

Sparks in the Dark

Olga Sunneborn Gudnadottir¹, Axel Gallén^{1*}, Giulia Ripellino¹,
Jochen J. Heinrich², Raazesh Sainudiin³ and Rebeca Gonzalez Suarez¹

¹ Department of Physics and Astronomy, Uppsala University, Läderhyggsvägen 1,
Uppsala, 752 37, Sweden

² Department of Physics, University of Oregon, 120 Willamette Hall, 1371 E 13th
Avenue, Eugene, Oregon, 97403, United States.

³ Department of Mathematics, Uppsala University, Läderhyggsvägen 1, Uppsala, 752 37,
Sweden.

★ axel.lars.gallen@cern.ch

Abstract

This study presents a novel method for the definition of signal regions in searches for new physics at collider experiments. By leveraging multi-dimensional histograms with precise arithmetic and utilizing the SparkDensityTree library, it is possible to identify high-density regions within the available phase space, potentially improving sensitivity to very small signals. Inspired by a search for dark mesons at the ATLAS experiment, CMS open data is used for this proof-of-concept intentionally targeting an already excluded signal. Signal regions are defined based on density estimates of signal and background. These preliminary regions align well with the physical properties of the signal while effectively rejecting background events.

Copyright attribution to authors.

This work is a submission to SciPost Physics.

License information to appear upon publication.

Publication information to appear upon publication.

Received Date

Accepted Date

Published Date

1

2 Contents

3	1 Introduction	2
4	2 Datasets and event selection	3
5	3 Discriminating variables	4
6	4 Method	5
7	5 Results	6
8	6 Conclusion and Outlook	11
9	References	12

10

11

1 Introduction

Collider experiments in high-energy physics often deal with large amounts of experimental data. The two general-purpose experiments at CERN's Large Hadron Collider (LHC), ATLAS and CMS, record about 10 PB of data per year. These data are then analysed for, e.g., consistency with different theoretical models, which involves both isolating a small signal from large background and data-driven corrections to background estimates. A pre-selection of data is performed based on the particles involved in the experimental signature of the signal. Subsequently, the resulting dataset is explored with the objective to create a phase-space region enriched in signal events. This enriched region allows for a statistical analysis that is sensitive to the signal. Optimising the region involves using theoretical knowledge of the signatures and kinematic behaviour of the signal and background processes to define new variables, and a tedious process of exploring the data using 1D or 2D histograms. Machine learning classifiers are also commonly used at this stage, which both hone in on the region without the same need for manual optimisation and utilise complex relationships between variables. The downside of these methods is that the interdependence of the variables is never made explicit, and the analysis becomes harder to understand than one defined in terms of intervals in each variable. This matters not only for the understanding of the individual physicist, but also matters for reinterpretations of the results. This paper is a proof-of-concept of a new method which has the potential to produce a more sensitive signal region in a shorter time than manual optimisation, while keeping the analysis and interpretability as simple as possible.

This work builds on multi-dimensional histograms as implemented in the SparkDensityTree library [1], following [2–4]. SparkDensityTree has the following advantages compared to other density estimation methods.

Firstly, unlike most density estimation methods, including various regularization and Bayesian methods based on the likelihood, the minimum distance estimate (MDE) returned by SparkDensityTree as a multidimensional histogram is guaranteed to be within an L_1 distance or integrated absolute distance bound from the unknown density f for any given sample size n , no matter what the underlying density f happens to be, i.e., any density in L_1 , the space of Lebesgue integrable functions, and is thus said to have *universal performance guarantees* [3].

Secondly, the method scales to arbitrarily large sample sizes in high dimensions due to a scalable implementation with sparse binary trees for representing the data. A sparse binary tree can represent only the existing data in its leaves by implicitly encoding the unrepresented leaves without data as zero akin to sparse vectors and matrices.

Thirdly, SparkDensityTree provides various further statistical insights from the MDE. In particular, calculating the coverage or highest density regions of the MDE histogram of the signal and background data allows for finding the region of phase space with the largest probability density in the signal and background. The highest density region covering the sample space for a given probability $1 - \alpha$, should have the smallest possible volume such that every point inside the region should have a probability density at least as large as every point outside the region. The method takes measured or simulated data for signal or background processes as input and returns the highest density region of its density estimate (MDE histogram). The signal region is given as a union of intervals, rectangles, cuboids and hyper-cuboids over the domain of the input variables.

The current proof-of-concept is largely inspired by a search for dark mesons in ATLAS data [5]. This search explores prompt decays of dark pions and dark rho mesons into standard model particles in LHC data. The data and simulation used in the following sections, as well as the selections applied, loosely follow the analysis. The full ATLAS analysis is

61 however quite complex, and so, comparing directly with this work is not completely pos-
 62 sible. The signal point chosen in this study has already been excluded by ATLAS [5], and
 63 so the data will be used as background.

64 2 Datasets and event selection

65 The dark sector signals searched by ATLAS [5] produce two final states: $t\bar{t}b$ and $t\bar{t}b\bar{b}$.
 66 The one lepton final state is the most sensitive and it is characterized by events with one
 67 lepton and from 6 to 8 jets where at least 4 are identified as coming from the decay of a
 68 b-quark.

69 The study uses 2.3 fb^{-1} of $\sqrt{s} = 13 \text{ TeV}$ proton–proton (pp) collision data collected by
 70 the CMS experiment [6] in 2015 to model the background to the dark meson signal. The
 71 analysed data correspond to the SingleElectron [7] and SingleMuon [8] datasets released
 72 on the CERN Open Data portal [9]. Only events in the list of validated runs [10] are
 73 retained for the study. A total of about 110 million single electron and 70 million single
 74 muon pp events are available for analysis.

75 The datasets are provided in the CMS miniAOD format, which contains high-level
 76 reconstructed objects that can be used for analysis [11]. This study is based on such re-
 77 constructed electrons, muons and jets. The data is accessed and processed using the CMS
 78 analysis code provided with the CMS open data [12]. Within this framework, jets are
 79 reconstructed using the anti- k_t algorithm [13] with a fixed radius parameter $R = 0.4$ and
 80 are tagged as containing a bottom hadron (b-tagged) based on the Combined Secondary
 81 Vertex (CSV) tagging algorithm. In a complete data analysis, the fact that b-tagging
 82 efficiencies can differ significantly between data and simulation, has to be taken into ac-
 83 count. This is however not critical for this study and so, no dedicated strategy has been
 84 implemented to address this.

85 The matrix element calculation for the dark meson production is performed at next-
 86 to-leading order (NLO) in QCD based on the model described in Ref. [14], using the corre-
 87 sponding UFO model files [15]. A signal sample is simulated using MadGraph5_aMC@NLO
 88 3.5.1 [16] interfaced with Pythia 8.306 [17] for the modeling of parton showering, hadroniza-
 89 tion and underlying event with the A14 set of tuned parameters [18] and the NNPDF2.3lo [19]
 90 set of parton distribution functions (PDF). Fast simulation of the detector is done with
 91 Delphes 3.5.0 [20] using the standard CMS detector card.

92 Within Delphes, jets are determined with the FastJet 3.3.4 [21] software package and
 93 the anti- k_t algorithm [13]. The default b -tagging of the CMS Delphes card is used to
 94 identify b -jets. The dark pion mass is set to $m_{\pi_D} = 500 \text{ GeV}$ and the dark rho mass
 95 to $m_{\rho_D} = 2 \text{ TeV}$. The signal cross-section is extracted from MadGraph5_aMC@NLO
 96 2.9.9 [16] and amounts to 18.9 fb . As previously mentioned, this signal point has been
 97 excluded by the ATLAS collaboration [5]. A total of 50k signal events are simulated and
 98 the sample is normalized to the integrated luminosity of the data sample.

99 Events are further selected for the study based on kinematic and quality criteria im-
 100 posed on the reconstructed leptons and jets. In the MC events, any electron or muon
 101 with transverse momentum $p_T > 28 \text{ GeV}$ is considered as a signal lepton. In data events,
 102 the signal lepton must additionally pass the *Tight* selection criteria [22, 23]. Only events
 103 containing exactly one signal lepton are retained for the study.

104 All jets are required to have a transverse momentum $p_T > 20 \text{ GeV}$ and to satisfy
 105 $|\eta| < 2.5$. In addition, any jet is required to have an angular distance $\Delta R > 0.4$ from
 106 the signal lepton in the event, in order to resolve any reconstruction ambiguities between
 107 the lepton and jets. If these requirements are not met, the jet is discarded. Events are

108 eventually required to have at least four jets, out of which at least two must be b -tagged.
 109 Events passing all requirements listed here are selected for analysis. A total of 120k
 110 and 6.47 events pass this baseline selection in data and signal respectively.

111 3 Discriminating variables

112 The method is demonstrated on four event-level quantities that are suitable as discrimi-
 113 nating variables. The first three are; $\Delta R(l, b_2)$, defined as the angle between the lepton
 114 and the second closest b -jet; $m_{bb\Delta R_{min}}$, defined as the invariant mass of the two b -jets in
 115 the event that are closest to each other; and H_T , defined as the scalar sum of the p_T of the
 116 jets in the event. The final variable is based on $R = 1.2$ jets reclustered from the $R = 0.4$
 117 jets using the anti- k_t algorithm with a fixed radius parameter of $R = 1.2$ [24]. The lepton
 118 is added to the $R = 0.4$ jet collection before the reclustering and the highest- p_T large- R jet
 119 containing the lepton is referred to as \mathbb{J}^{lep} while the highest- p_T fully hadronic large- R jet
 120 is referred to as \mathbb{J}^{had} . The sum of the masses of these two jets is used as a discriminating
 121 variable and is denoted by $m_{\mathbb{J}^{had}} + m_{\mathbb{J}^{lep}}$. Distributions of the discriminating variables in
 122 data and signal are shown in Fig. 1 for events passing the baseline selection described in
 the previous section.

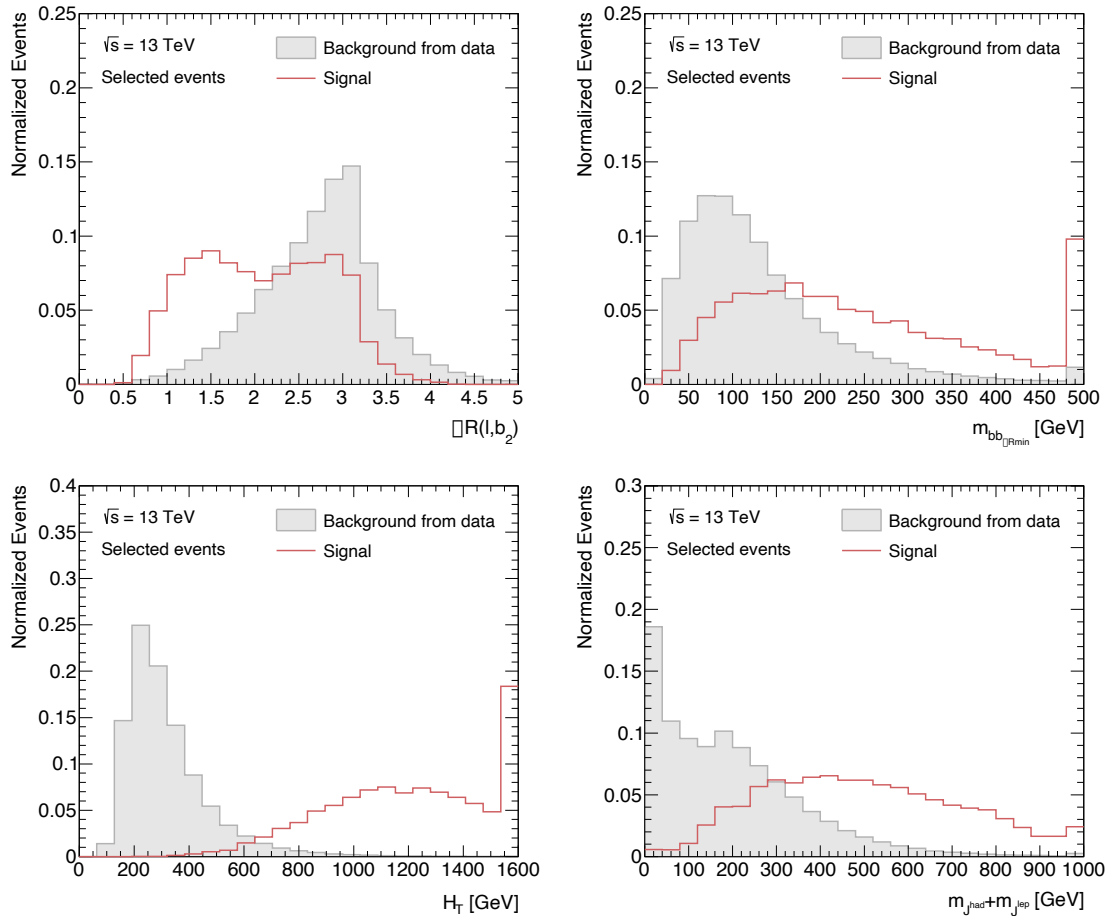


Figure 1: Distributions of the discriminating variables in data and signal for selected events, normalized to 1. (Top, Left): $\Delta R(l, b_2)$; (Top, Right): $m_{bb\Delta R_{min}}$; (Bottom, Left): H_T ; and (Bottom, Right): $m_{\mathbb{J}^{had}} + m_{\mathbb{J}^{lep}}$.

124 4 Method

125 The SparkDensityTree library is a library of statistical methods, with the base class being
 126 a multi-dimensional density estimator that for any sample generated from an unknown
 127 density returns an optimally smoothed histogram. The optimally smoothed histogram is
 128 taken to be the one that, per estimation, minimizes the L_1 distance to the true underlying
 129 distribution, using the MDE method. The statistical methods on these MDE histograms
 130 include arithmetic operations, conditional densities, coverage regions, and marginal den-
 131 sities. Recall that the marginal probability density of a subset of variables is obtained by
 132 integrating out all other variables in the joint probability density of all the variables.

133 The histogram object is represented as a binary tree in which each node represents a
 134 bisection of the phase space, and the leaves contain the event count in the finest resolution
 135 boxes thus obtained. The histogram construction begins with the definition of the root
 136 box, ideally the smallest hypercube containing all data points. From the root box, \mathbf{x}_ρ , the
 137 support is iteratively bisected until a stopping criterion is reached, as visualized in Fig. 2.
 138 The underlying tree structure [2] allows for giving each box in the splitting a unique
 139 address. The combination of the leaf address and the counts is defined as the label of the
 140 box, $(\rho v, \#\mathbf{x}_{\rho v})$.

141 The MDE histogram is described in [3], and is taken as the optimal density estimate
 142 in this work. This quantity, $f_n(x)$, is defined as following [3]:

$$f_n(x) = \frac{1}{n} \sum_{\rho v \in \mathbb{L}(s)} \frac{\#\mathbf{x}_{\rho v}}{\text{vol}(\mathbf{x}_{\rho v})} \mathbb{I}_{\mathbf{x}_{\rho v}}(x), \quad (1)$$

143 where the volume of a d -dimensional box is defined in detail in [2], \mathbb{L} corresponds to the
 144 full set of leafs, and n is the amount of data points in the root box s . This quantity is found
 145 by an adaptive search in sequentially coarser histograms, starting at the one obtained by
 146 the splitting.

147 The splitting is an inherently sequential process, but a distributed solution was devel-
 148 oped in [4,25]. This requires an initial splitting of the root box down to the finest resolution
 149 that might be needed instantaneously – possibly to the point that each leaf only has a
 150 count of one – and then merged again. This is accomplished by only representing the
 151 leaves with at least one data point using sparse binary trees.

152 In the distributed method, therefore, an additional step is added between the splitting
 153 and the MDE, which consists of merging the cells to a stopping criterion on the counts in
 154 each box, effectively representing the initial histogram for finding the MDE.

155 For a more in-depth explanation of the steps, the reader is referred to [2–4,25,26]. The
 156 procedure is sketched below:

157 **Stage 1:** Find the root box containing all the data points.

158 **Stage 2:** Define a stopping criterion for the splitting, such as a maximum box size. The root
 159 box is split until this criterion is reached, giving the *finest resolution histogram*. In
 160 this work, the finest splitting is determined by the stopping criterion that no leaf-box
 161 has any side length longer than the parameter **finestResSideLength**.

162 **Stage 3:** Merge leaves such that the counts are maximized, while not going higher than some
 163 limit **minimumCountLimit** and keeping the leaf depth as small as possible.

164 **Stage 4:** Starting from the histogram obtained in stage 3, find the optimally smoothed his-
 165 togram using MDE as described in [25,26].

166 Additionally, two user-defined parameters concerning the distributed aspect of the
 167 method are available: **numTrainingPartitions** and **sampleSizeHint**. Respectively, they
 168 correspond to how many times the training data is partitioned, related to distribution of
 169 work among computing nodes, and an initial guess of points connected to the size of the
 170 node batches [27].

171 The value of this method for data exploration in high-energy physics lies in the next
 172 step. When the MDE histogram is obtained, the highest density regions can be extracted
 173 by calculating the pdf coverage regions; and accordingly the highest and lowest density
 174 regions.

175 For simplicity, marginal densities are considered in this work, but the method can be
 176 extended to take the full density into account simultaneously.

177 The marginal densities for all unique pairs of the variables can be obtained from the
 178 4-dimensional MDE histogram. In this paper, $\binom{4}{2} = 6$ unique pairs of variables are chosen
 179 and these six combinations are what the highest density regions are computed from. This
 180 is done separately for signal and background. The signal and background highest density
 181 regions can be defined independently of each other, and can, crucially, be flipped around
 182 to allow for finding the least dense region in the background density. From here, the user
 183 has to consider the best ways to use these marginal densities, and an example is given
 184 below.

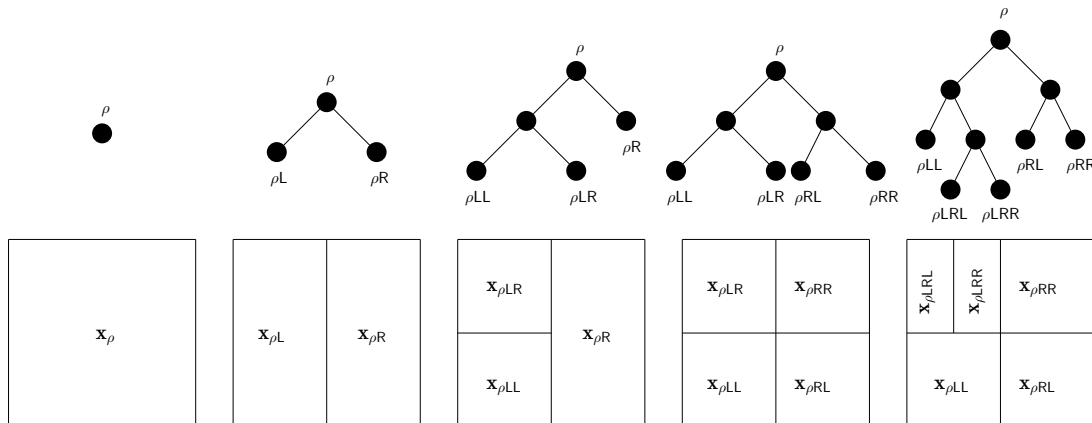


Figure 2: A sequence of splittings along the first widest coordinate, starting from the root box in two dimensions. Obtained from [2].

185 5 Results

186 The results presented in this work are documented in a Github repository [28]. All compu-
 187 tations for the upcoming results have been performed on Virtual Machines (VMs) hosted
 188 by Google as a part of a dataproc cluster. The cluster contains three VM instances, all
 189 of which run four Intel Skylake vCPUs and has 15 GB of RAM; all in order to utilize the
 190 distributed aspect of the method.

191 Figure 3 shows a comparison between a 2D frequency or count histogram of the data
 192 over a uniform grid and that over the optimally smoothed nonuniform partition corre-
 193 sponding to the MDE histogram of this method. All distributions considered in this work
 194 have been compared in this way to ensure sensible density estimates are returned by the
 195 method.

196 The density estimate is presented at three different highest density regions for back-
 197 ground in Fig. 4, and for signal in Fig. 5 for the $m_{\mathbb{J}had} + m_{\mathbb{J}lep}$ vs. $\Delta R(l, b_2)$ combination.

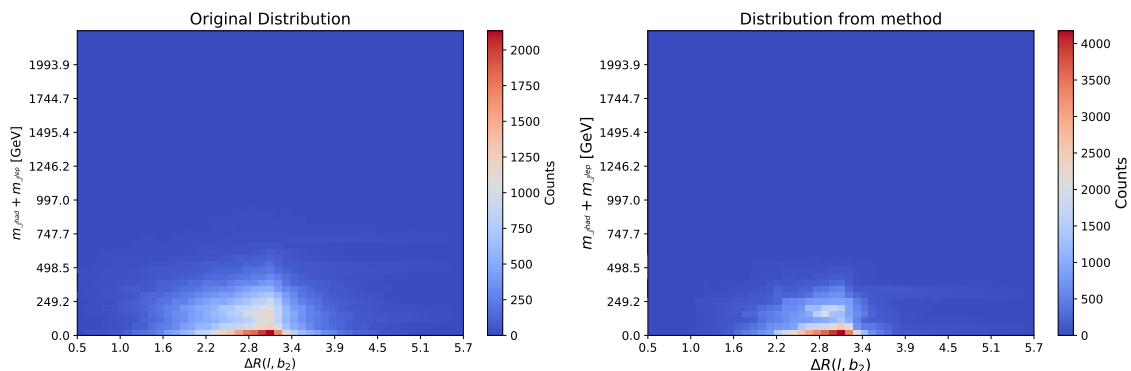


Figure 3: Comparison between a regular 2D histogram representation (Left) and the distribution obtained in this method (Right) of $m_{\mathbb{J}had} + m_{\mathbb{J}lep}$ vs. $\Delta R(l, b_2)$ in background from data.

198 Comparisons between signal and background distributions can also be made at different
 199 levels. Figure 6 shows the 3D and 2D combinations, together with the highest 50% density
 200 regions for $m_{\mathbb{J}had} + m_{\mathbb{J}lep}$ and H_T .

201 The density estimates for signal and background are combined to form $X_{\text{sig}} \otimes \bar{Y}_{\text{bkg}}$
 202 density regions, where X_{sig} indicates the $X\%$ highest signal density region and \bar{Y}_{bkg} indi-
 203 cates the complement of the $Y\%$ highest background density region. These combinations
 204 are used to design kinematic regions corresponding to the most dense signal and the least
 205 dense background. The regions are achieved from the $X\%$ highest signal density region
 206 and the $Y\%$ highest background density region using a bounding box around the density
 207 region in each pair of variables. From the bounding box, the sensitive interval of each
 208 variable is taken as the projection of the box onto that axis. The intersection of the signal
 209 interval and the complement of the background interval forms the final interval of inter-
 210 est for each variable pair. Each variable is associated with exactly three intervals from
 211 its participation in three variable pairs. In this work, the final region is defined by the
 212 union of these intervals in each variable. Three combinations are presented: $50\% \otimes \bar{50}\%$,
 213 $90\% \otimes \bar{20}\%$ and $90\% \otimes \bar{10}\%$. As an example, the obtained intervals for the $90\% \otimes \bar{10}\%$
 214 combination are:

$$\begin{aligned} \Delta R(l, b_2) &: [0.6, 1.1] & H_T &: [625, 2172] \text{ GeV} \\ m_{bb\Delta R_{min}} &: [312, 634] \text{ GeV} & m_{\mathbb{J}had} + m_{\mathbb{J}lep} &: [552, 996] \text{ GeV} \end{aligned}$$

215 When compared to the one-dimensional distributions in Fig. 1 it is clear that these corre-
 216 spond to regions with discrimination power between signal and background. While direct
 217 comparison with the ATLAS analysis [5] cannot be done, since for example $m_{\mathbb{J}had} + m_{\mathbb{J}lep}$
 218 is used as discriminant variable for the final fit, the intervals selected in this case are all
 219 fully contained in the signal region of the actual analysis.

220 The event selection corresponding to the intervals is applied to signal and data and
 221 the number of events passing the requirements are presented and compared in Table I.

222 The method results on less than one signal event on all tested scenarios and no back-
 223 ground events pass the selections in the most aggressive selection. Dark meson signals are
 224 usually very small, and unlikely to be accessible in 2.3 fb^{-1} of data. It is possible however
 225 to naively scale the 0.57 expected events in the $50\% \otimes \bar{50}\%$ scenario to, e.g. the full Run
 226 2 data set collected by ATLAS, containing 140 fb^{-1} , to more than 30 events, a reasonable
 227 signal for a new physics search.

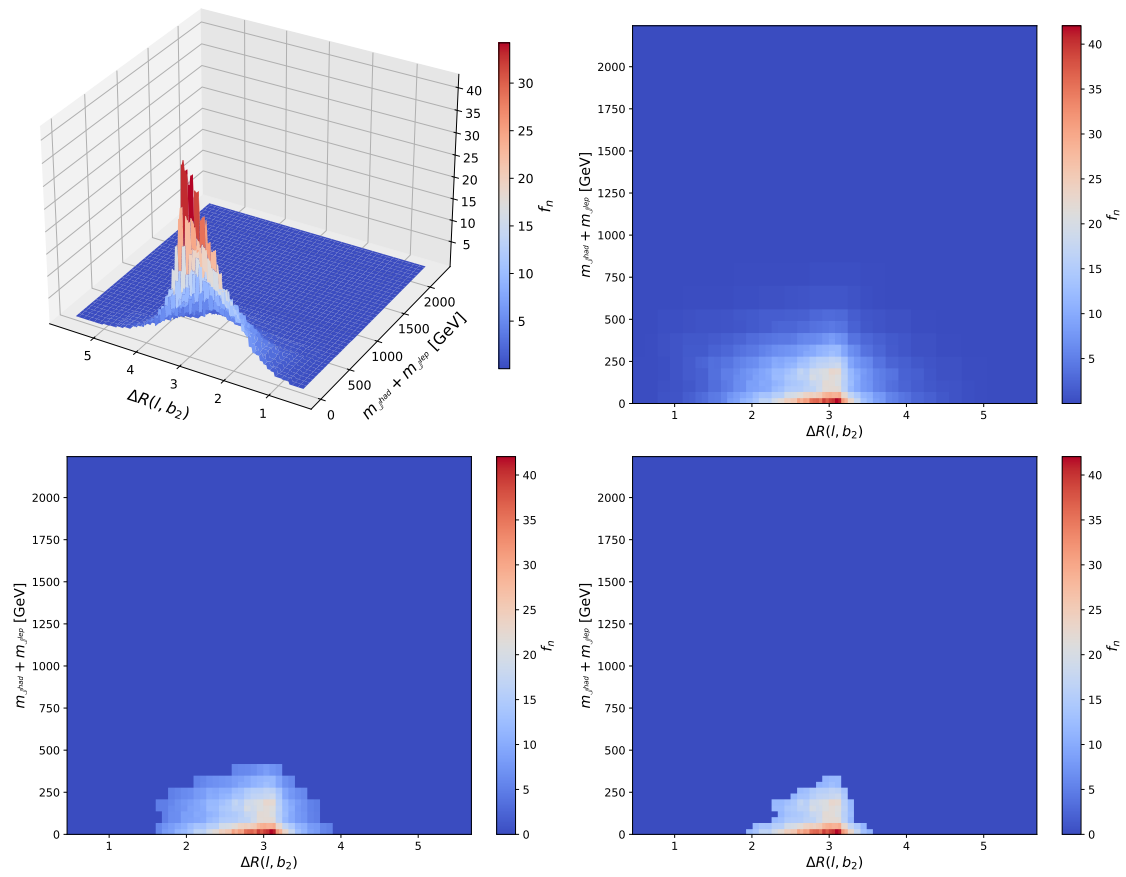


Figure 4: Full background density estimate of the $m_{\mathbb{J}had} + m_{\mathbb{J}lep}$ vs. $\Delta R(l, b_2)$ combination in 3D (Top Left) and 2D (Top Right) together with the highest 75% density region (Bottom Left) and the highest 50% density region (Bottom Right).

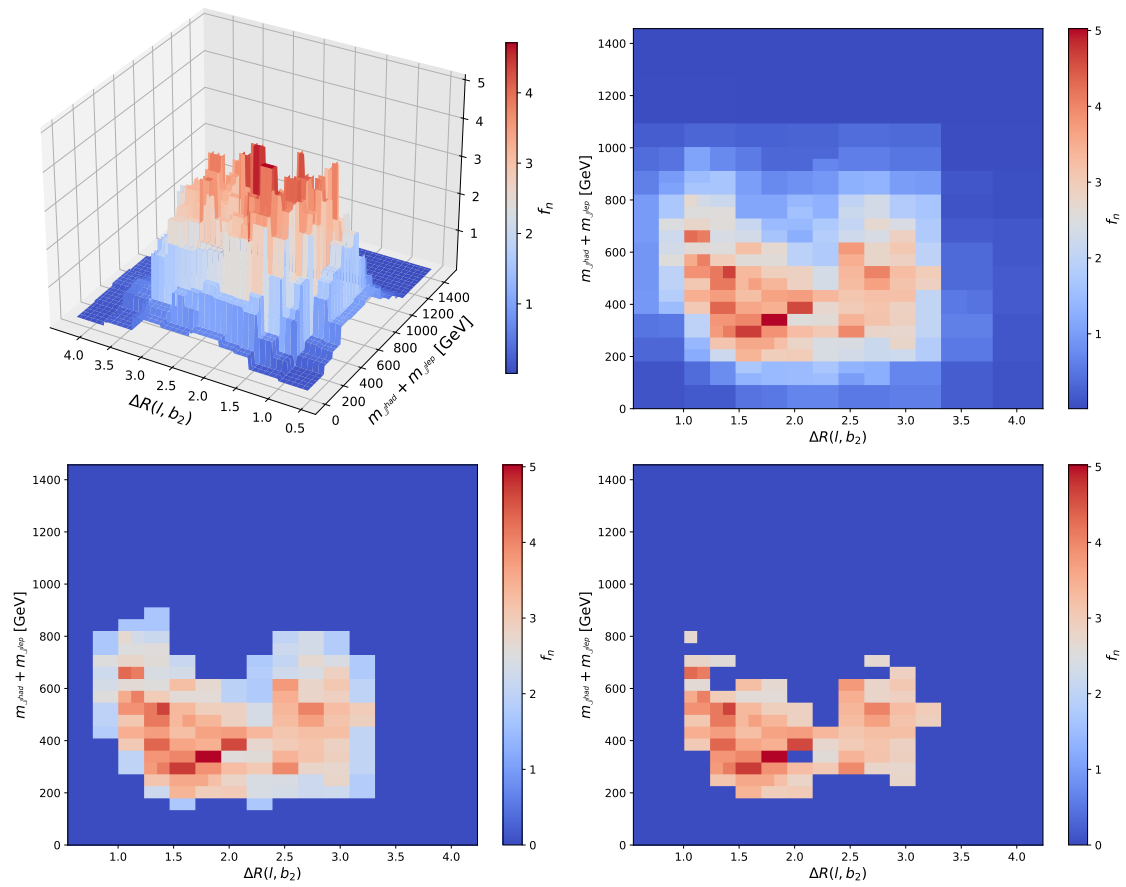


Figure 5: Full signal density estimate of the $m_{Jhad} + m_{Jlep}$ vs. $\Delta R(l, b_2)$ combination in 3D (Top Left) and 2D (Top Right) together with the highest 75% density region (Bottom Left) and the highest 50% density region (Bottom Right).

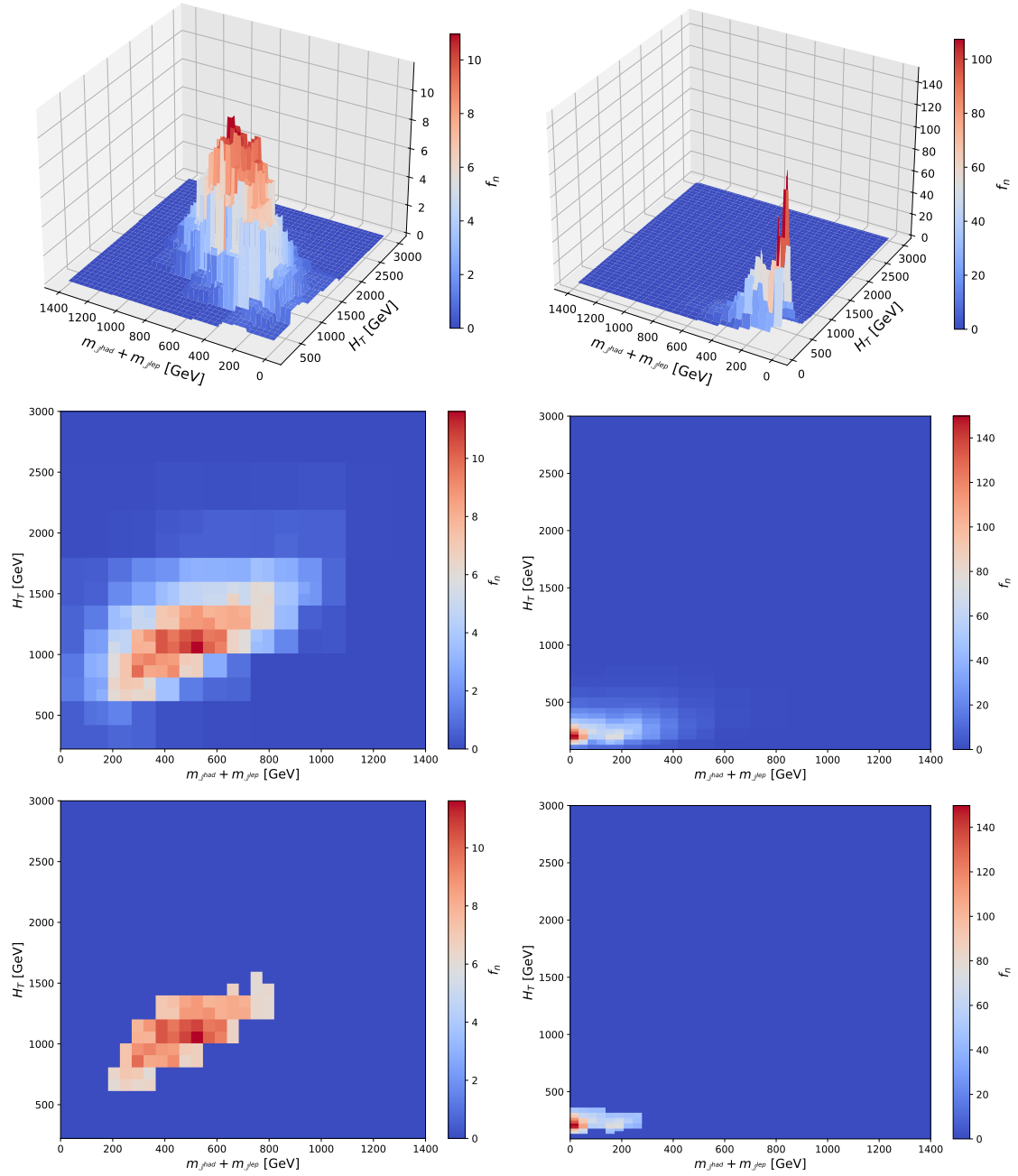


Figure 6: Full density estimate of the $m_{Jhad} + m_{Jlep}$ vs. H_T combination in 3D (Top), 2D (Middle) and the highest 50% density region (Bottom) for signal (Left) and background (Right).

Table I: Number of signal and background events passing the baseline analysis event selection and the selections derived from the density regions applied on top of the baseline. Relative numbers of events with respect to the baseline analysis selection are given within brackets. The signal numbers are normalized to the integrated luminosity of the dataset.

Selection	Signal	Background
Baseline	6.47 (100.00%)	123951 (100.00%)
50% \otimes $\overline{50}$ %	0.57 (8.74%)	364 (0.29%)
90% \otimes $\overline{20}$ %	0.30 (4.57%)	16 (0.01%)
90% \otimes $\overline{10}$ %	0.07 (1.11%)	0 (0.00%)

228 The method could further be developed to identify the highest density region directly in
 229 the 4D histogram, and then project this onto the four axes. The SparkDensityTree library
 230 allows for defining arithmetic on the histograms, and it might be possible to combine the
 231 signal and background histograms and find the densest region in, e.g., number of signal
 232 events divided by number of background events, or the difference between the histograms.

233 Finally, scalability is a very powerful aspect of this approach. This study did not delve
 234 into it, but as mentioned in [4, 25, 26], the original method has been tested on several
 235 terabytes of simulations, and great decreases in computational time can be seen with the
 236 increase of cores. In this work, for example, the time required (averaged over 5 runs) to
 237 run over the CMS open data was 42.9 seconds, while the time require to run over the
 238 signal was 26.9 seconds.

239 This is something of interest for the field of high-energy physics, as it would be straight-
 240 forward to run directly on the full collision datasets from the LHC.

241 6 Conclusion and Outlook

242 This paper introduces a scalable method, originally formulated in a purely mathemati-
 243 cal context, applied for the first time in a high-energy setting. The approach relies on
 244 optimally smoothed multi-dimensional histograms with universal performance guarantees
 245 through scalable sparse binary tree arithmetic, incorporated in the SparkDensityTree li-
 246 brary. It enables a rigorous definition of phase space regions enriched in signal, using
 247 multiple variables at a time. This method suggests promising avenues for the exploration
 248 of new physics phenomena at the LHC.

249 A large number of additional options is available from the SparkDensityTree library.
 250 This library contains several arithmetic operations and statistical methods (not covered
 251 here) that can be advantageous for studies on histograms, naturally interesting in a high-
 252 energy physics context.

253 Acknowledgements

254 The authors want to thank the conveners and members of the Heavy Quarks and Top
 255 Subgroup in ATLAS for the interest in this project and the discussions.

256 **Author contributions** The paper design and planning was done by OSG, and the method
 257 was developed by OSG and AG. Data preparation, monte carlo simulation, and analysis
 258 was done by GR and JJH. The manuscript was written by OSG, AG, GR, RS, and RGS.

259 Final editorial work done by GR and RGS. Figures are from AG and GR. All authors
260 participated in the discussion leading to this paper. The paper was approved by all
261 authors.

262 **Funding information** This research was partially supported by the project *AI4Research*
263 at Uppsala University. This material is based upon work supported by the Google Cloud
264 Research Credits program with the award GCP19980904. G. Ripellino is supported by the
265 Carl Trygger foundation (CTS 20:1169). J. Heinrich is supported by the Department of
266 Energy Office of Science Award DE-SC0017996. The Swedish Research Council supports
267 A. Gallén and R. Gonzalez Suarez (VR 2023-03403). R. Sainudiin is partially supported
268 by the Wallenberg AI, Autonomous Systems and Software Program funded by Knut and
269 Alice Wallenberg Foundation. O. Sunneborn Gudnadottir is partially supported by the
270 Centre for Interdisciplinary Mathematics (CIM) at Uppsala University.

271 References

- 272 [1] T. W. A. Sandstedt, J. Graner and R. Sainudiin, *SparkDensityTree: An Apache Spark*
273 *library for scalable density estimation, anomaly detection, and conditional density*
274 *regression with universal performance guarantees through distributed sparse binary*
275 *trees* (2023).
- 276 [2] W. T. J. Harlow, R. Sainudiin, *Mapped regular pavings*, Reliab. Comput. **16**, 252–282
277 (2012), doi:[10.1007/s42081-019-00054-y](https://doi.org/10.1007/s42081-019-00054-y).
- 278 [3] R. Sainudiin and G. Teng, *Minimum distance histograms with universal performance*
279 *guarantees*, Japanese J. Stat. Data Sci. **2**, 507–527 (2019), doi:[10.1007/s42081-019-](https://doi.org/10.1007/s42081-019-00054-y)
280 [00054-y](https://doi.org/10.1007/s42081-019-00054-y).
- 281 [4] R. Sainudiin, W. Tucker and T. Wiklund, *Scalable Multivariate Histograms* (2020),
282 [2012.14847](https://doi.org/10.1007/s42081-019-00054-y).
- 283 [5] G. Aad *et al.*, *Search for dark mesons decaying to top and bottom quarks in proton-*
284 *proton collisions at $\sqrt{s} = 13$ TeV with the ATLAS detector*, JHEP **09**, 005 (2024),
285 doi:[10.1007/JHEP09\(2024\)005](https://doi.org/10.1007/JHEP09(2024)005), [2405.20061](https://arxiv.org/abs/2405.20061).
- 286 [6] S. Chatrchyan *et al.*, *The CMS Experiment at the CERN LHC*, JINST **3**, S08004
287 (2008), doi:[10.1088/1748-0221/3/08/S08004](https://doi.org/10.1088/1748-0221/3/08/S08004).
- 288 [7] CMS Collaboration, *SingleElectron primary dataset in MINIAOD format from RunD*
289 *of 2015 (/SingleElectron/Run2015D-08Jun2016-v1/MINIAOD)*, CERN Open Data
290 Portal, doi:[10.7483/OPENDATA.CMS.29BN.FBTV](https://doi.org/10.7483/OPENDATA.CMS.29BN.FBTV) (2021).
- 291 [8] CMS Collaboration, *SingleMuon primary dataset in MINIAOD format from RunD of*
292 *2015 (/SingleMuon/Run2015D-08Jun2016-v1/MINIAOD)*, CERN Open Data Portal,
293 doi:[10.7483/OPENDATA.CMS.1LUB.Y1DH](https://doi.org/10.7483/OPENDATA.CMS.1LUB.Y1DH) (2021).
- 294 [9] CERN, *CERN Open Data*, <http://opendata.cern.ch> (2021).
- 295 [10] CMS Collaboration, *CMS list of validated runs for primary datasets of data taking*,
296 *CERN Open Data Portal*, <http://opendata.cern.ch/record/14210> (2021).
- 297 [11] CMS Collaboration, *Physics Objects, CMS Open Data Guide*, [https://](https://cms-opendata-guide.web.cern.ch/analysis/selection/objects/objects/)
298 cms-opendata-guide.web.cern.ch/analysis/selection/objects/objects/.

- 299 [12] CMS Collaboration, *CMS Open Data Guide*, <https://cms-opendata-guide.web.cern.ch>.
300
- 301 [13] M. Cacciari, G. P. Salam and G. Soyez, *The anti- k_t jet clustering algorithm*, JHEP
302 **04**, 063 (2008), doi:[10.1088/1126-6708/2008/04/063](https://doi.org/10.1088/1126-6708/2008/04/063), [0802.1189](https://arxiv.org/abs/0802.1189).
- 303 [14] G. D. Kribs, A. Martin, B. Ostdiek and T. Tong, *Dark Mesons at the LHC*, JHEP
304 **07**, 133 (2019), doi:[10.1007/JHEP07\(2019\)133](https://doi.org/10.1007/JHEP07(2019)133), [1809.10184](https://arxiv.org/abs/1809.10184).
- 305 [15] G. D. Kribs, A. Martin, B. Ostdiek and T. Tong, *HeavyDarkMesons*, <https://github.com/bostdiek/HeavyDarkMesons.git>.
306
- 307 [16] J. Alwall, R. Frederix, S. Frixione, V. Hirschi, F. Maltoni, O. Mattelaer, H. S. Shao,
308 T. Stelzer, P. Torrielli and M. Zaro, *The automated computation of tree-level and
309 next-to-leading order differential cross sections, and their matching to parton shower
310 simulations*, JHEP **07**, 079 (2014), doi:[10.1007/JHEP07\(2014\)079](https://doi.org/10.1007/JHEP07(2014)079), [1405.0301](https://arxiv.org/abs/1405.0301).
- 311 [17] T. Sjöstrand, S. Ask, J. R. Christiansen, R. Corke, N. Desai, P. Ilten, S. Mrenna,
312 S. Prestel, C. O. Rasmussen and P. Z. Skands, *An introduction to PYTHIA 8.2*,
313 Comput. Phys. Commun. **191**, 159 (2015), doi:[10.1016/j.cpc.2015.01.024](https://doi.org/10.1016/j.cpc.2015.01.024), [1410.3012](https://arxiv.org/abs/1410.3012).
- 314 [18] ATLAS Collaboration, *ATLAS Pythia 8 tunes to 7 TeV data*, Tech. rep., CERN,
315 Geneva (2014).
- 316 [19] R. D. Ball *et al.*, *Parton distributions with LHC data*, Nucl. Phys. B **867**, 244 (2013),
317 doi:[10.1016/j.nuclphysb.2012.10.003](https://doi.org/10.1016/j.nuclphysb.2012.10.003), [1207.1303](https://arxiv.org/abs/1207.1303).
- 318 [20] J. de Favereau, C. Delaere, P. Demin, A. Giammanco, V. Lemaître, A. Mertens
319 and M. Selvaggi, *DELPHES 3, A modular framework for fast simulation of a generic
320 collider experiment*, JHEP **02**, 057 (2014), doi:[10.1007/JHEP02\(2014\)057](https://doi.org/10.1007/JHEP02(2014)057), [1307.6346](https://arxiv.org/abs/1307.6346).
- 321 [21] M. Cacciari, G. P. Salam and G. Soyez, *FastJet User Manual*, Eur. Phys. J. C **72**,
322 1896 (2012), doi:[10.1140/epjc/s10052-012-1896-2](https://doi.org/10.1140/epjc/s10052-012-1896-2), [1111.6097](https://arxiv.org/abs/1111.6097).
- 323 [22] CMS Collaboration, *Electrons*, *CMS Open Data Guide*, [https://cms-opendata-guide.
324 web.cern.ch/analysis/selection/objects/muons/](https://cms-opendata-guide.web.cern.ch/analysis/selection/objects/muons/).
- 325 [23] CMS Collaboration, *Muons*, *CMS Open Data Guide*, [https://cms-opendata-guide.
326 web.cern.ch/analysis/selection/objects/muons/](https://cms-opendata-guide.web.cern.ch/analysis/selection/objects/muons/).
- 327 [24] B. Nachman, P. Nef, A. Schwartzman, M. Swiatlowski and C. Wanotayaroj, *Jets
328 from Jets: Re-clustering as a tool for large radius jet reconstruction and grooming at
329 the LHC*, JHEP **02**, 075 (2015), doi:[10.1007/JHEP02\(2015\)075](https://doi.org/10.1007/JHEP02(2015)075), [1407.2922](https://arxiv.org/abs/1407.2922).
- 330 [25] J. Graner, *Scalable algorithms in nonparametric computational statistics*, Master's
331 thesis, Uppsala University (2022).
- 332 [26] A. Sandstedt, *Scalable nonparametric L_1 density estimation via sparse subtree parti-
333 tioning*, Master's thesis, Uppsala University (2023).
- 334 [27] A. Sandstedt, J. Graner, T. Wiklund and R. Sainudiin, *SparkDensityTree-examples:
335 User-guide with examples for SparkDensityTree library, Version 1.0, License: Apache-
336 2.0*, <https://github.com/lamastex/SparkDensityTree-examples> (2023).
- 337 [28] O. Sunneborn Gudnadottir, A. Gallén, G. Ripellino and R. Gonzalez Suarez,
338 *SparksInTheDark*, <https://github.com/giuliaripellino/GOAR-ML-Project>.