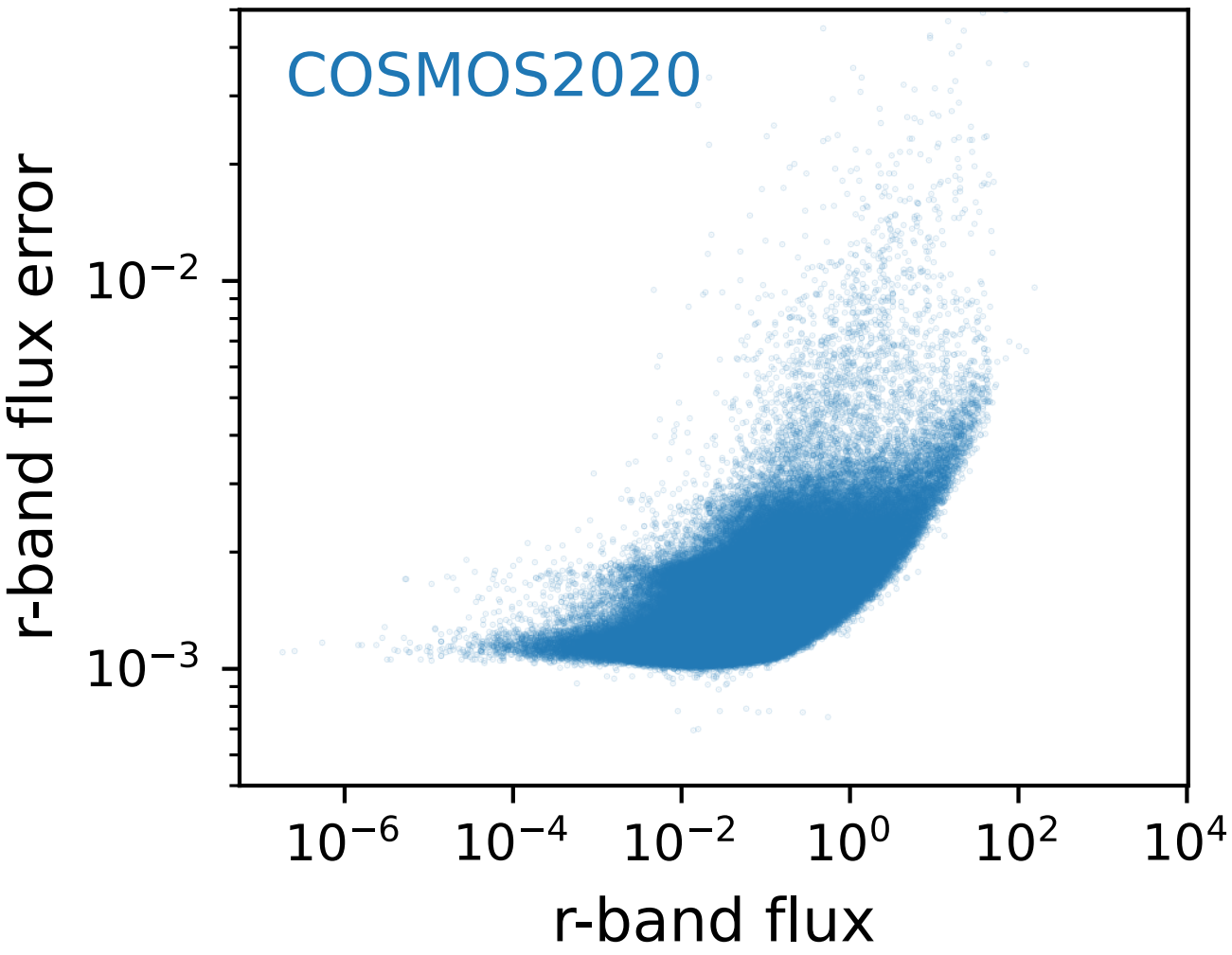
iid Gaussian noise
($t = T$)

What goes into the data model?

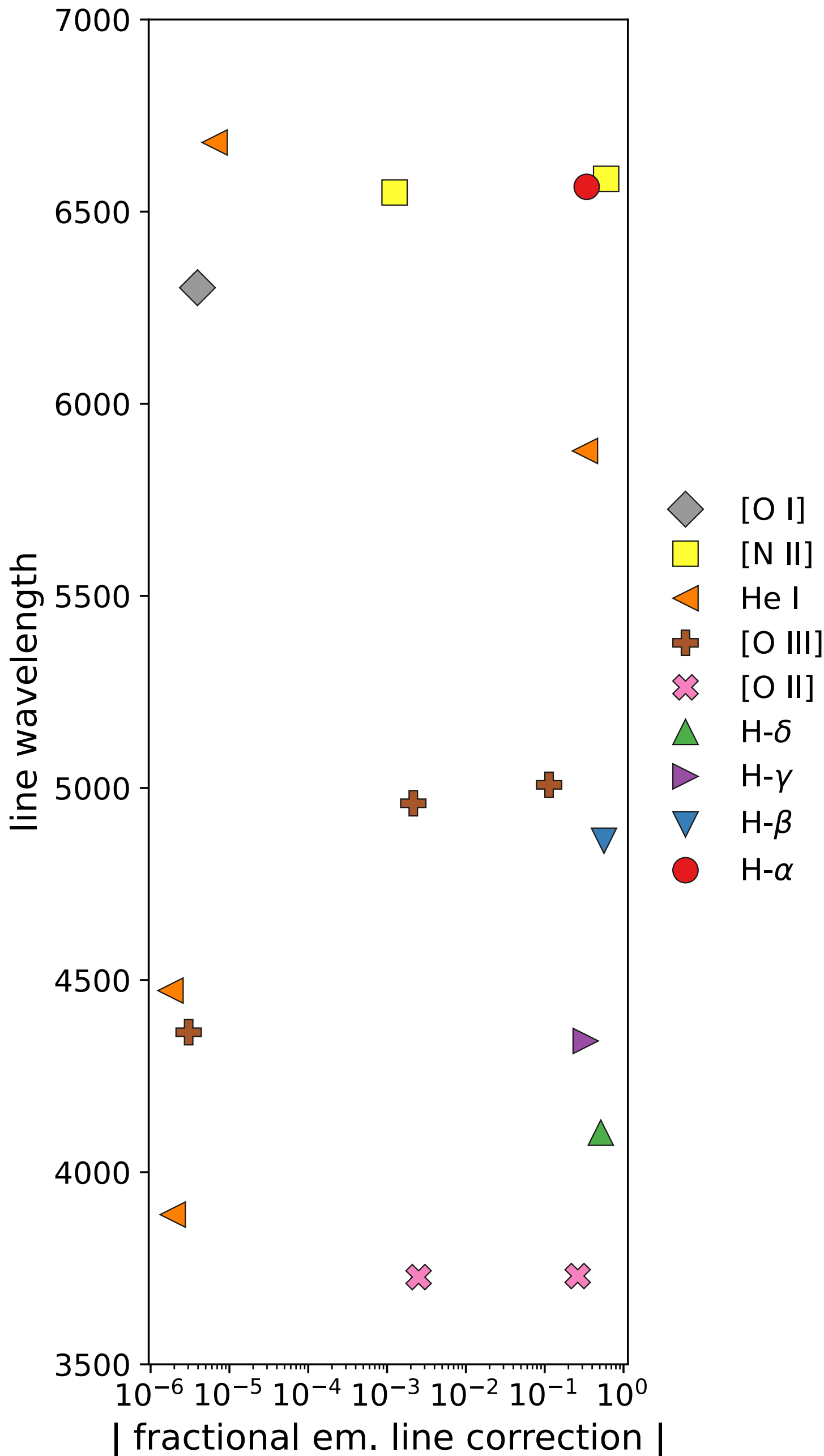
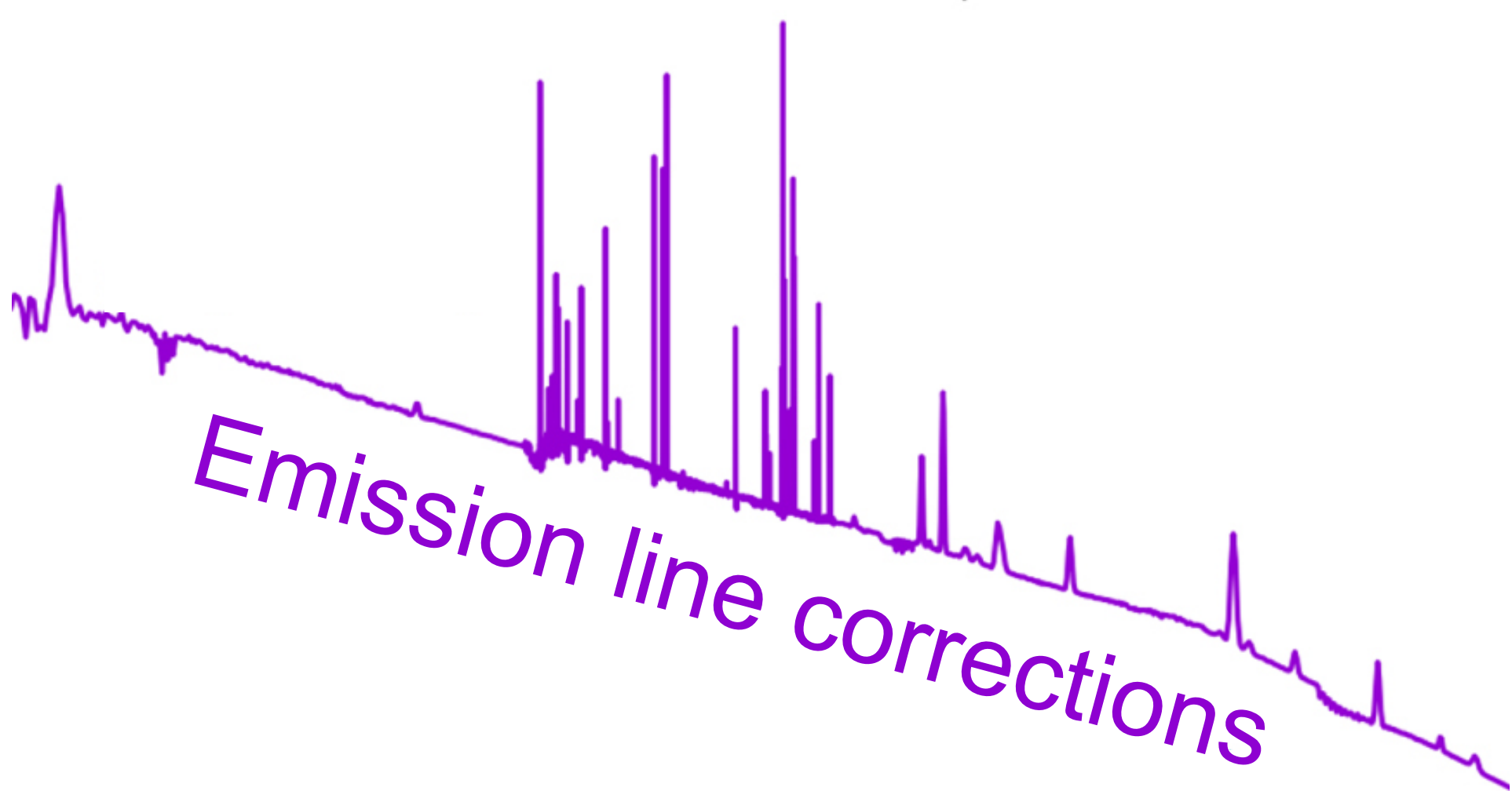
Uncertainty and calibration

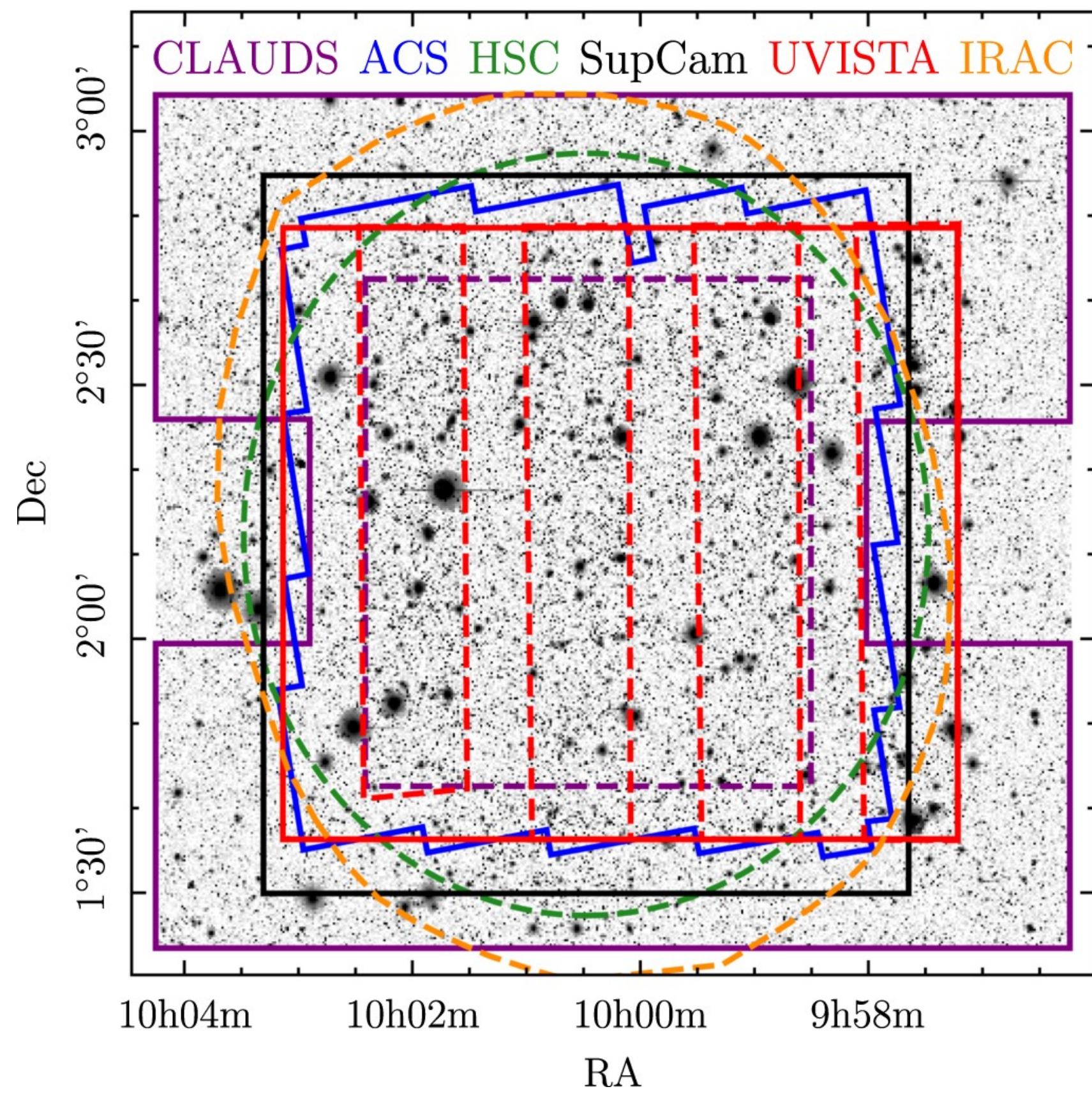


Zero-point corrections



Uncertainty model

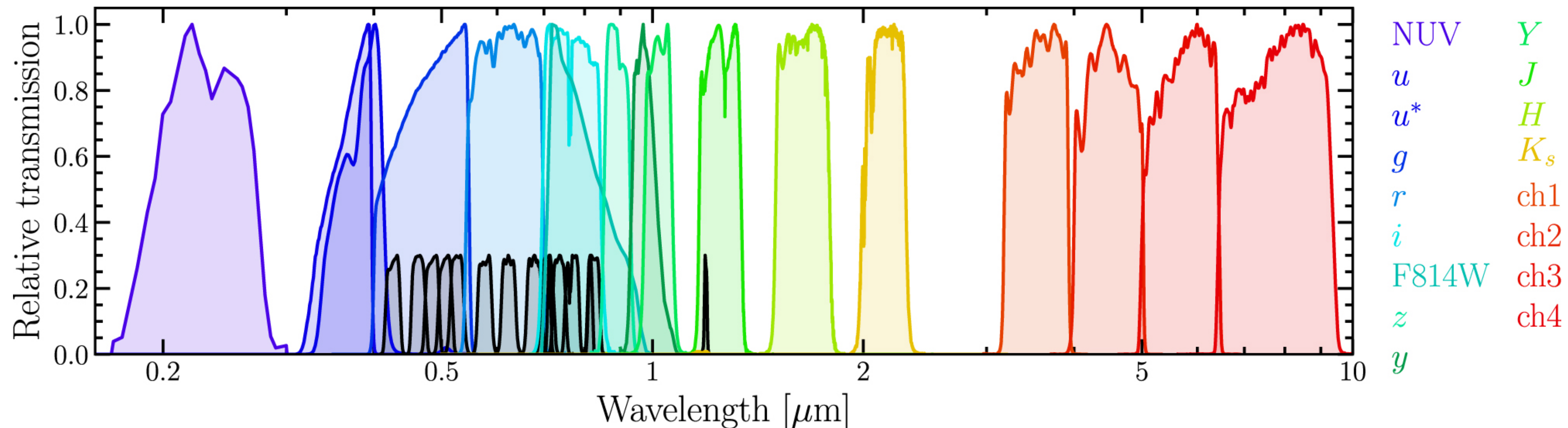




What data will we use?

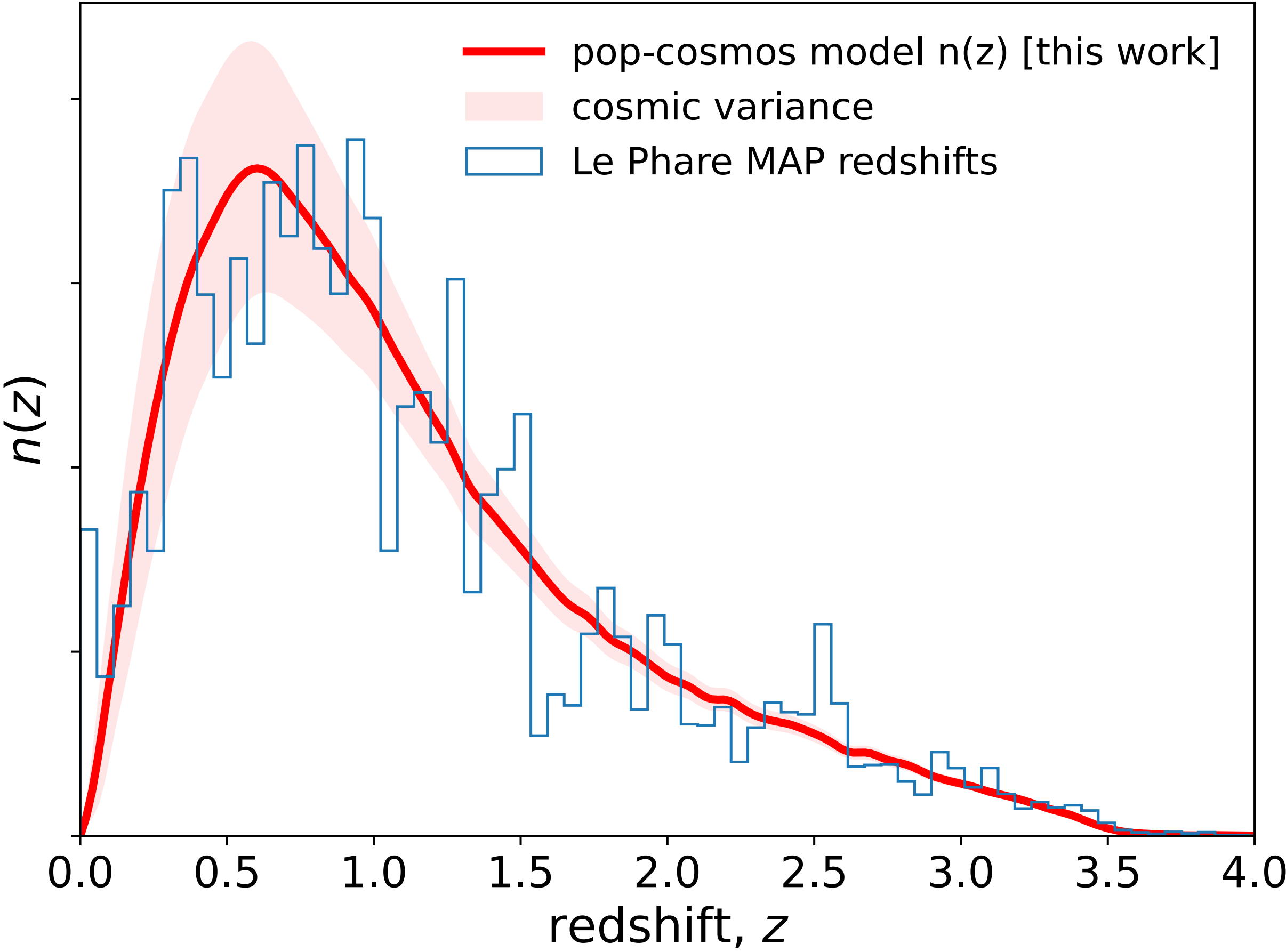
COSMOS2020 (Weaver et al. 2022)

- Has $\approx 420,000$ galaxies with *Spitzer* IRAC *Ch. 1* < 26
- Wide and narrow bands
- Coverage from near-UV to mid-IR

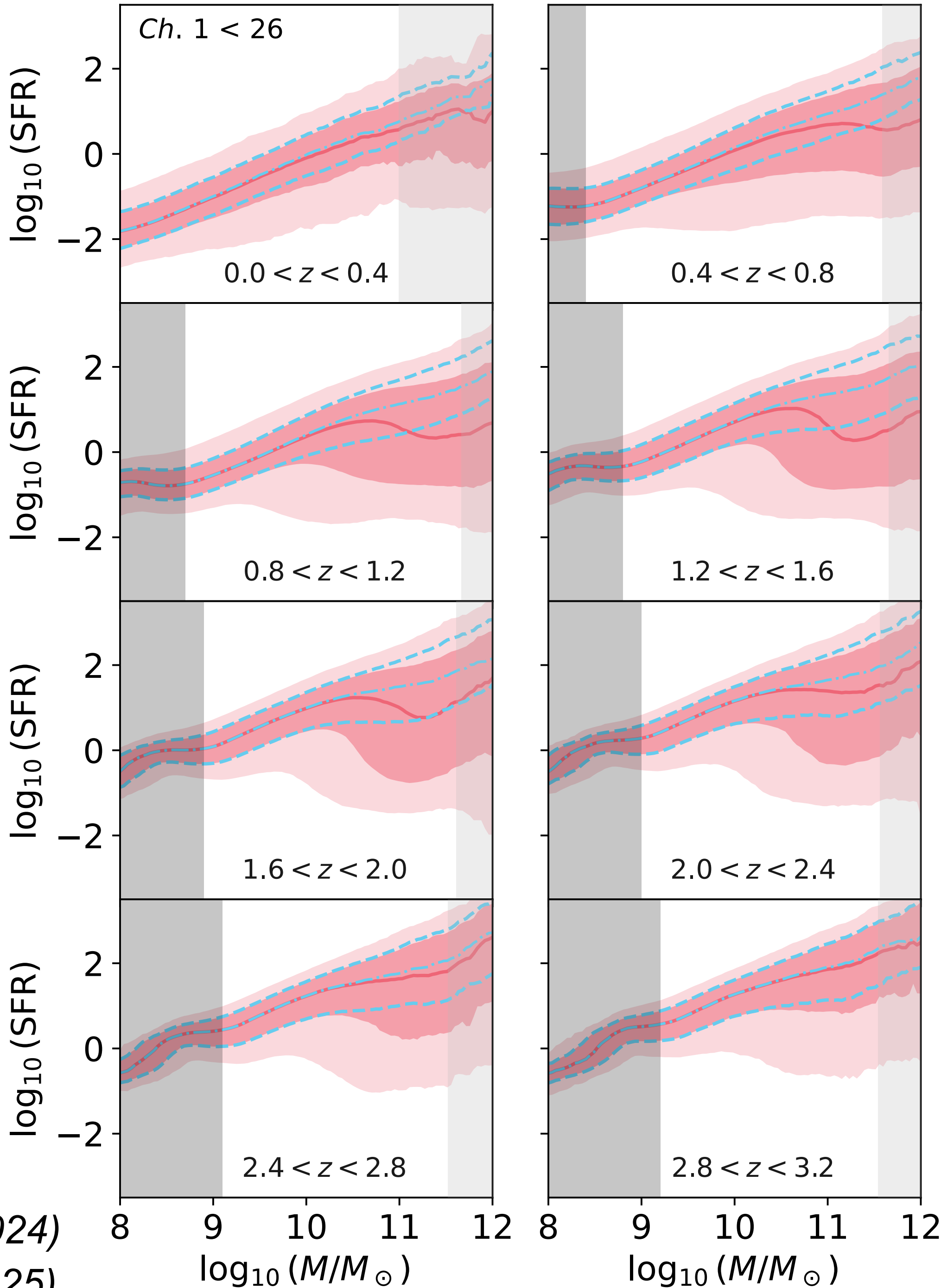


Trained Model

Population-level predictions...



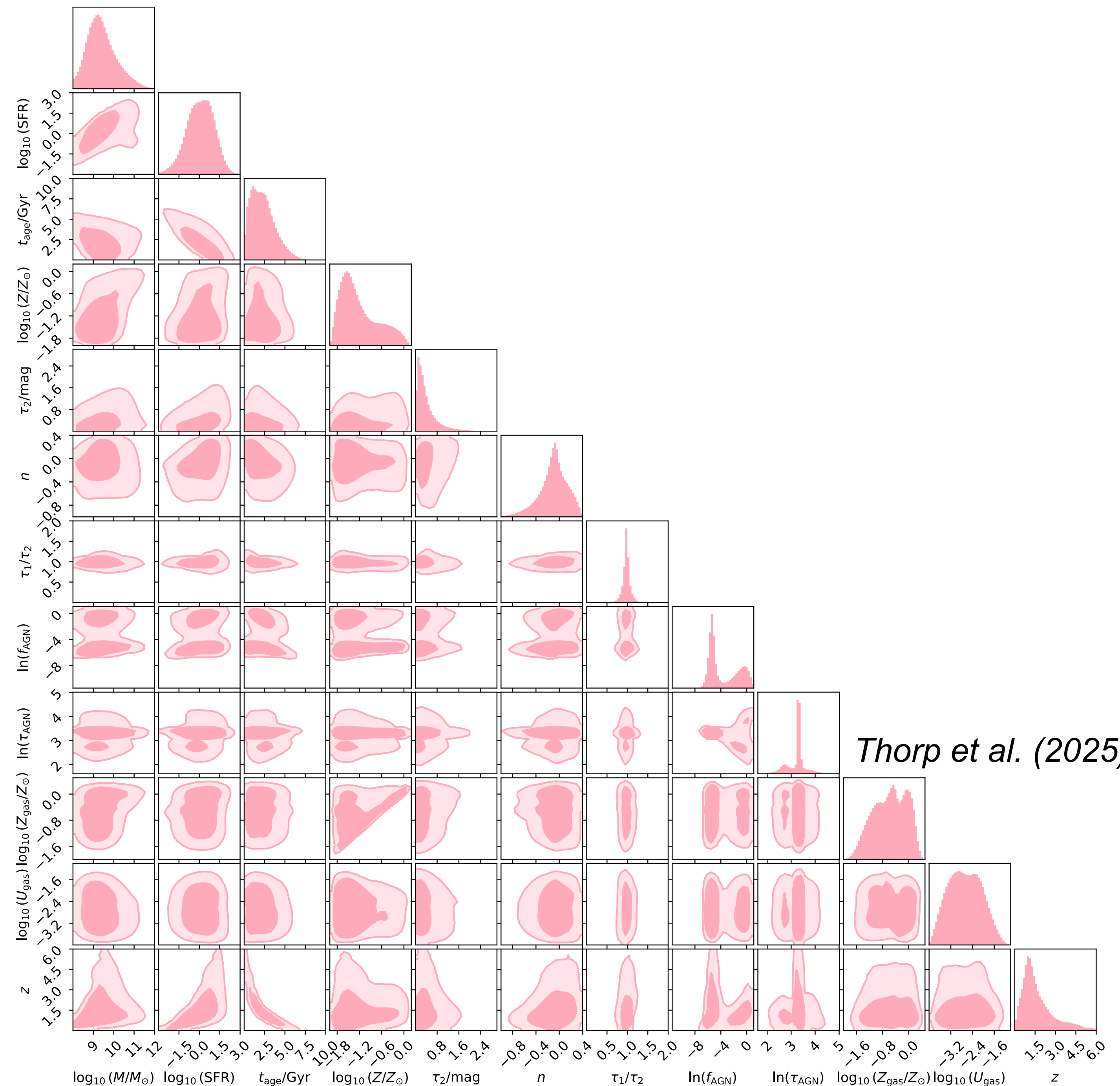
Alsing et al. (2024)
Thorp et al. (2025)

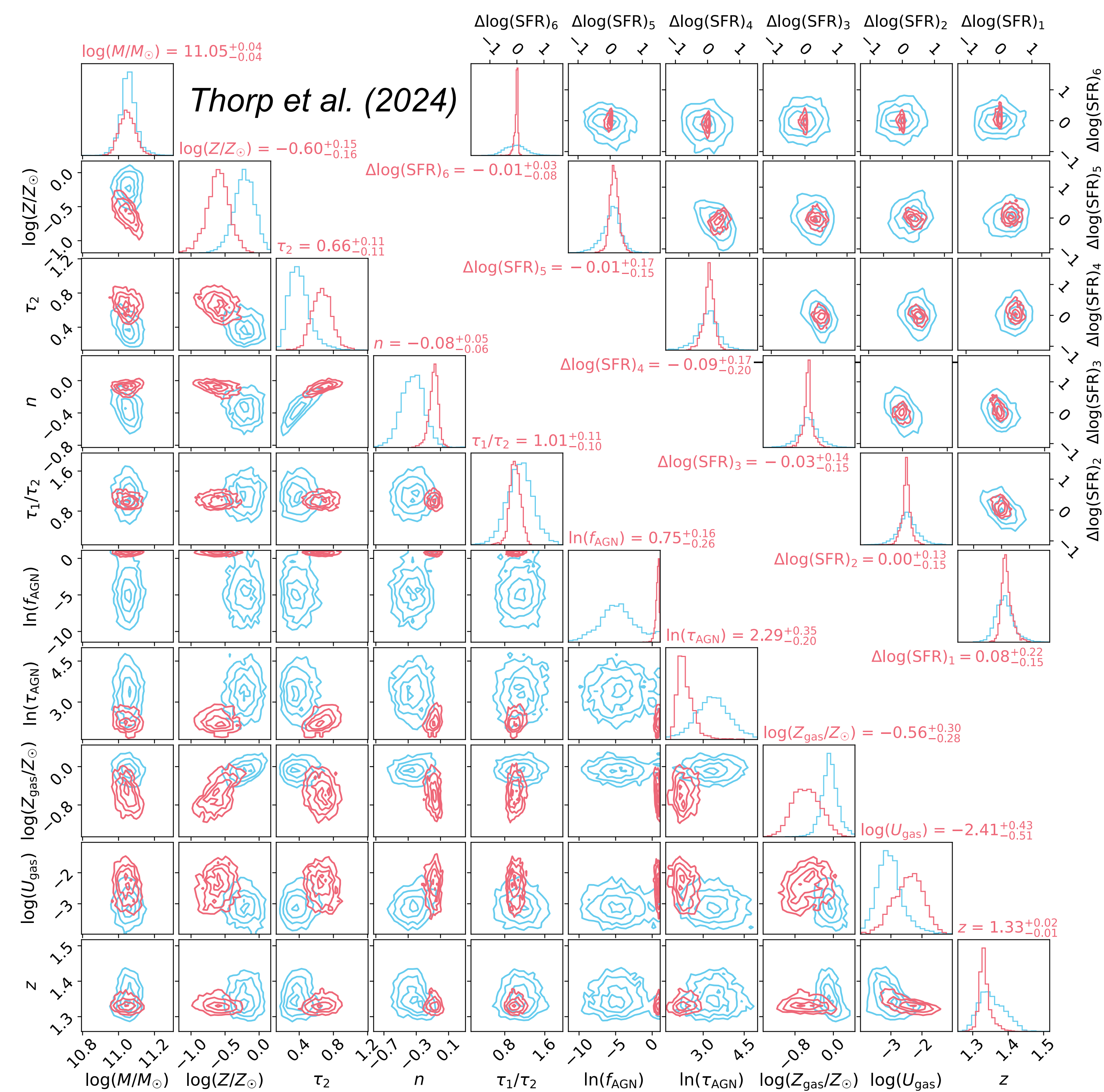


Trained Model

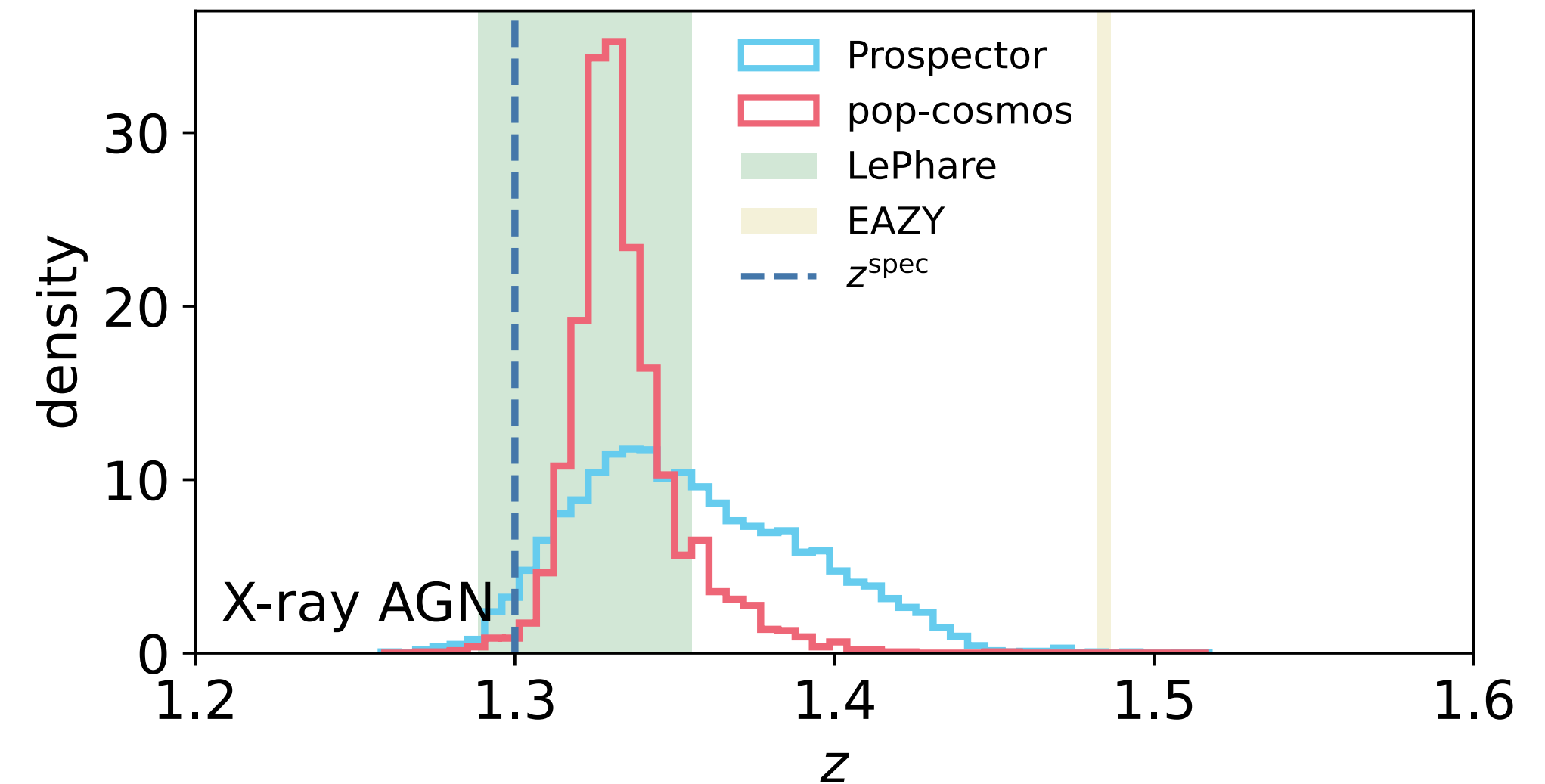
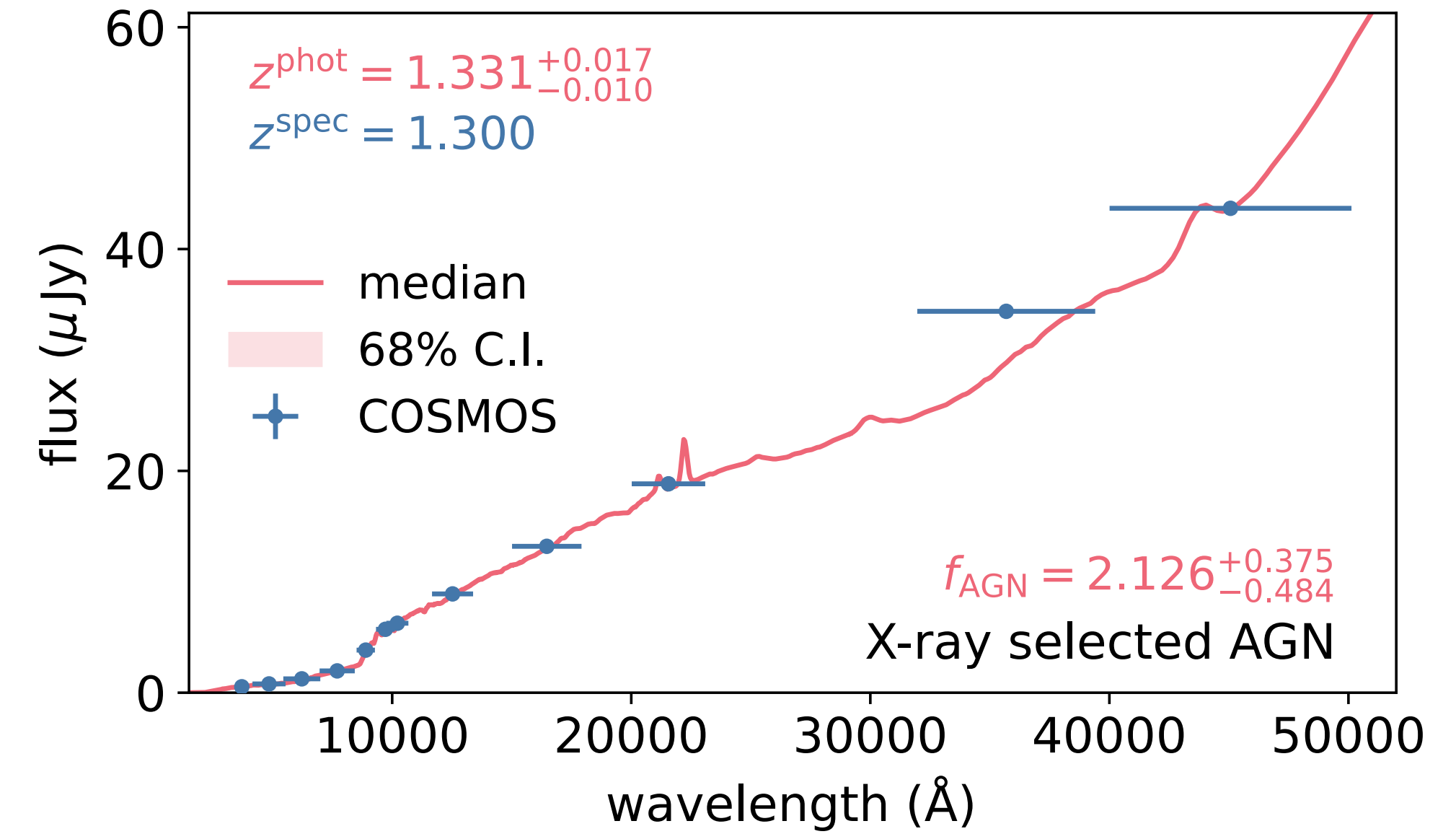
As a prior...

- The learned population distribution over SPS parameters can be used for downstream inference for individual galaxies
- Why? Some science applications need redshift/parameter inference for *specific galaxies*, not just population-level constraints



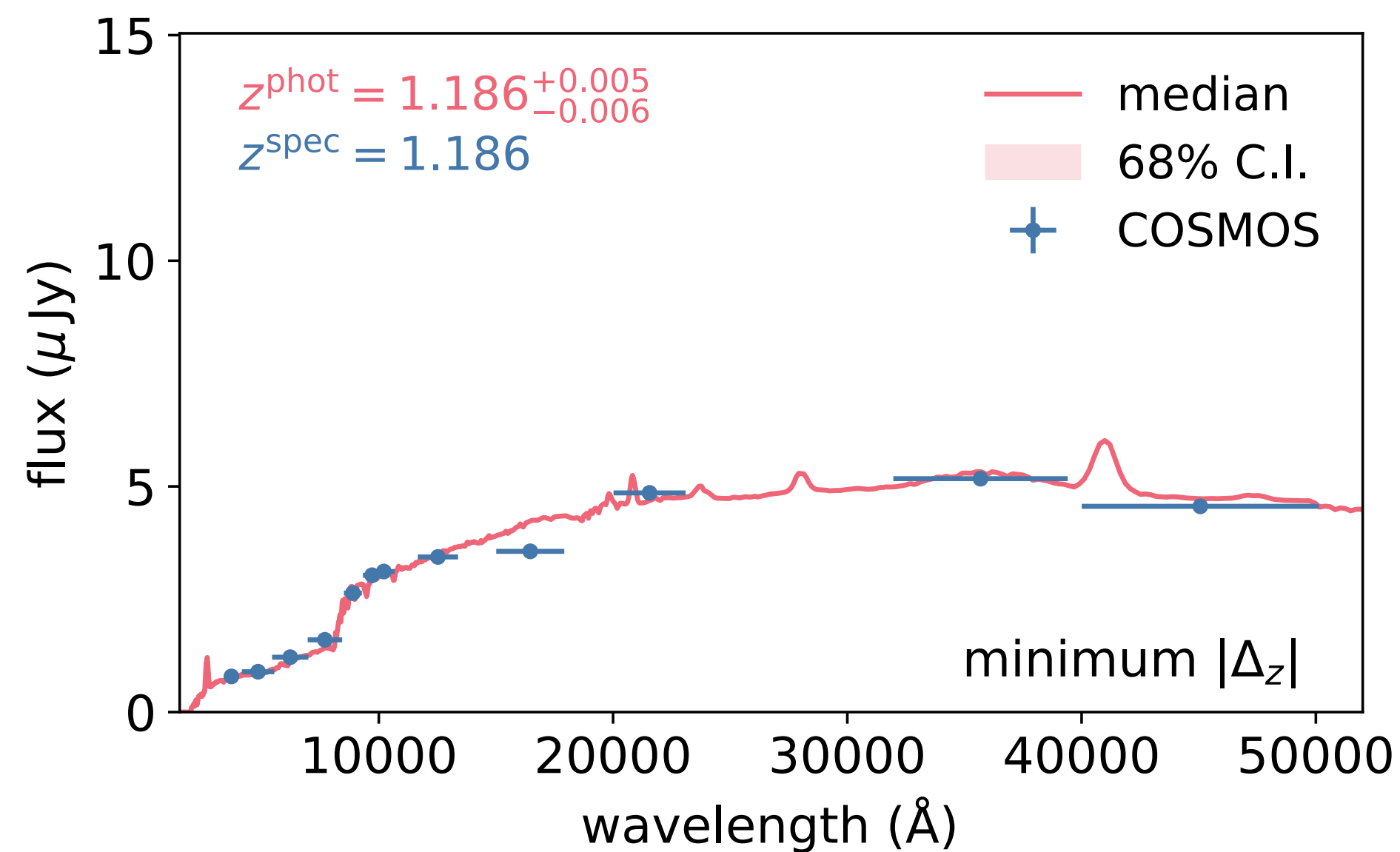
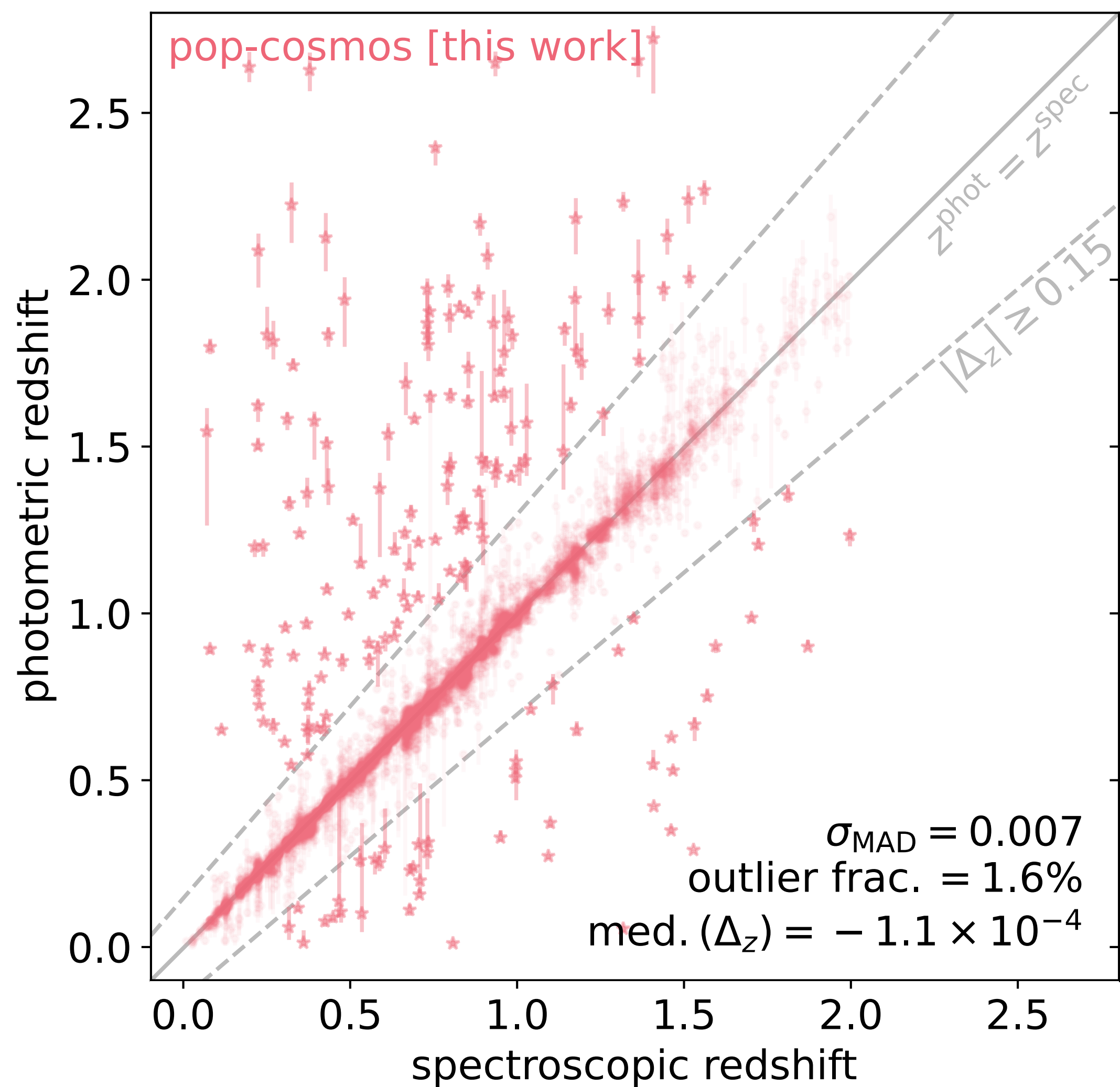


SED Fitting

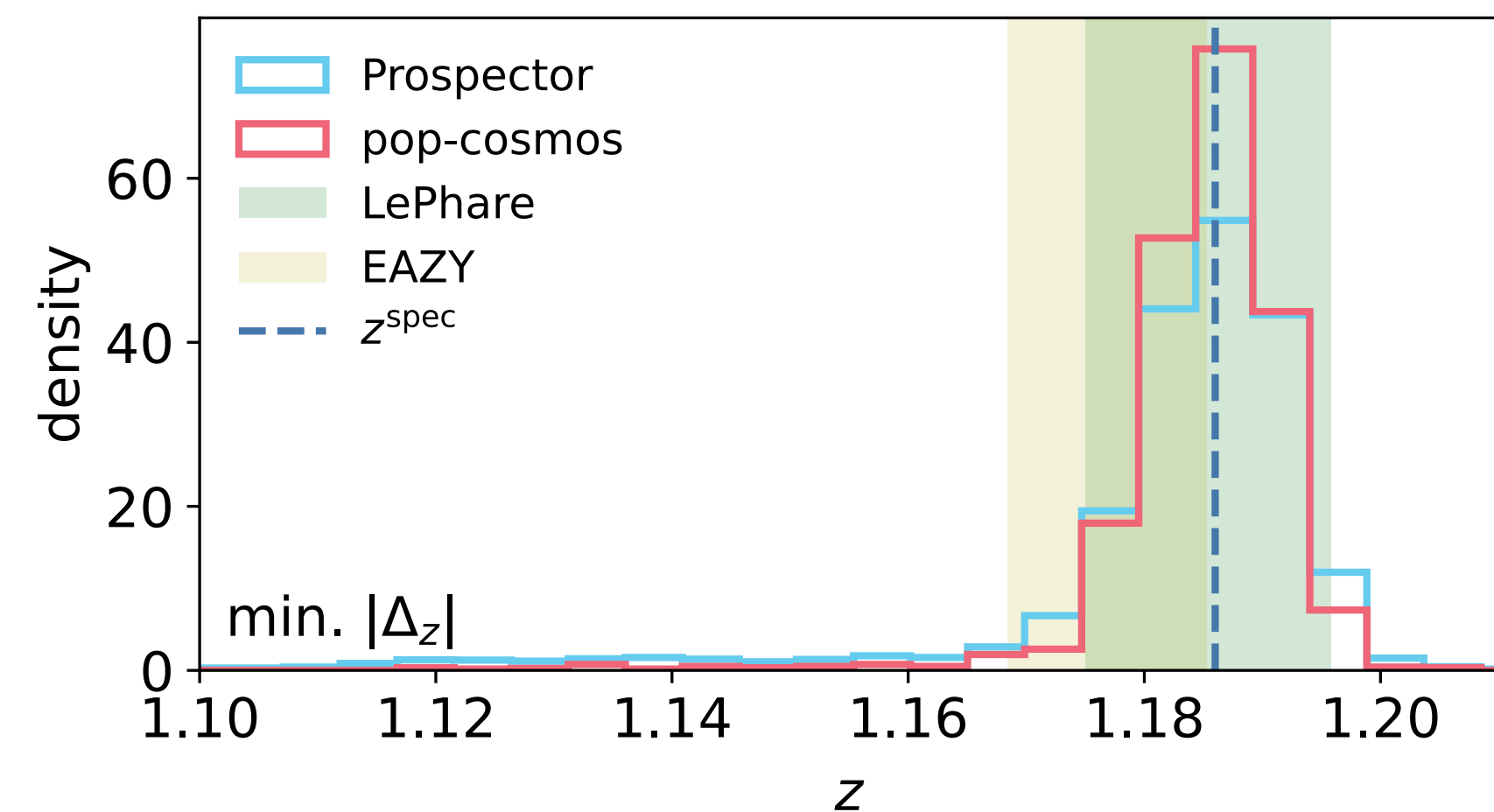


SED Fitting

Redshift inference compared to spectroscopy...

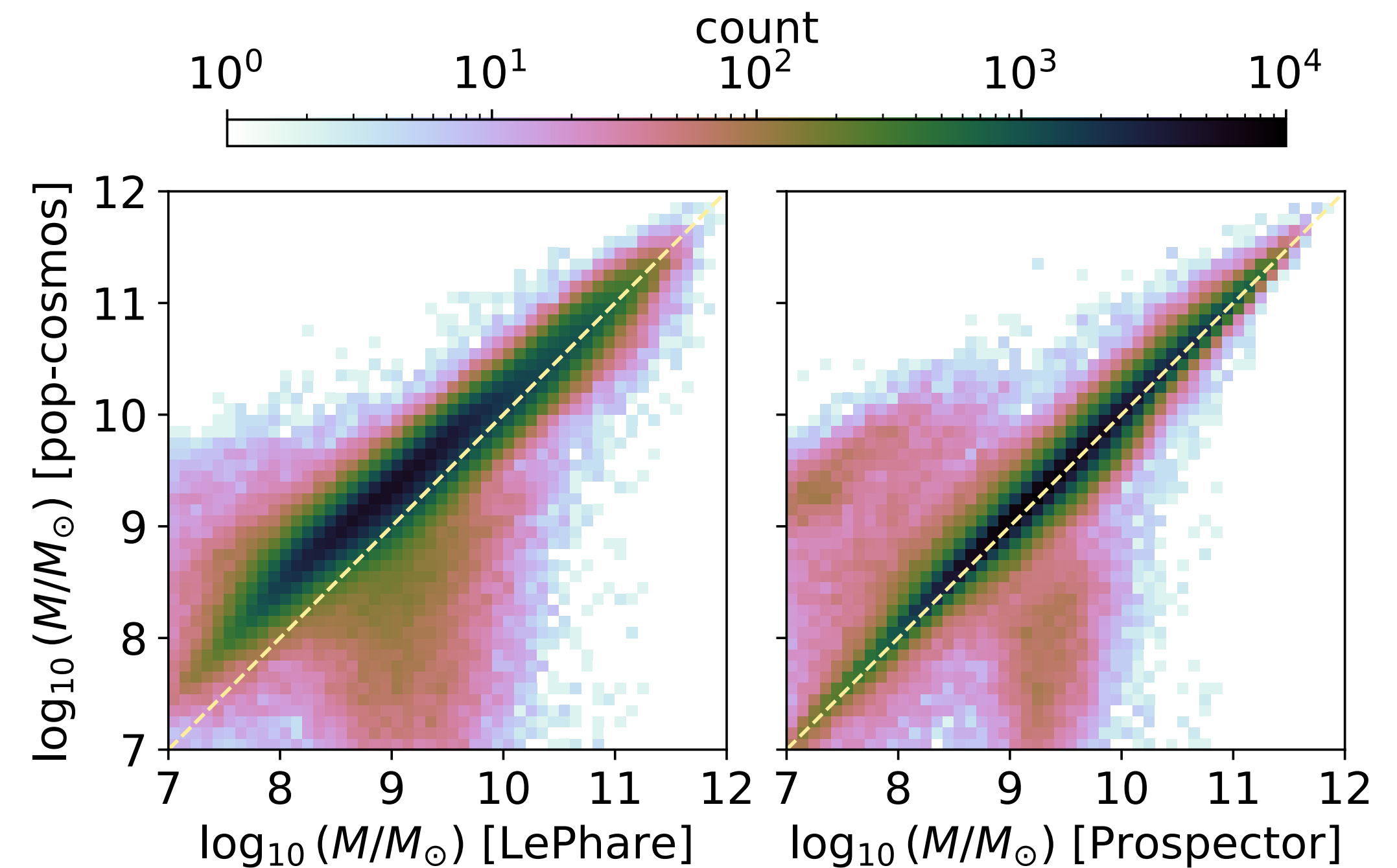
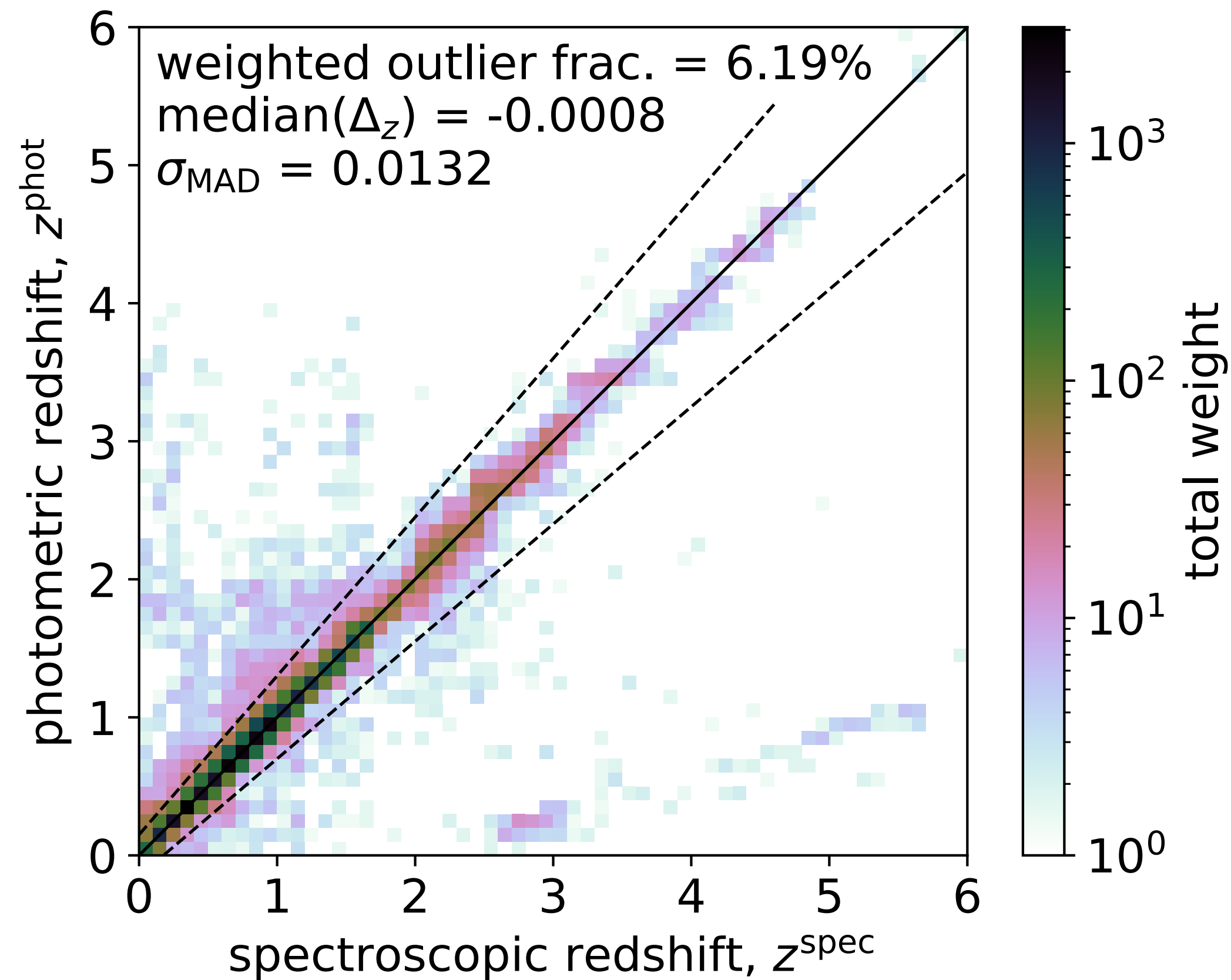


Thorp et al. (2024)



SED Fitting

Redshift and stellar mass inference...



Transferring to other surveys?

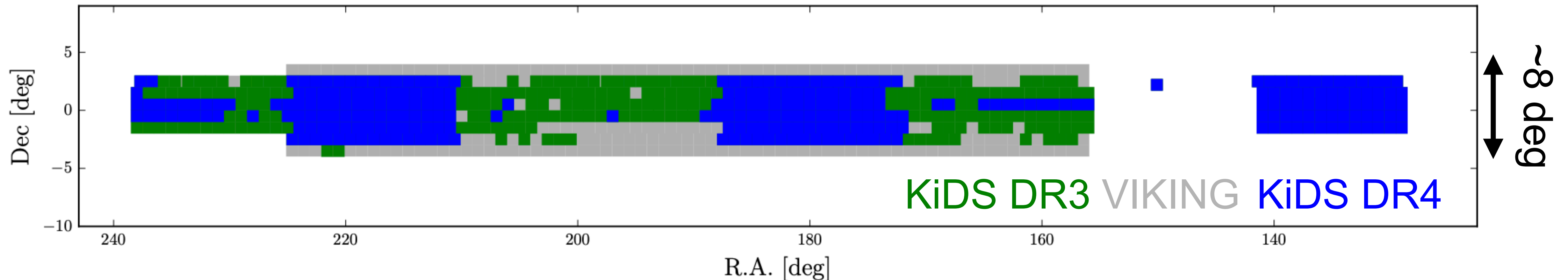
Kilo Degree Survey (KiDS; Kuijken et al. 2019)

- Has $\approx 70,000,000$ galaxies
- Broadband only: $ugriZYJHK_S$
- Much larger area ($1,000 \text{ deg}^2$)

$\sim 85 \text{ deg}$



KiDS-North



Transferring to other surveys?

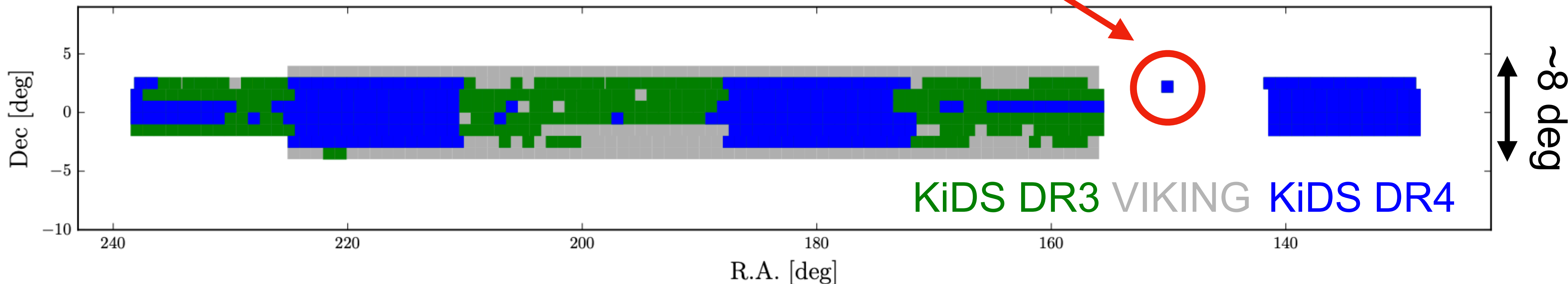
Kilo Degree Survey (KiDS; Kuijken et al. 2019)

- Has $\approx 70,000,000$ galaxies
- Broadband only: $ugriZYJHK_s$
- Much larger area ($1,000 \text{ deg}^2$)

COSMOS (our training data)

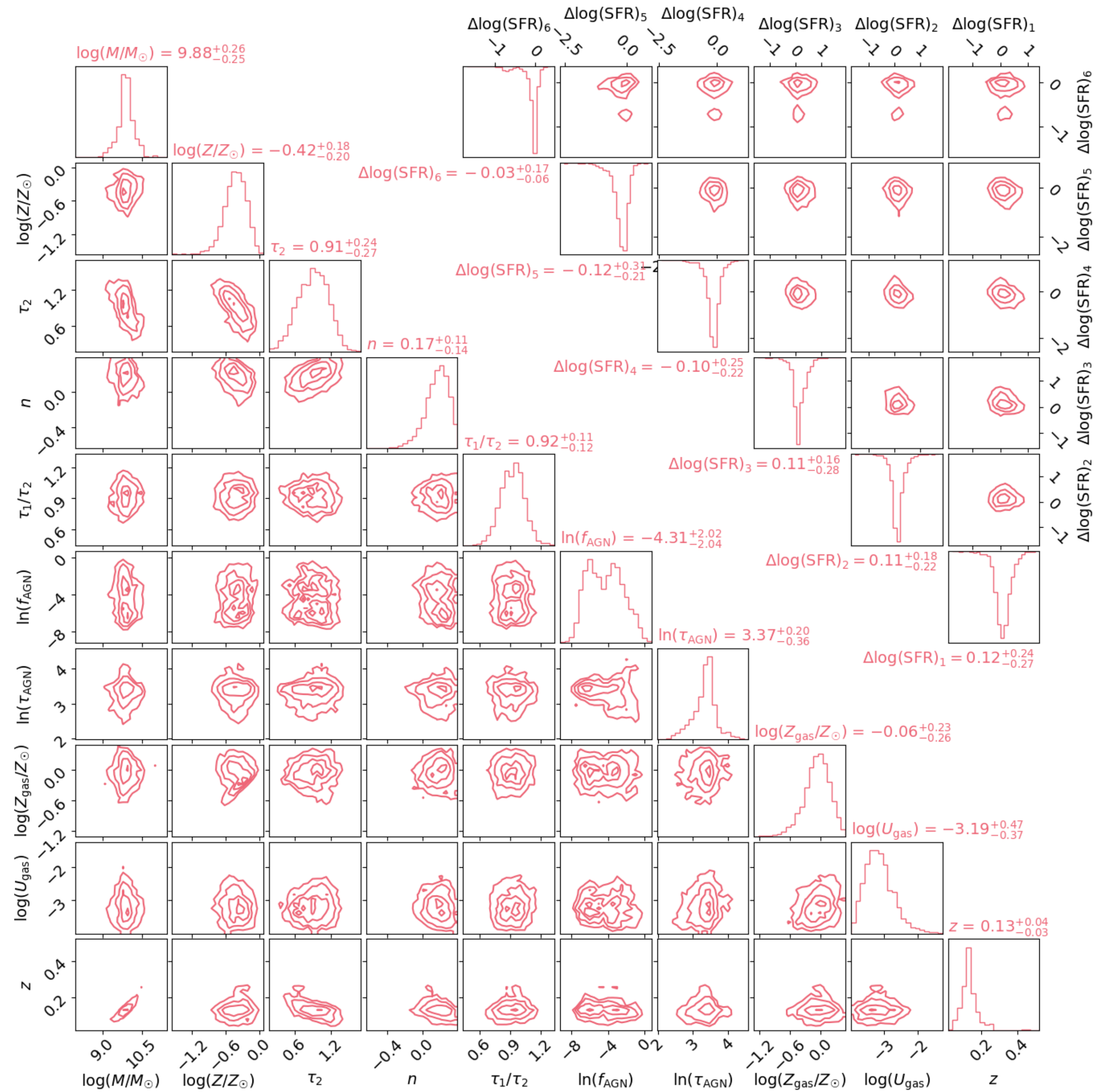
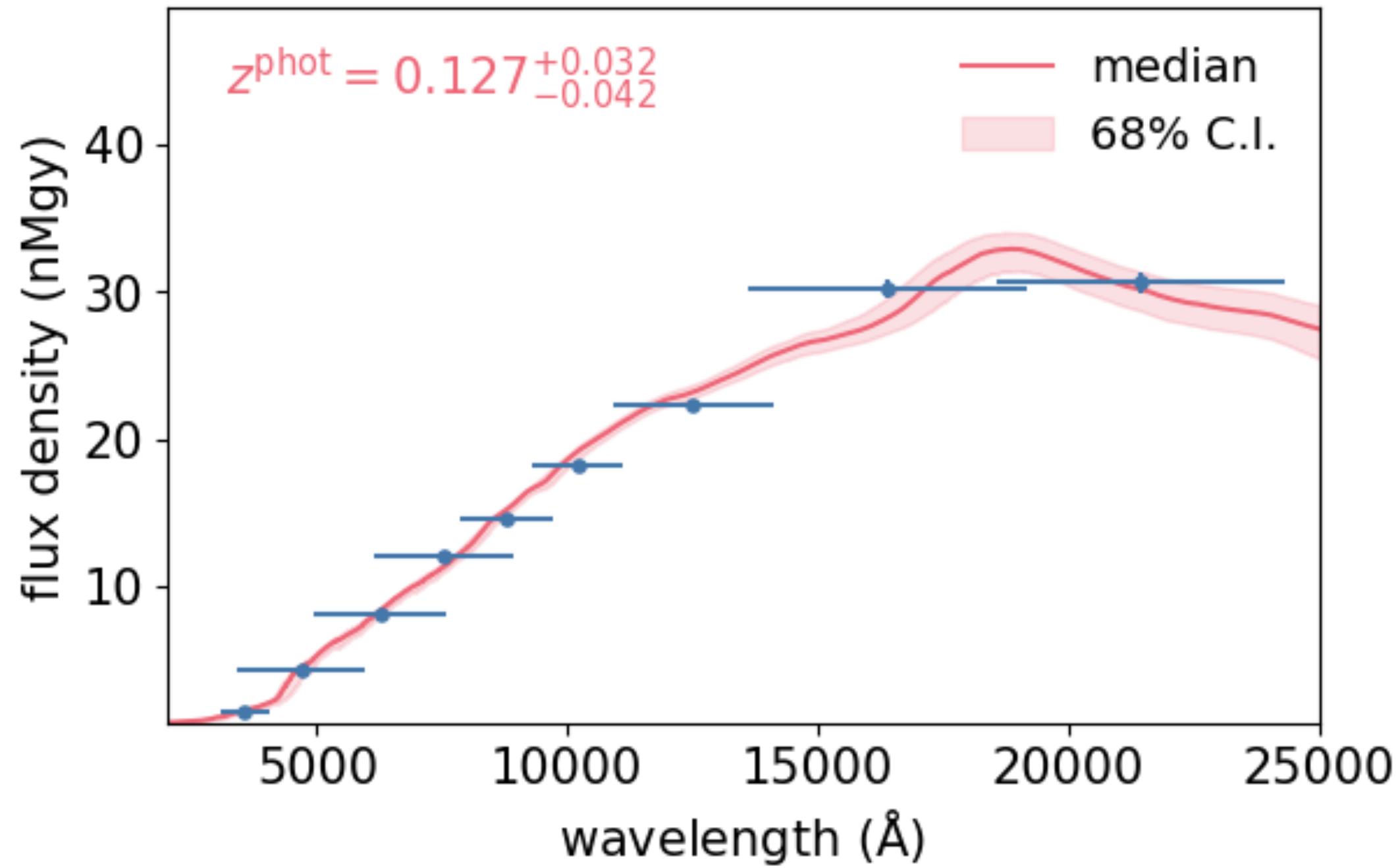
$\sim 85 \text{ deg}$

KiDS-North



SED Fitting

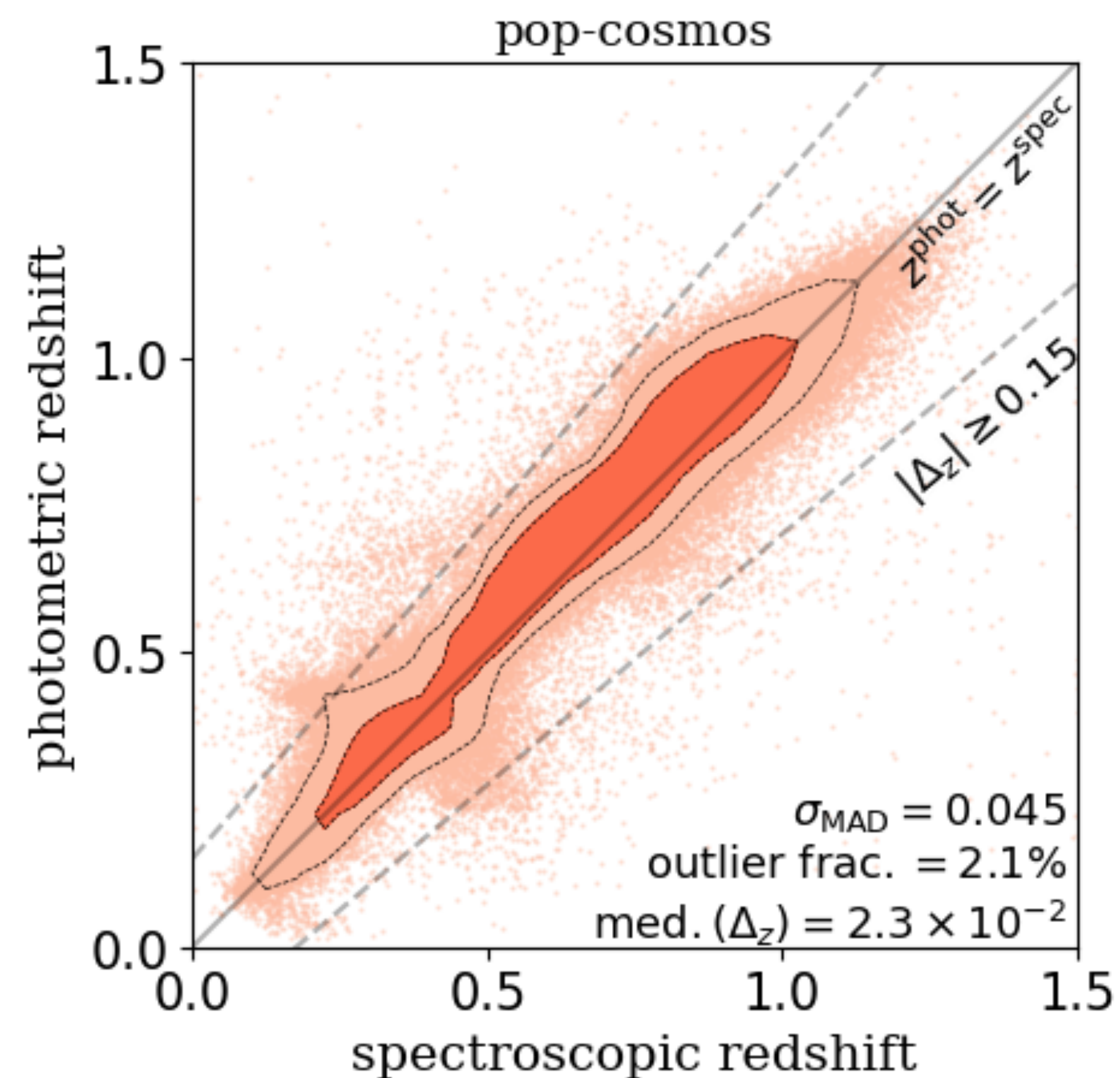
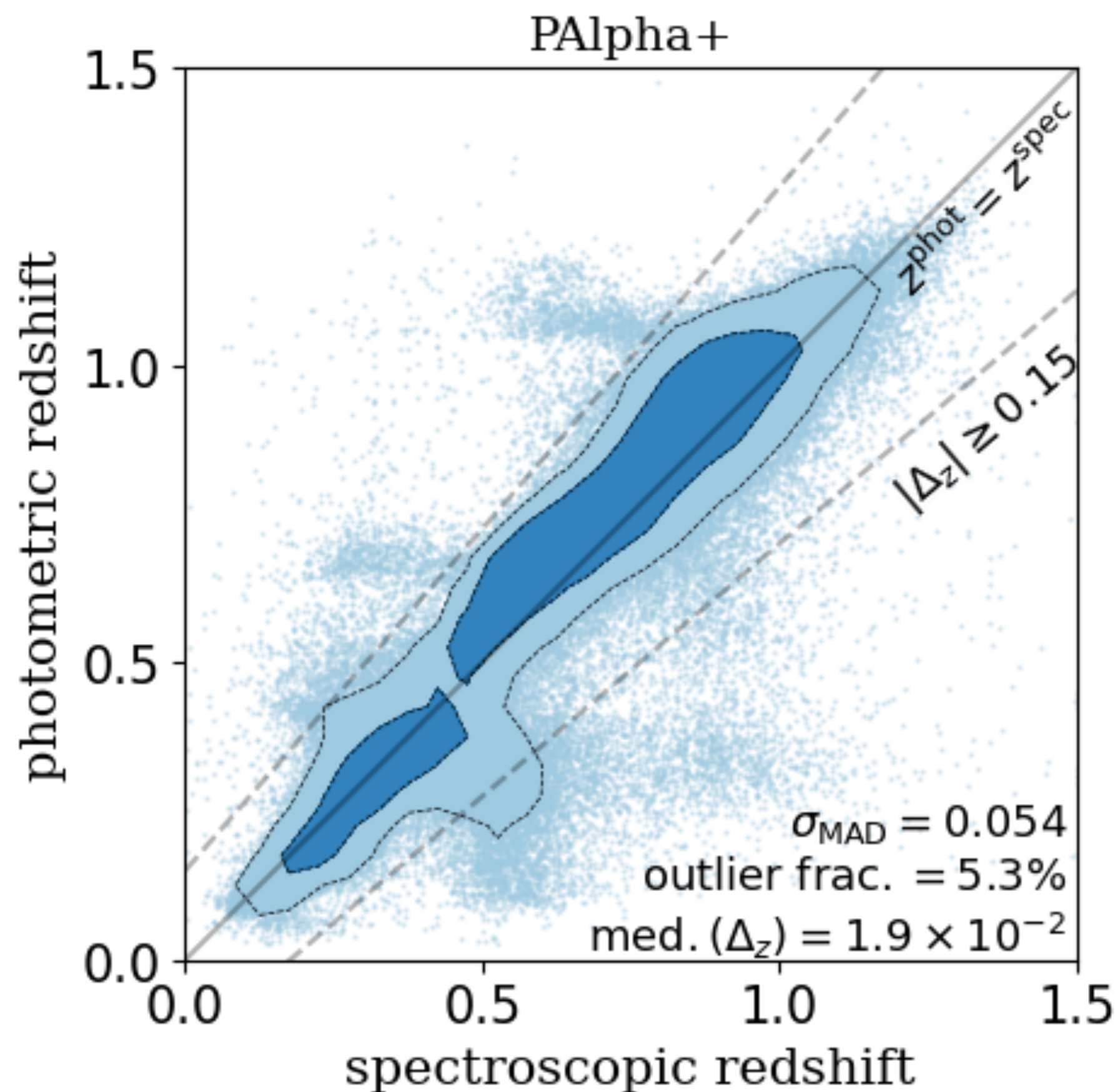
Kilo Degree Survey data...



SED Fitting

Kilo Degree Survey data...

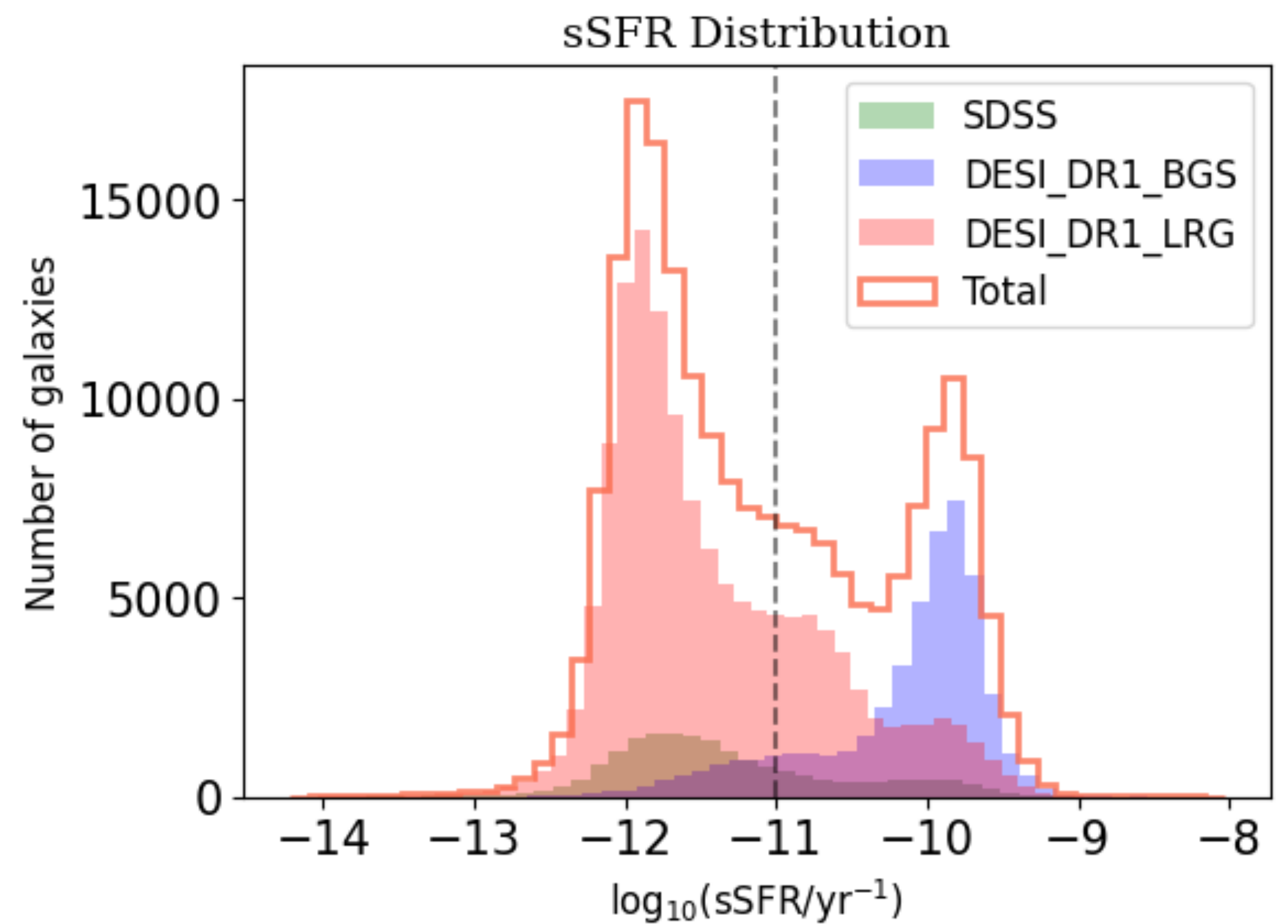
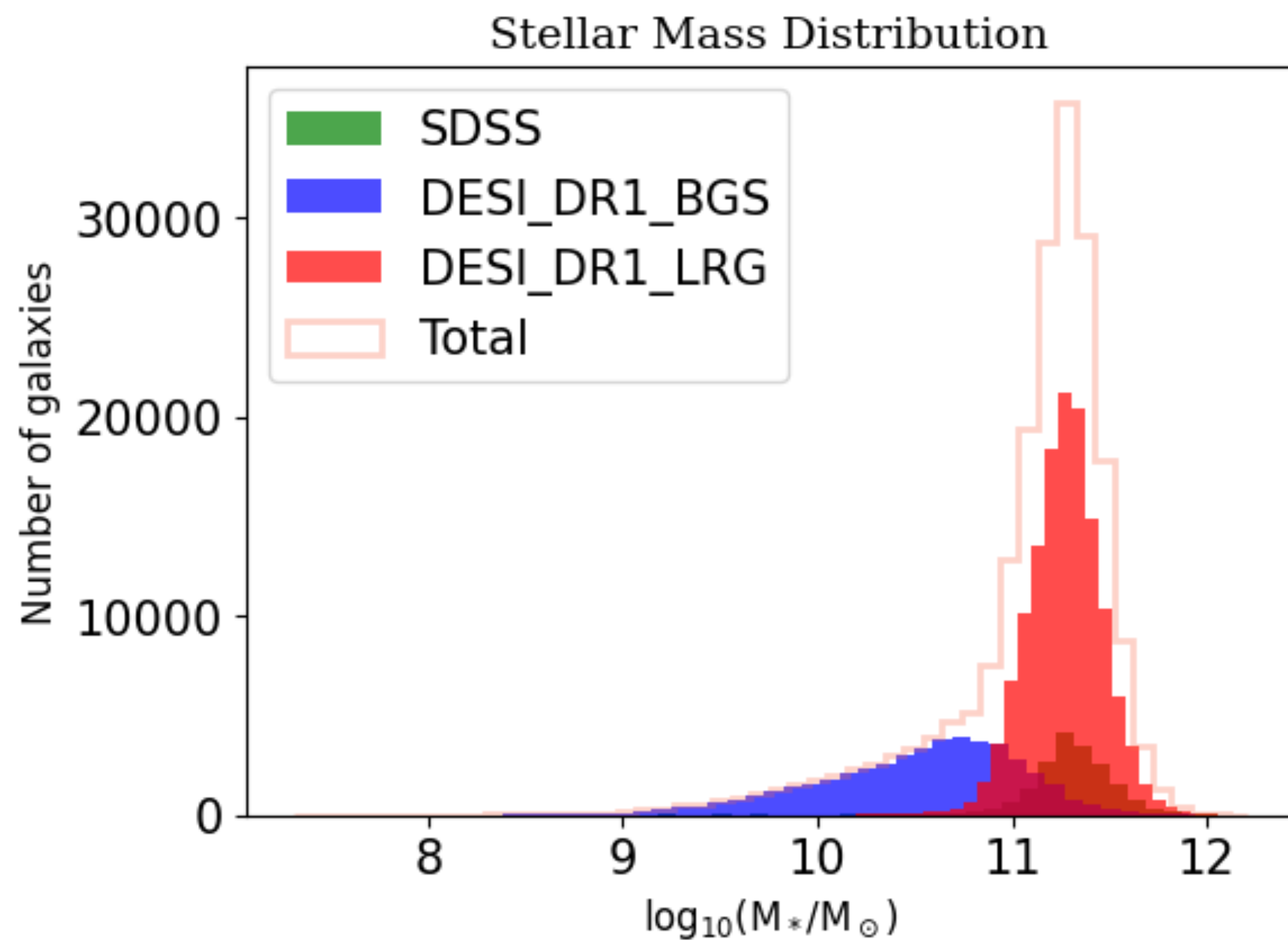
KiDS1000 cross-matched with CDFS, DESI DR1 (BGS+LRG), GAMA, SDSS, VVDS; $N = 202,326$ galaxies



SED Fitting

Kilo Degree Survey data...

KiDS1000 Galaxy Properties by Spectroscopic Sample



Scalability

- Currently using emulated SPS + diffusion prior + affine-invariant sampler
- Gets full MCMC chains with throughput of ~ 7 GPU-sec per galaxy (or ~ 0.7 GPU-sec per galaxy with a simpler analytic prior)
- Good enough for Stage III surveys; could be faster for Stage IV
- Bottleneck is flow/diffusion model log probability; room for improvement there, e.g. distillation/consistency models, ANPE, bespoke solvers (a cool idea from Luigi's highlight talk on Monday!)

Summary

- Data-driven prior over galaxy physical parameters, using a score-based diffusion model
- Train on the deepest, highest-fidelity data that we have
- Use the model as a prior over physical parameters when doing inference for new surveys
- Scaling to ~few millions of galaxies possible with emulation, parallelised MCMC, GPUs
- Bigger data with ANPE, distillation, bespoke solvers?



2506.12122