

# Bias Corrected Estimation of Paleointensity (BiCEP): An improved methodology for obtaining paleointensity estimates

Brendan Cych<sup>1</sup>, Matthias Morzfeld<sup>1</sup>, Lisa Tauxe<sup>1</sup>

<sup>1</sup>University of California, San Diego

## Key Points:

- Empirical evidence suggests that paleointensity estimates for non-ideal specimens are biased.
- BiCEP is a method for estimating paleointensity for ensembles of specimens, correcting for bias
- BiCEP produces accurate results when applied to data where the true field strength is known.

---

Corresponding author: Brendan Cych, [bcych@ucsd.edu](mailto:bcych@ucsd.edu)

## Abstract

The assumptions of paleointensity experiments are violated in many natural and archaeological materials, leading to Arai plots which do not appear linear and yield inaccurate paleointensity estimates, leading to bias in the result. Recently, paleomagnetists have adopted sets of “selection criteria” that exclude specimens with non linear Arai plots from the analysis, but there is little consensus in the paleomagnetic community on which set to use. In this paper, we present a statistical method we call Bias Corrected Estimation of Paleointensity (BiCEP), which assumes that the paleointensity recorded by each specimen is biased away from a true answer by an amount that is dependent a single metric of nonlinearity (the curvature parameter  $\vec{k}$ ) on the Arai plot. We can use this empirical relationship to estimate the recorded paleointensity for a specimen where  $\vec{k} = 0$ , i.e., a perfectly straight line. We apply the BiCEP method to a collection of 30 sites for which the true value of the original field is well constrained. Our method returns accurate estimates of paleointensity, with a higher level of accuracy and precision than the strict CCRIT selection criteria, and with higher accuracy and similar precision to the modified PICRIT03 criteria. The BiCEP method has a significant advantage over using these selection criteria because it achieves these accurate results without excluding large numbers of specimens from the analysis.

## Plain Language Summary

Paleomagnetists perform experiments on rocks and pottery sherds to estimate the strength of the ancient Earth’s magnetic field (the paleointensity) through time. These make assumptions which are frequently violated, leading to bias. Quantitative metrics (selection criteria) attempt to screen out ‘bad’ data. If a particular experiment fails the criteria, the results are ignored. However, there a lack of agreement as to which set of criteria are the most important and what is considered a failure. One of these criteria quantifies the deviation from the fundamental assumption of linearity of between the ancient and laboratory magnetizations. We present a new Bayesian method called Bias Corrected Estimation of Paleointensity (BiCEP), in which we assume that the estimated paleointensity depends on this deviation. We can then use this dependency to correct the paleointensity made on an ensemble of specimens with differing deviations from ideal behavior. This allows us to calculate accurate estimates of the ancient magnetic field, without ignoring results from non-ideal specimens. We test BiCEP on paleomagnetic data for which we the original field strength is well constrained. BiCEP recovers the field strength as precisely and slightly more accurately than the best performing set of selection criteria we tested.

## 1 Introduction

Estimates of the strength of the ancient Earth’s magnetic field are currently made by performing experiments that compare the natural remanent magnetization (NRM) acquired by a specimen while cooling in the Earth’s field, to a remanence known as thermal remanent magnetization (TRM) acquired by the specimen while cooling in a known laboratory field. Such experiments include the Königsberger-Thellier-Thellier (KTT) family of experiments (Königsberger, 1938; Thellier & Thellier, 1959), the Shaw family of experiments (Shaw, 1974), and the multi-specimen family of experiments (Hoffman et al., 1989), among others. All of these experimental families make assumptions about the relationship between the magnetic field and the remanent magnetization which may or may not be applicable (see the review by Tauxe & Yamazaki, 2015). In this paper, we will focus on the KTT family of experiments.

KTT family of experiments involve a double heating protocol in which a specimen is heated two or more times to a series of temperatures up to the Curie Temperature. At each temperature, the specimen is cooled in two different fields. This has the effect

of replacing the NRM with a TRM acquired in a known laboratory field. Data from KTT-type experiments are normally represented by the Arai diagram (Nagata et al., 1963), which plots the NRM magnetization remaining at each temperature step against the magnetization imparted in the laboratory (often referred to as partial TRM or pTRM). The ratio of these two magnetizations, as represented by the slope of the best fitting line to the Arai plot data, is generally taken to be the ratio of the two magnetizing fields (ancient,  $B_{anc}$  and laboratory,  $B_{lab}$ ).

KTT-type experiments rely on several assumptions which are frequently violated in paleointensity experiments. These include thermochemical alteration of specimens which may lead to the production of new magnetic minerals, and an assumption known as reciprocity, which requires that the blocking temperature (the temperature below which grains retain their magnetization after an external field is removed) is the same as the unblocking temperature (the temperature above which grains equilibrate with the external field).

The reciprocity assumption of Thellier and Thellier (1959) is fundamental to Néel’s theory for uniaxial single domain grains (Néel, 1949). Néel theory assumes that the electronic spins within magnetic grains are fully aligned, and that the alignment is in one of two directions along an energetically favorable ‘easy’ axis. In zero field, there is no preference for either direction, but in the presence of a field there is a slight preference for the direction along the easy axis with the smallest angle to the applied field. If the reciprocity assumption is met, then the energy required for the magnetization to change directions along the easy axis is always the same regardless of whether the specimen is cooled from higher temperature (blocking) or heated from room temperature (unblocking) and the two temperatures are identical.

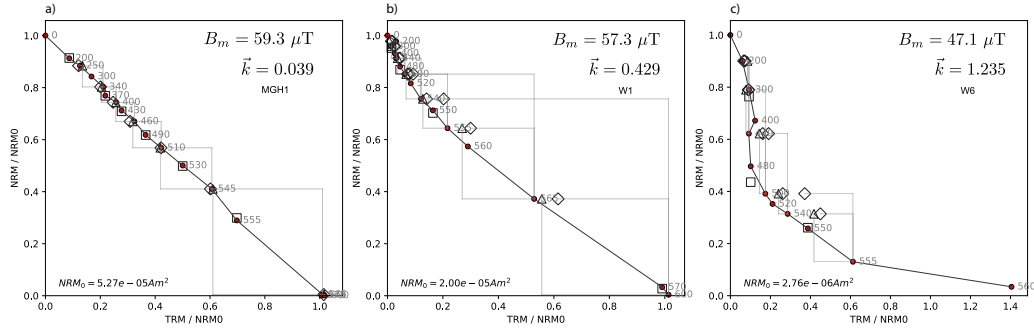
By assuming that electronic spins within magnetic grains are fully aligned, Néel theory fails to take into account a term in the magnetic energy of grains which causes deviations from full alignment, resulting in structures such as the vortex state of, e.g., Williams and Dunlop (1989). Although this effect is present in nearly all magnetic grains, it is insignificant over short length scales (10s of nm) and so uniaxial single domain theory may be a reasonable approximation for smaller, elongate grains. Specimens in paleointensity experiments contain mixtures of grains with different sizes and shapes and a specimen used for paleointensity is likely to include grains for which the applicability of single domain theory does not hold.

Failure of reciprocity and other fundamental assumptions embedded in the KTT family of experiments (laid out by e.g., Thellier & Thellier, 1959) provides a challenge for those analyzing paleointensity data. Paleomagnetists generally use a set of selection criteria which reject an intensity result if the NRM and pTRM data behave in a way which deviate from single domain theory (linear on the Arai plot, see Figure 1a) by more than some arbitrarily chosen threshold value. This is because data that contain a large proportion of non single domain-like grains or which otherwise violate the assumptions of the experiment are likely to give biased results (Tauxe et al., 2021). Selection criteria generally operate in a binary way, with specimens either being ‘accepted’ or ‘rejected’ from the estimation of the site mean, where ‘site’ is the collection of specimens assumed to have cooled in identical external magnetic fields (say, a lava flow or ceramic fragment).

Figure 1 gives a demonstration of biased results in specimens from prepared magnetite powders of increasing grain size that were magnetized in a 60  $\mu$ T field (Krása et al., 2003). If all assumptions of Thellier and Thellier (1959) were obeyed, we would expect the best fitting lines to data on Arai plots to give a range of values distributed closely about a mean of 60  $\mu$ T. As the grain size of the powder increases, the Arai plot becomes more curved and the best fitting line to the Arai plot yields a progressively lower intensity estimate. As all the paleointensities estimated from the curved plots are below the expected value, the estimate for the ensemble can be biased, with the high temperature segment having an even lower mean value, and the low temperature segment having a

high mean value. The data of Tauxe et al. (2021) also demonstrate downward curved Arai plots in natural samples are biased so this problem may effect many of the results compiled in paleointensity databases like the MagIC database (Tauxe et al., 2016) or PINT (Biggin, 2010).

The curvature of an Arai plot can be quantified using the curvature criterion ( $\vec{k}$ ) of Paterson (2011) (see also Paterson et al., 2014). Curvature is calculated using the reciprocal of the radius of a circle fit to scaled Arai plot data (see Section 2.2.1). While there is no theoretical basis for a circular fit (as opposed to the linear fit, which is firmly rooted in Néel theory), it is a useful approximation that we will use in this paper.



**Figure 1.** Arai plots from prepared magnetite powders given a TRM in a 60  $\mu\text{T}$  field (Krása et al., 2003). The curvature criterion,  $\vec{k}$  (Paterson, 2011) and specimen level paleointensity estimate  $B_m$  estimated from fitting a line to the entire Arai plot are plotted on the figure as text. The grain size of the magnetite powders increases from left to right. The coarser grains have non ideal domain state, leading to curved Arai plots and estimates of paleointensity which are biased to lower values than the expected 60  $\mu\text{T}$ . a) Nominal grain size of 23 nm. b) Mean grain size of 70 nm. c) Mean grain size of 12.1  $\mu\text{m}$ .

The practice of using binary (pass/fail) selection criteria is problematic for many reasons. Paleomagnetic specimens generally contain magnetic carriers which span a range of grain sizes and may or may not conform to the assumptions of the method. In addition, micromagnetic simulations (e.g., Williams and Dunlop (1989); Nagy et al. (2017)) demonstrate that the change in magnetic domain state with grain size is a continuum, and so one individual grain’s behaviour may be more or less ideal than any other’s. With binary pass/fail criteria, the distinction between ‘good’ and ‘bad’ data must be assessed with an arbitrary threshold value, which does not reflect the range of behaviors within both groups. Consequently there are a large number of selection criteria in common use (over 40 in Paterson et al., 2014), most of which have some empirical rationale, but there is little agreement on which set to use or their threshold values.

In this paper, we describe a new approach for paleointensity estimation that treats the quality of paleointensity data as a continuum. We assume that paleointensities become more biased as specimens’ magnetic behaviors become more non-ideal and their Arai plots become less linear. By allowing the data interpretation for specimens to be based on the shape of their Arai plots, we are able to obtain unbiased estimates of paleointensity without the need for many specimen level (binary) selection criteria. We call this method the ‘Bias Corrected Estimation of Paleointensity’ or BiCEP. In the next section, we develop a Bayesian approach to obtain accurate paleointensity estimates with realistic uncertainties, using  $k$  as a metric of bias, and show how to combine data at the site level. In Section 3 we compare results from the BiCEP method to those of more tra-

ditional selection criteria based approaches. We discuss the results in Section 4 and summarize our conclusions in Section 5. Accompanying this paper, we release a Graphical User Interface (GUI) which can apply the BiCEP method to MagIC formatted data. Links and instructions on how to access the code can be found in Appendix 6.3.

## 2 Methods

### 2.1 Accounting for bias in paleointensity experiments

Paleomagnetists determine the paleointensity for a site by performing a Thellier-type double heating experiment on multiple specimens from that site. According to the theory for single domain grains (assuming no alteration of the specimen during heating), the ratio of NRM lost to pTRM gained is the ratio of the ancient field to the laboratory field. If the specimen conforms to theory, the Arai plot data will fall along a line the slope of which is equal to the ratio of ancient to the laboratory field (see Figure 1a).

We expect that the field strength predicted by the slope of the line on the Arai plot for each specimen (here called  $B_m$ ) will be distributed about the true (expected) ancient field ( $B_{exp}$ ) at the site with a Gaussian distribution. However, rarely do a set of specimens from a site all produce linear Arai plots that are easily interpretable. For example, interpretation of data from specimens with magnetic grains exhibiting non single domain magnetic domain states produce non-linear Arai plots which violate the assumptions of the method (e.g., Dunlop & Özdemir, 2001). Fitting lines to the data on such Arai plots often produces estimates of paleointensity which are biased (see Figure 1c, Krása et al., 2003), which in turn would bias site level estimates.

Paleomagnetists generally deal with non-ideal data by using certain quantitative criteria chosen to eliminate results suffering from one or more pathologies (Paterson et al., 2014). If a particular criterion calculated for a specimen exceeds some threshold value, then the specimen is excluded from the analysis. In this paper, we present an alternative approach in which we allow for specimens to behave in a non-ideal (non-linear) fashion when considering how specimen intensity estimates are distributed about a site mean and weight the contribution of individual specimen estimates according to linearity. Under such a scheme, we start by predicting a bias for each specimen, and the specimens with the smallest predicted bias most strongly determine the paleointensity at that site. In this way, biased specimens do not strongly affect our site intensity estimate, as they are down-weighted, yet provide useful constraints on the uncertainty.

1991-1992 Eruption Site	10.1029/2005GC001141	lava flow	9.8	-104.3	1991	36.2	53
hw108	10.1016/j.pepi.2014.12.007	lava flow	19.9	-155.9	1859	39.3	23
hw123	10.1016/j.pepi.2014.12.007	lava flow	19.1	-155.7	1907	37.7	12
hw126	10.1016/j.pepi.2014.12.007	lava flow	19.7	-155.5	1935	36.4	13
hw128	10.1016/j.pepi.2014.12.007	lava flow	19.3	-155.9	1950	36.2	26
hw201	10.1016/j.pepi.2014.12.007	lava flow	19.4	-155.0	1990	35.2	12
hw226	10.1016/j.pepi.2014.12.007	lava flow	19.6	-155.5	1843	39.9	11
hw241	10.1016/j.pepi.2014.12.007	lava flow	19.5	-155.8	1960	36.0	18
BR06	10.1016/j.pepi.2007.10.002	brick	60.1	24.9	1906	49.7	3
P	10.1029/2010JB007844, 2011	lava flow	19.3	-102.1	1943	44.6	36
VM	10.1029/2010JB007844, 2013	lava flow	40.8	14.5	1944	43.8	18
BBQ	10.1029/93jb01160	submarine lava flow	9.8	-104.3	1990	36.2	12
rs25	10.1016/j.epsl.2009.12.022	synthetic	N/A	N/A	N/A	30.0	5
rs26	10.1016/j.epsl.2009.12.022	synthetic	N/A	N/A	N/A	60.0	5
rs27	10.1016/j.epsl.2009.12.022	synthetic	N/A	N/A	N/A	90.0	10
remag-rs61	10.1016/j.epsl.2011.08.024	synthetic	N/A	N/A	N/A	40.0	6
remag-rs62	10.1016/j.epsl.2011.08.025	synthetic	N/A	N/A	N/A	60.0	6
remag-rs63	10.1016/j.epsl.2011.08.024	synthetic	N/A	N/A	N/A	80.0	5
remag-rs78	10.1016/j.epsl.2011.08.025	synthetic	N/A	N/A	N/A	20.0	4
kf	10.1111/j.1365-246X.2012.05412.x	lava flow	65.7	-16.8	1984	52.0	3
Hawaii 1960 Flow	10.1046/j.1365-246X.2003.01909.x	lava flow	19.5	-155.8	1960	36.0	22
SW	10.1016/j.pepi.2008.03.006	lava flow	31.6	-130.6	1946	46.4	19
TS	10.1016/j.pepi.2008.03.006	lava flow	31.6	-130.6	1914	47.8	53
ET1	10.1016/j.epsl.2007.03.017	basaltic lava	37.8	15.0	1950	43.3	3
ET2	10.1016/j.epsl.2007.03.017	basaltic lava	37.8	15.0	1979	44.1	2
ET3	10.1016/j.epsl.2007.03.017	basaltic lava	37.8	15.0	1983	44.2	4
Synthetic60	10.1016/S1474-7065(03)00122-0	synthetic	N/A	N/A	N/A	60.0	7
LV	10.1029/2009JB006475	Lithic Clasts	-23.4	67.7	1993	24.0	45
MSH	10.1029/2009JB006475	Lithic Clasts	46.2	-122.2	1980	55.6	19
FreshTRM	10.1029/2018GC007946	remagnetized/synthetic	N/A	N/A	N/A	70.0	24

**Table 1.** Table of sites used for analysis in this study, including original study locations, latitude, longitude and year of magnetization (where applicable), expected field at that location ( $B_{exp}$ ), number of specimens used for analysis at that site  $M$ . Lat.: site latitude ( $^{\circ}$ N). Long. site longitude ( $^{\circ}$ E. N/A: Not Applicable (Synthetic)).

To predict the amount of bias a specimen is likely to have, we require a proxy for bias in paleointensity experiments. For this we use the curvature criterion  $\vec{k}$  of Paterson (2011) (see Section 2.2.1). There are several reasons that make this criterion a useful proxy for bias in paleointensity experiments:

- Specimens that are highly linear have, by definition, low values for  $|\vec{k}|$  and will generally give unbiased paleointensity estimates (e.g., Cromwell et al., 2015).
- By contrast, specimens with higher  $|\vec{k}|$  yield biased paleointensities, with the magnitude of the bias generally increasing with the magnitude of  $|\vec{k}|$  (e.g., Tauxe et al., 2021).
- $|\vec{k}|$  has an empirical correlation with magnetic grain size (Paterson, 2011).

To predict bias, we can use a method by which we minimize the misfit to a model assuming that  $B_m$  is linearly related to  $\vec{k}$  for all specimens. In other words, we say that:

$$B_m = B_{exp} + c\vec{k}_m + \epsilon \quad (1)$$

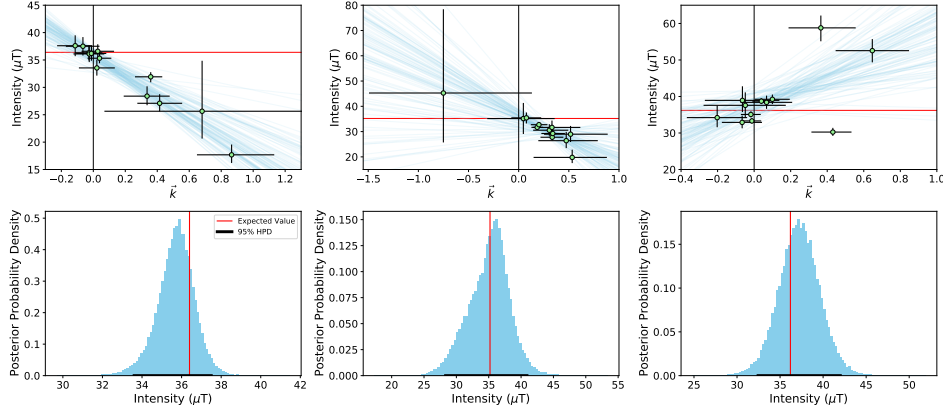
where  $m$  is an index reflecting the specimen number,  $\epsilon$  is an error term and  $B_{exp}$  is the true value of  $B$ . Effectively, our model just becomes a linear fit between the specimen estimate  $B_m$  and  $\vec{k}$ , the y-intercept of which is the true value of the field  $B_{exp}$  and  $c$  is a slope constant. While there is no theoretical justification (yet) for why  $B_m$  would be related to  $\vec{k}_m$ , although it has been observed empirically (Figure 1 and Tauxe et al. (2021)), a linear model is the simplest one to relate the two. We demonstrate in Section 3.3 that more complex models with a quadratic and cubic fit relating  $B_m$  to  $\vec{k}_m$  perform worse than the linear model when predicting the paleointensity for sites for which the paleointensity is well constrained (historical lava flows or laboratory remanences).

Arai plot curvature is not the sole cause of bias in paleointensity experiments. In some cases, specimens with Arai plots which do not have high  $|\vec{k}|$  but are still non linear (e.g., ‘zig-zagged’ as in, e.g., Yu et al., 2004), may still cause bias in paleointensity experiments. To counteract this, we use a Bayesian method of calculating  $\vec{k}_m$  and  $B_m$  which provides an uncertainty for both of these parameters. The benefit of this approach is that specimens whose Arai plots are not well fit by a line or an elliptical arc have less influence on the linear fit. Therefore, the specimens with the lowest uncertainty in  $\vec{k}$  are generally the most linear, and will have the most influence on the linear fit. Yet, for each specimen, there is a trade off between minimizing the circle fit at a specimen level and the linear fit between  $B_m$  and  $\vec{k}$  for specimens from the same site, an issue we will deal with in Section 2.2.3.

Figure 2 shows results from our method (detailed in Section 2.2) applied to several sites for which the true value of  $B_{anc}$  (here,  $B_{exp}$ ) is either calculated from the International Geomagnetic Reference Field (IGRF, Thébault et al., 2015) for historical flows, or known as the NRM is a laboratory TRM imparted to the specimens. Following Equation 1, the uncertainty in the intercept value of these linear fits gives us the uncertainty for our site value of  $B_{anc}$ . In this way, we can obtain an unbiased estimate of  $B_{anc}$  without relying on arbitrary binary (accept/reject) criteria to exclude specimen results.

In the following, we detail how the specimen level circle fit  $\vec{k}$  and site level paleointensity for unknown values for  $B$  (here called  $B_{anc}$ ) can be calculated. We then compare the efficacy of several different versions of our model to classical selection criteria. We do this using a data compilation from 30 sites updated from Paterson et al. (2014) and Tauxe et al. (2016) for which  $B_{exp}$  is well constrained (see Table 1 for details concerning the original publications of the data).





**Figure 2.** Example of results from the BiCEP method for several sites used as examples in this study. Lines (in blue) are fit to the values of  $B_m$  and  $\vec{k}$  for each specimen (blue dots, with uncertainties as black lines). The values of linear fits at  $\vec{k} = 0$  (blue histograms) provide an unbiased estimate of the expected paleointensity value at the site from the known field (red lines). a,d) hw126. b,e) hw201. c,f) BBQ. See Table 1 for sampling and citation details and Section 3 for comparison with the expected field values,  $B_{est}$ .

## 2.2 Statistical Methodology

### 2.2.1 Estimating curvature

Paterson (2011) proposed a least squares fit of circles in Arai plot data. The parameter  $\vec{k}$  of Paterson (2011) is defined as the reciprocal of the radius of a best-fitting circle through the data. It is positive if the circle center is to the upper right of the Arai plot data (upward facing bow, Figure 3a) and negative if the circle center is below and to the left of the Arai plot data (downward facing bow, Figure 3b).

Before fitting to the Arai plot data, Paterson (2011) scales the pTRMs by the maximum pTRM to ensure that the paleointensity data are independent of the laboratory field. For estimating  $\vec{k}$ , we also subtract the minimum remaining NRM ( $NRM_{min}$ ) for specimens for which full demagnetization has not been completed and we subtract the minimum pTRM ( $pTRM_{min}$ ) for specimens for which the low temperature steps were excluded from the analysis (e.g., because of viscous remanent magnetization). This modified form is termed  $\vec{k}'$  (Paterson et al., 2014).

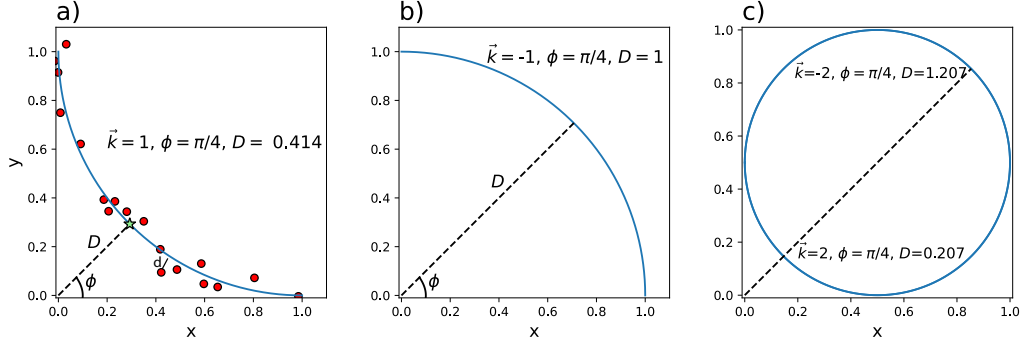
For the BiCEP method, we define two sets of data vectors  $x$  and  $y$ :

$$x_n = \frac{pTRM_n - pTRM_{min}}{pTRM_{max}}, \quad y_n = \frac{NRM_n - NRM_{min}}{NRM_0}, \quad (2)$$

where  $n$  is the index of the data point. Because scaling should be by the total (original) TRM (the NRM), we also exclude specimens whose  $NRM_{min}$  is more than 25% of the initial NRM. This is justified by the assumption that the experimenter did not carry out demagnetization to fully replace the NRM. Then, to fit a circle with center  $x_c, y_c$  and radius  $R$  to the data, we try to minimize the squared perpendicular distance  $d_n^2$  (Figure 3a) of all the  $n$  data points to the circle edge:

$$\sum_{n=1}^N d_n^2 \quad \text{where} \quad d_n^2 = (\sqrt{(x_n - x_c)^2 + (y_n - y_c)^2} - R)^2. \quad (3)$$





**Figure 3.** Example circles with different values for parameters  $\vec{k}$  and  $D$  with the same  $\phi$ , showing how these parameters define a circle. a) Positive  $\vec{k}$ . Red dots are example data, and the green star is the intersection of  $D, \phi$  with the circle edge (see text for definitions).  $d$  is the distance of an individual data point from the best-fit curve (blue). b) Negative  $\vec{k}$ . Note that in this case,  $\phi$  could take any value as the circle center is at the origin, making the definition of  $\phi$  meaningless in this case. c) Example showing how two sets of the parameters  $\vec{k}, \phi, D$  can describe the same circle.

In a total least squares fit, Equation 3 would be our objective function that we would minimize. To fit circles to the Arai plot using a Bayesian method, we use Bayes' formula (Equation 4). This formula allows us to assign a probability distribution to the values of different parameters (in this case,  $\vec{k}_m$  and  $B_m$ ), rather than just finding the 'best' value of the parameters. In a Bayesian context, we can simply assume that the data have some Gaussian noise distribution with some unknown standard deviation  $\sigma$  and apply Bayes' formula (e.g., Gelman et al., 2004):

$$P(\text{Parameters}|\text{Data}) = \frac{P(\text{Data}|\text{Parameters})P(\text{Parameters})}{P(\text{Data})}, \quad (4)$$

where the left hand side is the probability of the parameters given the data and the right hand side is the probability of the data given the parameters times the probability of the parameters, normalized by the probability of the data. In our case, the parameters are  $x_c, y_c, R$  and  $\sigma$  and our data are  $x$  and  $y$  so we rewrite Equation 4 as:

$$P(x_c, y_c, R, \sigma | x, y) = \frac{P(x, y | x_c, y_c, R, \sigma) P(x_c, y_c, R, \sigma)}{P(x, y)}. \quad (5)$$

The term  $P(x, y | x_c, y_c, R, \sigma)$  is known as the "likelihood" and is based on the probability of generating the observed data from a given set of parameters using the assumed Gaussian distribution. The term  $P(x_c, y_c, R, \sigma)$  is known as the "prior" and is a probability distribution for values of  $x_c, y_c, R$  and  $\sigma$  we consider to be reasonable before we see any data. We consider the priors on these parameters to be independent of one another, so we could rewrite this as  $P(x_c)P(y_c)P(R)P(\sigma)$ . The term  $P(x, y)$  is known as the "evidence", and is simply a normalizing constant that makes the "posterior" probability distribution,  $P(x_c, y_c, R, \sigma | x, y)$ , integrate to 1. In our application, we can simplify the relationship by ignoring the normalization. Furthermore, we can say from the definition of the Gaussian distribution that:

$$P(x, y | x_c, y_c, R, \sigma) = \left( \frac{1}{2\pi\sigma^2} \right)^{N/2} \exp \left( -\sum_{n=1}^N \frac{d_n^2}{\sigma^2} \right). \quad (6)$$

Now we have an expression for our posterior probability distribution:

$$P(x_c, y_c, R, \sigma | x, y) \propto \left( \frac{1}{2\pi\sigma^2} \right)^{N/2} \exp \left( \sum_{n=1}^N -\frac{d_n^2}{\sigma^2} \right) P(x_c, y_c, R) P(\sigma). \quad (7)$$

Because the actual noise distribution of the Arai plot data is quite complicated (Paterson et al., 2012), we do not know the value of  $\sigma$ , so we use the uninformative prior  $P(\sigma) \propto \frac{1}{\sigma}$ ; in other words, the smaller  $\sigma$ , the more likely the result. We can then substitute this prior into Equation 7 and integrate out  $\sigma$  to obtain:

$$P(x_c, y_c, R | x, y) \propto \left( \sum_{n=1}^N d_n^2 \right)^{-N/2} P(R, x_c, y_c) \quad (8)$$

where  $N$  is the total number of measurements considered.

The set of parameters  $x_c, y_c$  and  $R$  is not easy to solve for, because Equation 3 has multiple local minima (see Chernov and Lesort (2005) for a more detailed discussion). Consider the simple case of a specimen with a linear Arai plot; in even this simplest case, there are four minima, as both  $R$  and  $x_c, y_c$  will be either positive or negative and very large. To avoid this complexity, we can use instead a change of parameters similar to that of Chernov and Lesort (2005) which Paterson (2011) used as a basis for the circle fitting protocol. Based on this, we define a set of three new parameters which avoid the problem of multiple minima.

Firstly, we require a point on the Arai plot which can be related to a unimodal distribution. We know that linear data will plot along the edge of a circle (the tangent), so if we draw a line from the origin toward the center  $(x_c, y_c)$  (not shown), this will touch the edge of the circle at some distance  $D$  (green star in Figure 3a). The angle to the horizontal of this line we call  $\phi$  and we can directly estimate the  $\vec{k}$  parameter of Paterson (2011) using Equations 9,10,11. We can then establish equations for transforming between these two sets of parameters (see Appendix 6.1 for a more detailed derivation):

$$x_c = \left( D + \frac{1}{\vec{k}} \right) \cos(\phi), \quad (9)$$

$$y_c = \left( D + \frac{1}{\vec{k}} \right) \sin(\phi), \quad (10)$$

$$R = \frac{1}{|\vec{k}|}. \quad (11)$$

Despite this transformation, the circle fitting equation can still have multiple minima, even with  $\vec{k}, D, \phi$  as our parameters, as the line connecting the origin to the horizontal touches the circle edge in two locations (see Figure 3c). However, we can use prior distributions to avoid this.

Chernov and Lesort (2005) define a function of the data  $d_{max}$  to define the region of possible values for  $\vec{k}$ :

$$d_{max} = \max_{i,j} \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2} \quad (12)$$

Additionally, we define distance from the origin to the centroid of the data,  $d_{cent}$ :

$$d_{cent} = \sqrt{\bar{x}^2 + \bar{y}^2} \quad (13)$$

Using this function, we can assume that  $D < 2d_{cent}$  and  $|\vec{k}| < N/d_{max}$  and can define priors for our parameters:

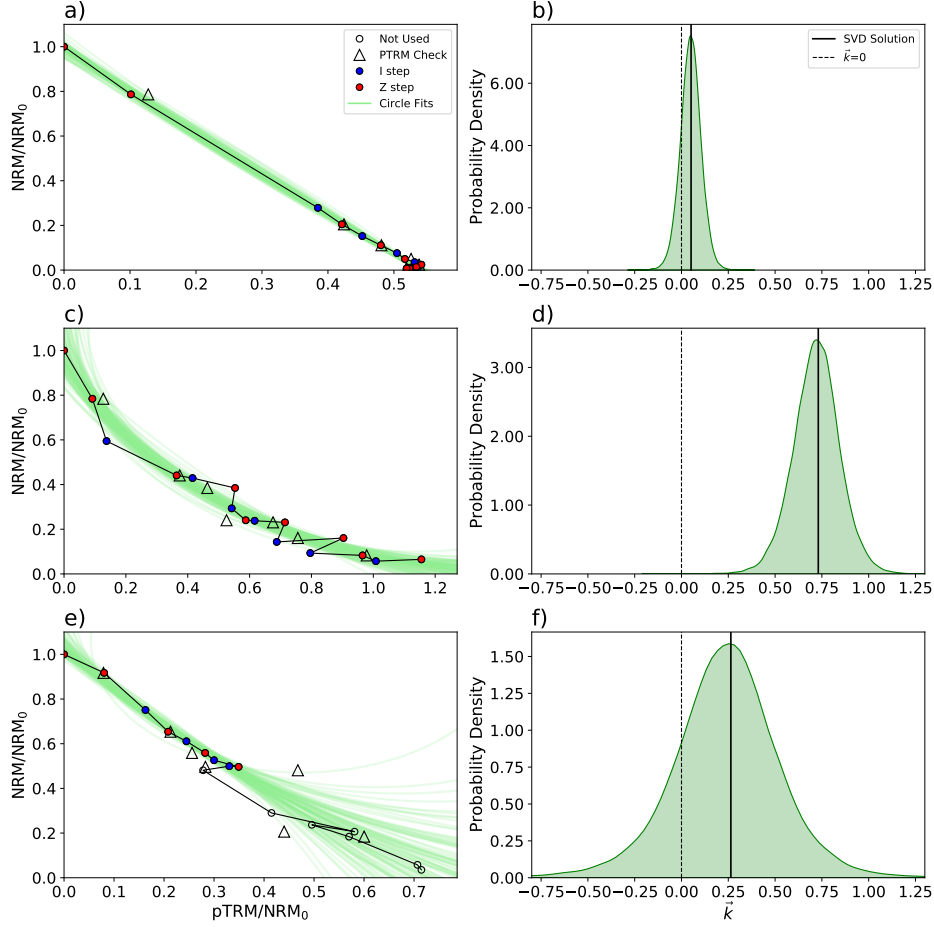
$$P(D) \sim \text{Uniform}(0, 2d_{cent}), \quad (14)$$

$$P(\phi) \sim \text{Uniform}(0, \pi), \quad (15)$$

and

$$P(\vec{k}) \sim \text{Uniform}(-N/d_{max}, N/d_{max}). \quad (16)$$

Using these priors gives us a posterior with a single maximum in most cases, which makes the problem much easier to solve computationally.



**Figure 4.** Examples of circle fits to Arai plots (left column) and approximate probability densities of  $\vec{k}$  (right column). Dashed lines in the left hand plots are the tangents to circles with the median values for  $\phi$  and  $D$ . We use tangents to the circle to get an estimate for  $B_m$  as outlined in Section 2.2.2. Triangles in a), c), e) are repeated lower temperature steps (pTRM checks) that indicate alteration of magnetic minerals during the experiment when offset from the original measurements (red dots). a) Specimen hw126a1. A fit to a straight line yields a precise  $\vec{k}$  distribution with a maximum close to zero (b). c) Specimen hw126a7. A curved Arai plot with a high amount of scatter/zigzagging (left) results in a higher uncertainty in the value of  $\vec{k}$  (d). e) Specimen hw126a6. Arai plot for a specimen that underwent thermochemical alteration at high temperature. A circle fit to just the low temperature steps results in a high uncertainty in the value of  $\vec{k}$  (f).

We can now apply a Bayesian approach to estimate  $\vec{k}$  for all temperature steps for a given specimen  $m$ . It is frequently useful to choose a subset of the temperature steps (e.g., if there is evidence for multiple components of the NRM or heating related alteration, as detected by repeated lower temperature pTRM steps). When using a subset of steps, we scale by the maximum pTRM for all temperature steps and the NRM at room temperature; in this way we can predict the curvature for the part of the Arai plot that is missing. This means that interpretations based on a fraction of the Arai plot will have large uncertainties in the value of  $\vec{k}$ . Therefore, our circle fit can prioritize interpretations using the largest fraction of the NRM.

Figures 4a,c,e show circle fits sampled from the posterior distributions for specimens from site hw126 (site level results shown in Figure 2a). The probability densities of all the  $\vec{k}$  values for each specimen are plotted in Figures 4b,d,f. The plot demonstrates how a straight Arai plot (Figure 4a) produces a narrow posterior about  $\vec{k} = 0$  (Figure 4b), while a curved one (Figure 4c) produces a posterior which does not contain  $\vec{k} = 0$  (Figure 4d). In the example with failed pTRM checks at higher temperatures (offset triangles in Figure 4e), we exclude the data points represented by open circles and use a linear segment with only a portion of the results, the posterior distribution of  $\vec{k}$  has a larger uncertainty on the value, translating to a larger uncertainty in the bias for that specimen.

### 2.2.2 Obtaining a specimen level paleointensity estimate

Analogous to the case in which paleointensity estimates are made using the slope of a fitted line to the Arai plot data, we can obtain a similar “slope” value for a circular arc fit to the data. Consider the case in which the edge of the circle forms an exact line ( $\vec{k}=0$ , see Figure 4a). In this case, the slope of the line can be given by the tangent to the circle at the point where it intersects a line drawn from the origin (0,0) to the circle center (Figure 3a). In other words, the “slope” of the Arai plot can be estimated as  $\cot \phi$ , which gives the tangent to the circle. We can then turn this into an intensity estimate  $B_m$  using the formula:

$$B_m = \frac{B_{lab} \cot(\phi)}{\text{pTRM}_{max}}, \quad (17)$$

where  $B_{lab}$  is the laboratory field used to impart a pTRM to the specimen.

We now have a way of obtaining estimates for  $B_m$  and  $\vec{k}_m$  for each specimen. We use the methodology laid out in Sections 2.2.1 and 2.2.2 to plot the median value of the posterior for these parameters (with error bars) in Figure 5a, and examples of circle fits in Figures 5c, e, g. For specimens with values of  $\vec{k}$  that are approximately 0 (Figure 5g), the  $B_m$  values are quite accurate. There appears to be a bias for specimens with large  $\vec{k}$ , with the amount of bias increasing as  $\vec{k}$  increases. In this example, large positive values of  $\vec{k}$  lead to a large underestimates of  $B_m$  while negative values of  $\vec{k}$  lead to overestimates of  $B_m$  (although small in this example).

### 2.2.3 Obtaining a site level paleointensity estimate

The main problem with the method presented thus far is that we still do not have a way of obtaining an estimate for  $B_{anc}$ , the unknown value at the site level. However, in Figure 5a there appears to be a dependence between  $\vec{k}_m$  and  $B_m$  as suggested earlier, with most of the specimens showing a quasi-linear relationship (the only exception being the point labeled e) whose Arai plot is shown in Figure 5e) and suggests there is a great deal of uncertainty in the value of  $\vec{k}$  itself. Because of this, we can modify our model slightly by imposing the extra restriction that  $B_m$  must be linearly dependent on  $\vec{k}_m$  (with noise) using Equation 1 (substituting  $B_{anc}$  for the unknown value of  $B_{est}$ ).

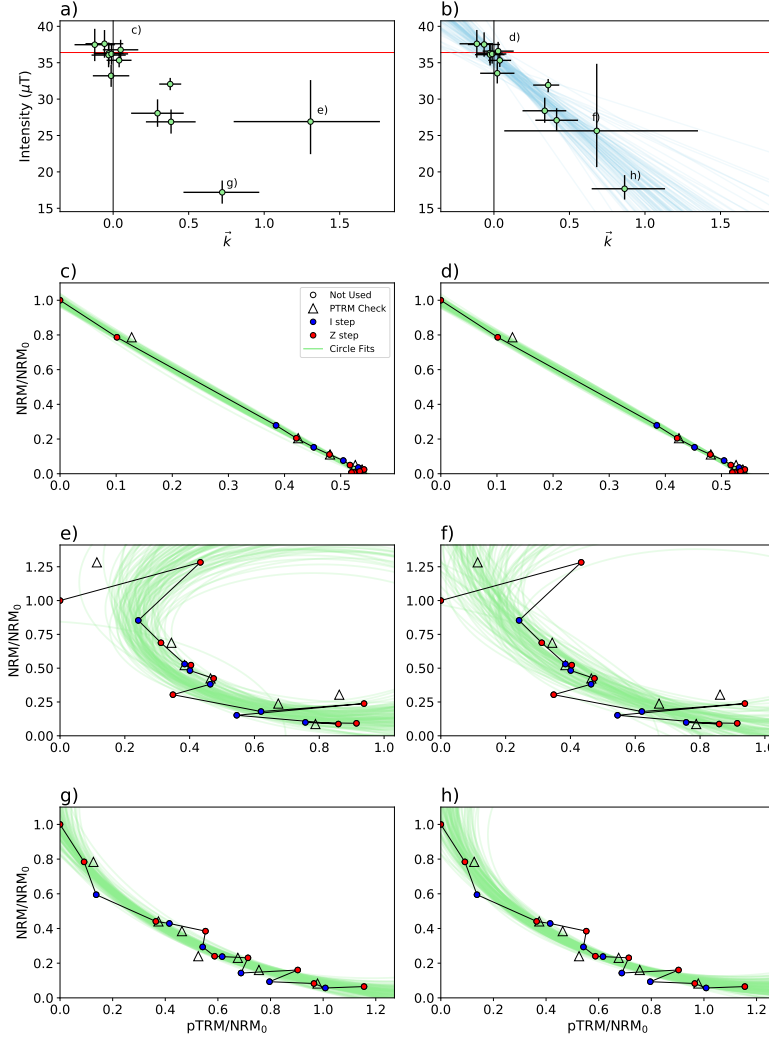
Previous papers have assumed that  $B_{anc}$  for selected specimens follows a Gaussian distribution and we can also make this assumption here. In the following, we will show how this modification can shift results from specimens that are offset from the linear relationship toward the line (as in the point labeled ‘f’ in Figure 5b) and produce models (shown as blue lines) that estimate all of our  $B_m$ . We can then use the resulting models to estimate the probability distribution for  $B_m$  as:

$$P(B_m|k_m, B_{anc}, \sigma_{site}, c) = \frac{1}{\sqrt{2\pi\sigma_{site}^2}} \exp\left(-\frac{(B_{anc} + c\vec{k}_m - B_m)^2}{2\sigma_{site}^2}\right). \quad (18)$$

Now we can combine our expressions for  $B_m$  and  $\vec{k}_m$  (Equations 17, Sections 2.2.1 and 2.2.2) with the new constraint of a linear relationship between  $B_m$  and  $\vec{k}_m$  (Equation 18). This allows us to obtain an expression for the site level intensity estimate  $B_{anc}$ :

$$P(B_{anc}, \sigma_{site}, c, B_m, k_m, D_m|x_m, y_m) \propto P(x_m, y_m|k_m, D_m, B_m)P(B_m|k_m, B_{anc}, \sigma_{site}, c)P(B_{anc}, \sigma_{site}, c)P(D_m, k_m). \quad (19)$$

Equation 19 may look complicated, but we defined each of the terms already. The benefit of this treatment is that we can obtain  $P(x_m, y_m|k_m, D_m, B_m)$  from our circle fitting in Equation 8 (see also Appendix 6.1). We defined  $P(B_m|k_m, B_{anc}, \sigma_{site}, c)$  in Equation 18. The values of  $\vec{k}$  and  $B_m$  for each specimen needed to fit both of these terms. This means that specimens with large scatter in their Arai plots (those which have Arai plots that are not fit well by a line or a circle) are more strongly affected by the site level fit  $B_{anc}$ , and therefore by the specimens with more linear (or circular) Arai plots. Conversely, those specimens with a small uncertainty in  $\vec{k}$  or  $B_m$  are tightly constrained by the Arai plot fit and so have more control over the fit at the site level.



**Figure 5.** Examples demonstrating how the predicted  $\vec{k}$  and  $B_m$  for each specimen are modified for a site by using a hierarchical model (Equation 14). The left column shows samples from the posterior for an “unpooled” model where we estimate  $B_m$  and  $\vec{k}_m$  independently. The right column shows samples from the posterior for the BiCEP method where we assume a linear relationship between  $B_m$  and  $\vec{k}_m$ . a) Red horizontal line is  $B_{exp}$  (hw126, see Table 1). 95% credible intervals for  $\vec{k}_m$  and  $B_m$  are plotted using black error bars, with the medians as green points. b) Representative MCMC samples from the posterior distribution are plotted as blue lines assuming that the individual specimen values  $B_m$  follow the relationship stated in Equation 14. Note that the higher curvature specimens with large uncertainty in  $\vec{k}$  follow a linear trend away from  $B_{exp}$ . c),e),g): [Symbols same as in Figure 4.] Arai plots of particular specimens are shown with circle fits sampled from the posterior of the unpooled model shown in a) and plotted in green. In d), f), h), same specimens as in c), e), g) but using the posterior of the BiCEP model in b). Note that there is little change in the specimen in d) for which a close fit to the data is possible, but in f) and h) the curvature (and intensity) of the specimen are modified to fit the line better.

The other two terms on the right side of Equation 19 ( $P(B_{anc}, \sigma_{site}, c)P(D_m, k_m)$ ), are priors.  $P(D_m, k_m)$  were defined in Equations 14 and 16 respectively. Now, we need to define priors for  $P(B_{anc}, \sigma_{site}, c)$ . For this purpose, we use a poorly constrained prior for the slope,  $c$ , where  $P(c) \propto 1$ . Although this is not a probability distribution, the resulting posterior distribution for  $B_{anc}$  is always a real probability distribution if the number of specimens is greater than one. We use a uniform prior between 0 and 250  $\mu\text{T}$  for  $P(B_{anc})$  as intensity values can never be negative and in databases such as the MagIC database (Tauxe et al., 2016) or the PINT database of (Biggin, 2010) rarely (if ever) exceed 250  $\mu\text{T}$ . For  $P(\sigma_{site})$  we use a normal distribution with zero mean and standard deviation of 5  $\mu\text{T}$ , truncated to always be positive.

Figure 5b shows our median estimates for  $B_m$  and  $\vec{k}_m$  after applying the linear restriction. Here, there is a tradeoff between fitting the Arai plot data with the circle, and fitting the linear trend at a site level. The effect of the linear fitting is apparent when compared to estimating  $\vec{k}_m$  and  $B_m$  for each specimen in isolation, which is shown in Figure 5a. With the linear restriction, the  $\vec{k}$  and  $B_m$  of specimens are “pulled” closer to a linear trend by modifying the Arai plot fits; specimens with more uncertain  $\vec{k}_m$  are more strongly affected (e.g., specimen labeled e) and f) in Figure 5a and b). The specimens with highly linear Arai plots (for which we have small uncertainty in  $\vec{k}_m$ ), the circle fits (see g and h) are mostly unchanged. Despite this modification of the circle fits to the Arai plots by the linear model, the circle fits to those specimens do not look unreasonable.

### 2.3 Metrics of success

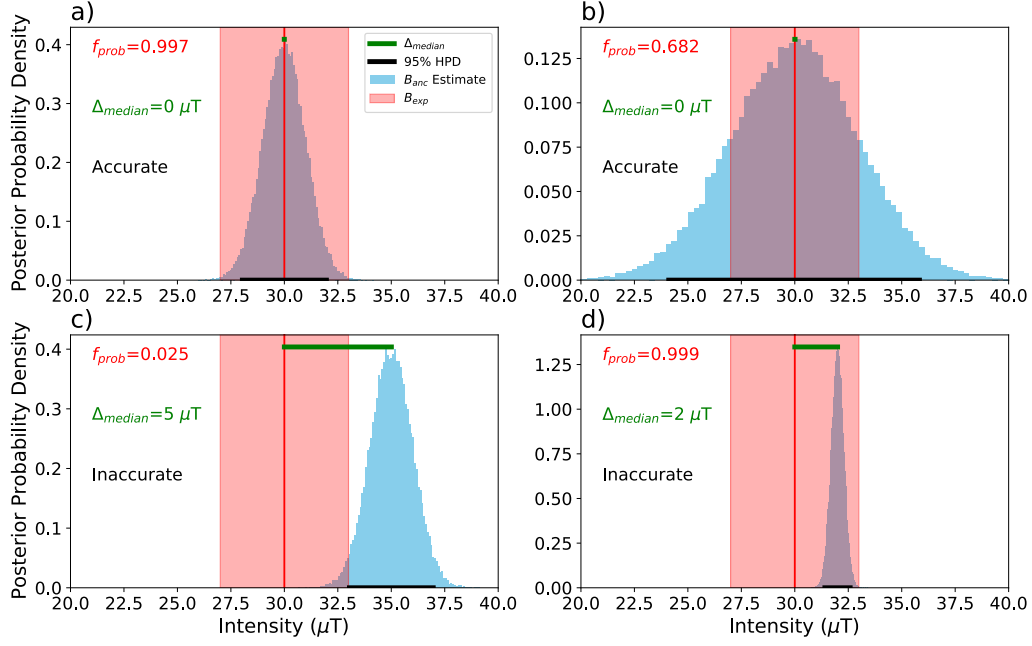
In order to ‘ground-truth’ the method, we rely on a compilation of paleointensity data updated from that of Paterson et al. (2014) and Tauxe et al. (2016). This compilation has data from 30 sites for which  $B_{anc}$  is well constrained (hence we use  $B_{exp}$ ), either through the IGRF, or because the specimens were given TRMs in a known lab field before the Thellier experiment. A list of sites used here is given in Table 1. Instead of choosing a range of temperatures for each site, we simply use every temperature on the Arai plot for all specimens.

Because we have to estimate multiple parameters for each specimen, our method involves a high dimensional optimization problem. Therefore, we generate the estimates for  $B_{anc}$  for a given site using a Markov chain Monte Carlo (MCMC) method which approximates the posterior distribution by generating pseudosamples from it (see Appendix 6.2 for details). MCMC techniques are frequently used to solve high dimensional problems of this kind.

For each site, we quantify the effectiveness of the BiCEP method using several metrics,  $f_{prob}$ ,  $\Delta_{median}$  (see Figure 6 for graphical representation),  $\bar{f}_{prob}$ , and  $n_{acc}$ :

1.  $f_{prob}$ : We report the median value of our posterior distribution and the 2.5<sup>th</sup> and 97.5<sup>th</sup> percentile of the Monte Carlo sample (95% credible interval) as error bars. To quantify the effectiveness of our method, we look at the proportion of the posterior distribution that lies within 3  $\mu\text{T}$  of the expected value of  $B$  ( $B_{exp}$ ) and call this proportion  $f_{prob}$ .
2.  $\bar{f}_{prob}$ : the mean value of  $f_{prob}$  over all sites included in the study. A value of 1 is the best possible value and means all our results are accurate and precise to better than 3  $\mu\text{T}$ .
3.  $\Delta_{median}$ : the difference (in  $\mu\text{T}$ ) between the median value of the MCMC sample (see Section 6.2 for explanation) and  $B_{exp}$ . The median value of  $\Delta_{median}$  is  $\tilde{\Delta}_{median}$ . Values of  $\tilde{\Delta}_{median}$  close to zero are best.
4.  $n_{acc}$ : the number of sites for which  $B_{exp}$  lies within our 95% credible interval. A related parameter,  $f_{acc}$  is the fraction of results that are accurate ( $n_{acc}/M$ ), where





**Figure 6.** Examples of accuracy and precision metrics used in this study with simulated Gaussian distributions of  $B_{anc}$  for illustration. a) An accurate and precise estimate, b) An accurate but imprecise estimate, c) An inaccurate and imprecise estimate. d) A slightly inaccurate and highly precise estimate. Accuracy check used for  $n_{acc}$  checks whether the black line intersects the expected value ( $B_{exp}$ ).  $f_{prob}$  is the area of the blue histogram that lies within the red shaded area.  $\Delta_{median}$  is the length of the green line.

$M$  is the number at the site level. We expect this number to be 0.95 in ideal circumstances.

We use these metrics to compare the BiCEP results to those obtained by two different sets of selection criteria: CCRIT (Cromwell et al., 2015) and Paterson’s modified PICRIT03 criteria (Paterson et al., 2014) (here called PICRITMOD) without the curvature parameter  $\vec{k}$ . Most sets of commonly used selection criteria rely on an assumption of a Gaussian probability distribution for the site level estimate  $B_{anc}$ , which allows us to calculate these same metrics. For this analysis, we exclude sites that contain fewer than three specimens, or those for which have one or fewer specimens that meet any of the selection criteria in CCRIT or PICRITMOD. We discuss the results of this comparison in Section 3.1.

## 2.4 Width of prior and order of fit

Here we consider several alternative contingent models in order to explore our choices for  $P(\sigma_{site})$  and assumptions about the relationship of  $B_m$  and  $\vec{k}$ . In addition to using a standard deviation of  $5 \mu T$  for  $P(\sigma_{site})$ , we use standard deviations of  $10 \mu T$  and  $20 \mu T$ . The effect of this is hard to conceptualize, but wider priors will prioritize fitting circles to the individual specimens over fitting the linear relationship between  $B_m$  and  $\vec{k}_m$  at a site level. The practical effect of this is wider posteriors for sites where the number of specimens is small.

So far, we have assumed *a priori* that  $B_m$  is linearly dependent on  $\vec{k}_m$ . Because there is no theoretical reason why this should be the case, we test models for which the relationship between  $B_m$  and  $\vec{k}_m$  is described by a quadratic polynomial and a cubic polynomial. We would expect a higher order model to more closely fit the individual  $\vec{k}_m$  and  $B_m$  values, but with a loss of precision due to the more complicated model.

Results for our method, as well as for two sets of selection criteria, are given in Table 2. For each model, we calculate  $\bar{f}_{prob}$ ,  $\tilde{\Delta}_{median}$  and  $f_{acc}$  for comparison. In this table, our models are named for the value of the standard deviation of  $P(\sigma_{site})$  as well as the order of the fit. Our preferred model is referred to as “Linear 5  $\mu$ T”, and this is the model used in this paper where otherwise unspecified.

## 2.5 MCMC sampler diagnostics

MCMC samplers are only ever an approximation of the posterior distribution, and the number of Monte Carlo samples needed to make an accurate approximation is not the same for every site, or every run of the sampler. To determine whether we are accurately sampling the posterior distribution, we look at three diagnostics which are also described in Appendix 6.2:

1.  $\hat{R}$ : (Gelman & Rubin, 1992) quantifies convergence between chains in the MCMC method. This parameter is required to be between 1.1 and 0.9 for the sampler to converge.
2.  $n_{eff}$ : the effective MCMC sample size. We are using 30,000 Monte Carlo samples and  $n_{eff}$  should be large ( $> 1000$ ) to have a good representation of our parameters.
3.  $f_{div}$ : the proportion of divergent transitions  $f_{div}$  in the MCMC sample. This should ideally be zero, but it does not appear to cause large problems for the estimate of  $B_{anc}$  if it is non zero (see Section 6.2).

The diagnostics  $n_{eff}$  and  $\hat{R}$  are produced for each of our parameters (each of our  $B_m$ ,  $\vec{k}_m$ ,  $D_m$  and  $B_{anc}$ ,  $\sigma_{site}$ ). When reporting these values, we look at the worst value of  $\hat{R}$  (furthest from unity) and the value of  $n_{eff}$  for  $B_{anc}$ . If  $\hat{R} > 1.1$ , we replace the distribution on  $B_{anc}$  with a uniform distribution between 0 and 250  $\mu$ T (the prior). The results of the MCMC sampler are presented in Section 3.4.

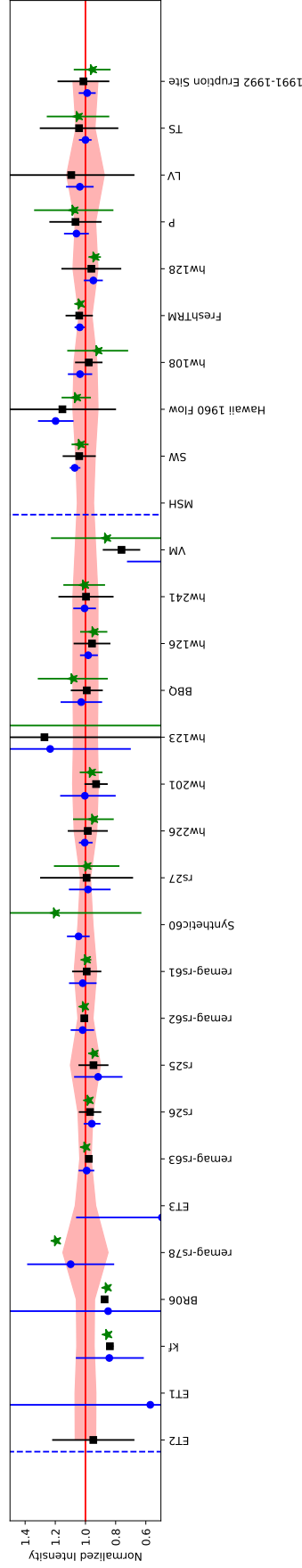
## 3 Results

### 3.1 Comparison of BiCEP to CCRIT and PICRITMOD

In this section, we compare the BiCEP to the CCRIT and PICRITMOD sets of selection criteria (see Section 2.3). The full set of results for all sites can be seen in Figure 7, and are summarized in Table 3.

Figure 7 shows the 95% credible intervals for each method, normalized by the expected value at the site. The median values of our results are generally similar to those found by CCRIT and PICRITMOD. Results from BiCEP were the most accurate, followed by CCRIT and PICRITMOD. BiCEP and PICRITMOD achieved similar levels of precision, with CCRIT’s results being slightly less precise. This indicates that our method achieves its goal of eliminating the bias introduced by poor quality specimens at least as well as selection criteria without having to exclude those specimens from the analysis based on arbitrarily chosen threshold values.

Sites in Figure 7 are sorted by the number of specimens used by BiCEP for the analysis. Unique to our method, sites with low numbers of specimens ( $M$ ) have wide credible intervals and sites with high  $M$  have narrow credible intervals, so the estimate of



**Figure 7.** Paleointensity estimates for our collection of sites (Table 1) using our method (blue circles) with 95% confidence interval compared to results using CCRIT (red diamonds) and Paterson's modified PICRITMOD (green stars). A dashed blue line indicates a site where the sampler failed with  $\hat{R} > 1.1$ , so the prior distribution (a uniform distribution between 0 and  $250 \mu\text{T}$ ) was used. The results are normalized to the expected field value for each site (black line, by definition 1) and the pink shaded region represents  $\pm 3 \mu\text{T}$ . Sites are ordered by the number of specimens used by our method for paleointensity analysis in that site. Results are summarized in Table 3

Model Name	$\bar{f}_{prob}$	$\tilde{\Delta}_{median}$ ( $\mu\text{T}$ )	$f_{acc}$
Linear, 5 $\mu\text{T}$	0.67	1.5	0.80
Linear, 10 $\mu\text{T}$	0.65	1.4	0.80
Linear, 20 $\mu\text{T}$	0.65	1.5	0.80
Quadratic, 5 $\mu\text{T}$	0.60	1.6	0.75
Quadratic, 10 $\mu\text{T}$	0.59	1.7	0.80
Quadratic, 20 $\mu\text{T}$	0.59	1.6	0.80
Cubic, 5 $\mu\text{T}$	0.47	2.3	0.85
Cubic, 10 $\mu\text{T}$	0.40	3.3	0.80
Cubic 20 $\mu\text{T}$	0.40	3.9	0.90
CCRIT	0.59	1.7	0.90
PICRIT (Modified)	0.67	2.1	0.80

**Table 2.** Results comparing the models used in this study to results using CCRIT (Cromwell et al., 2015) and PICRITMOD (Paterson et al., 2014). See details in text and Figure 6 for explanations of the different parameters presented here. Results are sorted by the number of specimens in the site used to make the estimate using our method.

$B_{anc}$  becomes more precise as more specimens are measured. This is because calculating the credible interval for a  $B_{anc}$  is more similar to calculating the standard error of the mean than the site level standard deviation, which is done for CCRIT and PICRITMOD.

The increasing precision on  $B_{anc}$  leads to some sites with high  $M$  having estimates of  $B_{anc}$  which are seemingly too precise. These estimates are still generally only a few  $\mu\text{T}$  away from the expected value, however, and we discuss potential reasons for this in Section 4.3.

Our increased level of accuracy and precision relative to CCRIT is demonstrated in our metrics (Table 2). We see that BiCEP has a higher  $\bar{f}_{prob}$  (0.67 vs 0.59) and  $\tilde{\Delta}_{median}$  (about 1.7  $\mu\text{T}$  vs 1.5  $\mu\text{T}$ ), although a difference of 0.2  $\mu\text{T}$  is not particularly significant. The sites which are overly precise are reflected in BiCEP’s lower  $f_{acc}$  (0.80 vs 0.90). PICRITMOD achieved a similar  $\bar{f}_{prob}$  to our method (also 0.67). This is probably because it is a looser set of criteria, which means there are more interpretations to pick from to optimize the standard deviation.

### 3.2 Width of the prior

To investigate the role of the prior distribution ( $P(\sigma_{site})$ ), we apply the BiCEP method on the data compilation using a variety of values for its standard deviation (see Table 2). The main effect of varying  $\sigma_{site}$  is that for smaller values, the estimates of  $B_m$  and  $\vec{k}_m$  for specimens are “pulled” closer to the line being fitted at a site level (see Figure 5a,b). For our estimate of  $B_{anc}$ , this means that sites with fewer specimens will be more precise, as it is unlikely that specimen  $B_m$  will deviate strongly from the mean. For sites with many specimens, there is little effect as  $\sigma_{site}$  is well constrained by the data.

From Table 2, we see that changes to  $P(\sigma_{site})$  seem to have little influence on the effectiveness of the model, as all our  $f_{acc}$  values are the same for our linear model regardless of the prior distribution used. We can also see graphically in Figure 7 that our precision is low for these sites. Because of this, we favor the version of the model with a 5  $\mu\text{T}$  standard deviation on  $P(\sigma_{site})$ , as models with higher standard deviations reduce precision without capturing any more sites within their 95% credible intervals.

**Table 3.** Results for each site, including the value of  $f_{prob}$  using our method (BiCEP, superscript ‘B’) and two different sets of selection criteria, (CCRIT, superscript ‘C’ and PICRITMOD, superscript ‘P’).

Site Name	$B_{exp}$	$M^B$	$B_{min}^B$	$B_{med}^B$	$B_{max}^B$	$f_{prob}^B$	$B_{med}^C$	$B_{min}^C$	$B_{max}^C$	$f_{prob}^C$	$B_{med}^P$	$B_{min}^P$	$B_{max}^P$	$f_{prob}^P$
1991-1992 Eruption Site	36.2	53	33.9	35.8	37.6	1.00	36.7	30.7	42.7	0.67	34.6	30.4	38.8	0.73
hw108	39.3	23	37.8	40.7	43.7	0.86	38.4	35.1	41.8	0.89	36.1	28.4	43.8	0.43
hw123	37.7	12	26.6	46.5	62.9	0.15	48.0	-10.7	106.6	0.08	63.4	9.56	117.3	0.06
hw126	36.4	13	33.6	35.8	37.5	0.98	34.8	30.6	39.1	0.73	34.4	31.3	37.5	0.74
hw128	36.2	26	32.2	34.3	36.4	0.86	34.8	27.8	41.8	0.57	34.0	32.7	35.3	0.89
hw201	35.2	12	28.3	35.3	40.9	0.68	32.7	30.2	35.2	0.66	33.9	31.4	36.3	0.91
hw226	39.9	11	38.2	40.1	41.4	1.00	39.3	34.2	44.4	0.75	37.8	32.6	43.0	0.61
hw241	36.0	18	33.7	36.2	38.7	0.98	35.9	29.5	42.3	0.65	36.3	31.5	41.1	0.79
BR06	49.7	3	16.5	42.2	85.6	0.11	43.4	43.4	43.4	0.00	42.6	42.6	42.6	0.00
P	44.6	36	43.9	47.3	50.7	0.58	47.5	40.1	55.0	0.45	48.1	36.6	59.5	0.34
VM	43.8	18	4.0	18.7	31.5	0.00	33.3	28.1	38.5	0.00	37.7	21.9	53.6	0.22
BBQ	36.2	12	32.4	37.2	42.0	0.76	35.9	32.2	39.5	0.90	39.2	31.0	47.4	0.43
rs25	30.0	5	22.8	27.5	32.1	0.60	28.4	25.6	31.2	0.84	28.3	27.7	29.0	1.00
rs26	60.0	5	54.4	57.5	60.4	0.66	58.2	54.0	62.3	0.70	58.6	58.6	58.7	1.00
rs27	90.0	10	75.7	88.5	99.4	0.43	89.3	62.1	116.6	0.17	89.3	70.2	108.3	0.25
remag-rs61	40.0	6	37.3	40.7	44.1	0.91	39.7	36.0	43.3	0.89	39.8	38.7	40.9	1.00
remag-rs62	60.0	6	56.9	61.1	65.6	0.82	60.5	60.5	60.5	N/A	60.5	60.4	60.7	1.00
remag-rs63	80.0	5	75.8	79.4	83.2	0.90	78.2	78.0	78.5	1.0	79.9	79.5	80.1	1.00
remag-rs78	20.0	4	16.3	22.0	27.6	0.66	N/A	N/A	N/A	N/A	23.9	23.9	23.9	0.00
kf	52.0	3	32.2	43.8	55.0	0.05	43.6	42.6	44.5	0.00	44.4	44.3	44.5	0.00
Hawaii 1960 Flow	36.0	22	39.0	43.1	47.1	0.02	41.5	28.9	54.1	0.26	38.2	34.9	41.5	0.69
SW	46.4	19	48.2	49.7	51.0	0.33	48.3	43.5	53.2	0.65	48.1	45.7	50.4	0.87
TS	47.8	53	46.1	47.8	49.6	1.00	49.8	37.7	62.0	0.36	50.2	40.5	59.8	0.42
ET1	43.3	3	3.7	24.7	106.3	0.04	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
ET2	44.1	2	6.2	125.0	243.7	0.02	42.0	30.0	53.6	0.36	N/A	N/A	N/A	N/A
ET3	44.2	4	2.3	21.8	46.7	0.03	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
Synthetic60	60.0	7	58.7	62.8	67.0	0.55	N/A	N/A	N/A	N/A	72.0	38.1	106.0	0.11
LV	24.0	45	22.8	24.9	27.0	0.98	26.3	16.3	36.2	0.41	N/A	N/A	N/A	N/A
MSH	55.6	19	6.1	124.6	243.9	0.02	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
FreshTRM	70.0	24	70.7	72.6	74.7	0.65	72.9	66.9	78.9	0.49	72.5	72.0	73.1	0.96

### 3.3 Order of polynomial fit

The results for our test sites (Table 2) demonstrate that increasing the order of the polynomial fit significantly decreases the precision of the estimate as demonstrated by reduced values of  $\hat{f}_{prob}$ . This is expected as there are more parameters to be estimated with the same number of data. The level of accuracy is also reduced, with  $\hat{\Delta}_{median}$  increasing as the order of the fit increases. For this reason, we assume a linear relationship between  $B_m$  and  $\vec{k}_m$ .

### 3.4 Sampler Diagnostics

Site Name	Worst $\hat{R}$	$n_{eff}$	$f_{div}$
1991-1992 Eruption Site	1.00	59741	0.00
hw108	1.00	77959	0.00
hw123	1.01	11687	0.00
hw126	1.00	36130	0.01
hw128	1.00	78978	0.00
hw201	1.00	10641	0.01
hw226	1.00	7139	0.05
hw241	1.00	66565	0.00
BR06	1.01	451	0.00
P	1.00	62252	0.00
VM	1.05	1447	0.00
BBQ	1.00	63082	0.00
rs25	1.00	5614	0.00
rs26	1.00	11866	0.00
rs27	1.00	22211	0.00
remag-rs61	1.00	26746	0.00
remag-rs62	1.00	16916	0.00
remag-rs63	1.00	3788	0.00
remag-rs78	1.00	12388	0.00
kf	1.02	2712	0.00
Hawaii 1960 Flow	1.00	60184	0.00
SW	1.00	36390	0.00
TS	1.00	56518	0.00
ET1	1.01	995	0.00
ET2	6.93	6	0.03
ET3	1.01	424	0.00
Synthetic60	1.00	36572	0.01
LV	1.02	5931	0.08
MSH	2.78	24	0.45
FreshTRM	1.00	81007	0.00

**Table 4.** Sampler diagnostics (see Section 2.5 for an explanation of each diagnostic) for each site using the BiCEP method.

The sampler diagnostics for each site are given in Table 4. Indicators of poor MCMC sampler performance (worst  $\hat{R} > 1.1$ , low  $n_{eff}$ , high  $f_{div}$ ) tend to occur at sites with four or fewer specimens, or sites where the Arai plots are extremely scattered and the sampler struggles to fit them. This indicates that to get a strongly reproducible answer from this method, paleomagnetists ought to measure five or more specimens per site. In practice, most studies already do this in order to have enough specimens that pass the cho-

sen selection criteria, yet many specimens may be excluded from analysis. Here, we can use all of the specimens measured so there may be no additional burden.

### 3.5 Summary of Results

After testing all of our contingent models, we prefer the model which assumes the relationship between  $B_m$  and  $\vec{k}_m$  is linear, and which uses a  $5 \mu\text{T}$  standard deviation on  $P(\sigma_{site})$ . This model has a higher level of accuracy than the PICRIT and CCRIT sets of selection criteria and higher precision than CCRIT. Our precision increases for sites for which the number of specimens is large, similarly to calculating the standard error of the mean when using selection criteria. Unlike selection criteria, the BiCEP method does not require exclusion of large numbers of specimens to obtain an accurate result, which leads us to prefer it over those methods.

## 4 Discussion

### 4.1 Advantages of BiCEP compared to selection criteria

BiCEP has significant advantages over classical selection criteria. Firstly, we obtain estimates for all sites with at least three specimens, including some which do not contain any specimens that pass classical selection criteria, see Figure 7. In most cases, our estimates have give similar or higher accuracy than the selection criteria (evidenced by  $\hat{\Delta}_{median}$  and Figure 7), and this is accomplished while only excluding specimens from the analysis which were not fully demagnetized.

Secondly, the increasing precision of our paleointensity estimate as the number of specimens increases allows for an improved workflow when compared to classical paleointensity criteria. Instead of needing a minimum number of specimens to pass our selection criteria, we can keep measuring specimens until we reach a desired level of precision. We discuss this workflow in more detail in Section 4.2. The property of increasing precision with number of specimens is inherent to Bayesian models and can also be found in the method of Kosareva et al. (2020), although this method does not include the bias correction found in our method.

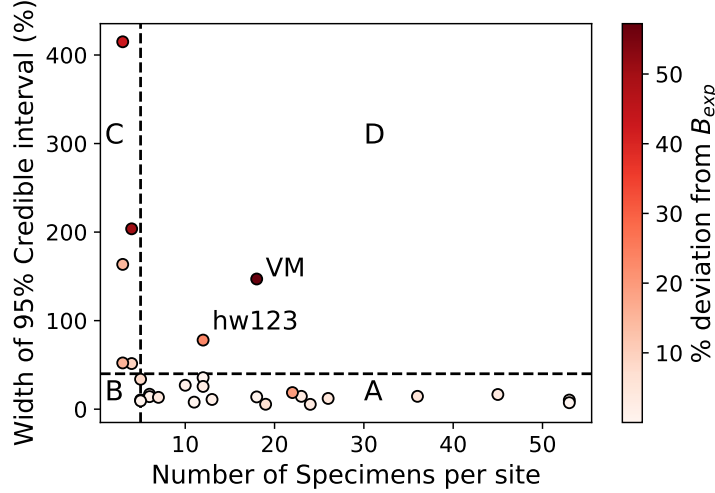
Thirdly, our method propagates the uncertainties from a specimen to the site level. Specimens with more scattered (or non linear, or non circular) Arai plots will have less influence over the specimen mean than those with highly linear Arai plots. In addition to this, the BiCEP method foregoes the need for criteria which are concerned with the length of the line on the Arai plot used to make an interpretation, like the NRM Fraction (e.g., FRAC of Shaar & Tauxe, 2013). Using a set of temperatures with small FRAC will cause an increase in the uncertainty in  $\vec{k}$  (see Figure 4e, f), which will cause this specimen to have less effect on the estimate of  $B_{anc}$ , without excluding it from the analysis entirely. We discuss this further in Section 4.4.

### 4.2 Workflow with BiCEP

Figure 8 plots precision (here expressed as the full width of the 95% credible interval as a percentage of the median) against the number of specimens per site ( $M$ ). The use of the 95% credible interval in BiCEP differs from the use of the standard error of the mean with selection criteria in that relying on the standard error generally leads to overly precise measurements in the frequent case of low numbers of specimens ( $M < 5$ ).

We have divided Figure 8a and b into four regions (A-D). Region A has high precision with many specimens. Region B has high precision (better than 40%, which for a Gaussian distribution would be equivalent to a standard deviation of  $\pm 10\%$ ) with few





**Figure 8.** Plot of precision of estimates using the Linear, BiCEP against the number of specimens per site for all sites where  $\hat{R} < 1.1$ . Colors indicate the deviation of the median value of the estimate from the expected site value ( $B_{exp}$ ) as a percentage. The horizontal dashed line indicates a value of 40% for the full width of the 95% credible interval, which for a Gaussian distribution would correspond to a standard deviation of  $\pm 10\%$ . Suggested workflow for sites in regions: A) Accept the site. B) Continue measuring if better precision is desired, if not, accept the site as is. C) Continue measuring specimens, as improved precision is likely. D) Stop measuring the site as further effort is likely to be futile.

specimens ( $M < 5$ ). Region C has low precision but a low number of specimens ( $M < 5$ ) and Region D has similarly poor precision but a large number of specimens ( $M > 5$ ). With few exceptions (VM, hw123), the BiCEP method allows for increasing precision in the estimate of  $B_{anc}$  as the number of specimens increases. For sites where the median value for the intensity is low, it may be more appropriate to define these regions using an absolute width for the intensity bound rather than a percentage.

Because all of the sites considered here have known values for the site intensity, we can also consider accuracy. From Figure 8, we see that almost all sites with an estimated precision better than 40% have median values within 20% of  $B_{exp}$ , as might be expected. (color of dots reflect the deviation from  $B_{exp}$ ). Considering the region in which a particular site plots leads to a workflow based on the likelihood of success. For sites in Region C, obtaining more specimens is likely to result in improved precision and accuracy. If after measuring more specimens, a site may move from region C into region D (e.g., sites hw123 and VM) and further effort is likely to be futile. If a site moves from region C into regions A or B, the site may be acceptable, depending on the desired precision. If a higher level of precision is desired (better than 40%), increasing  $M$  is likely to be successful.

### 4.3 Overly precise estimates of $B_{anc}$

The BiCEP method has a lower  $f_{acc}$  than CCRIT, despite having a similar degree of accuracy when using a metric like  $\Delta_{median}$ . The reason for this is that the increasing precision on the BiCEP estimate leads to estimates which are highly precise when  $M$  is large. This is the case shown in Figure 6d.

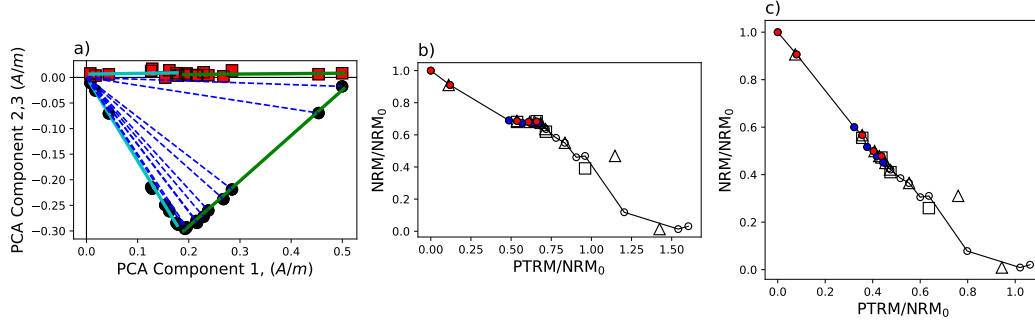
Labeling sites with extremely high precision in the estimate may be misleading, as we have not taken into account uncertainties in the value of the expected fields at the sites in this study. For example, using differences between the observed directions and the IGRF, Yamamoto and Hoshi (2008) quoted the expected value at the site “SW” as  $46.0 \pm 2.6 \mu\text{T}$ , which is just consistent with the 95% credible interval for our specimen ( $48.2\text{--}49.7 \mu\text{T}$ ). Because of this, we prefer to use  $\bar{f}_{prob}$  as a metric of how well a model performs as it allows for a few  $\mu\text{T}$  of uncertainty in the expected field value. Additionally, Yamamoto and Yamaoka (2018) suggested that the IZZI-Thellier results for sites SW and TS may be biased slightly high due to acquisition of a thermo-chemical remanent magnetization (TCRM), which is not detectable by our method. Yamamoto et al. (2003) also invoke a TCRM mechanism to explain the paleointensity overestimate for the Hawaii 1960 Flow, which is another of their sites for which we overestimate the expected intensity (see Table 3). We note that Cromwell et al. (2015) also sampled the 1960 flow (hw241 which targeted the fine grained flow top) and all selection criteria resulted in accurate results, with BiCEP producing the tightest confidence interval.

#### 4.4 Exclusion of measurement level data

It is frequently possible to improve the accuracy and precision of results by finding the ‘best’ set of temperature steps to use in the intensity interpretation. Two situations frequently occur for which this might be justified. The first is the case in which thermochemical alteration occurs at high temperature (e.g., Figure 4e). For such specimens, the low temperature measurements can be used to make a paleointensity estimate (colored dots in the figure). Figures 4e and f show how our method can be used on a reduced range of temperature steps on the Arai plot at the cost of precision. The plot of circle fits (green lines in Figure 4e) demonstrates that the Arai plot interpretations are poorly constrained and can continue in any direction after the last temperature step chosen. This results in a higher uncertainty in the curvature associated with this (Figure 4f). The second case in which a portion of the data could be excluded from the calculation, would be when the magnetization has multiple components (Figure 9a). In such a case, a paleointensity estimate can only be made using the small range of temperature steps that correspond to the characteristic component. We currently do not have an objective method to choose which set of temperature steps on the Arai plot to use. We suggest that decisions about which data points to include should not be made based on the original in-field or zero field Arai plot measurements (dots in the Arai plots), but rather exclusively on deviating pTRM checks (triangles in, e.g., Figure 4e) or other indicators of alteration for the first case and on the directions of the magnetization vector (it must trend to the origin and be well defined) in the second case, e.g., Figure 9a.

#### 4.5 Application to multi-component magnetizations

We test an application of the BiCEP method on data with multi-component directions as shown in Figure 9a using the data of Lisé-Pronovost et al. (2020). The data are from Scottish firebricks which were used in a foundry in Australia. The date and location of firing are both well constrained, hence we have a reasonably well constrained value for  $B_{exp}$ . The bricks all contained a low temperature component associated with the Australian field. Some also displayed a high temperature component associated with the original firing in Scotland as shown in Figure 9a. Lisé-Pronovost et al. (2020) already have interpretations which separate these components in the original study. To account for the change in direction of the NRM, we subtract the high temperature component from the low temperature component, and then add the magnitude of these values to the magnitude of the low temperature component (see Figure 9 for a graphical explanation). The vector subtraction is necessary for the low-temperature component as we need a total TRM ( $p\text{TRM}_{max}$ ) to scale by in order to penalize the result for shorter components. We then proceed to use the BiCEP method as previously described, using the



**Figure 9.** a) Example of vector endpoint diagram for specimen FB2-B1 from Lisé-Pronovost et al. (2020). The magnetization is rotated so that the principal component of the TRM direction for all steps lies along the x axis. Green line fit to the low temperature component and cyan line fit to the high temperature component. b) Arai plot and c) “corrected” Arai plot for a specimen from the data shown in b). NRM values for the low temperature component (filled circles) are calculated by taking the magnitude of the vector endpoint (blue dashed lines in the vector endpoint diagram in a). In b), these NRM values are calculated by vector subtracting the high temperature component (cyan line), taking the magnitude of our new NRM vectors (distance along green line), and adding the magnitude of the low temperature component (length of cyan line).

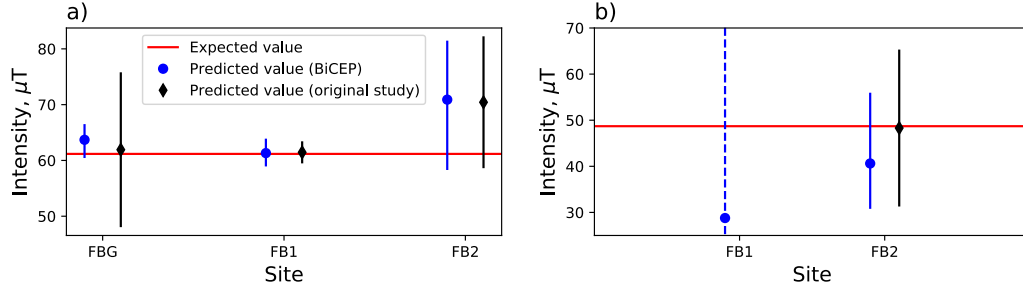
original interpretations for the different components. For the sake of simplicity, we do not perform the magnetomineralogical change (MMC) correction (Valet et al., 1996). We also do not apply the corrections for anisotropy of TRM or cooling rate with these data, as they appeared to be negligible. Of course these could be applied in the usual fashion if necessary.

We display the results from multi-component remanences in Figure 10. We find that for the low temperature, Australian field, component (Figure 10a), our estimates for all firebricks contain the expected answer ( $61.17 \mu T$ ) within the 95% credible interval. Our interpretation for site FBG is slightly less accurate than the original analysis but with much higher precision. This difference is likely caused by not applying the MMC correction, as the specimens at this site were mostly of good quality, with none being excluded from the original analysis. (Figure 10b) behaves differently.

The sampler does not converge for site FB1, indicating too few specimens in the analysis. For site FB2, we have a result that is less accurate, but more precise than in the original study. The lack of MMC correction may contribute to the decreased accuracy in this example, whereas the reduced precision is likely caused by the smaller length of the interpretation on the Arai plot, leading to a higher uncertainty in the curvature for that specimen. Our results for this study demonstrate that BiCEP will be effective for obtaining precise estimates for components which represent most of the magnetization, and ineffective for components which have small NRM fraction.

#### 4.6 Implications for bias in curved Arai plots

The success of our method demonstrates that Arai plot “curvature” or sagging does lead to a progressive bias in paleointensity estimation which increases as the amount of curvature increases as described by Tauxe et al. (2021) and strongly suggested by the data of Krása et al. (2003) (see Figure 1). Our estimates are made by using the tangent to a circle fit rather than fitting a line to part of the data, so one might expect them to



**Figure 10.** Expected and predicted intensities on the data of Lisé-Pronovost et al. (2020) using BiCEP (blue circles) and the method used in the original study (black diamonds). a) Results for the low temperature component (Australia, expected field value 61.17  $\mu\text{T}$ ) for each firebrick. b) Results for the high temperature component (Scotland, expected field value 48.3  $\mu\text{T}$ ), where this component was present. The dashed blue line indicates that the MCMC sampler failed to converge for site FB1.

be biased. However, it has been demonstrated by e.g. the data of (Krása et al., 2003) that fitting lines to the high temperature or low temperature slope of Arai plots yields even more biased results. Fitting a tangent gives a result more similar to the best fitting line to all data, or the total TRM, both of which exhibit a similar progressive bias for curved Arai plots. The bias seen generally underestimates paleointensity with higher (positive) curvature, but this is not the case for all sites, some of which exhibit the opposite trend.

The assumption of a quasi-linear dependence between the specimen level paleointensities and the curvature of the Arai plot does not have any theoretical basis. We stress that this relationship only needs to be loosely followed for our method to work. In cases where there does not appear to be a strong linear relationship between  $B_m$  and  $\vec{k}_m$  (e.g. in Figure 2c), an accurate paleointensity estimate is still possible if there are enough specimens with low  $|\vec{k}|$ , as the intercept of the linear fit is still well constrained even if the slope is not. Conversely, if there are few specimens with high  $|\vec{k}|$  and there is a poor linear relationship, then both the slope and intercept are poorly constrained, resulting in a huge uncertainty in  $B_{anc}$ .

## 5 Conclusions

- We present a new Bayesian method (BiCEP) which accounts for bias in paleointensity estimates in specimens.
- Instead of excluding specimens from the paleointensity analysis in the traditional (binary) selection criteria based approach, our method predicts an amount of bias for each specimen, using the curvature of the Arai plot as a metric of non-linearity and a predictor of bias. In this way, the BiCEP method is quite different from the recently published Bayesian approach of Kosareva et al. (2020).
- When tested on a compilation of sites for which an approximate paleointensity is known *a priori*, our method is more accurate than two commonly used sets of selection criteria, and has a similar level of accuracy to the modified PICRIT criteria of Paterson et al. (2014).
- Our method generates some slightly inaccurate paleointensity estimates with high levels of precision, but these can generally be explained with inaccuracies in the expected field (see Section 4.3).

- The BiCEP method handles uncertainties in a different way than using classical selection criteria, as the uncertainty in site level estimates decreases as the number of specimens increases, but this uncertainty remains high when the number of specimens is low due to inclusion of prior information. The Bayesian uncertainties are in this way more similar to the ‘extended error bars’ in the Thellier\_GUI auto-interpreter of Shaar and Tauxe (2013).
- We propose a workflow in which sites are accepted and measurement of specimens can cease once a desired level of confidence in the site level estimate has been reached. Sites which do not reach this level of confidence after measuring several ( $> 5$ ) specimens likely do not contain useful information and can be discarded.

## Data Availability Statement

Data used in this paper may be found in the MagIC database at: <https://earthref.org/MagIC/17104/0326fdaa-4bcf-44f3-989d-0116b9a2fb75> for review and will be available to the public at <https://earthref.org/MagIC/17104> on publication.

## 6 Appendix

### 6.1 Change of variables

In Section 2.2.1 we mention that we need to use a change of variables to get from our original circle fitting parameters  $R, x_c, y_c$  to our new set of parameters  $\vec{k}, D, \phi$ . We can use the Jacobian of the parameter change to get the new formula for the posterior probability under our new parameters:

$$P(D, \phi, \vec{k}|x, y) = P(x_c, y_c, R|x, y) \left| \frac{\partial(x_c, y_c, R)}{\partial(D, \phi, \vec{k})} \right|. \quad (20)$$

We can evaluate this Jacobian as:

$$\left| \frac{\partial(x_c, y_c, R)}{\partial(D, \phi, \vec{k})} \right| = \left| \frac{\vec{k}}{|\vec{k}|^3} \left( D + \frac{1}{\vec{k}} \right) (\cos \phi + \sin \phi) \right|. \quad (21)$$

So our posterior looks like:

$$P(D, \phi, \vec{k}|x, y) \propto \left( \sum_{n=1}^N \sqrt{\left( \left( D + \frac{1}{\vec{k}} \cos \theta \right) - x_n \right)^2 + \left( \left( D + \frac{1}{\vec{k}} \sin \theta \right) - y_n \right)^2 - \frac{1}{|\vec{k}|}} \right)^{-N/2} \left| \frac{\vec{k}}{|\vec{k}|^3} \left( D + \frac{1}{\vec{k}} \right) (\cos \phi + \sin \phi) \right| P(\vec{k}, \phi, D). \quad (22)$$

### 6.2 Markov chain Monte Carlo sampling

The Markov chain Monte Carlo (MCMC) sampling method generates a set of samples from the posterior probability distribution of  $B_{anc}$  which allows us to approximate it. We use the python bindings for the Stan software package (<http://mc-stan.org>) to generate these samples which provides diagnostic information and runs relatively quickly. For each site we run four Markov chains and generate 30,000 samples of  $B_{anc}$  in each chain. We discard the first half of the chain as ‘burn in’ for a total of 60,000 samples.

Stan provides several diagnostics that tell us whether we have successfully sampled the posterior distribution. These include the  $\hat{R}$  score (Gelman & Rubin, 1992) which tells

us about the convergence between chains, and is required to be between 1.1 and 0.9 which is necessary for convergence, the effective sample size,  $n_{eff}$  which should be large ( $> 1000$ ) for a good sample and the number of divergent transitions ( $f_{div}$ ) which should be zero in ideal cases. In most cases our results display high degrees of convergence with  $\hat{R}$  close to 1 and high effective sample sizes. Some sites included divergent transitions in small numbers. These seem to occur at a specimen level for specimens where the posterior distribution of one of the circle parameters is long-tailed. In theory this can mean the posterior was inefficiently sampled, but because these specimens generally have large uncertainties on their  $\hat{k}$  parameter, the final results do not change, even under a change of parameters. The sampler struggled to converge, with  $\hat{R} > 1.1$  for several sites with very few specimens, where once again the distributions are extremely long tailed. The sampler also did not converge for site MSH, where the Arai plots were so non linear, with few points, that BiCEP struggled to fit circles to them. We consider these sites to have “failed” using our method (grade of ‘D’ in Figure 8) and use the prior distribution on  $B_{anc}$  (uniform between 0 and 250  $\mu$ T) as an estimate of their intensity. We calculate the  $\hat{R}$  furthest from unity, the  $n_{eff}$  for  $B_{anc}$  and the proportion of divergent samples  $f_{div}$  for our model.

### 6.3 Code and GUI

We present a simple GUI that can perform the BiCEP method on data in the MagIC format. The code uses Jupyter notebooks and can be found at ([http://github.com/bcych/BiCEP\\_GUI](http://github.com/bcych/BiCEP_GUI)) and contains a readme file detailing how to use the notebook. The GUI can also be accessed at the Earthref JupyterHub site (<http://jupyterhub.earthref.org>). To access the GUI this way:

- Sign up to Earthref at (<http://earthref.org>)
- Navigate to the Earthref JupyterHub site at (<http://jupyterhub.earthref.org>)
- Open and run all the cells in the “BiCEP GUI - Setup.ipynb” notebook.
- Upload MagIC formatted “sites”, “samples”, “specimens” and “measurements” files to the BiCEP\_GUI directory in JupyterHub. These can be formatted using pmag\_gui. (Tauxe et al., 2016).
- Open the BiCEP GUI notebook and press the “App Mode” button.

For more detailed instructions, read the included readme file at the github site.

### Acknowledgments

We are deeply grateful for the advice given by Andrew Roberts, David Heslop and Joseph Wilson. This research was supported in part by NSF Grants EAR1547263 and EAR1827263 to LT. We are also grateful to Agnes Lisé-Pronovost for sharing her measurement level data.

### References

- Biggin, A. (2010). Paleointensity database updated and upgraded. *EOS*, 91, 15.
- Chernov, N., & Lesort, C. (2005). Least squares fitting of circles. *Journal of Mathematical Imaging and Vision*, 23(3), 239–252. doi: 10.1007/s10851-005-0482-8
- Cromwell, G., Tauxe, L., Staudigel, H., & Ron, H. (2015). Paleointensity estimates from historic and modern hawaiian lava flows using glassy basalt as a primary source material. *Phys. Earth Planet. Int.*, 241, 44–56. doi: 10.1016/j.pepi.2014.12.007
- Dunlop, D., & Özdemir, O. (2001). Beyond Néel’s theories: thermal demagnetization of narrow-band partial thermoremanent magnetization. *Phys. Earth Planet. Int.*, 126, 43–57.



- Gelman, A., Carlin, J. B., Stern, H. S., & Rubin, D. B. (2004). *Bayesian data analysis* (Second ed.). Chapman & Hall/CRC, Boca Raton, FL.
- Gelman, A., & Rubin, D. B. (1992, 11). Inference from iterative simulation using multiple sequences. *Statist. Sci.*, 7(4), 457–472. Retrieved from <https://doi.org/10.1214/ss/1177011136> doi: 10.1214/ss/1177011136
- Hoffman, K. A., Constantine, V. L., & Morse, D. L. (1989). Determination of absolute palaeointensity using a multi-specimen procedure. *Nature*, 339, 295–297.
- Königsberger, J. G. (1938). Natural residual magnetism of eruptive rocks. *Terrestrial Magnetism and Atmospheric Electricity*, 43(3), 299–320. doi: 10.1029/TE043i003p00299
- Kosareva, L. R., Kuzina, D. M., Nurgaliev, D. K., Sitdikov, A. G., Luneva, O. V., Khasanov, D. I., ... Spassov, S. (2020). Archaeomagnetic investigations in Bolgar (Tatarstan). *Stud. Geophys. Geod.*, 64(2), 255–292. doi: 10.1007/s11200-019-0493-3
- Krásá, D., Heunemann, C., Leonhardt, R., & Petersen, N. (2003). Experimental procedure to detect multidomain remanence during thellier–thellier experiments. *Phys. Chem Earth (A/B/C)*, 28(16), 681–687. (Paleo, Rock and Environmental Magnetism 2002) doi: 10.1016/S1474-7065(03)00122-0
- Lisé-Pronovost, A., Mallett, T., & Herries, A. I. R. (2020). Archaeointensity of nineteenth-century scottish firebricks from a foundry in melbourne, australia: comparisons with field models and magnetic observatory data. *Geological Society, London, Special Publications*, 497(1), 27–45. Retrieved from <https://sp.lyellcollection.org/content/497/1/27> doi: 10.1144/SP497-2019-72
- Nagata, T., Arai, Y., & Momose, K. (1963). Secular variation of the geomagnetic total force during the last 5000 years. *J. Geophys. Res.*, 68(18), 5277–5281. doi: 10.1029/j.2156-2202.1963.tb00005.x
- Nagy, L., Williams, W., Muxworthy, A. R., Fabian, K., Almeida, T. P., Conbhuí, P. Ó., & Shcherbakov, V. P. (2017). Stability of equidimensional pseudo-single-domain magnetite over billion-year timescales. *Proc. Natl. Acad. Sci. U.S.A.*, 114(39), 10356–10360. doi: 10.1073/pnas.1708344114
- Néel, L. (1949). Théorie du traînage magnétique des ferromagnétiques en grains fins avec applications aux terres cuites. *Ann. géophys.*, 5, 99–136.
- Paterson, G. A. (2011). A simple test for the presence of multidomain behavior during paleointensity experiments. *J. Geophys. Res.*, 116. doi: 10.1029/2011JB008369
- Paterson, G. A., Biggin, A. J., Yamamoto, Y., & Pan, Y. (2012). Towards the robust selection of Thellier-type paleointensity data: The influence of experimental noise. *Geochem. Geophys. Geosyst.*, 13(5). doi: 10.1029/2012GC004046
- Paterson, G. A., Tauxe, L., Biggin, A. J., Shaar, R., & Jonestrask, L. C. (2014). On improving the selection of thellier-type paleointensity data. *Geochem. Geophys. Geosyst.*, 15(4), 1180–1192. doi: 10.1002/2013GC005135
- Shaar, R., & Tauxe, L. (2013). Thellier\_gui: An integrated tool for analyzing paleointensity data from thellier-type experiments. *Geochem. Geophys. Geosyst.*, 14, 677–692. doi: 10.1002/ggge.20062
- Shaw, J. (1974). A new method of determining the magnitude of the paleomagnetic field application to 5 historic lavas and five archeological samples. *Geophys. J. R. astr. Soc.*, 39, 133–141.
- Tauxe, L., Santos, C., Cych, B., Zhao, X., Roberts, A., Nagy, L., & Williams, W. (2021). Understanding non-ideal paleointensity recording in igneous rocks: Insights from aging experiments on lava samples and the causes and consequences of 'fragile' curvature in arai plots. *Geochem. Geophys. Geosyst.*, 22, e2020GC009423. doi: 10.1029/2020GC009423
- Tauxe, L., Shaar, R., Jonestrask, L., Swanson-Hysell, N. L., Minnett, R., Koppers, A. a. P., ... Fairchild, L. (2016). PmagPy: Software package for paleomag-



- netic data analysis and a bridge to the magnetism information consortium (MagIC) database. *Geochem., Geophys., Geosyst.*, 17(6), 2450–2463. doi: 10.1002/2016GC006307
- Tauxe, L., & Yamazaki, T. (2015). Paleointensities. In M. Kono (Ed.), *Geomagnetism* (2nd Edition ed., Vol. 5, p. 461–509). Elsevier.
- Thébault, E., Finlay, C. C., Beggan, C. D., Alken, P., Aubert, J., Barrois, O., ... Zvereva, T. (2015). International Geomagnetic Reference Field: the 12th generation. *Earth Planets Space*, 67(1), 79. doi: 10.1186/s40623-015-0228-9
- Thellier, E., & Thellier, O. (1959). Sur l'intensité du champ magnétique terrestre dans le passé historique et géologique. *Ann. Geophys.*, 15, 285.
- Valet, J.-P., Brassart, J., Le Meur, I., Soler, V., Quidelleur, X., Tric, E., & Gillot, P.-Y. (1996). Absolute paleointensity and magnetomineralogical changes. *Journal of Geophysical Research: Solid Earth*, 101(B11), 25029–25044. Retrieved from <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/96JB02115> doi: <https://doi.org/10.1029/96JB02115>
- Williams, W., & Dunlop, D. J. (1989). Three-dimensional micromagnetic modelling of ferromagnetic domain structure. *Nature*, 337, 634–637.
- Yamamoto, Y., & Hoshi, H. (2008). Paleomagnetic and rock magnetic studies of the sakurajima 1914 and 1946 andesitic lavas from japan: A comparison of the ltd-dht shaw and thellier paleointensity methods. *Phys. Earth and Planet. Inter.*, 167, 118–143.
- Yamamoto, Y., Tsunakawa, H., & Shibuya, H. (2003). Palaeointensity study of the hawaiian 1960 lava: implications for possible causes of erroneously high intensities. *Geophys J Int*, 153(1), 263–276.
- Yamamoto, Y., & Yamaoka, R. (2018). Paleointensity study on the Holocene surface lavas on the Island of Hawaii using the Tsunakawa-Shaw method. *Front. Earth Sci.*, 6. doi: 10.3389/feart.2018.00048
- Yu, Y., Tauxe, L., & Genevey, A. (2004). Toward an optimal geomagnetic field intensity determination technique. *Geochem., Geophys., Geosyst.*, 5(2). doi: 10.1029/2003GC000630