# Controlling Systematics in Large-Scale Structure Surveys

Noah J. Weaverdyck
University of Michigan
*nweaverd@umich.edu*

Lawrence Berkeley National Laboratory

UNIVERSITY OF MICHIGAN

# Outline

- Background

- Mitigation methods: insights from a common framework

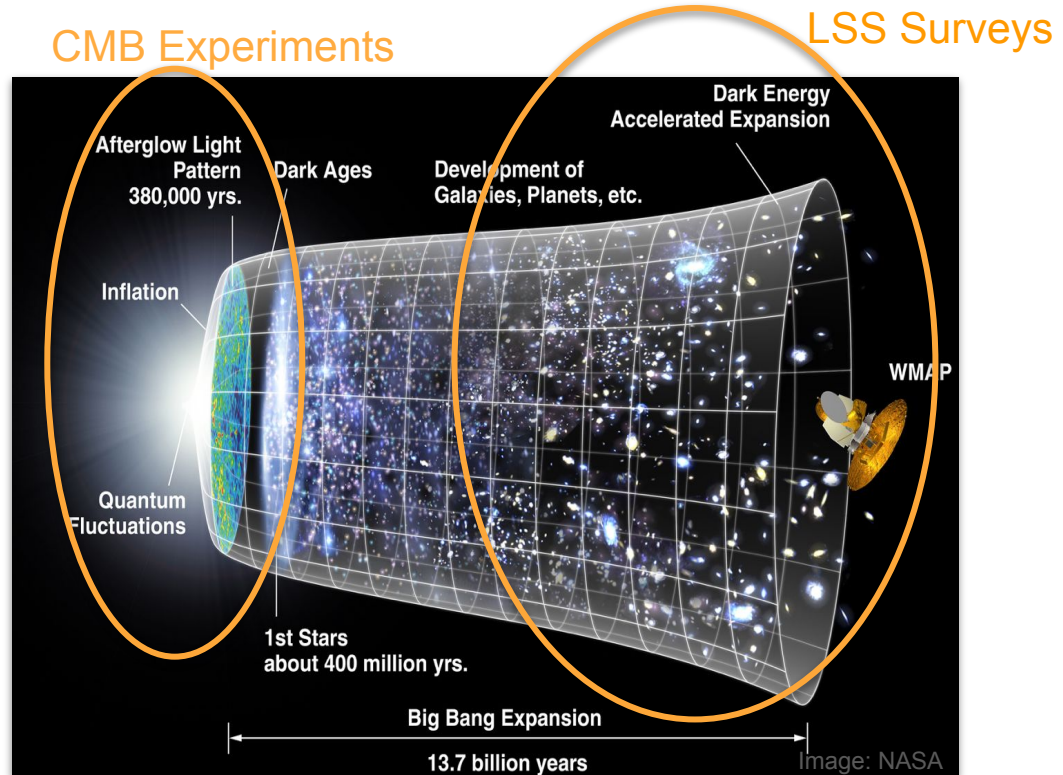- Simulated comparison

- Outlook

Largely based on
Weaverdyck & Huterer (2007.14499)

Some other work not covered in this talk:

- Rapid and generic systematics testing via importance sampling

- Small-scale modeling challenges for constraining inflation via the spectral runnings

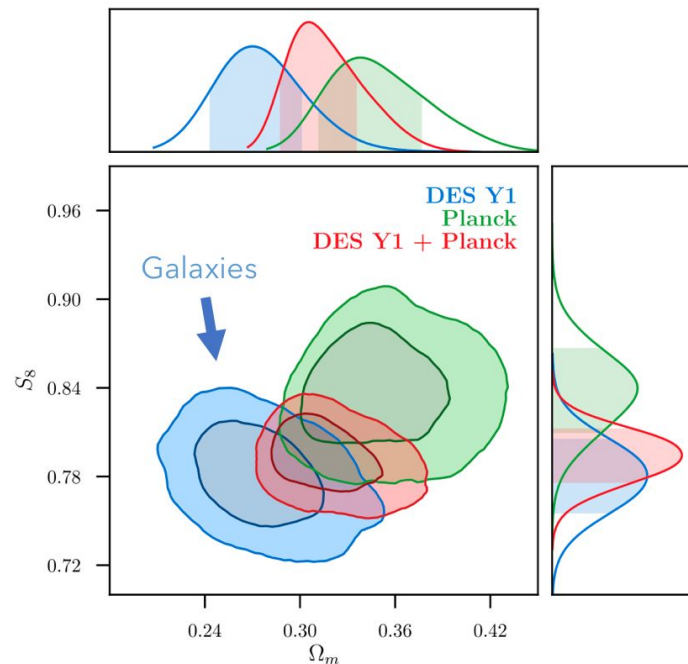- Vetting MCMC samplers for cosmological inference and model testing

# Large-scale structure (LSS) surveys

- Map "late-time" density fluctuations

- Complement primordial fluctuations from CMB

- Probe *expansion history* and *growth of structure*; dark energy, neutrino mass, primordial non-Gaussianity



CMB Experiments

LSS Surveys

Afterglow Light Pattern 380,000 yrs.

Dark Ages

Development of Galaxies, Planets, etc.

Dark Energy Accelerated Expansion

Inflation

WMAP

Quantum Fluctuations

1st Stars about 400 million yrs.

Big Bang Expansion

13.7 billion years
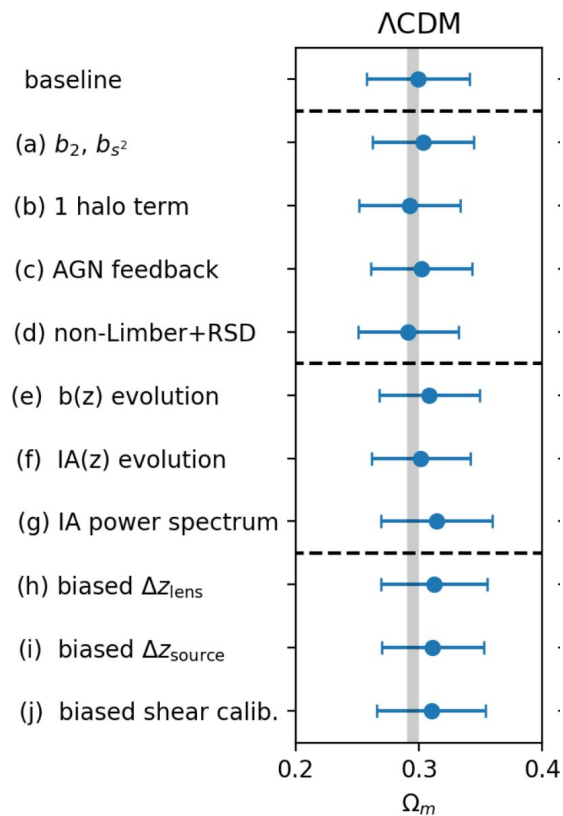
Image: NASA

# LSS surveys

- Primary observables:
  - Galaxy number density → galaxy clustering
  - Galaxy shapes → weak lensing

- Now competitive with CMB

- LSST, DESI, Roman, SphereX...
  **Large** number densities → **small** statistical error
  - Control of systematics paramount to discover new physics



DES Collab (1708.01530)

# LSS systematics

- Galaxy bias

- Small-scale modeling (non-linear Pk)

- Intrinsic alignments

- Photo-z errors

- **_Spatial systematics_**
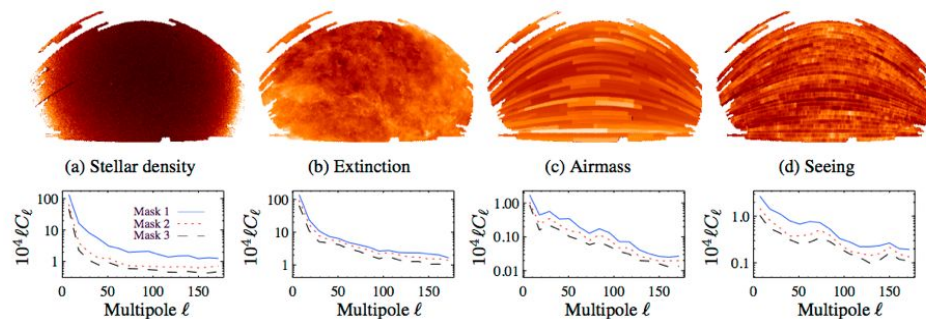  - Modify selection function: map-level



ΛCDM

baseline

(a) $b_2$, $b_{s^2}$

(b) 1 halo term

(c) AGN feedback

(d) non-Limber+RSD

(e) b(z) evolution

(f) IA(z) evolution

(g) IA power spectrum

(h) biased $\Delta z_{lens}$

(i) biased $\Delta z_{source}$

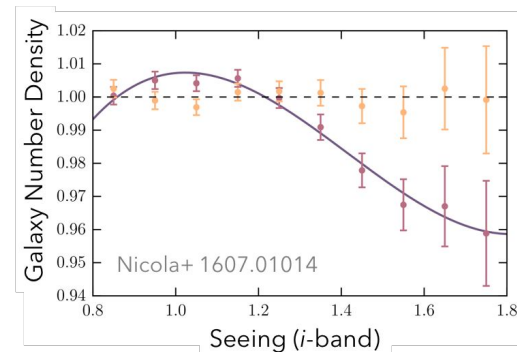(j) biased shear calib.

$\Omega_m$

Krause et al. (DES) 1706.09359

# Spatial systematics

Observed galaxy field ≠ truth

- Astrophysical (stellar contamination, dust extinction, ...)
- Observing conditions (seeing, sky brightness, ...),
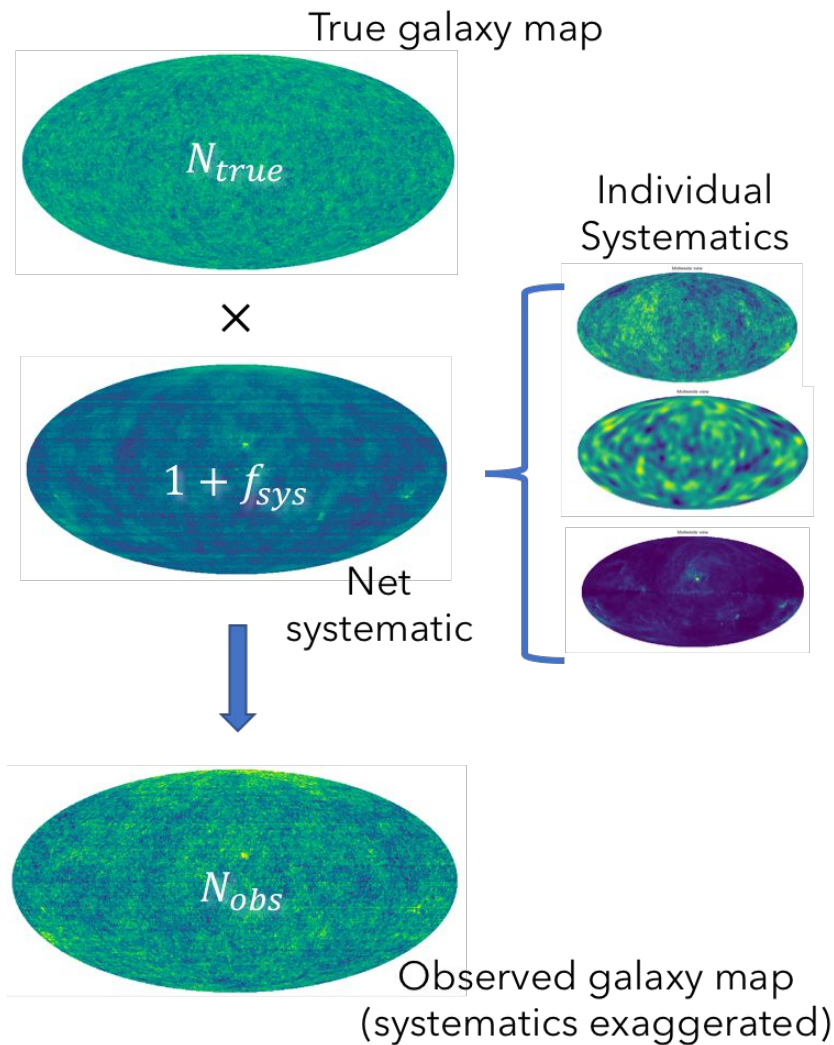- Instrumental (flux calibration, source detection algorithms, ...)



(a) Stellar density   (b) Extinction   (c) Airmass   (d) Seeing
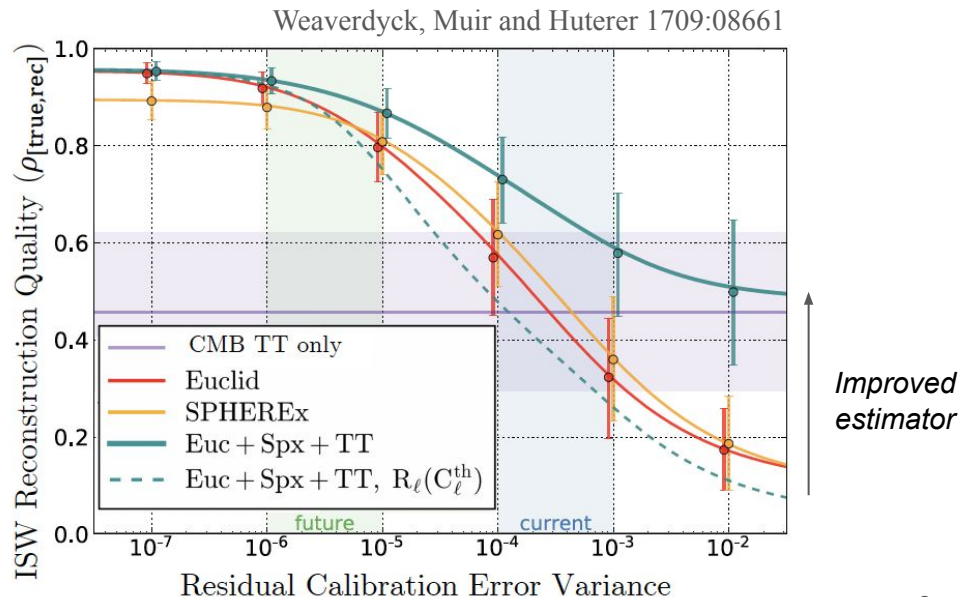
Leistedt et al. 1306.0005

Nicola+ 1607.01014

# Spatial systematics

- Spatially dependent screen ($f_{sys}$) modulates observed galaxy density

- Result: density maps biased! (and 2-pt functions, 3-pt, …)

True galaxy map

$N_{true}$

$\times$

$1 + f_{sys}$

Net systematic

Individual Systematics

$N_{obs}$

Observed galaxy map (systematics exaggerated)

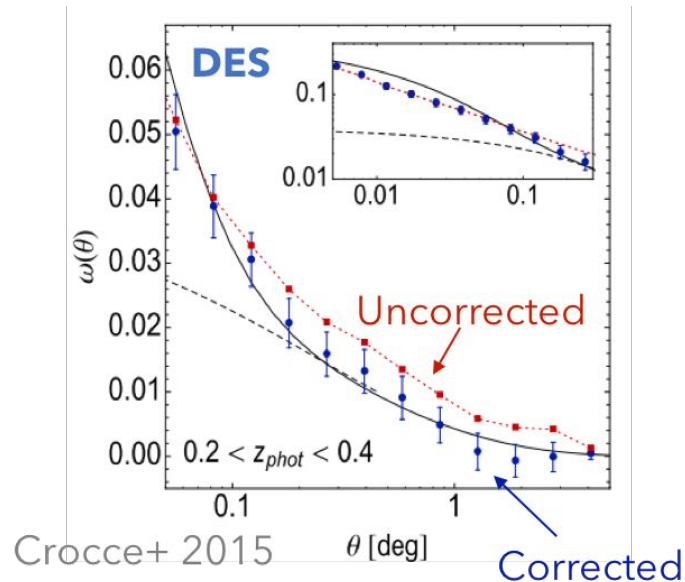# Spatial systematics: impact on ISW

- First PhD project:
  assess control needed for accurate Integrated Sachs Wolfe effect (ISW)

- Leading contribution to CMB at large scales, important for DE/MG

- Infer from x-correlation with LSS

- Optimization, improved estimator for upcoming LSS surveys



Weaverdyck, Muir and Huterer 1709:08661

*Improved estimator*

# How to control systematics?

- Most common: use *systematic templates*, which trace potential contamination
  - Mask extreme regions
  - Estimate and correct for contamination
    (also: Balrog, Obiwan)

- Effects can be large
  - E.g. ELG and QSO densities in DESI imaging: ~10% variation after aggressive masking (Kitanidis et al. 1911.05714)

- Approaches varied, mostly ad hoc
  - Weaverdyck & Huterer (2007.14499):
    compare common methods, establish interpretive framework, improvements



Crocce+ 2015

# How to control systematics?

Prominent methods investigated:

- Mode (De)Projection (e.g. HSC, SDSS QSOs)

- Multiple Linear Regression (e.g. KiDS LRGs, CFHTLenS)

- Template Subtraction (e.g. BOSS LRGs)

- DES-Y1 weighting (DES LRGs)

- "E.Net"                          } New

- "Forward Selection"

All be reformulated as forms of **regression**

# Mode (De)Projection

$$\delta_{\rm obs} \approx \delta_{\rm true} + \alpha t$$
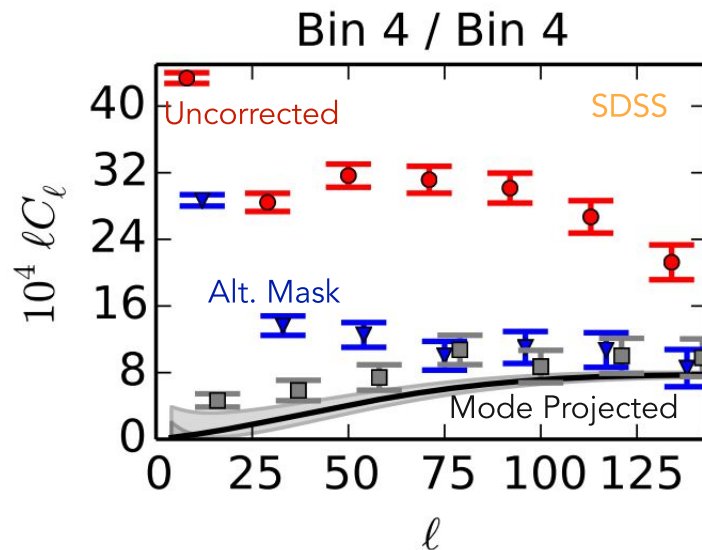
Observed overdensity map

Template map

- Template map *t:*
  Marginalize over **additive** contaminant to overdensity

- *pseudo*-Cl version developed by Elsner+ 2016
  - Avoids $N_{pix}$ x $N_{pix}$ inversion

- Expanded to spin-2 fields
  - Public code *NaMaster* for LSST
    (Alonso+ 2018)

- Equivalent to OLS regression + step to debias *Cl*

Leistedt+ 2015



Bin 4 / Bin 4

Uncorrected     SDSS

Alt. Mask

Mode Projected

$10^4\ \ell C_\ell$

$\ell$

# Mode (De)Projection

$$\delta_{\text{obs}} \approx \delta_{\text{true}} + \alpha t$$

*Multiple systematic templates:*

$$t \to T \quad (\text{N}_{\text{pix}} \times \text{N}_{\text{tpl}})$$

**MP for Pseudo-Cl**

$$\hat{\delta} = \boldsymbol{F}\delta_{\text{obs}}$$

$$= \left[ \lim_{\beta \to \infty} \left( I + \beta t t^{\dagger} \right)^{-1} \right] \delta_{\text{obs}}$$

$$= \left[ I - t(t^{\dagger}t)^{-1}t^{\dagger} \right] \delta_{\text{obs}}$$

$$\hat{\delta} = \delta_{\text{obs}} - t\hat{\alpha}$$

Map estimate

MP estimate of contamination coefficient $\alpha$
Is MLE, *assuming*:

$$\delta \sim \mathcal{N}(0, \, \sigma^2 I)$$

i.e. $\quad \hat{\alpha} = \text{argmin}_{\alpha} ||\delta_{\text{obs}} - T\alpha||^2$

$$y = X\beta + \epsilon$$
$$\hat{\beta} = (X^{\dagger}X)^{-1}X^{\dagger}y$$

*OLS to predict y from X*

$$y = X(X^{\dagger}X)^{-1}X^{\dagger}y + \hat{\epsilon}$$

$$\delta_{\text{obs}} = T[T^{\dagger}T]^{-1}T^{\dagger}\delta_{\text{obs}} + \hat{\delta}$$

$$\hat{\alpha}$$

Actually care about residuals and their clustering

12

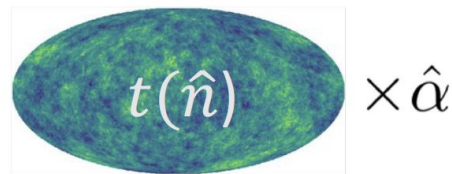$$\delta_{\text{obs}} \approx \delta_{\text{true}} + \alpha t$$

## Template Subtraction

*Decompose template, fit each harmonic*



$$\times \hat{\alpha}_4$$
$$\vdots$$
$$+$$
$$\times \hat{\alpha}_{20}$$
$$\vdots$$
$$+$$
$$\times \hat{\alpha}_{200}$$

$$\sum_m t_{\ell m} Y_{\ell m}(\hat{n})$$

$$\hat{\alpha}_\ell = \frac{\tilde{C}_\ell^{td}}{C_\ell^{tt}}$$

## (PCL) Mode Projection

*Fit full template map*
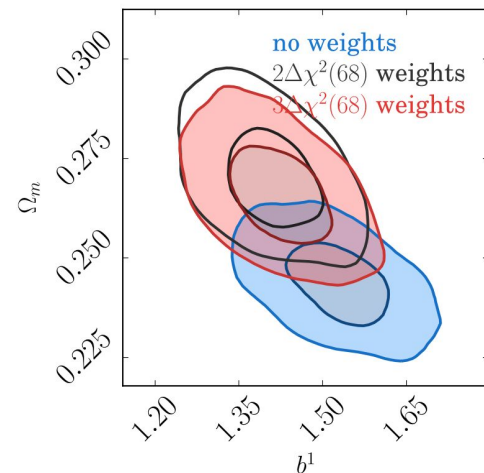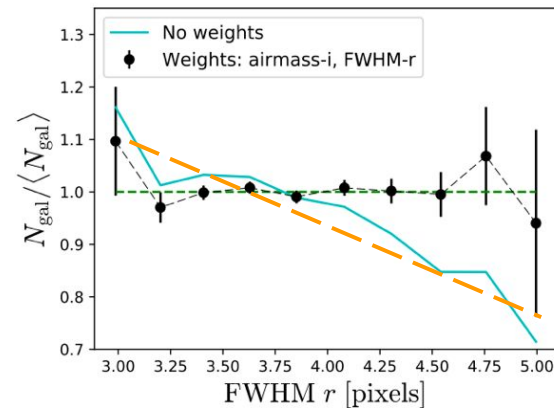
$$t(\hat{n}) \times \hat{\alpha}$$

$$\hat{\alpha} = \frac{\sum_\ell (2\ell + 1)\tilde{C}_\ell^{td}}{\sum_\ell (2\ell + 1)C_\ell^{tt}}$$

13

# "Weights" method (DES-Y1)

- Series of 1D, binned regressions on each template, iteratively reweight galaxies

- Pros vs OLS methods:
  - **Covariance** from mocks,
  - **Significance** threshold to control overfitting

- Cons vs OLS methods:
  - Only detect marginal relationships
  - Computation and time intensive (~1 day)

# Elastic Net Weighting

- Regression extension: form of regularization (Zou & Hastie 2005)
- Incorporate <u>template selection</u>, operate in full-D space

$$\hat{\alpha} = \mathrm{argmin}_{\alpha} \left( ||\delta_{\mathrm{obs}} - T\alpha||^2 + \lambda_1||\alpha||_1 + \lambda_2||\alpha||_2^2 \right)$$

OLS penalty      Sparsity prior (LASSO)      Regularization (Ridge)

*In terms of Maximum Posterior Estimate, equivalent to:*      *Gaussian Likelihood*      *Laplace prior on coefficients*      *Gaussian prior on coefficients*
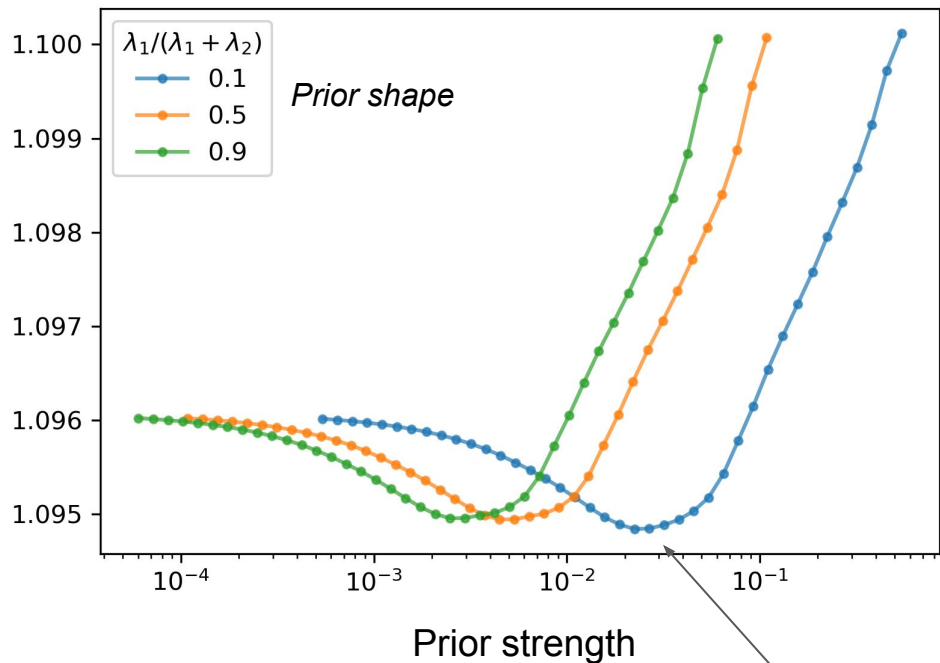
In practice, select $\{\lambda_1, \lambda_2\}$ through cross-validation
(trained on subsets of the data)

# Elastic Net Weighting



Use all templates
(OLS)
High variance

$N_{tpl} = 0$
(no cleaning)
High bias

Average
mean
squared
error on test

$\lambda_1/(\lambda_1 + \lambda_2)$
0.1
0.5
0.9

*Prior shape*

Prior strength

Let data
determine
effective number
of templates

*Optimal hyperparameters*

*Also apply multiplicative correction*

# Multiplicative Correction

$$1 + \delta_{\mathrm{obs}} = (1 + \delta_{\mathrm{true}})(1 + f_{\mathrm{sys}})\gamma$$

$$\delta_{\mathrm{obs}} \approx \underline{\delta_{\mathrm{true}} + f_{\mathrm{sys}}} + \underline{\delta_{\mathrm{true}} f_{\mathrm{sys}}}$$

- Additive estimates (MP, EN, OLS...) leave residual scatter in map
  - Contaminant to small-scale power

- Remove with simple multiplicative correction

$$\hat{\delta} = \frac{\delta_{\mathrm{obs}} - \hat{f}_{\mathrm{sys}}}{1 + \hat{f}_{\mathrm{sys}}}$$

*Next → compare methods on simulation*
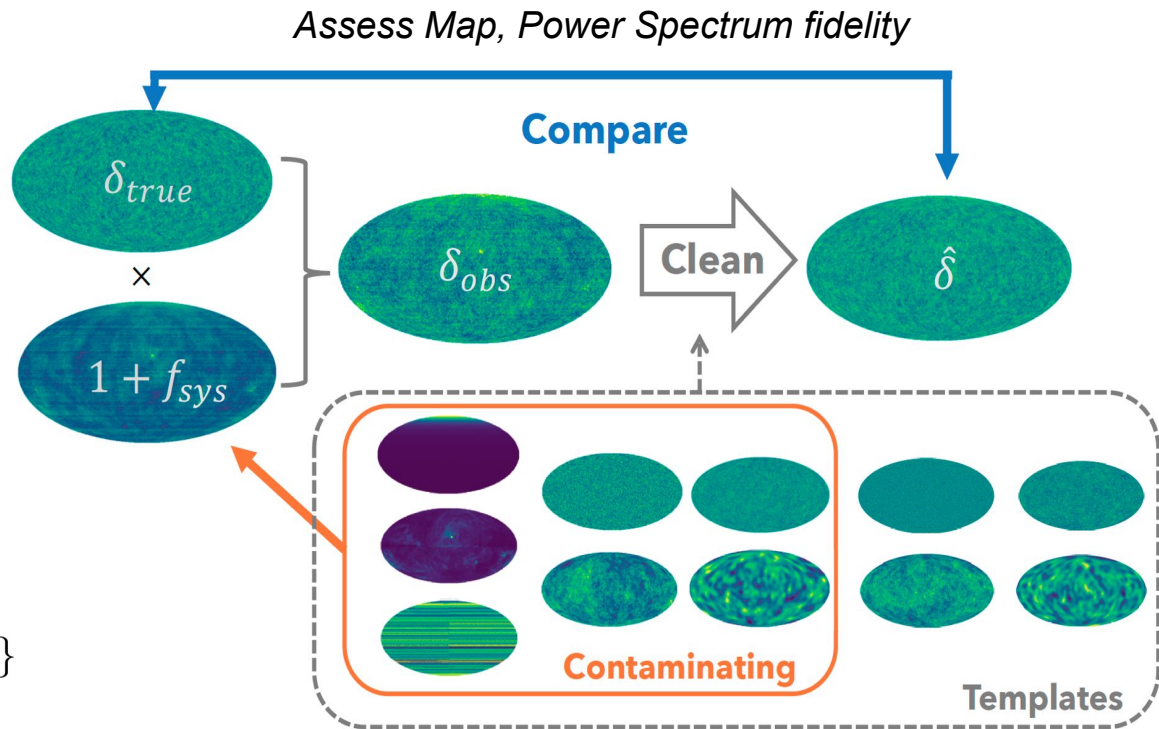


*Each point = map pixel*

17

# Simulation Pipeline

- DES-Y5
- 5 z-bins
- Results not strongly sensitive to survey specs

Templates:

- Gaussian realizations

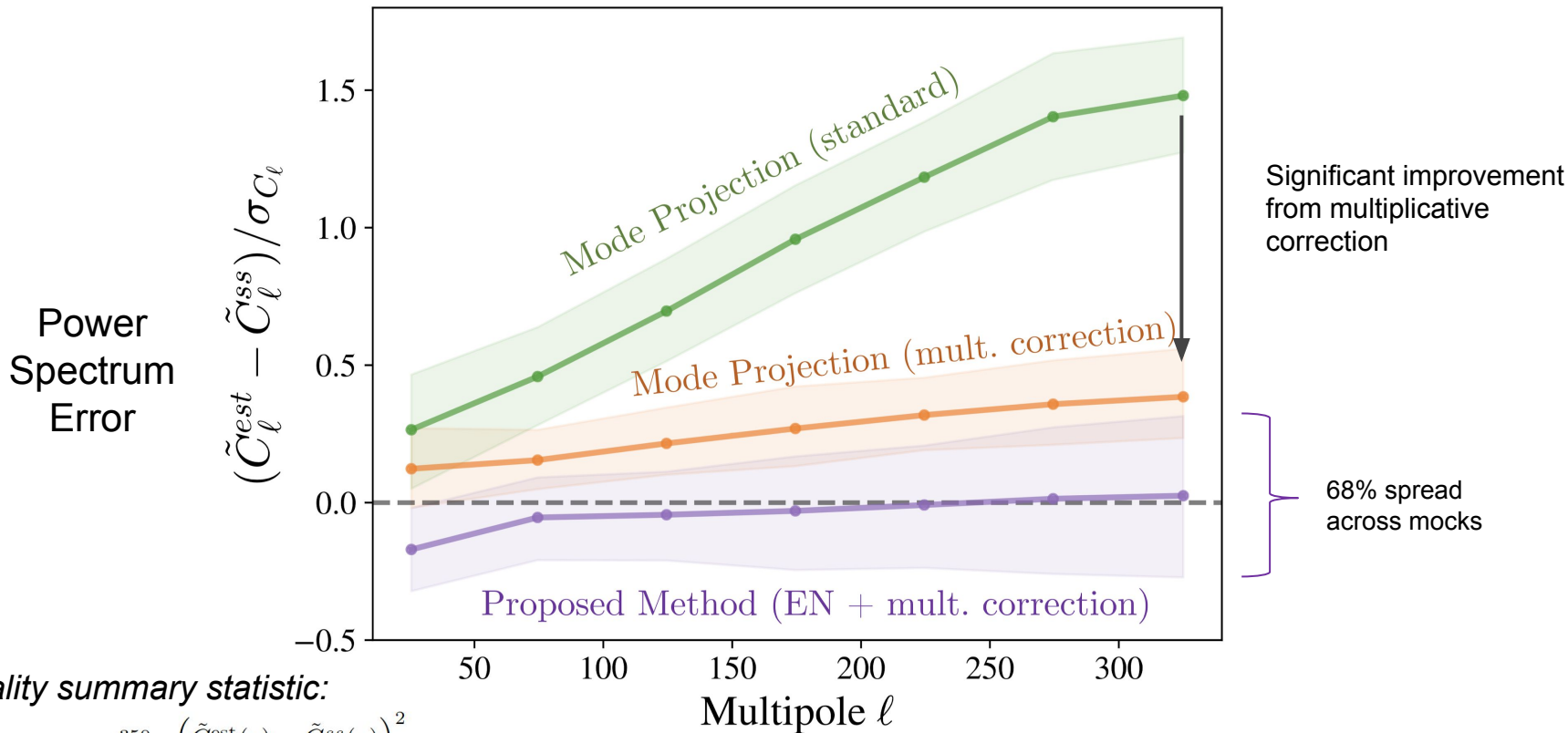$$C_\ell \propto (\ell + 1)^{-p} \qquad p \in \{0, 1, 2\}$$

- Static (Dust, scanning strategy, etc)



*Assess Map, Power Spectrum fidelity*

**Compare**

$\delta_{true}$ × $1 + f_{sys}$ → $\delta_{obs}$ **Clean** → $\hat{\delta}$

**Contaminating**

**Templates**

Note: Methods applicable to any contaminated signal with templates. Here galaxy clustering, with signal = galaxy overdensity.
Generically:  $\delta_{true} \rightarrow s, \quad \delta_{obs} \rightarrow d_{obs}$
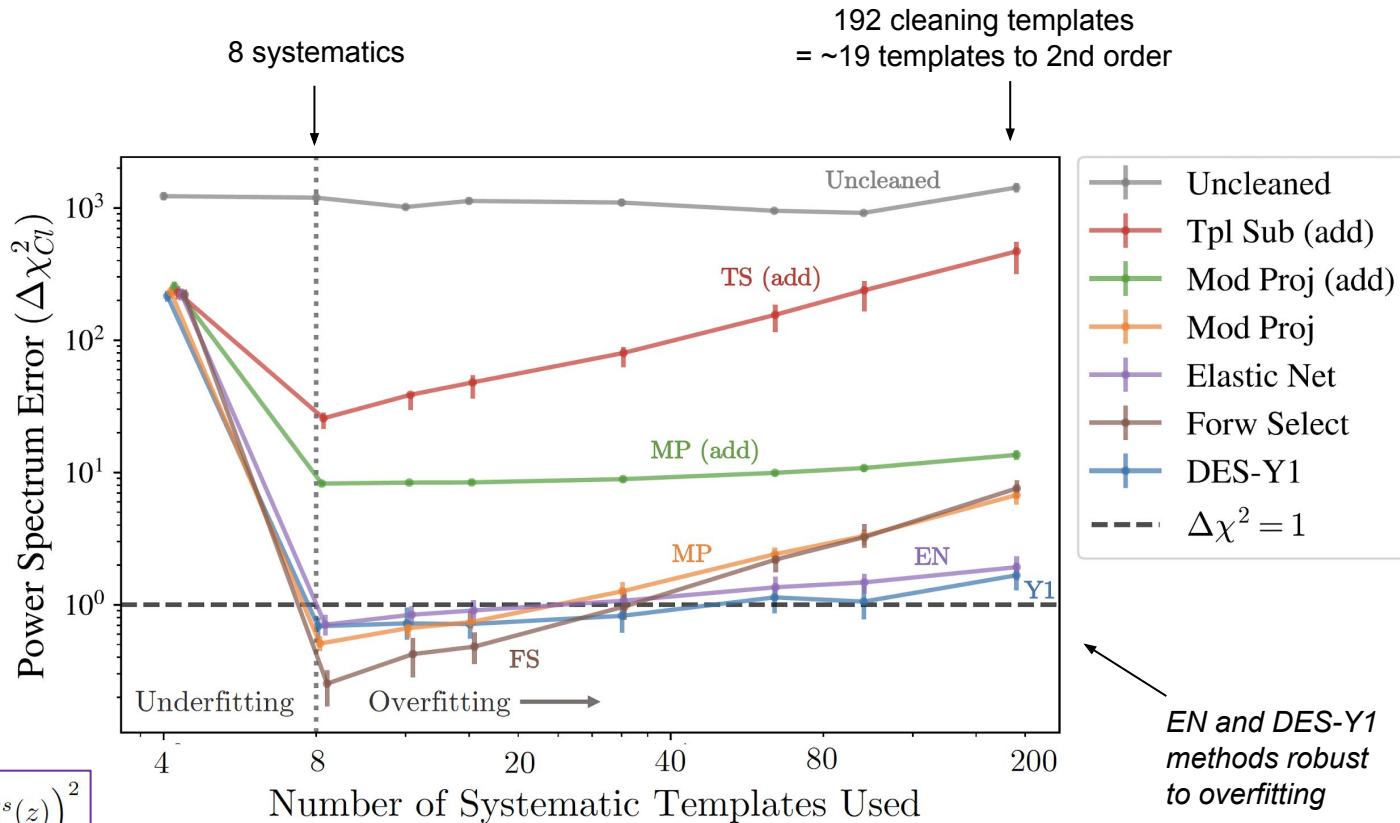
# Importance of Multiplicative Correction



Power Spectrum Error

Significant improvement from multiplicative correction

68% spread across mocks

*Quality summary statistic:*

$$\Delta\chi^2_{C_\ell} = \sum_{z\,\text{bins}} \sum_{\ell=\ell_{\min}}^{350} \frac{\left(\tilde{C}_\ell^{\text{est}}(z) - \tilde{C}_\ell^{ss}(z)\right)^2}{\sigma^2_{C_\ell^{ss}(z)}},$$
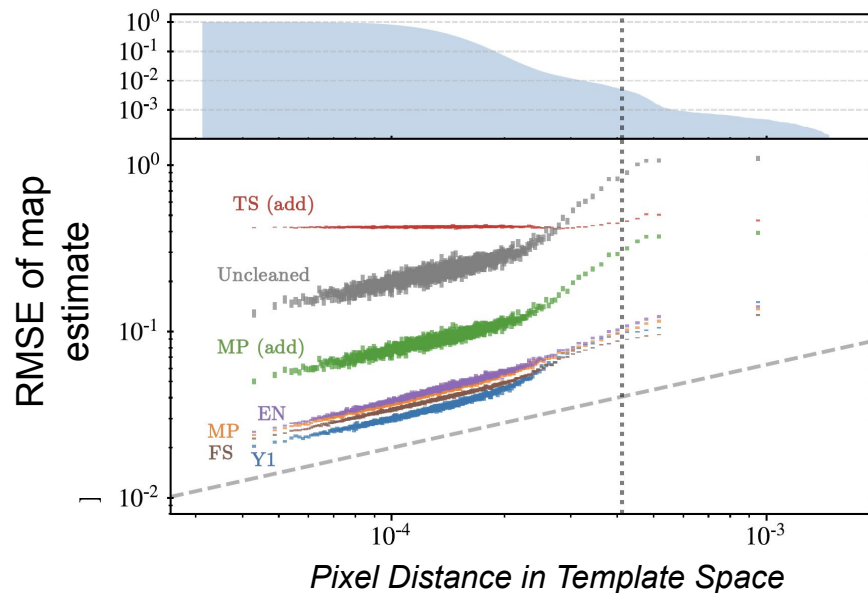
# Importance of Data-driven Template Selection



$$\Delta\chi^2_{\mathrm{C}_\ell} = \sum_{z\mathrm{bins}} \sum_{\ell=\ell_{\min}}^{350} \frac{\left(\tilde{C}^{\mathrm{est}}_\ell(z) - \tilde{C}^{ss}_\ell(z)\right)^2}{\sigma^2_{C^{ss}_\ell(z)}},$$

# Further development

- Mask optimization with template map-statistics

- Scale-optimized cleaning
  - Harmonic prewhitening
  - Maximize S/N for cosmology

- Systematics mitigation for primordial non-Gaussianity ($f_{NL}$)
  - Key target of LSS
  - Cleaning large scales *crucial* (e.g. Castorina et al 2019)

# Outlook

- Common framework unleashes new, powerful tools for systematics mitigation
    - Supervised learning/regression with *residuals* and clustering as signal of interest

- Corrections at *both* map and 2-pt function level

- Mask is important
  → rapid mitigation enables iteration

- Template selection should be *data-driven*
    - Self-calibrated sparsity + shrinkage priors work well!

**Thank you!**



*R. Hahn*

# MP Assumptions on Noise

- True clustering signal = regression "noise"

  Only optimal if clustering signal

  1) Gaussian

  2) Diagonal

  3) Flat

  Can estimate $\alpha$ in pixel space or harmonic space $\quad \hat{\alpha} = [T^\dagger T]^{-1} T^\dagger \delta_{\text{obs}}$

*Diagonalize and optimally weight in harmonic space*

$$\hat{\alpha} = \frac{\sum_{\ell=0}^{\ell_{\max}} (2\ell+1) \tilde{C}_\ell^{\,td}/C_\ell^{ss}}{\sum_{\ell=0}^{\ell_{\max}} (2\ell+1) \tilde{C}_\ell^{\,tt}/C_\ell^{ss}}.$$
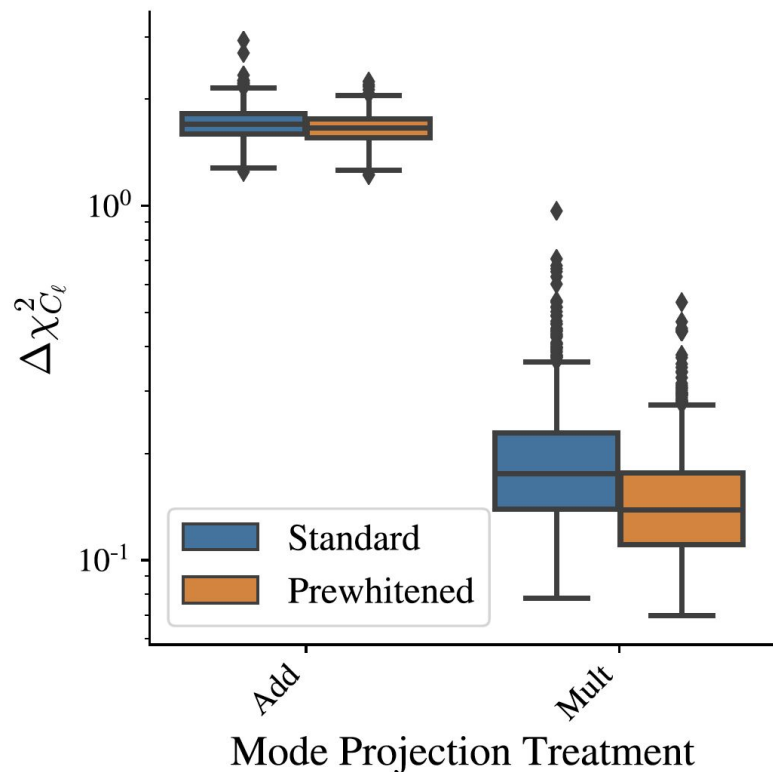
|  | **Pixel Space** | **Harmonic Space** |
|---|---|---|
| **Data** | $\delta_{obs}(\hat{n}_i)$ | $[\delta_{obs}]_{\ell m}$ |
| **Dims of $T$** | $N_{pix} \times N_{tpl}$ (real) | $N_{\ell m} \times N_{tpl}$ (complex) |
| **Regression Noise (additive)** | $\delta(\hat{n}_i)$ | $\delta_{\ell m}$ |
| **Gaussian** | Approx. (~lognormal) | Yes |
| **Diagonal** | No | Yes |
| **Flat** | Yes | No |

$(d_{\text{obs}})'_{\ell m} = (d_{\text{obs}})_{\ell m}/\sqrt{C_\ell^{ss}}.$

$(t_i)'_{\ell m} = (t_i)_{\ell m}/\sqrt{C_\ell^{ss}},$

24

Impact of pixel covariance
*Minor compared to methodological differences*.

No method particularly susceptible
to Gaussian assumption