

Implementing a 3D histogram version of the Energy-Test in ROOT

April 11, 2016

E. O. Cohen^a, I. D. Reid^b, E. Piassetzky^a,
cohen.erez7@gmail.com, ivan.reid@brunel.ac.uk, eip@tauphy.tau.ac.il

^a*School of Physics and Astronomy, Tel Aviv University, Tel Aviv 69978, Israel*

^b*College of Engineering, Design and Physical Sciences, Brunel University London, Uxbridge UB8 3PH, UK*

Abstract

Comparing simulation and data histograms is of interest in nuclear and particle physics experiments; however, the leading three-dimensional histogram comparison tool available in ROOT, the 3D Kolmogorov-Smirnov test, exhibits shortcomings. Throughout the following, we present and discuss the implementation of an alternative comparison test for three-dimensional histograms, based on the Energy-Test by Aslan and Zech.

The software package can be found at <http://www-nuclear.tau.ac.il/~ecohen/>.

Keywords: Two-Sample Test, Goodness of Fit, Energy-Test, ROOT

1. Introduction

Goodness of Fit (GoF) comparisons are a recurrent task when analyzing nuclear physics and high-energy experiments. Particularly common are GoF comparisons between histograms of data and Monte Carlo (MC) simulation. Such comparisons typically serve to determine whether the data and an MC sample are consistent with being generated from the same parent distribution. Often multiple MC sets with different parameters are generated, and GoF comparisons are needed to determine which best describes the data (The null-hypothesis (distributions are the same) is well defined, and it is

important to obtain appropriate GoF methods to check its validity).

One-dimensional comparison methods are well known in the literature. Some are designed for histogrammed data comparison (e.g, the χ^2 test), while others are intended for discrete data application (e.g, the Kolmogorov-Smirnov (KS) test), though also applicable to histogrammed data provided that the binning effects are considered.

GoF using the KS test (and other existing cumulative tests) is problematic for comparing multidimensional data, as it relies on the ordering of the data to

25 obtain the Cumulative Distribution Function (CDF) 53 the ROOT 3D-KS test approach those of the discrete
26 and because of the large number of distinct ways of or- 54 data ordered in one dimension along each coordinate
27 dering the data in space ($2^d - 1$ in d -dimensional space). 55 separately. In extreme cases this can lead to false pos-
28 Another disadvantage of multidimensional GoF tests is 56 itives as histograms with similar projections onto the
29 the lack of metric invariance, which leads to an undesir- 57 axes are compared (see e.g. [5] for the 2D case).
30 able high sensitivity of the comparison on a scale factor
31 - or the number of bins in the histogrammed case.

32 1.1. Histograms comparisons in ROOT

33 ROOT is the most widely used data analysis tool in 61
34 high-energy physics experiments [1]. The major existing 62
35 method for comparing 3-dimensional (3D) histograms 63
36 in ROOT is the **Kolmogorov-Smirnov Test** (the KS 64
37 test). ROOT also implements a 3D version of the χ^2 65
38 test, though due to exceptionally inferior performance 66
39 in previous 2D investigations [4, 5], it was not consid- 67
40 ered in this work. The 3D extension of the KS test 68
41 is complicated by the problem of ordering the data to 69
42 build the CDF. In addressing this, ROOT computes six 70
43 CDFs for each histogram, accumulating the binned data 71
44 raster-wise, in all distinct possible patterns, so that the 72
45 comparison yields six maximum differences to which the 73
46 Kolmogorov function is applied to the averages, return- 74
47 ing the null hypothesis probability (i.e., that the two 75
48 histograms represent selections from the same distribu- 76
49 tion). However, at finer histogram binning, the order 77
50 in which the binned data are accumulated approaches 78
51 the order of the discrete data in the slowest varying 79
52 dimension [4]. Consequently, the CDFs generated by 80

58 1.2. An alternative 3D test

59 The Energy Test (ETest), first proposed by Aslan
60 and Zech [2, 3], can serve as a powerful and robust tool
61 for multidimensional data comparison. Although this
62 test was originally designed for discrete data, apply-
63 ing it to histogrammed or clustered data may expedite
64 calculations [3].

65 The ETest is a two-sample test, in which the null
66 hypothesis to be examined is that both samples origi-
67 nate from the same distribution. The ETest can also
68 be considered as a standard GoF test, if there is an
69 MC sample large compared to a data sample. In this
70 case, the null hypothesis is that the data follow the
71 parent distribution of the MC sample. The difference
72 between these two cases is important for obtaining the
73 distribution of the ETest statistic. For model-dependent
74 calculations, a large number of MC samples can be gen-
75 erated and compared with the data to accumulate a
76 distribution of the Energy-Test statistic; however, in
77 the case of two-samples originating from real experi-
78 ments, this might not be possible. The only solution
79 in this case is to perform the test multiple times using
80 bootstrap samples of the data.

Reid et al. [4] have implemented a version of the ETest for 2-dimensional histogrammed data within the ROOT framework, provided some evaluations of its performance [4], and presented some of its advantages over χ^2 -2D and KS-2D ROOT implementations. A revisit of the 2D histogrammed implementation of the ETest was introduced in 2012 to a wider audience, together with comparisons to available 2D tests (χ^2 and KS) [5].

In this work we follow [5] and introduce a 3D histogrammed implementation of the ETest, as well as demonstrate some of its performances.

2. The Energy-Test

Consider a sample of **Data (D)** and **MC** points in a d -dimensional space, consisting of n_D and n_{MC} charges, $\{\mathbf{x}_i^D\}$ and $\{\mathbf{x}_j^{MC}\}$, respectively. The hypothesis that they arise from the same parent distribution is to be examined.

If **D (MC)** represents a system of positive (negative) point charges $1/n_D$ ($-1/n_{MC}$), then, in the limit of $n_D \rightarrow \infty$ and $n_{MC} \rightarrow \infty$, the total electrostatic energy (for a $1/r$ potential) of the two samples will reach a minimum when both samples have the same distribution. The ETest generalizes this concept.

2.1. The test statistic

The ETest statistic consists of three terms, corresponding to the self-energies of the two samples, **D** and

MC, and the interaction energy between the two samples, $\Phi = \Phi_D + \Phi_{MC} + \Phi_{DMC}$, where

$$\left\{ \begin{array}{l} \Phi_D = \frac{1}{n_D^2} \sum_{i=2}^{n_D} \sum_{j=1}^{i-1} \psi(|\mathbf{x}_i^D - \mathbf{x}_j^D|) \\ \Phi_{MC} = \frac{1}{n_{MC}^2} \sum_{i=2}^{n_{MC}} \sum_{j=1}^{i-1} \psi(|\mathbf{x}_i^{MC} - \mathbf{x}_j^{MC}|) \\ \Phi_{DMC} = -\frac{1}{n_D n_{MC}} \sum_{i=1}^{n_D} \sum_{j=1}^{n_{MC}} \psi(|\mathbf{x}_i^D - \mathbf{x}_j^{MC}|) \end{array} \right.$$

and ψ is a continuous, monotonically-decreasing function of the Euclidean distance r between the charges. Following [5], we choose to use $\psi = -\ln(r + \epsilon)$, rather than the electrostatic potential $1/r$, since it renders a scale-invariant function for the test, and offers better rejection powers against alternatives to the null hypothesis¹. The value of the cutoff parameter ϵ is not critical so long as it is of the order of the mean distance between points at the densest region of the sample distributions.

2.2. Implementation of a 3D histogrammed version of the ETest in ROOT

The ETest was implemented as a compiled ROOT macro for equally-binned ($N \times N \times N$) histograms. Aslan and Zech [3] suggest that the ranges of the data can be normalized, to equalize the relative scales of the x , y , and z -coordinates. We found that for our specific application a similar normalization is not necessary. Underflow and overflow bins (with indices 0 and $N+1$,

¹if all axes are scaled identically.

124 respectively, in ROOT notation) can be included with 137 must be treated individually, i.e., when bin (i_1, i_2, i_3) is
 125 nominal widths of $1/N$ below or above the histogram 138 being compared to bin (i_1, i_2, i_3)
 126 limits, selected by a user input parameter. 139 We assume the original points are randomly dis-

The choice of the number of bins chosen can be 140 tributed within the bin limits, and take the average
 based on statistical methods proposed in the literature. 141 distance between pairs of random points in a unit
 The authors found the Freedman-Diaconis rule to work 142 cube to calculate an effective cutoff ϵ . This value is
 well in practice [6]. In this approach the bin size is 143 $\langle r \rangle = 0.66170...^3$ [7], so we use $\epsilon = \langle r \rangle / N$ as the cutoff
 chosen by 144 distance. See below the sensitivity study to the cutoff

$$\text{bin size} = 2 n(x)^{-1/3} \text{IRQ}(x),$$

127 where $n(x)$ is the number of observations in the sample 147 of k points within a given bin by the weight $k^2/2$ rather
 128 x , and $\text{IRQ}(x)$ is the interquartile distance². For the 148 than the rigorous $k(k-1)/2$, to ensure that compar-
 129 example of 135,000 points uniformly distributed in a 149 isons between identical histograms return exactly zero
 130 unit cube, this results in $N \sim 50$ bins in each direction. 150 analytically.

131 To summarize, the implementation of the three
 132 Histograms neglect intrabin positional information 152 terms in the energy sum when comparing two
 133 as all points within a given bin are assigned a single 153 $N \times N \times N$ ROOT histograms, **hD** representing the
 134 position, i.e., the bin centre. Unlike the discrete case, 154 **data** and **hMC** representing the **Monte-Carlo expecta-**
 135 the self-energy between points in the same bin must 155 **tion**, with total content n_D and n_{MC} , respectively, is
 136 be taken into account. This means that the $r = 0$ case 156 given by:

²The interquartile distance, sometimes also referred to as the midspread, is the difference between the upper and lower quartiles.
³ $\langle r \rangle = \frac{1}{105} (4 + 17\sqrt{2} - 6\sqrt{3} + 21 \sinh^{-1} 1 + 42 \ln(2 + \sqrt{3}) - 7\pi)$

$$\left\{ \begin{array}{l}
\Phi_{\mathbf{D}} = \frac{1}{n_{\mathbf{D}}^2} \sum_{d_1=0}^{N+1} \sum_{d_2=0}^{N+1} \sum_{d_3=0}^{N+1} D_{d_1, d_2, d_3} \left(\begin{array}{l}
\sum_{d'_1=0}^{d_1-1} \sum_{d'_2=0}^{N+1} \sum_{d'_3=0}^{N+1} D_{d'_1, d'_2, d'_3} \psi_{d_1, d_2, d_3}^{d'_1, d'_2, d'_3} \\
+ \sum_{d'_2=0}^{d_2-1} \sum_{d'_3=0}^{N+1} D_{d_1, d'_2, d'_3} \psi_{d_1, d_2, d_3}^{d_1, d'_2, d'_3} \\
+ \sum_{d'_3=0}^{d_3-1} D_{d_1, d_2, d'_3} \psi_{d_1, d_2, d_3}^{d_1, d_2, d'_3} \\
+ 0.5 D_{d_1, d_2, d_3} \mathcal{D}_0
\end{array} \right), \\
\Phi_{\mathbf{MC}} = \frac{1}{n_{\mathbf{MC}}^2} \sum_{m_1=0}^{N+1} \sum_{m_2=0}^{N+1} \sum_{m_3=0}^{N+1} MC_{m_1, m_2, m_3} \left(\begin{array}{l}
\sum_{m'_1=0}^{m_1-1} \sum_{m'_2=0}^{N+1} \sum_{m'_3=0}^{N+1} MC_{m'_1, m'_2, m'_3} \psi_{m_1, m_2, m_3}^{m'_1, m'_2, m'_3} \\
+ \sum_{m'_2=0}^{m_2-1} \sum_{m'_3=0}^{N+1} MC_{m_1, m'_2, m'_3} \psi_{m_1, m_2, m_3}^{m_1, m'_2, m'_3} \\
+ \sum_{m'_3=0}^{m_3-1} MC_{m_1, m_2, m'_3} \psi_{m_1, m_2, m_3}^{m_1, m_2, m'_3} \\
+ 0.5 MC_{m_1, m_2, m_3} \mathcal{D}_0
\end{array} \right), \\
\Phi_{\mathbf{DMC}} = -\frac{1}{n_{\mathbf{D}} n_{\mathbf{MC}}} \sum_{d_1=0}^{N+1} \sum_{d_2=0}^{N+1} \sum_{d_3=0}^{N+1} D_{d_1, d_2, d_3} \sum_{m_1=0}^{N+1} \sum_{m_2=0}^{N+1} \sum_{m_3=0}^{N+1} MC_{m_1, m_2, m_3} \psi_{d_1, d_2, d_3}^{m_1, m_2, m_3},
\end{array} \right.$$

where

$$\psi_{i_1, i_2, i_3}^{j_1, j_2, j_3} = \begin{cases} \mathcal{D}_0 = -\ln(\langle r \rangle / N), & \text{if } i_1 = j_1, i_2 = j_2, i_3 = j_3, \\ -\frac{1}{2} \ln \left[\frac{(i_1 - j_1)^2 + (i_2 - j_2)^2 + (i_3 - j_3)^2}{N^2} \right], & \text{otherwise,} \end{cases}$$

and D_{d_1, d_2, d_3} , MC_{m_1, m_2, m_3} are the contents of individual bins within the histograms.

2.3. Computation speed

The computation time complexity of the test statistic is $\mathcal{O}(n^2)$, and in terms of histogram dimensions $\mathcal{O}(N^6)$. In order to minimize computation time, time-consuming operations were eliminated by the following:

1. Allocating local arrays holding the histogram data to enable pointer indexing rather than using the time-consuming `GetCellContents()` method when retrieving bin counts.
2. Constructing a local array to hold the potential

function $\psi_{i_1, i_2, i_3}^{j_1, j_2, j_3}$.

3. Skipping computations involving empty bins.

Table 1 shows the time expenditure for comparisons between histogram pairs filled with 10^6 randomly uniformly distributed points with various binning. The comparison of data samples with distribution of equally spaced points is meant for testing, and not to describe a real application. Despite attempts to reduce calculation time, the time expenditure for fine binning ($N \geq 50$) is very large, and time-reduction programming should

178 be further studied to address this issue. We also note 197
 179 that ROOT experiences frequent memory crashes for 198
 180 3-dimensional arrays with large sizes ($N > 60$), due 199
 181 to the fixed (and finite) memory size allocated on the 200
 182 stack. To address this, allocated variables were put 201
 183 in the heap so as to manually emulate 3D arrays. All 202
 184 calculations reported in this work were performed on a 203
 185 3 GHz Intel Core i7 processor (8 GB 1600 MHz DDR3 204
 186 memory) using ROOT version 5.34/21.

Table 1: Comparison time for 10^6 points histograms of various binning with the ROOT 3D-KS test and the ETest.

Histograms Size	ROOT 3D-KS	ETest
$10 \times 10 \times 10$	< 10 ms	< 10 ms
$30 \times 30 \times 30$	< 10 ms	5.3 s
$50 \times 50 \times 50$	30 ms	150 s
$100 \times 100 \times 100$	320 ms	2×10^4 s

187 *2.4. Testing resolving power*

188 The ability of a test to discriminate against non- 216
 189 conforming data, usually referred to as the *power* of the 217
 190 test, serves as a measure for the test capability to reject 218
 191 incompatible data sets based on selected criterion. De- 219
 192 termining the power is possible only if a confidence level 220
 193 for accepting the test result is established. A traditional 221
 194 criterion is a confidence level of 95% $CL_{95\%}$. 222

195 In order to test our implementation of the 3D ETest, 223
 196 two reference sets were generated: (a) A unit cube filled 224

with a constant distribution (no statistical fluctuations)
 of 37 points in each one of a $30 \times 30 \times 30$ bins, and (b)
 a continually re-generated sample of 1,000,000 points
 randomly and uniformly distributed in the unit cube.
 10,000 tests were performed against these references
 using samples of 1,000,000 random points. The first
 sample served as a reference for a one-sample GoF test
 that can determine the consistency with the assumption
 of a constant distribution, and the second for a two-
 sample comparison test to determine if both resulted
 from the same parent distribution.

Fig. 1 shows the resulting test statistic distributions.

The values for $CL_{95\%}$ are 2.2×10^{-6} for a constant
 parent and 4.1×10^{-6} for comparison between uniform
 random distributions.

212 *2.5. Gaussian contamination*

The test for sensitivity to contamination was con-
 ducted by the following [5]. The comparisons de-
 scribed above in Section 2.4 were repeated 1,800 times
 with 1,000,000 points, but where $n = 0 - 20\%$ of the
 points from each sample were replaced by a trivariate
 $\mathcal{N}(\mu = 0.5, \sigma = 0.1)$ Gaussian distribution. The ETest
 discrimination power was determined as the fraction of
 comparison below the corresponding $CL_{95\%}$. Results
 are presented in Fig. 2 and Table 2. As expected, for 0%
 contamination the result is consistent with the choice of
 $CL_{95\%}$, which clearly rejects distributions with $n > 1\%$
 contamination. The ETest exhibits superior perfor-

225 mance than the 3D KS test.

Table 2: Discrimination power of the ETest and the ROOT 3D-KS ($30 \times 30 \times 30$ binning), as a function of the contamination. See text for details.

Gaussian Contamination	ETest power	ROOT 3D-KS power
0%	0.044	0.0
0.01%	0.051	0.0
0.1%	0.129	0.0
0.7%	1.0	0.0
1%	1.0	0.0
1.3%	1.0	0.0
3%	1.0	0.010
5%	1.0	0.260
10%	1.0	0.942
15%	1.0	0.999
20%	1.0	1.0

226 *2.6. Binning effects*

227 To study the effects of the different number of bins on
 228 the ETest resolving power, a set of 1,000,000 points uni-
 229 formly distributed inside the unit cube was compared to
 230 3,000 similar sets, each contaminated by a fixed fraction
 231 of $n = 0.1\%$ Gaussian distributed $\mathcal{N}(\mu = 0.5, \sigma = 0.1)$
 232 points. The discrimination power for different binning
 233 (for $N = 10^3, 20^3, 30^3, 40^3$ and 50^3) is reported in Table
 234 3. As expected, the discrimination power is improved

235 with finer binning, though not drastically.

Table 3: ETest 95% confidence level for comparison between two sets of 1,000,000 uniform random distributed points and contamination of $n = 0.1\%$ as a function of the number of bins.

Histogram binning	ETest $CL_{95\%}$	ETest power
$10 \times 10 \times 10$	3.35×10^{-6}	0.11
$20 \times 20 \times 20$	4.05×10^{-6}	0.11
$30 \times 30 \times 30$	4.10×10^{-6}	0.13
$40 \times 40 \times 40$	4.65×10^{-6}	0.12
$50 \times 50 \times 50$	4.85×10^{-6}	0.14

236 *2.7. Cutoff parameter impact*

237 To study the effects of different cutoff parameters
 238 values on the ETest results, the comparisons described
 239 in section 2.4 were repeated 3,000 times using cutoff
 240 parameters $\langle r \rangle$ in the range 0.1 – 1.0. The Gaussian
 241 contamination was fixed at $n = 0.1\%$ Gaussian distri-
 242 bution $\mathcal{N}(\mu = 0.5, \sigma = 0.1)$.

243 Figure 3 shows results from this study. As expected,
 244 the choice of the cutoff parameter is not critical if its or-
 245 der of magnitude equals the mean intra-points distance
 246 in the densest distributions region.

247 *2.8. Displacement sensitivity*

248 The sensitivity of the tests to a shift in the position
 249 of a histogrammed sample was investigated by compar-
 250 ing 1,000 pairs of 135,000 trivariate $\mathcal{N}(\mu = 0.5, \sigma = 0.1)$

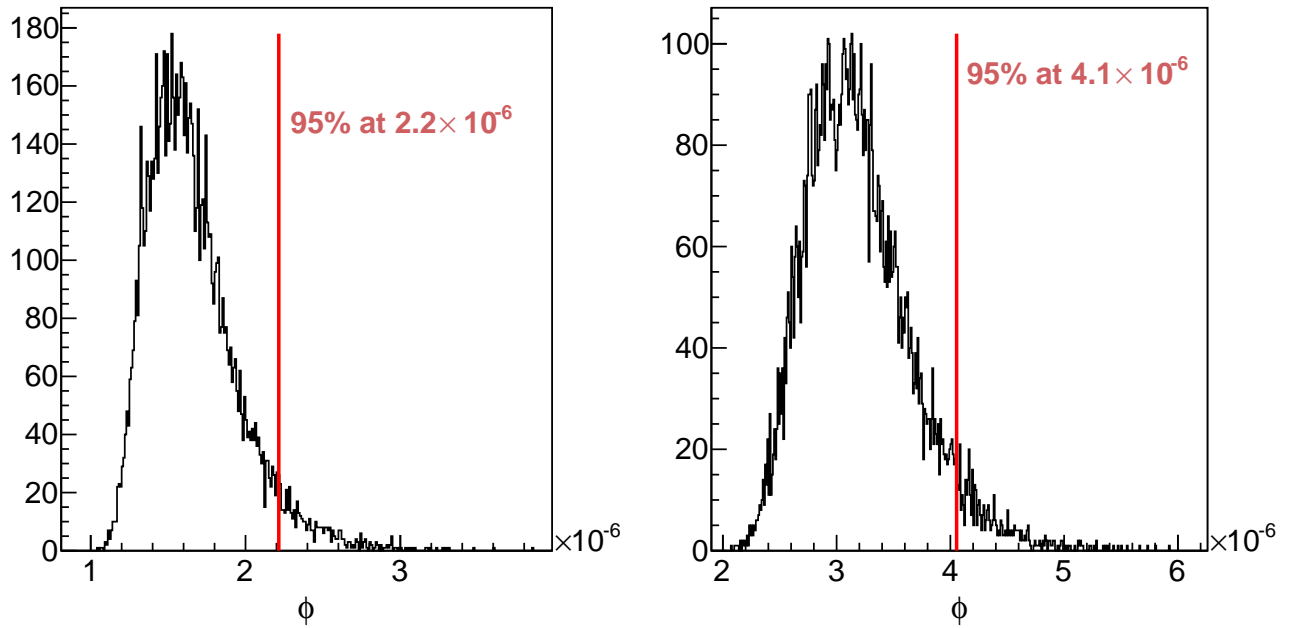


Figure 1: Distribution of the 3D ETest statistic. Compared are 10,000 sets of 1,000,000 randomly distributed points in the unit cube to a constant distribution and to a second uniform distribution, with $30 \times 30 \times 30$ bins.

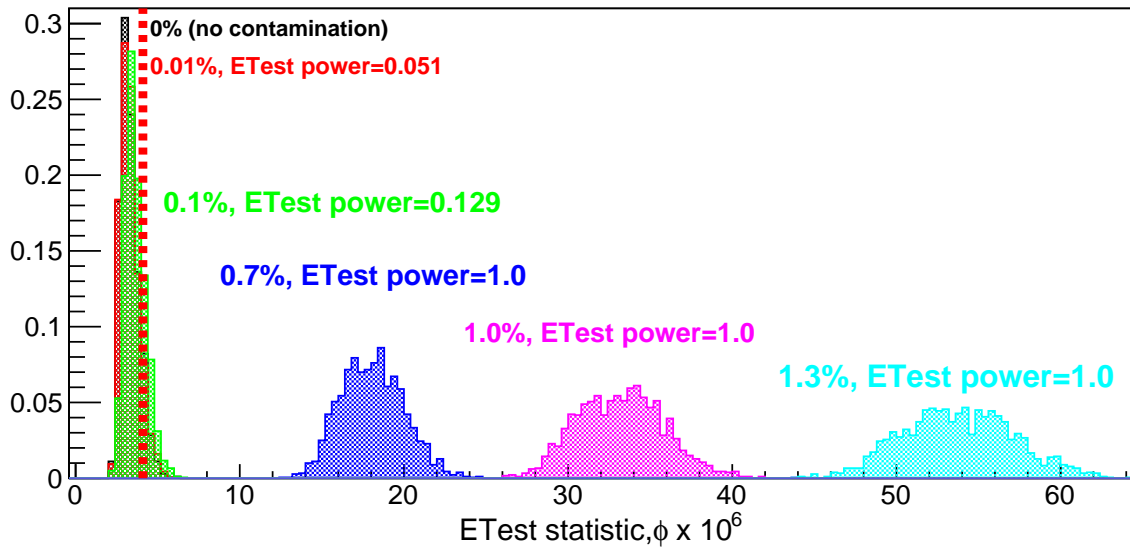


Figure 2: Same as Fig. 1 right, with one sample contaminated by $n = 0, 0.01, 0.1, 0.7, 1,$ and 1.3% trivariate $\mathcal{N}(\mu = 0.5, \sigma = 0.1)$ Gaussian distribution. The red dotted line indicates the $CL_{95\%}$.

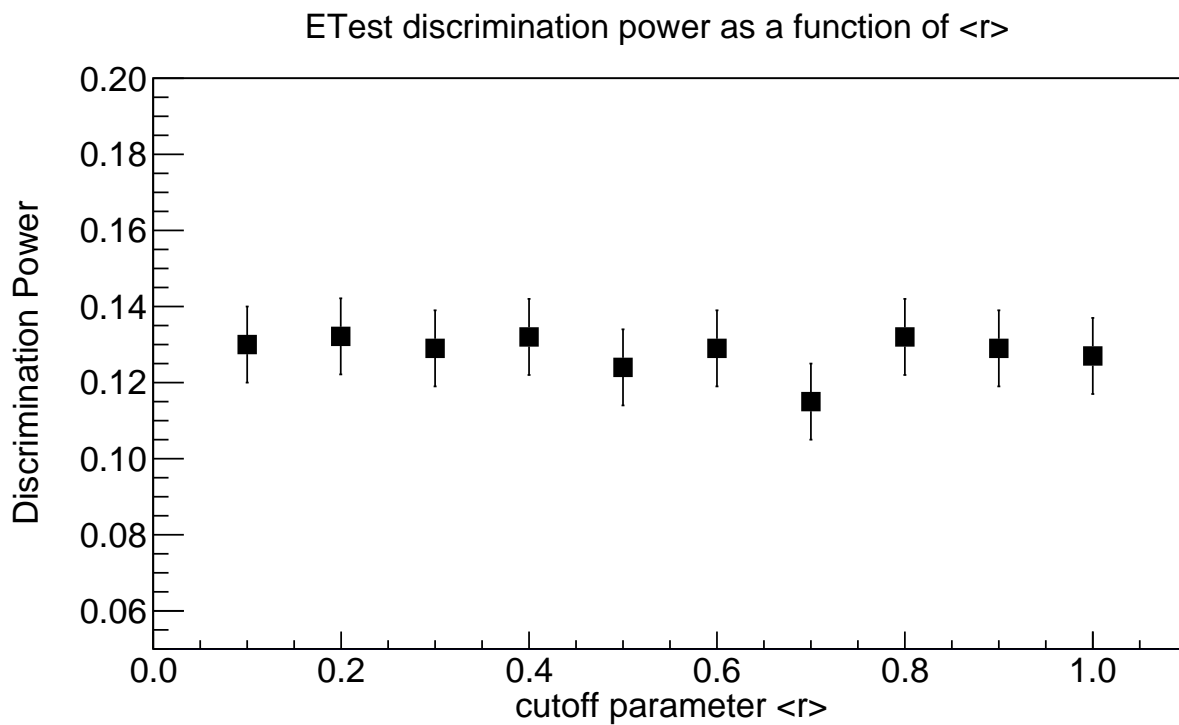


Figure 3: ETest discrimination power for different cutoff parameters $\langle r \rangle$ in the range 0.1 – 1.0, with $30 \times 30 \times 30$ bins. Compared are sets of 1,000,000 uniformly distributed points inside the unit cube and contamination of 0.1% against a uniform reference.

distributed points and $30 \times 30 \times 30$ bins. The second distribution was shifted away $(0.5, 0.5, 0.5)$ by several values (δx). For the histogrammed Energy-Test, CL_{95} was taken from the test metric distribution obtained from 10,000 pair-wise comparisons at $\delta x = 0$, which yielded a value of 1.95×10^{-5} (Figure 4); The selection criterion for the ROOT 3D-KS tests was a 5% acceptance level. The calculated powers for the tests are given in Table 4. The histogrammed ETest provides significantly better rejection than the ROOT 3D-KS test, approaching full rejection at $\delta x = 0.002$ (about 6% of bin size), compared to $\delta x = 0.2$ for the 3D-KS test.

Table 4: Discrimination power of the ETest and the ROOT 3D-KS test for various δx displacements between trivariate $N(\mu = 0.5, \sigma = 0.1)$.

δx	ETest power	ROOT 3D-KS power
0.0001	0.150	0.0
0.0005	0.337	0.0
0.0007	0.477	0.0
0.001	0.910	0.0
0.002	0.999	0.0
0.003	1.0	0.0
0.004	1.0	0.002
0.1	1.0	0.350
0.15	1.0	0.790
0.2	1.0	1.0

3. Conclusions

A new implementation of the Energy Test of Aslan and Zech, for performing GoF comparisons between three-dimensional histograms, was introduced and investigated. The software package can be found at <http://www-nuclear.tau.ac.il/~ecohen/>.

Concluding this investigation, we show that the histogrammed ETest is superior to the only available ROOT **Kolmogorov-Smirnov Test**, for comparing synthetic data sets.

The main reason for this seems to be the fact that the histogrammed ETest is a global test that compares each pair of bins in the histograms, while the ROOT 3D-KS is sensitive to neighborhood variations, dependent on the way in which the CDFs are built.

The disadvantage of the histogrammed ETest is that its calculations are time consuming, especially with fine binnings. For moderately-sized histograms the penalty is slight, particularly if the time taken to construct the histograms is also considered.

An upgraded version of the 3D ETest, which also includes an un-binned test option, is planned for implementation in ROOT in the near future.

4. Acknowledgments

This work was supported by the United States-Israel Binational Science Foundation, as well as the Science and Technology Facilities Council, UK. Additionally, Erez O. Cohen would like to acknowledge the support

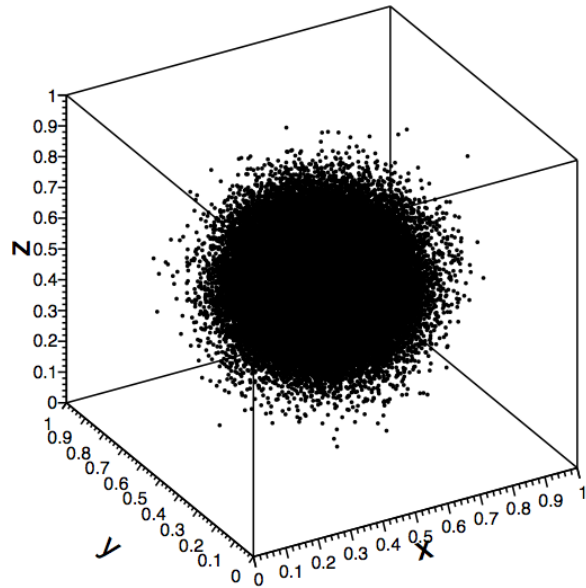
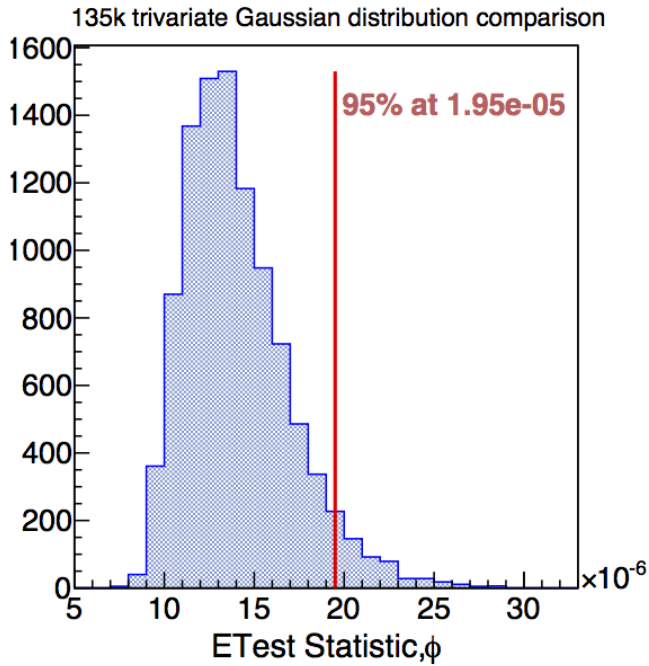


Figure 4: The distribution of results of the histogrammed ETest, comparing 10,000 pairs of histograms, each consisting of 135,000 points drawn from a trivariate $N(\mu = 0.5, \sigma = 0.1)$ distribution (shown on the right) in the unit cube.

291 of the Azrieli Foundation.

292 **References**

293 [1] <http://root.cern.ch>.

294 [2] B. Aslan, G. Zech, Nucl. Instr. Meth. A 537 (2005)
295 626.

296 [3] B. Aslan, G. Zech, J. Statist. Comput. Simul. 75
297 (2004) 109.

298 [4] I. D. Reid, R. H. C. Lopes and P. R. Hobson,
299 CERN-CMS-NOTE-2008-023.

300 [5] I. D. Reid, R. H. C. Lopes and P. R. Hobson, J.
301 Phys. Conf. Ser. **368**, 012046 (2012).

302 [6] Freedman, D., Diaconis, P., “On the histogram as
303 a density estimator: L2 theory”, <http://dx.doi.org/10.1007/BF01025868>.

304 [7] Weisstein, Eric W. “Cube Line Picking.” From
305 MathWorld— <http://mathworld.wolfram.com/CubeLinePicking.html>.

306 [8] Philip, J. “probability distributions of distribu-
307 tions between two random points in a box”, <https://people.kth.se/~johanph/habc.pdf>.