

## BAYESIAN INFERENCE OF CMB GRAVITATIONAL LENSING

ETHAN ANDERES<sup>1</sup>

Department of Statistics, University of California, Davis CA 95616, USA.

BENJAMIN D. WANDELT<sup>2</sup> AND GUILHEM LAVAUX<sup>2</sup>

Sorbonne Universités, UPMC Univ Paris 06 & CNRS, UMR7095, Institut d’Astrophysique de Paris, F-75014, Paris, France

*Draft version October 21, 2015*

### ABSTRACT

The Planck satellite, along with several ground based telescopes, have mapped the cosmic microwave background (CMB) at sufficient resolution and signal-to-noise so as to allow a detection of the subtle distortions due to the gravitational influence of the intervening matter distribution. A natural modeling approach is to write a Bayesian hierarchical model for the lensed CMB in terms of the unlensed CMB and the lensing potential. So far there has been no feasible algorithm for inferring the posterior distribution of the lensing potential from the lensed CMB map. We propose a solution that allows efficient Markov Chain Monte Carlo sampling from the joint posterior of the lensing potential and the unlensed CMB map using the Hamiltonian Monte Carlo technique. The main conceptual step in the solution is a re-parameterization of CMB lensing in terms of the lensed CMB and the “inverse lensing” potential. We demonstrate a fast implementation on simulated data including noise and a sky cut, that uses a further acceleration based on a very mild approximation of the inverse lensing potential. We find that the resulting Markov Chain has short correlation lengths and excellent convergence properties, making it promising for application to high resolution CMB data sets of the future.

*Keywords:* gravitational lensing: weak – methods: statistical – cosmology: cosmic background radiation – cosmology: dark matter

### 1. INTRODUCTION

Over the past few years, data from ground based telescopes (ACT, SPT, Polarbear) and the Planck satellite have resulted in an unprecedented detection of weak gravitational lensing of the cosmic microwave background (CMB) (Das et al. 2011; Engelen et al. 2012; Planck Collaboration 2014; The Polarbear Collaboration: P. A. R. Ade et al. 2014; Planck Collaboration 2015). Upcoming high resolution, high signal-to-noise experiments are poised to make the gravitational lensing distortion a powerful probe of cosmology, dark matter, and neutrino physics. The state-of-the-art estimator of CMB gravitational lensing, the quadratic estimator developed by Hu and Okamoto (Hu 2001; Hu and Okamoto 2002), works in part through a delicate cancellation of terms in an infinite Taylor expansion of the lensing effect on the CMB. The effect of this cancellation is particularly sensitive to foreground contaminants and sky masking, which if not fully accounted for, limits the statistical inferential power of this new data.

Possibly the most promising alternative to the quadratic estimator is Bayesian lensing. It has been known for some time that the quadratic estimator is suboptimal for high signal-to-noise, high resolution experiments and that a full Bayesian treatment can overcome this limitation (Hirata and Seljak 2003a,b). Indeed, Bayesian techniques applied to the lensed CMB observations have the potential to drastically change the way lensing is estimated and used for inference. Current frequentist estimators of the unknown lensing potential treat the unlensed CMB as a source of shape noise which is marginalized out. Conversely, a Bayesian lensing posterior treats the lensing potential *and* the unlensed CMB as joint unknowns, whereby obtaining scientific constrains jointly rather than marginally. Moreover, the posterior distribution is easier to interpret and sequentially update with additional data. From the geometry of weak lensing, most of the lensing power comes from matter at high redshift  $z \sim 2$ . At these distances the matter distribution on large scales is well approximated by Gaussian density fluctuations. In addition, the unlensed CMB is, at present, indistinguishable from an isotropic Gaussian random field. From a statistical perspective, this is a perfect scenario for Bayesian methods in that both the observations and the unknown lensing potential are *physically predicted* to be Gaussian random fields.

Physicists have known, for some time, that Bayesian methods could potentially provide next-generation lensing estimates. In their seminal review Lewis and Challinor (2006) discuss the possibility of obtaining posterior draws from the lensing potential and the unlensed CMB jointly. However, they acknowledge the main obstacle for naive Gibbs implementations:

“... given a particular lensing potential the delensed sky is given essentially by a delta function. This means that naive Gibbs iterations will not converge within a reasonable time. At the time of writing there are no known practical methods for sampling from the full posterior distribution.”

In this paper we show that, indeed, there does exist a practical way to obtain Gibbs iterations which converge quickly. The solution is through a re-parameterization of CMB lensing problem. Instead of treating the lensing potential as

unknown we work with inverse-lensing or what we call anti-lensing. Surprisingly, the slowness of naive Gibbs translates to fast convergence of the re-parameterized Gibbs chain.

In Section 3 we motivate our re-parameterization by analyzing a simple two parameter statistical problem. The concepts are then applied to the Bayesian lensing problem in Section 4. The two conditional distributions in our Gibbs implementation are discussed in Section 5 and Section 6. We finish with some simulation examples in Section 7.

All the code presented in this paper is written in the language *Julia* (Bezanson et al. 2012) and is publicly available through the on-line repository <https://github.com/EthanAnderes/BayesianCmbLensing>.

## 2. WEAK LENSING PRIMER AND A BAYESIAN CHALLENGE

The effect of weak lensing is to simply remap the CMB, preserving surface brightness. Up to leading order, the remapping displacements are given by  $\nabla\phi$ , where  $\phi$  denotes the lensing potential and is the planar projection of the three dimensional gravitational potential (see Dodelson (2003), for example). Therefore the lensed CMB can be written  $T(x + \nabla\phi(x))$  where  $T(x)$  denotes the unlensed CMB temperature fluctuations and  $x$  represents an observational direction on the unit sphere. For this paper we will be focusing on the small angle limit so that  $x$  is assumed to vary in a small patch of  $\mathbb{R}^2$ . However, we do not expect the fast convergence properties of our algorithm to be sensitive to the small angle approximation and the methodology presented here should hold for a full treatment on the sphere. The lensed CMB is observed with additive noise (denoted  $n(x)$ ) to result in data of the form

$$\text{data}(x) = T(x + \nabla\phi(x)) + n(x). \quad (1)$$

The goal of weak lensing surveys is to use the data in (1) to estimate  $\phi$ ,  $T$  and possibly the spectral densities of  $T$  and  $\phi$ .

A natural approach to develop a Bayesian lensing estimator is to generate posterior samples through a Gibbs algorithm which iteratively samples from the two conditionals:  $P(T|\phi, \text{data})$  and  $P(\phi|T, \text{data})$ . Sampling from  $P(T|\phi, \text{data})$  is simply a Gaussian random field prediction problem since conditioning on  $\phi$  models the data as

$$\text{data}(x) = T(\underbrace{x + \nabla\phi(x)}_{\text{known obs locations}}) + n(x).$$

In other words, the data is a noisy version of  $T$  observed on an irregular grid. Conversely, when sampling from  $P(\phi|T, \text{data})$  the data is of the form

$$\text{data}(x) = \underbrace{T}_{\text{known}}(x + \nabla\phi(x)) + n(x).$$

To see how one might approximate this conditional notice first that the CMB field  $T(x)$  is very smooth. Indeed, Silk damping predicts a exponentially decaying power spectrum. Therefore a linear Taylor approximation,  $\text{data}(x) \approx T(x) + \nabla T(x) \cdot \nabla\phi(x) + n(x)$ , may be useful. In fact, the derivation of the quadratic estimator explicitly uses this linear approximation. If one is willing to use this linear approximation then the conditional  $P(\phi|T, \text{data})$  is simply a Bayesian regression problem since  $T$  (and thus  $\nabla T$ ) are both known with a Gaussian prior on  $\nabla\phi$ .

Unfortunately, the structure of both of these conditionals make the Gibbs very slow to converge. The case is exacerbated in the situation when noise level is small. For example, in the second conditional, if  $T$  is known and fixed, the extent of the likely  $\phi$ 's under  $P(\phi|T, \text{data})$  is very small compared to the likely  $\phi$ 's under  $P(\phi, T|\text{data})$ . This suggests a highly dependent posterior  $P(\phi, T|\text{data})$ .

## 3. TWO PARAMETER ANALOGY

To motivate our solution to the Bayesian lensing problem we start with a simple two parameter statistical problem. This system has two unknown parameters  $t, \varphi$  with a single data point given by

$$\text{data} = t + \varphi + n$$

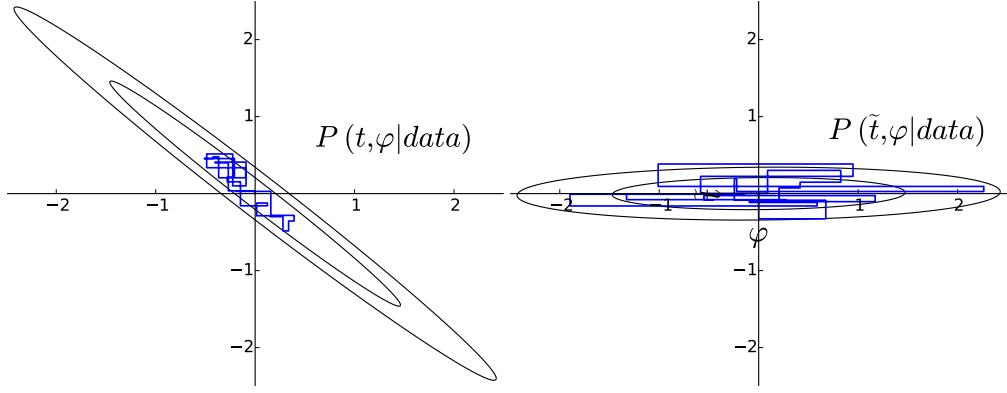
where  $n$  denotes additive noise. In the Bayesian setting, the posterior distribution is computed as

$$P(t, \varphi|\text{data}) \propto P(\text{data}|t, \varphi)P(t, \varphi) \quad (2)$$

where  $P(\text{data}|t, \varphi)$  denotes the likelihood of the data given  $t, \varphi$  and  $P(t, \varphi)$  denotes the prior on  $t, \varphi$ . The Gibbs sampler is a widely used algorithm for generating (asymptotic) samples from  $P(t, \varphi|\text{data})$ . The algorithm generates a Markov chain of parameter values  $(t^1, \varphi^1), (t^2, \varphi^2), \dots$  generated by iteratively sampling from the conditional distributions:

$$\begin{aligned} t^{i+1} &\sim P(t|\varphi^i, \text{data}) \\ \varphi^{i+1} &\sim P(\varphi|t^{i+1}, \text{data}). \end{aligned}$$

A useful heuristic for determining the convergence rate of a Gibbs chain is the extent to which the two parameters  $t$  and  $\varphi$  are dependent in  $P(t, \varphi|\text{data})$ . A highly dependent posterior  $P(t, \varphi|\text{data})$  leads to a slow Gibbs chain, near independence leads to a fast Gibbs chain. Indeed, exact independence gives a sample of the posterior after one Gibbs step. A technique for accelerating the convergence of a Gibbs sampler is to find a re-parameterization of  $t$  and  $\varphi$  in a way which makes the posterior less dependent. In the remainder of this section we discuss a specific re-parameterization which, by analogy, can be applied to Bayesian lensing.



**Figure 1.** *Left:* density contours of the **ancillary** chain  $P(t, \varphi|\text{data})$  with 20 steps of a Gibbs sampler. *Right:* density contours of the **sufficient** chain  $P(\tilde{t}, \varphi|\text{data})$  with 20 steps of a Gibbs sampler. This illustrates the general heuristic that a slowly converging ancillary chain translates to a quickly converging the sufficient chain.

The relevant situation for Bayesian lensing is the case that  $t$  and  $\varphi$  are highly negatively correlated in  $P(t, \varphi|\text{data})$ . This motivates re-parameterizing  $(t, \varphi)$  to  $(\tilde{t}, \varphi)$  where  $\tilde{t} \equiv t + \varphi$  so that

$$\text{data} = \tilde{t} + n.$$

In the statistics literature,  $(t, \varphi)$  has been referred to as an **ancillary parameterization** whereas  $(\tilde{t}, \varphi)$  is referred to as a **sufficient parameterization**. We note that the terms *ancillary* and *sufficient* parameterization have been used interchangeably with the nomenclature *non-centered* and *centered* parameterizations, respectively, in the statistics literature Roberts et al. (2003); Gelfand et al. (1995); Papaspiliopoulos and Roberts (2008); Papaspiliopoulos et al. (2007); Yu and Meng (2011). Figure 1 illustrates the difference between an ancillary versus sufficient posterior distribution for our simple two parameter model. The left plot shows the posterior density contours for the ancillary parameterization  $(t, \varphi)$ , along with 20 steps of a Gibbs sampler. Conversely, the right plot shows the posterior density contours for the sufficient chain  $(\tilde{t}, \varphi)$  with 20 Gibbs steps. Notice that negative correlation in the ancillary parameterization manifests in near independence for the sufficient chain. Indeed, the slower the ancillary chain the faster the sufficient chain and vice-versa.

#### 4. ANCILLARY VERSUS SUFFICIENT PARAMETERS FOR THE LENSED CMB

The ancillary parameterization presented in the previous section is analogous to the lensed CMB problem as follows

$$\text{data}(x) = T(x + \nabla\phi(x)) + n(x) \quad \text{analogous to} \quad \text{data} = \tilde{T} + n$$

where the unlensed CMB temperature field  $T$  and the lensing potential  $\phi$  are the two unknown parameters. As was discussed in Section 2 the Gibbs chain based on the ancillary parameters  $T(x)$  and  $\phi(x)$  is exceedingly slow. This clearly motivates the following re-parameterization to sufficient parameters for the lensed CMB problem

$$\text{data}(x) = \tilde{T}(x) + n(x) \quad \text{analogous to} \quad \text{data} = \tilde{t} + n$$

where now  $\tilde{T}$  denotes the lensed CMB temperature field with no noise or beam. The sufficient chain then proceeds as

$$\tilde{T}^{i+1} \sim P(\tilde{T}|\phi^i, \text{data}) \tag{3}$$

$$\phi^{i+1} \sim P(\phi|\tilde{T}^{i+1}, \text{data}). \tag{4}$$

In Section 6 we adapt an iterative message passing algorithm, originally developed in Elsner and Wandelt (2013); Jasche and Lavaux (2015), for Wiener filtering and sampling from (3). In Section 5 we derive a Hamiltonian Markov Chain algorithm to sample from (4). Our Hamiltonian Markov Chain algorithm relies on an approximation—motivated again by the two parameter system—we call *anti-lensing*.

##### 4.1. Anti-lensing approximation

In the two parameter analogy from Section 3, the relation between the sufficient parameter  $\tilde{t}$  and the ancillary parameter  $t$  is given by  $\tilde{t} - \varphi = t$ . The corresponding relation for CMB lensing we refer to as *anti-lensing*:

$$\tilde{T}\left(\underbrace{x - \nabla\phi(x)}_{\text{anti-lensing}}\right) \approx T(x). \tag{5}$$

We distinguish between *inverse lensing* and *anti-lensing*. Inverse lensing denotes the true coordinate displacement which, when applied to  $\tilde{T}$ , recovers the unlensed  $T$ . Conversely, anti-lensing is given by  $-\nabla\phi$  and approximates inverse lensing.

To examine the difference between anti-lensing and inverse lensing notice that an extra divergence-free potential is needed to model the inverse lensing displacement field. Indeed, let  $f(x) := x + \nabla\phi(x)$  denote the lensing map. With this notation we have

$$\tilde{T}(x) = T(f(x)) \quad \text{and} \quad T(x) = \tilde{T}(f^{-1}(x))$$

where  $f^{-1}$  is the inverse lensing map that satisfies  $x = f^{-1}(f(x))$ . Now let  $d(x)$  denote the displacement vector field for inverse lensing so that  $f^{-1}(x) = x + d(x)$ . Therefore  $x = f^{-1}(f(x)) = f(x) + d(f(x))$ . In particular  $d(f(x)) = x - f(x) = -\nabla\phi(x)$  which gives

$$d(x) = -\nabla\phi(f^{-1}(x)).$$

This implies that the inverse lensing displacement is modeled as a warped version of the curl-free vector field  $-\nabla\phi$  (warped by  $f^{-1}$ ). This warping introduces a non-zero divergence-free term (just as lensing adds non-zero B-mode power in the CMB polarization).

To illustrate the expected magnitudes of the divergence-free and curl-free terms, start with a Helmholtz decomposition of the inverse lensing displacement:  $d(x) = -\nabla\phi^{\text{inv}}(x) - \nabla^\perp\psi^{\text{inv}}(x)$ , where  $\nabla^\perp \equiv (-\frac{\partial}{\partial y}, \frac{\partial}{\partial x})$  and  $\psi^{\text{inv}}$  denotes a stream function potential which models a field rotation so that

$$\tilde{T}\left(\underbrace{x - \nabla\phi^{\text{inv}}(x) - \nabla^\perp\psi^{\text{inv}}(x)}_{\text{inverse lensing}}\right) = T(x).$$

Due to the fact that the expected size of the inverse lensing displacement  $d(x)$  is on the order of arcmin but the correlation length scale of  $\phi$  is on the order of degrees we claim that  $-\nabla\phi(f^{-1}(x))$  is well approximated by  $-\nabla\phi(x)$ . In particular, the divergence-free term  $-\nabla^\perp\psi^{\text{inv}}$  is small and

$$-\nabla\phi \approx -\nabla\phi^{\text{inv}} \approx -\nabla\phi^{\text{inv}} - \nabla^\perp\psi^{\text{inv}} = d. \quad (6)$$

Figure 2 illustrates the magnitudes of the above terms. The anti-lensing potential  $-\phi$  is shown (upper-left) along with the corresponding inverse lensing potential  $-\phi^{\text{inv}}$  (upper-right). The difference  $\phi - \phi^{\text{inv}}$  is also shown (bottom left) along with the stream function  $-\psi^{\text{inv}}$  (bottom-right). Clearly, the magnitude of the difference  $\phi^{\text{inv}} - \phi$  and  $-\psi^{\text{inv}}$  is sub-dominant to estimation error expected in current lensing experimental conditions.

## 5. HAMILTONIAN MONTE CARLO SAMPLER FOR $P(\phi|\tilde{T}, \text{data})$

The Hamiltonian Monte Carlo (HMC) algorithm is an iterative sampling algorithm designed to mitigate the low-acceptance rate of the Metropolis-Hastings algorithm when working in high dimension. A nice review of HMC can be found in Neal (2011). For applications of HMC in cosmology see Hajian (2007); Taylor et al. (2008); Elsner and Wandelt (2010); Jasche et al. (2010); Jasche and Wandelt (2012, 2013a,b). In the present case we utilize the HMC algorithm to produce samples of  $\phi$  from  $P(\phi|\tilde{T}, \text{data})$ . The key to making HMC work for lensing is to parameterize  $\phi$  in terms of its Fourier transform. One can then utilize Claim 1, presented below, to efficiently compute the gradient of the log conditional density of  $P(\phi|\tilde{T}, \text{data})$ , which is a necessary computation for the HMC algorithm.

*Notation:* Throughout the remainder of this paper, the Fourier transform of any function  $f(x)$  will be denoted by  $f_l$  or  $f_k$  so that  $f_l = \int_{\mathbb{R}^2} e^{-ix \cdot l} f(x) \frac{dx}{2\pi}$  and  $f(x) = \int_{\mathbb{R}^2} e^{ix \cdot l} f_l \frac{dl}{2\pi}$  where  $l \in \mathbb{R}^2$  is a two dimensional frequency vector and  $x \in \mathbb{R}^2$  is a two dimensional spatial coordinate.

To describe the HMC algorithm let  $\phi$  denote the concatenation of the real and imaginary parts of  $\phi_l$  as  $l$  ranges through discrete frequencies  $l$  ranging up to a pre-specified  $|l|_{\max}$  (but excluding half of the Fourier frequencies due to the Hermitian symmetry associated with the Fourier transform of a real field). Note that  $\phi$  is a vector of real numbers.

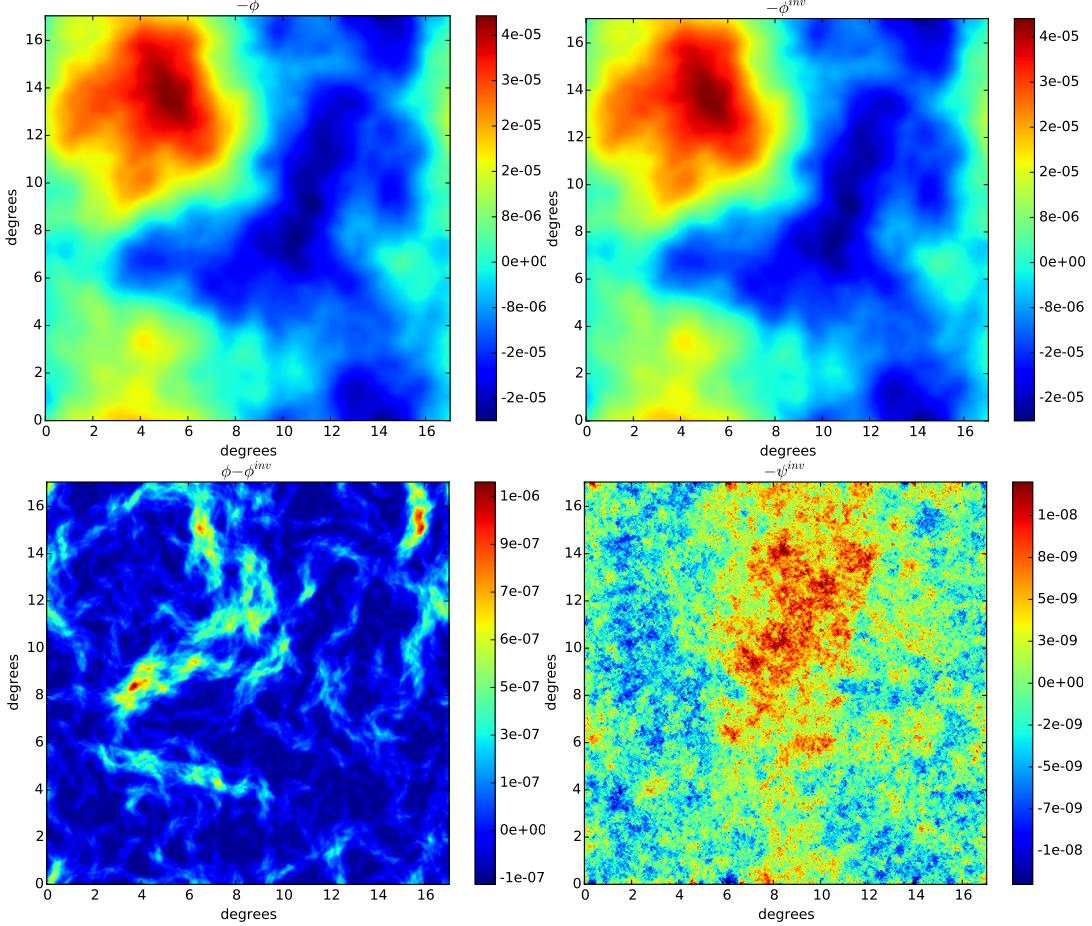
Let  $P(\phi|\tilde{T}, \text{data})$  denote the density of  $\phi$  given  $\tilde{T}$  and the data. Let  $\mathbf{p}$  denote a ‘momentum’ vector and  $\mathbf{m}$  denote a ‘mass’ vector, which are both the same length as  $\phi$ . The Hamiltonian is a function of  $\phi$  and  $\mathbf{p}$  and is defined as follows

$$H(\phi, \mathbf{p}) := -\log P(\phi|\tilde{T}, \text{data}) + \sum_k \frac{\mathbf{p}_k^2}{2\mathbf{m}_k}.$$

This Hamiltonian generates a time-dependent evolution of  $\phi$  and  $\mathbf{p}$  given by

$$\begin{aligned} \frac{d\phi^t}{dt} &= \nabla_{\mathbf{p}} H(\phi^t, \mathbf{p}^t) \\ \frac{d\mathbf{p}^t}{dt} &= -\nabla_{\phi} H(\phi^t, \mathbf{p}^t). \end{aligned}$$

The HMC is a discrete version of this time-dynamic equation, using a leapfrog method, which produces a Markov chain  $(\phi_1, \mathbf{p}_1), (\phi_2, \mathbf{p}_2), \dots$  where the  $i^{\text{th}}$  iteration is given by Algorithm 1 below.



**Figure 2.** The difference between anti-lensing and inverse lensing. *Upper left:* anti-lensing potential  $-\phi$ . *Upper right:* The inverse lensing potential  $-\phi^{inv}$ . *Bottom left:* The difference  $\phi^{inv} - \phi$ . *Bottom right:* The inverse lensing stream function  $-\psi^{inv}$ .

---

**Algorithm 1**  $i^{\text{th}}$  step of the Hamiltonian Markov Chain

---

- 1: Set  $\phi^0 := \phi_{i-1}$  and simulate  $\mathbf{p}^0 \sim \mathcal{N}(0, \Lambda_m)$  where  $\Lambda_m$  is diagonal with  $\text{diag}(\Lambda_m) = \mathbf{m}$ .
- 2: Recursively compute  $\phi^{k\epsilon}$  and  $\mathbf{p}^{k\epsilon}$  for  $k = 1, \dots, n$  using the following equations:

$$\begin{aligned}\phi^{t+\epsilon} &:= \phi^t + \epsilon \Lambda_m^{-1} \left[ \mathbf{p}^t - \frac{\epsilon}{2} \nabla_\phi H(\phi^t, \mathbf{p}^t) \right], \\ \mathbf{p}^{t+\epsilon} &:= \mathbf{p}^t - \frac{\epsilon}{2} \left[ \nabla_\phi H(\phi^t, \mathbf{p}^t) + \nabla_\phi H(\phi^{t+\epsilon}, \mathbf{p}^t) \right].\end{aligned}$$

- 3: Simulate  $u \sim \mathcal{U}(0, 1)$ , and define  $p := \min \left( 1, e^{-H(\phi^{n\epsilon}, \mathbf{p}^{n\epsilon})} / e^{-H(\phi^0, \mathbf{p}^0)} \right)$ .
  - 4: If  $u < p$ , set  $\phi_i := \phi^{n\epsilon}$ , otherwise set  $\phi_i := \phi_{i-1}$ .
- 

The HMC algorithm is notoriously sensitive to tuning parameters. The prevailing wisdom (see Neal (2011) page 22, for example) is that one should set  $\mathbf{m}$  to match the reciprocal of the posterior variance of  $\phi$ . For the simulation presented in Section 7 we simply set  $\mathbf{m}_l^{-1}$  to be nearly proportional to  $C_l^{\phi\phi}$  with a slight attenuation at low wavenumber. In particular, we set  $\mathbf{m}_l^{-1} := 2 \times 10^2 \left[ \frac{3}{4} + \frac{1}{4} \tanh\left(\frac{|l|-1500}{200}\right) \right] C_l^{\phi\phi} \delta_0$  where  $\delta_0$  denotes the constant obtained by evaluating the discrete Dirac delta function  $\delta_l$  at  $l = 0$  (see the Appendix for a detailed discussion). This choice was motivated by the fact that the high frequency terms  $\phi_l$  are not well constrained by the posterior distribution which results in a posterior variance closely matching  $C_l^{\phi\phi}$ . The remaining parameters of Algorithm 1 are set to  $n = 30$  and  $\epsilon = 2 \times 10^{-3} u$  where  $u$  is a uniform  $(0, 1)$  random variable sampled anew at each pass of Algorithm 1 (the use of random  $\epsilon$  is designed to avoid resonant frequencies, as advocated in Taylor et al. (2008)).

The key difficulty in using Algorithm 1 is the computation of the  $\nabla_\phi H(\phi^t, \mathbf{p}^t)$ , or equivalently the computation of  $\nabla_\phi \log P(\phi | \tilde{T}, \text{data})$ . The number of frequencies is extremely large and therefore, any slow computation of the gradients will present a serious bottleneck. The follow claim shows that the gradient of the log density of  $P(\phi | \tilde{T}, \text{data})$ ,

with respect to the Fourier basis of  $\phi$ , can be computed quickly with Fourier and inverse Fourier transforms.

**Claim 1.** *Under the anti-lensing approximation (5) for any nonzero frequency vector  $l \equiv (l_1, l_2) \in \mathbb{R}^2$*

$$\frac{\partial}{\partial \phi_l} \log P(\phi | \tilde{T}, \text{data}) \propto -\frac{\phi_l}{C_l^{\phi\phi}} - \sum_{q=1,2} il_q \int_{\mathbb{R}^2} e^{-ix \cdot l} A^q(x) B(x) \frac{dx}{2\pi} \quad (7)$$

where  $\phi_l = \text{re}\phi_l + i\text{im}\phi_l$ ,  $\frac{\partial}{\partial \phi_l} \equiv \frac{\partial}{\partial \text{re}\phi_l} + i\frac{\partial}{\partial \text{im}\phi_l}$  and

$$B_l \equiv \frac{1}{C_l^{TT}} \int e^{-ix \cdot l} \tilde{T}(x - \nabla\phi(x)) \frac{dx}{2\pi} \quad (8)$$

$$A^q(x) \equiv \frac{\partial \tilde{T}}{\partial x_q}(x - \nabla\phi(x)). \quad (9)$$

An important fact used in the derivation of (7) is that the lensing and anti-lensing operator is invertible. For example, if  $\phi(x)$  and  $\tilde{T}(x)$  are known at all pixel locations  $x$ , then it is possible to perfectly reconstruct  $T(x)$ . This implies that the anti-lensing operation (which is a linear action on the CMB) can be represented as an infinitesimal permutation matrix. Therefore, the determinant of the anti-lensing operator  $\det(d\tilde{T}^\phi/d\tilde{T})$  equals 1, where  $\tilde{T}^\phi(x) \equiv \tilde{T}(x - \nabla\phi(x))$ . Now, to compute the likelihood surface  $P(\phi | \tilde{T}, \text{data})$  as a function of  $\phi$  we obtain the following formula:

$$\begin{aligned} P(\phi | \tilde{T}, \text{data}) &= P(\phi | \tilde{T}) \\ &\propto P(\tilde{T} | \phi) P(\phi) \\ &= \underbrace{|\det(d\tilde{T}^\phi/d\tilde{T})|}_{=1} P(\tilde{T}^\phi | \phi) P(\phi) \end{aligned}$$

where  $P(\tilde{T}^\phi | \phi)$  represents the likelihood that  $\tilde{T}^\phi$  is statistically unlensed by  $\phi$ . In other words,  $P(\tilde{T}^\phi | \phi)$  measures the likelihood that  $\tilde{T}^\phi$  is an isotropic Gaussian random field with spectral density  $C_l^{TT}$ . This explains the following characterization of the log likelihood of  $\phi$  given  $\tilde{T}$  and the data:

$$\log P(\phi | \tilde{T}, \text{data}) = c - \frac{1}{2} \int_{\mathbb{R}^2} \left[ \frac{|\tilde{T}_k^\phi|^2}{C_k^{TT}} + \frac{|\phi_k|^2}{C_k^{\phi\phi}} \right] dk \quad (10)$$

where  $c$  is a constant which does not depend on  $\phi$ . The remaining details of the derivation of Claim 1 is left to the appendix.

It is instructive to compare the gradient calculation (7) with the quadratic estimate of  $\phi$  developed in Hu (2001); Hu and Okamoto (2002). The quadratic estimate, applied to observations of the form  $\tilde{T}(x) + n(x)$ , is given by

$$\hat{\phi}_l = -N_l \sum_{q=1,2} il_q \int_{\mathbb{R}^2} e^{-ix \cdot l} A^q(x) B(x) \frac{dx}{2\pi} \quad (11)$$

where  $B_l \equiv [C_l^{\tilde{T}\tilde{T}} + C_l^{nn}]^{-1} [\tilde{T}_l + n_l]$ ,  $A_l^q \equiv il_q [C_l^{TT}] [C_l^{\tilde{T}\tilde{T}} + C_l^{nn}]^{-1} [\tilde{T}_l + n_l]$  and

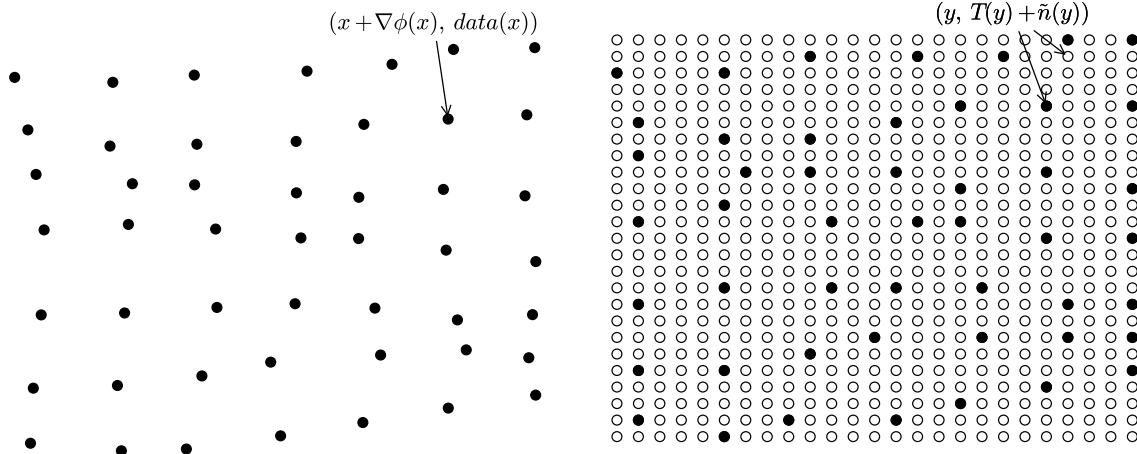
$$N_l^{-1} \equiv \frac{1}{2} \int_{\mathbb{R}^2} \frac{(l \cdot (k+l) C_{k+l}^{TT} - l \cdot k C_k^{TT})^2}{(C_{k+l}^{\tilde{T}\tilde{T}} + C_{k+l}^{nn})(C_k^{\tilde{T}\tilde{T}} + C_k^{nn})} \frac{dk}{(2\pi)^2}.$$

The term  $N_l$  is radially symmetric in frequency  $l$  and corresponds to a normalization which makes the quadratic estimate unbiased up to first order. After substituting  $C_l^{\tilde{T}\tilde{T}} + C_l^{nn} \rightarrow C_l^{TT}$  and  $n_l \rightarrow 0$  in the formula for the quadratic estimate, one obtains

$$\frac{\partial}{\partial \phi_l} \log P(\phi | \tilde{T}, \text{data}) \Big|_{\phi=0} = \frac{\hat{\phi}_l}{N_l}.$$

Indeed, an approximate Newton step, using (7), is an accurate approximation to the quadratic estimate  $\hat{\phi}_l$ . This illustrates how the parameterization  $(\tilde{T}, \phi)$  results in Gibbs iterations which make drastic moves, on the order of the size of the quadratic estimate.

One of the features of the quadratic estimate is that the fast Fourier transform (FFT) and inverse fast Fourier transform (IFFT) can be used to compute  $\hat{\phi}_l$  for all frequencies  $l$ . Naively computing the quadratic form of the quadratic estimate requires  $O(n^2)$  flops, rather than the  $O(n \log n)$  flops obtained by the FFT/IFFT method, where  $n$  denotes the number of pixels. We note that Claim 1 establishes that the gradient computation inherits a similar FFT/IFFT characterization to compute  $\frac{\partial}{\partial \phi_l} \log P(\phi | \tilde{T}, \text{data})$  at all frequencies  $l$ , in  $O(n \log n)$  flops. Since this gradient



**Figure 3.** This graphic illustrates how knowledge of  $\phi(x)$ , used when sampling from the Gibbs step  $P(\tilde{T}|\phi, \text{data})$ , converts white noise corrupted gridded observations of the lensed CMB into masked observations of the unlensed CMB on a more dense grid. The data associated with the original grid, indexed by  $x$ , is moved via advection to the lensed grid  $x + \nabla\phi(x)$  seen at left. The right panel shows the lensed grid embedded into a higher resolution grid. The data, indexed by the dense regular grid on the right panel, is of the form  $T(y) + \tilde{n}(y)$  where  $T(y)$  denotes the unlensed CMB and  $\tilde{n}(y)$  denotes the noise. The unobserved locations, represented by open circles, are characterized with an infinite variance for  $\tilde{n}(y)$ .

computation needs to be embedded in a Hamiltonian Markov step within a Gibbs Chain, the computation efficiency gained by the FFT/IFFT is absolutely crucial.

## 6. ITERATIVE MESSAGE PASSING ALGORITHM FOR $P(\tilde{T}|\phi, \text{DATA})$

There are two natural ways to model the lensed CMB  $\tilde{T}$ . If one marginalizes out  $\phi$ , then  $\tilde{T}$  is modeled as a *non-Gaussian* but isotropic random field. Conversely, if one conditions on  $\phi$  the field  $\tilde{T}$  is modeled as a *non-isotropic* but Gaussian random field. The latter case is relevant for sampling from  $P(\tilde{T}|\phi, \text{data})$  which is, therefore, simply a Gaussian conditional simulation problem. Unfortunately, the non-isotropic (indeed, non-stationary) nature of the conditional distribution of  $\tilde{T}$  presents serious computational challenges. In what follows we utilize a new iterative algorithm developed in Elsner and Wandelt (2013) for Gaussian conditional expectation when the signal is diagonalized in harmonic space that the noise is diagonalized in pixel space. The method we present here is similar to the Gibbs sampling adaptation of Jasche and Lavaux (2015).

Start by transforming each pixel location  $x$  by the lensing operation  $x + \nabla\phi(x)$ , while simultaneously preserving the data associated with that pixel. This effectively de-lenses  $\text{data}(x) = T(x + \nabla\phi(x)) + n(x)$  but produces observations on an irregular grid. In particular, one may switch to the lensed coordinates  $y = x + \nabla\phi(x)$  so that

$$\underbrace{(x + \nabla\phi(x), \text{data}(x))}_{\text{(pixel, data) tuple}} = (y, T(y) + \tilde{n}(y))$$

where  $\tilde{n}(x + \nabla\phi(x)) = n(x)$ . Now the data  $(y, T(y) + \tilde{n}(y))$  is arranged on an irregular grid in  $y$ . This irregular grid is then embedded into a high resolution regular grid by nearest neighbor interpolation. The points  $y$  which do not get assigned an observation  $T(y) + \tilde{n}(y)$  under the interpolation we consider to be masked. Figure 3 illustrates this situation. The left hand plot shows the irregularly sampled data  $(x + \nabla\phi(x), \text{data}(x))$  and the right hand plot shows the grid embedding. The filled dots represent observations of  $T(y) + \tilde{n}(y)$  whereas the empty dots correspond to a masked observation of  $T(y)$ . Finally we extend the definition of  $\tilde{n}(y)$  to have infinite variance over the masked region, whereby producing data  $T(y) + \tilde{n}(y)$  over a dense regular grid in  $y$ .

As an intermediate step in producing a sample from  $P(\tilde{T}|\phi, \text{data})$  we produce a conditional sample of  $T(y)$  given the observations  $T(y) + \tilde{n}(y)$ . The difficulty of this step is that  $\tilde{n}(y)$  is non-homogeneous noise—from the masking and any inhomogeneity in  $n(x)$ —and therefore it is not decorrelated by the Fourier transform. To handle this situation we adapted a new method for Gaussian conditional expectation developed in Elsner and Wandelt (2013). This method works particularly well for observations with large amounts of irregular masking, as in our case. The algorithm utilizes a messenger field which effectively behaves as a latent—signal plus white noise—model which is amenable to Gibbs sampling (Jasche and Lavaux 2015).

The delensing algorithm described in this paper requires the capability to do fast constrained realization of non-lensed CMB. The embedding illustrated in Figure 3 means that each constrained realization needs to be computed on a mask with a great deal of structure on the scale of the pixels of dense grid. Several algorithms exist in the literature to solve the general problem of constrained Gaussian random field on a given mask and power spectrum:

the conjugate gradient method (Wandelt et al. 2004; Eriksen et al. 2004), the multiscale conjugate gradient method (Smith et al. 2007), the multigrid method (Seljebotn et al. 2014), the Messenger algorithm (Elsner and Wandelt 2013) and its variant the Gibbs-Messenger (Jasche and Lavaux 2015).

Since the lensing potential changes from each iteration to the next, every constrained realization of non-lensed CMB needs to be computed for a different set of active points in the embedding grid. This effectively means that the solution is done for a different mask at every iteration. This rules out linear solvers that require expensive pre-computations, e.g. of pre-conditioners, that depend on the coefficient matrix of the system since that depends on the mask. We also require exact acceptance and higher speed than direct or standard iterative methods. This reduces the possibilities to either the Messenger algorithm or the Gibbs-Messenger.

The Gibbs-Messenger generates very fast constrained realizations that converge to the correct distribution in a statistical sense without iterating to a numerical solution (thus obviating the need to specify a cooling schedule in Algorithm 2) for the price of losing independence between subsequent samples. In contrast the Messenger algorithm simulates independent constrained realisations, but requires iteration of the linear system. To ensure a numerically accurate solution implies a conservative choice of cooling schedule in Algorithm 2 (though this is still much faster than the alternatives described in the previous paragraph).

We adopt a hybrid approach where we occasionally generate an independent sample using the Messenger algorithm, and then generate many quick samples using the Gibbs-Messenger approach. The detailed choices for the cooling schedule and the number of samples between full Messenger solutions will be described in Section 7.

---

**Algorithm 2** Algorithm for sampling from  $P(T|\text{data})$  where  $\text{data}(y) = T(y) + \tilde{n}(y)$ 


---

- 1: Set cooling schedule  $\lambda_1 \geq \dots \geq \lambda_n$  where  $\lambda_n = 1$ .
- 2: Decompose  $\tilde{n}(y)$  into a homogeneous part with variance  $\bar{\sigma}^2$  and a non-homogeneous part with variance  $\tilde{\sigma}^2(y)$  so that

$$\text{var}(\tilde{n}(y)) = \bar{\sigma}^2 + \tilde{\sigma}^2(y)$$

where  $\tilde{\sigma}^2(y) = \infty$  on all masked pixels  $y$ . Notice that the spectral density of the homogeneous part is given by  $\bar{\sigma}^2 dy$  where  $dy$  denotes the pixel grid area.

- 3: Initialize the fields  $M(y)$  and  $T(y)$  to be zero at all pixel locations  $y$ .
- 4: Recursively update fields  $M(y)$  and  $T(y)$  by iterating the following steps for  $j = 1, \dots, n$ :

- Simulate a mean zero Gaussian random field  $Z(y)$  which is independent across pixels and with pointwise variance  $(\frac{1}{\lambda_j \bar{\sigma}^2} + \frac{1}{\tilde{\sigma}^2(y)})^{-1}$ .
- Update  $M(y) \leftarrow \text{data}(y) \frac{\lambda_j \bar{\sigma}^2}{\lambda_j \bar{\sigma}^2 + \tilde{\sigma}^2(y)} + T(y) \frac{\tilde{\sigma}^2(y)}{\lambda_j \bar{\sigma}^2 + \tilde{\sigma}^2(y)} + Z(y)$ .
- Simulate a mean zero Gaussian random field,  $W(y)$ , with spectral density  $\langle W_l W_{l'}^* \rangle = \delta_{l-l'} (\frac{1}{C_l^{TT}} + \frac{1}{\lambda_j \bar{\sigma}^2 dy})^{-1}$
- Update  $T_l \leftarrow M_l \frac{C_l^{TT}}{C_l^{TT} + \lambda_j \bar{\sigma}^2 dy} + W_l$ .

- 5: Return  $T(x)$ .
- 

The following algorithm describes the use of Algorithm 2 to produce a sample from  $P(\tilde{T}|\phi, \text{data})$ .

---

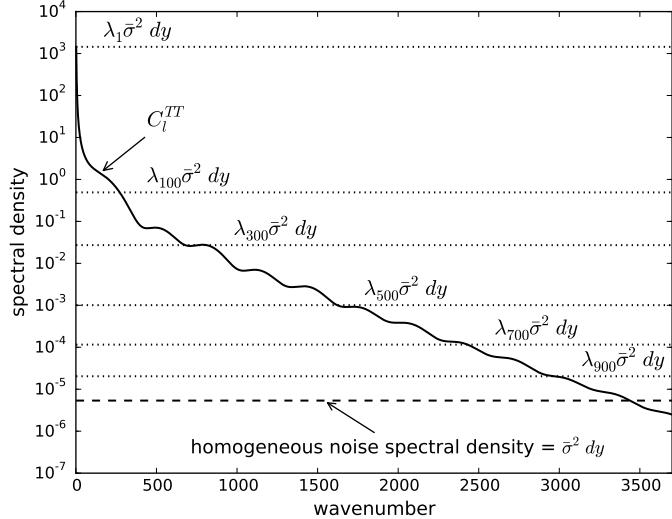
**Algorithm 3** Sampling from  $P(\tilde{T}|\phi, \text{data})$ 


---

- 1: Embedded the pixel/data pairs  $(x, \text{data}(x))$  into observations of the form  $(y, T(y) + \tilde{n}(y))$  where  $y$  ranges over a high resolution regular grid as illustrated in the right plot of Figure 3.
  - 2: Use Algorithm 2 to produce a sample  $T \sim P(T|T + \tilde{n})$ .
  - 3: Return  $\tilde{T}(x) = T(x + \nabla\phi(x))$ .
- 

At present, algorithms 2 and 3 are designed for the situation that the pixels are sufficiently small compared to the magnitude of  $\nabla\phi(x)$  and the noise is approximately white on these scales. Indeed, our goal is to explore the low noise and small beam experimental conditions where the quadratic estimate is known to be suboptimal (see Hirata and Seljak (2003a,b)). That being said, the only change needed to incorporate other experimental details in the Bayesian lensing methodology presented here, including foreground contaminants, is how one samples from  $P(\tilde{T}|\phi, \text{data})$ . Algorithms 2 and 3 take advantage of the special lensed-grid structure of the data when conditioning on  $\phi$  to accomplish this goal. Adding different/new experimental details to the data will still result in a Gaussian constrained realization problem. We acknowledged that more complicated modeling of the data will most certainly introduce additional computational challenges. However, we consider these computational challenges to be sub-dominant to the fundamental bottleneck for Bayesian lensing which was the extremely slow mixing time of the original Gibbs formulation. Moreover, the structure of the new Gibbs formulation isolates all experimental details to the Gibbs step (4) where conditioning on the lensing potential  $\phi$  results in a classic Gaussian constrained realization problem.

## 7. SIMULATION EXAMPLE



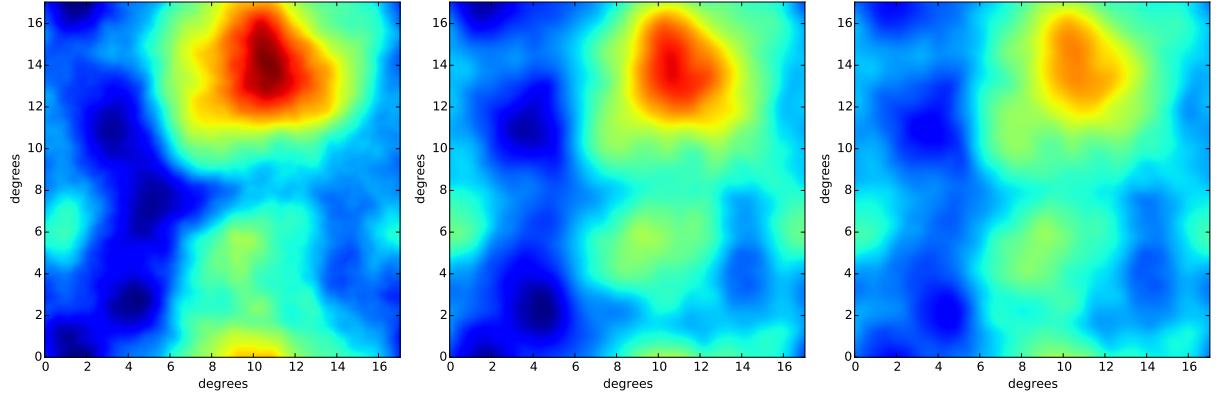
**Figure 4.** This figure shows values of the cooling parameter  $\lambda_j$  for  $j = 1, 100, 300, 500, 700, 900$  used in Algorithm 2. This schedule is applied every 100<sup>th</sup> step in the Gibbs algorithm. In Algorithm 2, the value of  $\lambda_j \bar{\sigma}^2 dy$  serves as the spectral density of artificial additive white noise in the latent field  $M(x)$ . Therefore setting  $\lambda_j$  greater than 1, encourages fast mixing of  $T_l$  at all frequency vectors  $l$  such that  $C_l^{TT} \lesssim \lambda_j \bar{\sigma}^2 dy$ . The cooling schedule shown above is an attempt to let  $\lambda_j$  approach 1 in such a way as to encourage all frequency vectors up to  $|l|_{\max}$  to mix quickly.

In this section we present a simulation to illustrate the methodology presented above. The simulated lensing potential used in this section, shown at left in Figure 5, is generated on a flat sky with periodic boundary conditions. The data, shown upper-left in Figure 6, is generated on 2 arcmin pixels with independent additive noise and masking. The noise level is set to  $8.0 \mu K$  arcmin and the masking covers approximately 10% of the pixels. The parameters of the Bayesian lensing procedure are the Fourier modes of  $\tilde{T}$  and  $\phi$ . For the lensing potential we set  $|l|_{\max}$  to 460. For this sky coverage the scale-resolution in Fourier space  $\Delta l = 21$  yields 1500 unknown Fourier coefficients for  $\phi_l$ . The  $|l|_{\max}$  of 2700 for the unlensed temperature  $T$  is set in Algorithm 2 and corresponds to half of the Nyquist limit at 2 arcmin pixels. CAMB (Lewis et al. 2000) is used to generate the fiducial theoretical power spectra,  $C_l^{TT}$  and  $C_l^{\phi\phi}$ , for simulating the data and for the Bayesian priors on  $T$  and  $\phi$  (both fields are assumed to be Gaussian). The fiducial cosmology used in CAMB is based on a flat, power law CDM cosmological model, with baryon density  $\Omega_b = 0.044$ ; cold dark matter density  $\Omega_{cdm} = 0.21$ ; cosmological constant density  $\Omega_\Lambda = 0.74$ ; Hubble parameter  $h = 0.71$  in units of  $100 \text{ km s}^{-1} \text{Mpc}^{-1}$ ; primordial scalar fluctuation amplitude  $A_s(k = 0.002 \text{ Mpc}^{-1}) = 2.45 \times 10^{-9}$ ; scalar spectral index  $n_s(k = 0.002 \text{ Mpc}^{-1}) = 0.96$ ; primordial helium abundance  $Y_P = 0.24$ ; and reionization optical depth  $\tau_r = 0.088$ .

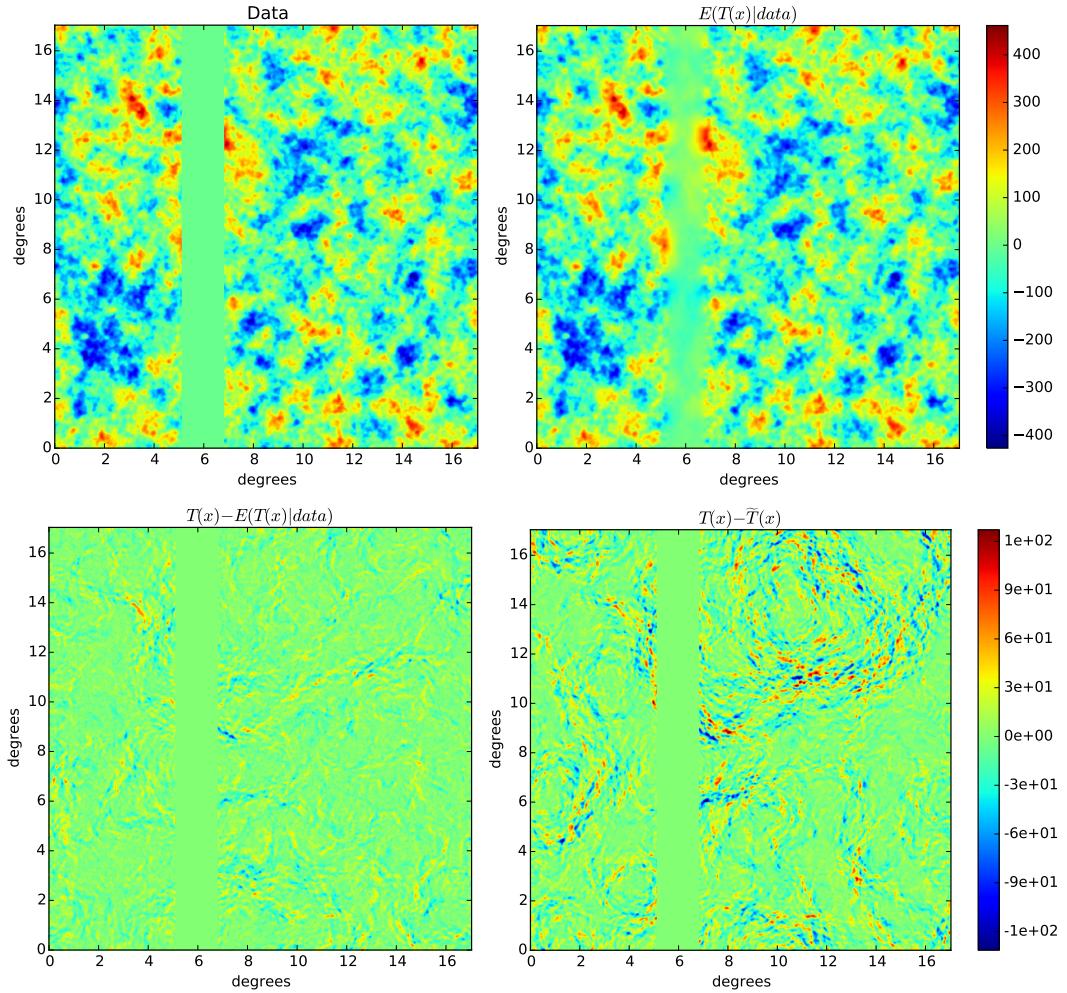
We ran 10 parallel Gibbs chains for a total of 2500 steps. The timings for each Gibbs iteration averaged approximately 200 seconds using a Dual Intel Xeon E5-2690 2.90GHz processor. Each chain was initially warmed up by replacing the HMC draws in the first 5 iterations with a gradient ascent. A burn-in of approximately 550 runs were discarded and the remaining runs were thinned by 100. The result is a total of 200 posterior samples. The cooling schedule for the iterative message passing algorithm was selected by numerical experimentation. Most of the Gibbs iterations set the cooling terms  $(\lambda_1, \dots, \lambda_{400}) \equiv (1, \dots, 1)$  in Algorithm 2. However, we did find it advantageous to periodically run a nontrivial 1000-step cooling schedule every 100<sup>th</sup> pass of the Gibbs algorithm. This nontrivial cooling schedule for  $\lambda_j$  is plotted in Figure 4 and is set in an attempt to encourage fast mixing of the Fourier modes  $T_l$  up to  $|l|_{\max}$ .

The best Bayesian estimate of  $\phi(x)$  corresponds to the posterior mean  $E(\phi(x)|\text{data})$ . This quantity is approximated by the average of the 200 draws from the Gibbs chain and is shown in the middle plot of Figure 5. The right plot of Figure 5 shows the quadratic estimate of  $\phi(x)$  for comparison. However, due to the difficulty when using the quadratic estimate in the presence of sky cuts, the quadratic estimate shown uses all of the data—including the pixels which are masked—in producing the estimate of  $\phi$ . In general, one can see good agreement with  $E(\phi(x)|\text{data})$  and  $\phi(x)$ . Indeed, the effect of masking is visually undetectable as compared to the quadratic estimate. To get a better visualization of the individual draws from the posterior, the left plot in Figure 7 shows a horizontal cross section of the posterior draws of  $\phi(x)$  taken at vertical degree mark  $12.7^\circ$ .

The Gibbs methodology presented here yields samples of the lensed CMB  $\tilde{T}(x)$  and the lensing potential  $\phi(x)$  conditional on the data. Moreover, as a byproduct of Algorithm 3, we also obtain samples of the unlensed  $T(x)$  given the data. By averaging 200 draws from the Gibbs chain one can construct an approximation to  $E(T(x)|\text{data})$ , shown in the upper-right plot of Figure 6. In the bottom-left plot of Figure 6 we show the difference  $T(x) - E(T(x)|\text{data})$ . When compared to the nominal difference between lensed and unlensed CMB  $T(x) - \tilde{T}(x)$ , shown bottom-right in Figure 6, one can see that the Gibbs methodology is successful at delensing the observed CMB. To get a better visualization of the individual draws from the posterior, the right plot in Figure 7 shows a horizontal cross section of the posterior



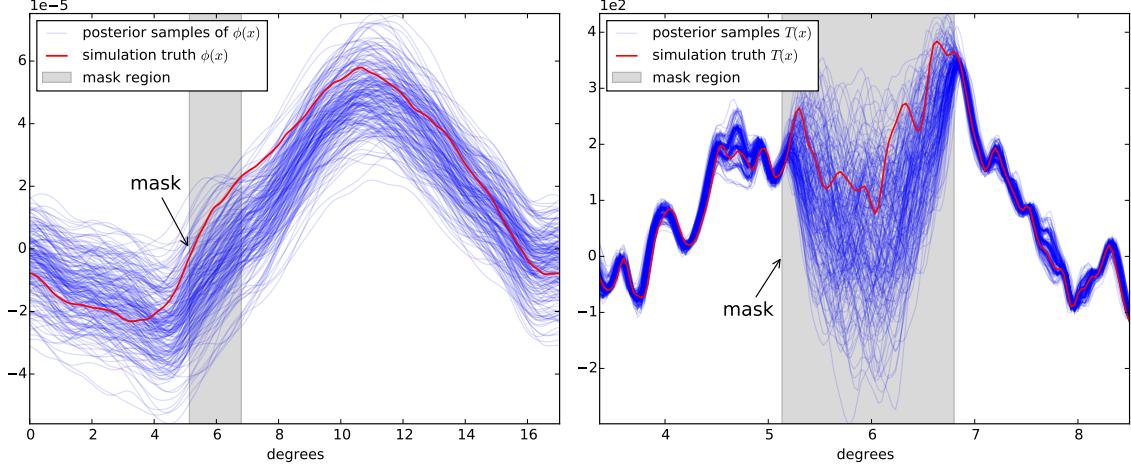
**Figure 5.** *Left:* simulation truth  $\phi(x)$ . *Middle:* posterior mean  $E(\phi(x)|\text{data})$ . *Right:* The quadratic estimate. To avoid difficulties associated with masked data when using the quadratic estimate, the estimate shown at right is applied to the full data set with the masked region removed. In contrast, the data used for the Bayesian methodology is masked as shown in Figure 6.



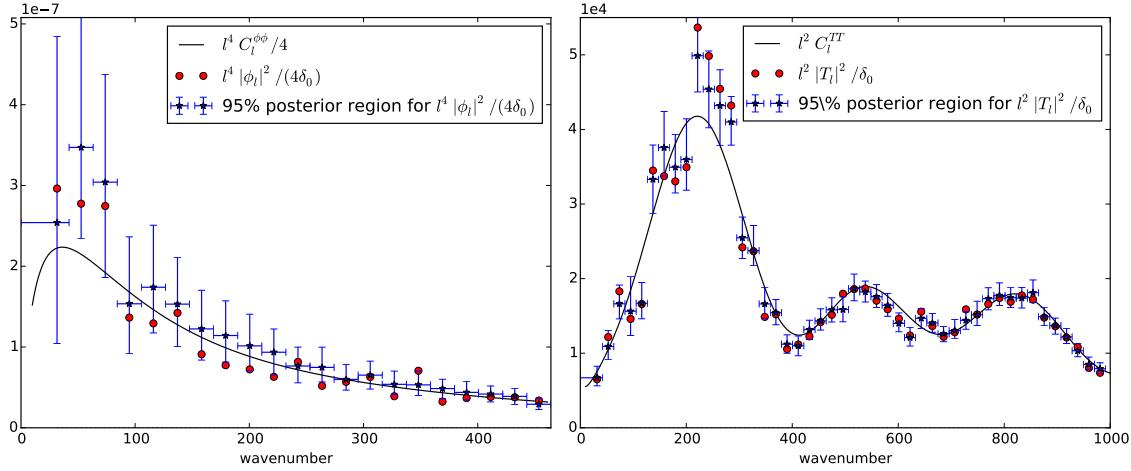
**Figure 6.** *Upper left:* simulated lensed CMB data with masking and additive white noise (at level  $8.0 \mu K$  arcmin). *Upper right:* posterior mean  $E(T(x)|\text{data})$ . *Lower left:* This plot shows  $T(x) - E(T(x)|\text{data})$  and probes the ability of the Bayesian methodology to delense the observations. This plot should be compared with the nominal difference between the simulation truth unlensed CMB and the lensed CMB,  $T(x) - \tilde{T}(x)$ , shown *bottom right*.

draws of  $T(x)$  taken at vertical degree mark  $12.7^\circ$  and is magnified near the masking region for better visual inspection.

In Figures 8 and 9 we summarize the posterior draws for  $\phi$  and  $T$  in the Fourier domain. The left plot of Figure 8 shows 95% posterior regions for  $l^4|\phi_l|^2/(4\delta_0)$  averaged over  $l$  in wavenumber bins. For comparison the simulation true



**Figure 7.** Here we plot one dimensional slices of the posterior draws from  $P(\phi(x)|\text{data})$  and  $P(T(x)|\text{data})$ . For reference, these slices are taken from the horizontal regions at vertical degree mark  $12.7^\circ$  in Figures 6 and 5.



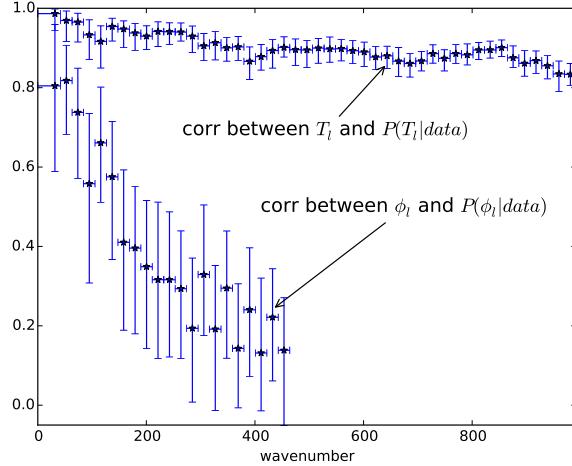
**Figure 8.** Estimates of  $l^4|\phi_l|^2/4$  and  $l^2|T_l|^2$  (shown in blue), scaled to the units of the corresponding spectral density. The red dots show  $l^4|\phi_l|^2/4$  and  $l^2|T_l|^2$  for the simulation truth (similarly scaled). The discrepancy between the red dots and the spectral densities, shown in black, is exclusively due to cosmic variance. The confidence bars show 95% probability regions from the posterior distributions  $P(l^4|\phi_l|^2/4|\text{data})$  and  $P(l^2|T_l|^2|\text{data})$ .

values of  $l^4|\phi_l|^2/(4\delta_0)$  are shown in red and the spectral density  $l^4C_l^{\phi\phi}/4$  is plotted in the black solid line. The same quantities are shown for  $l^2|T_l|^2/(\delta_0)$  in the right plot of Figure 8. Finally, in Figure 9 we show the empirical cross correlation of the posterior draws over wavenumber bins between the simulation truth and the posterior samples of  $\phi_l$  and  $T_l$ . In particular, samples were generated from  $\frac{1}{c} \sum_{l \in \Delta l} \phi_l^{\text{sim}} \phi_l^*$  and  $\frac{1}{c} \sum_{l \in \Delta l} T_l^{\text{sim}} T_l^*$  where  $\Delta l$  is a frequency wavenumber bin,  $\phi_l$  and  $T_l$  are the simulation truth,  $\phi_l^{\text{sim}}$  and  $T_l^{\text{sim}}$  are sampled from the Gibbs algorithm presented here and  $c$  is a normalization constant which transforms to a correlation scale. Notice that the plotted correlations trend to 0 for larger wavenumber. This is what one would expect since larger wavenumber have correspondingly smaller signal-to-noise ratio which causes the posterior to revert back to the prior.

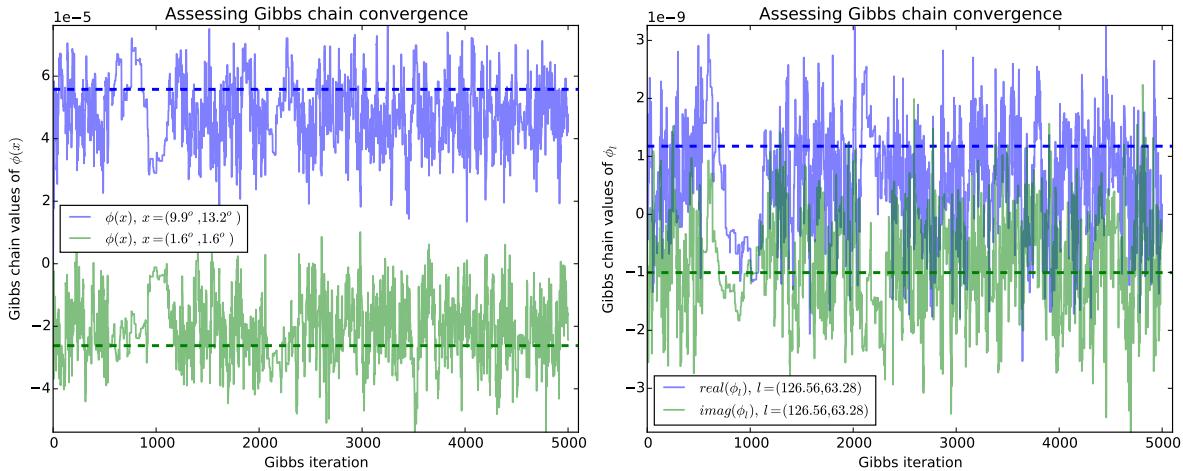
In Figure 10 we show the Gibbs chain correlation length scale and the speed of mixing for different statistics of the lensing potential. The left plot shows the Gibbs chain for  $\phi(x)$  where  $x = (9.9^\circ, 13.2^\circ)$  and  $x = (1.6^\circ, 1.6^\circ)$  in the same degree coordinates given in Figures 6 and 5. The right plot shows the real and imaginary parts of  $\phi_l$  where the frequency vector  $l$  is set to  $(126.56, 63.28)$ . Each dashed line represents the corresponding simulation truth parameters. These plots suggest that the Gibbs chain is mixing well and that the correlation length scale is small enough so that thinning by 100 is sufficient to yield relatively uncorrelated samples.

## 8. CONCLUDING REMARKS

In this paper we construct a prototype algorithm which establishes that it is possible to construct a fast Gibbs sampler of the Bayesian posterior for the unknown lensing potential and the de-noised CMB temperature map. This prototype solves one of the fundamental obstacles in a Gibbs implementation of the Bayesian lensing problem: the



**Figure 9.** This plot summarizes the correlation, in  $\Delta l$  wavenumber bins, between the simulation truth and their corresponding posterior samples. In particular, samples were generated from  $\frac{1}{c} \sum_{l \in \Delta l} \phi_l^{\text{sim}} \phi_l^*$  and  $\frac{1}{c} \sum_{l \in \Delta l} T_l^{\text{sim}} T_l^*$  where  $\Delta l$  is a frequency wavenumber bin,  $\phi_l$  and  $T_l$  are the simulation truth,  $\phi_l^{\text{sim}}$  and  $T_l^{\text{sim}}$  is sampled from the Gibbs algorithm presented here and  $c$  is a normalization constant which transforms to a correlation scale. Recall that  $|l|_{\max}$  for the lensing potential is  $\sim 460$  and is  $\sim 2700$  for the unlensed CMB, which explains why the correlation for  $\phi_l$  only extends to 460. Notice that the plotted correlations trend to 0 for larger wavenumber. This is what one would expect from the Bayesian posterior. Indeed, the signal-to-noise ratio is lower at larger wavenumber. This causes the Bayesian posterior to revert to the prior, which will be uncorrelated with the simulation truth.



**Figure 10.** This plot illustrates the Gibbs chain correlation length scale and the speed of mixing for different statistics of the lensing potential. The left plot shows the Gibbs chain for  $\phi(x)$  where  $x = (9.9^\circ, 13.2^\circ)$  and  $x = (1.6^\circ, 1.6^\circ)$  in the same degree coordinates given in Figures 6 and 5. The right plot shows the real and imaginary parts of  $\phi_l$  where the frequency vector  $l$  is set to  $(126.56, 63.28)$ . Each dashed line represents the corresponding simulation truth parameters. Recall that the algorithm is initialized with a zero lensing potential.

naive parameterization  $(T, \phi)$  is extremely slow. We identify the ancillary and sufficient parametrization duality for this problem and notice that the slowness of the Gibbs chain for the ancillary parametrization  $(T, \phi)$  translates to a fast chain for the sufficient parametrization  $(\tilde{T}, \phi)$ . This observation is one of the main contributions of this paper. The second contribution is the use of the anti-lensing approximation along with Claim 1 which makes feasible the development of a Hamiltonian Markov Chain algorithm for sampling from  $P(\phi|\tilde{T})$ . Without the Fourier transform characterization in Claim 1 the HMC would be computational prohibitive. The third contribution of this paper is to recognize that a new messenger algorithm Elsner and Wandelt (2013); Jasche and Lavaux (2015) can be adapted for high resolution conditional Gaussian sampling under the irregular sampling scenario needed for  $P(\tilde{T}|\phi, \text{data})$ .

Notice that both sampling steps  $P(\phi|\tilde{T})$  and  $P(\tilde{T}|\phi, \text{data})$  in our algorithm utilize a high resolution embedding for  $\tilde{T}$ . This high resolution embedding is most likely the dominant bottleneck for scaling the current prototype implementation presented here. In this paragraph we discuss what is needed to avoid using this embedding for scaling up this algorithm. When sampling from the conditional  $P(\phi|\tilde{T})$ , the main challenge is to compute  $A^q(x)$  and  $B(x)$ , as defined in Claim 1. Within the HMC algorithm, a proposed lensing potential  $\phi$  changes iteratively. At each iteration one requires a new computation of  $A^q(x)$  and  $B(x)$ . In our prototype, a spline interpolation performs the task of fast anti-lensing required

for  $A^q(x)$  and  $B(x)$ . It is an open problem how to compute this fast anti-lensing without the need for a high resolution  $\tilde{T}$ . Simulating from  $P(\tilde{T}|\phi, \text{data})$  also requires a high resolution embedding in our prototype. This simply expresses the fact that given  $\phi$  the field  $\tilde{T}$  is modeled as a non-stationary random field. To circumvent this difficulty we transform to lensed coordinates as illustrated in Figure 3. The challenge when avoiding this high resolution embedding, then, is to directly generate conditional simulations of the non-stationary  $\tilde{T}$  given  $\text{data}(x) = \tilde{T}(x) + n(x)$  and the lensing potential  $\phi(x)$ .

## ACKNOWLEDGMENTS

The authors wish to thank an anonymous referee for constructive comments and suggestions. BW acknowledges funding through his Chaire d’Excellence from the Agence Nationale de la Recherche (ANR-10-CEXC-004-01). This work has been done within the Labex ILP (reference ANR-10-LABX-63) part of the Idex SUPER, and received financial state aid managed by the Agence Nationale de la Recherche, as part of the programme Investissements d’avenir under the reference ANR-11-IDEX-0004-02. EA acknowledges grant support from NSF CAREER DMS-1252795.

## REFERENCES

- S. Das et al., Physical Review Letters **107**, 021301 (2011).  
 V. Engelen et al., The Astrophysical Journal **756**, 142 (2012).  
 Planck Collaboration, A&A **571**, A17 (2014), 1303.5077.  
 The Polarbear Collaboration: P. A. R. Ade, Y. Akiba, A. E. Anthony, K. Arnold, M. Atlas, D. Barron, D. Boettger, J. Borrill, S. Chapman, Y. Chinone, et al., ApJ **794**, 171 (2014), 1403.2369.  
 Planck Collaboration, ArXiv e-prints (2015), 1502.01591.  
 W. Hu, The Astrophysical Journal Letters **557**, L79 (2001).  
 W. Hu and T. Okamoto, The Astrophysical Journal **574**, 566 (2002).  
 C. M. Hirata and U. c. v. Seljak, Phys. Rev. D **67**, 043001 (2003a), URL <http://link.aps.org/doi/10.1103/PhysRevD.67.043001>.  
 C. M. Hirata and U. c. v. Seljak, Phys. Rev. D **68**, 083002 (2003b), URL <http://link.aps.org/doi/10.1103/PhysRevD.68.083002>.  
 A. Lewis and A. Challinor, Physics Reports **429**, 1 (2006), ISSN 0370-1573, URL <http://www.sciencedirect.com/science/article/pii/S0370157306000810>.  
 J. Bezanson, S. Karpinski, V. Shah, and A. Edelman, arXiv preprint arXiv:1209.5145 (2012).  
 S. Dodelson, *Modern cosmology* (Academic press, 2003).  
 G. Roberts, O. Papaspiliopoulos, and M. Sköld, in *Bayesian Statistics 7: Proceedings of the Seventh Valencia International Meeting* (Oxford University Press, USA, 2003), p. 307.  
 A. Gelfand, S. Sahu, and B. Carlin, Biometrika **82**, 479 (1995).  
 O. Papaspiliopoulos and G. Roberts, The Annals of Statistics pp. 95–117 (2008).  
 O. Papaspiliopoulos, G. Roberts, and M. Sköld, Statistical Science pp. 59–73 (2007).  
 Y. Yu and X.-L. Meng, Journal of Computational and Graphical Statistics **20**, 531 (2011).  
 F. Elsner and B. D. Wandelt, A&A **549**, A111 (2013), 1210.4931.  
 J. Jasche and G. Lavaux, MNRAS **447**, 1204 (2015), 1402.1763.  
 R. M. Neal, Handbook of Markov Chain Monte Carlo **2** (2011).  
 A. Hajian, Phys. Rev. D **75**, 083525 (2007), URL <http://link.aps.org/doi/10.1103/PhysRevD.75.083525>.  
 J. F. Taylor, M. A. J. Ashdown, and M. P. Hobson, Monthly Notices of the Royal Astronomical Society **389**, 1284 (2008).  
 F. Elsner and B. Wandelt, The Astrophysical Journal **724**, 1262 (2010).  
 J. Jasche, F. S. Kitaura, C. Li, and T. A. Enßlin, Monthly Notices of the Royal Astronomical Society **409**, 355 (2010), 0911.2498.  
 J. Jasche and B. Wandelt, Monthly Notices of the Royal Astronomical Society **425**, 1042 (2012), 1106.2757.  
 J. Jasche and B. Wandelt, Monthly Notices of the Royal Astronomical Society p. stt449 (2013a).  
 J. Jasche and B. Wandelt, The Astrophysical Journal **779**, 15 (2013b).  
 B. D. Wandelt, D. L. Larson, and A. Lakshminarayanan, Phys. Rev. D **70**, 083511 (2004), astro-ph/0310080.  
 H. K. Eriksen, I. J. O’Dwyer, J. B. Jewell, B. D. Wandelt, D. L. Larson, K. M. Górska, S. Levin, A. J. Banday, and P. B. Lilje, ApJS **155**, 227 (2004), astro-ph/0407028.  
 K. M. Smith, O. Zahn, and O. Doré, Phys. Rev. D **76**, 043510 (2007), 0705.3980.  
 D. S. Seljebotn, K.-A. Mardal, J. B. Jewell, H. K. Eriksen, and P. Bull, ApJS **210**, 24 (2014), 1308.5299.  
 A. Lewis, A. Challinor, and A. Lasenby, Astrophys. J. **538**, 473 (2000), astro-ph/9911177.

## APPENDIX

Before we proceed to the proofs we briefly discuss notation. First, we do not differentiate, notationally, a random field with periodic boundary conditions on  $(-L/2, L/2]^2$  and the case where  $L \rightarrow \infty$  so that the Fourier series  $\sum_{l \in \frac{2\pi}{L}\mathbb{Z}} e^{ix \cdot l} f_l \frac{2\pi/L}{2\pi}$  converges to the continuous Fourier transform  $\int_{\mathbb{R}^2} e^{ix \cdot l} f_l \frac{dl}{2\pi}$ . For example, at times we will refer to an infinitesimal area element  $dl$  or  $dk$  in Fourier space, which simply equals  $(2\pi/L)^2$  for large  $L$ . In this case,  $\delta_l$  denotes a discrete Dirac delta function which we equate with  $1/dl$  when  $l = 0$  and zero otherwise. Indeed, throughout the paper it will be convenient to use  $\delta_0 \equiv \delta_l|_{l=0}$  to denote the constant  $1/dl$ . Secondly, for any function  $f(x)$  let  $f^\phi(x) = f(x - \nabla\phi(x))$  denote anti-lensing of  $f$  and  $f_l^\phi$  denote the Fourier transform of  $f^\phi(x)$ .

*Proof of Claim 1.* Since  $\tilde{T}$  is sufficient for the unknown  $\phi$  we have that

$$P(\phi|\tilde{T}, \text{data}) = P(\phi|\tilde{T}) \propto P(\tilde{T}|\phi)P(\phi).$$

Since  $\phi(x)$  is an isotropic random field with spectral density  $C_l^{\phi\phi}$  we have that  $E(\phi_l \phi_{l'}^*) = \delta_{l-l'} C_l^{\phi\phi}$ . Therefore  $E(\phi_l \phi_l^*) = \delta_0 C_l^{\phi\phi}$  and  $E(\phi_l \phi_l) = 0$  implies that the random variables  $\text{re}\phi_l$ ,  $\text{im}\phi_l$  are independent  $\mathcal{N}(0, \frac{1}{2}\delta_0 C_l^{\phi\phi})$  for each fixed  $l$ . Moreovoer  $\phi(x)$  takes values in  $\mathbb{R}$  so that  $\phi_l = \phi_{-l}^*$ . This implies that  $\phi_l$  are independent random variables over all  $l$  which are restricted to the Hermitian half of the Fourier grid, denoted  $\mathbb{H}$  here. In particular, if we exclude

the zero frequency  $l = 0$  we get

$$\log P(\phi) - c_1 = -\frac{1}{2} \sum_{k \in \mathbb{H} \setminus \{0\}} \left[ \frac{(\operatorname{re}\phi_k)^2}{\frac{1}{2}\delta_0 C_k^{\phi\phi}} + \frac{(\operatorname{im}\phi_k)^2}{\frac{1}{2}\delta_0 C_k^{\phi\phi}} \right] = -\frac{1}{2} \int_{\mathbb{R}^2} \frac{|\phi_k|^2}{C_k^{\phi\phi}} dk \quad (1)$$

$$\log P(\tilde{T}|\phi) - c_2 = -\frac{1}{2} \sum_{k \in \mathbb{H} \setminus \{0\}} \left[ \frac{(\operatorname{re}\tilde{T}_k^\phi)^2}{\frac{1}{2}\delta_0 C_k^{TT}} + \frac{(\operatorname{im}\tilde{T}_k^\phi)^2}{\frac{1}{2}\delta_0 C_k^{TT}} \right] = -\frac{1}{2} \int_{\mathbb{R}^2} \frac{|\tilde{T}_k^\phi|^2}{C_k^{TT}} dk \quad (2)$$

where  $c_1$  and  $c_2$  are constants and  $\tilde{T}^\phi(x) \equiv \tilde{T}(x - \nabla\phi(x))$ .

Taking derivatives in (1) gives

$$\frac{\partial}{\partial\phi_l} \log P(\phi) = -2(dl) \frac{\phi_l}{C_l^{\phi\phi}}. \quad (3)$$

Taking derivatives in (2) gives

$$\frac{\partial}{\partial \operatorname{re}\phi_l} \log P(\tilde{T}|\phi) = -\operatorname{re} \int_{\mathbb{R}^2} \frac{\partial \tilde{T}_k^\phi}{\partial \operatorname{re}\phi_l} \frac{\tilde{T}_k^{\phi*}}{C_k^{TT}} dk \quad (4)$$

$$\frac{\partial}{\partial \operatorname{im}\phi_l} \log P(\tilde{T}|\phi) = -\operatorname{re} \int_{\mathbb{R}^2} \frac{\partial \tilde{T}_k^\phi}{\partial \operatorname{im}\phi_l} \frac{\tilde{T}_k^{\phi*}}{C_k^{TT}} dk. \quad (5)$$

Taking linear combinations of the two equalities in Lemma 1 below we get

$$\frac{\partial \tilde{T}_k^\phi}{\partial \operatorname{re}\phi_l} = \frac{1}{2} \frac{\partial \tilde{T}_k^\phi}{\partial \phi_l} + \frac{1}{2} \frac{\partial \tilde{T}_k^\phi}{\partial \phi_l^*} = \frac{dk}{2\pi} \sum_{q=1,2} il_q \left\{ [(\nabla^q \tilde{T})^\phi]_{k-l} - [(\nabla^q \tilde{T})^\phi]_{k+l} \right\} \quad (6)$$

$$\frac{\partial \tilde{T}_k^\phi}{\partial \operatorname{im}\phi_l} = \frac{-i}{2} \frac{\partial \tilde{T}_k^\phi}{\partial \phi_l} + \frac{i}{2} \frac{\partial \tilde{T}_k^\phi}{\partial \phi_l^*} = \frac{dk}{2\pi} \sum_{q=1,2} l_q \left\{ -[(\nabla^q \tilde{T})^\phi]_{k-l} - [(\nabla^q \tilde{T})^\phi]_{k+l} \right\}. \quad (7)$$

Now the above two equations establish, by Lemma 2 below, that both integrals  $\int_{\mathbb{R}^2} \frac{\partial \tilde{T}_k^\phi}{\partial \operatorname{re}\phi_l} \frac{\tilde{T}_k^{\phi*}}{C_k^{TT}} dk$  and  $\int_{\mathbb{R}^2} \frac{\partial \tilde{T}_k^\phi}{\partial \operatorname{im}\phi_l} \frac{\tilde{T}_k^{\phi*}}{C_k^{TT}} dk$  are real which implies

$$\begin{aligned} \frac{\partial}{\partial\phi_l} \log P(\tilde{T}|\phi) &= - \int_{\mathbb{R}^2} \frac{\partial \tilde{T}_k^\phi}{\partial \phi_l} \frac{\tilde{T}_k^{\phi*}}{C_k^{TT}} dk \\ &= -\frac{dk}{\pi} \sum_{q=1,2} il_q \int_{\mathbb{R}^2} [(\nabla^q \tilde{T})^\phi]_{k+l} \frac{\tilde{T}_k^{\phi*}}{C_k^{TT}} dk \\ &= -i2(dk) \sum_{q=1,2} l_q \int_{\mathbb{R}^2} [(\nabla^q \tilde{T})^\phi]_{k+l} \frac{\tilde{T}_k^{\phi*}}{C_k^{TT}} \frac{dk}{2\pi} \\ &= -i2(dk) \sum_{q=1,2} l_q \int_{\mathbb{R}^2} e^{-ix \cdot l} A^q(x) B(x) \frac{dx}{2\pi}, \quad \text{by Lemma 3 below} \end{aligned}$$

where  $A^q(x) \equiv (\nabla^q \tilde{T})^\phi(x)$  and  $B_k \equiv (\tilde{T}_k^\phi)^*/C_k^{TT}$ . □

### Lemma 1.

$$\frac{\partial \tilde{T}_k^\phi}{\partial \phi_l} = \frac{dk}{\pi} \sum_{q=1,2} il_q [(\nabla^q \tilde{T})^\phi]_{k+l} \quad (8)$$

$$\frac{\partial \tilde{T}_k^\phi}{\partial \phi_l^*} = \frac{dk}{\pi} \sum_{q=1,2} -il_q [(\nabla^q \tilde{T})^\phi]_{k-l} \quad (9)$$

where  $\nabla^q \tilde{T} \equiv \frac{\partial \tilde{T}}{\partial x_q}$ .

*Proof.* First notice

$$\frac{\partial}{\partial \operatorname{re}\phi_l} \frac{\partial \phi(x)}{\partial x_q} = \int_{\mathbb{R}^2} i k_q e^{ix \cdot k} \frac{\partial \phi_k}{\partial \operatorname{re}\phi_l} \frac{dk}{2\pi} = [il_q e^{ix \cdot l} - il_q e^{-ix \cdot l}] \frac{dk}{2\pi} \quad (10)$$

$$\frac{\partial}{\partial \operatorname{im}\phi_l} \frac{\partial \phi(x)}{\partial x_q} = \int_{\mathbb{R}^2} i k_q e^{ix \cdot k} \frac{\partial \phi_k}{\partial \operatorname{im}\phi_l} \frac{dk}{2\pi} = [-l_q e^{ix \cdot l} - l_q e^{-ix \cdot l}] \frac{dk}{2\pi}. \quad (11)$$

This implies

$$\begin{aligned} \frac{\partial \tilde{T}_k^\phi}{\partial \phi_l} &= \frac{\partial}{\partial \phi_l} \int_{\mathbb{R}^2} e^{-ix \cdot k} \tilde{T}(x - \nabla \phi(x)) \frac{dx}{2\pi} \\ &= \sum_{q=1,2} \int_{\mathbb{R}^2} e^{-ix \cdot k} \nabla^q \tilde{T}(x - \nabla \phi(x)) \left[ -\frac{\partial}{\partial \operatorname{re}\phi_l} \frac{\partial \phi(x)}{\partial x_q} - i \frac{\partial}{\partial \operatorname{im}\phi_l} \frac{\partial \phi(x)}{\partial x_q} \right] \frac{dx}{2\pi} \\ &= \sum_{q=1,2} \frac{il_q dk}{\pi} \int_{\mathbb{R}^2} e^{-ix \cdot (k+l)} \nabla^q \tilde{T}(x - \nabla \phi(x)) \frac{dx}{2\pi}, \quad \text{by (10) and (11)} \\ &= \sum_{q=1,2} \frac{il_q dk}{\pi} [(\nabla^q \tilde{T})^\phi]_{k+l} \end{aligned} \quad (12)$$

Similarly

$$\frac{\partial \tilde{T}_k^\phi}{\partial \phi_l^*} = \sum_{q=1,2} \frac{-il_q dk}{\pi} [(\nabla^q \tilde{T})^\phi]_{k-l}. \quad (13)$$

□

**Lemma 2.** If  $A(x)$  and  $B(x)$  are real scalar fields then the two integrals,  $\int_{\mathbb{R}^2} i \{A_{k-l} - A_{k+l}\} B_k^* dk$  and  $\int_{\mathbb{R}^2} \{A_{k-l} + A_{k+l}\} B_k^* dk$ , are both real numbers.

*Proof.* By a simple change of variables it is clear that  $\int_{\mathbb{R}^2} (i \{A_{k-l} - A_{k+l}\} B_k^*)^* dk = \int_{\mathbb{R}^2} i \{A_{k'-l} - A_{k'+l}\} B_{k'}^* dk'$  and  $\int_{\mathbb{R}^2} (\{A_{k-l} + A_{k+l}\} B_k^*)^* dk = \int_{\mathbb{R}^2} \{A_{k'-l} + A_{k'+l}\} B_{k'}^* dk'$ . □

The following lemma is equivalent to the so-called Convolution Theorem. We state it here for reference.

**Lemma 3.** If  $A(x)$  and  $B(x)$  are real scalar fields then  $\int_{\mathbb{R}^2} A_{k+l} B_k^* \frac{dk}{2\pi} = \int_{\mathbb{R}^2} e^{-ix \cdot l} A(x) B(x) \frac{dx}{2\pi}$ .