

# Testing for Non-Gaussianity on Cosmic Microwave Background Radiation: A Review

Domenico Marinucci

*Abstract.* Cosmic microwave background (CMB) radiation can be viewed as a snapshot of the Universe 13 billion years ago, when it had 0.002% of its current age. A flood of data on CMB is becoming available thanks to satellite and balloon-borne missions, and a number of statistical issues have been raised consequently. A very relevant issue is the characterization of the statistical distribution of CMB and, in particular, procedures to test the assumption that the generating random field is Gaussian. Gaussianity tests are of fundamental importance both to validate statistical inference procedures and to discriminate between competing scenarios for the Big Bang dynamics. Several procedures have been proposed in the cosmological literature. This article is an attempt to provide a brief survey of developments in this area.

*Key words and phrases:* Cosmic microwave background radiation, Gaussianity, spherical random fields.

## 1. INTRODUCTION AND MOTIVATIONS

One hundred years ago we did not know how stars generate energy, the age of the Universe was thought to be only millions of years, and our Milky Way galaxy was the only galaxy known. Today, we know that we live in an evolving and expanding Universe comprising billions of galaxies, all held together by dark matter. With the hot Big Bang model, we can trace the evolution of the Universe from the hot soup of quarks and leptons that existed a fraction of a second after the beginning to the formation of galaxies a few billion years later, and finally to the Universe we see today, 13 billion years after the Big Bang, with its clusters of galaxies, superclusters, voids, and great walls. . . . As we enter the 21st century, a flood of observations is testing this paradigm (Turner and Tyson, 1999).

It has been extensively argued that a golden age is beginning for cosmology. This is due to an impressive combination of theoretical developments and experimental advances. In particular, a fundamental probe for the physics of primordial epochs is the analysis of cosmic microwave background (CMB) radiation, whose existence was theoretically predicted in the late 1940s and then observationally confirmed by Penzias and Wilson in 1964. Loosely speaking, CMB can be viewed as a snapshot of the Universe at the time of so-called *recombination*, which is reckoned to have occurred approximately 300,000 years after the Big Bang (i.e., at about 0.002% of the current age). The basic mechanism for CMB formation can be sketched as follows. There is now overwhelming evidence that the Universe evolved from a hot and dense state, where the temperature was so high that matter was completely ionized, that is, electrons were not bound in atoms and thus were scattered continuously with photons; during this stage, the Universe was therefore completely opaque to light. As the expansion went on, the temperature decreased steadily; the energy of photons hence dropped until it reached a critical value where it was no longer enough to keep the electrons and protons apart. Thus electrons were captured into orbits to form stable atomic nuclei, and photons started to move freely; as

---

*Domenico Marinucci is Associate Professor of Statistics, Department of Mathematics, University of Rome "Tor Vergata," 00133 Rome, Italy (e-mail: marinucc@mat.uniroma2.it).*

a first approximation, we can assume they reached the present time without any further interaction with matter. This stream of photons that permeates the Universe and reaches the Earth from every direction is the CMB; in this sense, the latter can be viewed as a signature of the distribution of matter and radiation in the very early Universe, and as such it is expected to yield very tight constraints on physical models for the Big Bang.

The photons are distributed over the various frequencies according to a *Planckian curve*, which characterizes the distribution of radiation emitted by a black body in thermal equilibrium. Thus, the amount of radiation at any given frequency uniquely identifies a corresponding temperature, which we can easily derive from CMB observations. This temperature has a mean value of about 2.73 K. Fluctuations around this value are labelled anisotropies in the cosmological literature; this terminology is to some extent misleading, because the CMB random field is indeed assumed to be statistically isotropic, that is, invariant in law under rotations, and the observed random fluctuations are just viewed as sample realizations of this random field. The magnitude of these fluctuations is on the order of  $10^{-5}$  times their mean value. Their existence was theoretically predicted in the late 1960s by a number of scientists (including the well-known Soviet physicist Andrei Sakharov) and first confirmed by the National Aeronautics and Space Administration (NASA) satellite mission COBE in the early 1990s (Smoot et al., 1992). Many, if not most, of the greatest challenges of current cosmological research relate to the analysis of the distributional properties of the random field of CMB fluctuations. The dependence structure of the field can be expressed in terms of the angular correlation function and its harmonic transform, the angular power spectrum. These statistics are expected to yield very tight constraints on fundamental issues such as the large scale geometry of the Universe (in turn determining its ultimate fate—perpetual expansion or recollapse into a Big Crunch), the existence and nature of nonordinary (baryonic) matter, the existence and nature of vacuum energy, the nature of physical interactions at the highest energies (Great Unification theories) and the shape and mechanisms of primordial fluctuations that have led to currently observed large scale structure. The ongoing NASA satellite mission WMAP (<http://map.gsfc.nasa.gov>) and the forthcoming European Space Agency (ESA) mission Planck (<http://astro.estec.esa.nl/index.php?project=PLANCK>) will probe CMB to an unprecedented accuracy, providing a flood of data on the order of several billion observations to address these issues. Preliminary

data sets are already available from balloon-borne missions, the most important so far being MAXIMA (<http://cosmology.berkeley.edu/group/cmb>) and BOOMERanG (<http://oberon.roma1.infn.it/boomerang>).

Of course, the correlation function and the angular power spectrum completely determine the distribution of a spherical random field if the field is Gaussian. By itself, this provides strong motivation for Gaussianity tests. There is, however, much more important motivation, concerning the mechanisms that underlie the dynamics of the Universe in the tiny fractions of seconds which followed the Big Bang. There are alternative models for the latter, the most popular of which is the inflationary scenario (see Peebles, 1993; Peacock, 1999), which predicts a Gaussian distribution for the density fluctuations. Alternative cosmological theories, such as topological defects or nonstandard inflationary models, predict other types of behavior, so Gaussianity checks are expected to play a crucial role in discriminating between competing scenarios. Non-Gaussianities may also have a nonphysical origin, that is, they might be generated by systematic errors in the CMB map, such as noise which has not been properly removed, contamination from the galaxy or distortions in the optics of the telescope. Of course, as mentioned above, a proper understanding of the distribution of fluctuations is also instrumental for making correct inference about the physical constants which can be estimated from CMB radiation. For these reasons, recent years have witnessed an enormous amount of attention, in the astrophysical and cosmological literature, to testing for Gaussianity of spherical random fields.

Our aim in this article is to provide a review of this literature, trying to point out directions for future research from the point of view of mathematical and computational statistics. An impressive amount of brilliant ideas has originated in the cosmological community. However, they mostly have been developed with surprisingly little connection to the existing probabilistic/statistical literature. An enormous amount of work is needed to assess their performance, not only from the point of view of analytic properties, but often even in terms of Monte Carlo simulations for size and power. The latter task only very recently has become feasible: the construction of numerical algorithms to simulate non-Gaussian spherical random fields of some physical relevance is by itself a very important topic for research. Our plan for this work is as follows. In Section 2, we review some basic facts and notation on spherical random fields; particular attention is spent on the discussion of the two possible representations

of the random field, that is, in real space or in harmonic space. Section 3 addresses methods and procedures which focus on the morphological properties of Gaussian fields, such as the distribution of their maxima, minima and saddle points, their local curvature and the length of the boundary between *hot* and *cold* regions. Section 4 is concerned with harmonic space methods, that is, those procedures based on harmonic transforms of higher order correlation functions and related ideas. Section 5 draws some conclusions and points out directions for future research. To make references more easily available for the statistical community, we have provided the web site locations of several preprints; in the cosmological community, such web sites are by far the most popular instruments for spreading scientific methods and results.

**2. SPHERICAL RANDOM FIELDS**

As stated in the previous section, CMB can be represented as a random field  $T(\theta, \varphi)$  indexed by the unit sphere  $S^2$ , that is, for each azimuth  $0 \leq \theta \leq \pi$  and elongation  $0 \leq \varphi < 2\pi$ ,  $T(\theta, \varphi)$  is a real random variable defined on some probability space. For notational simplicity, we assume that  $T(\theta, \varphi)$  has zero mean; physical motivations imply that it has finite variance and it is mean square continuous and isotropic, that is, its covariance is invariant with respect to the group of rotations. For such fields, it is useful to introduce the spherical harmonics, defined by

$$Y_{lm}(\theta, \varphi) = \begin{cases} \sqrt{\frac{2l+1}{4\pi} \frac{(l-m)!}{(l+m)!}} \cdot P_{lm}(\cos\theta) \exp(im\varphi), & \text{for } m > 0, \\ (-1)^m Y_{l,-m}^*(\theta, \varphi), & \text{for } m < 0, \end{cases}$$

where the asterisk denotes complex conjugation and  $P_{lm}(\cos\theta)$  denotes the associated Legendre polynomial of degree  $l, m$ , that is,

$$P_{lm}(x) = (-1)^m (1-x^2)^{m/2} \frac{d^m}{dx^m} P_l(x),$$

$$P_l(x) = \frac{1}{2^l l!} \frac{d^l}{dx^l} (x^2-1)^l,$$

$$m = 0, 1, 2, \dots, l, \quad l = 1, 2, 3, \dots$$

A detailed discussion of the properties of the spherical harmonics can be found in Liboff (1998, Chapter 9) and in Varshalovich, Moskalev and Khersonskii (1988, Chapter 5). The following spectral representation holds

in the mean square sense (see, e.g., Hannan, 1970; Wong, 1971; Leonenko, 1999):

$$(1) \quad T(\theta, \varphi) = \sum_{l=1}^{\infty} \sum_{m=-l}^l a_{lm} Y_{lm}(\theta, \varphi).$$

In (1), the triangular array  $\{a_{lm}\}$  represents a set of random coefficients, which can be obtained from  $T(\theta, \varphi)$  through the inversion formula

$$(2) \quad a_{lm} = \int_{-\pi}^{\pi} \int_0^{\pi} T(\theta, \varphi) Y_{lm}^*(\theta, \varphi) \sin\theta \, d\theta \, d\varphi,$$

$$m = 0, \pm 1, \dots, \pm l, \quad l = 1, 2, \dots$$

These coefficients are zero mean and uncorrelated; hence, if  $T(\theta, \varphi)$  is Gaussian, they have a complex Gaussian distribution, and they are independent over  $l$  and  $m \geq 0$  [although  $a_{l,-m} = (-1)^m a_{lm}^*$ ], with variance  $E|a_{lm}|^2 = C_l, m = 0, \pm 1, \dots, \pm l$ . The index  $l$  is usually labeled a multipole and, in principle, it runs from 1 to infinity. In any realistic experiment, however, there is an upper limit (which we denote by  $L$ ) on the multipoles we may observe, depending on the resolution of the experiment and the presence of noise;  $L$  is reckoned to be on the order of 500/800 for WMAP and 1500/2000 for Planck.

The sequence  $\{C_l\}$  denotes the angular power spectrum. We always assume that  $C_l$  is strictly positive, for all values of  $l$ . The angular power spectrum completely identifies the correlation structure of a spherical random field. Denote by  $\psi = \psi(\theta, \varphi; \theta', \varphi')$  the angle between  $(\theta, \varphi)$  and  $(\theta', \varphi')$ , and denote by  $\Gamma(\psi) = ET(\theta, \varphi)T(\theta', \varphi')$  the angular correlation function. We have

$$\Gamma(\psi) = \sum_{l=1}^{\infty} \sum_{l'=1}^{\infty} \sum_{m=-l}^l \sum_{m'=-l'}^{l'} E a_{lm} \bar{a}_{l'm'} Y_{lm}(\theta, \varphi) \cdot \bar{Y}_{l'm'}(\theta', \varphi')$$

$$= \sum_{l=1}^{\infty} C_l \sum_{m=-l}^l Y_{lm}(\theta, \varphi) \bar{Y}_{lm}(\theta', \varphi')$$

$$= \frac{1}{4\pi} \sum_{l=1}^{\infty} (2l+1) C_l P_l(\cos\psi),$$

where we have used the identity (Wong, 1971)

$$\sum_{m=-l}^l Y_{lm}(\theta, \varphi) \bar{Y}_{lm}(\theta', \varphi') = \frac{2l+1}{4\pi} P_l(\cos\psi).$$

Thus the angular spectrum can be viewed as the Legendre transform of the angular correlation function,

in obvious analogy with well-known results from the Fourier analysis of random processes. A natural estimator for  $\widehat{C}_l$  is

$$(3) \quad \widehat{C}_l = \frac{1}{2l+1} \sum_{m=-l}^l |a_{lm}|^2, \quad l = 1, 2, \dots,$$

which is clearly unbiased and consistent for  $l \rightarrow \infty$ ; see also Wasserman et al. (2001) and Miller, Nichol, Genovese and Wasserman (2002).

The preceding framework is appropriate for missions which provide a full-sky map of CMB radiation. This is true for the NASA satellite WMAP, for which data first were released on February 14, 2003 and much more data are anticipated over the next three years, and for the ESA satellite mission Planck, which is due to be launched in 2007 and will provide high resolution data for the following two or three years. Most of the data currently available, however, are provided by balloon-borne experiments, such as MAXIMA, DASI, BOOMERanG and ARCHEOPS. These experiments observe only a small fraction of the sky, a “patch” of a few degrees width (approximately square regions with sides on the order of 10–30°). It is thus not surprising that many articles so far have adopted the so-called flat sky approximation of CMB. The latter can be explained as follows. The observed regions can be considered to be small enough to be approximated by a plane tangent to the sphere. The CMB field is thus defined on an ordinary subset of  $R^2$  and the spherical harmonic representation is replaced by standard Fourier methods. The two approaches are usually considered broadly analogous and, indeed, many equivalence results in the limit where the scale goes to zero can be proved. Note, however, that some differences persist for any finite scale; for instance, the spherical harmonic coefficients are exactly orthogonal for a full-sky map, whereas Fourier transforms are not if we exclude the degenerate case of a white noise field.

From the practical point of view, it is also important to note that, although the parameter space is typically assumed to be continuous, the observations provided by any experiment are pixelized, that is, the observed field is discretized. In other words, only  $T(\vartheta_p, \varphi_p)$  on a finite grid  $p = 1, 2, \dots, N_p$  is observed. A rigorous treatment of these issues should include correct allowance for so-called aliasing effects, but, for brevity’s sake, we do not address these issues herein.

Although the literature on Gaussianity tests in a cosmological environment is fairly recent, it comprises such a high number of proposals that we cannot

aim at completeness. Our effort has been to divide these contributions into two main classes: (1) methods that search for non-Gaussianity by looking directly at  $T(\theta, \varphi)$  and (2) methods that search for non-Gaussianity by looking at the spherical harmonic coefficients  $a_{lm}$ . We label these approaches real-space and harmonic space methods, respectively, and we start our review with the former class.

### 3. REAL-SPACE METHODS

We start by considering methods which search for departures from Gaussianity in real space, that is, directly on the observed field  $T(\cdot)$ . Early attempts in this direction were made by Bond and Efstathiou (1987) and Vittorio and Juszkiewicz (1987): the aim of their work was to characterize the expected shape of “hot spots” (e.g., connected regions where the temperature is above some threshold level) and their correlation structure under Gaussianity. This approach was generalized in Novikov and Jørgensen (1996), where a detailed analysis of the morphological properties of two-dimensional Gaussian fields is provided. These authors were able to derive an explicit expression for the density of saddle points, the expected value and the variance of a Gaussian field around two neighboring maxima and the expected value of the length of a cluster of peaks. It was then suggested that these quantities could be exploited to test for Gaussianity, but no specific testing procedure was advocated.

More recently, Heavens and Sheth (1999) and Heavens and Gupta (2000) produced exact predictions for the correlation function of extrema (maxima and minima) of Gaussian random fields. More precisely, they defined a six-dimensional vector  $\mathbb{T}(\theta, \varphi)$  by

$$\mathbb{T}(\theta, \varphi) = (T, T_\theta, T_\varphi, T_{\theta\theta}, T_{\theta\varphi}, T_{\varphi\varphi})',$$

where

$$T_a = \frac{\partial T}{\partial a}, \quad a = \theta, \varphi$$

and

$$T_{ab} = \frac{\partial^2 T}{\partial a \partial b}, \quad a, b = \theta, \varphi.$$

Explicit expressions for these derivatives can be obtained from (1), which yields, for instance,

$$\frac{\partial T(\theta, \varphi)}{\partial \theta} = \sum_{l=1}^{\infty} \sum_{m=-l}^l a_{lm} \frac{\partial Y_{lm}(\theta, \varphi)}{\partial \theta},$$

with similar expressions for the other partial derivatives. Heavens and Gupta (2000) then considered the 12-dimensional Gaussian zero mean vector  $\mathbb{T} = (\mathbb{T}(\theta_1, \varphi_1)', \mathbb{T}(\theta_2, \varphi_2)')'$  for the covariance matrix for

which they provided an explicit expression. The correlation function of maxima can then be obtained by integrating numerically the joint density function, subject to the constraint that the two points are maxima, that is,

$$T_\theta = T_\varphi = 0, \quad T_{\theta\theta}, T_{\varphi\varphi} < 0, \quad T_{\theta\theta}T_{\varphi\varphi} - T_{\theta\varphi}^2 > 0.$$

A heuristic testing procedure was then suggested that compares the shape of the expected correlation function among maxima with the sample correlation function evaluated on realizations of Gaussian and non-Gaussian fields. Some exploratory analysis was then performed on the power of the test by applying this procedure over two non-Gaussian maps generated by *cosmic strings* models.

More recently, a related approach was advocated by Doré, Colombi and Bouchet (2003) (hereafter DCB). Choose first a threshold  $\nu$  and divide the sphere into two parts: hot regions, where the random field  $T$  passes the threshold, and cold regions, where  $T < \nu$ . The hot region is also called the excursion set of the field  $T$  over the threshold  $\nu$  (Adler, 1981) or the spherical measure of excess of level  $\nu$  (Leonenko, 1999, page 164). The insight behind the DCB approach is to consider the probability distribution not only of extrema, but of all points in the maps that exceed  $\nu$ , classified according to their local curvature and for varying values of the threshold. The relative abundances of the three types of points are then measured and compared with their Gaussian expected values, which were analytically derived. More precisely, and maintaining our previous notation, we consider the Hessian matrix  $\{T_{ab}\}_{a,b=\theta,\varphi}$ , which is real and symmetric: denote by  $\lambda_1$  and  $\lambda_2$  its eigenvalues, which are well known to be real. We call *hills* the points where both eigenvalues are negative, *lakes* the points where they are both positive and *saddles* the points where the two have opposite signs. It is obvious that the local maxima are a proper subset of the hills and the local minima are a proper subset of the lakes, whereas both the observed and the expected proportions of the three types of points must sum to unity. DCB then provided an expression for the probability that a point above a given threshold belongs to any of these three classes. We do not report their explicit result here, but we note its two most important features.

1. The probability depends on a single unknown parameter, related to the shape of the spectral density.
2. The asymptotic value of the probability for  $\nu \rightarrow \infty$  is independent of any nuisance parameter, that is, identical for any Gaussian field.

These nice properties are not shared, in general, by other statistics that focus only on extrema.

DCB then proposed a testing procedure based on these results. Write  $p(v) = (p_h(v), p_l(v))'$ , where  $p_h(v)$  and  $p_l(v)$  are the observed relative frequencies of hills and lakes among pixels above the threshold  $\nu$ . As discussed before, their expected value  $Ep(v)$  is known analytically under Gaussianity, up to a single nuisance parameter to be estimated from the data. Then consider a grid of threshold values  $\nu_i$ ,  $i = 1, 2, \dots, k$ , and compute the chi-square statistic

$$(4) \quad \sum_{i=1}^k \{p(\nu_i) - Ep(\nu_i)\}' \Omega^{-1}(\nu_i) \{p(\nu_i) - Ep(\nu_i)\}.$$

Here, the variance-covariance matrix  $\Omega(\nu_i)$  is defined by

$$\Omega(\nu_i) = E\{p(\nu_i) - E(p(\nu_i))\}\{p(\nu_i) - E(p(\nu_i))\}'.$$

In principle,  $\Omega(\nu_i)$  can be derived analytically; the computations are hard, however, and the authors found it more convenient to resort to Monte Carlo simulations. It was then argued that the distribution of (4) should be well approximated by a chi square with  $k$  degrees of freedom. Although some exploratory simulation results in the article showed a good agreement with this assumption in the Gaussian case, note that the addends in (4) are not independent; indeed, they become closer and closer for higher  $k$ , so that some caution is needed with this approximation. With this caveat in mind, some exploratory analysis suggests that this method can be a useful probe of non-Gaussianity. In particular, DCB considered a linear mixture between a Gaussian map and a non-Gaussian map simulated according to a cosmic string model; that is, they considered a map formed by

$$T(\theta, \varphi) = \beta T^G(\theta, \varphi) + (1 - \beta) T^{\text{CS}}(\theta, \varphi), \quad \beta \in [0, 1],$$

where  $T^G(\theta, \varphi)$  and  $T^{\text{CS}}(\theta, \varphi)$  represent, respectively, a Gaussian and a cosmic string field. DCB were then able to show that a significant value of the chi-square test (for  $\alpha = 1\%$ ) is obtained at  $\beta = 0.7$ , even for an experiment that observes less than 1% of the full sky. As expected, the  $p$  value of the test rises steadily as  $\beta$  decreases to zero. These results refer, however, to a single map; a proper Monte Carlo evaluation of power has not yet been provided for this as has been for related methods.

Martinez-González et al. (2000) considered two different methods to test for Gaussianity in real space. On one hand, they analyzed the curvature and the

eccentricity of regions around a local maximum of the field. Simulation results, however, suggest that such a procedure might have very little power against non-Gaussian alternatives. Alternatively, they focused on the partition function, introduced in the cosmological literature by Diego et al. (1999). The partition function is defined as follows. Divide the CMB map into a grid of cells, each of size  $\delta$ , and denote by  $N_{\text{boxes}}(\delta)$  the number of such cells needed to cover the map; these cells are labelled by  $i$ . Define a measure  $\mu_i(\delta)$  on each cell, for instance, by  $\mu_i(\delta) = \sum_{p \in \text{cell}_i} |T(\vartheta_p, \varphi_p)|$ . The partition function is then defined as

$$Z(q, \delta) = \sum_{i=1}^{N_{\text{boxes}}(\delta)} \mu_i(\delta)^q.$$

The idea is that variations over  $\delta$  should be sensitive to the dependence structure of the map, whereas fluctuations over  $q$  should help to detect skewness, kurtosis and other non-Gaussian features of the marginal distributions. Again, some exploratory analysis is provided on simulated maps, but a full Monte Carlo study is not yet available.

Among real-space methods, probably the most popular set of non-Gaussianity tests relates to Minkowski functionals, which were introduced into the cosmologists' toolbox by Mecke, Buchert and Wagner (1994). We mainly follow the presentation of Winitzki and Kosowsky (1998) and Novikov, Schmalzing and Mukhanov (2000); for a probabilistic reference, see Stoyan, Kendall and Mecke (1987). To analyze a spherical map in terms of Minkowski functionals, we consider again the excursion sets, that is, the map subsets which exceed a given threshold value. The threshold is labelled  $v$ , as before, and it is treated as an independent variable on which these functionals depend. The three functionals of interest then are, up to irrelevant constant factors, the following.

1. *Area*:  $M_0(v)$  is the total area of all hot regions, that is, points on a sphere where  $T(\theta, \varphi) > v$ .
2. *Boundary length*:  $M_1(v)$  is proportional to the total length of the boundary between cold and hot regions.
3. *Euler characteristic*:  $M_2(v)$ , a purely topological quantity, counts the number of isolated hot regions minus the number of isolated cold regions.

The rationale behind these statistics can be explained from mathematical results in Hadwiger (1959). The idea is to use a completeness theorem to characterize the morphological properties of a spherical map. Here,

by morphological we mean the properties which are invariant under translations and rotations, and which are additive. Denote by  $K^d$  the class of convex, compact sets in  $R^d$ , and consider the class of maps which satisfy the properties of motion invariance and additivity, that is,

$$T(gK) = T(K),$$

$$T(K_1 \cup K_2) + T(K_1 \cap K_2) = T(K_1) + T(K_2),$$

where  $g$  belongs to the groups of rigid motions (i.e., rotations and translations) on  $R^d$  and  $K, K_1, K_2, K_1 \cup K_2$  are assumed to be convex sets. Then  $T$  can be expressed as a linear combination of Minkowski functionals:

$$T(K) = \sum_{i=0}^d \alpha_i M_i(K), \quad \alpha_i \in R.$$

This can be interpreted by stating that all the morphological information of a convex body is contained in the Minkowski functionals (Winitzki and Kosowsky, 1998).

For Gaussian random fields, the expected values of the three Minkowski functionals can be derived in analytic form. Assuming for simplicity that units are chosen to give a unit variance, and denoting by  $\Phi(\cdot)$  the cumulative distribution function of a standard Gaussian variate, they are

$$(5) \quad EM_0(v) = 1 - \Phi(v),$$

$$(6) \quad EM_1(v) = \frac{\sqrt{\tau}}{8} \exp\left(-\frac{v^2}{2}\right),$$

$$(7) \quad EM_2(v) = \frac{\tau}{\sqrt{8\pi^3}} v \exp\left(-\frac{v^2}{2}\right),$$

where

$$\tau = \sum_{l=1}^{\infty} (2l+1) C_l \frac{l(l+1)}{2}.$$

An immediate consequence of (5)–(7) is that, although the expected value of the first Minkowski functional is invariant with respect to the dependence structure of  $T(\cdot)$ , for the second and third Minkowski functionals this is not the case and calibration for a given angular power spectrum  $C_l$  is needed. In practice, the latter is unknown and must be estimated from the data. Moreover, even for the first Minkowski functional, knowledge of the angular power spectrum is required for a Monte Carlo evaluation of its sample variance.

This is clearly an important drawback of these methods, and significant effort has been undertaken to provide at least some crude upper bound for the functionals' variance (Winitzki and Kosowsky, 1998).

We note also that if we define, as is usually done, the empirical distribution function by

$$F_{N_p}(v) = \frac{1}{N_p} \sum_{p=1}^{N_p} \mathbb{I}(T(\theta_p, \varphi_p) \leq v),$$

then

$$M_0(v) = 1 - F_{N_p}(v).$$

Therefore, it seems possible to derive some approximation to the asymptotic distribution of the first Minkowski functional by taking into account the rich literature on empirical processes. For instance, because the CMB temperature field is known to be characterized by a slow decay of autocovariances (long range dependence), it seems potentially useful to expand the first Minkowski functional into Hermite polynomials and derive uniform approximations for its limiting behavior; see, for instance, Dehling and Taqqu (1989) and Doukhan, Lang and Surgailis (2002). Some attention must be paid, however, to the interpretation of the indices: in the usual framework, an increase in the number of observations  $N_p$  yields nice ergodic theorems and convergence results. Here, as  $N_p$  grows, we are simply sampling, on a finer grid, the same mean-square continuous field: although, of course, as the information grows, no ergodic property is applicable, so some care is needed when appealing to asymptotic results. This is a typical problem for any attempt at an asymptotic theory with CMB statistical analysis.

Many articles have applied Minkowski functionals to experimental data. References include Kogut et al. (1996) for COBE data, Polenta et al. (2002) for data from the BOOMERanG experiment, and Komatsu et al. (2003) for the first release from the WMAP mission; none of these works has detected significant non-Gaussianities. Note, however, that the most interesting (and popular) non-Gaussian models predict deviations from Gaussianity far too low to be detectable from currently available data, so it is safe to wait at least until the Planck release (expected in 2008/2009) before any definite conclusion is drawn.

The Minkowski functionals, and the other approaches so far considered, can be performed without any assumption on the distribution of the random field under the alternative. Less popular, but still widely diffused among cosmologists are approaches where

a specific alternative is taken into account. Phillips and Kogut (2001), for instance, considered a class of pseudo-likelihood ratio tests defined as follows. Let  $S = (S_1, \dots, S_n)'$  be a set of statistics, for instance, the values of one among the Minkowski functionals, evaluated at threshold  $v_i$ ,  $i = 1, \dots, n$ . Denote its variance-covariance matrix under the assumption  $H$  by

$$\Omega_H = E_H\{S - E_H(S)\}\{S - E_H(S)\}'$$

and consider the normalized (pseudo) chi-square statistic defined by

$$\chi_H^2 = \{S - E_H(S)\}' \Omega_H^{-1} \{S - E_H(S)\}.$$

A pseudo-likelihood ratio test for the two competing hypotheses  $H_0$  vs.  $H_1$  is then obtained by considering

$$\text{LR}(H_0, H_1) = \left\{ \frac{\det(\Omega_{H_0})}{\det(\Omega_{H_1})} \right\}^{1/2} \exp\left\{ -\frac{1}{2}(\chi_{H_1}^2 - \chi_{H_0}^2) \right\}.$$

Under non-Gaussian circumstances, the expected values  $\Omega_{H_1}$  and  $E_{H_1}(S)$  are usually very hard to derive analytically, and their numerical values obtained from simulations depend very heavily on the chosen specifications. As such, this approach is highly model dependent and very demanding from a computational point of view. Phillips and Kogut (2001) took  $H_0$  to represent a standard Gaussian model, joint with some pre-specified angular power spectra, and took  $H_1$  to entail a so-called *texture model* alternative. They then focused on the third Minkowski functional and on the correlation function for extrema under the two competing alternatives. These authors did not attempt to derive threshold values for  $\text{LR}(H_0, H_1)$  by a Monte Carlo experiment, but rather they performed some exploratory analysis in which they claim a detection every time a value higher than unity is obtained. The results are interesting, but it is difficult to assess the size and power of this procedure and its robustness against departures from the joint assumptions in  $H_0$  and  $H_1$ .

To conclude this section, we address one more practical issue. As argued extensively, CMB can be characterized as a random field on the unit sphere. In real-world experiments, however, the data are provided as a time series of signal plus noise. If we denote by  $(\theta_t, \varphi_t)$  the pixel observed at time  $t$  and we take  $e_t$  to represent a zero mean, covariance stationary noise sequence, we have that the time series of the data can be written as  $\mathbf{d} = \mathbf{T} + \mathbf{e}$ , where  $\mathbf{d} = (d_1, \dots, d_N)'$ ,  $\mathbf{T} = (T(\theta_1, \varphi_1), \dots, T(\theta_N, \varphi_N))'$  and  $\mathbf{e} = (e_1, \dots, e_N)$ ,  $N$  denoting the total number of observations. The series  $\mathbf{T}$  and  $\mathbf{e}$  are assumed to be uncorrelated; it is customary to label their corresponding covariance matrices as  $\mathbf{S} = E\mathbf{T}\mathbf{T}'$ ,  $\mathbf{N} = E\mathbf{e}\mathbf{e}'$ , and  $\mathbf{C} = \mathbf{S} + \mathbf{N} = E\mathbf{d}\mathbf{d}'$ .

Of course, we are interested in the Gaussianity of the signal, which can be disguised by the distribution properties of noise. To overcome this problem, Wu et al. (2001) used principal components as follows. Assume that  $\mathbf{S}$  and  $\mathbf{N}$  are known, the former depending on the angular power spectrum of the field, the latter on the spectral density of noise. In reality, in both cases, there are at least many parameters to be estimated (for noise, see, e.g., Natoli et al., 2002), but we neglect this issue as a first order approximation. Then we form  $\widehat{\mathbf{d}}_i = h'_i \mathbf{d}$ , where for  $i = 1, 2, \dots, N$ ,  $h_i$  represents the eigenvector that corresponds to the  $i$ th highest eigenvalue of the matrix  $\mathbf{N}^{-1/2} \mathbf{C} \mathbf{N}^{-1/2}$ . Wu et al. (2001) applied these ideas to data from the balloon experiment MAXIMA, and focused only on eigenvalues larger than unity, thus considering only 639 modes out of almost 6000 observations. It is then possible to analyze the new sample  $\widehat{\mathbf{d}}_i$ ,  $i = 1, 2, \dots, 639$ . By construction, the elements of the sample are uncorrelated, that is, independent under Gaussianity. On these observations we can then perform standard tests for independent and identically distributed (i.i.d.) samples, for instance, a Kolmogorov–Smirnov statistic. By the same rationale, Wu et al. (2001) considered other related procedures, such as Wiener filtering of the CMB signal from a noisy time series prior to implementing Gaussianity tests.

Of course, the effect of noise is an important issue for any procedure which is meant to be useful in practical applications. It is important to be cautious, however, when performing linear transformations on data prior to implementing a Gaussianity test. These transformations, in fact, can erase, to various degrees, non-Gaussian components due to the central limit theorem. We return to this issue when we discuss harmonic space methods in the following section.

#### 4. HARMONIC SPACE METHODS

Probably the single most popular approach to the detection of non-Gaussian signatures relates to the *bispectrum* of the CMB random field. The bispectrum can be viewed as the harmonic transform of the three-point correlation function, much as the angular power spectrum is the Legendre transform of the (two-point) angular correlation function. More precisely, write  $\Psi_i = (\theta_i, \varphi_i)$  for  $i = 1, 2, 3$ . We have

$$\begin{aligned}
 & ET(\Psi_1)T(\Psi_2)T(\Psi_3) \\
 (8) \quad &= \sum_{l_1, l_2, l_3} \sum_{m_1, m_2, m_3} B_{l_1 m_1 l_2 m_2 l_3 m_3} Y_{l_1 m_1}(\Psi_1) \\
 &\quad \cdot Y_{l_2 m_2}(\Psi_2) Y_{l_3 m_3}(\Psi_3),
 \end{aligned}$$

where the bispectrum is given by

$$(9) \quad B_{l_1 m_1 l_2 m_2 l_3 m_3} = E(a_{l_1 m_1} a_{l_2 m_2} a_{l_3 m_3}).$$

By symmetry of the Gaussian distribution, both (8) and (9) are clearly equal to zero for centered fields. Moreover, the assumption that the CMB random field is statistically isotropic entails a number of constraints on  $B_{l_1 m_1 l_2 m_2 l_3 m_3}$ . Indeed, the latter must ensure that the three-point correlation function on the left-hand side of (8) remains unchanged if the three directions  $\Psi_1, \Psi_2$  and  $\Psi_3$  are rotated by the same angle. Careful choices of the orientations allow us to prove that the angular bispectrum of an isotropic field can be nonzero only if the following conditions are met:

- 4(a)  $l_1, l_2$  and  $l_3$  satisfy the triangle rule,  $l_i \leq l_j + l_k$ , for all choices of  $i, j, k = 1, 2, 3$ .
- 4(b)  $l_1 + l_2 + l_3 = \text{even}$ .
- 4(c)  $m_1 + m_2 + m_3 = 0$ .

More generally, it can be shown that a necessary and sufficient condition for  $B_{l_1 m_1 l_2 m_2 l_3 m_3}$  to represent the angular bispectrum of an isotropic random field is that there exists a real symmetric function of  $l_1, l_2, l_3$ , which we denote  $b_{l_1 l_2 l_3}$ , such that we have the identity

$$(10) \quad B_{l_1 m_1 l_2 m_2 l_3 m_3} = \mathcal{G}_{l_1 l_2 l_3}^{m_1 m_2 m_3} b_{l_1 l_2 l_3}.$$

In (10) we used the Gaunt integral  $\mathcal{G}_{l_1 l_2 l_3}^{m_1 m_2 m_3}$ , defined as

$$\begin{aligned}
 \mathcal{G}_{l_1 l_2 l_3}^{m_1 m_2 m_3} &= \int_{S^2} Y_{l_1 m_1}(\Psi) Y_{l_2 m_2}(\Psi) Y_{l_3 m_3}(\Psi) d\Psi \\
 &= \left( \frac{(2l_1 + 1)(2l_2 + 1)(2l_3 + 1)}{4\pi} \right)^{1/2} \\
 &\quad \cdot \begin{pmatrix} l_1 & l_2 & l_3 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} l_1 & l_2 & l_3 \\ m_1 & m_2 & m_3 \end{pmatrix},
 \end{aligned}$$

where the  $\begin{pmatrix} l_1 & l_2 & l_3 \\ m_1 & m_2 & m_3 \end{pmatrix}$  represent the ‘‘Wigner 3j symbols’’; see Varshalovich, Moskalev and Khersonskii (1988, Chapter 8). It can be shown that the Gaunt integral is identically equal to zero unless conditions 4(a)–(c) are fulfilled.

In practice, of course, the bispectrum is not observable. The statistics which are usually considered are the (observed) angle averaged bispectrum, defined as

$$\begin{aligned}
 \widehat{B}_{l_1 l_2 l_3} &= \sum_{m_1=-l_1}^{l_1} \sum_{m_2=-l_2}^{l_2} \sum_{m_3=-l_3}^{l_3} \begin{pmatrix} l_1 & l_2 & l_3 \\ m_1 & m_2 & m_3 \end{pmatrix} \\
 &\quad \cdot (a_{l_1 m_1} a_{l_2 m_2} a_{l_3 m_3}),
 \end{aligned}$$



and the (observed) reduced bispectrum, defined as

$$\left( \frac{(2l_1 + 1)(2l_2 + 1)(2l_3 + 1)}{4\pi} \right)^{1/2} \cdot \begin{pmatrix} l_1 & l_2 & l_3 \\ 0 & 0 & 0 \end{pmatrix} \widehat{b}_{l_1 l_2 l_3} = \widehat{B}_{l_1 l_2 l_3},$$

where the following useful identity is used:

$$\sum_{m_1=-l_1}^{l_1} \sum_{m_2=-l_2}^{l_2} \sum_{m_3=-l_3}^{l_3} \begin{pmatrix} l_1 & l_2 & l_3 \\ m_1 & m_2 & m_3 \end{pmatrix}^2 = 1.$$

The bispectrum has a number of nice characteristics, which suggest it is potentially important as a probe of non-Gaussianity. In particular, for some important non-Gaussian models, it has been possible to evaluate analytically the form of the bispectrum up to finitely many parameters to be estimated, so that a model fit approach can be used to test these types of non-Gaussianity. A particularly relevant contribution in this area is the article by Komatsu and Spergel (2001), hereafter denoted KS. As in much other related work, in this article the amount of non-Gaussianity is measured by a parameter  $f_{NL}$ , such that  $f_{NL} \times 10^{-5}$  is the ratio between the non-Gaussian and the Gaussian components of the map (which are not independent, though). For a specific, physically motivated non-Gaussian model, KS were then able to derive an explicit formula for the bispectrum in terms of parameters which can be estimated from the angular power spectrum and  $f_{NL}$ ; we label this predicted bispectrum  $B_{l_1 l_2 l_3}^P$ . In an attempt to provide a much more faithful representation of realistic experiments, they also provided an analytic expression for the bispectrum for various kinds of noise components which can partially contaminate the cosmological signal, such as, for instance, instrumental noise, dust and foreground sources (astrophysical objects covering the CMB); we label this component  $B_{l_1 l_2 l_3}^N$ . They went on to suggest the regression model

$$(11) \quad \widehat{B}_{l_1 l_2 l_3} = A_1 B_{l_1 l_2 l_3}^P + A_2 B_{l_1 l_2 l_3}^N + e_{l_1 l_2 l_3},$$

where  $e_{l_1 l_2 l_3}$  denotes a regression residual. Provided the assumptions on the non-Gaussian components are correct, we have  $A_1 = A_2 = 1$ . The idea is then to implement a regression-based test of the latter assumption, performed on the weighted least square estimates

$$(12) \quad \widehat{A}_1, \widehat{A}_2 = \arg \min_{A_1, A_2} \left\{ \sum_{l_1, l_2, l_3 \in \beta(l_1, l_2, l_3)} \left( \widehat{B}_{l_1 l_2 l_3} - A_1 B_{l_1 l_2 l_3}^P - A_2 B_{l_1 l_2 l_3}^N \right)^2 / (w_{l_1 l_2 l_3}) \right\},$$

where the weights  $w_{l_1 l_2 l_3}$  are based on an approximate estimate of the variance of the residual term in (11); here, we write  $\beta(l_1, l_2, l_3)$  for the set of  $(l_1, l_2, l_3)$  which fulfills conditions 4(a) and (b). As we mentioned in the Introduction, in the analysis of real data the highest observable multipole must be considered equal to a finite number  $L$ , which depends on the resolution of the experiment. From physical arguments, KS argued that even for an ideal experiment, no cosmological information should be expected over the value  $L \sim 3000$ .

By an approximate evaluation of the Fisher matrix for this regression model and by numerical implementation for a given set of parameters, KS concluded that a high resolution experiment like Planck may be able to distinguish non-Gaussian features down to  $f_{NL} \sim 5$ . This result rests on a number of specific assumptions and approximations, and it is an important research topic to validate it either theoretically or by means of a proper Monte Carlo experiment. Also of great importance is the analysis of the robustness properties of this model-fit approach in the presence of some misspecification either in  $B_{l_1 l_2 l_3}^P$  or  $B_{l_1 l_2 l_3}^N$  and taking into account the extra variance due to estimated parameters. From a computational point of view, even a single estimate of the parameters  $\widehat{A}_1$  and  $\widehat{A}_2$  is very demanding, because  $O(L^3)$  bispectrum values need to be calculated: a Monte Carlo study is completely out of reach on a standard computer with the current techniques, because the required CPU time is on the order of several years. Indeed, when the WMAP data were made available in February, 2003, KS decided not to adopt this methodology, because they found it numerically infeasible. They thus implemented a Minkowski functionals approach on real data (see the previous section) and a more heuristic method in harmonic space that has statistical properties that call for further discussion (Komatsu et al., 2003). The main result is that available data suggest that  $f_{NL}$  should not be larger than 150. Much tighter constraints are needed, however, to bridge the gap between theory and observations, because theoretical physicists expect values of  $f_{NL}$  on the order of unity or below.

The approach by KS can be considered a parametric attempt to identify non-Gaussian components. As such, it is clearly very powerful against specific departures from the Gaussian assumption. On the other hand, this approach relies heavily on specific conditions on the physical mechanisms to generate non-Gaussianities and on a number of other details, including the values of many unknown physical parameters. It is thus also important to consider methods which impose little

or no assumptions on the form of the non-Gaussian alternative. The bispectrum has been widely adopted in this context as well. Note first that, under Gaussianity, the distribution of  $a_{lm}/C_l^{1/2}$  does not depend on any nuisance parameter. Hence, the bispectrum can be easily made model-independent, namely we can focus on the normalized bispectrum defined by

$$I_{l_1 l_2 l_3} = \frac{\widehat{B}_{l_1 l_2 l_3}}{\sqrt{\widehat{C}_{l_1} \widehat{C}_{l_2} \widehat{C}_{l_3}}}.$$

It is not difficult to show that  $\text{Var}\{I_{l_1 l_2 l_3}\} = \Delta_{l_1 l_2 l_3}$ , where  $\Delta_{l_1 l_2 l_3} = 1, 2, 6$  for the cases where the three indexes are all different, two of them are the same, and all three of them are equal, respectively. In practice, of course,  $I_{l_1 l_2 l_3}$  is infeasible and it is usually replaced by the observable statistic

$$\widehat{I}_{l_1 l_2 l_3} = \frac{\widehat{B}_{l_1 l_2 l_3}}{\sqrt{\widehat{C}_{l_1} \widehat{C}_{l_2} \widehat{C}_{l_3}}};$$

see (3). A possible way to test for Gaussianity is then to compare the absolute value of  $\widehat{I}_{l_1 l_2 l_3}$  with threshold levels in the Gaussian case obtained, for instance, by Monte Carlo simulations. This idea was pursued, for instance, by Komatsu and Spergel (2001), who performed  $i = 1, 2, \dots, 50,000$  realizations of Gaussian fields and evaluated  $\widehat{I}_{l_1 l_2 l_3, i}$  for each realization. Then the normalized bispectrum  $\widehat{I}_{l_1 l_2 l_3, \text{COBE}}$  also was estimated from the COBE data, thus obtaining the empirical  $p$  values,

$$P_{l_1 l_2 l_3} \equiv \frac{\sum_{i=1}^{50,000} \mathbb{I}\{\widehat{I}_{l_1 l_2 l_3, i} > \widehat{I}_{l_1 l_2 l_3, \text{COBE}}\}}{50,000}.$$

Note that this approach is computationally feasible because the resolution of COBE was approximately only  $7^\circ$ , so the highest observed multipole was approximately  $L \sim 25$  (rather than 500/750 for WMAP or 1500/2000 for Planck). No detection of non-Gaussianity was claimed, but again it is difficult to assess the power of this procedure, especially for the low resolution experiments for which it is computationally feasible.

Other authors (e.g., Ferreira, Magueijo and Górski, 1998) have used chi-square statistics such as

$$(13) \quad \left[ \frac{2}{L} \right] \sum_{l \in \beta(l, l, l)} \frac{\widehat{I}_{lll}^2}{\Delta_{lll}},$$

where  $[\cdot]$  denotes the integer part:  $[2/L]$  is just the number of elements in the sum, as  $l \in \beta(l, l, l)$  if and only if it is even. Ferreira, Magueijo and Górski (1998)

argued that better statistical properties, in terms of both size and power, can be obtained by resorting to the related statistic

$$\left[ \frac{2}{L} \right] \sum_{l \in \beta(l, l, l)} \{-2 \log \Phi_l(\widehat{I}_{lll}) + c_l\}^2,$$

where  $\Phi_l(\cdot)$  approximates the distribution function of  $\widehat{I}_{lll}$  and the  $c_l$  are chosen to make the summands mean zero. For a low number of multipoles, threshold values can again be derived by Monte Carlo simulations. Note that we consider only bispectrum coefficients with  $l_1 = l_2 = l_3$ , an assumption that can be relaxed only at the cost of an enormous increase in computational complexity. This was the first method to detect a significant non-Gaussianity on experimental data, more precisely, on the observations provided by the pioneering experiment COBE (Ferreira, Magueijo and Górski, 1998). This finding was especially remarkable, because the data from COBE are much less informative than other data from subsequent experiments (as mentioned before, the observations are considered reliable only up to multipoles on the order of 20, much less than the few hundreds allowed for instance by WMAP) and the resulting level of non-Gaussianity was therefore several orders of magnitude larger than expected from the most popular physical scenarios. In fact, this result was not confirmed by alternative experiments and many subsequent articles, for instance, Banday, Zaroubi and Górski (2000), have argued that the detection of a non-Gaussianity in such a case was merely due to systematic errors in the data.

Winitzki and Wu (2000) entertained an analytic investigation into the properties of the bispectrum at small scales, that is, as  $l \rightarrow \infty$ . Their results are presented in terms of the flat sky approximation and Fourier, rather than Legendre, transforms, but for our purposes they can be summarized as follows. For high  $l$ , the sample bispectra  $\widehat{B}_{l_1 l_2 l_3}$  of a Gaussian field are themselves approximately Gaussian and independent variables for different values of at least one of the three indexes  $l_1, l_2, l_3$ . Therefore, one can treat them as a sample of i.i.d. Gaussian observations and perform standard Gaussianity tests based on skewness, kurtosis or joint higher order moments. The article by Winitzki and Wu is analytically sophisticated and presents some exploratory analysis in non-Gaussian circumstances; the evidence provided on the power of this procedure, however, is somewhat mixed and certainly deserves more investigation.

The joint distribution of the spherical harmonic coefficients for a Gaussian map is also the starting point

for an approach proposed by Marinucci and Piccioni (2004) and further developed by Hansen, Marinucci, Natoli and Vittorio (2002) and Hansen, Marinucci and Vittorio (2003). The idea can be explained as follows. Let us drop, for brevity's sake, the  $m = 0$  term and consider for each  $l$ , the sequences

$$\zeta_{lm} = \frac{2|a_{lm}|^2}{C_l}, \quad \eta_{lm} = \frac{\operatorname{Re}(a_{lm})}{\operatorname{Im}(a_{lm})}, \quad m = 1, 2, \dots, l.$$

Under Gaussianity, it is not difficult to see that  $\{\zeta_{lm}\}$  is a sequence of i.i.d. exponential random variables, whereas  $\{\eta_{lm}\}$  is an i.i.d. Cauchy sequence. The two sequences represent the random (normalized) amplitudes and phases of  $a_{lm}$ , respectively. It is easy to implement a Smirnov transformation and make both sequences uniform on  $[0, 1]$ :

$$u_{lm} = 1 - \exp(-\zeta_{lm}), \quad v_{lm} = \frac{1}{2} + \frac{1}{\pi} \arctan(\eta_{lm}).$$

It is then natural to propose (possibly multivariate) empirical processes techniques, for instance, considering

$$\begin{aligned} & K_L^c(\alpha_1, \dots, \alpha_k; r) \\ &= \frac{1}{\sqrt{L}} \sum_{l=1}^{[Lr]} \left\{ F_l^c(\alpha_1, \dots, \alpha_k) - \prod_{j=1}^k \alpha_j \right\}, \quad c = \eta, \zeta, \end{aligned}$$

where

$$\begin{aligned} & F_l^c(\alpha_1, \dots, \alpha_k) \\ &= \frac{1}{l-k+1} \sum_{m=1}^{l-k+1} \mathbb{I}(c_{lm} \leq \alpha_1, \dots, c_{l,m+k-1} \leq \alpha_k), \\ & \quad c = \eta, \zeta. \end{aligned}$$

The asymptotic behavior of  $K_L^c(\alpha_1, \dots, \alpha_k; r)$  can be established rigorously under Gaussianity, and it is immediate to use such results to propose Cramér–Von Mises or Kolmogorov–Smirnov tests. As for other harmonic space methods, a nice feature of such tests is that they provide some information not only on departures from Gaussianity, but also on its location in the multipole space. There are of course alternative methods by which the same procedure can be implemented. For instance, there is no need to restrict attention only to adjacent  $m$ 's; indeed it is useful also to consider the joint distribution of  $\eta_{lm}$ ,  $\zeta_{lm}$  at different multipoles (i.e., for different values of  $l$ ). It is also worth noticing that  $\zeta_{lm}$  is infeasible, because the angular power spectrum is unknown. It is, however, possible to obtain a feasible statistic by replacing  $C_l$  with  $\widehat{C}_l$ . The effect of introducing an estimated angular power

spectrum is not negligible, however, in the sense that both the asymptotic expected value and the variance of  $K_L^c(\alpha_1, \dots, \alpha_k; r)$  are affected (see Marinucci and Piccioni, 2004, for  $k = 1$ , and Hansen et al., 2002, for general  $k$ ). The power properties of these tests also have been investigated: the general outcome is that the power of the test is usually (much) larger for  $k = 3$  (especially) and  $k = 2$  than for the univariate case ( $k \geq 4$  is almost infeasible from a computational point of view). A possible heuristic explanation is as follows. For some non-Gaussian models, the marginal distributions of the  $a_{lm}$ 's can be close to Gaussian, especially for high values of  $l$ , because of central limit theorem-like arguments. Note that the spherical harmonic coefficients are linear transforms of the observed map; see (2). This is basically the same sort of problem that we mentioned in the previous section for filtered maps in real space. On the other hand, the joint assumption that the  $a_{lm}$  are Gaussian and independent uniquely identifies a Gaussian random field. Therefore, it is natural to expect that procedures for probing this joint assumption should have satisfactory power properties against a variety of alternatives: see Hansen et al. (2002, 2003) for more discussion on this point. The same references also address a number of practical issues, such as the presence of noise and gaps in the map. The idea to probe non-Gaussianity by investigating the dependence structure of the random phases  $\{\eta_{lm}\}$  also was considered by Chiang, Naselsky and Coles (2001) and Chiang, Naselsky, Verkhodanov and Way (2003).

Many other articles have focused on harmonic methods to detect non-Gaussian features in the CMB. A very technical contribution by Hu (2001) discussed the general form of the four-point correlation function and its harmonic transform, the trispectrum. Particular care is devoted to the determination of an explicit form for the trispectrum and the analysis of its signal-to-noise properties for forthcoming experiments. A major issue in this area is the possibility of making such procedures computationally feasible, given the enormous amount of data provided by satellite missions and the difficulties in implementing higher order harmonic transforms.

## 5. FINAL REMARKS

No definite conclusion can be drawn presently on the relative performance of harmonic space and real-space methods. It is difficult to perform comparisons, because the various procedures require quite different inputs; for instance, all real-space methods require a

pretty accurate knowledge of the angular power spectrum to derive threshold values under the null hypothesis of Gaussianity. This is not the case for most harmonic space methods, which can be performed on normalized coefficients and are thus free of nuisance parameters under the null. A proper comparison, therefore, requires some preliminary assumptions on the accuracy of our estimates for the sequences of  $C_l$ . On the other hand, it is known that harmonic space methods are less robust in the presence of partially observed sky maps, that is, when there are gaps, for instance, due to the presence of the Milky Way and other foreground sources. Again, a number of assumptions are needed here on the possibility of filtering out these spurious effects in realistic experiments. Another practical issue relates to the effect of instrumental noise, the statistical properties of which are by themselves a very important area for theoretical and applied research. More importantly, the power of the different procedures clearly depends heavily on the nature of the non-Gaussian alternative; there are models where non-Gaussianities produce stronger effects in pixel space and other models where the reverse is the case (see Hansen et al., 2002, 2003). In general, there is no doubt that any empirical analysis must resort to a combination of different methods to produce reliable results.

With these caveats in mind, it is still possible to have some preliminary indication of the relative performance of the various procedures. Among real-space methods, the first Minkowski functional seems to represent a valuable option, because it is less sensitive to the specification of the correlation structure and presents good power properties against a wide variety of alternatives (Cabella et al., 2003). Among harmonic space methods, there is some evidence that the bispectrum may have little power against generic alternatives (Gangui, Pogosian and Winitzki, 2002; Cabella et al., 2003). On the other hand, this statistic is at present the most promising as a match filter for specific non-Gaussian models (Komatsu and Spergel, 2001, Section 4). Indeed, the widely popular nonstandard inflationary models predict the ratio of the Gaussian to the non-Gaussian component to be on the order of  $f_{NL} \times 10^{-5}$ , with  $f_{NL} \sim O(1)$ ; for such tiny signals, there are basically no alternatives to a model fit of the non-Gaussian part, and such model fit will most probably require the bispectrum or some related statistics.

It is important to stress that harmonic analysis is by no means the unique approach for investigating non-Gaussianity in a dual space. In particular, a number of recent articles have focused on the expansion of the

CMB field into spherical wavelets (see, e.g., Barreiro et al., 2000; Cayón et al., 2001, 2003; Martínez-González et al., 2002; Aghanim, Kunz, Castro and Forni, 2003). The first aim of these articles was the extension to spherical random fields of the most common wavelet bases, such as the Haar basis or the Mexican Hat; optimality properties were also discussed. Gaussianity tests can then be implemented by looking at the skewness and kurtosis of the spherical wavelet random coefficients. This approach is extremely interesting and promising; the area of wavelet analysis for CMB certainly deserves a lot of attention from the statistical point of view, not only as far as Gaussianity is concerned but also in connection with other statistical issues. In this article, however, we avoided a full discussion of the wavelet formalism on spherical surfaces for brevity's sake. We refer to a very recent work by Aghanim et al. (2003) for a nice introduction to this area of research and some numerical comparisons with other procedures we have discussed.

We conclude by noting that a major issue, on the border between physics and computational statistics, is the numerical generation of non-Gaussian maps. Of course, it is not hard to make a map non-Gaussian, for instance, by simply squaring and centering its value pixel by pixel. However, this is of little, if any, practical significance. There have been interesting attempts to produce non-Gaussian maps with a given spectrum and bispectrum; see, for instance, Contaldi and Magueijo (2001). The resulting maps, however, do not correspond strictly to any viable physical alternative. Algorithms to implement non-Gaussian full-sky maps with physical meaning are currently being studied (Bartolo, Matarrese and Riotto, 2002). In the near future, such maps will provide an important benchmark for comparison of different testing procedures.

#### ACKNOWLEDGMENTS

I am grateful to A. Balbi, F. K. Hansen, P. Natoli and N. Vittorio for many useful discussions. Research supported by MURST.

#### REFERENCES

- ADLER, R. (1981). *The Geometry of Random Fields*. Wiley, New York.
- AGHANIM, N., KUNZ, M., CASTRO, P. G. and FORNI, O. (2003). Non-Gaussianity: Comparing wavelet and Fourier based methods. Preprint. Available at <http://it.arxiv.org> as astro-ph/0301220.

- BANDAY, A. J., ZAROUBI, S. and GÓRSKI, K. M. (2000). On the non-Gaussianity observed in the COBE–DMR sky maps. *Astrophysical J.* **533** 575–587. Available at <http://it.arxiv.org> as astro-ph/9908070.
- BARREIRO, R. B. et al. (2000). Testing the Gaussianity of the COBE–DMR data with spherical wavelets. *Monthly Notices of the Royal Astronomical Society* **318** 475–481. Available at <http://it.arxiv.org> as astro-ph/0004202.
- BARTOLO, N., MATARRESE, S. and RIOTTO, A. (2002). Non-Gaussianity from inflation. *Phys. Rev. D* **65** 103505. Available at <http://it.arxiv.org> as hep-ph/0112261.
- BOND, J. R. and EFSTATHIOU, G. (1987). The statistics of cosmic background radiation fluctuations. *Monthly Notices of the Royal Astronomical Society* **226** 655–687.
- CABELLA, P. et al. (2003). The relative performance of pixel and harmonic space methods to search for non-Gaussianity in CMB: A Monte Carlo study. Unpublished manuscript.
- CAYÓN, L. et al. (2001). Spherical Mexican hat wavelet: An application to detect non-Gaussianity in the COBE–DMR maps. *Monthly Notices of the Royal Astronomical Society* **326** 1243–1248.
- CAYÓN, L., MARTÍNEZ-GONZÁLEZ, E., ARGÜESO, F., BANDAY, A. J. and GÓRSKI, K. M. (2003). COBE–DMR constraints on the non-linear coupling parameter: A wavelet based method. *Monthly Notices of the Royal Astronomical Society* **339** 1189–1194.
- CHIANG, L.-Y., NASELSKY, P. and COLES, P. (2001). Phase-mapping as a powerful diagnostic of primordial non-Gaussianity. Available at <http://it.arxiv.org> as astro-ph/0208235.
- CHIANG, L.-Y., NASELSKY, P. D., VERKHODANOV, O. V. and WAY, M. J. (2003). Non-Gaussianity of the derived maps from the first-year WMAP data. *Astrophysical J. Letters* **590** L65–L68. Available at <http://it.arxiv.org> as astro-ph/0303643.
- CONTALDI, C. R. and MAGUEIJO, J. (2001). Generating non-Gaussian maps with a given power spectrum and bispectrum. *Phys. Rev. D* **63** 103512. Available at <http://it.arxiv.org> as astro-ph/0101512.
- DEHLING, H. and TAQUU, M. S. (1989). The empirical process of some long-range dependent sequences with an application to  $U$ -statistics. *Ann. Statist.* **17** 1767–1783.
- DIEGO, J. M. et al. (1999). Partition function based analysis of cosmic microwave background maps. *Monthly Notices of the Royal Astronomical Society* **306** 427–436.
- DORÉ, O., COLOMBI, S. and BOUCHET, F. R. (2003). Probing CMB non-Gaussianity using local curvature. *Monthly Notices of the Royal Astronomical Society* **344** 905–916. Available at <http://it.arxiv.org> as astro-ph/0202135.
- DOUKHAN, P., LANG, G. and SURGAILIS, D. (2002). Asymptotics of weighted empirical processes of linear fields with long-range dependence. *Ann. Inst. H. Poincaré Probab. Statist.* **38** 879–896.
- FERREIRA, P. G., MAGUEIJO, J. and GÓRSKI, K. M. (1998). Evidence for non-Gaussianity in the COBE DMR four-year sky maps. *Astrophysical J. Letters* **503** L1. Available at <http://it.arxiv.org> as astro-ph/9803256.
- GANGUI, A., POGOSIAN, L. and WINITZKI, S. (2002). Cosmic string signatures in anisotropies of the cosmic microwave background. *New Astronomy Reviews* **46** 681–691. Available at <http://it.arxiv.org> as astro-ph/0112145.
- HADWIGER, H. (1959). Normale Körper im euklidischen Raum und ihre topologischen und metrischen Eigenschaften. *Math. Z.* **71** 124–140.
- HANNAN, E. J. (1970). *Multiple Time Series*. Wiley, New York.
- HANSEN, F. K., MARINUCCI, D., NATOLI, P. and VITTORIO, N. (2002). Testing for non-Gaussianity of the cosmic microwave background in harmonic space: An empirical process approach. *Phys. Rev. D* **66** 063002. Available at <http://it.arxiv.org> as astro-ph/0206501.
- HANSEN, F. K., MARINUCCI, D. and VITTORIO, N. (2003). The extended empirical process test for non-Gaussianity in the CMB, with an application to non-Gaussian inflationary models. *Phys. Rev. D* **67** 123004. Available at <http://it.arxiv.org> as astro-ph/0302202.
- HEAVENS, A. F. and GUPTA, S. (2000). Full-sky correlations of peaks in the microwave background. *Monthly Notices of the Royal Astronomical Society* **324** 960–968.
- HEAVENS, A. F. and SHETH, R. K. (1999). The correlation of peaks in the microwave background. *Monthly Notices of the Royal Astronomical Society* **310** 1062–1070.
- HU, W. (2001). Angular trispectrum of the cosmic microwave background. *Phys. Rev. D* **64** 083005. Available at <http://it.arxiv.org> as astro-ph/0105117.
- KOGUT, A. et al. (1996). Tests for non-Gaussian statistics in the DMR four-year sky maps. *Astrophysical J. Letters* **464** L29.
- KOMATSU, E. et al. (2003). First-year Wilkinson microwave anisotropy probe (WMAP) observations: Tests of Gaussianity. *Astrophysical J. Supplement Series* **148** 119–134. Available at <http://it.arxiv.org> as astro-ph/0302223.
- KOMATSU, E. and SPERGEL, D. N. (2001). Acoustic signatures in the primary microwave background bispectrum. *Phys. Rev. D* **63** 063002. Available at <http://it.arxiv.org> as astro-ph/0005036.
- KOMATSU, E., WANDEL, B. D., SPERGEL, D. N., BANDAY, A. J. and GÓRSKI, K. M. (2002). Measurement of the cosmic microwave background bispectrum on the COBE DMR sky maps. *Astrophysical J.* **566** 19–29. Available at <http://it.arxiv.org> as astro-ph/0107605.
- LIBOFF, R. L. (1998). *Introductory Quantum Mechanics*, 3rd ed. Addison–Wesley, Reading, MA.
- LEONENKO, N. N. (1999). *Limit Theorems for Random Fields with Singular Spectrum*. Kluwer, Dordrecht.
- LUO, X. (1994). The angular bispectrum of the cosmic microwave background. *Astrophysical J. Letters* **427** L71–L74. Available at <http://it.arxiv.org> as astro-ph/9312004.
- MARINUCCI, D. and PICCIONI, M. (2004). The empirical process on Gaussian spherical harmonics. *Ann. Statist.* **32** 1261–1288.
- MARTINEZ-GONZÁLEZ, E. et al. (2000). Tests of Gaussianity of CMB maps. *Astrophysical Letters and Communications* **37** 335–340. Available at <http://it.arxiv.org> as astro-ph/0010330.
- MARTINEZ-GONZÁLEZ, E., GALLEGOS, J. E., ARGÜESO, F., CAYÓN, L. and SANZ, J. L. (2002). The performance of spherical wavelets to detect non-Gaussianity in the cosmic microwave background sky. *Monthly Notices of the Royal Astronomical Society* **336** 22–32. Available at <http://it.arxiv.org> as astro-ph/0111284.
- MECKE, K. R., BUCHERT, T. and WAGNER, H. (1994). Robust morphological measures for large-scale structure in the Universe. *Astronomy and Astrophysics* **288** 697–704.

- MILLER, C. J., NICHOL, R. C., GENOVESE, C. and WASSERMAN, L. (2002). A nonparametric analysis of the cosmic microwave background power spectrum. *Astrophysical J. Letters* **565** L67–L70. Available at <http://it.arxiv.org> as astro-ph/0112049.
- NATOLI, P., MARINUCCI, D., CABELLA, P., DE GASPERIS, G. and VITTORIO, N. (2002). Non-iterative methods to estimate the in-flight noise properties of CMB detectors. *Astronomy and Astrophysics* **383** 1100–1112.
- NOVIKOV, D. I. and JØRGENSEN, H. E. (1996). A theoretical investigation of the topology of the cosmic microwave background anisotropy on the scale  $\sim 1$  degree. *Astrophysical J.* **471** 521–541.
- NOVIKOV, D., SCHMALZING, J. and MUKHANOV, V. F. (2000). On non-Gaussianity in the cosmic microwave background. *Astronomy and Astrophysics* **364** 17–25. Available at <http://it.arxiv.org> as astro-ph/0006097.
- PEACOCK, J. A. (1999). *Cosmological Physics*. Cambridge Univ. Press.
- PEEBLES, P. J. E. (1993). *Principles of Physical Cosmology*. Princeton Univ. Press.
- PHILLIPS, N. G. and KOGUT, A. (2001). Statistical power, the bispectrum and the search for non-Gaussianity in the cosmic microwave background anisotropy. *Astrophysical J.* **548** 540–549. Available at <http://it.arxiv.org> as astro-ph/0010333.
- POLENTA, G. et al. (2002). Search for non-Gaussian signals in the BOOMERANG maps: Pixel-space analysis. *Astrophysical J. Letters* **572** L27–L31.
- SMOOT, G. F. et al. (1992). Structure in the COBE differential microwave radiometer first-year maps. *Astrophysical J. Letters* **396** L1–L5.
- STOYAN, D., KENDALL, W. S. and MECKE, J. (1987). *Stochastic Geometry and Its Applications*. Wiley, Chichester.
- TURNER, M. S. and TYSON, J. A. (1999). Cosmology at the millennium. *Rev. Modern Phys.* **71** S145–S164. Available at <http://it.arxiv.org> as astro-ph/9901113.
- VARSHALOVICH, D. A., MOSKALEV, A. N. and KHERSONSKII, V. K. (1988). *Quantum Theory of Angular Momentum*. World Scientific, Singapore.
- VITTORIO, N. and JUSZKIEWICZ, R. (1987). Hot spots in the microwave sky. *Astrophysical J. Letters* **314** L29–L32.
- WASSERMAN, L. et al. (2001). Nonparametric inference in astrophysics. Preprint. Available at <http://it.arxiv.org> as astro-ph/0112050.
- WINITZKI, S. and KOSOWSKY, A. (1998). Minkowski functional description of microwave background Gaussianity. *New Astronomy* **3** 75–100. Available at <http://it.arxiv.org> as astro-ph/9710164.
- WINITZKI, S. and WU, J. H. P. (2000). Inter-scale correlations as measures of CMB Gaussianity. Preprint. Available at <http://it.arxiv.org> as astro-ph/0007213.
- WONG, E. (1971). *Stochastic Processes in Information and Dynamical Systems*. McGraw–Hill, New York.
- WU, J. H. P. et al. (2001). Tests for Gaussianity of the MAXIMA-1 cosmic microwave background map. *Phys. Rev. Lett.* **87** 251303. Available at <http://it.arxiv.org> as astro-ph/0104248.