

## CHARACTERIZING SIGNAL LOSS, ERROR, AND BIAS IN THE 21 CM REIONIZATION POWER SPECTRUM: A REVISED STUDY OF PAPER-64

CARINA CHENG<sup>1,◊</sup>, AARON R. PARSONS<sup>1,2</sup>, MATTHEW KOLOPANIS<sup>3</sup>, DANIEL C. JACOBS<sup>3</sup>, ADRIAN LIU<sup>1,4,†</sup>, SAUL A. KOHN<sup>5</sup>,  
 JAMES E. AGUIRRE<sup>5</sup>, JONATHAN C. POBER<sup>6</sup>, ZAKI S. ALI<sup>1</sup>, GIANNI BERNARDI<sup>7,8,9</sup>, RICHARD F. BRADLEY<sup>10,11,12</sup>, CHRIS L.  
 CARILLI<sup>13,14</sup>, DAVID R. DEBOER<sup>2</sup>, MATTHEW R. DEXTER<sup>2</sup>, JOSHUA S. DILLON<sup>1,\*</sup>, PAT KLIMA<sup>11</sup>, DAVID H. E. MACMAHON<sup>2</sup>,  
 DAVID F. MOORE<sup>5</sup>, CHUNEETA D. NUNHOKEE<sup>8</sup>, WILLIAM P. WALBRUGH<sup>7</sup>, ANDRE WALKER<sup>7</sup>

<sup>1</sup>Astronomy Dept., U. California, Berkeley, CA

<sup>2</sup>Radio Astronomy Lab., U. California, Berkeley CA

<sup>3</sup>School of Earth and Space Exploration, Arizona State U., Tempe AZ

<sup>4</sup>Berkeley Center for Cosmological Physics, Berkeley, CA

<sup>5</sup>Dept. of Physics and Astronomy, U. Penn., Philadelphia PA

<sup>6</sup>Dept. of Physics, Brown University, Providence RI

<sup>7</sup>Square Kilometer Array, S. Africa, Cape Town South Africa

<sup>8</sup>Dept. of Physics and Electronics, Rhodes University, South Africa

<sup>9</sup>INAF-Instituto di Radioastronomia, Bologna Italy

<sup>10</sup>Dept. of Electrical and Computer Engineering, U. Virginia, Charlottesville VA

<sup>11</sup>National Radio Astronomy Obs., Charlottesville VA

<sup>12</sup>Dept. of Astronomy, U. Virginia, Charlottesville VA

<sup>13</sup>National Radio Astronomy Obs., Socorro NM

<sup>14</sup>Cavendish Lab., Cambridge UK

### ABSTRACT

The Epoch of Reionization (EoR) is an uncharted era in our Universe’s history during which the first stars and galaxies led to the ionization of neutral hydrogen in the intergalactic medium. There are many experiments investigating the EoR by tracing the 21 cm line of neutral hydrogen, a signal which is very faint and difficult to isolate. With a new generation of instruments and a statistical power spectrum detection in our foreseeable future, it has become increasingly important to develop techniques that help maximize sensitivity while validating results. Additionally, it is imperative to understand the trade-offs between different analysis methods and their effects on power spectrum estimates. In this paper, we focus on three major themes — signal loss, power spectrum error estimation, and bias in measurements. We describe techniques that affect these themes using both toy models and data taken by the 64-element configuration of the Donald C. Backer Precision Array for Probing the Epoch of Reionization (PAPER). In particular, we highlight how detailed investigations of these themes have led to a revised, higher 21 cm power spectrum upper limit from PAPER-64. While we only focus on a subset of PAPER-64 data in this paper, revised power spectrum limits from the PAPER experiment are presented in a companion paper by Kolopanis et al. (*in prep.*) and supersede results from previously published PAPER analyses.

### 1. INTRODUCTION

By about one billion years after the Big Bang ( $z \sim 6$ ), the first stars and galaxies are thought to have ionized all the neutral hydrogen that dominated the baryonic matter content in the Universe. This transition period, during which the first luminous structures formed from

gravitational collapse and began to emit intense radiation that ionized the cold neutral gas into a plasma, is known as the Epoch of Reionization (EoR). The EoR is a relatively unexplored era in our *cosmic dawn*, which spans the birth of the first stars to the full reionization of the intergalactic medium (IGM). Its history encodes important information regarding the nature of the first galaxies and the processes of structure formation. Direct measurements of the EoR would unlock powerful characteristics about the IGM, revealing connections between the matter distribution exhibited via cosmic microwave background (CMB) studies and the highly structured web of galaxies we observe today (for a review, see

◊ccheng@berkeley.edu

†Hubble Fellow

\*NSF AAPF Fellow

Barkana & Loeb (2001), Furlanetto et al. (2006) and Loeb & Furlanetto (2013)).

One promising technique to probe the EoR is to target the 21 cm wavelength signal that is emitted and absorbed by neutral hydrogen via its spin-flip transition (Furlanetto et al. 2006; Barkana & Loeb 2008; Morales & Wyithe 2010; Pritchard & Loeb 2010; Pritchard & Loeb 2012). This technique is powerful because it can be observed both spatially and as a function of redshift — that is, the wavelength of the signal reaching our telescopes can be directly mapped to a distance from where the emission originated before stretching out as it traveled through expanding space. Hence, 21 cm tomography offers a unique window into both the spatial and temporal evolution of ionization, temperature, and density fluctuations.

In addition to the first tentative detection of the EoR signal made by the Experiment to Detect the Global EoR Signature (EDGES; Bowman et al. 2018; Bowman & Rogers 2010), there are several radio telescope experiments that have succeeded in using the 21 cm signal from hydrogen to place constraints on the brightness of the signal. Examples of experiments investigating the mean brightness temperature of the EoR relative to the CMB are the Large Aperture Experiment to Detect the Dark Ages (LEDA; Bernardi et al. 2016), the Dark Ages Radio Explorer (DARE; Burns et al. 2012), the Sonda Cosmológica de las Islas para la Detección de Hidrógeno NeutroSciHi (SCI-HI; Voytek et al. 2014), the Broadband Instrument for Global HyDrOgen Reionisation Signal (BIGHORNS; Sokolowski et al. 2015), and the Shaped Antenna measurement of the background RAdio Spectrum (SARAS; Patra et al. 2015). Radio interferometers which seek to measure statistical power spectra include the Giant Metre-wave Radio Telescope (GMRT; Paciga et al. 2013a), the LOw Frequency ARray (LOFAR; van Haarlem et al. 2013), the Murchison Widefield Array (MWA; Tingay et al. 2013), the 21 Centimeter Array (21CMA; Peterson 2004; Wu 2009), the Square Kilometre Array (SKA; Koopmans et al. 2015), and PAPER (Parsons et al. 2010). The Hydrogen Epoch of Reionization Array (HERA), which is currently being built, is a next-generation instrument that aims to combine lessons learned from previous experiments and is forecast to be able to make a high-significance power spectrum detection with an eventual 350 elements using current analysis techniques (Pober et al. 2014; Liu & Parsons 2016; Dillon & Parsons 2016; DeBoer et al. 2017).

The major challenge that faces all 21 cm experiments is isolating a small signal that is buried underneath foregrounds and instrumental systematics that are, when combined, four to five orders of magnitude brighter (e.g., Santos et al. 2005; Ali et al. 2008; de Oliveira-Costa et al. 2008; Jelić et al. 2008; Bernardi et al. 2009, 2010; Ghosh et al. 2011; Pober et al. 2013; Bernardi et al. 2013; Dillon et al. 2014; Kohn et al. 2016). A clean measurement therefore requires an intimate understanding

of the instrument and a rigorous study of data analysis choices. With continual progress being made in the field and HERA on the horizon, it is becoming increasingly important to understand how the methods we choose interact with each other to affect power spectrum results. More specifically, it is imperative to develop techniques and tests that ensure the accuracy and reliability of a potential EoR detection. In this paper, we discuss three topics (signal loss, error estimation, and bias) that are essential to investigate for a robust 21 cm power spectrum analysis. We also highlight four power spectrum techniques (fringe-rate filtering, weighting, bootstrapping, jackknife testing) and their trade-offs, potential pitfalls, and connections to the themes. We first approach the themes from a broad perspective, and then perform a detailed case study using data from the 64-element configuration of PAPER. In this study we use a subset of PAPER-64 data to illustrate our revised analysis methods, while a companion paper, Kolopanis et al. (*in prep.*), builds off of the methods in this paper to present revised PAPER-64 results for multiple redshifts and baseline types.

Finally, this paper adds to the growing foundations of lessons which have been documented, for example, in Paciga et al. (2013b), Patil et al. (2016), and Jacobs et al. (2016), by the GMRT, LOFAR, and MWA projects respectively. These lessons are imperative as the community as a whole moves towards higher sensitivities and potential EoR detections.

This paper is organized into two main halves. In Section 2 we introduce the three themes of our focus, using toy models to develop intuition for each one. In Section 3 we present a case study into each theme using data from the PAPER-64 array, highlighting key changes from the methods used in the previously published result in Ali et al. (2015), henceforth known as A15, which have led to a revised PAPER-64 power spectrum result (Kolopanis et al. (*in prep.*)). We conclude in Section 4.

## 2. POWER SPECTRUM THEMES AND TECHNIQUES

There are many choices a 21 cm data analyst must consider. How can time-ordered measurements be combined? How can the variance of the data be estimated? In what way(s) can the data be weighted to suppress contaminated modes while not destroying an EoR signal? How can a statistically significant detection of a signal be properly identified? Many common techniques, such as averaging data, weighting, bootstrapping, and jackknife testing, address these issues but harbor additional trade-offs. For example, an aggressive filtering method may succeed in eliminating interfering systematics but comes at the cost of losing some EoR signal. A chosen weighting scheme may theoretically maximize sensitivity but fail to suppress foregrounds in practice.

Though there are many data analysis choices, measuring the statistical 21 cm power spectrum ultimately requires robust methods for determining accurate confi-

dence intervals and rigorous techniques to identify and control systematics. In this paper, we focus on three 21 cm power spectrum themes that encapsulate this goal and discuss four techniques that interplay with each other and impact the themes. We will give brief definitions now, and build intuition for each theme in the sections to follow.

- **Signal Loss** (Section 2.1): Signal loss refers to attenuation of the target *cosmological* signal in a power spectrum estimate. Certain analysis techniques can cause this loss, and if the amount of loss is not quantified accurately, it could lead to false non-detections and overly aggressive upper limits. Determining whether an analysis pipeline is lossy, and estimating the amount of loss if so, has subtle challenges but is necessary to ensure the accuracy of any result.
- **Error Estimation** (Section 2.2): Confidence intervals on a 21 cm power spectrum result determine the difference between a detection and a null result, which have two very different implications. Additionally, accurate error estimation is crucial for the comparison of results to theoretical models. Errors can be estimated in a variety of ways, and we will discuss a few of them.
- **Bias** (Section 2.3): There are several possible sources of power offset in a visibility measurement that can show up as a detection in a power spectrum, such as bias from noise and foregrounds. In particular, a successful EoR detection would also imitate a bias. Proving that a bias is an EoR detection may be the most difficult challenge for future 21 cm analyses, as it is crucial to be able to distinguish a detection of foreground leakage, for example, from that of EoR. In this paper we will highlight some sources of bias, discuss ways to mitigate their effects, and describe example tests that a true EoR detection must pass.

The following techniques each have advantages when it comes to maximizing sensitivity and understanding systematics in data. However, some have limitations, and we will discuss circumstances in which there are trade-offs. We choose to focus on these four techniques because they represent major steps in PAPER’s power spectrum pipeline, with several of them also being standard steps in general 21 cm analyses.

- **Fringe-rate filtering**: Fringe-rate filtering is an averaging scheme for time-ordered data (Parsons et al. 2016). Broadly, a fringe-rate filter averages visibilities in time to produce a smaller number of more sensitive independent samples. However, such a filter also affects the presence of foregrounds and systematics. We explain the trade-offs of filtering in more detail in Section 2.1.3.

- **Weighting**: A dataset can be weighted to emphasize certain features and minimize others. One particular flavor of weighting employed by previous PAPER analyses is inverse covariance weighting in frequency, which is a generalized version of inverse variance weighting that also takes into account frequency correlations (Liu & Tegmark 2011; Dillon et al. 2013; Liu et al. 2014a; Liu et al. 2014b; Dillon et al. 2014; Dillon et al. 2015). Using such a technique enables the down-weighting of contaminant modes that obey a different covariance structure from that of cosmological modes. However, a challenge of inverse covariance weighting is in estimating a covariance matrix that is closest to the true covariance of the data; the discrepancy between the two can have large impacts on signal loss. We investigate the impact of different types of weighting on signal loss in Section 2.1.
- **Bootstrapping**: In addition to using theoretical models for covariance matrices and theoretical error estimation methods, bootstrapping is one way to estimate errors. Namely, bootstrapping is a useful method for estimating errors of a dataset from itself (Andrae 2010). By randomly drawing many subsamples of the data, we obtain a sense of its inherent variance, though there are subtleties to consider such as the independence of values in a dataset. We explore this potential pitfall of bootstrapping in Section 2.2.
- **Jackknife testing**: A resampling technique useful for estimating bias, jackknives can be taken along different dimensions of a dataset to cross-validate results. In particular, null tests can be used to verify whether results are free of systematics, as done with CO power spectra (Keating et al. 2016) and CMB measurements (see e.g. Ade et al. 2008; Chiang et al. 2010; Bischoff et al. 2011; Das et al. 2011a; Araujo et al. 2012; Crites et al. 2015; BICEP2 Collaboration et al. 2016; Ade et al. 2017; Sherwin et al. 2017). An EoR detection must pass both jackknife and null tests, which we highlight in Section 2.3.2.

In the next three subsections, we study each theme in depth, focusing on how power spectrum technique trade-offs affect each. We use toy data models to develop intuition into why certain analysis choices may be appealing and discuss ways in which they are limited. We highlight problems that can arise regarding each theme and offer suggestions to mitigate the issues. Ultimately, we show that rigorous investigations into signal loss, error estimation, and bias must be performed for robust 21 cm results.

### 2.1. Signal Loss

Signal loss can arise in a variety of ways in an analysis pipeline, such as by fitting a polynomial during spectral

calibration, applying a delay domain filter, or deriving weights from data and applying them to itself. Here we focus on signal loss associated with the use of an empirically estimated covariance matrix with the “optimal quadratic estimator” formalism. This loss was significantly under-estimated in the A15 analysis.

### 2.1.1. The Quadratic Estimator Method

The goal of power spectrum analysis is to produce an unbiased estimator of the EoR power spectrum in the presence of both noise and foreground emission. Prior to power spectrum estimation, the data will often have been prepared to have minimal foregrounds by some method of subtraction, so this foreground emission may appear either directly (because it was not subtracted) or as a residual of some subtraction process not in the power spectrum domain. If an accurate estimate of the total covariance of the data is known, including both the desired signal and any contaminants, then the “optimal quadratic estimator” formalism provides a method of producing a minimum variance, unbiased estimator of the desired signal, as shown in Liu & Tegmark (2011), Dillon et al. (2013), Liu et al. (2014a), Liu et al. (2014b), Trott et al. (2012), Dillon et al. (2014), Dillon et al. (2015), Switzer et al. (2015), and Trott et al. (2016). We provide a brief summary of the optimal quadratic estimator and explain its advantages below.

The measured visibilities for a single baseline in Jy are arranged as a data vector,  $\mathbf{x}$ . It has length  $N_t N_f$ , where  $N_t$  is the number of time integrations and  $N_f$  is the number of frequency channels. The covariance of the data is given by

$$\mathbf{C} \equiv \langle \mathbf{x} \mathbf{x}^\dagger \rangle = \mathbf{S} + \mathbf{U} \quad (1)$$

where the average over an ensemble of data realizations produces the true covariance, and we further assume it may be written as the sum of the desired cosmological signal  $\mathbf{S}$  and other terms  $\mathbf{U}$ .

We are interested in estimating the three-dimensional power spectrum of the EoR. Visibilities are measurements of the Fourier transform of the sky along two spatial dimensions (using the flat-sky approximation), and since we are interested in three-dimensional Fourier modes we only need to take one Fourier transform of our visibilities along the line-of-sight dimension. We consider band powers  $P^\alpha$  of the power spectrum of  $\mathbf{x}$  over some range in cosmological  $\mathbf{k}$ , where  $\alpha$  indexes a waveband in  $k_\parallel$  (a cosmological wavenumber  $k_\parallel$  is the Fourier dual to frequency under the delay approximation (Parsons et al. 2012b), which is a good approximation for the short baselines that PAPER analyzes). The fundamental dependence of the covariance on the power spectrum band powers  $P^\alpha$  is encoded as

$$\mathbf{S} = \sum_{\alpha} P^\alpha \frac{\partial \mathbf{C}}{\partial P^\alpha} \equiv \sum_{\alpha} P^\alpha \mathbf{Q}^\alpha \quad (2)$$

where we define  $\frac{\partial \mathbf{C}}{\partial P^\alpha} \equiv \mathbf{Q}^\alpha$ .

The optimal quadratic estimator prescription is then to compute

$$\hat{P}^\alpha = \sum_{\beta} (\mathbf{F}^{-1})^{\alpha\beta} (\hat{q}^\beta - \hat{b}^\beta) \quad (3)$$

where  $\mathbf{F}$  is the Fisher matrix

$$F^{\alpha\beta} \equiv \frac{1}{2} \text{tr} \left( \mathbf{C}^{-1} \mathbf{Q}^\alpha \mathbf{C}^{-1} \mathbf{Q}^\beta \right) \quad (4)$$

and

$$\hat{q}^\alpha = \frac{1}{2} \mathbf{x}^\dagger \mathbf{C}^{-1} \mathbf{Q}^\alpha \mathbf{C}^{-1} \mathbf{x} \quad (5)$$

and

$$\hat{b}^\alpha = \frac{1}{2} \text{tr} \left( \mathbf{U} \mathbf{C}^{-1} \mathbf{Q}^\alpha \mathbf{C}^{-1} \right). \quad (6)$$

This estimator is the minimum variance (smallest error bar) estimate of the power spectrum subject to the constraint that it is also unbiased; that is, the ensemble average of the estimator is equal to its true value

$$\langle \hat{P}^\alpha \rangle = P^\alpha \quad (7)$$

(Tegmark 1997; Bond et al. 1998).

It is clear that in order to obtain an unbiased estimate in the presence of contaminants, the estimator defined by Equation (3) must be capable of “suppressing” or “removing” their effects. By construction, the subtraction of the residual foreground and noise bias accomplishes this, removing any additive bias. However, the  $\mathbf{C}^{-1}$  piece of Equation (5) also has the effect of suppressing residual foregrounds and noise, in both the additive bias and any contributions the residuals may have to the variance. The way in which the optimal estimator accomplishes this is illustrated with a toy model in Appendix A. In this appendix, we show that the effect of the weighting in Equation 5 is to project out the modes of  $\mathbf{U}$  with a different covariance structure than  $\mathbf{S}$  in the power spectrum estimate, and the effect of Equation 6 is to subtract out the remaining bias. Similar effects for a realistic model of the EoR and foregrounds are shown in Liu & Tegmark (2011).

The toy model in Appendix A also illustrates that if the covariance structure of the contaminants is sufficiently different from the desired power spectrum, then the linear bias term may be expected to be quite small, and it is only necessary to know  $\mathbf{C}$  and  $\mathbf{Q}^\alpha$ , but not  $\mathbf{U}$ . Since the foregrounds are expected to be strongly correlated between frequencies whereas the EoR is not, we expect different covariance structures and therefore a small linear bias. Moreover, because the linear bias is always positive and there is no multiplicative bias, the quadratic-only term will always produce an estimate which is *high* relative to the true value, and which can conservatively be interpreted as an upper limit. These considerations, and the difficulty of obtaining an estimate for  $\mathbf{U}$ , motivate the neglect of the linear bias in the rest of this analysis.

Motivated by the desire to retain the advantageous behavior of suppressing contributions of  $\mathbf{U}$  to estimates



of the EoR power spectrum, we note that it is possible to define a modified version of the quadratic estimator where Equation (5) is replaced by

$$\hat{q}^\alpha = \frac{1}{2} \mathbf{x}^\dagger \mathbf{R} \mathbf{Q}^\alpha \mathbf{R} \mathbf{x} \quad (8)$$

where  $\mathbf{R}$  is a weighting matrix chosen by the data analyst. For example, inverse covariance weighting (the optimal form of QE) would set  $\mathbf{R} \equiv \mathbf{C}^{-1}$  and a uniform-weighted case would use  $\mathbf{R} \equiv \mathbf{I}$ , the identity matrix. The matrix  $\mathbf{Q}^\alpha$  encodes the dependence of the covariance on the power spectrum parameters  $\partial \mathbf{C} / \partial P^\alpha$  but in practice also does other things, including implementing a transform of the frequency domain visibilities to  $\mathbf{k}$ -space, taking into account cosmological scalings, and converting the visibilities from Jy to Kelvin.

With an appropriate normalization matrix  $\mathbf{M}$ , the quantity

$$\hat{\mathbf{P}} = \mathbf{M} \hat{\mathbf{q}} \quad (9)$$

is a sensible estimate of the true power spectrum  $\mathbf{P}$ . To ensure that  $\mathbf{M}$  correctly normalizes our power spectrum, one may take the expectation value of Equation (9) to obtain

$$\begin{aligned} \langle \hat{P}^\alpha \rangle &= \frac{1}{2} \sum_{\beta \gamma} M^{\alpha \gamma} \text{tr}(\mathbf{R} \mathbf{Q}^\gamma \mathbf{R} \mathbf{Q}^\beta) P^\beta + \frac{1}{2} \sum_{\gamma} \text{tr}(\mathbf{U} \mathbf{R} \mathbf{Q}^\gamma \mathbf{R}) \\ &\equiv \sum_{\beta} W^{\alpha \beta} P^\beta + \frac{1}{2} \sum_{\gamma} \text{tr}(\mathbf{U} \mathbf{R} \mathbf{Q}^\gamma \mathbf{R}), \end{aligned} \quad (10)$$

where  $W^{\alpha \beta}$  are elements of a window function matrix. Let's consider the first term of this expression (again, we are assuming that the linear bias term is significantly suppressed; and if this is not the case, we are simply assuming that we are setting a conservative upper limit). If  $\mathbf{W}$  ends up being the identity matrix for our choices of  $\mathbf{R}$  and  $\mathbf{M}$ , then we recover Equation (7) for the first term, and we have an estimator that has no multiplicative matrix bias. However, Equation (7) is a rather restrictive condition, and it is possible to violate it and still have a sensible (and correctly normalized) power spectrum estimate. In particular, as long as the rows of  $\mathbf{W}$  sum to unity, our power spectrum will be correctly normalized. Beyond this, the data analyst has a choice for  $\mathbf{M}$ . For simplicity in this section we choose  $\mathbf{M}$  to be diagonal, although we discuss other cases for the analysis of PAPER-64 data as explained in Section 3.3. In a preview of what is to come, we also stress that the derivation that leads to Equation (10) assumes that  $\mathbf{R}$  and  $\mathbf{x}$  are not correlated. If this assumption is violated, a simple application of the (now incorrect) formulae in this section can result in an improperly normalized power spectrum estimator that does not conserve power, i.e., one that has signal loss.

Given the advantages of inverse covariance weighting, a question arises of how one goes about estimating  $\mathbf{C}$ . One method is to empirically derive it from the data  $\mathbf{x}$  itself. Similar types of weightings that are based on variance information in data are done in Chang et al.

(2010) and Switzer et al. (2015). In previous PAPER analyses, one time-averages the data to obtain:

$$\hat{\mathbf{C}} \equiv \langle \mathbf{x} \mathbf{x}^\dagger \rangle_t \approx \langle \mathbf{x} \mathbf{x}^\dagger \rangle, \quad (11)$$

assuming  $\langle \mathbf{x} \rangle_t = 0$  (a reasonable assumption since fringes average to 0 over a sufficient amount of time), where  $\langle \rangle_t$  denotes a finite average over time. The weighting matrix for our empirically estimated inverse covariance weighting is then  $\mathbf{R} \equiv \hat{\mathbf{C}}^{-1}$ , where we use a hat symbol to distinguish the empirical covariance from the true covariance  $\mathbf{C}$ .

In the next three sections, we use toy models to investigate the effects of weighting matrices on signal loss by experimenting with different matrices  $\mathbf{R}$  and examining their impact on the resulting power spectrum estimates  $\hat{\mathbf{P}}$ . Our goal in experimenting with weighting is to suppress foregrounds and investigate EoR losses associated with it. We note that we purposely take a thorough and pedagogical approach to describing the toy model examples given in the next few sections. The specifics of how signal loss appears in PAPER's analysis is later described in Section 3.1.

As a brief preview, we summarize our findings in the following sections here:

- If the covariance matrix is estimated from the data, a strong correlation between the estimated modes and the data will in general produce an estimate of the signal power spectrum which is strongly biased *low* relative to the true value. In this context, this is what we call “signal loss” (Section 2.1.2).
- The effect of the bias is worsened when the number of independent samples used to estimate the covariance matrix is reduced (Section 2.1.3).
- The rate at which empirical eigenvectors converge to their true forms depends on sample variance in the empirical estimate and the steepness (the degeneracy) of the empirical eigenspectrum. In general, larger sample variances and flatter eigenspectra (more degeneracies) lead to more loss (Section 2.1.3).
- Knowing these things, there are some simple ways of altering the empirical covariance matrix to decouple it from the data and produce unbiased power spectrum estimates (Section 2.1.4).

### 2.1.2. Toy Model: Inverse Covariance Weighting

Using a toy model, we will now build intuition into how weighting by the inverse of the empirically estimated covariance,  $\hat{\mathbf{C}}^{-1}$ , can give rise to signal loss. We construct a simple dataset that contains visibility data with 100 time integrations and 20 frequency channels. This model represents realistic dimensions of about an hour of PAPER data which might be used for a power

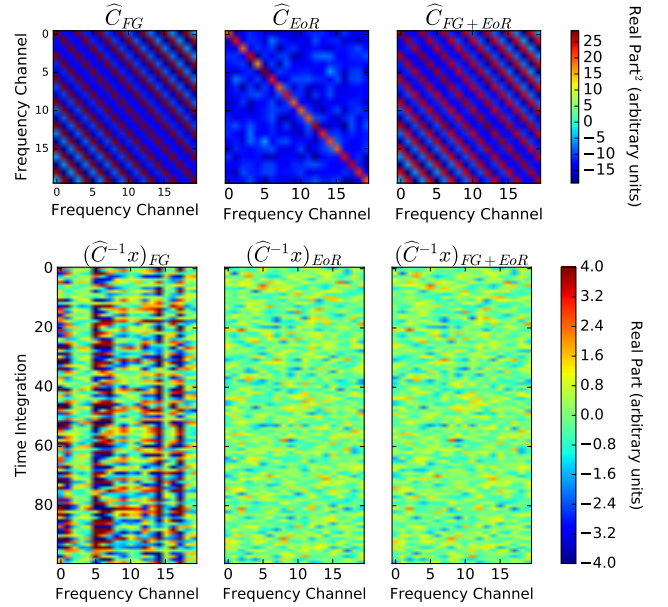


**Figure 1.** Our toy model dataset to which we apply different weighting schemes to in order to investigate signal loss. We model a mock foreground-only visibility with a sinusoid signal that varies smoothly in time and frequency. We model a mock visibility of an EoR signal as a random Gaussian signal. We add the two together to form  $\mathbf{x} = \mathbf{x}_{\text{FG}} + \mathbf{x}_{\text{EoR}}$ . Real parts are shown here.

spectrum analysis. For PAPER-64 (both the A15 analysis and our new analysis) we use  $\sim 8$  hours of data (with channel widths of 0.5 MHz and integration times of 43 seconds), but here we scale it down with no loss of generality.

We create mock visibilities,  $\mathbf{x}$ , and assume a non-tracking, drift-scan observation. Hence, flat spectrum sources (away from zenith) lead to measured visibilities which oscillate in time and frequency. We therefore form a mock visibility measurement of a bright foreground signal,  $\mathbf{x}_{\text{FG}}$ , as a complex sinusoid that varies smoothly in time and frequency, a simplistic but realistic representation of a single bright source. We also create a mock visibility measurement of an EoR signal  $\mathbf{x}_{\text{EoR}}$  as a complex, Gaussian random signal. A more realistic EoR signal would have a sloped power spectrum in  $p(k)$  (instead of flat, as in the case of white noise), which could be simulated by introducing frequency correlations into the mock EoR signal. However, here we treat all  $k$ 's separately, so a simplistic white noise approximation can be used. Our combined data vector is then  $\mathbf{x} = \mathbf{x}_{\text{FG}} + \mathbf{x}_{\text{EoR}}$ , to which we apply different weighting schemes throughout Section 2.1. The three data components are shown in Figure 1.

We compute the power spectrum of our toy model dataset  $\mathbf{x}$  using Equations 8 and 9, with  $\mathbf{R} \equiv \hat{\mathbf{C}}^{-1}$ . Figure 2 shows the estimated covariances of our toy model datasets along with the  $\hat{\mathbf{C}}^{-1}$  weighted data. The foreground sinusoid is clearly visible in  $\hat{\mathbf{C}}_{\text{FG}}$ . The power spectrum result is shown in green in the left plot of Figure 3. Also plotted in the figure are the uniform-weighted ( $\mathbf{R} \equiv \mathbf{I}$ ) power spectrum of the individual components  $\mathbf{x}_{\text{FG}}$  (blue) and  $\mathbf{x}_{\text{EoR}}$  (red). As shown, our  $\hat{\mathbf{C}}^{-1}$  weighted result successfully suppresses foregrounds, demonstrated in Figure 3 by the missing foreground peak in the weighted power spectrum estimate (green). It is also evident that our result fails to recover the EoR signal — it exhibits the correct shape, but the ampli-



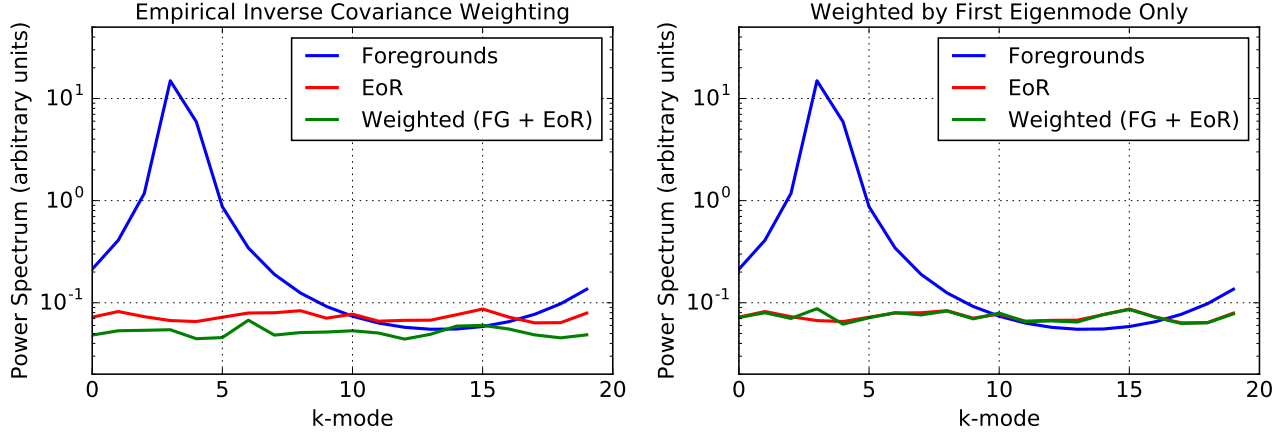
**Figure 2.** The estimated covariance matrices (top row) and inverse covariance-weighted data (bottom row) for FG only (left), EoR only (middle), and FG + EoR (right). Real parts are shown here.

tude level is slightly low. It is this behavior which we describe as signal loss.

As discussed in Section 2.1.1, this behavior is *not* expected in the case that we were to use a true  $\mathbf{C}^{-1}$  weighting. Rather, we would obtain the behavior shown in the toy model in Appendix A<sup>1</sup>, with suppression of the foreground mode resulting in a nearly unbiased estimate of the power spectrum. The key difference is that since  $\hat{\mathbf{C}}$  is estimated from the data, its eigenvectors and eigenvalues are strongly coupled to the particular data realization that was used to compute it, and this coupling leads to loss.

For the case of a mode which can be safely assumed to be predominantly a foreground, its presence in the true covariance matrix will result in the desired suppression via a kind of projection (Equations A3 and A4); whether or not it is strongly correlated with the actual data vector is irrelevant. However, in the case of an empirically estimated covariance matrix, the eigenmodes of  $\hat{\mathbf{C}}_{\text{EoR}}$  will both be incorrect and can be correlated with the data. If these incorrect eigenmodes are not correlated with the data, it will lead to non-minimum variance estimates but will not produce the suppression of the power spectrum amplitude as seen in the left plot of Figure 3. As shown mathematically in Appendix B, however, if  $\hat{\mathbf{C}}_{\text{EoR}}$  is correlated with the data vector  $\mathbf{x}$ , there is a kind of projection of power in the *non*-foreground modes from the resulting power

<sup>1</sup> Note that there the true covariance matrix is also the sum of a diagonal portion describing the signal, and a single mode describing the contaminant (similar to Figure 2).



**Figure 3.** Resulting power spectrum estimates for the toy model simulation described in Section 2.1.2 — foregrounds only (blue), EoR only (red), and the weighted FG + EoR dataset (green). The power spectrum of the foregrounds peaks at a  $k$ -mode based on the frequency of the sinusoid used to create the mock FG signal. In the two panels, we compare using empirically estimated inverse covariance weighting where  $\mathbf{C}$  is derived from the data (left), and projecting out the zeroth eigenmode only (right). In the former case, signal loss arises from the coupling of the eigenmodes of  $\hat{\mathbf{C}}$  to the data. There is negligible signal loss when all eigenmodes besides the foreground one are no longer correlated with the data.

spectrum estimate, thus producing an estimate that is biased low. In short, *if the covariance is computed from the data itself, it carries the risk of over-fitting information in the data and introducing a multiplicative bias (per  $k$ ) to estimates of the signal.*

The danger of an empirically estimated covariance matrix comes mostly from not being able to describe the EoR-dominated eigenmodes of  $\mathbf{C}$  accurately, for which the EoR signal is brighter than foregrounds. In such a case, the coupling between these modes to the data realization leads to the overfitting and subtraction of the EoR signal. More specifically, the coupling between the estimated covariance and the data is anti-correlated in nature (which is explained in more detail in Section 3.1.1 and Appendix B), which leads to loss. Mis-estimating  $\mathbf{C}$  for EoR-dominated eigenmodes is therefore more harmful than for FG-dominated modes, and since the lowest-valued eigenmodes of an eigenspectrum are typically EoR-dominated, using this part of the spectrum for weighting is most dangerous.

Armed with this information, we can tweak the covariance in a simple way to suppress foregrounds and yield minimal signal loss. Recall that our toy model foreground can be perfectly described by a single eigenmode. Using the full dataset’s (foreground plus EoR signal) empirical covariance, we can project out the zeroth eigenmode and then take the remaining covariance to be the identity matrix. This decouples the covariance from the data for the EoR modes. The resulting power spectrum estimate for this case is shown in the right plot of Figure 3. In this case we recover the EoR signal, demonstrating that if we can disentangle the foreground-dominated modes and EoR-dominated modes, we can suppress foregrounds with negligible signal loss.

Altering  $\hat{\mathbf{C}}$  as such is one specific example of a regularization method for this toy model, in which we are changing  $\hat{\mathbf{C}}$  in a way that changes its coupling to the

data realization. There are several other simple ways to regularize  $\hat{\mathbf{C}}$ , and we will discuss some in Section 2.1.4.

### 2.1.3. Toy Model: Fringe-Rate Filtering

We have shown how signal loss can arise due to the coupling of EoR-dominated eigenmodes to the data. We will next show how this effect is exacerbated by reducing the total number of independent samples in a dataset.

A fringe-rate filter is an analysis technique designed to maximize sensitivity by integrating in time (Parsons et al. 2016). Rather than a traditional box-car average in time, a time domain filter can be designed to up-weight temporal modes consistent with the sidereal motion on the sky, while down-weighting modes which are noise-like.

Because fringe-rate filtering is analogous to averaging in time, it comes at the cost of reducing the total number of independent samples in the data. To mimic this filter, we average every four time integrations of our toy model dataset together, yielding 25 independent samples in time (Figure 4). We choose these numbers so that the total number of independent samples is similar to the number of frequency channels — hence our matrices will still be full rank.

With fringe-rate filtering resulting in fewer independent modes, it becomes more difficult for the empirical covariance to estimate the true covariance matrix of the fringe-rate filtered data. We can quantify this effect by evaluating a convergence metric  $\varepsilon(\hat{\mathbf{C}})$  for the empirical covariance, which we define as

$$\varepsilon(\hat{\mathbf{C}}) \equiv \sqrt{\frac{\sum_{ij} (\hat{C}_{ij} - C_{ij})^2}{\sum_{ij} C_{ij}^2}}, \quad (12)$$

where  $\mathbf{C}$  is the true covariance matrix. To compute this metric, we draw different numbers of realizations (different draws of Gaussian noise) of our toy model EoR



measurement,  $\mathbf{x}_{\text{EoR}}$ , and take their ensemble average. We then compare this to the “true” covariance, which in our simulation is set to be the empirical covariance after a large number (500) of realizations. As shown in Figure 5, we perform this computation for a range of total independent ensemble realizations (horizontal axis) and number of independent samples in the data following fringe-rate filtering (different colors). With more independent time samples (i.e., more realizations) in the data, one converges to the true fringe-rate filtered covariance more quickly.

The situation here with using a finite number of time samples to estimate our covariance is analogous to a problem faced in galaxy surveys, where the non-linear covariance of the matter power spectrum is estimated using a large — but finite — number of expensive simulations. There, the limited number of independent simulations results in inaccuracies in estimated covariance matrices (Dodelson & Schneider 2013; Taylor & Joachimi 2014), which in turn result in biases in the final parameter constraints (Hartlap et al. 2007). In our case, the empirically estimated covariances are used for estimating the power spectrum, and as we discussed in the previous section (and will argue more thoroughly in Section 3.1.1 and Appendix B), couplings between these covariances and the data can lead to power spectrum estimates that are biased *low*—which is precisely signal loss. In future work, it will be fruitful to investigate whether advanced techniques from the galaxy survey literature for estimating accurate covariance matrices can be successfully adapted for 21 cm cosmology. These techniques include the imposition of sparsity priors (Padmanabhan et al. 2016), the fitting of theoretically motivated parametric forms (Pearson & Samushia 2016), covariance tapering (Paz & Sánchez 2015), marginalization over the true covariance (Sellentin & Heavens 2016), and shrinkage methods (Pope & Szapudi 2008; Joachimi 2017).

The overall convergence of the covariance is important, but also noteworthy is the fact that different eigenvectors converge to their true forms at different rates. This is illustrated by Figure 6, which shows the convergence of eigenvectors in an empirical estimate of a covariance matrix. For this particular toy model, we construct a covariance whose true form combines the same mock foreground from the previous toy models with an EoR component that is modeled as a diagonal matrix with eigenvalues spanning one order of magnitude (more specifically, we construct the EoR covariance as a diagonal matrix in the Fourier domain, where the signal is expected to be uncorrelated; its Fourier transform is then the true covariance of the EoR in the frequency domain, or  $\mathbf{C}_{\text{EoR}}$ ). For different numbers of realizations, we draw random EoR signals that are consistent with  $\mathbf{C}_{\text{EoR}}$ , add them to the mock foreground data, and compute the combined empirical covariance by averaging over the realizations. The eigenvectors of this empirical covariance are then compared to the true eigenvectors  $\hat{\mathbf{v}}$ , where we

use as a convergence metric  $\varepsilon(\hat{\mathbf{v}})$ , defined as:

$$\varepsilon(\hat{\mathbf{v}}) \equiv \sqrt{\sum_i^{N_f} |\mathbf{v} - \hat{\mathbf{v}}|_i^2}, \quad (13)$$

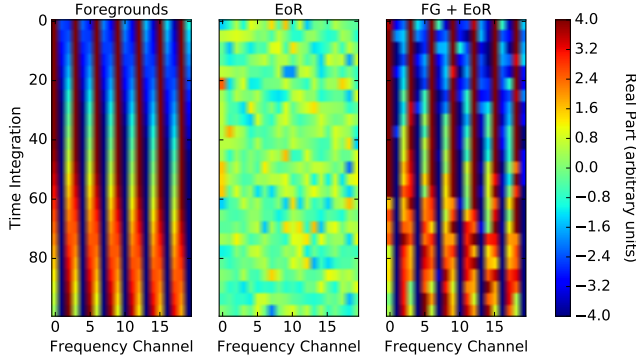
where  $N_f$  is the number of frequencies (20) in the mock data. The eigenmode convergence curves in Figure 6 are ranked ordered by eigenvalue, such that “Eigenmode #0” illustrates the convergence of the eigenvector with the largest eigenvalue, “Eigenmode #1” for the second largest eigenvalue, and so on. We see that the zeroth eigenmode — the mode describing the foreground signal — is quickest to converge.

Our numerical test reveals that the convergence rates of empirical eigenvectors is related to the sample variance in our empirical estimate. In general, computing an empirical covariance from a finite ensemble average means that the empirical eigenmodes have sample variances. Consider first a limiting case where all eigenvalues are equal. In such a scenario, any linear combination of eigenvectors is also an eigenvector, and thus there is no sensible way to define the convergence of eigenvectors. In our current test, aside from the zeroth mode, the eigenvalues have similar values but are not precisely equal. Hence, there is a well-defined set of eigenvectors to converge to. However, due to the sample variance of our empirical covariance estimate, there may be accidental degeneracies between modes, where some modes are mixing and swapping with others. Therefore, the steeper an eigenspectrum, the easier it is for the eigenmodes to decouple from each other and approach their true forms. A particularly drastic example of this can be seen in the behavior of mode 0 (the foreground mode), whose eigenvalue differs enough from the others that it is able to converge reasonably quickly despite substantial sample variance in our empirical covariance estimate. To break degeneracies in the remaining modes, however, requires many more realizations.

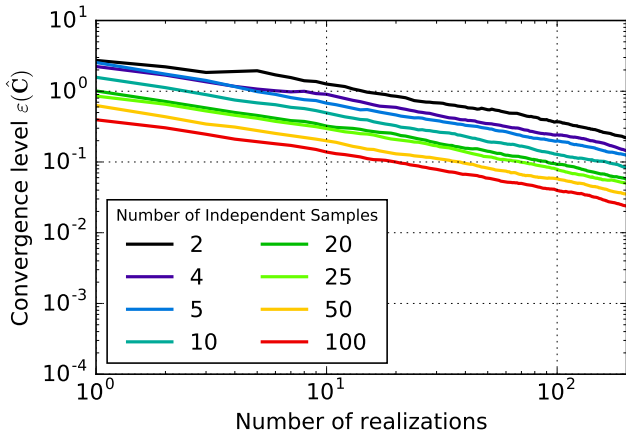
While the connection between the rate of convergence of an empirical eigenvector with the steepness of an eigenspectrum is interesting, it is also important to note that regardless of convergence rate, any mode that is coupled to the data is susceptible to signal loss. The true eigenvectors are not correlated with the data realizations; thus, if our empirical eigenvectors are converged, there will not be any signal loss. However, an unconverged eigenvector estimate will retain some memory of the data realizations used in its generation, leading to signal loss.

In the toy models throughout Section 2.1, we exploit the fact that the strongest eigenmode (highest eigenvalue mode) is dominated by foregrounds in order to purposely incur signal loss for that mode. Even for the case of real PAPER data (Section 3.1), we make the assumption that the strongest eigenmodes are likely the most contaminated by foregrounds. However, in general, foregrounds need not be restricted to the strongest eigenmodes, and as we have seen, it is really the degen-





**Figure 4.** Our “fringe-rate filtered” (time-averaged) toy model dataset. We average every four samples together, yielding 25 independent samples in time. Real parts are shown here.

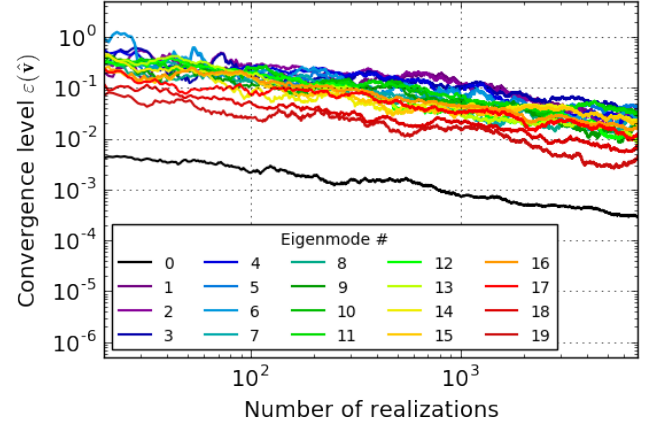


**Figure 5.** The convergence level, as defined by Equation (12), of empirically estimated covariances of mock EoR signals with different numbers of independent samples. In red, the mock EoR signal is comprised entirely of independent samples. Subsequent colors show time-averaged signals. As the number of realizations increases, we see that the empirical covariances approach the true covariances. With more independent samples, the quicker an empirical covariance converges (i.e., the quicker it decouples from the data), and the less signal loss we would expect to result.

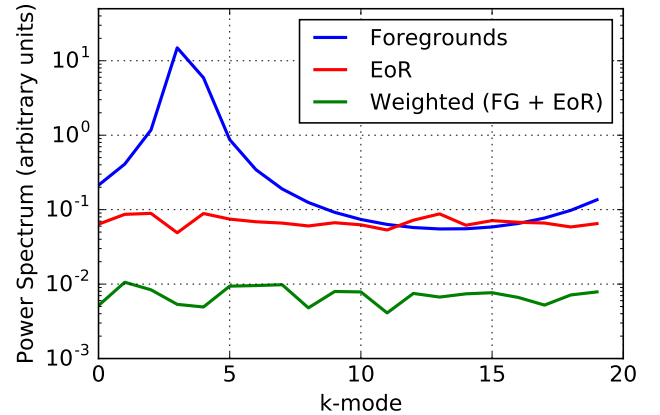
eracies between modes that determines how quickly they converge, and hence how much signal loss can result.

With Figures 5 and 6 establishing the connection between convergence rates (of empirical covariances and eigenvectors) and number of realizations, we now turn back to our fringe-rate filtered toy model. The power spectrum results for this model are shown in Figure 7, and as expected there is a much larger amount of signal loss for this time-averaged dataset since we do a worse job estimating the true covariance. In addition, as a result of having fewer independent samples, we obtain an estimate with more scatter. This is evident by noticing that the green curve in Figure 7 fails to trace the shape of the uniform-weighted EoR power spectrum (red).

Using our toy model, we have seen that a sensitivity-driven analysis technique like fringe-rate filtering has trade-offs of signal loss and noisier estimates when using data-estimated covariance matrices. Longer integrations



**Figure 6.** The convergence level, as defined by Equation (13), of empirically estimated eigenvectors for different numbers of mock data realizations. The colors span from the 0th eigenmode (has the highest eigenvalue) to the 19th eigenmode (has the lowest eigenvalue), where they are ordered by eigenvalue in descending order. This figure shows that the zeroth eigenmode converges the quickest, implying that eigenvectors with eigenvalues that are substantially different than the rest (the FG-dominated mode has a much higher eigenvalue than the EoR modes) are able to converge to the true eigenvectors the quickest. On the other hand, eigenmodes 1-19 have similar eigenvalues and are slower to converge because of degeneracies between them.



**Figure 7.** Resulting power spectrum estimate for the “fringe-rate filtered” (time-averaged) toy model simulation — foregrounds only (blue), EoR only (red), and the weighted FG + EoR dataset (green). We use empirically estimated inverse covariance weighting where  $\mathbf{C}$  is computed from the data. There is a larger amount of signal loss than for the non-averaged data, a consequence of weighting by eigenmodes that are more strongly coupled to the data due to there being fewer independent modes in the data.

increase sensitivity but reduce the number of independent samples, resulting in eigenmodes correlated with the data that can overfit signal greatly. We note that a fringe-rate filter does have a range of benefits, many described in Parsons et al. (2016), so it can still be advantageous to use one despite the trade-offs.

#### 2.1.4. Toy Model: Other Weighting Options

In Section 2.1.2 we showed one example of how altering  $\hat{\mathbf{C}}$  can make the difference between nearly zero and some signal loss. We will now use our toy model to describe several other ways to tailor  $\hat{\mathbf{C}}$  in order to minimize signal loss. We choose four independent regularization methods to highlight in this section, which have been chosen due to their simplicity in implementation and straightforward interpretations. We illustrate the resulting power spectra for the different cases in Figure 8. These examples are not meant to be taken as suggested analysis methods but rather as illustrative cases.

As a first test, we model the covariance matrix of EoR as a proof of concept that if perfect models are known, signal loss can be avoided. We know that our simulated EoR signal should have a covariance matrix that mimics the identity matrix, with its variance encoded along the diagonal. We model  $\mathbf{C}_{\text{EoR}}$  as such (i.e., the identity), instead of computing it based on  $\mathbf{x}_{\text{EoR}}$  itself. Next, we add  $\mathbf{C}_{\text{EoR}} + \hat{\mathbf{C}}_{\text{FG}}$  (where  $\hat{\mathbf{C}}_{\text{FG}} = \langle \mathbf{x}_{\text{FG}} \mathbf{x}_{\text{FG}}^\dagger \rangle_t$ ) to obtain a final  $\hat{\mathbf{C}}$  to use in weighting. In Figure 8 (upper left), we see that there is negligible signal loss. This is because by modeling  $\mathbf{C}_{\text{EoR}}$ , we avoid over-fitting EoR fluctuations in the data that our model doesn't know about (but an empirically derived  $\hat{\mathbf{C}}_{\text{EoR}}$  would). In practice such a weighting option is not feasible, as it is difficult to model  $\mathbf{C}_{\text{EoR}}$ , and  $\hat{\mathbf{C}}_{\text{FG}}$  is unknown because we do not know how to separate out the foregrounds from the EoR in our data.

The second panel (top right) in Figure 8 uses a regularization method of setting  $\hat{\mathbf{C}} \equiv \hat{\mathbf{C}} + \gamma \mathbf{I}$ , where  $\gamma = 5$  (an arbitrary strength of  $\mathbf{I}$  for the purpose of this toy model). By adding the identity matrix, element-wise, we are weighting the diagonal elements of the estimated covariance matrix more heavily than those off-diagonal. Since the identity component does not know anything about the data realization, it alters the covariance to be less coupled to the data. Although there is negligible signal loss using this regularization, the small green peak at the third  $k$ -mode represents residual foregrounds that still exist since the shapes encoded in the off-diagonal frequency correlations of the covariance matrix were deemed not as prominent as the diagonal elements using this weighting scheme.

The third panel (bottom left) in Figure 8 minimizes signal loss a different way — by only using the first three eigenmodes of the estimated covariance. Recalling that our toy model foregrounds can be described entirely by the zeroth eigenmode, this method intentionally projects out the highest-valued modes only by replacing all but the three highest weights in the eigenspectrum with 1's (equal weights). Again, avoiding the over-fitting of EoR-dominated modes which are coupled to the data results in negligible signal loss. However, we do not perfectly recover the shape of the EoR power spectrum because we lost information when ignoring the relative weights of most of the modes. While this case is illuminating for the toy model, in practice it is not obvious which

eigenmodes are FG or EoR dominated (and they could be mixed as well), so determining which subset of modes to down-weight is not trivial. We experiment with this idea using PAPER data in Section 3.1.3.

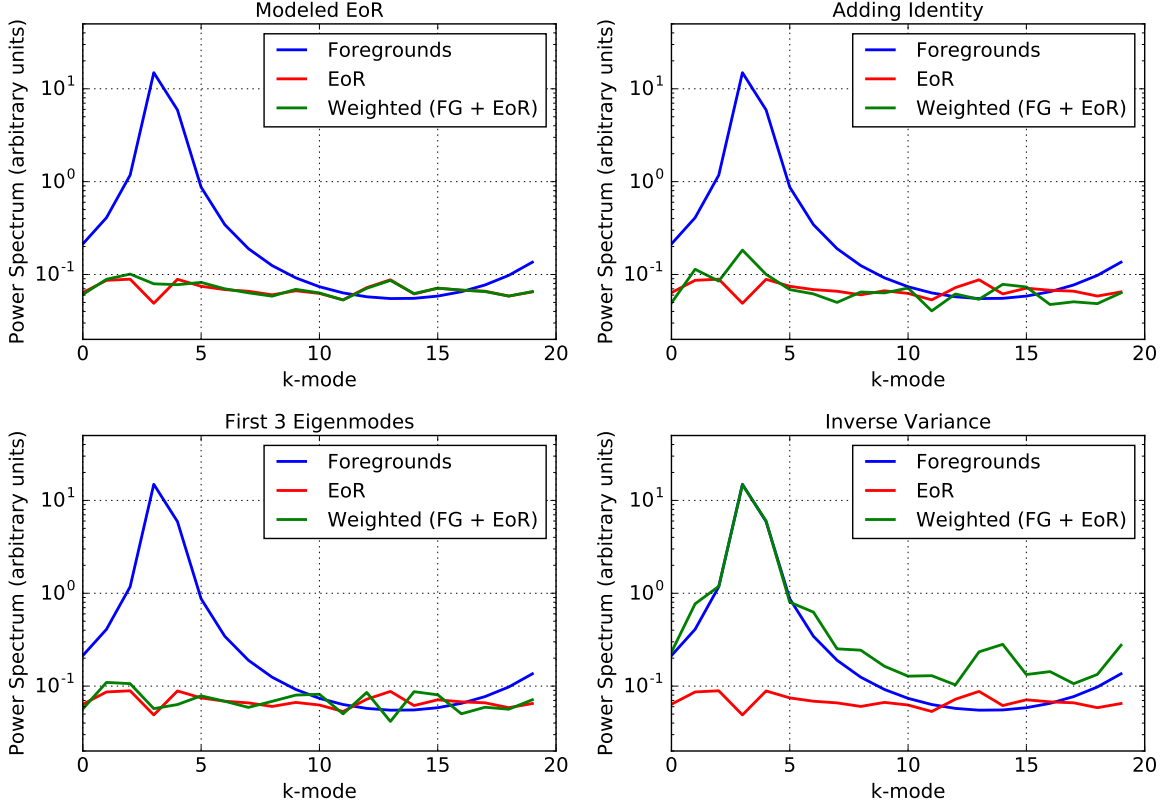
The last regularization scheme we are highlighting here is setting  $\hat{\mathbf{C}} \equiv \hat{\mathbf{C}} \circ \mathbf{I}$  (element-wise multiplication), or inverse variance weighting (i.e., keeping only the diagonal elements of  $\hat{\mathbf{C}}$ ). In the bottom right panel of Figure 8, we see that this method does not down-weight the foregrounds at all — this regularization altered  $\hat{\mathbf{C}}$  in a way where it is no longer coupled to *any* of the empirically estimated eigenmodes, including the FG-dominated one. For this toy model, our foregrounds are spread out in frequency and therefore have non-negligible frequency-frequency correlations. Multiplying by the identity, element-wise, results in a diagonal matrix, meaning we do not have any correlation information. Because of this, we do a poor job suppressing the foreground. But because we decoupled the whole eigenspectrum from the data, we also avoid signal loss. Although this method did not successfully recover the EoR signal for this particular simulation, it is important that we show that there are many options for estimating a covariance matrix, and some may down-weight certain eigenmodes more effectively than others based on the spectral nature of the components in a dataset.

In summary, we have a choice of how to weight 21 cm data. Ideally, we want to down-weight bright foregrounds without removing the underlying cosmological signal. However, there are trade-offs between the weighting method used, its foreground-removal effectiveness, the number of independent signal samples in a dataset, and the amount of resulting signal loss.

## 2.2. Error Estimation

Our second major 21 cm power spectrum theme is error estimation, as we desire robust methods for determining accurate confidence intervals for our measurements. Two popular ways of calculating errors on a power spectrum measurement are calculating the variance of power spectrum results, and computing a theoretical error estimate based on an instrument's system temperature and observational parameters. In a perfect world, both methods would match up. However, in practice the two do not always agree due to a number of factors, including possible non-Gaussianities in the noise properties of our instruments and possible systematics in the data.

A third option which acts as a middle ground between purely theoretical and purely empirical errors is using Gaussian error. This involves the assumption of Gaussianity, but allows the variance of the power spectrum estimator to be written as a function of the two-point estimator, or covariance. One could empirically calculate the covariance and then propagate it into an analytic expression to compute the errors, making this method fall somewhere between being fully empirical and fully modeled (see Das et al. (2011b) for an example of its



**Figure 8.** Resulting power spectra estimates for our “fringe-rate filtered” (time-averaged) toy model simulation — foregrounds only (blue), EoR only (red), and the weighted FG + EoR dataset (green). We show four alternate weighting options that each minimize signal loss, including modeling the covariance matrix of EoR (upper left), regularizing  $\hat{\mathbf{C}}$  by adding an identity matrix to it (upper right), using only the first three eigenmodes of  $\hat{\mathbf{C}}$  (lower left), and keeping only the diagonal elements of  $\hat{\mathbf{C}}$  (lower right). The first case (upper left) is not feasible in practice since we do not know  $\mathbf{C}_{\text{FG}}$  and  $\mathbf{C}_{\text{EoR}}$  like we do in the toy model.

implementation).

For PAPER’s analysis, we choose a data-driven method of error estimation that does not rely on assumptions of Gaussianity. Namely, we compute error bars that have been derived from the inherent variance of our measurements. A common technique used to do this is bootstrapping. For pedagogical purposes, we first define the technique of bootstrapping and then illustrate one of its pitfalls through a toy model.

Bootstrapping uses sampling with replacement to estimate a posterior distribution. For example, measurements (like power spectra) can be made from different samples of data. Each of these measurements is a different realization drawn from some underlying distribution, and realizations are correlated with each other to a degree set by the fraction of sampled points that are held in common between them. Through the process of re-sampling and averaging along different axes of a dataset, such as along baselines or times, we can estimate error bars for our results which represent the underlying distribution of values that are allowed by our measurements (Efron & Tibshirani 1994; Andrae 2010).

Suppose we have  $N$  different measurements targeting the same quantity (for example,  $N$  power spectrum measurements). Bootstrapping means that we form  $N_{\text{boot}}$

(often a large number) bootstraps, where each bootstrap is a random selection of the  $N$  measurements. Bootstraps each have dimensions of  $N$ , and the values populated into each bootstrap are drawn from the original set of measurements with replacement (i.e., every  $n^{\text{th}}$  slot in  $N$  is filled randomly for each bootstrap). Next we take the mean of each bootstrap to collapse it from an array of length  $N$  to a single number (we are interested in the mean statistic here, but any function of interest can be applied to each bootstrap as long as it’s the same function for each one). The error (on the mean) is then computed as the standard deviation across all bootstraps.

We must be careful in distinguishing  $N_{\text{boot}}$ , the number of bootstraps, from  $N$ , the number of samples, or elements, or values, that comprise a bootstrap. In the toy models presented in this section,  $N_{\text{boot}}$  is typically large, and the standard deviation across bootstraps (the error we are computing) converges for large  $N_{\text{boot}}$ . Typically  $N$  is a straightforward value to set that just depends on the experiment. However, we will illustrate one case in which it is not simply the number of samples along the axis that is being re-sampled. More specifically, we will see that  $N$  depends on sample independence and may not always be straightforward to approximate.

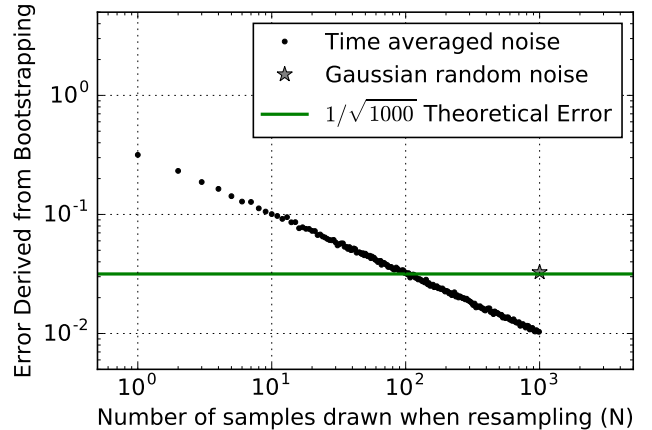
For our toy model, suppose we have a Gaussian random signal dataset of length  $N = 1000$  and unity variance (zero mean). This could represent 1000 power spectrum measurements, for which we are interested in its error. We predict that the error on the mean should obey  $1/\sqrt{N}$ , where  $N$  is the number of samples.

We next form 500 bootstraps ( $N_{\text{boot}} = 500$ ). To create each bootstrap, we draw  $N$  samples, with replacement, of the original data, and take the mean over the  $N$  samples. The standard deviation over the 500 bootstraps gives an error estimate for our dataset. This error is indicated by the gray star in Figure 9 and matches our theoretical prediction (green).

One major caveat of bootstrapping arises when working with correlated data. If, for example, a dataset has many repeated values inside it, this would be reflected in each bootstrap. The same value would be present multiple times within a bootstrap and also be present between bootstraps, purely because it has a more likely chance of being drawn if there are repeats of itself. Therefore, bootstrapping correlated data results in a smaller variation between bootstraps, and hence, under-estimates errors. The use of a fringe-rate filter, which averages data in time to increase sensitivity, is one example which leads to a reduction in the number of independent samples, creating a situation in which errors can be under-estimated. We will now show this effect using our toy model.

Going back to our toy model, we apply a sliding boxcar average to 10 samples at a time, thus reducing the number of independent data samples to  $N/10 = 100$ . Bootstrapping this time-averaged noise, using the same method as described earlier (drawing  $N = 1000$  elements per bootstrap sample), under-estimates the error by a factor of  $\sim 3$  (black points in Figure 9, at  $N = 1000$ ). This occurs because we are drawing more samples than independent ones available, and thus some samples are repeated multiple times in all bootstraps, leading to less variation between the bootstraps. In fact, the error derived from bootstrapping is a strong function of the number of elements that are drawn per bootstrap (Figure 9, black points), and we can both under-estimate the error by drawing too many or over-estimate it by drawing too few. However, since we know that we have 100 independent samples in this toy model, the error associated with drawing  $N = 100$  samples with replacement does match the theoretical prediction as expected (the black points cross the green line at  $N = 100$  in Figure 9).

This example highlights the importance of understanding how analysis techniques (such as fringe-rate filtering) can affect a common statistical procedure like bootstrapping. Bootstrapping as a means of estimating power spectrum errors from real fringe-rate filtered data requires knowledge of the number of independent samples, which is not always a trivial task. For example, computing the effective number of independent samples of fringe-rate filtered data is not as simple as counting



**Figure 9.** Error estimation from bootstrapping as a function of the number of elements drawn per bootstrap when sampling with replacement. The star represents the standard deviation of  $N_{\text{boot}} = 500$  bootstraps, each created by drawing 1000 elements (with replacement) from a length 1000 array of a Gaussian random signal. The black points correspond to time-averaged data (correlated data) which has 100 independent samples. They illustrate how errors can be under-estimated if drawing more elements than there are independent samples in the data. The estimated errors match up with the theoretical prediction only at  $N = 100$ .

the number of averages performed. Down-sampling a time-averaged signal is straightforward using a boxcar average, but non-trivial with a more complicated convolution function that has long tails. Hence, we do not recommend bootstrapping unless the number of independent samples along the axis that is being re-sampled is well-determined. In Section 3.2.2, we explain how we under-estimated errors in the A15 analysis of PAPER and how our bootstrapping procedure has now changed to avoid the over-sampling of correlated data.

In summary, bootstrapping can be an effective and straightforward way to estimate errors of a dataset. However, we have illustrated a situation in which bootstrapping can lead to under-estimated errors and therefore under-estimated power spectrum limits. We have shown that bootstrapped error depends strongly on the number of elements drawn in a bootstrap sample. Estimated errors can drop to arbitrarily small values when the number of elements drawn exceeds the effective number of independent elements. While bootstrapping is convenient because it provides a way to estimate errors from the data itself, one must assess whether certain analysis choices have compromised the method and whether a variation or an avoidance of traditional re-sampling could be preferred instead.

### 2.3. Bias

In a 21 cm power spectrum, detections could be the EoR signal, but they could also be attributed to other sources of bias. Connecting a detection to EoR as opposed to noise or foreground bias is a key challenge of future 21 cm data analyses (e.g. Petrovic & Oh 2011). In this section we will discuss possible sources of bias in a measurement, as well as techniques that can help



mitigate their effects. We will also present a series of tests in a pedagogical fashion which we suggest be used to help evaluate deep limits and/or detections.

### 2.3.1. Foreground and Noise Bias

In Section 2.1, we discussed signal loss as a form of multiplicative bias to estimates of the signal. Foregrounds are another type of bias, but an additive instead of a multiplicative one. Foreground bias is perhaps one of the main factors limiting 21 cm results, as foreground signals lie  $\sim 4$ -5 orders of magnitude above the cosmological signal. Though there are many techniques proposed for removing foregrounds (see e.g. Vedantham et al. 2012; Chapman et al. 2012; Parsons et al. 2012a; Parsons et al. 2012b; Dillon et al. 2013; Wang et al. 2013; Parsons et al. 2014; Liu et al. 2014a; Wolz et al. 2014; Liu et al. 2014b; Dillon et al. 2015; Pofer et al. 2016b; Trott et al. 2016), most experiments currently remain limited by residuals rather than noise, especially at low  $k$ .

One common method to isolate and filter foregrounds is to exploit their behavior in  $k$ -space. For a particular baseline length, there is a maximum delay imposed on sources attached to the sky, which corresponds to the light-crossing time between two antennas in a baseline. For longer baselines, this value increases, producing what is known as “the wedge” (Datta et al. 2010; Parsons et al. 2012b; Vedantham et al. 2012; Pofer et al. 2013; Thyagarajan et al. 2013; Liu et al. 2014a,b; Patil et al. 2017). The wedge describes a region in  $k$ -space contaminated by smooth spectrum foregrounds, bounded by baseline length (which is proportional to  $k_{\perp}$ ) and delay (which is proportional to  $k_{\parallel}$ ). Properties of the wedge can be used to isolate and avoid foregrounds, as done by A15, Parsons et al. (2014), Dillon et al. (2014), Dillon et al. (2015), Jacobs et al. (2015), Beardsley et al. (2016), and Trott et al. (2016).

Although smooth-spectrum foregrounds preferentially show up at low delay, or low  $k$ -modes, their isolation within the wedge is not perfect. In deep measurements, power spectrum measurements at  $k_{\parallel}$  values beyond the delay associated with the length of a baseline are often still contaminated at a low level. This leakage, particularly at low  $k$ ’s, can be attributed to convolution kernels associated with Fourier-transforming visibilities into delay-space. In other words, smooth-spectrum foregrounds appear as  $\delta$ -functions in delay-space, convolved by the Fourier transform of the source spectrum, the signal chain, and the antenna response, all of which could smear out the foregrounds and cause leakage outside the wedge (e.g. Ewall-Wice et al. 2017; Kerrigan et al. 2018).

There are analysis techniques to mitigate the effects of foreground leakage and prevent information from low  $k$ ’s from spreading to high  $k$  values. For example, narrow window functions in delay-space can be used to minimize the leakage from a particular  $k$  value into other ones (Liu et al. 2014b). In other words, one can construct an estimator using QE that forces a window function to have a minimum response to low  $k$  values. The win-

dow function used in A15 is constructed in such a way, specifically to prevent foregrounds that live at low  $k$ ’s from contaminating higher  $k$ -modes (see Section 3.3).

Determining the source of positive non-EoR detections at higher  $k$ ’s is more difficult. In previous power spectrum results, these detections have been explained as instrumental systematics, particularly time-variable cross talk, RFI, cable reflections, and calibration errors (A15; Parsons et al. 2014; Dillon et al. 2014; Beardsley et al. 2016; Patil et al. 2017). In the next section, we will present some tests that can help distinguish these excesses from that of EoR.

In addition to foreground bias, noise can also be responsible for positive power spectrum detections if thermal noise is multiplied by itself. Every 21 cm visibility measurement contains thermal noise that is comprised of receiver and sky noise. We expect this noise to be independent between antennas and thus we can beat it down (increase sensitivity) by integrating longer, using more baselines, etc. However, the squaring of noise can occur when cross-multiplying visibilities, which is shown by the two copies of  $\mathbf{x}$  in Equation (8). If both copies of  $\mathbf{x}$  come from the same baseline and time, it can result in power spectrum measurements that are higher than those predicted by the thermal noise of the instrument. One way to avoid this type of noise bias is to avoid cross-multiplying data from the same baselines or days. This ensures that the two quantities that go into a measurement have separate noises that don’t correlate with each other. We also note that if the noise level is known, this type of bias can be subtracted off, though this procedure is argued to be dangerous (Dillon et al. 2014; Parsons et al. 2014).

Another type of noise bias can stem from the spurious cross-coupling of signals between antennas. This excess is known as instrumental crosstalk and is an inadvertent correlation between two independent measurements via a coupled signal path. Crosstalk appears as a constant phase bias in time in visibilities, and it varies slowly compared to the typical fringe-rates of sources. Because it is slow-varying, crosstalk can be suppressed using time-averages or fringe-rate filters. However, there remains a possibility that power spectrum detections that aren’t the cosmological signal are caused by residual, low-level crosstalk which survived any suppression techniques.

### 2.3.2. Jackknife Tests

We now approach the difficult task of tracing excesses to foreground, noise, and EoR biases through a discussion of useful jackknife tests. Again, we first approach this topic pedagogically as an introduction to the related PAPER discussion in Section 3.3.

The jackknife is a resampling technique in which a statistic (i.e., power spectrum) is computed in subsets of the data (Quenouille 1949; Tukey 1958). These subsets are then compared to reveal systematics. In this section we define two main tests — the null test and the traditional jackknife — and explain how a power spectrum detection must pass each. We then highlight

how these tests can be used to help distinguish between different sources of bias.

- **Null Test:** A null test is a type of jackknife test that removes the astronomical signal from data in order to investigate underlying systematics (see Keating et al. (2016) for examples from intensity mapping that are closely related to our current application). For example, one can divide data into two subsets by separating odd and even Julian dates, or the first half of the observing season from the second. Subtracting the two removes signal that is common to both subsets, including foregrounds and the EoR signal. The resulting power spectrum should be consistent with thermal noise estimates; if it is not, it suggests the presence of a systematic that differs from one of the data subsets to the other (i.e., doesn't get subtracted perfectly).
- **Traditional Jackknife:** In a broader sense, it is important to perform many jackknife tests in order to instill confidence in a final result. A stable result must be steadfast throughout all jackknives no matter how the data is sliced. Jackknives can be taken along several different axes — for example, one could start with a full dataset, and compute a new power spectrum every time as a day of data is removed, or a baseline is removed. This type of jackknife would reveal bias present only at certain LSTs (such as a foreground source), for example, or misbehaving baselines.

While the null test hunts for deviations from thermal noise and the jackknife tests for deviations in subsamples, they are both closely related. We can highlight the connection between the two using a toy model dataset.

Suppose we have four independent measurements made along two different axes. As an example, we construct  $\mathbf{x}_{1a}$ ,  $\mathbf{x}_{1b}$ ,  $\mathbf{x}_{2a}$ , and  $\mathbf{x}_{2b}$ , where the numbers symbolize two different days of data and the letters represent two different baselines. Each of the measurements have dimensions of 100 time integrations and 20 frequency channels. They each have separate thermal noises constructed as a Gaussian random signal for each, and identical EoR signals.

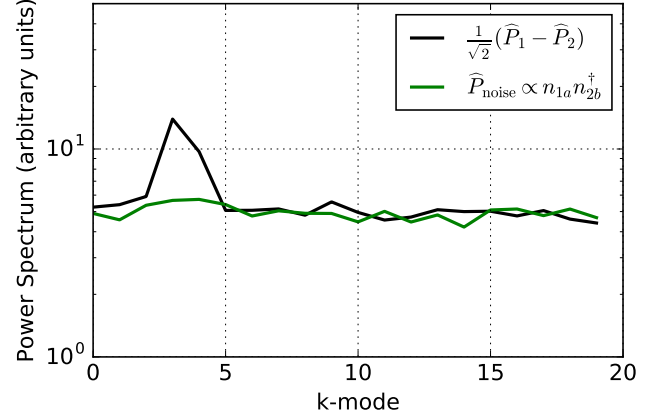
To mimic the presence of a systematic, we add a toy sinusoid foreground, similar to the one used in Section 2.1.2, to only  $\mathbf{x}_{2a}$  and  $\mathbf{x}_{2b}$ . This represents a foreground signal present in only the second day of data, but not the first. Mathematically, if  $\mathbf{n}$  is noise,  $\mathbf{e}$  is the EoR signal, and  $\mathbf{f}$  is the foreground signal, the four measurements can be written as:

$$\mathbf{x}_{1a} = \mathbf{n}_{1a} + \mathbf{e} \quad (14)$$

$$\mathbf{x}_{1b} = \mathbf{n}_{1b} + \mathbf{e} \quad (15)$$

$$\mathbf{x}_{2a} = \mathbf{n}_{2a} + \mathbf{e} + \mathbf{f} \quad (16)$$

$$\mathbf{x}_{2b} = \mathbf{n}_{2b} + \mathbf{e} + \mathbf{f}. \quad (17)$$



**Figure 10.** A null jackknife test shown as the power spectrum difference between two measurements (black), compared to the power spectrum of noise alone (green). Because the null test is not consistent with noise, it suggests the presence of a systematic in either  $\mathbf{x}_1$  or  $\mathbf{x}_2$ . Null tests of clean measurements should be consistent with thermal noise.

We now take a jackknife along the day-axis, forming separate power spectrum estimates for each day:

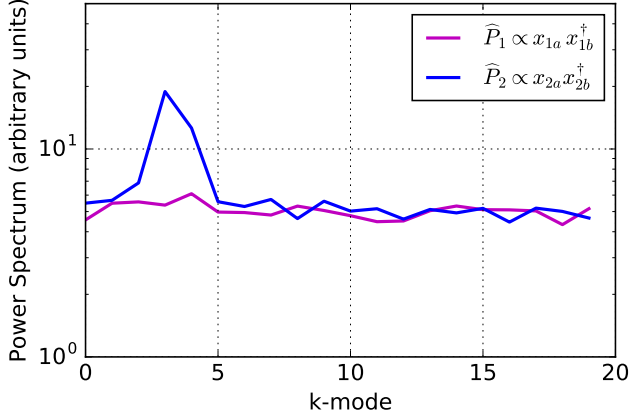
$$\hat{\mathbf{P}}_1 \propto \mathbf{x}_{1a} \mathbf{x}_{1b}^\dagger \quad (18)$$

$$\hat{\mathbf{P}}_2 \propto \mathbf{x}_{2a} \mathbf{x}_{2b}^\dagger. \quad (19)$$

We do not perform a time-average or apply a fringe-rate filter to this toy model, since we are interested only in what jackknife tests can tell us about biases. For the same reason, we use a weighting matrix of  $\mathbf{I}$  for power spectrum estimation to avoid signal loss.

To construct a null test, we difference the two power spectra, with the result shown in Figure 10 (black) along with the power spectrum of noise only (green). Subtracting the two estimates removes sky signal that should ideally be present in both jackknives. However, we see a clear difference between the null test and the power spectrum of noise. This signifies a non-EoR bias that is only present in either  $\mathbf{x}_1$  or  $\mathbf{x}_2$ , but not both.

While the null test is useful for testing noise properties and the uniformity of a dataset, jackknives are useful in pinpointing which data subsets are contaminated by biases and which are not; in our toy model we see that the bias exists only in  $\mathbf{x}_2$  (Figure 11). If foreground or noise biases exist in a dataset, jackknives can tease them out and provide insight into possible sources. For example, if jackknives along the time-axis reveal a bias present at a certain LST, a likely explanation would be excess foreground emission from a radio source in the sky at that time. A jackknife test involving data before and after the application of a fringe-rate filter can reveal whether crosstalk noise bias is successfully suppressed with the filter, or if similar-shaped detections in both power spectra suggest otherwise. There are many other jackknife axes of which we will not go into detail here, including baseline, frequency, and polarization. Ultimately, an EoR detection should persist through them all and



**Figure 11.** Power spectrum estimates for  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , two jackknives of the toy model. They suggest the presence of a systematic in  $\mathbf{x}_2$  only, illustrating how jackknives can be used to tease out excesses. Clean measurements should remain consistent despite the jackknife taken.

a clean measurement should exhibit noise-like null spectra.

In this section we have highlighted how null tests and jackknife tests are key for determining the nature of a power spectrum detection. In Section 3.3 we perform some examples of these tests on PAPER-64 data in order to show that our excesses are not EoR and to identify their likely cause.

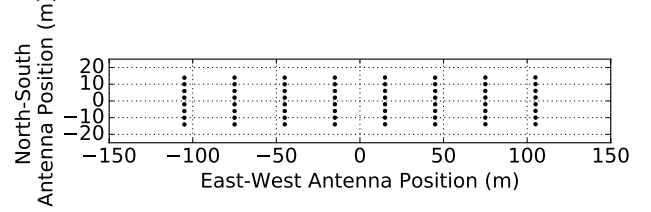
### 3. DEMONSTRATION IN PAPER-64 DATA

In the previous sections we have discussed three overarching 21 cm power spectrum themes — signal loss, error estimation, and bias. Understanding the subtleties and trade-offs involved in each is necessary for an accurate and robust understanding of a power spectrum result.

We now apply these lessons to data from the PAPER experiment, using a subset of the PAPER-64 dataset from A15 in order to illustrate our revised analysis pipeline.

As a brief review, PAPER is a dedicated 21 cm experiment located in the Karoo Desert in South Africa. The PAPER-64 configuration consists of 64 dual-polarization drift-scan elements that are arranged in a grid layout. For our case study, we focus solely on Stokes I estimated data (Moore et al. 2013) from PAPER’s 30 m East/West baselines (Figure 12). All data is compressed, calibrated (using self-calibration and redundant calibration), delay-filtered (to remove foregrounds inside the wedge), LST-binned, and fringe-rate filtered. For detailed information about the backend system of PAPER-64, its observations, and data reduction pipeline, we refer the reader to Parsons et al. (2010) and A15. We note that all data processing steps are identical to those in A15 until after the LST-binning step in Figure 3 of A15.

The previously best published 21 cm upper limit result



**Figure 12.** The PAPER-64 antenna layout. We use only 10 of the 30 m East/West baselines for the analysis in this paper (i.e., a subset of the shortest horizontal spacings).

from A15 placed a  $2\sigma$  upper limit on  $\Delta^2(k)$ , defined as

$$\Delta^2(k) = \frac{k^3}{2\pi^2} \hat{P}(k), \quad (20)$$

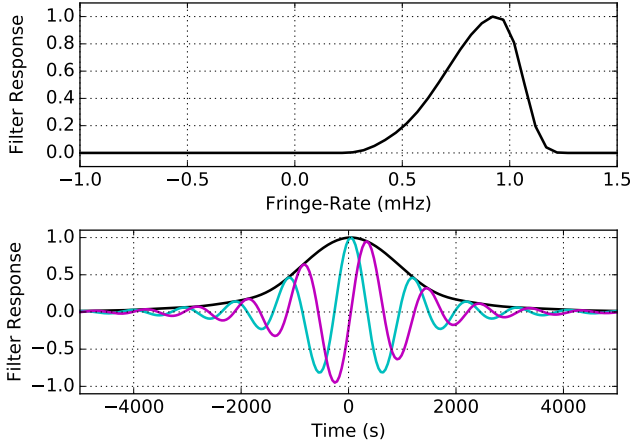
of  $(22.4 \text{ mK})^2$  in the range  $0.15 < k < 0.5 h \text{ Mpc}^{-1}$  at  $z = 8.4$ . The need to revise this limit stems mostly from previously under-estimated signal loss and under-estimated error bars, both of which we address in the following sections.

For the analysis in this paper, we use 8.1 hours of LST, namely an RA range of 0.5-8.6 hours (A15 uses a slightly longer RA range of 0-8.6 hours; we found that some early LSTs were more severely foreground contaminated). We also use only 10 baselines, a subset of the 51 total East/West baselines used in A15, in order to illustrate our revised methods. All power spectrum results are produced for a center frequency of 151 MHz using a width of 10 MHz (20 channels), identical to the analysis in A15. In the case study in this paper, we only use one baseline type instead of the three as in A15, but Kolopanis et al. (*in prep.*) uses the full dataset presented in A15 to revise the result and place limits on the EoR at multiple redshifts.

The most significant changes from A15 occur in our revised power spectrum analysis, which is explained in the rest of this paper, but we also note that the applied fringe-rate filter is also slightly different. In A15, the applied filter was not equivalent to the optimal fringe-rate filter (which is designed to maximize power spectrum sensitivity). Instead, the optimal filter was degraded slightly by widening it in fringe-rate space. This was chosen in order to increase the number of independent modes and reduce signal loss, though as we will explain in the next section, this signal loss was still under-estimated. With the development of a new, robust method for assessing signal loss, we choose to use the optimal filter in order to maximize sensitivity. This filter is computed for a fiducial 30 m baseline at 150 MHz, the center frequency in our band. The filter in both the fringe-rate domain and time domain is shown in Figure 13.

#### 3.1. PAPER-64: Signal Loss

In Section 2.1, we showed how signal loss arises when weighting data using information from the data itself.



**Figure 13.** Top: the normalized optimal power-spectrum sensitivity weighting in fringe-rate space for our fiducial baseline and Stokes I polarization beam. Bottom: the time domain convolution kernel corresponding to the top panel. Real and imaginary components are illustrated in cyan and magenta, respectively, with the absolute amplitude in black. The fringe-rate filter acts as an integration in time, increasing sensitivity but reducing the number of independent samples in the dataset.

Here we describe a methodology for estimating the amount of signal loss caused by a particular power spectrum estimator when applied to a particular dataset. The exact amount of signal loss will depend on the specific realizations of the signals present in the data and is not something we can directly compute. In this work, as in A15, we inject simulated cosmological signals into our data and test the recovery of those signals (an approach also taken by Masui et al. (2013)). As we will see, correlations between the injected signals and the data are significant complicating factors which were previously not taken into account.

We present our PAPER-64 signal loss investigation in three parts: we first give an overview of our injection framework and illustrate how different power spectrum components affect loss. We next describe our methodology in practice and detail how we map our simulations into a posterior for the EoR signal. Finally, we build off of Section 2.1 by experimenting with different regularization schemes on PAPER data in order to minimize loss. Throughout each section, we also highlight major differences from the signal loss computation used in A15, which previously under-estimated losses.

### 3.1.1. Signal Loss Methodology

In short, our method consists of adding an EoR-like signal into the data and then measuring how much of this injected signal would be detectable given any attenuation of this signal by the (lossy) data analysis pipeline. To capture the full statistical likelihood of signal loss, one requires a quick way to generate many realizations of simulated 21 cm signal visibilities. Here we use the same method as in A15, where mock Gaussian noise visibilities (mock EoR signals) are filtered in time using an optimal fringe-rate filter to retain only “sky-like” modes.

Since the optimal filter has a shape that matches the rate of the sidereal motion of the sky, this transforms the Gaussian noise into a measurement that PAPER could make. This signal is then added to the data.<sup>2</sup>

Suppose that  $\mathbf{e}$  is the mock injected EoR signal (at some amplitude level), and  $\mathbf{x}$  is our data. We define  $\mathbf{r}$  to be the data plus the EoR signal:

$$\mathbf{r} = \mathbf{x} + \mathbf{e}. \quad (21)$$

We are interested in quantifying how much variance in  $\mathbf{e}$  is lost after weighting  $\mathbf{r}$  and estimating the power spectrum according to QE formalism. We investigate this by comparing two quantities we call the input power spectrum and output power spectrum:  $\hat{P}_{\text{in}}$  and  $\hat{P}_{\text{out}}$ , estimated using QE as

$$\hat{P}_{\text{in}}^{\alpha} \equiv \mathbf{M}_{\text{in}}^{\alpha} \mathbf{e}^{\dagger} \mathbf{I} \mathbf{Q}^{\alpha} \mathbf{I} \mathbf{e} \quad (22)$$

and

$$\begin{aligned} \hat{P}_{\text{out}}^{\alpha} &\equiv \hat{\mathbf{P}}_{\mathbf{r}}^{\alpha} \\ &= \mathbf{M}_{\mathbf{r}}^{\alpha} \mathbf{r}^{\dagger} \mathbf{R}_{\mathbf{r}} \mathbf{Q}^{\alpha} \mathbf{R}_{\mathbf{r}} \mathbf{r}, \end{aligned} \quad (23)$$

where, for illustrative purposes and notational simplicity, we have written these equations with scalar normalizations  $\mathbf{M}$ , even though for our numerical results we choose a diagonal matrix normalization using  $\mathbf{M}$  as in Equation (9).

The quantity  $\hat{P}_{\text{in}}$ , defined by Equation (22), is a uniformly weighted estimator of the power spectrum of  $\mathbf{e}$ . Since there is no noise contribution to  $\mathbf{e}$ , it can be considered the power spectrum of this particular realization of the EoR; alternatively, it can be viewed as the true power spectrum of the injected signal up to cosmic variance fluctuations. The role of  $\hat{P}_{\text{in}}$  in our analysis is to serve as a reference for the power spectrum that would be measured if there were no signal loss or other systematics. This is then to be compared to  $\hat{P}_{\text{out}}$ , which approximates the (lossy) power spectrum estimate that is output by our analysis pipeline prior to any signal loss adjustments. In principle, one could compute the true  $\hat{P}_{\text{out}}$  using end-to-end simulations of the instrument and data analysis pipeline. In practice, however, such simulations may not accurately reflect real-life systematics and foregrounds. To overcome this obstacle, one can make the assumption that since the EoR signal is expected to be small, the data vector  $\mathbf{x}$  itself is our best model of these contaminants. Making this assumption, the injected EoR signal  $\mathbf{e}$  takes on the role of the true

<sup>2</sup> One specific change from A15 is that we add this simulated signal into the analysis pipeline before the final fringe-rate filter is applied to the data. Previously, the addition was done after that final fringe-rate filter step. This change results in an increased estimate of signal loss, likely due to the use of the fringe-rate filter as a simulator. However, this pipeline difference, while significant, is not the dominant reason why signal loss was under-estimated in A15 (the dominant reason is explained in the main text in Section 3.1.1).



EOIR signal, and the sum of  $\mathbf{x}$  and  $\mathbf{e}$  (Equation (21)) replaces  $\mathbf{x}$  as our model of the measured data. Therefore,  $\hat{P}_{\text{out}}$  can be directly compared to the measurement that we make.

Under this injection framework, we can begin to see explicitly why there can be large signal loss. Expanding out Equation (23),  $\hat{P}_{\text{out}}$  becomes:

$$\begin{aligned}\hat{P}_{\text{out}}^\alpha &= M_r^\alpha (\mathbf{x} + \mathbf{e})^\dagger \mathbf{R}_r \mathbf{Q}^\alpha \mathbf{R}_r (\mathbf{x} + \mathbf{e}) \\ &= M_a^\alpha \mathbf{x}^\dagger \mathbf{R}_r \mathbf{Q}^\alpha \mathbf{R}_r \mathbf{x} + M_b^\alpha \mathbf{e}^\dagger \mathbf{R}_r \mathbf{Q}^\alpha \mathbf{R}_r \mathbf{e} \\ &\quad + M_c^\alpha \mathbf{x}^\dagger \mathbf{R}_r \mathbf{Q}^\alpha \mathbf{R}_r \mathbf{e} + M_d^\alpha \mathbf{e}^\dagger \mathbf{R}_r \mathbf{Q}^\alpha \mathbf{R}_r \mathbf{x}. \quad (24)\end{aligned}$$

Assuming  $\mathbf{R}_r$  is symmetric, the two cross-terms (terms with one copy of  $\mathbf{e}$  and one copy of  $\mathbf{x}$ ) can be summed together as:

$$\begin{aligned}\hat{P}_{\text{out}}^\alpha &= M_a^\alpha \mathbf{x}^\dagger \mathbf{R}_r \mathbf{Q}^\alpha \mathbf{R}_r \mathbf{x} + M_b^\alpha \mathbf{e}^\dagger \mathbf{R}_r \mathbf{Q}^\alpha \mathbf{R}_r \mathbf{e} \\ &\quad + 2M_c^\alpha \mathbf{x}^\dagger \mathbf{R}_r \mathbf{Q}^\alpha \mathbf{R}_r \mathbf{e}. \quad (25)\end{aligned}$$

In order to investigate the effect of each of these terms on signal loss, all three components are plotted in Figure 14 for two cases: empirically estimated inverse covariance weighting ( $\mathbf{R}_r \equiv \hat{\mathbf{C}}_r^{-1}$ ) and uniform weighting ( $\mathbf{R}_r \equiv \mathbf{I}$ ). We will now examine the behavior of this equation in three different regimes of the injected signal - very weak (left ends of the  $P_{\text{in}}$  axes in Figure 14), very strong (right ends), and in between (middle portions).

**Small injection:** In this regime, the cross-terms (red) behave as noise averaged over a finite number of samples. Output values are Gaussian distributed around zero, spanning a range of values set by the injection level. This is because  $\hat{\mathbf{R}}_r$  is dominated by the data  $\mathbf{x}$ , avoiding correlations with  $\mathbf{e}$  that can lead to solely negative power (explained further below). In fact, for the uniformly weighted case, the cross-term  $M_c^\alpha \mathbf{x}^\dagger \mathbf{I} \mathbf{Q}^\alpha \mathbf{I} \mathbf{e}$  is well modeled as a symmetric distribution with zero mean and width  $\sqrt{\hat{P}_e} \sqrt{\hat{P}_x}$ . We also note that in this regime,  $\hat{P}_r$  (black) approaches the data-only power spectrum value (gray) as expected.

**Large injection:** When the injected signal is much larger than the measured power spectrum, the data-only components can be neglected as many orders of magnitude smaller. We include a description of this regime for completeness in our discussion, but note that the upper limits that we compute are typically not determined by simulations in this regime (i.e., in using an empirical weighting scheme we've assumed the data to be dominated by foregrounds rather than the cosmological signal). However, it is useful as a check of our system in a relatively simple case. As we can see from Figure 14, the cross-terms (red) are small in comparison to the signal-only term (green). Here only does the signal-only term used in A15 dominate the total power output. We again see that, in the empirical inverse covariance weighted case, the cross-terms behave as noise (positive and negative fluctuations around zero mean).

This is for the same reason as at small injections — here  $\hat{\mathbf{C}}_r$  is dominated by the signal  $\mathbf{e}$ . The cross-correlation can again be modeled as a symmetric distribution of zero mean and width  $\sqrt{\hat{P}_e} \sqrt{\hat{P}_x}$ .

**In between:** When the injected signal is of a similar amplitude to the data by itself, the situation becomes less straightforward. We see that the weighted injected power spectrum component mirrors the input power indicating little loss (i.e., the green curve follows the dotted black line), eventually departing from unity when the injected amplitude is well above the level of the data power spectrum. However, in this regime the cross-term (red) has nearly the same amplitude, but with a negative sign. As explained below, this negativity is the result of cross-correlating inverse covariance weighted terms. This negative component drives down the  $\hat{P}_{\text{out}}$  estimator (black). We note that in A15, signal loss was computed by only looking at the second term in Equation (25) (green), which incorrectly implies no loss at the data-only power spectrum level. Ignoring the effect of the negative power from the cross-terms is the main reason for under-estimating power spectrum limits in A15.

The source of the strong negative cross-term is not immediately obvious, however it is an explainable effect. When  $\mathbf{R}_r$  is taken to be  $\hat{\mathbf{C}}_r^{-1}$ , the third term of Equation (25) is a cross-correlation between  $\hat{\mathbf{C}}_r^{-1} \mathbf{x}$  and  $\hat{\mathbf{C}}_r^{-1} \mathbf{e}$ . As shown in Switzer et al. (2015), this cross-correlation term is non-zero, and in fact negative in expectation. This negative cross-term power arises from a coupling between the inverse of  $\hat{\mathbf{C}}_r$  and  $\mathbf{x}$ . Intuitively, we can see this by expanding the empirical covariance of  $\mathbf{r} = \mathbf{x} + \mathbf{e}$ :

$$\begin{aligned}\hat{\mathbf{C}}_r &= \langle \mathbf{r} \mathbf{r}^\dagger \rangle_t \\ &= \langle \mathbf{x} \mathbf{x}^\dagger \rangle_t + \langle \mathbf{x} \mathbf{e}^\dagger \rangle_t + \langle \mathbf{e} \mathbf{x}^\dagger \rangle_t + \langle \mathbf{e} \mathbf{e}^\dagger \rangle_t, \quad (26)\end{aligned}$$

where we can neglect the first term because  $\mathbf{x}$  is small (i.e., the large negative cross-term power in the left panel of Figure 14 occurs when the injected amplitude surpasses the level of the data-only power spectrum). Without loss of generality, we will assume an eigenbasis of  $\mathbf{e}$ , so that  $\langle \mathbf{e} \mathbf{e}^\dagger \rangle_t$  is diagonal. The middle two terms, however, can have power in their off-diagonal terms due to the fact that, when averaging over a finite ensemble,  $\langle \mathbf{x} \mathbf{e}^\dagger \rangle_t$  is not zero. As shown in Appendix C of Parsons et al. (2014), to leading order the inversion of a diagonal-dominant matrix like  $\hat{\mathbf{C}}_r$  (from  $\langle \mathbf{e} \mathbf{e}^\dagger \rangle_t$ ) with smaller off-diagonal terms results in a new diagonal-dominant matrix with negative off-diagonal terms. These off-diagonal terms depend on both  $\mathbf{x}$  and  $\mathbf{e}$ . Then, when  $\hat{\mathbf{C}}_r^{-1}$  is multiplied into  $\mathbf{x}$ , the result is a vector that is similar to  $\mathbf{x}$  but contains a residual correlation to  $\mathbf{e}$  from the off-diagonal components of  $\hat{\mathbf{C}}_r^{-1}$ . The correlation is negative because the product  $\hat{\mathbf{C}}_r^{-1} \mathbf{x}$  effectively squares the  $\mathbf{x}$ -dependence of the off-diagonal terms in  $\hat{\mathbf{C}}_r^{-1}$  while re-



**Figure 14.** Illustration of the power spectrum amplitude of different power spectrum terms as a function of injected EoR power level summed into the data. Left: The empirically estimated inverse covariance weighted case used in A15. The output power spectrum ( $\hat{P}_{\text{out}} = \hat{P}_r$ ) is shown in black, and the individual terms in Equation (25) are shown in blue, red, and green. The dotted diagonal black line indicates perfect 1:1 input-to-output mapping (no signal loss). The details of the simulation used to generate the figure is explained in Section 3.1.2; here we sample a larger  $P_{\text{in}}$  range and fit smooth polynomials to our data points to make an illustrative example. The gray horizontal line is the power spectrum value of data alone,  $\hat{P}_x$  (it does not depend on injected power). The green signal-signal component is the term used in A15 to estimate signal loss. It is significantly higher than  $\hat{P}_r$  (black) when the cross-terms (red) are large and negative (black = green + red + blue). In the regime where cross-correlations between injection and data are not dominant (small and large  $P_{\text{in}}$ ), the cross-terms have a noise-like term with width  $\sqrt{\hat{P}_e} \sqrt{\hat{P}_x}$ . However, at power levels comparable to the data (the middle region), the cross-terms can produce large, negative signal due to couplings between  $\mathbf{x}$  and  $\mathbf{e}$  which affect  $\hat{\mathbf{C}}_r$ . This causes the difference between the green curve (which exhibits negligible loss at the data-only power spectrum value) and the black curve (which exhibits  $\sim 4$  orders of magnitude of loss). Right: The same power spectrum terms illustrated for the uniform weighted case.

taining the negative sign that arose from the inversion of a diagonal-dominant matrix.

**In general:** Another way to phrase the shortcoming of the empirical inverse covariance estimator (which is also discussed in Appendix B) is that it is not properly normalized. Signal loss due to couplings between the data and its weightings arise because our unnormalized quadratic estimator from Equation (8) ceases to be a quadratic quantity, and instead contains higher order powers of the data. However, the normalization matrix  $\mathbf{M}$  is derived assuming that the unnormalized estimator is quadratic in the data. The power spectrum estimate will therefore be incorrectly normalized, which manifests as signal loss. We leave a full analytic solution for  $\mathbf{M}$  for future work, since our simulations already capture the full phenomenology of signal loss and have the added benefit of being more easily generalizable in the face of non-Gaussian systematics.

### 3.1.2. Signal Loss in Practice

We now shift our attention towards computing upper limits on the EoR signal for the fringe-rate filtered PAPER-64 dataset in a way that accounts for signal loss. While our methodology outlined below is independent of weighting scheme, here we demonstrate the computation using empirically estimated inverse covariance weighting ( $\mathbf{R} \equiv \hat{\mathbf{C}}^{-1}$ ), the weighting scheme used in A15 which leads to substantial loss. With this weighting, our expressions for  $\hat{P}_{\text{in}}$  and  $\hat{P}_{\text{out}}$  become:

$$\hat{P}_{\text{in}}^\alpha = \mathbf{M}_{\text{in}}^\alpha \mathbf{e}^\dagger \mathbf{I} \mathbf{Q}^\alpha \mathbf{I} \mathbf{e} \quad (27)$$

$$\hat{P}_{\text{out}}^\alpha = \mathbf{M}_r^\alpha \mathbf{r}^\dagger \hat{\mathbf{C}}_r^{-1} \mathbf{Q}^\alpha \hat{\mathbf{C}}_r^{-1} \mathbf{r}. \quad (28)$$

One issue to address is how one incorporates the randomness of  $\hat{P}_{\text{out}}$  into our signal loss corrections. A different realization of the mock EoR signal is injected with each bootstrap run, causing the output to vary in three

ways — there is noise variation from the bootstraps, there is cosmic variation from generating multiple realizations of the mock EoR signal, and there is a variation caused by whether the injected signal looks more or less “like” the data (i.e., how much coupling there is, which affects how much loss results).

For each injection level, the true  $P_{\text{in}}$  is simply the average of our bootstrapped estimates  $\hat{P}_{\text{in}}$ , since  $\hat{P}_{\text{in},\alpha}$  is by construction an unbiased estimator. Phrased in the context of Bayes’ rule, we wish to find the posterior probability distribution  $p(P_{\text{in}}|\hat{P}_{\text{out}})$ , which is the probability of  $P_{\text{in}}$  given the uncorrected/measured power spectrum estimate  $\hat{P}_{\text{out}}$ . Bayes’ rule relates the posterior, which we don’t know, to the likelihood, which we can forward model. In other words,

$$p(P_{\text{in}}|\hat{P}_{\text{out}}) \propto \mathcal{L}(\hat{P}_{\text{out}}|P_{\text{in}})p(P_{\text{in}}), \quad (29)$$

where  $\mathcal{L}$  is the likelihood function defined as the distribution of data plus injection ( $\hat{P}_{\text{out}}$ ) given the injection  $P_{\text{in}}$ . We construct this distribution by fixing  $P_{\text{in}}$  and simulating our analysis pipeline for many realizations of the injected EoR signal consistent with this power spectrum. The resulting distribution is normalized such that the sum over  $\hat{P}_{\text{out}}$  is unity, and the whole process is then repeated for a different value of  $P_{\text{in}}$ .

The implementation details of the injection process require some more detailed explanation. In our code, we add a new realization of EoR to each independent bootstrap of data (see Section 3.2.2 for a description of PAPER’s bootstrapping routine) with the goal of simultaneously capturing cosmic variance, noise variance, and signal loss. To limit computing time we perform 20 realizations of each  $P_{\text{in}}$  level. We also run 50 total EoR injection levels, yielding  $P_{\text{in}}$  values that range from  $\sim 10^5 \text{ mK}^2 (h^{-1} \text{ Mpc})^3$  to  $\sim 10^{11} \text{ mK}^2 (h^{-1} \text{ Mpc})^3$ , resulting in a total of 1000 data points on our  $P_{\text{in}}$  vs.  $\hat{P}_{\text{out}}$  grid.

Going forward, we treat every  $k$ -value separately in order to determine an upper limit on the EoR signal per  $k$ . We bin our simulation outputs along the  $P_{\text{in}}$  axis (one bin per injection level) and, since they are well-approximated by a Gaussian distribution in our numerical results, we smooth the distribution of  $\hat{P}_{\text{out}}$  values by fitting Gaussians for each bin based on its mean and variance (and normalize them). Stitching all of them together results in a 2-dimensional transfer function — the likelihood function in Bayes’ rule, namely  $\mathcal{L}(\hat{P}_{\text{out}}|P_{\text{in}})$ . We then have a choice for our prior,  $p(P_{\text{in}})$ , and we choose to invoke a Jeffreys prior (Jaynes 1968) because it is a true uninformative prior. For a derivation and more details about the Jeffreys prior used in our analysis, see Appendix C.

Finally, our transfer functions are shown in Figure 15 for both the weighted (left) and unweighted (right) cases. Our bootstrapped power spectrum outputs are shown as black points and the colored heat-map overlaid on top is the likelihood function modified by our prior.

Although we only show figures for one  $k$ -value, we note that the shape of the transfer curve is similar for all  $k$ ’s. We then invoke Bayes’ interpretation and re-interpret it as the posterior  $p(P_{\text{in}}|\hat{P}_{\text{out}})$  where we recall that  $\hat{P}_{\text{out}}$  is a model of our data. To do this we make a horizontal cut across at the data value  $\hat{P}_x$  (setting  $\hat{P}_{\text{out}} = \hat{P}_x$ ), shown by the gray solid line, to yield a posterior distribution for the signal. We normalize this final distribution and compute the 95% confidence interval (an upper limit on EoR).

By-eye inspection of the transfer function in Figure 15 gives a sense of what the signal loss result should be. The power spectrum value of our data,  $\hat{P}_x$  is marked by the solid gray horizontal lines. From the left plot (empirically estimated inverse covariance weighting), one can eyeball that a data value of  $10^5 \text{ mK}^2 (h^{-1} \text{ Mpc})^3$ , for example, would map approximately to an upper limit of  $\sim 10^9 \text{ mK}^2 (h^{-1} \text{ Mpc})^3$ , implying a signal loss factor of  $\sim 10^4$ . For the uniform-weighted case (right plot), we see no loss at a data value of  $\sim 10^7 \text{ mK}^2 (h^{-1} \text{ Mpc})^3$ .

The loss-corrected power spectrum limit for empirically estimated inverse covariance weighted PAPER-64 data is shown in Figure 16 (solid red), which we can compare to the original lossy result (dashed red). Post-signal loss estimation, the power spectrum limits are higher than both the theoretical noise level (green) and uniform-weighted power spectrum (which is shown three ways: black and gray points are positive and negative power spectrum values, respectively, with  $2\sigma$  error bars from bootstrapping, the solid blue is the upper limit on the EoR signal using the full signal injection framework, and the shaded gray is the power spectrum values with thermal noise errors). We elaborate on this point in the next section, as well as investigate alternate weighting schemes to inverse covariance weighting, with the goal of finding one that balances the aggressiveness of down-weighting contaminants and minimizing the loss of the EoR signal.

### 3.1.3. Minimizing Signal Loss

With a signal loss formalism established, we now have the capability of experimenting with different weighting options for  $\mathbf{R}$ . Our goal here is to choose a weighting method that successfully down-weights foregrounds and systematics in our data without generating large amounts of signal loss as we have seen with the inverse covariance estimator. We have found that the balance between the two is a delicate one and requires a careful understanding and altering of empirical covariances.

We saw in Section 2.1.4 how limiting the number of down-weighted eigenmodes (i.e., flattening out part of the eigenspectrum and effectively decoupling the lowest-valued eigenmodes, which are typically EoR-dominated, from the data) can help minimize signal loss. We experiment with this idea on PAPER-64 data, dialing the number of modes that are down-weighted from zero (which is equivalent to identity-weighting, or the uniform-weighted case) to 21 (which is the full inverse



**Figure 15.** Signal loss transfer functions showing the relationship of  $P_{\text{in}}$  and  $\hat{P}_{\text{out}}$ , as defined by Equations (22) and (23). Power spectra values (black points) are generated for 20 realizations of  $\mathbf{e}$  per signal injection level. Since our  $\hat{P}_{\text{out}}$  values are well-approximated by a Gaussian distribution, we fit Gaussians to each injection level based on the mean and variance of the simulation outputs. This entire likelihood function is then multiplied by a Jeffreys prior for  $p(P_{\text{in}})$ , with the final result shown as the colored heat-maps on top of the points. Two cases are displayed: empirically estimated inverse covariance weighted PAPER-64 data (left) and uniform-weighted data (right). The dotted black diagonal lines mark a perfect unity mapping, and the solid gray horizontal line denotes the power spectrum value of the data  $\hat{P}_x$ , from which a posterior distribution for the signal is extracted. From these plots, it is clear that the weighted case results in  $\sim 4$  orders of magnitude of signal loss at the data-only power spectrum value, whereas the uniform-weighted case does not exhibit loss. The general shape of these transfer functions are also shown by the black curves in Figure 14 for comparison.

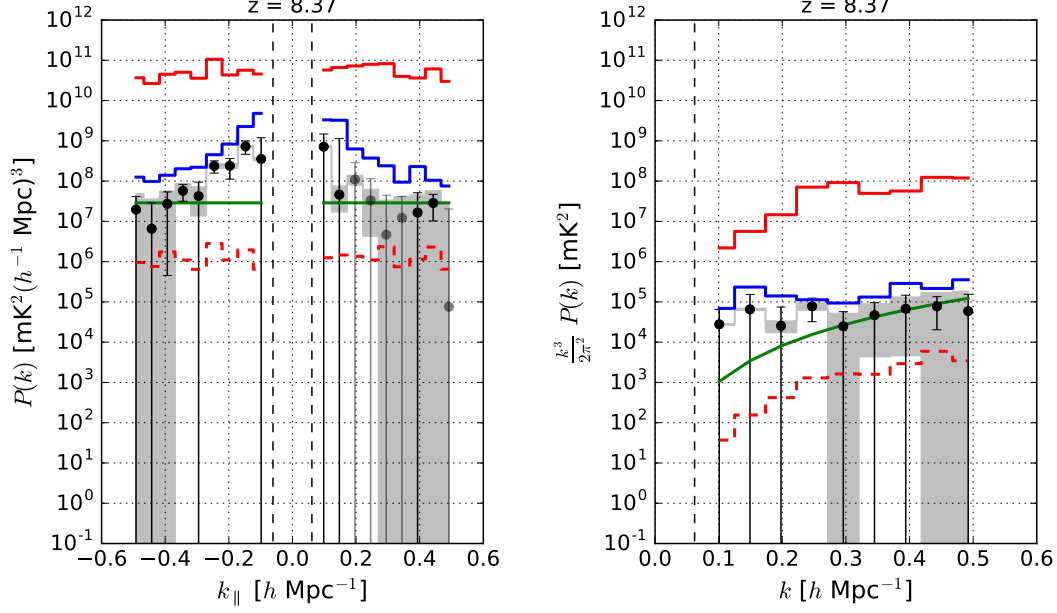
covariance estimator). The power spectrum results for one  $k$ -value, both before and after signal loss estimation, are shown in the top panel in Figure 17. We see that the amount of signal loss increases as weighting becomes more aggressive (dashed red). In other words, more EoR-dominated fluctuations are being overfit and subtracted as more modes are down-weighted. We also find that the power spectrum upper limit, post-signal loss estimation, increases with the number of down-weighted modes (solid red). The more modes we use in down-weighting, the stronger the coupling between the weighting and the data, and the greater the error we have in estimating the power spectrum. Switzer et al. (2013) took a similar approach in determining the optimal number of modes to down-weight in GBT data, finding similar trends and noting that removing too few modes is limited by residual foregrounds and removing too many modes is limited by large error bars and signal loss.

Optimistically, we expect there to be a “sweet spot” as we dial our regularization knob; a level of regularization where weighting is beneficial compared to uniform weighting (blue). In other words, we would like a weighting scheme that down-weights eigenmodes that predominantly describe foreground modes, but not EoR modes. We see in Figure 17 that this occurs roughly when only

the  $\sim 2$ -3 highest-valued eigenmodes are down-weighted and the rest are given equal weights (though for the case shown, weighting does not actually outperform uniform weighting). For a similar discussion on projecting out modes (zeroing out eigenmodes, rather than just ignoring their relative weightings as we do in this study), see Switzer et al. (2013).

We also saw in Section 2.1.4 how adding the identity matrix to the empirical covariance can minimize signal loss. We experiment with this idea as well, shown in the bottom panel of Figure 17. The dashed red and solid red lines represent power spectrum limits pre and post-signal loss estimation, respectively, as a function of the strength of  $\mathbf{I}$  that is added to  $\hat{\mathbf{C}}$ , quantified as a percentage of  $\text{Tr}(\hat{\mathbf{C}})\mathbf{I}$  added to  $\hat{\mathbf{C}}$ . We parameterize this “regularization strength” parameter as  $\gamma$ , namely  $\hat{\mathbf{C}} \equiv \hat{\mathbf{C}} + \gamma \text{Tr}(\hat{\mathbf{C}})\mathbf{I}$ . From this plot we see that only a small percentage of  $\text{Tr}(\hat{\mathbf{C}})$  is needed to significantly reduce loss. We expect that as the strength of  $\mathbf{I}$  is increased (going to the left), both the red curves will approach the uniform-weighted case. We also notice that the post-signal loss limit hovers around the uniform-weighted limit for a large range of regularization strengths and while an overall trend from high-to-low signal loss is seen as the strength increases, there does not appear to be a clear “minimum” that produces





**Figure 16.** A power spectrum of a subset of PAPER-64 data illustrating the use of empirical inverse covariance weighting. The solid red curve is the  $2\sigma$  upper limit on the EoR signal estimated from our signal injection framework using empirical inverse covariance weighting. Shown for comparison is the lossy limit prior to signal loss estimation (dashed red). The theoretical  $2\sigma$  thermal noise level prediction based on observational parameters is in green, whose calculation is detailed in Section 3.2.1. Additionally, the power spectrum result for the uniform weighted case is shown in three different ways: power spectrum values (black and gray points as positive and negative values, respectively, with  $2\sigma$  error bars from bootstrapping), the  $2\sigma$  upper limit on the EoR signal using our full signal injection framework (solid blue), and the measured power spectrum values with  $2\sigma$  thermal noise errors (gray shaded regions). The vertical dashed black lines signify the horizon limit for this analysis using 30 m baselines. In this example, we see that the lossy power spectrum limit is  $\sim 4$  orders of magnitude too low when using empirical inverse covariance weighting.

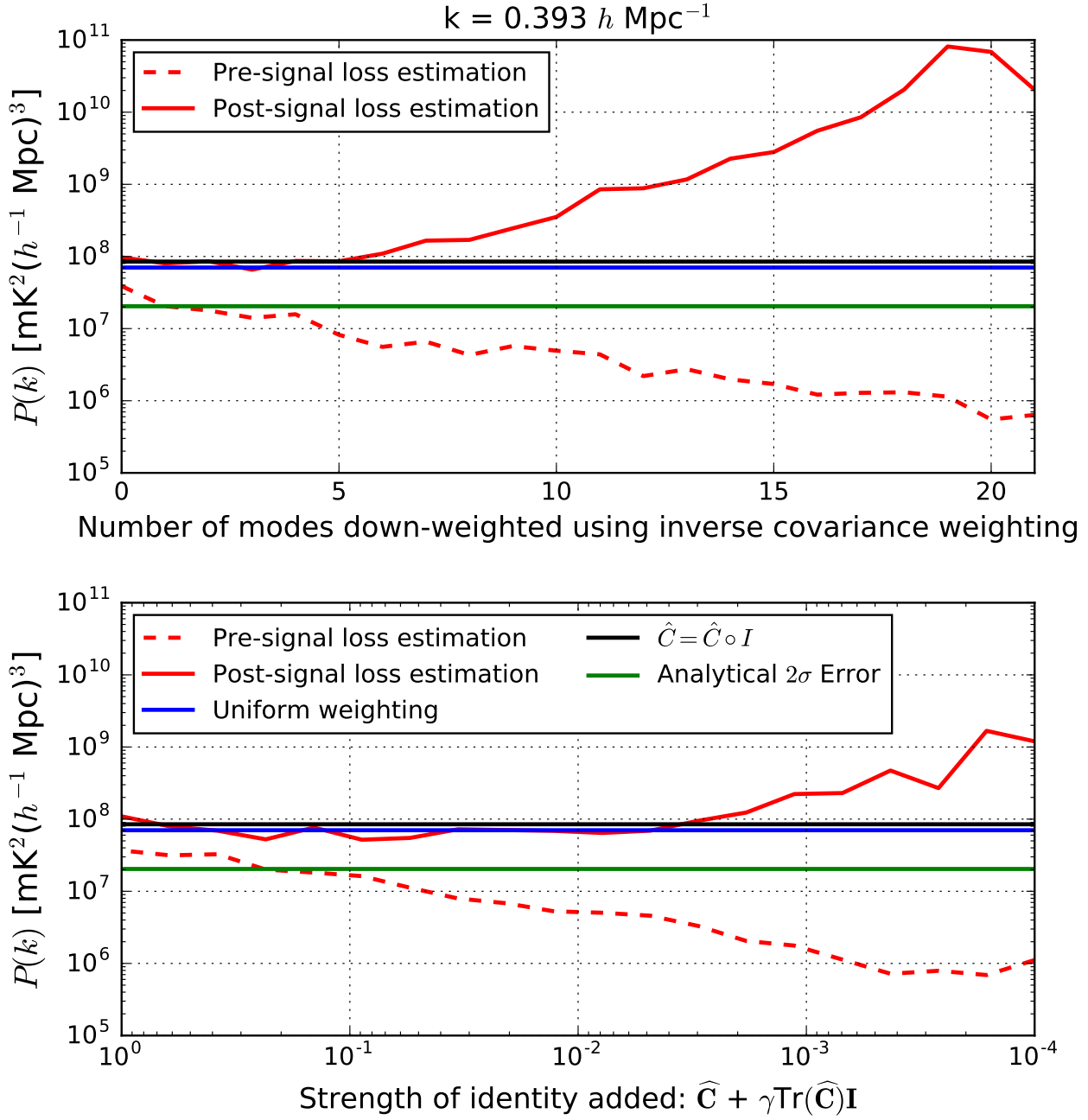
the least loss.

In addition to our thermal noise prediction (green) and uniform-weighted power spectrum limit (blue), one additional horizontal line is shown in Figure 17 in both panels and represents a third regularization technique. This line (black) denotes the power spectrum value, post-signal loss estimation, for inverse variance weighting (multiplying an identity matrix element-wise to  $\hat{\mathbf{C}}$ ). This result is single-valued and not a function of the horizontal axis. We see that all three regularization schemes shown (solid red top panel, solid red bottom panel, black) perform similarly at their best (i.e., when  $\sim 2$ -3 eigenmodes are down-weighted in the case of the top panel's solid red curve). However, for the remainder of this paper, we choose to use the weighting option of  $\hat{\mathbf{C}} \equiv \hat{\mathbf{C}} + 0.09 \text{Tr}(\hat{\mathbf{C}})\mathbf{I}$ , or  $\gamma = 0.09$ , which we will denote as  $\hat{\mathbf{C}}_{\text{eff}}$ . We choose this weighting scheme merely as a simple example of regularizing PAPER-64 covariances, noting that the power spectrum upper limit remains roughly constant for a broad range of values of  $\gamma$ . We also note that all of the regularizations described here perform similarly to the uniformly weighted case. Thus, in Kolopanis et al. (*in prep.*), uniform weightings are used to produce straightforward power spectrum limits that do not suffer from loss.

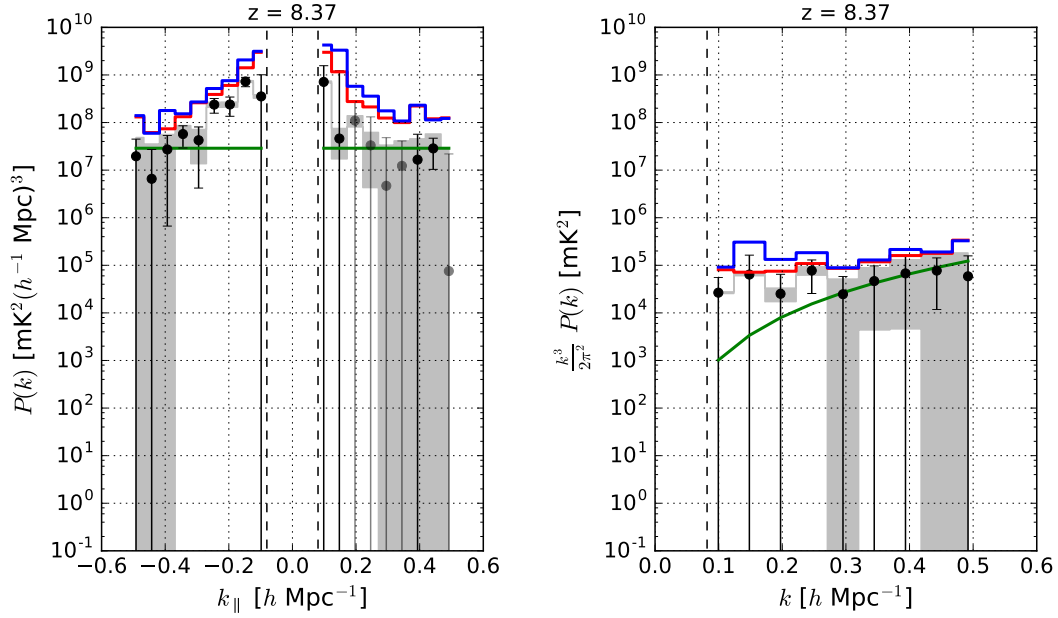
The power spectrum result for our subset of PAPER-64 data (using only one baseline separation type, 10

baselines, and  $\hat{\mathbf{C}}_{\text{eff}}$ ) is shown in Figure 18. Again, the solid red curve represents our upper limit on the EoR signal using the full signal injection framework. The uniform weighted case is shown as the black and gray points, which correspond to positive and negative power spectrum values respectively (with  $2\sigma$  errors bars from bootstrapping). It is also shown as an upper limit using the signal injection framework (solid blue), which is interestingly larger than the errors computed from bootstrapping, likely because the full injection framework takes into account additional sample variance whereas the bootstrapped errors do not. Finally, the gray shaded regions combine the measured uniform weighted power spectrum values with thermal noise errors. We show this power spectrum result as one example of how a simple regularization of an empirical covariance matrix can minimize signal loss, though we also note that this weighting does not produce more stringent limits than the uniform weighted case.

In this section we have shown three simple ways of regularizing  $\hat{\mathbf{C}}$  to minimize signal loss using PAPER-64 data. There are many other weighting schemes that we leave for consideration in future work. For example, one could estimate  $\hat{\mathbf{C}}$  using information from different subsets of baselines. For redundant arrays this might mean calculating  $\hat{\mathbf{C}}$  from a different but similar baseline type, such as the  $\sim 30$  m diagonal PAPER baselines (instead



**Figure 17.** Power spectra  $2\sigma$  upper limits for  $k = 0.393 h \text{ Mpc}^{-1}$  for fringe-rate filtered PAPER-64 data. Top: Values are shown before (dashed red) and after (solid red) signal loss estimation via our signal injection framework as a function of number of eigenmodes of  $\hat{\mathbf{C}}$  that are down-weighted. This regularization knob is tuned from 0 modes on the left (i.e., unweighted) to 21 modes on the right (i.e., the full inverse covariance estimator).  $\sim 4$  orders of magnitude of signal loss results when using empirically estimated inverse covariance weighting. Bottom: Power spectrum upper limits before (dashed red) and after (solid red) signal loss estimation as a function of identity added to the empirical covariance. This regularization knob is tuned from  $\gamma = 10^{-4}$  on the right (i.e., very little regularization) to  $\gamma = 1$  on the left (see main text for the definition of  $\gamma$ ). Also plotted in both panels for comparison are  $2\sigma$  power spectrum upper limits for the uniform-weighted case (blue) and inverse variance weighted case (black); both are after signal loss estimation. Finally, a theoretical prediction for noise ( $2\sigma$  error) is plotted as green. In the PAPER-64 analysis in this paper, we choose to use a regularization scheme of  $\hat{\mathbf{C}}_{\text{eff}} \equiv 0.09 \text{Tr}(\hat{\mathbf{C}})\mathbf{I} + \hat{\mathbf{C}}$  ( $\gamma = 0.09$ ) as a simple example of regularization that minimizes loss, and note that the power spectrum limits using this type of regularization are roughly constant across a large range of values of  $\gamma$ .



**Figure 18.** A power spectrum of a subset of PAPER-64 data illustrating the use of  $\hat{\mathbf{C}}_{\text{eff}}$  to minimize signal loss. The solid red curve is the  $2\sigma$  upper limit on the EoR signal estimated from our signal injection framework. The theoretical  $2\sigma$  thermal noise level prediction based on observational parameters is in green. Additionally, the power spectrum result for the uniform weighted case is shown in three different ways: power spectrum values (black and gray points as positive and negative values, respectively, with  $2\sigma$  error bars from bootstrapping), the  $2\sigma$  upper limit on the EoR signal using our full signal injection framework (solid blue), and the measured power spectrum values with  $2\sigma$  thermal noise errors (gray shaded regions). The vertical dashed black lines signify the horizon limit for this analysis using 30 m baselines. This power spectrum result does not use the full dataset’s sensitivity as in A15 and Kolopanis et al. (*in prep.*), though we include all analysis changes which have mostly stemmed from revisions regarding signal loss, bootstrapping, and the theoretical error computation. We see that the regularization scheme used here produces limits similar to the unweighted limits.

of the horizontal E/W ones). Alternatively, covariances could be estimated from all baselines except the two being cross-multiplied when forming a power spectrum estimate. This method was used in [Parsons et al. \(2014\)](#) (a similar method was also used in [Dillon et al. \(2015\)](#)) in order to avoid suppressing the 21 cm signal, and it is worth noting that the PAPER-32 results are likely less impacted from the issue of signal loss underestimation because of this very reason (however, they are affected by the error estimation issues described in Section 3.2, so we also regard those results as suspect and superseded by those of Kolopanis et al. (*in prep.*)).

Another possible way to regularize  $\hat{\mathbf{C}}$  is to use information from different ranges of LST. For example, one could calculate  $\hat{\mathbf{C}}$  with data from LSTs where foregrounds are stronger (earlier or later LSTs than the “foreground-quiet” range typically used in forming power spectra) — doing so may yield a better description of the foregrounds that we desire to down-weight, especially if residual foreground chromaticity is instrumental in origin and stable in time. Fundamentally, each of these examples are similar in that they rely on a computation of  $\hat{\mathbf{C}}$  from data that is similar but not exactly the same as the data that is being down-weighted. Ideally this would be effective in down-weighting shared contaminants yet avoid signal loss from over-fitting EoR modes in the power spectrum dataset itself.

In Section 3.1, we have detailed several aspects of signal loss in PAPER-64: how the loss arises, how it can be estimated from an injection framework, and ways it can be minimized. We again emphasize that these lessons learned about signal loss are largely responsible for shaping our revised analysis of PAPER data. In the remainder of this paper, we will transition to other new aspects of our analysis, framed within the context of error estimation and (non-EoR) bias in PAPER-64.

### 3.2. PAPER-64: Error Estimation

In this section we discuss the ways in which we estimate errors for PAPER-64 power spectra. We first walk through an expression for a theoretical error estimation (of thermal noise) based on observational parameters. Although a theoretical model often differs from true errors, it is helpful to understand the ideal case and the factors that affect its sensitivity. Additionally, we build on the lessons learned about bootstrapping in Section 2.2 to revise our bootstrapping method as applied to PAPER-64 data in order to compute accurate errors from the data itself.

In particular, we highlight major changes in both our sensitivity calculation and bootstrapping method that differ from the [A15](#) analysis of PAPER-64. While we do not discuss the changes within the context of PAPER-32, it is worth noting that the power spectrum results in [Parsons et al. \(2014\)](#) are affected by the same issues.

#### 3.2.1. Theoretical Error Estimation

Re-analysis of the PAPER-64 data included a detailed study using several independently generated noise simulations. What we found was that these simulations all agreed but were discrepant with the previous analytic sensitivity calculations. The analytic calculation is only an approximation; however, the differences were large enough (factors of 10 in some cases) to warrant a careful investigation. The analytic calculation attempts to combine a large number of pieces of information in an approximate way, and when re-considering some of the approximations, we have found there to be large effects. What follows here is an accounting of the differences which have been discovered. Our revised theoretical error estimation, which is plotted as the solid green curve in many of the previous power spectrum plots, is computed with these changes accounted for.

The noise prediction  $n(k)$  ([Parsons et al. 2012a](#); [Pober et al. 2013](#)) for a power spectral analysis of interferometric 21 cm data, in temperature-units, is:

$$N(k) = \frac{X^2 Y \Omega_{\text{eff}} T_{\text{sys}}^2}{\sqrt{2 N_{\text{lst}} N_{\text{seps}} t_{\text{int}} N_{\text{days}} N_{\text{bls}} N_{\text{pols}}}}. \quad (30)$$

We will now explain each factor in Equation (30) and highlight key differences from the numbers used in [A15](#).

- $X^2 Y$ : Conversion factors from observing coordinates (angles on the sky and frequency) to cosmological coordinates (co-moving distances). For  $z = 8.4$ ,  $X^2 Y = 5 \times 10^{11} h^{-3} \text{ Mpc}^3 \text{ str}^{-1} \text{ GHz}^{-1}$ .
- $\Omega_{\text{eff}}$ : The effective primary beam area in steradians ([Parsons et al. 2010](#); [Pober et al. 2012](#)). The effective beam area changes with the application of a fringe-rate filter, since different parts of the beam are up-weighted and down-weighted. Using numbers from Table 1 in [Parsons et al. \(2016\)](#),  $\Omega_{\text{eff}} = 0.74^2 / 0.24$  for an optimal fringe-rate filter and the PAPER primary beam.
- $T_{\text{sys}}$ : The system temperature is set by:

$$T_{\text{sys}} = 180 \left( \frac{\nu}{0.18} \right)^{-2.55} + T_{\text{rcvr}}, \quad (31)$$

where  $\nu$  are frequencies in GHz ([Thompson et al. 2001](#)). We use a receiver temperature of 144 K, yielding  $T_{\text{sys}} = 431 \text{ K}$  at 150 MHz. This is lower than the  $T_{\text{sys}}$  of 500 K used in [A15](#) because of several small miscalculation errors that were identified<sup>3</sup>.

- $\sqrt{2}$ : This factor in the denominator of the sensitivity equation comes from taking the real part of the power spectrum estimates after cross-multiplying

<sup>3</sup> For example, there was a missing a square root in going from a variance to a standard deviation.



independent “even” and “odd” visibility measurements (this cross-multiplication is done principally to avoid a noise bias). In A15, a factor of 2 was mistakenly used.

- $N_{\text{lst}}$ : The number of independent LST bins that go into a power spectrum estimation. The sensitivity scales as the square root because we integrate incoherently over time. For PAPER-64,  $N_{\text{lst}} = 8$ .
- $N_{\text{seps}}$ : The number of baseline separation types (where baselines of a unique separation type have the same orientation and length) averaged incoherently in a final power spectrum estimate. For the analysis in this paper, we only use one type of baseline (PAPER’s 30 m East/West baselines). However, both the updated limits in Kolopanis et al. (*in prep.*) and the sensitivity prediction in Figure 19 use three separation types ( $N_{\text{seps}} = 3$ ) to match A15.
- $t_{\text{int}}$ : Length of an independent integration of the data. It is crucial to adapt this number if filtering is applied along the time axis (i.e., a fringe-rate filter). We compute the effective integration time of our fringe-rate filtered data by scaling the original integration time  $t_i$  using the following:

$$t_{\text{int}} = t_i \int \frac{1 df}{w^2(f) df}, \quad (32)$$

where  $t_i = 43$  seconds,  $t_{\text{int}}$  is the fringe-rate filtered integration time,  $w$  is the fringe-rate profile, and the integral is taken over all fringe-rates. For PAPER-64, this number is  $t_{\text{int}} = 3857$  s.

- $N_{\text{days}}$ : The total number of days of data analyzed. In A15, this number was set to 135. However, because we divide our data in half (to form “even” and “odd” datasets, or  $N_{\text{datasets}} = 2$ ), this number should reflect the number of days in each individual dataset instead of the total. Additionally, this number should be adjusted to reflect the actual number of cross-multiplications that occur between datasets (“even” with “odd” and “odd” with “even”, but not “odd” with “odd” or “even” with “even”, for reasons explained in Section 3.3.1). Finally, because our LST coverage is not 100% complete (it doesn’t overlap for every single day), we incorporate a root-mean-square statistic in computing a realistic value of  $N_{\text{days}}$ . Our expression therefore becomes:

$$N_{\text{days}} = \sqrt{\langle N_i^2 \rangle} \sqrt{(N_{\text{datasets}}^2 - N_{\text{datasets}})} \quad (33)$$

where  $i$  indexes LST and frequency channel over all datasets (Jacobs et al. 2015). For PAPER-64, our revised estimate of  $N_{\text{days}}$  is  $\sim 47$  days.

- $N_{\text{bls}}$ : The number of baselines contributing to the sensitivity of a power spectrum estimate.

In A15, this number was the total number of 30 m East/West baselines used in the analysis. However, using the total number of baselines ( $N_{\text{bls\_total}} = 51$ ) neglects the fact that the A15 analysis averages baselines into groups for computational speed-up when cross-multiplying data. Our revised estimate for the parameter is:

$$N_{\text{bls}} = \frac{N_{\text{bls\_total}}}{N_{\text{gps}}} \sqrt{\frac{N_{\text{gps}}^2 - N_{\text{gps}}}{2}}, \quad (34)$$

where, in the A15 analysis,  $N_{\text{gps}} = 5$ . Each baseline group averages down linearly as the number of baselines entering the group ( $N_{\text{bls\_total}}/N_{\text{gps}}$ ) and then as the square root of the number of cross-multiplied pairs ( $\sqrt{\frac{N_{\text{gps}}^2 - N_{\text{gps}}}{2}}$ ). A revised A15 analysis should therefore use  $N_{\text{bls}} \sim 32$  instead of 51, and this change is taken into account in Figure 19. However, the analysis in this paper and in Kolopanis et al. (*in prep.*) no longer averages baselines into groups ( $N_{\text{gps}} = 1$ ). For the subset of data presented in this paper,  $N_{\text{bls}} = 10$ .

- $N_{\text{pols}}$ : The number of polarizations averaged together. For the case of Stokes I,  $N_{\text{pols}} = 2$ .

An additional factor of  $\sqrt{2}$  is gained in sensitivity when folding together positive and negative  $k$ ’s to form  $\Delta^2(k)$ .

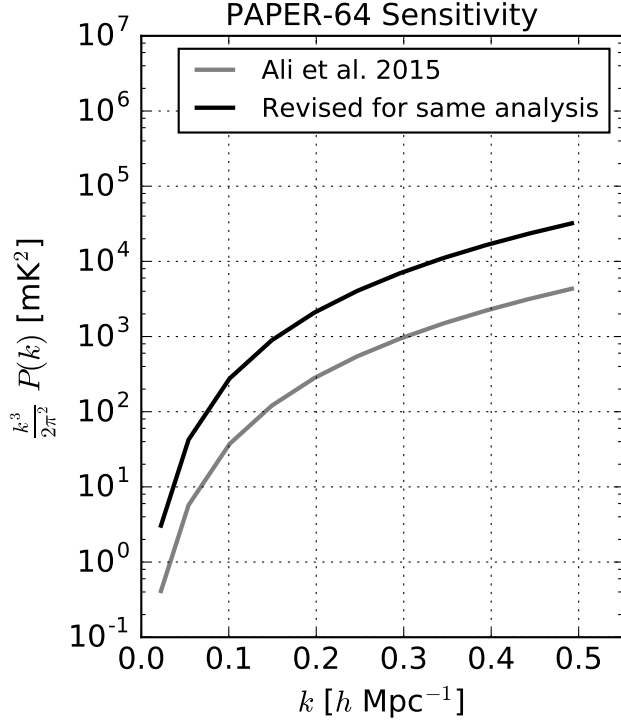
Our revised sensitivity estimate for the A15 analysis of PAPER-64 is shown in Figure 19. Together, the revised parameters yield a decrease in sensitivity (higher noise floor) by a factor of  $\sim 7$  in mK<sup>2</sup>.

To verify our thermal noise prediction, we form power spectra estimates using a pure noise simulation. We create Gaussian random noise assuming a constant  $T_{\text{rcvr}}$  (translated into  $T_{\text{sys}}$  via Equation (31)) but accounting for the true  $N_{\text{days}}$  as determined by LST sampling counts for each time and frequency in the LST-binned data. We convert  $T_{\text{sys}}$  into a root-mean-square variance statistic using:

$$T_{\text{rms}} = \frac{T_{\text{sys}}}{\sqrt{\Delta\nu \Delta t N_{\text{days}} N_{\text{pols}}}}, \quad (35)$$

where  $\Delta\nu$  is the channel spacing,  $\Delta t$  is the integration time,  $N_{\text{days}}$  is the number of daily counts for a particular time and frequency that went into our LST-binned set, and  $N_{\text{pols}}$  is the number of polarizations (2 for Stokes I). This temperature sets the variance of the Gaussian random noise.

Power spectrum results for the noise simulation, which uses our full power spectrum pipeline, are shown in Figure 20. We highlight that the bootstrapped data (black and gray points, with  $2\sigma$  error bars) and thermal noise prediction (solid green) show good agreement, as bootstrapping provides an accurate estimate of the noise variance. However, the limits from the full signal loss framework (weighted and unweighted in red and blue,



**Figure 19.** An updated prediction for the thermal noise level of PAPER-64 data (black) is shown in comparison to previously published sensitivity limits (gray), both computed for the parameters and methods used in A15. Major factors that contribute to the discrepancy are  $\Omega_{\text{eff}}$ ,  $N_{\text{days}}$  and  $N_{\text{bls}}$ , as in Equation (30) and described in Section 3.2.1, which when combined decreases our sensitivity (higher noise floor) by a factor of  $\sim 7$  in  $\text{mK}^2$ .

respectively) are inflated, likely due to the additional inclusion of sample variance that comes from the EoR simulations.

### 3.2.2. Bootstrapping

We bootstrap PAPER-64 power spectra in order to determine confidence intervals for our results. In this section, we highlight four specific changes in the way we estimate errors since A15, the first of which builds off of the lesson we have learned previously about bootstrapping independent samples.

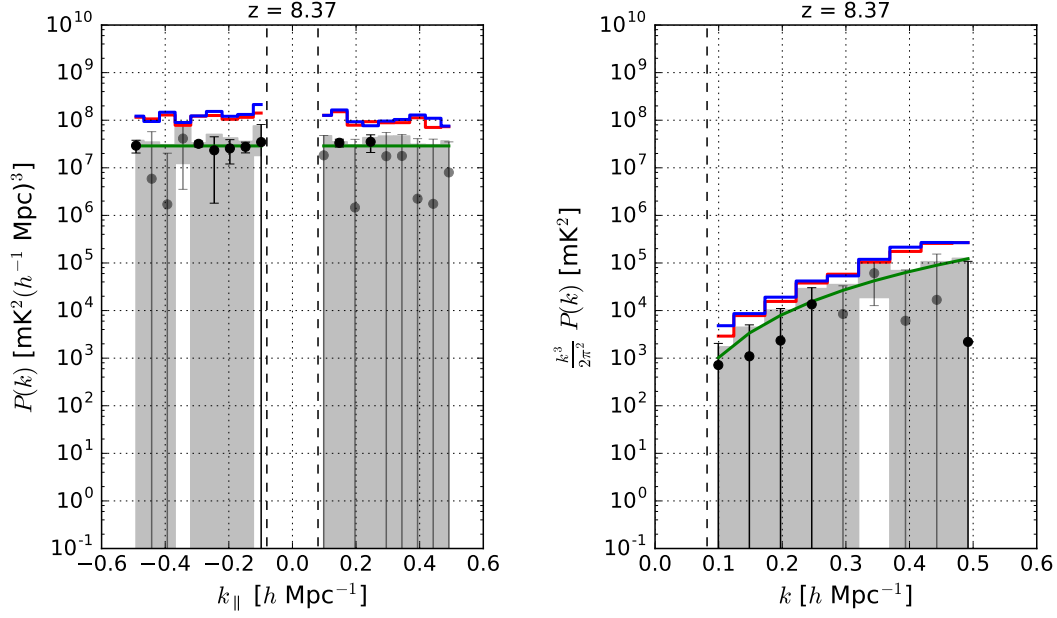
- As discussed in Section 2.2, bootstrapping is only a valid way of estimating errors if a dataset is comprised of independent samples, or the number of independent samples is well known. The PAPER-64 pipeline outputs 20 bootstraps (over baselines), each a 2-dimensional power spectrum that is a function of  $k$  and time.

In A15, a second round of bootstrapping occurred over the time axis. A total of 400 bootstraps were created in this step ( $N_{\text{boot}} = 400$ ), each comprised of randomly selected values sampled with replacement along the time axis. More specifically, each of these bootstraps contained the same number of values as the number of time integrations (which,

at  $\sim 700$ , greatly exceeds the approximate number of independent samples after fringe-rate filtering). Means were then taken of the values in each bootstrap. Finally, power spectrum limits were computed by taking the mean and standard deviation over all the bootstraps. We emphasize again that in this previous analysis, the number of elements sampled per bootstrap greatly exceeded the number of independent LST samples, under-estimating errors. A random draw of 700 measurements from this dataset has many repeated values, and the variance between hundreds ( $N_{\text{boot}}$ ) of these random samples is smaller than the true underlying variance of the data.

Given our new understanding of the sensitivity of bootstraps to the number of elements sampled, we have removed the second bootstrapping step along time entirely and now simply bootstrap over the baseline axis. Power spectrum  $2\sigma$  errors (computed from bootstrap variances) with this bootstrapping change for fringe-rate filtered noise are shown in Figure 21. The estimates are uniformly weighted in order to disentangle the effects of bootstrapping from signal loss. As shown in the figure, when more elements are drawn for each bootstrap than the number of independent samples (by over-sampling elements along the time axis), repeated values begin to crop up and the apparent variation between bootstraps drops, resulting in limits (gray) below the predicted noise level (green). Using the revised bootstrapping method, where bootstrapping only occurs over the baseline axis, the limits (black) are shown to better agree with the analytic prediction for noise. While Figure 21 implies that errors are under-estimated by a factor of  $\sim 5$  in  $\text{mK}^2$  for the noise simulation, in practice this factor is lower for the case of real data (a factor of  $\sim 3$  in  $\text{mK}^2$  instead), possibly due to the data being less correlated in time than the fringe-rate filtered noise in the simulation.

- A second change to our bootstrapping procedure is that we now bootstrap over baseline cross-products, instead of the baselines themselves. In the previous analysis, baselines were bootstrapped prior to forming cross power spectra, and using this particular ordering of operations (bootstrapping, then cross-multiplication) yields variances that have been found to disagree with predicted errors from bootstrapping using simulations. On the contrary, bootstrapping over cross power spectra ensures that we are estimating the variance of our quantity of interest (i.e., the power spectrum). This change, while fundamental in retaining the integrity of the bootstrapping method in general, alters the resulting power spectrum errors by factors of  $< 2$  in practice.



**Figure 20.** The power spectrum for a noise simulation that mimics the noise level of a subset of PAPER-64 data, where the solid red curve is the  $2\sigma$  upper limit on the EoR signal estimated from our signal injection framework using  $\hat{C}_{\text{eff}}$ . The theoretical  $2\sigma$  thermal noise level prediction based on observational parameters (calculated by Equation (30)) is in green. Additionally, the power spectrum result for the uniform weighted case is shown in three different ways: power spectrum values (black and gray points as positive and negative values, respectively, with  $2\sigma$  error bars from bootstrapping), the  $2\sigma$  upper limit on the EoR signal using our full signal injection framework (solid blue), and the measured power spectrum values with  $2\sigma$  thermal noise errors (gray shaded regions). The vertical dashed black lines signify the horizon limit for this analysis using 30 m baselines. We highlight that the bootstrapped data points and thermal noise prediction show good agreement, while the limits from the full injection framework (red and blue) are inflated due to the additional inclusion of sample variance that comes from the injection simulations.

- In A15, individual baselines were divided into 5 independent groups, where no baselines were repeated in each group. Then, baselines within each group were averaged together, and the groups were cross-multiplied to form power spectra. This grouping method was used to reduce computational time, however upon closer examination it has been found that the initial grouping introduces an element of randomness into the final measurements — more specifically, the power spectrum value fluctuates depending on how baselines are assigned into their initial groups. Our new approach removes this element of randomness at the cost of computational expense, as we now perform all baseline cross-products.
- Finally, the last change from the A15 method is that our power spectrum points (previously computed as the mean of all bootstraps), are now computed as the power spectrum estimate resulting from not bootstrapping at all. More specifically, we compute one estimate without sampling, and this estimate is propagated through our signal loss computation (this estimate is  $\hat{\mathbf{P}}_x$ ). The difference between taking the mean of the bootstrapped values and using the estimate from the no-bootstrapping case is small, but doing the lat-

ter ensures that we are forming results that reflect the estimate preferred by all our data.

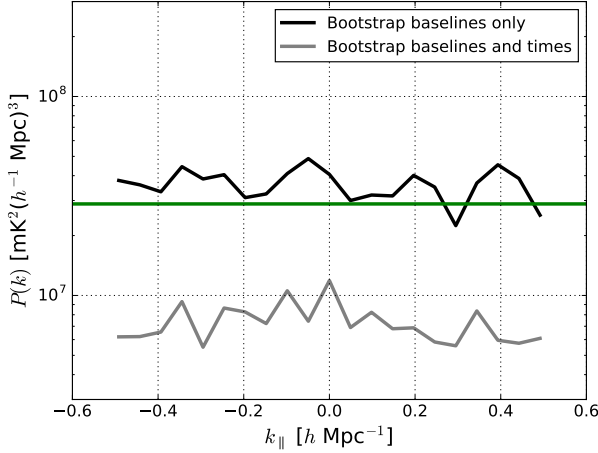
In summary, we have learned several lessons regarding bootstrapping and have revised our analysis procedure in order to determine error bars that correctly reflect the variance in our power spectrum estimates. Additionally, we have carried out substantial testing with these changes in place to verify that the computed errors for our noise simulations are accurate and agree with thermal noise predictions.

### 3.3. PAPER-64: Bias

In Section 2.3 we highlighted some common sources of bias that can show up as power spectrum detections and imitate an EoR signal. We discussed the importance of using jackknife and null tests for instilling confidence in an EoR detection, as well as for identifying other sources of biases. Here we demonstrate methods used by PAPER-64 to mitigate foreground and noise bias and we perform null tests in order to characterize the stability and implications of our results.

#### 3.3.1. Mitigating Bias

We briefly discuss one way we mitigate foreground leakage in a power spectrum estimate, and two ways we suppress noise biases. These methods are not novel



**Figure 21.**  $2\sigma$  power spectrum errors (from bootstrap variances) for a noise simulation (computed via Equation (35) using PAPER-64 observing parameters) using two different bootstrapping methods. The noise is fringe-rate filtered and a weighting matrix of  $\mathbf{I}$  (uniform-weighted) is used in order to disentangle the effects of bootstrapping from signal loss. The bootstrapping method used in A15 is shown in gray, where bootstrapping occurs along both the baseline and time axes. This under-estimates errors by sampling more values than independent ones in the dataset (fringe-rate filtering reduces the number of independent samples along time). We use the method illustrated by the black curve in our updated analysis, where bootstrapping only occurs along the baseline axis. We find that these revised limits better agree with the  $2\sigma$  analytic prediction for noise (green).

to this analysis but here we frame them in the context of minimizing false (non-EoR) detections.

Tailoring window functions is one way to suppress foreground biases (similar discussions to the following one are in Liu et al. (2014b) and A15). As alluded to in Section 2.1, we have a choice for the normalization matrix  $\mathbf{M}$  in Equation (9). For the analysis of PAPER-64 data, we compute  $\mathbf{M}$  using the matrix  $\mathbf{G}$  (which would be the Fisher matrix if  $\mathbf{R} \equiv \mathbf{C}^{-1}$ ), defined as:

$$\mathbf{G}^{\alpha\beta} = \frac{1}{2} \text{tr}[\mathbf{R}\mathbf{Q}^\alpha \mathbf{R}\mathbf{Q}^\beta] \quad (36)$$

where  $\mathbf{R}$  is the data-weighting matrix and  $\alpha$  and  $\beta$  are wavebands in  $k_\parallel$ . We take the Cholesky decomposition of  $\mathbf{G}$ , decomposing it into two lower triangular matrices (which is possible since  $\mathbf{G}$  is Hermitian):

$$\mathbf{G} = \mathbf{L}\mathbf{L}^\dagger. \quad (37)$$

Next, we construct  $\mathbf{M}$ :

$$\mathbf{M} = \mathbf{D}\mathbf{L}^{-1} \quad (38)$$

where  $\mathbf{D}$  is a diagonal matrix. In doing so, our window function, defined as  $\mathbf{W} = \mathbf{M}\mathbf{G}$  (see Equation (10)), becomes

$$\mathbf{W} = \mathbf{D}\mathbf{L}^\dagger. \quad (39)$$

Because of the nature of the lower triangular matrix, this window function has the property of preventing the

leakage of foreground power from low- $k$  to high- $k$  modes. Specifically, we order the elements in  $\mathbf{G}$  in such a way so that power can leak from high- $k$  modes to low- $k$  modes, but not vice versa. Since most foreground power shows up at low- $k$ 's, this method ensures a window function that retains clean, noise-dominated measurements while minimizing the contamination of foreground bias. This tailored window function was used in the A15 analysis, however throughout this paper, we use a diagonal  $\mathbf{M}$  for simplicity.

In addition to mitigating foreground bias at high- $k$ 's, two other sources of bias that we actively suppress in the PAPER-64 analysis are noise bias associated with the squaring of thermal noise and noise bias from crosstalk. In order to avoid the former, we filter out certain cross-multiplications when forming  $\hat{q}$  in Equation (8). Namely, the PAPER-64 dataset is divided into two halves: even Julian dates and odd Julian dates. Our data vectors are then  $\mathbf{x}_{\text{even},1}$  for the “even” dataset and baseline 1,  $\mathbf{x}_{\text{odd},1}$  for the “odd” dataset and baseline 1, etc. We only form  $\hat{q}$  when the two copies of  $\mathbf{x}$  come from different groups and baselines, never multiplying “baseline 1” with “baseline 1”, for example, in order to prevent the squaring of the same thermal noise.

To mitigate crosstalk bias, which appears as a static bias in time, we apply a fringe-rate filter that suppresses fringe-rates of zero. Figure 13 shows that the filter response is zero for such static signals. The effect of filtering out zero fringe-rates on power spectrum results is shown in A15. Most notably, even without accounting for signal loss, the crosstalk bias at all  $k$ 's is very strong compared to the removed case.

### 3.3.2. Jackknife/Null Tests

As shown in Figure 18, our illustrative PAPER-64 power spectrum shows biases above the predicted noise level, particularly at low- $k$  values. As discussed in Section 2.3.1, this bias is most likely attributable to foreground leakage.

Here we perform three null tests on PAPER-64 data that aim to isolate systematics in the data and verify that our biases are not attributable to EoR. Similar to in Section 2.3.2, we take jackknives along different axes of the dataset to produce multiple power spectra. We then difference them (i.e., the null test) to tease out excess variances.

The three results are shown in Figure 22. Each test displays the differenced power spectrum between two halves of a jackknife, where the plotted points are the differenced power spectrum values, and the plotted errors are the bootstrapped errors of the two dataset halves added in quadrature. The expected thermal noise level (gray shaded regions) is the thermal noise of each dataset added in quadrature as well. Constructing the tests as such ensures that we are probing whether the variances of each dataset differ by an amount consistent with the thermal noise. We use uniform weightings for all tests.

We take jackknives along three different axes:



- **Baselines:** We split our dataset into two halves, where each contains half of the total baselines used in the analysis. No baselines are repeated between the two datasets.
- **Sidereal Hour:** We split our dataset into two halves based on LST, namely the first half (LSTs 0.5-4.5 hours) and second half (LSTs 4.5-8.6 hours).
- **Day:** We split our dataset into even and odd Julian dates. We form power spectra for each separately, allowing the cross-multiplication of “even” with “even”, for example, for this null test only. If the same sky signal is in both the “even” and “odd” datasets, we expect it to cancel out.

In investigating Figure 22, we focus on three main possibilities — whether the data points and error bars are consistent with thermal noise (“passing”), whether the error bars are consistent with zero but not consistent with thermal noise (“passing but has an additional variance”), or whether the error bars are not consistent with zero at all (“failing”). We examine each case in the context of our results below.

Firstly, all three null tests display data points that lie within the thermal noise gray band for  $k > \pm 0.2 h \text{ Mpc}^{-1}$ . In addition, all three null tests show error bars consistent with the thermal noise level for those same  $k$ ’s. This implies that the two jackknife halves do not differ by an amount greater than the thermal noise (i.e., the baselines making up the two jackknives either do not contain bias, or contain similar amounts of bias; we suspect it is the former though more thorough jackknives along this axis are needed to make this conclusion). We deem these as “passing” null tests for the specific jackknives taken (again, dividing up the data in a different way along the same axes may not yield the same results, so more thorough testing is needed to be sure).

The second null test possibility (error bars consistent with zero but not with thermal noise) is displayed by the  $k$ -values just outside the horizon ( $k \sim \pm 0.1 h \text{ Mpc}^{-1}$ ) for all three tests. This indicates an additive noise component that is increasing our errors. More specifically, although we expect each cross-multiplication that is used in power spectrum estimation to have independent noise, there is still the possibility of a noise-foreground coupling term that can introduce power. This is because cross-multiplications produce four additive terms — a signal-squared term (where “signal” includes both foregrounds and EoR), two cross-terms between the signal and noise, and one noise-only term. When differencing two power spectra (each with their own four terms), we expect the signal-only term to subtract out for a “passing” null test, and we expect the noise-only terms to be consistent with thermal noise. While the cross-terms have a mathematical expectation value of zero, in practice we are limited by our number of samples (90 cross-products for this analysis times  $\sim 8$  independent LST samples). Combined with the fact that the

foregrounds are so bright, the finite ensemble of the couplings can introduce extra variance that varies with foreground strength. It is therefore not surprising that this effect is largest at  $k$ -values just outside the horizon, where we expect foregrounds to be brightest post-delay filtering.

Lastly, a third null test result is an error bar not consistent with zero. This is the case for the LST null test at  $k \sim -0.15 h \text{ Mpc}^{-1}$  and  $k \sim -0.2 h \text{ Mpc}^{-1}$ . In such a case, the two jackknife halves differ by an amount greater than the thermal noise (i.e., the data point is not in the thermal noise band), yet they are each constrained tightly by individual error bars that when combined, are not consistent with zero. This result implies that there exists a low level bias that is LST-dependent, and likely caused by residual foregrounds that vary in LST. Again, it is not surprising that this type of bias occurs near the horizon limit.

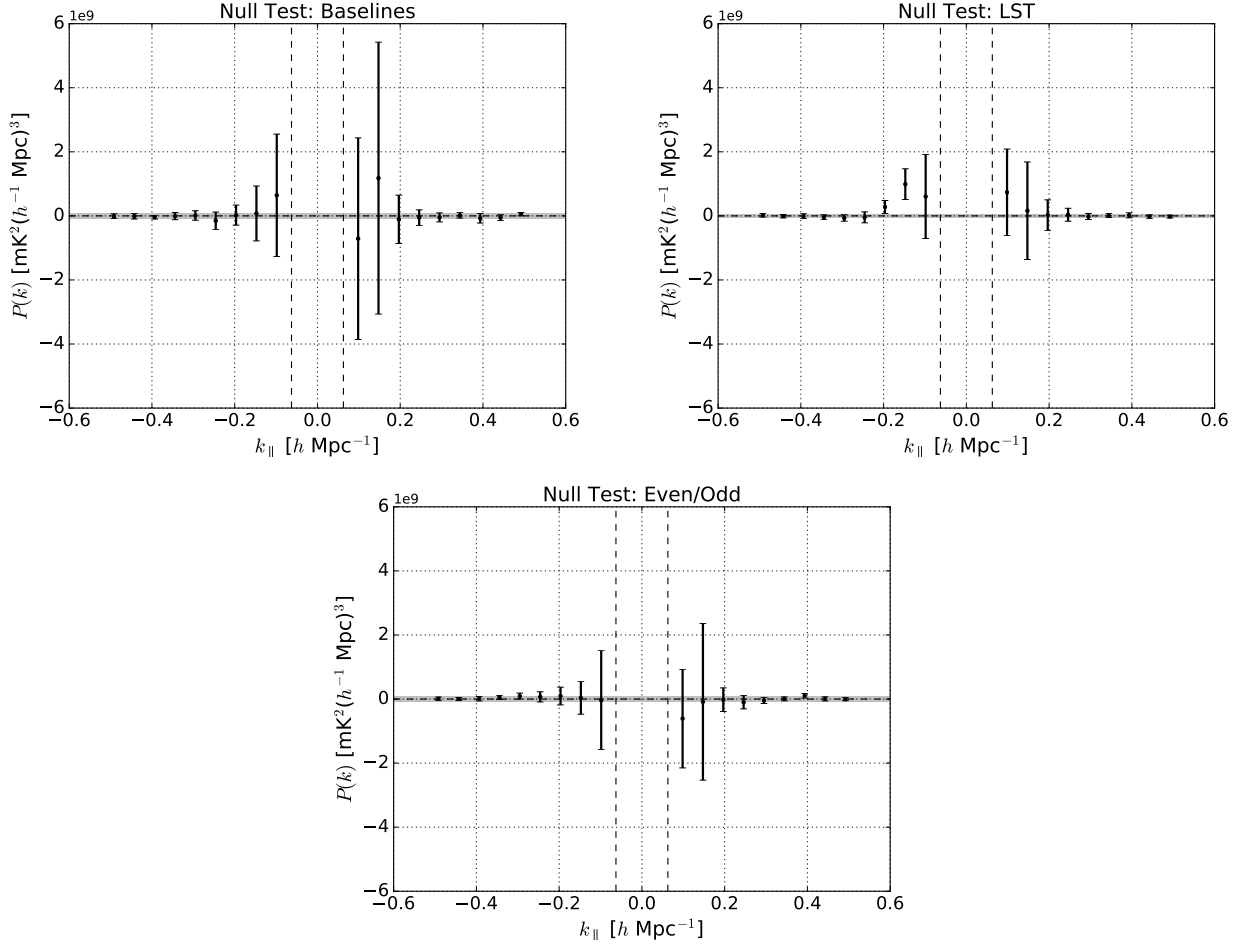
In this section we have presented the first jackknife and null tests from the PAPER experiment. Unsurprisingly, they imply that our measurements are biased by foregrounds, and not the EoR signal (a clean detection of EoR would have passed all three tests). While simple, these tests outline a framework that can be used by future measurements. The 21 cm community is beginning to recognize the importance of these types of tests (Poher et al. 2016a) in characterizing power spectra at the EoR level, and it is clear that future results will require more substantial and thorough investigations of this nature.

#### 4. CONCLUSION

Although current 21 cm published power spectrum upper limits lie several orders of magnitude above predicted EoR levels, ongoing analyses of deeper sensitivity datasets from PAPER, MWA, and LOFAR, as well as next generation instruments like HERA, are expected to continue to push towards EoR sensitivities. As the field progresses towards a detection, we have shown that it is crucial for future analyses to have a rigorous understanding of signal loss in an analysis pipeline, be able to accurately and robustly calculate both power spectrum and theoretical errors, and consistently undergo a comprehensive set of jackknife and null tests.

In particular, in this paper we have investigated the subtleties and tradeoffs of common 21 cm power spectrum techniques on signal loss, error estimation, and bias, which can be summarized as follows:

- Substantial signal loss can result when weighting data using empirically estimated covariances due to couplings with the data realizations (Section 2.1). Loss of the 21 cm signal is especially significant the fewer number of independent modes that exist in the data. Hence, there exists a trade-off between sensitivity driven time-averaging techniques such as fringe-rate filtering and signal loss when using empirically estimated covariances.



**Figure 22.** Differenced power spectrum results (with  $2\sigma$  bootstrapped errors) for three null tests, where a jackknife is taken along the baseline axis (top left), LST axis (top right), and even/odd Julian date axis (bottom). The results shown are unweighted (no signal loss), where the power spectrum values plotted are computed from the difference between two power spectra produced on either side of the jackknife axis. The gray shaded region in each plot is the estimated  $2\sigma$  theoretical noise limit given the parameters of each test. We find that there are no significant systematics for  $k > \pm 0.2 h \text{ Mpc}^{-1}$  for all three tests. However, we find that all tests exhibit an extra variance at  $k$ -values near the horizon ( $k \sim \pm 0.1 h \text{ Mpc}^{-1}$ ), likely due to foreground-noise coupling terms when foregrounds are brightest. Additionally, we find that the LST null test is not fully consistent with zero, implying a bias that is LST dependent and likely caused by varying foregrounds.

- Signal injection and recovery simulations can be used to quantify signal loss (Section 3.1). However, a signal-only simulation (i.e., comparing a uniformly weighted vs. weighted power spectrum of EoR only) can under-estimate loss by failing to account for correlations between the data and signal which can be large and negative (Section 3.1.1).
- Errors that are estimated via bootstrapping can be under-estimated if samples in the dataset are significantly correlated (Section 2.2). However, if the number of independent samples in a dataset is well-determined, bootstrapping is a simple and accurate way of estimating errors.
- Meaningful null tests are vital to validate an EoR detection (Section 2.3.2). Similarly, performing jackknife tests along multiple axes of a dataset is

necessary for confidence in an EoR detection and can also be used to tease out systematics.

As a consequence of our investigations, we have also used a subset of PAPER-64 data to make a new power spectrum analysis. This serves as an illustrative example of using a signal injection framework, correctly computing errors via bootstrapping, accurately estimating thermal noise, and implementing jackknife tests. Our revised PAPER-64 limits are presented in Kolopanis et al. (*in prep.*), which supersede all previously published PAPER limits. The main reasons for a previously under-estimated limit and ways in which our new analysis differs can be summarized by the following:

- Signal loss, previously found to be  $< 2\%$  in A15, was under-estimated by a factor of  $> 1000$  for the case of empirically estimated inverse covariance weighting. Using a regularized covariance weight-

ing method can minimize loss (Section 3.1.3), however, because a regularized weighting method is not as aggressive as the former, it produces limits that are still higher than the lossy empirical inverse covariance limits. Under-estimated signal loss therefore represents the bulk of our revision. This revision is similar in nature to the re-analysis of results from the GMRT (Paciga et al. 2013a) which were also revised from new signal loss calculations associated with their singular value decomposition foreground filter.

- Power spectrum errors, originally computed by bootstrapping, were under-estimated for the data by a factor of  $\sim 2$  in mK due to oversampling data whose effective number of independent samples was reduced from fringe-rate filtering (Section 3.2.2). Several other errors were also found regarding error estimation, though with smaller effects.
- Several factors used in an analytic expression to predict the noise-level in PAPER-64 data were revised, yielding a decrease in predicted sensitivity level by a factor of  $\sim 3$  in mK (Section 3.2.1). We note that our sensitivity prediction is revised by a factor less than our overall power spectrum result, implying that if taken at face value, the theoretical prediction for noise in A15 was too high for its data points.

The future of 21 cm cosmology is exciting, as new experiments have sensitivities that expect to reach and surpass EoR levels, improved foreground mitigation and removal strategies are being developed, and simulations

are being designed to better understand instruments. On the power spectrum analysis side, robust signal loss simulations, precise error calculations, and comprehensive jackknife tests will play critical roles in accurate 21 cm results. With strong foundations being established now, it is safe to say that we can expect to learn much about reionization and our early Universe in the coming years.

## 5. ACKNOWLEDGEMENTS

CC would like to acknowledge the UC Berkeley Chancellor's Fellowship and National Science Foundation Graduate Research Fellowship (Division of Graduate Education award 1106400). She would also like to thank Phil Bull, Bryna Hazelton, Miguel Morales, and Eric Switzer for helpful discussions. PAPER and HERA are supported by grants from the National Science Foundation (awards 1440343, and 1636646). ARP, DCJ, and JEA would also like to acknowledge NSF support (awards 1352519, 1401708, and 1455151, respectively). AL acknowledges support for this work by NASA through Hubble Fellowship grant #HST-HF2-51363.001-A awarded by the Space Telescope Science Institute, which is operated by the Association of Universities for Research in Astronomy, Inc., for NASA, under contract NAS5-26555. SAK is supported by the University of Pennsylvania School of Arts and Sciences Dissertation Completion Fellowship. JSD acknowledges NSF AAPF award 1701536. GB acknowledges support from the Royal Society and the Newton Fund under grant NA150184. This work is based on research supported in part by the National Research Foundation of South Africa (award 103424). We graciously thank SKA-SA for site infrastructure and observing support.

## APPENDIX

### A. A TOY MODEL FOR INVERSE COVARIANCE WEIGHTING

We illustrate the way in which the quadratic estimator, Equations (3), (5), and (6), estimate the power spectrum of EoR in the presence of contamination by other terms by considering two simpler, but similar, cases of estimating the variance of data. We have specifically chosen toy models where the data covariance is diagonal, as indeed we expect the EoR signal to be. We assume we have  $N$  data points  $\Delta_i$  which are the sum of a desired signal  $\sigma_i$  and an undesired contaminant  $v_i$

$$\Delta_i = \sigma_i + v_i \quad (\text{A1})$$

with

$$\langle \sigma_i \rangle = 0; \langle \sigma_i^2 \rangle = s; \text{ and } \langle \sigma \sigma^T \rangle = s \mathbf{I}_{N \times N} \equiv \mathbf{S}, \quad (\text{A2})$$

where we wish to estimate  $s$ . The contaminant in this case has a similar structure (as the EoR) for its covariance, and is assumed uncorrelated with the signal

$$\langle v_i \rangle = 0; \langle v_i^2 \rangle = u; \langle v v^T \rangle = u \mathbf{I}_{N \times N} \equiv \mathbf{U}; \text{ and } \langle \sigma_i v_j \rangle = 0. \quad (\text{A3})$$

With the covariance matrix given by  $\mathbf{C} = \mathbf{S} + \mathbf{U}$ , the estimator for  $s$  using only the quadratic part of Equation 3 is

$$\hat{s} = \frac{\Delta^T \Delta}{N} \quad (\text{A4})$$

and its expectation is

$$\langle \hat{s} \rangle = s + u. \quad (\text{A5})$$

Thus, *when the covariance structure of the contaminant is identical to the signal* ( $\frac{\partial \mathbf{S}}{\partial s} = \frac{\partial \mathbf{U}}{\partial u} = \frac{\partial \mathbf{C}}{\partial s}$ ), the information available to the quadratic portion of the estimator to distinguish between the two is degenerate, and knowledge only of  $\mathbf{C}$  and  $\frac{\partial \mathbf{C}}{\partial s}$  is inadequate. In order to obtain an unbiased estimate of  $s$ , one must also use knowledge of  $\mathbf{U}$ . Indeed, computing the linear bias from Equation (6), one finds  $b = u$ .

Now consider a case, chosen to be very similar to the toy model in 2.1.2, in which the data again have an additive contaminant, now given by

$$\Delta_i = \sigma_i + v m_i \quad (\text{A6})$$

where the properties of  $\sigma_i$  are as before, but now  $v$  is a random variable and  $m_i$  is a fixed function of  $i$  with

$$\langle v \rangle = 0; \langle v^2 \rangle = u; \langle v \mathbf{v}^T \rangle = u \mathbf{m} \mathbf{m}^T \equiv \mathbf{U}; \mathbf{m}^T \mathbf{m} = 1; \text{ and } \langle \sigma_i v \rangle = 0. \quad (\text{A7})$$

Here  $\mathbf{m}$  represents a mode which is correlated across many data points (i.e., we are assuming  $\mathbf{U}$  need *not* be diagonal), with amplitude given by  $v$ . The normalization of  $\mathbf{m}$  is a matter of convention, and can be absorbed in the variance  $u$ ; the choice above will be convenient for understanding the limiting case  $u \gg s$ .

We can calculate the quadratic portion of the estimator explicitly by using the Sherman-Morrison identity to invert the covariance matrix. Defining

$$\xi \equiv \frac{u/s}{1 + u/s}, \quad (\text{A8})$$

we have

$$\mathbf{C}^{-1} = \frac{1}{s} (\mathbf{I} - \xi \mathbf{m} \mathbf{m}^T) \quad (\text{A9})$$

and

$$\hat{s} = \frac{\Delta^T (\mathbf{I} + (\xi^2 - 2\xi) \mathbf{m} \mathbf{m}^T) \Delta}{N + \xi^2 - 2\xi} \quad (\text{A10})$$

with expectation

$$\langle \hat{s} \rangle = s + \frac{1 - 2\xi + \xi^2}{N + -2\xi + \xi^2} u. \quad (\text{A11})$$

It is worth observing immediately that there is no multiplicative bias on  $s$ , and that the additive bias is strictly  $< u/N$ .

An instructive limit is  $u \gg s$ ,  $\xi \rightarrow 1$ , in which case the virtue of weighting by  $\mathbf{C}^{-1}$  becomes clearer, as it becomes

$$\mathbf{C}^{-1} = \frac{1}{s} (\mathbf{I} - \mathbf{m} \mathbf{m}^T) \quad (\text{A12})$$

where  $\mathbf{I} - \mathbf{m} \mathbf{m}^T$  is the projection operator, projecting out  $\mathbf{m}$  from any vector it acts on, and further, the linear bias tends to 0 as  $\xi \rightarrow 1$  (i.e., the projection is “perfect” and not “undone” by the Fisher matrix normalization).

This is the ideal case for the inverse covariance weighting performed in the PAPER analysis, where removal of contamination with a known covariance can be suppressed by a kind of projection of the offending modes. But even in this case, it is worth pointing out that the estimator still has a linear bias for finite  $u$ . We have also assumed that the contaminating mode  $\mathbf{m}$  is known perfectly; the next appendix takes up the case where the modes are estimated from the data.

## B. A TOY MODEL FOR SIGNAL LOSS

In this Appendix, we examine a toy model for signal loss. Our goal is to derive an analytic formula for power spectrum signal loss. While this model does not apply generally to all the scenarios presented in this paper, it provides some analytic intuition for how the coupling between data and an empirical covariance can result in signal loss.

The minimum-variance quadratic estimator  $\hat{P}^\alpha$  for the  $\alpha$ th bandpower of the power spectrum is given by

$$\hat{P}^\alpha = \frac{1}{2\mathbf{F}^{\alpha\alpha}} \mathbf{x}^t \mathbf{C}^{-1} \mathbf{Q}^\alpha \mathbf{C}^{-1} \mathbf{x}, \quad (\text{B13})$$

where

$$F^{\alpha\alpha} \equiv \frac{1}{2} \text{tr} (\mathbf{C}^{-1} \mathbf{Q}^\alpha \mathbf{C}^{-1} \mathbf{Q}^\alpha) \quad (\text{B14})$$

is the  $\alpha$ th diagonal element of the Fisher matrix. For this section only, with no loss of generality, we assume that the data  $\mathbf{x}$  are real. We also assume for simplicity that  $\mathbf{x}$  is the data from a single instant in time, so that it is of length  $N_f$ , where  $N_f$  is the number of frequency channels.

In our case, we do not have *a priori* knowledge of the covariance matrix. Thus, we deviate from the true minimum-variance quadratic estimator and replace  $\mathbf{C}$  with  $\hat{\mathbf{C}}$ , its data-derived approximation. Our estimator then becomes

$$\hat{P}_{\text{loss}}^\alpha = \frac{1}{2\hat{\mathbf{F}}^{\alpha\alpha}} \mathbf{x}^t \hat{\mathbf{C}}^{-1} \mathbf{Q}^\alpha \hat{\mathbf{C}}^{-1} \mathbf{x}, \quad (\text{B15})$$



where

$$\hat{F}^{\alpha\alpha} \equiv \frac{1}{2} \text{tr} \left( \hat{\mathbf{C}}^{-1} \mathbf{Q}^\alpha \hat{\mathbf{C}}^{-1} \mathbf{Q}^\alpha \right), \quad (\text{B16})$$

with the label “loss” to foreshadow the fact that this will be an estimator with signal loss (i.e., a multiplicative bias of less than unity). We will now provide an explicit demonstration of this by modeling the estimated covariance as

$$\hat{\mathbf{C}} = (1 - \eta) \mathbf{C} + \eta \mathbf{x} \mathbf{x}^t, \quad (\text{B17})$$

where  $\eta$  is a parameter quantifying our success at estimating the true covariance matrix. If  $\eta = 0$ , our covariance estimate has perfectly modeled the true covariance and  $\hat{\mathbf{C}} = \mathbf{C}$ . On the other hand, if  $\eta = 1$ , then our covariance estimate is based purely on the one realization of the covariance that is our actual data, and we would expect a high level of overfitting and signal loss.

Our strategy for computing the signal loss will be to insert Equation (B17) into Equation (B15) and to express the resulting estimator  $\hat{P}_{\text{loss}}^\alpha$  in terms of  $\hat{P}^\alpha$ . We begin by expressing  $\hat{\mathbf{C}}^{-1}$  in terms of  $\mathbf{C}^{-1}$  using the Woodbury identity so that

$$\hat{\mathbf{C}}^{-1} = \frac{\mathbf{C}^{-1}}{1 - \eta} \left[ \mathbf{I} - \frac{\eta \mathbf{x} \mathbf{x}^t \mathbf{C}^{-1}}{1 + \eta(g - 1)} \right], \quad (\text{B18})$$

where we have defined  $g \equiv \mathbf{x}^t \mathbf{C}^{-1} \mathbf{x}$ . Inserting this into our Fisher estimate we have

$$\hat{F}^{\alpha\alpha} = \frac{F^{\alpha\alpha}}{(1 - \eta)^2} \left[ 1 - \frac{\eta}{1 + \eta(g - 1)} \frac{h^{\alpha\alpha}}{F^{\alpha\alpha}} + \frac{1}{2} \left( \frac{\eta}{1 + \eta(g - 1)} \right)^2 \frac{(h^\alpha)^2}{F^{\alpha\alpha}} \right], \quad (\text{B19})$$

where  $h^\alpha \equiv \mathbf{x}^t \mathbf{C}^{-1} \mathbf{Q}^\alpha \mathbf{C}^{-1} \mathbf{x}$  and  $h^{\alpha\alpha} \equiv \mathbf{x}^t \mathbf{C}^{-1} \mathbf{Q}^\alpha \mathbf{C}^{-1} \mathbf{Q}^\alpha \mathbf{C}^{-1} \mathbf{x}$ . Note that  $g$ ,  $h^\alpha$ , and  $h^{\alpha\alpha}$  are all random variables, since they depend on  $\mathbf{x}$ . Inserting these expressions into our estimator gives

$$\hat{P}_{\text{loss}}^\alpha = \frac{1}{2} \frac{h^\alpha}{F^{\alpha\alpha}} \left[ 1 - \frac{\eta g}{1 + \eta(g - 1)} \right]^2 \left[ 1 - \frac{\eta}{1 + \eta(g - 1)} \frac{h^{\alpha\alpha}}{F^{\alpha\alpha}} + \frac{1}{2} \left( \frac{\eta}{1 + \eta(g - 1)} \right)^2 \frac{(h^\alpha)^2}{F^{\alpha\alpha}} \right]^{-1}. \quad (\text{B20})$$

Both for the purposes of analytical tractability and to provide intuition, we expand this expression to leading order in  $\eta$ . This approximates the limiting case where the covariance  $\hat{\mathbf{C}}$  is close to the ideal and the lossy covariance is a small perturbation. The result is

$$\hat{P}_{\text{loss}}^\alpha \approx \frac{1}{2} \frac{h^\alpha}{F^{\alpha\alpha}} \left[ 1 - \eta \left( g - \frac{h^{\alpha\alpha}}{F^{\alpha\alpha}} \right) \right]. \quad (\text{B21})$$

Taking the ensemble average of both sides and noting that the true power spectrum  $P^\alpha$  is equal to  $\langle h^\alpha \rangle / 2F^{\alpha\alpha}$ , we obtain

$$\langle \hat{P}_{\text{loss}}^\alpha \rangle \approx (1 - \eta N_f) P^\alpha + 4\eta \frac{\text{tr}(\mathbf{C}^{-1} \mathbf{Q}^\alpha \mathbf{C}^{-1} \mathbf{Q}^\alpha \mathbf{C}^{-1} \mathbf{Q}^\alpha)}{[\text{tr}(\mathbf{C}^{-1} \mathbf{Q}^\alpha \mathbf{C}^{-1} \mathbf{Q}^\alpha)]^2} \approx (1 - \eta N_f) P^\alpha, \quad (\text{B22})$$

where recall that  $N_f$  is the length of  $\mathbf{x}$ , or the number of frequency channels. In the last step we dropped the final term, since it scales as  $\eta P^\alpha$  (without the factor of  $N$ ) and is therefore typically small compared to the terms that have been retained.

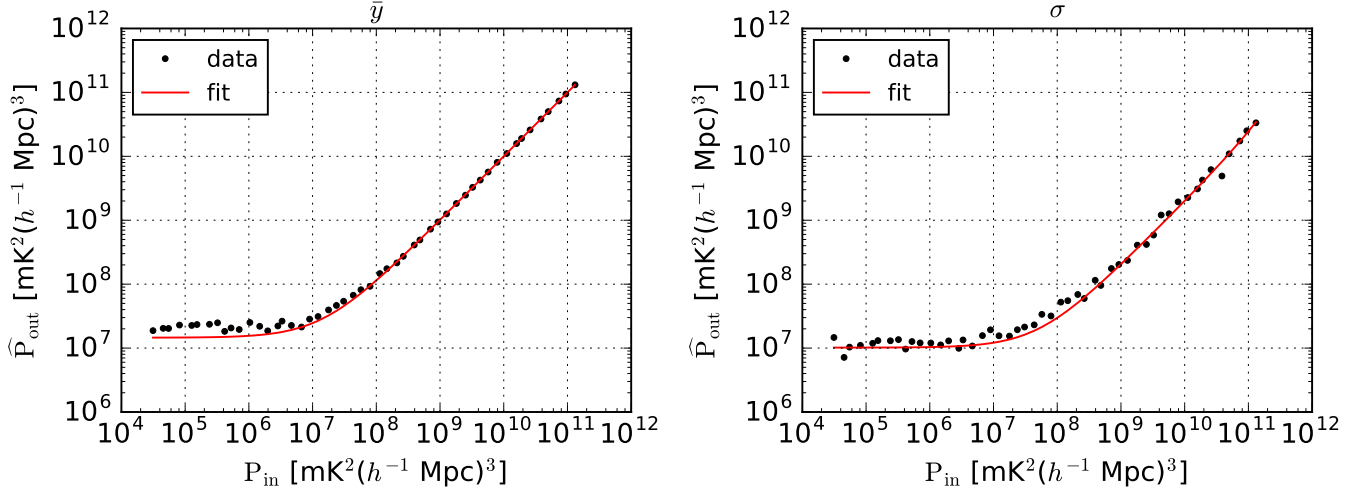
Recalling that  $P^\alpha$  is the *true* power spectrum, one sees that when the covariance in the optimal quadratic estimator is naively replaced by an empirical covariance, the resulting power spectrum estimate is biased low, i.e., there is signal loss. This occurs because of couplings between  $\hat{\mathbf{C}}$  and  $\mathbf{x}$ , which formally means that what was originally a quadratic estimator is no longer quadratic, but contains higher-order correlations. This violates the assumptions implicit in the derivation of  $F^{\alpha\alpha}$  as the normalization factor for converting unnormalized bandpowers  $\frac{1}{2} \mathbf{x}^t \mathbf{C}^{-1} \mathbf{Q}^\alpha \mathbf{C}^{-1} \mathbf{x}$  into properly normalized power spectrum estimates, where the unnormalized bandpowers are assumed to be two-point (i.e., quadratic) statistics (Liu & Tegmark 2011). The result is an improperly normalized—and thus lossy—power spectrum estimate.

### C. DERIVATION OF THE JEFFREYS PRIOR

The Jeffreys prior is an objective, non-informative prior distribution for a parameter space using Bayesian probability (Jaynes 1968). For the signal injection framework outlined in Section 3.1.2, we wish to compute the prior  $p(P_{\text{in}})$ , or the probability density of the power spectrum of the EoR signal.

The Jeffreys prior is defined as:

$$p(P_{\text{in}}) \propto \sqrt{\left\langle \left( \frac{\partial \mathcal{L}}{\partial P_{\text{in}}} \right)^2 \right\rangle}, \quad (\text{C23})$$



**Figure C1.** An illustrative example (for the PAPER-64 analysis using uniform weighting and  $k = 0.393 h \text{ Mpc}^{-1}$ ) of how the mean of  $P_{\text{out}}$  (left) and standard deviation of  $P_{\text{out}}$  (right) behave as a function of  $P_{\text{in}}$ . Polynomials are fit to each (red) to describe how  $\bar{y}$  and  $\sigma$  evolve with  $x$  (injection level), respectively, for the computation of the Jeffreys prior as defined in Equation (C27). The polynomial fits for this example are  $y = (-5.1 \times 10^{-15})x^2 + x + (1.5 \times 10^7)$  and  $y = (5.0 \times 10^{-13})x^2 + 0.2x + 10^7$  for  $\bar{y}$  and  $\sigma$ , respectively.

where

$$\mathcal{L} = \ln p(\hat{P}_{\text{out}} | P_{\text{in}}), \quad (\text{C24})$$

recalling that in our framework  $P_{\text{in}}$  is the power spectrum of the EoR signal (uniformly weighted), and  $\hat{P}_{\text{out}}$  is the weighted output power spectrum of the data plus EoR.

For a single injection amplitude, our bootstrapped  $\hat{P}_{\text{out}}$  values are well-approximated by a Gaussian distribution. Simplifying our notation so that  $x = P_{\text{in}}$  and  $y = \hat{P}_{\text{out}}$ :

$$p(y|x) = \frac{1}{\sigma(x)\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{y-\bar{y}(x)}{\sigma}\right)^2}, \quad (\text{C25})$$

where  $\sigma$  is the standard deviation of  $\hat{P}_{\text{out}}$  and  $\bar{y}$  is the mean of  $\hat{P}_{\text{out}}$ , and they are both functions of  $P_{\text{in}}$ . Using Equations (C25) and (C24), the quantity inside the expectation value of Equation (C23) becomes:

$$\begin{aligned} \left(\frac{\partial \mathcal{L}}{\partial x}\right)^2 &= \frac{1}{\sigma^2} \left(\frac{\partial \sigma}{\partial x}\right)^2 - \left(\frac{2(y-\bar{y})}{\sigma^3}\right) \frac{\partial \sigma}{\partial x} \frac{\partial \bar{y}}{\partial x} - \left(\frac{2(y-\bar{y})^2}{\sigma^4}\right) \left(\frac{\partial \sigma}{\partial x}\right)^2 \\ &\quad + \left(\frac{(y-\bar{y})^2}{\sigma^4}\right) \left(\frac{\partial \bar{y}}{\partial x}\right)^2 + \left(\frac{2(y-\bar{y})^3}{\sigma^5}\right) \frac{\partial \sigma}{\partial x} \frac{\partial \bar{y}}{\partial x} + \left(\frac{(y-\bar{y})^4}{\sigma^6}\right) \left(\frac{\partial \sigma}{\partial x}\right)^2. \end{aligned} \quad (\text{C26})$$

Taking the expectation value then removes all terms with odd powers of  $(y-\bar{y})$  because those Gaussian moments evaluate to zero. Additionally, the second moment can be simplified since  $\langle (y-\bar{y})^2 \rangle = \sigma^2$  and the fourth moment can be simplified since  $\langle (y-\bar{y})^4 \rangle = 3\sigma^4$ . Finally, after some additional simplification the Jeffreys prior becomes:

$$p(x) \propto \sqrt{\frac{1}{\sigma^2} \left( 2 \left( \frac{\partial \sigma}{\partial x} \right)^2 + \left( \frac{\partial \bar{y}}{\partial x} \right)^2 \right)}. \quad (\text{C27})$$

When we simulate our full injection framework as in Section 3.1.2, we sample 50  $P_{\text{in}}$  values that range from  $\sim 10^5 \text{ mK}^2 (h^{-1} \text{ Mpc})^3$  to  $\sim 10^{11} \text{ mK}^2 (h^{-1} \text{ Mpc})^3$ , and we note that the prior is set to zero outside those regions. For the injections that we do sample, we can simply fit analytic functions to the mean and standard deviations of  $\hat{P}_{\text{out}}$  ( $\bar{y}$  and  $\sigma$ ) as functions of  $P_{\text{in}}$ . An example of the typical shape of these functions for the PAPER-64 analysis is shown in Figure C1, though in practice we fit solutions for every  $k$ -value and simulation independently.

We also show the typical shape of the Jeffreys prior used in our analysis in Figure C2, as computed by Equation (C27). Most noticeably, it is not constant with  $P_{\text{in}}$ , meaning a uniform prior, which is often used for simplicity, is informative in our application. Therefore, due to its objective nature we choose to use a Jeffreys prior in our analysis, multiplying our likelihood functions by Equation (C27) before computing posterior distributions.



**Figure C2.** An example of the typical Jeffreys prior shape for the PAPER-64 analysis as computed by Equation (C27) (black). We smooth the prior using a sliding boxcar average over every 5 injection levels (red). Most noticeably, the Jeffreys prior is not constant with  $P_{\text{in}}$ , meaning a uniform prior would be an informative prior.

## REFERENCES

- Ade, P., et al. 2008, *Astrophysical Journal*, 674, 22
- Ade, P. A. R., et al. 2017, *PhRvD*, 96, 102003
- Ali, S. S., Bharadwaj, S., & Chengalur, J. N. 2008, *MNRAS*, 385, 2166
- Ali, Z. S., et al. 2015, *ApJ*, 809, 61
- Andrae, R. 2010, *ArXiv e-prints*
- Araujo, D., et al. 2012, *The Astrophysical Journal*, 760, 145
- Barkana, R., & Loeb, A. 2001, *PhR*, 349, 125
- . 2008, *Monthly Notices of the Royal Astronomical Society*, 384, 1069
- Beardsley, A. P., et al. 2016, *The Astrophysical Journal*, 833, 102
- Bernardi, G., et al. 2009, *A&A*, 500, 965
- . 2010, *A&A*, 522, A67+
- Bernardi, G., et al. 2013, *The Astrophysical Journal*, 771, 105
- Bernardi, G., et al. 2016, *MNRAS*, 461, 2847
- BICEP2 Collaboration et al. 2016, *ApJ*, 833, 228
- Bischoff, C., et al. 2011, *The Astrophysical Journal*, 741, 111
- Bond, J. R., Jaffe, A. H., & Knox, L. 1998, *PhRvD*, 57, 2117
- Bowman, J. D., & Rogers, A. E. E. 2010, *Nature*, 468, 796
- Bowman, J. D., Rogers, A. E. E., Monsalve, R. A., Mozdzen, T. J., & Mahesh, N. 2018, *Nature*, 555, 67
- Burns, J. O., et al. 2012, *Advances in Space Research*, 49, 433
- Chang, T.-C., Pen, U.-L., Bandura, K., & Peterson, J. B. 2010, *Nature*, 466, 463
- Chapman, E., et al. 2012, *Monthly Notices of the Royal Astronomical Society*, 423, 2518
- Chiang, H. C., et al. 2010, *The Astrophysical Journal*, 711, 1123
- Crites, A. T., et al. 2015, *The Astrophysical Journal*, 805, 36
- Das, S., et al. 2011a, *Physical Review Letters*, 107, 021301
- . 2011b, *ApJ*, 729, 62
- Datta, A., Bowman, J. D., & Carilli, C. L. 2010, *The Astrophysical Journal*, 724, 526
- de Oliveira-Costa, A., Tegmark, M., Gaensler, B. M., Jonas, J., Landecker, T. L., & Reich, P. 2008, *MNRAS*, 388, 247
- DeBoer, D. R., et al. 2017, *Publications of the Astronomical Society of the Pacific*, 129, 045001
- Dillon, J. S., Liu, A., & Tegmark, M. 2013, *PhRvD*, 87, 043005
- Dillon, J. S., & Parsons, A. R. 2016, *The Astrophysical Journal*, 826, 181
- Dillon, J. S., et al. 2014, *Phys. Rev. D*, 89, 023002
- . 2015, *Phys. Rev. D*, 91, 123011
- Dodelson, S., & Schneider, M. D. 2013, *PhRvD*, 88, 063537
- Efron, B., & Tibshirani, R. 1994, *An Introduction to the Bootstrap*, Chapman & Hall/CRC Monographs on Statistics & Applied Probability (Taylor & Francis)
- Ewall-Wice, A., Dillon, J. S., Liu, A., & Hewitt, J. 2017, *MNRAS*, 470, 1849
- Furlanetto, S. R., Oh, S. P., & Briggs, F. H. 2006, *PhR*, 433, 181
- Ghosh, A., Bharadwaj, S., Ali, S. S., & Chengalur, J. N. 2011, *MNRAS*, 418, 2584
- Hartlap, J., Simon, P., & Schneider, P. 2007, *A&A*, 464, 399
- Jacobs, D. C., et al. 2015, *ApJ*, 801, 51
- Jacobs, D. C., et al. 2016, *ApJ*, 825, 114
- Jaynes, E. 1968, *IEEE Transactions on Systems Science and Cybernetics*, 4, 227
- Jelić, V., et al. 2008, *MNRAS*, 389, 1319
- Joachimi, B. 2017, *MNRAS*, 466, L83
- Keating, G. K., Marrone, D. P., Bower, G. C., Leitch, E., Carlstrom, J. E., & DeBoer, D. R. 2016, *The Astrophysical Journal*, 830, 34
- Kerrigan, J., et al. 2018, *ArXiv e-prints*
- Kohn, S. A., et al. 2016, *ApJ*, 823, 88
- Koopmans, L., et al. 2015, *Advancing Astrophysics with the Square Kilometre Array (AASKA14)*, 1
- Liu, A., & Parsons, A. R. 2016, *Monthly Notices of the Royal Astronomical Society*, 457, 1864
- Liu, A., Parsons, A. R., & Trott, C. M. 2014a, *PhRvD*, 90, 023018
- . 2014b, *PhRvD*, 90, 023019
- Liu, A., & Tegmark, M. 2011, *Phys. Rev. D*, 83, 103006
- Loeb, A., & Furlanetto, S. 2013, *The First Galaxies in the Universe* (Princeton University Press)
- Masui, K. W., et al. 2013, *ApJL*, 763, L20
- Moore, D. F., Aguirre, J. E., Parsons, A. R., Jacobs, D. C., & Poher, J. C. 2013, *The Astrophysical Journal*, 769, 154
- Morales, M. F., & Wyithe, J. S. B. 2010, *ARA&A*, 48, 127
- Paciga, G., et al. 2013a, *MNRAS*
- . 2013b, *MNRAS*, 433, 639
- Padmanabhan, N., White, M., Zhou, H. H., & O’Connell, R. 2016, *MNRAS*, 460, 1567
- Parsons, A., Poher, J., McQuinn, M., Jacobs, D., & Aguirre, J. 2012a, *ApJ*, 753, 81

- Parsons, A. R., Liu, A., Ali, Z. S., & Cheng, C. 2016, *ApJ*, 820, 51
- Parsons, A. R., Pober, J. C., Aguirre, J. E., Carilli, C. L., Jacobs, D. C., & Moore, D. F. 2012b, *ApJ*, 756, 165
- Parsons, A. R., et al. 2010, *AJ*, 139, 1468
- . 2014, *ApJ*, 788, 106
- Patil, A. H., et al. 2016, *MNRAS*, 463, 4317
- . 2017, *ApJ*, 838, 65
- Patra, N., Subrahmanyam, R., Sethi, S., Udaya Shankar, N., & Raghunathan, A. 2015, *ApJ*, 801, 138
- Paz, D. J., & Sánchez, A. G. 2015, *MNRAS*, 454, 4326
- Pearson, D. W., & Samushia, L. 2016, *MNRAS*, 457, 993
- Peterson, U.-L. P. X.-P. W. J. 2004, *ArXiv Astrophysics e-prints*
- Petrovic, N., & Oh, S. P. 2011, *MNRAS*, 413, 2103
- Pober, J. C., Greig, B., & Mesinger, A. 2016a, *MNRAS*, 463, L56
- Pober, J. C., et al. 2012, *AJ*, 143, 53
- Pober, J. C., et al. 2013, *The Astrophysical Journal Letters*, 768, L36
- Pober, J. C., et al. 2013, *AJ*, 145, 65
- . 2014, *ApJ*, 782, 66
- . 2016b, *ApJ*, 819, 8
- Pope, A. C., & Szapudi, I. 2008, *MNRAS*, 389, 766
- Pritchard, J. R., & Loeb, A. 2010, *PhRvD*, 82, 023006
- . 2012, *Reports on Progress in Physics*, 75, 086901
- Quenouille, M. H. 1949, *Ann. Math. Statist.*, 20, 355
- Santos, M. G., Cooray, A., & Knox, L. 2005, *ApJ*, 625, 575
- Sellentin, E., & Heavens, A. F. 2016, *MNRAS*, 456, L132
- Sherwin, B. D., et al. 2017, *Phys. Rev. D*, 95, 123529
- Sokolowski, M., et al. 2015, *PASA*, 32, e004
- Switzer, E. R., Chang, T.-C., Masui, K. W., Pen, U.-L., & Voytek, T. C. 2015, *ApJ*, 815, 51
- Switzer, E. R., et al. 2013, *MNRAS*, 434, L46
- Taylor, A., & Joachimi, B. 2014, *MNRAS*, 442, 2728
- Tegmark, M. 1997, *PhRvD*, 55, 5895
- Thompson, A. R., Moran, J. M., & Swenson, Jr., G. W. 2001, *Interferometry and Synthesis in Radio Astronomy*, 2nd Edition
- Thyagarajan, N., et al. 2013, *ApJ*, 776, 6
- Tingay, S. J., et al. 2013, *PASA*, 30, 7
- Trott, C. M., Wayth, R. B., & Tingay, S. J. 2012, *ApJ*, 757, 101
- Trott, C. M., et al. 2016, *The Astrophysical Journal*, 818, 139
- Tukey. 1958, *Ann. Math. Statist.*, 29, 614
- van Haarlem, M. P., et al. 2013, *A&A*, 556, A2
- Vedantham, H., Shankar, N. U., & Subrahmanyam, R. 2012, *The Astrophysical Journal*, 745, 176
- Voytek, T. C., Natarajan, A., Jáuregui García, J. M., Peterson, J. B., & López-Cruz, O. 2014, *ApJL*, 782, L9
- Wang, J., et al. 2013, *The Astrophysical Journal*, 763, 90
- Wolz, L., Abdalla, F. B., Blake, C., Shaw, J. R., Chapman, E., & Rawlings, S. 2014, *Monthly Notices of the Royal Astronomical Society*, 441, 3271
- Wu, X. 2009, in *Bulletin of the American Astronomical Society*, Vol. 41, American Astronomical Society Meeting Abstracts #213, 474