

1 **Key words.** Clustering, Path Length, Consensus, N-Dimensional, Line of Sight

2 **HIGH DIMENSIONAL CLUSTER ANALYSIS USING PATH LENGTHS***

3 KEVIN L. MCILHANY [†] AND STEPHEN WIGGINS [‡]

4 **Abstract.** A hierarchical scheme for clustering data is presented which applies to spaces with a high number of dimensions
5 ($N_D > 3$). The data set is first reduced to a smaller set of partitions (multi-dimensional bins). Multiple clustering techniques
6 are used, including spectral clustering, however, new techniques are also introduced based on the path length between partitions
7 that are connected to one another. A Line-Of-Sight algorithm is also developed for clustering. A test bank of 12 data sets with
8 varying properties is used to expose the strengths and weaknesses of each technique. Finally, a robust clustering technique is
9 discussed based on reaching a consensus among the multiple approaches, overcoming the weaknesses found individually.

10 **1. Introduction.** Clustering is a fundamental technique and methodology in data analysis and machine
11 learning. The explosion of the field of data science has, consequently, led to an expansion in how this notion
12 is applied. In this respect, it would be more appropriate to refer to clustering as data organization, which
13 would encompass the ideas of 1) data reduction, 2) data identification, 3) data clustering, and 4) data
14 grouping.

15 Data reduction is the process of converting raw data into a form that is more amenable for the application
16 of a specific analytical and/or computational methodology. Data identification is the process of analysing
17 trends or distributions within the data. Data clustering is the process of associating data through proximity,
18 similarity, or dissimilarity. Data grouping refers to breaking down data into groups according to a criterion
19 that is appropriate for the specific application under consideration.

20 The literature on clustering is extensive and it is beyond the scope of this paper to provide an adequate
21 review of this topic. The following papers [Jain et al. \(1999\)](#); [Ng et al. \(2002\)](#); [Barbakh et al. \(2009\)](#); [Jain \(2010\)](#) provide background on the clustering methods in this paper and the book [Kaufman and Rousseeuw \(2009\)](#) provides a broad overview of clustering methodologies, as well as their numerical implementation.

22 There is no single algorithm that realizes all four of these aspects of data organization. The approach
23 to this problem pursued in this paper is to develop a hierarchical scheme leading to a cluster analysis that
24 encompasses the issues raised above and adapts to high dimensional spaces.

25 The data analysis scheme presented in this paper uses a blend of traditional data analysis via a multi-
26 variate histogram along with standard clustering techniques, such as k-means, k-medoids and spectral clus-
27 tering. By binning the data onto a multi-dimensional grid, data is partitioned into regions on the grid which
28 may be connected or separated depending on the character of the data set. Data reduction is realized by
29 only retaining bins that have a population above a user selected threshold. The resulting multidimensional
30 bins are referred to as partitions. The passage to partitions is the data reduction step.

31 Data identification is the process of assigning known data distributions (parent) to an entangled set of
32 data. Typical examples are found in the literature of Bayesian analysis [Binder \(1978\)](#); [Kass and Raftery \(1995\)](#), however, this pursuit dates farther back to earlier attempts to understand how to distinguish data
33 from two or more distributions with overlapping tails. In more difficult scenarios, several distributions might
34 overlap within the peak regions, changing the problem to the identification of subdomains of the mixed
35 versus non-mixed distributions.

36 Data clustering traditionally refers to assigning data to subsets based on the proximity of data to one
37 another. The goals of the field of data clustering have expanded from this definition, taking on some of the
38 other roles identified here. For the purposes of this study, the term clustering will refer to both the overall
39 techniques applied as well as the specific property a set has when its members are close to one another when
40 appropriate. In the broadest sense, a cluster is simply a label given to data to identify common features.

*Submitted to the editors JUN-2018.

Funding: SW acknowledges the support of ONR Grant No. N00014-01-1-0769 and EPSRC Grant No. EP/P021123/1. KM acknowledges the support by ONR grants: N00014-15-WX-01814, N00014-16-WX-01705 and N00014-17-WX-01705 as well as the Kinnear Fellowship from the USNA Foundation.

[†]Physics Department, U. S. Naval Academy, 572c Holloway Rd, Annapolis, MD 21402-5002, USA (mcilhany@usna.edu).

[‡]School of Mathematics, University of Bristol, University Walk, Bristol, BS8 1TW, United Kingdom, (s.wiggins@bristol.ac.uk).

44 Data grouping is the process of assigning labels to data, without regard for proximity or parent distributions.
 45 An example might be to segregate a class of thirty 2nd grade children into five subgroups before
 46 entering a museum for a tour. How the larger group is broken apart is unimportant, merely that the larger
 47 group is distributed into smaller groups.

48 In this study, standard clustering techniques are applied such as k-means, k-medoids and spectral clustering,
 49 along with new path-based approaches. After data reduction, data within partitions may be connected
 50 in regions where a path length can be calculated along the grid of partitions between any two data. Several
 51 new clustering algorithms have been developed using the path length. Further, if two partitions are visible
 52 to each other by a Line-Of-Sight criteria, the relationship between them is given additional significance.
 53 These ideas are used, in conjunction with standard clustering techniques, to construct 26 different clustering
 54 algorithms.

55 This paper presents five new variations of approaches to data clustering:

- 56 1. Data reduction is achieved by segmenting the data set into partitions.
- 57 2. Data clustering is sought using path lengths as a distance metric.
- 58 3. Data clustering is achieved using a Line-Of Sight criteria.
- 59 4. Spectral clustering is sought using alternatives to the graph Laplacian and the eigenspace formed.
- 60 5. Final cluster assignment is accomplished using a consensus among multiple clustering techniques.

61 An analysis configuration is the set of choices made that determines how a study is performed. The three
 62 most important choices are which clustering techniques out of the 26 available to use; what variables are
 63 used to describe the data, where each variable is a dimension in the data space; the number of bins chosen
 64 along each dimension. Changes to the resolution of how the data space is partitioned may lead to changes
 65 in a datum's cluster assignment. For each choice of clustering technique, variables used (dimensions) and
 66 resolution (binning), each datum is assigned to a cluster. When data consistently cluster in one arrangement
 67 across multiple analysis configurations, the data is assigned robustly to its cluster. To determine a *robust*
 68 clustering assignment, a polling technique is used to arrive at a consensus amongst the clustering algorithms.
 69 While any one technique has faults, the consensus of techniques overcomes any one failure mode, giving the
 70 best all-round identification [Strehl and Ghosh \(2003\)](#).

71 This paper is organized as follows, sections 2 and 3 define the basic component used in this study, the
 72 partition. Section 4 shows the calculations of several values used throughout the analysis. Section 5 discusses
 73 a Line-Of-Sight criteria. Section 6 outlines the strategy taken for this study and it lists the comprehensive
 74 set of arrays calculated that are needed for the suite of algorithms. This section also introduces a test-bank
 75 of data sets used for clustering. Section 7 presents each algorithm, with details left for the appendix. Section
 76 8 shows the results for each clustering algorithm, discussing the strengths and weaknesses of each approach.
 77 Section 9 introduces the approach to robust clustering, employing multiple techniques and how a consensus is
 78 reached. Section 10 concludes with suggestions for extending this suite of clustering techniques. Throughout
 79 this paper, matrices and vectors are shown in bold face, while components are given subscripts.

2. Reduction of Data to Partitions. In this study, *data* refer to collection of real values forming
 a vector, $\mathbf{x} = \{x \in \mathbb{R}^{N_D}\}$, residing in a data space of dimension, N_D , whose elements total N . Along
 each dimension of the data space, the data is coarsely delineated into a set of bins, $\{b_i \in 1 \dots N_{B,i}\}$ where
 $i = 1 \dots N_D$ and $N_{B,i}$ is the number of bins per dimension. For each datum, the collection of indices form
 a bin address vector, $\mathbf{b} = \{b \in \mathbb{N}^{N_D}\}$ giving the unique location of a bin within the data space. Each bin
 is given a unique index, \tilde{k} , serialized by the expression given below. Within each bin, multiple data may
 reside, where $w_{\tilde{k}}$ is the number of elements in each bin (population).

$$(2.1) \quad \tilde{k} = \sum_{i=1}^{N_D} \left[(b_i - 1) \prod_{q=0}^{i-1} N_{B,q} \right] + 1, \quad \text{where } N_{B,0} = 1.$$

80 The maximal value the single index, \tilde{k} , can take is the total number of possible bins in the data space, given
 81 by the product of the number of bins, $\prod_{i=1}^{N_D} N_{B,i}$. Even though a data set may be large ($\approx 10^9$), the number
 82 of possible bins can be much larger. Consider the case with a billion data points and a data space of 12
 83 dimensions, each using 10 bins (very coarse), yielding 10^{12} possible bins. Depending on how the data is
 84 distributed, most likely the data will reside in small groupings within the data space, leaving much of the

85 domain sparse.

3. Reduction of Partitions to Clusters. The data has been reduced to a set of bins, $\mathcal{D}_k = \{\tilde{k}, w \in \mathbb{N}^2\}$ identified by an index and a population of only those bins containing data. The number of bins maybe be further reduced based on the population of the bins. Low density bins can be excluded from further study by either setting a threshold (Θ_{pop}) on the minimal number of data per bin, or by setting a threshold (Θ_{perc}) based on the cumulative percentage of the low density bins with respect to the total population of all the data.

$$(3.1) \quad \tilde{\mathcal{D}} = \left\{ \tilde{k}, w \in \mathcal{D}_k \mid [w_{\tilde{k}} > \Theta_{pop}] \text{ or } [\mathcal{F}(\tilde{k}) > \Theta_{perc}N] \right\}, \quad \text{where } \mathcal{F}(\tilde{k}) = \sum_{k'=1}^{\tilde{k}} w_{k'}, \text{ with } w_{k'} = \text{sort}(w_{\tilde{k}}).$$

86 The set $\tilde{\mathcal{D}}$ contains the bins of data which will be considered for clustering. The bin index, \tilde{k} , is mapped to
87 a sequential list of indices, $k = 1 \dots N_P$, where the total number of bins under consideration, N_P , will be
88 referred to as *partitions*, with the vector of populations, $\mathbf{w} = \{w_k \in \mathbb{N}\}$, for each partition addressed by k ,
89 and the partition data space given by, $\mathcal{D}_P = \{k, w \in \mathbb{N}^2\}$. All calculations for this study are performed on
90 the partition data space, \mathcal{D}_P , which represents the integer-based grid of bin locations. The complimentary
91 data space of either empty or low population partitions is given by $\mathcal{D}^o = \{k \in \mathcal{D}_k \mid k \notin \mathcal{D}_P\}$.

92 *Clusters* are subsets of data grouped based on a common feature. Cluster algorithms use a *criteria*
93 to delineate data, which are then *gathered* by some mechanism and then assigned to clusters. Traditional
94 definitions rely on proximity of data to one another, yet clustering can also be defined as a simple grouping
95 of the data, which could be based alphabetically, by income, or some property that is difficult to map
96 numerically such as an objects shape. Proximity alone can fail to cluster data appropriately when considering
97 data distributed along tails of distributions far from a centroid, such as a horseshoe. By altering the definition
98 of “proximity” to include distance measures such as path length, clustering can still be viewed as a local
99 grouping. This paper explores multiple clustering algorithms to later sort the clustering assignments into
100 groupings reached by *consensus*.

101 **4. Intermediary Calculations.** Several calculations are common to multiple techniques which re-
102 quire only the partition bin address vector. These low level calculations define geometrical features of how
103 the partitions are related to one another. Calculations between two partitions form matrices indexed by
104 $[k, \ell]$. Specific algorithms for each calculation can be found in the supplemental material online. The dis-
105 tances calculated here fall into two broad categories; *path lengths*, where the distance measured is between
106 partitions connected to one another, and *global*, where a connection is not required. Among path lengths,
107 two further distinctions are made; *stepwise*, where the distance is the sum of values from one partition to
108 the next, and *pathwise*, where the distance is the sum of values added from the start of the path to the
109 current partition for each step taken. The block of equations shown here are described in the following

$$\begin{aligned} \Delta \mathbf{b}_i &= \mathbf{b}_{i,k} - \mathbf{b}_{i,\ell} & \Delta \mathbf{R} &= \sqrt{\sum_{i=1}^{N_D} \Delta \mathbf{b}_i^2} \\ \text{NN1} &= \mathcal{I} \circ \Delta \mathbf{R}, & \mathcal{I} &\equiv \begin{cases} 0 & |\Delta \mathbf{b}_i| > 1, \quad \text{for any } i \\ 1 & |\Delta \mathbf{b}_i| \leq 1, \quad \text{for all } i \end{cases} \\ \text{L2} &= \sum_{j=k}^{\ell} \text{NN1}_{j,j+1} \text{ (stepwise)} & \mathbf{L2_T} &= \min \left[\sum_{j=k}^{\ell} \text{NN1}_{j,j+1} \right] \\ \text{L1} &= \sum_{i=1}^{N_D} |\Delta \mathbf{b}_i| & \Sigma \mathbf{L1} &= \sum_{j=k}^{\ell} \mathbf{L1}_{kj} \text{ (pathwise)} \\ \Sigma \mathbf{L1}_{\min} &= \min \left[\sum_{j=k}^{\ell} \mathbf{L1}_{kj} \right] & \Sigma \mathbf{L1_T} &= \sum_{j=k}^{\ell} \mathbf{L1}_{T,kj} \\ \Sigma \mathbf{L1}_{\text{VAR}} &= \sum_{j=k}^{\ell} |\Sigma \mathbf{L1}_{kj} - \Sigma \mathbf{L1}_{T,kj}|^2 & \Delta \mathbf{w} &= \mathbf{w}_k - \mathbf{w}_{\ell} \\ \mathbf{w} \mathbf{w}^{\top} &= \mathbf{w} \otimes \mathbf{w} - \text{diag}(\mathbf{w}) \end{aligned}$$

110 text. 111 The Euclidean distance is calculated between all partitions in \mathcal{D}_P . First, the difference between two par-
112 titions bin address' are calculated for each component, $\Delta \mathbf{b}_i$. The distance, $\Delta \mathbf{R}$, is then calculated from
113 the sum over $\Delta \mathbf{b}_i^2$. The first nearest neighbor matrix, NN1 , defines the distance between any two bins
114 that are in contact with one another. Two partitions are in contact with one another if there exists no bin
115 address component difference greater than one in magnitude, leading to the interpretation that they share
116 a common geometric feature; a point, line, area, etc... The matrix, NN1 , is the adjacency matrix weighted
117 by Euclidean distance, $\Delta \mathbf{R}$. As each partition is a unit hypercube, the distances range from $\{1 \dots \sqrt{N_D}\}$.

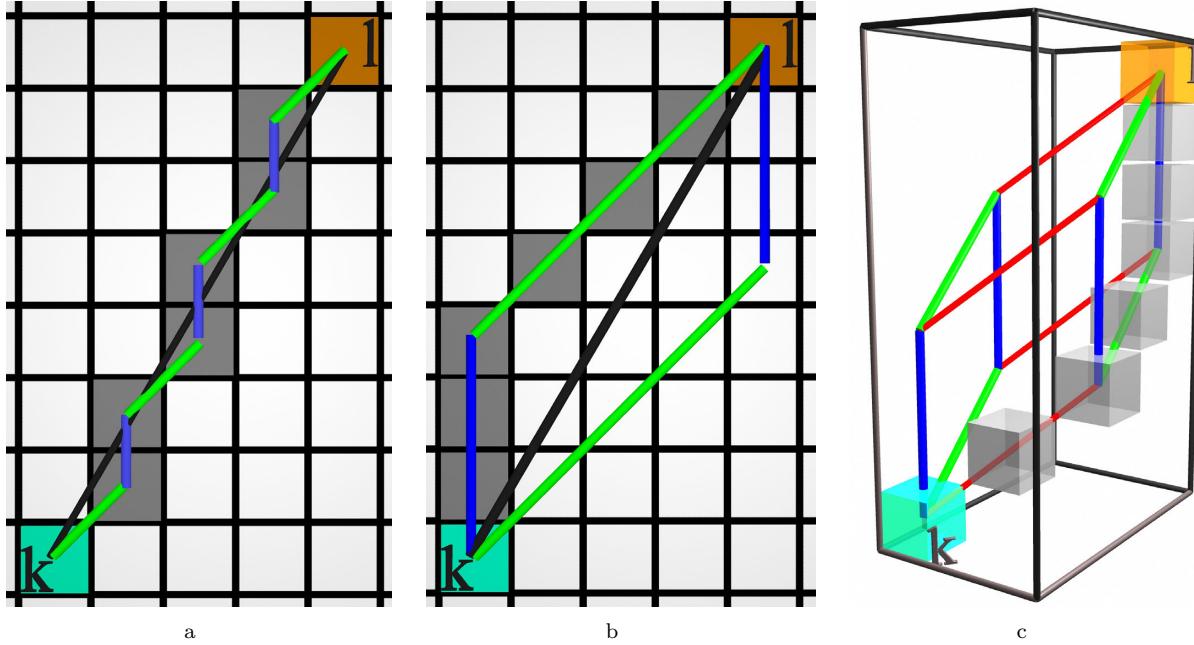


FIG. 1. Illustrating various paths from the k^{th} partition (cyan) to the ℓ^{th} (orange). Panel 1a) shows the direct (true) path and the nearest neighbor steps required to follow that path. Panel 1b) shows that the same path length can be taken along either of the edges of a parallelogram, bounded by the convex hull formed by $[k, \ell]$. The last panel 1c) shows how to extend the idea into higher dimensions, where the red segments take steps of $\sqrt{3}$, the green steps are $\sqrt{2}$ and the blue are unit steps. The grey partitions represent; 1a) the true path, 1b) the summed L1 minimal path and 1c) the summed L1 maximal path. The summed L1 path values increase monotonically from the minimal path to the maximal path, which in panel 1b) is from the upper edge of the parallelogram to the lower edge, respectively.

118 The path length, $\mathbf{L2}$, is the distance between any two partitions taken by stepping from one partition
 119 to another through $NN1$ steps, summing $\Delta\mathbf{R}$ along the path stepwise, where the initial partition is k ,
 120 interim partitions, j , up to the final partition, ℓ . Partitions are *connected* when a path is found, and for
 121 partitions having no connecting path, the path length is set to ∞ . The number of steps taken between any
 122 two partitions is the Path Count, \mathbf{PC} .

123 In order to find if two partitions meet a Line-Of-Sight (LOS) criteria, only paths that fall within the
 124 convex hull formed between the two partitions are considered. For each connecting path found, six values are
 125 calculated to determine the LOS criteria. The true path length, $\mathbf{L2_T}$, assumes a straight path exists between
 126 two partitions, giving the stepwise length formed taking the *least* number of $NN1$ steps with the smallest
 127 $\mathbf{L2}$ values possible. The Summed L1 length, $\Sigma\mathbf{L1}$, is the summation of the pathwise $\mathbf{L1}$ distances taken
 128 from the initial partition to each subsequent partition along a path. The Minimal Summed L1 path is the
 129 unique path with the least possible $\Sigma\mathbf{L1}_{\min}$, while the True Summed L1 distance is the $\Sigma\mathbf{L1}_T$ taken along
 130 the straight path established earlier. Finally, variance of the squared difference along a path, $\Sigma\mathbf{L1}_{\text{VAR}}$, is
 131 taken between the Summed L1 norm to the True Summed L1 norm along each step of a path. From these
 132 values, the true path is found which tests the LOS criteria.

133 The following calculations are performed before any paths are sought as they do not require knowledge
 134 of the exact path found, merely the endpoints which give the dimensions of the convex hull containing
 135 the two partitions $[k, \ell]$; $\Delta\mathbf{b}_i$, $NN1$, $\Delta\mathbf{R}$, $\mathbf{L2_T}$, $\mathbf{L1}$, $\Sigma\mathbf{L1}_{\min}$ and $\Sigma\mathbf{L1}_T$. The remaining three employ
 136 Dijkstra's algorithm (1959) and variations of it to calculate the shortest path lengths of: $\mathbf{L2}$ (stepwise), $\Sigma\mathbf{L1}$
 137 (pathwise), and $\Sigma\mathbf{L1}_{\text{VAR}}$ (pathwise).

138 Several calculations require that the partitions be weighted by the product of the populations of $[k, \ell]$,
 139 leading to, $\mathbf{w}\mathbf{w}^\top$, the outer product formed from the population vector taken with itself minus the weights
 140 along the diagonal to account for the self-weighting within a single partition. Further, the difference between
 141 two partitions populations is also needed, leading to the matrix, $\Delta\mathbf{w}$.

TABLE 1

Table of steps, coordinates and the distance given by the summed L1 path, $\Sigma\mathbf{L1}$, for the two examples in Figs. 1b, 1c. The top three rows represent the number of steps taken along each dimension from one partition to the next. The next three rows are the coordinates of the partitions along the paths taken, remembering that the path starts with the k^{th} partition at the origin. The final row is the summed L1 path distance. The axes are labeled where x_i is from the longest dimension to the shortest of the convex hull.

dims	Steps Taken Along Each Dimension - k^{th} partition is the origin															
	2D Case (min path)				2D Case (max path)				3D Case (min path)				3D Case (max path)			
x_1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
x_2	0	0	0	1	1	1	1	1	0	0	0	0	1	1	0	0
x_3													1	1	0	0
	Coordinates Of Path Along Each Dimension - k^{th} partition is the origin															
	2D Case (min path)				2D Case (max path)				3D Case (min path)				3D Case (max path)			
x_1	1	2	3	4	5	6	7	1	2	3	4	5	6	7	1	2
x_2	0	0	0	1	2	3	4	1	2	3	4	4	4	4	1	2
x_3															1	2
$\Sigma\mathbf{L1}$	1	3	6	11	18	27	38	2	6	12	20	29	39	50	1	3
															3	9
															17	27
															38	50
															63	

5. **Line-of-Sight (LOS) Criteria .** A Line Of Sight (LOS) criteria is introduced in this paper as a means to cluster data which gives additional significance to data while being independent of proximity. This approach assumes that data within a convex region of other data are likely to be associated together. When seeking the LOS criteria, the data space is divided into two broad subdomains, those partitions filled with sufficient data above threshold, and those partitions containing little or no data (\approx empty space). Within the filled regions of space, Dijkstra's algorithm is employed to find which partitions are connected to each other via a path and to measure the path length. Partitions of the empty set, \mathcal{D}^o , are viewed as obstacles to paths within \mathcal{D}_k . By analogy, the empty set serves to prevent LOS just as walls prevent continuity in vision.

The criteria used to establish a Line-Of-Sight (LOS) between two partitions relies on the pathwise summation of L1 distances along the path taken from $[k, \ell]$. This distance has the property that when traversing a grid from $[k, \ell]$, the distance calculated is different than when returning from $[\ell, k]$. The asymmetry of this measure proves useful in determining the LOS condition. Figure 1 illustrates how the distance is asymmetric with regard to the path taken. Three conditions must be met if two partitions are LOS:

1. A path must exist between $[k, \ell]$ that does not exit the convex hull, requiring $\mathbf{L2} = \mathbf{L2_T}$.
2. The path found must take a direct path between $[k, \ell]$, requiring that $\Sigma\mathbf{L1} \leq \Sigma\mathbf{L1_T}$.
3. The path found must follow the direct path, requiring $\Sigma\mathbf{L1_{VAR}} \leq (\mathbf{P_C}/2)^2$.

Dijkstra's algorithm finds the minimal path taken between two points on a grid given an adjacency matrix, $\mathbf{NN1}$, giving the path length, $\mathbf{L2}$. Figure 1 illustrates that multiple paths have the same value for $\mathbf{L2}$, forming a parallelotope of possible paths each with the same minimal value. The true path length can be found simply by summing $\mathbf{NN1}$ steps along the edge of the parallelotope, giving $\mathbf{L2_T}$. If $\mathbf{L2}$ exceeds $\mathbf{L2_T}$, the path found by Dijkstra has left the convex hull, leading to the first criteria. To find the pathwise $\Sigma\mathbf{L1}$ value, the adjacency matrix is altered, taking the row from $\mathbf{L1}$ for the k^{th} partition and multiplying by every row of the logical matrix, $\mathbf{NN1} > 0$; in this way, the adjacency matrix presented to Dijkstra's algorithm in a second round finds the minimal pathwise value from $[k, \ell]$. For an open region (no obstacles), the minimal summed L1, $\Sigma\mathbf{L1_{min}}$, path is along the edge of the parallelotope. Table 1 shows the calculation of $\Sigma\mathbf{L1}$ for both the minimal value as well as the maximal value from $[k, \ell]$, explaining why $\Sigma\mathbf{L1}$ is asymmetric. The example shown assumes that all paths between $[k, \ell]$ are possible, ie. there are no obstructions. From the table, the sum of L1 steps is easier to calculate when viewed along each dimension. The minimal path takes the smallest size steps first, then increasing in size until the final partition is reached. In this manner, the sum of steps along each dimension is simply the summation of $1/2 n(n + 1)$ for each dimension along the convex hull. When traversing the maximal path, the largest steps are taken first, proceeding to the smallest steps taken last. In this case, a similar summation occurs, however, an additional amount for each dimension is added because the path starts farther away from the initial point, which then adds successively along the path. The $\Sigma\mathbf{L1}$ matrix will be asymmetric as the minimal path from $[k, \ell]$ will not be the same from $[\ell, k]$. When following the minimal path according to Dijkstra from $[\ell, k]$, the path found follows the far

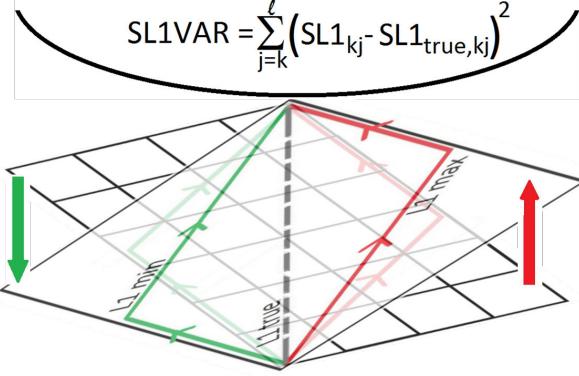


FIG. 2. Figure illustrates the two paths of $\Sigma L1_{min}$ and $\Sigma L1_{max}$, which in 2D lie within a plane. The average between them gives the straight line “true” value, $\Sigma L1_T$. When all paths within the L2 convex hull are available, Dijkstra’s algorithm will seek the $\Sigma L1_{min}$ path when going from $[k, \ell]$ and will seek the $\Sigma L1_{max}$ path going from $[\ell, k]$. This asymmetric behavior is exploited in order to find the straight line path by applying Dijkstra on a third pass, where the pathwise value applied is the squared difference between the $\Sigma L1$ and $\Sigma L1_T$, as is shown above the plane as a parabolic bowl.

177 opposite edge of the parallelotope, which is equivalent to the maximal $\Sigma L1$ from $[k, \ell]$. The two extrema of
 178 the parallelotope between $[k, \ell]$ represent the farthest paths that can be taken within the convex hull while
 179 always moving closer to the endpoint. The average between the two values of $\Sigma L1_{min}$ and $\Sigma L1_{max}$ is the
 180 true path $\Sigma L1_T$.

181 When sufficient obstacles force the paths from $[k, \ell]$ as well as $[\ell, k]$ to be along the same side of the
 182 parallelotope with respect to the true path, the path found “turns a corner” in order to reach the final
 183 partition. In this case, one of the two values, $\Sigma L1_{k,\ell}$ or $\Sigma L1_{\ell,k}$ will exceed the true path summed L1,
 184 $\Sigma L1_T$, leading to the second criteria.

185 The last criteria uses the results from the second application of Dijkstra, now, attempting to find a path
 186 that minimizes the variance of $\Sigma L1$ with respect to $\Sigma L1_T$. Calculating the value $(\Sigma L1 - \Sigma L1_T)^2$ then
 187 copying the k^{th} row of this matrix and multiplying it by every row of the logical matrix, $NN1 > 0$, a new
 188 adjacency matrix is formed and applied using Dijkstra’s algorithm for the third time. At each step, the
 189 minimal summed path variance gives the most direct path from $[k, \ell]$, finally giving the path that is LOS
 190 between the two partitions, illustrated in Fig. 2. The smallest error that can exist is when a path is found
 191 that is one step off of the true path near the middle of the path. In this case, the error is the difference
 192 between $1/2 n(n+1)$ and $1/2 (n-1)(n)$, where $n = P_C/2$, leading to the third criteria for LOS.

193 **6. Strategy.** This study employs 26 different clustering techniques to a bank of 12 representative test
 194 cases. The data sets forming the test bank were comprised of various shapes, both connected and disconnected
 195 as well as point clouds in both 2D and 3D. In each of the point clouds, four gaussian distributions were placed
 196 near one another, with three densely populated regions and a fourth low density gaussian which spans the
 197 domain. The point clouds were further varied by creating one case in 2D and 3D where the dense gaussians
 198 are clearly separated, and another two cases in 2D and 3D where the three gaussians overlap. Figure 3
 199 illustrates the test banks used, in this order: *L*, *Plus1*, *Plus2*, *Concentric1*, *Concentric2*, *Flame1*, *Flame2*,
 200 *Flame3*, *Data2D-1*, *Data2D-2*, *Data3D-1*, *Data3D-2*. Table 2 lists the test bank set as well as the features
 201 sought to examine in each case. The first test is the simple *L* as discussed in section 7.5. The *Plus1* and *Plus2*
 202 cases are extensions to the *L* case where symmetry is employed, testing how algorithms respond to symmetry
 203 as well as an open region (*Plus2*). *Concentric1* and *Concentric2* test how the routines respond to curved
 204 domains with symmetry and whether the domain is connected or not. *Flame1*, *Flame2* and *Flame3* test
 205 how asymmetry is dealt with as well as connected versus disconnected regions. *Flame3* also tests how well
 206 “tendrils” or filamentary data is handled. As a test of a 2D point cloud, *Data2D-1* and *Data2D-2* test how
 207 well four gaussian point clouds can be clustered for the case of three separated clusters, *Data2D-1*, and three
 208 close-by clusters, *Data2D-2*, where the fourth gaussian is evenly distributed across the domain simulating
 209 noise present in the data. *Data3D-1* and *Data3D-2* show the point cloud in 3D of four gaussian distributions
 210 similar in definition to the 2D cases, where in the first case are three disconnected ellipsoidal distributions

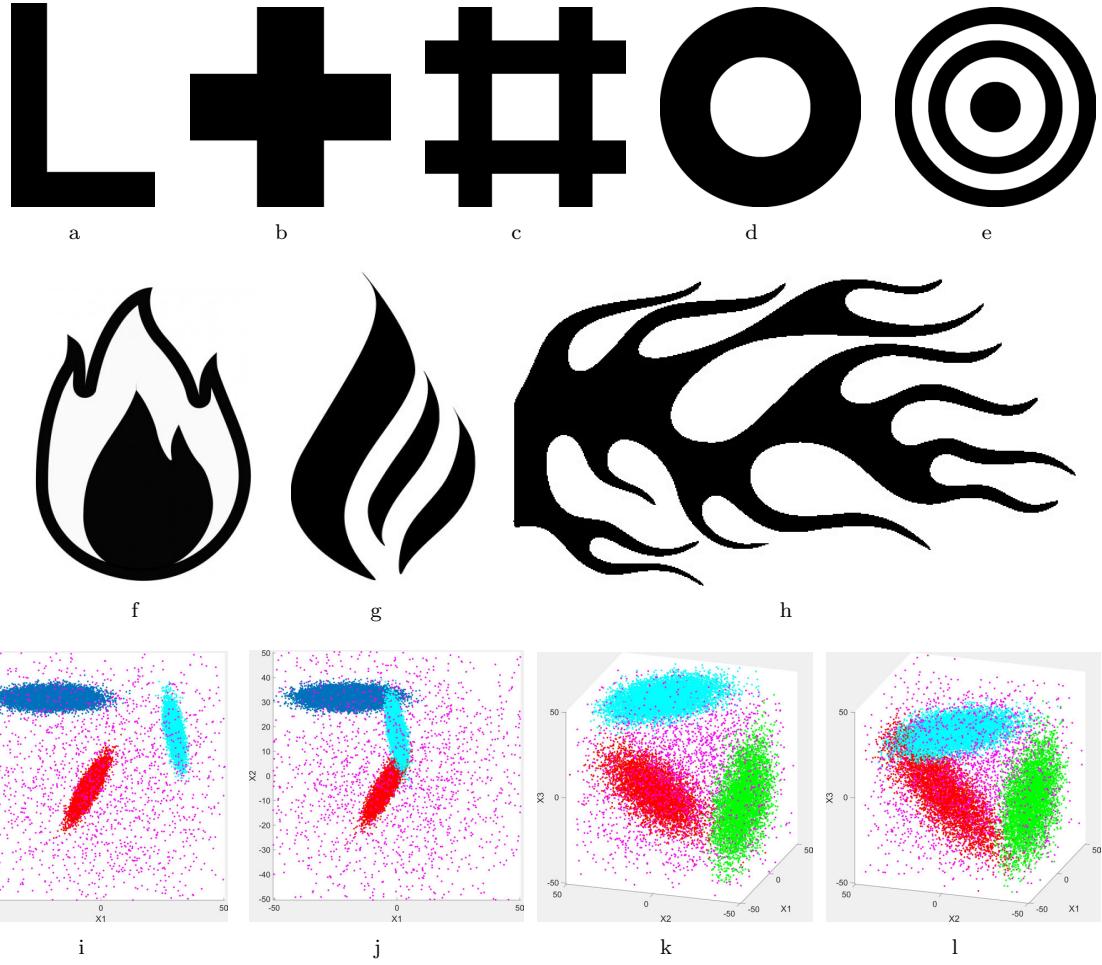


FIG. 3. Test bank of 12 shapes: L (a), Plus-1 (b), Plus-2 (c), Concentric-1 (d), Concentric-2 (e), Flame-1 (f), Flame-2 (g), Flame-3 (h), Data2D-1 (i), Data2D-2 (j), Data3D-1 (k) and Data3D-2 (l).

TABLE 2
Test bank data sets.

Labels	#	dim	size (pixels/pts)	Test Bank Data Sets					
				connected	symmetry	plateau	filamentary	overlap	noise
L	1	2D	1200x1200	✓	X	✓	X	X	X
Plus1	2	2D	1200x1200	✓	✓	✓	X	X	X
Plus2	3	2D	1200x1200	✓	✓	✓	X	X	X
Concentric1	4	2D	1200x1200	✓	✓	✓	X	X	X
Concentric2	5	2D	1200x1200	X	✓	✓	X	X	X
Flame1	6	2D	1200x1200	✓	X	✓	✓	X	X
Flame2	7	2D	1200x1200	X	X	✓	X	X	X
Flame3	8	2D	1200x1200	✓	X	✓	✓	X	X
Data2D-1	9	2D	200,000	X	X	X	✓	X	✓
Data2D-2	10	2D	200,000	✓	X	X	✓	✓	✓
Data3D-1	11	3D	200,000	X	X	X	✓	X	✓
Data3D-2	12	3D	200,000	✓	X	X	✓	✓	✓

211 with a fourth acting as noise, while in the second case shows the same three ellipsoidal distributions moved
 212 closer to one another such that two of the tails overlap. For all cases other than the point clouds, the data
 213 is derived from an image, where a binary set of points is established for all 8-bit grey-scale values above 100
 214 (1) or below (0). The image sizes when possible are 1200x1200, unless the aspect ratio prevented that exact
 215 size. The point clouds are based on four distributions with a summed value of 200,000 points.

216 Along with the 26 clustering algorithms applied, four additional cluster assignments are derived from
 217 consensus among the 26, leading to 30 differing cluster assignments per test case for a total of 360 figures
 218 showing the clustering results. These results are supplied as supplemental figures and can be found on the
 website. A sampling of these results is shown in Sec. 8.

TABLE 3

Clustering techniques for 26 algorithms highlighting requirements, pros and cons in each case. Some algorithms required partitions to be connected in order to search for clustering within the connected region. Clusters can be feature driven or can even the distribution of cluster assignments (balanced, group). LOS is a criteria for some clustering, which in turn can help identify data distributions. Finally, some algorithms treat isolated partitions on equal footing with larger connected subsets, making the clustering sensitive to these smaller subsets, interpreted as noise. Checks indicate a feature is used, “X” indicates the feature is not required whereas a “-” indicates the parameter is not applicable to the technique, finally “” indicates that population weighting could be applied to the technique or not - for the results shown in this study, weights were applied to spectral algorithms making the Laplacian sensitive to the populations of the partitions.*

Labels	#	Clustering Algorithms									
		Connected required	Proximity	Weights	Sensitive to Noise	Balanced	LOS criteria	Gathering method	Laplacian type	Eigen-vectors	Fixed k guess
KMEANS	1	X	ΔR	ww^T	X	X	X	weighted	-	-	✓
KMEDOIDS	2	X	ΔR	ww^T	X	X	X	weighted	-	-	✓
MAXGLOB	3	X	ΔR	Δw	X	X	X	slopes	-	-	-
MAXPATHL	4	✓	L_2	Δw	X	X	X	slopes	-	-	-
CONN	5	✓	X	X	X	X	X	-	-	-	-
LOS-MAXVIS	6	✓	$L_2, \Sigma L_1$	*	X	X	✓	max vis.	-	-	-
LOS-MUTUAL	7	✓	$L_2, \Sigma L_1$	*	X	X	✓	mutual vis.	-	-	-
SPECTRAL01	8	✓	X	*	✓	X	X	2D histo	NN1	1,2	-
SPECTRAL02	9	✓	X	*	✓	X	X	kmeans	NN1	1,2	✓
SPECTRAL03	10	✓	X	*	✓	X	X	kmedoids	NN1	1,2	✓
SPECTRAL04	11	✓	X	*	✓	✓	X	2D histo	NN1	2,3	-
SPECTRAL05	12	✓	X	*	✓	✓	X	kmeans	NN1	2,3	✓
SPECTRAL06	13	✓	X	*	✓	✓	X	kmedoids	NN1	2,3	✓
SPECTRAL07	14	✓	X	*	✓	X	✓	2D histo	LOS	1,2	-
SPECTRAL08	15	✓	X	*	✓	X	✓	kmeans	LOS	1,2	✓
SPECTRAL09	16	✓	X	*	✓	X	✓	kmedoids	LOS	1,2	✓
SPECTRAL10	17	✓	X	*	✓	✓	✓	2D histo	LOS	2,3	-
SPECTRAL11	18	✓	X	*	✓	✓	✓	kmeans	LOS	2,3	✓
SPECTRAL12	19	✓	X	*	✓	✓	✓	kmedoids	LOS	2,3	✓
SPECTRAL13	20	X	ΔR	*	X	X	X	2D histo	RAD	1,2	-
SPECTRAL14	21	X	ΔR	*	X	X	X	kmeans	RAD	1,2	✓
SPECTRAL15	22	X	ΔR	*	X	X	X	kmedoids	RAD	1,2	✓
SPECTRAL16	23	X	ΔR	*	X	✓	X	2D histo	RAD	2,3	-
SPECTRAL17	24	X	ΔR	*	X	✓	X	kmeans	RAD	2,3	✓
SPECTRAL18	25	X	ΔR	*	X	✓	X	kmedoids	RAD	2,3	✓
LMH-POS	26	X	X	X	X	X	X	-	-	-	-

219

220 **7. Clustering Algorithms.** This section discusses the clustering algorithms used in this paper. Some
 221 techniques are standard approaches, but several are variations on existing techniques with new methods.
 222 The new approaches involve treating the data in terms of partitions with populations of data serving as
 223 weights to the partitions. Also new, the distance metric used is changed from a traditional L2-norm to a
 224 path length along a grid of partitions. Along with investigating path length based clustering, a Line-Of-Sight
 225 criteria is also developed. An alternative approach of spectral clustering is also used, utilizing a different set
 226 of eigenvectors to establish clusters, and alternatives to the traditional Laplacian operator are used as well.
 227 Once all twenty-six clustering techniques are used to assign a cluster identity, an overall cluster identity is
 228 given to each data based on the consensus of the set of techniques, with four algorithms employed differing
 229 in degrees of consensus reached.

230 Table 3 lists the twenty-six techniques used. Each technique attempts to cluster data according to
 231 features present in the data. The table lists those features which best suit each technique. As the chief data
 232 reduction scheme here is to partition the data into multi-dimensional bins, the clustering is performed over

233 the weighted partitions on a grid. Features indicated in the table are; require partitions to be *Connected*
234 in order to cluster, *Proximity* uses distance as a criteria, *Weights* indicates populations affect the result,
235 *Sensitivity to Noise* indicates some methods fail to find structure within larger connected subsets in the
236 presence of noisy data, *Balanced* indicates methods which evenly divide partitions into clusters, *LOS* criteria
237 is required for some, *Gathering* indicates the method used to gather partitions for clustering, *Laplacian*
238 indicates which type of Laplacian is used for spectral algorithms, *Eigenvectors* indicate which modes are
239 used in gathering, and *Fixed k* requires an initial guess as to the number of clusters.

240 **7.1. K-Means and K-Medoids Clustering - KMEANS, KMEDOIDS.** K-means is a well estab-
241 lished clustering technique Forgy (1965); Lloyd (1982), seeking from a data set, the lowest possible distance
242 from individual data to a set of possible mean positions of the data, indicative of clusters. Over several
243 passes, the cluster definitions are altered to minimize the distance from each datum to clusters found. An
244 initial guess of the number of clusters to seek is required. K-means has been discussed thoroughly in the com-
245 munity for its strengths and weaknesses Jain *et al.* (1999). K-medoids has been proposed to overcome many
246 of the shortcomings of k-means and is similarly well-established in the community Kaufman and Rousseeuw
247 (2009). In both cases, an initial guess (k) as to the number of clusters sought is required which can be
248 problematic when the actual number of clusters does not match the guessed value. Further, both techniques
249 perform at $\mathcal{O}(N_D)$, which for large datasets are costly. Progress in improvements to speed have been made
250 to both techniques Park and Jun (2009); Razavi Zadegan *et al.* (2013), yet remain costly in high dimensions
251 for large data sets. By shifting the analysis from individual datum to partitions with weights, the k-means
252 and k-medoids algorithms are adjusted to accommodate the weighted bins. All calculations for distance
253 between two partitions are multiplied by the weight of each partition, $\mathbf{w}\mathbf{w}^\top$, and any centroid calculation is
254 treated as a weighted value.

255 **7.2. Maxima Clustering - Global And Path Length.** In this study, data has been reduced to
256 a set of partitions with a population assigned for each. The two schemes, MAXGLOB and MAXPATHL,
257 assign data to clusters based on how close a partition is to a significant nearby maxima among the partitions.
258 Treating the weights of the partitions as the height of a multi-dimensional map, the significance of a nearby
259 maxima is determined by calculating the slopes between any two partitions, where the slope is the ratio
260 of the weight difference, $\Delta\mathbf{w}$, to the distance between any two partitions. In the *global* case, the distance
261 used is the Euclidean distance, $\Delta\mathbf{R}$, and for the *path length* case, the distance used is the path length,
262 **L2**. MAXGLOB seeks to assign clusters between partitions that are not required to be connected, while
263 MAXPATHL requires a connection. Initially, local maxima among the partitions are found which are then
264 categorized into three types: lone peaks, ridges and plateaus. Once the maxima are classified, a peak and
265 all of the partitions associated with it are then assigned a cluster identification number, where the slopes
266 and distances from partition to peak are contributing factors in determining which peaks associate with
267 partitions. Definitions of local maxima, peaks and slopes as well as details of the algorithms for these two
268 techniques are included in the supplemental material online.

269 **7.3. Clustering via Connection - CONN.** In cases where local clusters of partitions are sparsely
270 found within the data space, a simple clustering algorithm is to determine which partitions are connected to
271 one another using first nearest neighbor steps, **NN1**. Section 4 discusses path lengths calculated from one
272 partition to another where those with a finite value are *connected*. A logical value is set between any two
273 connected partitions creating the matrix **CONN**. A unique cluster ID is assigned for each connected set of
274 partitions.

275 **7.4. Clustering by Line-Of-Sight - LOS.** Clustering by Line-Of-Sight is motivated by the idea that
276 two data within a convex region of a subset of the data have a higher chance of being correlated than data
277 outside that convex region. Considering a set of data comprised of various types of distributions, it is possible
278 for overlapping regions to form, where the tail of one distribution mingles with the tail of another. In the
279 worst case scenario, peaks of two differing distributions may overlap. Further, distributions may also form
280 along curved paths, where the peak may be far from the tails. Clustering via **CONN** will associate all data
281 in these distributions, however, checking whether two data lie within a convex hull more closely associates
282 those data with one another. The Line-Of-Sight criteria from Sec. 5 determines which partitions are convex
283 to one another. As examples, Figs. 3i-3l illustrate several distributions which have both convex regions

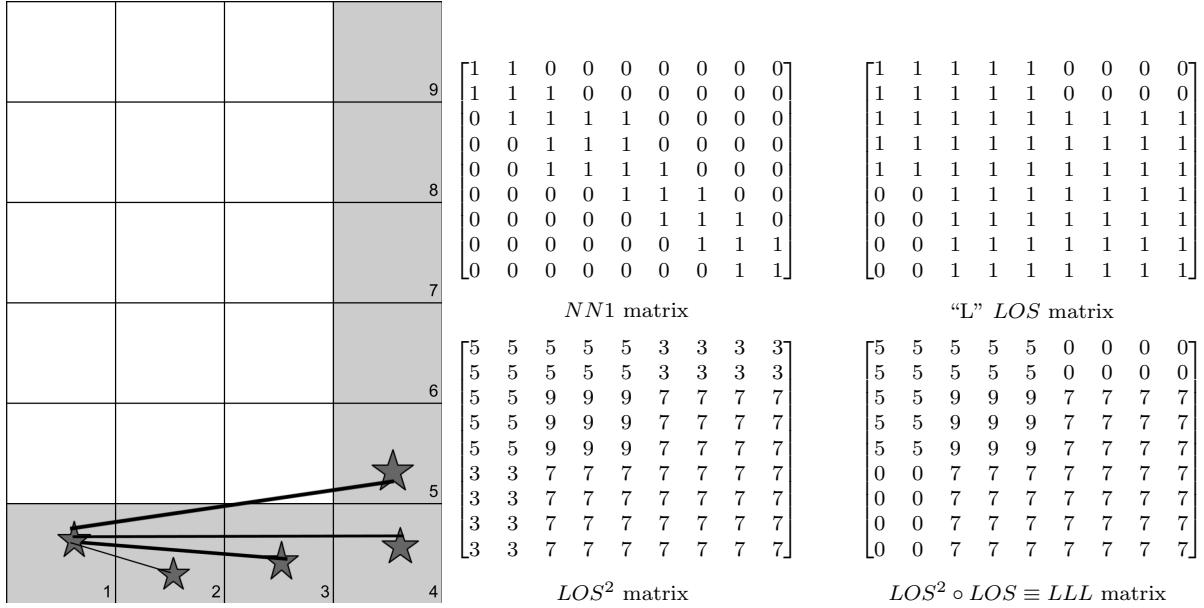


FIG. 4. Simple example to illustrate the ideas behind LOS clustering. This “L” shaped domain has nine populated bins. Starting from the bottom left to right then and moving upwards, the bins are numbered initially by rasterizing the domain, then contracting the bin indices to simply number from #1...9. Beginning with bin #1, the line-of-sight bins are indicated by the starred (*) bins, noting that bin #5 is considered LOS although a line connecting bin centers is not possible. This “corridor” condition means that for a long hallway, all partitions facing inward are also LOS. Matrices calculated for the simple example, the **NN1** (upper left) **LOS** (upper right), **LOS²** (lower left), **LOS² o LOS ≡ LLL** (lower right).

as well as overlapping tails of distributions. In this discussion, the term *visibility* refers to the number of partitions that are LOS to a specific partition. A detailed discussion of the algorithms used to form clusters based on the LOS criteria is provided by the supplemental material online.

The **LOS** matrix is formed where each row represents a partition and each column represents all other partitions where a logical value indicates whether the two are LOS, making the **LOS** matrix symmetric. Squaring the **LOS** matrix, **LOS²**, gives a matrix whose values along each row tally the number of partitions which are mutually LOS to one another. For the L example given next, in the first row, the last three partitions are not LOS to the first, yet they share three partitions that are LOS in common. In order to eliminate the entries in **LOS²** that are not present in the **LOS** matrix, a Hadamard product is taken between **LOS** and **LOS²** yielding a third matrix, **LLL**. To form clusters from the information in **LLL**, a gathering process finds partitions that meet one of two cluster criteria; *maximal visibility* finds those partitions that share a high value of visibility and are connected to one another, and *greatest mutual visibility* finds the largest sets of partitions with a common value, regardless of how high in value is their visibility.

7.5. Simple Example: L. A simple example serves to demonstrate how these matrices interact with one another. Consider a small distribution of partitions forming a 6x4 grid connected to each other in an “L” configuration as shown in Fig. 4. The serialization of partitions given in Eqn. 3.1 for the inverted L gives the partitions along the bottom row $k = 1\dots 4$, then along the vertical side $k = 5\dots 9$. In this case, there are only nine partitions connected to each other, requiring a 9x9 matrix to represent the information. As each partition is connected to all of the other partitions, the **CONN** matrix is full, with values of one. The **NN1** matrix reflects which partitions share a common geometrical feature. The **LOS**, **LOS²** and **LLL** matrices show which partitions are visible to each other. Note that partition five is visible to partition one, meaning that partitions can see the edges of one another. From the matrices shown, partitions (3,4,5) form a cluster with the maximal visibility, followed by partitions (6,7,8,9) then (1,2) (LOS-MAXVIS). Partitions (3,4,5,6,7,8,9) form a cluster with the highest mutual visibility followed by (1,2) with the lowest (LOS-MUTUAL).

309 **7.5.1. LOS Clustering With Maximal Visibility - LOS-MAXVIS.** The **LOS** matrix contains
 310 for each row the logical status of which partitions are LOS to the current partition. Further, the **LLL**
 311 matrix shows the number of mutually visible partitions within LOS of the current. From the **LLL** matrix,
 312 two values can be used to determine clustering using LOS. The highest value in the **LLL** matrix indicates
 313 which partitions are within LOS of the most other partitions. These highest valued **LLL** partitions have the
 314 *maximal visibility*, LOS-MAXVIS, of the set of partitions that are LOS. An example would be any partition
 315 that is located at an intersection of several distributions of partitions. Consider the test cases: *L* and *Plus1*,
 316 where the corner of the *L* and the center of the *Plus1* will have maximal visibility. The clusters formed in
 317 this manner find intersections and corners of data distributions preferentially, leading to *data identification*
 318 of the entangled portions of data sets arising from multiple distributions present.

319 Clustering by LOS-MAXVIS is achieved by forming a histogram from the visibility values of **LLL**, shown
 320 in Fig. 5 for the *Data3D2* test case. The horizontal axis indicates the visibility while the vertical axis is the
 321 number of partitions sharing a common visibility value. Starting from the maximal value of the visibility,
 322 a cluster is formed by taking all partitions sharing the maximum or nearby, defined by including all bins in
 323 the histogram starting from the leftmost until a minimum in the bins is reached. In the case of the simple
 324 *L*, the most visible partitions are the corner partitions with values **LLL** = 9. Of the set of partitions found,
 325 a cluster is assigned to the largest connected group of partitions. As each cluster is identified, the partitions
 326 are excluded from further searches by removing the rows and columns from **LLL** of the cluster found then
 327 recalculating the histogram. Further clusters are then identified by taking partitions associated with the
 328 next highest visibility bin in the histogram, beginning where the last set left off, and including all partitions
 329 with successively lower visibilities until the next minimum in the bins is reached. This process continues
 330 until the set of partitions is fully associated with clusters.

331 **7.5.2. LOS Clustering With Mutual Visibility - LOS-MUTUAL.** The **LLL** matrix can alterna-
 332 tively be used to cluster partitions with the *highest mutual visibility* (LOS-MUTUAL) by selecting clusters
 333 with the most common shared **LLL** value instead of the maximal value. In this manner, clusters are formed
 334 around partitions that can mutually see each other the most. From the same **LLL** histogram, starting from
 335 the bin with the most frequent visibility, a cluster is formed by seeking the minima on both sides of the peak
 336 in the histogram nearest the most populated bin. Once the lower and upper bins are found, all partitions
 337 which have any visibility values in **LLL** within this range are clustered together. Identifies partitions are
 338 removed from further searches and the process is repeated until all partitions are identified. LOS-MUTUAL
 339 clustering finds the largest set of partitions that are LOS to each other first, then searches for the next
 340 largest set of partitions that do not include the first set and so on. In the case of the simple *L*, the highest
 341 mutually visible partitions are the partitions forming the long arm of the *L*, with values **LLL** = 7. For the
 342 *Data3D2* case, all partitions with a visibility between 1700 up to 3000 are included in the first cluster found.
 343 As before, once a cluster is found, the partitions are removed from further searches. Clusters formed in this
 344 manner find full data distributions first, associating tails over mixed regions with the largest distributions
 345 first, giving an alternative to the *data identification* offered by LOS-MAXVIS.

346 **7.6. Spectral Clustering.** Spectral clustering [Chung (1997); von Luxburg (2007)] represents data
 347 as a graph, where data become vertices and relationships between data points are represented by edges
 348 and weights in the graph. The eigenmodes of the graph are sought solving Helmholtz equation from a
 349 Laplacian chosen based on the edge weights. This analysis uses the **NN1** matrix, the **LOS** matrix as well
 350 as a radial basis function to form the graphs. The Laplacian operator is a matrix formed by setting the
 351 degree of the vertex along the diagonal with off diagonal elements set to a negative weight factor. The off
 352 diagonal components are formed from either the summation of all nearest neighbors (**NN1**), the sum of
 353 all LOS partitions (**LOS**) or a radial factor related to the distance squared to all other partitions (**RAD**).
 354 **RAD** is chosen as a gaussian with a large sigma equal to the maximal distance of ΔR . Clustering with
 355 **NN1** seeks clusters as partitions connected to one another, using **LOS** seeks similar clustering for partitions
 356 connected through visibility, while **RAD** seeks clusters of partitions over a region, regardless of connection.
 357 In all cases, the analysis that follows is similar. The eigenvectors are calculated for the Laplacian, where
 358 the lowest two eigenvectors are typically used to define a *new data space* using each eigenvector as a basis.
 359 The partitions are then mapped to the eigenspace and clusters within the space are sought using novel 2D
 360 clustering techniques, either KMEANS, KMEDOIDS or a simple 2D histogram over the domain.

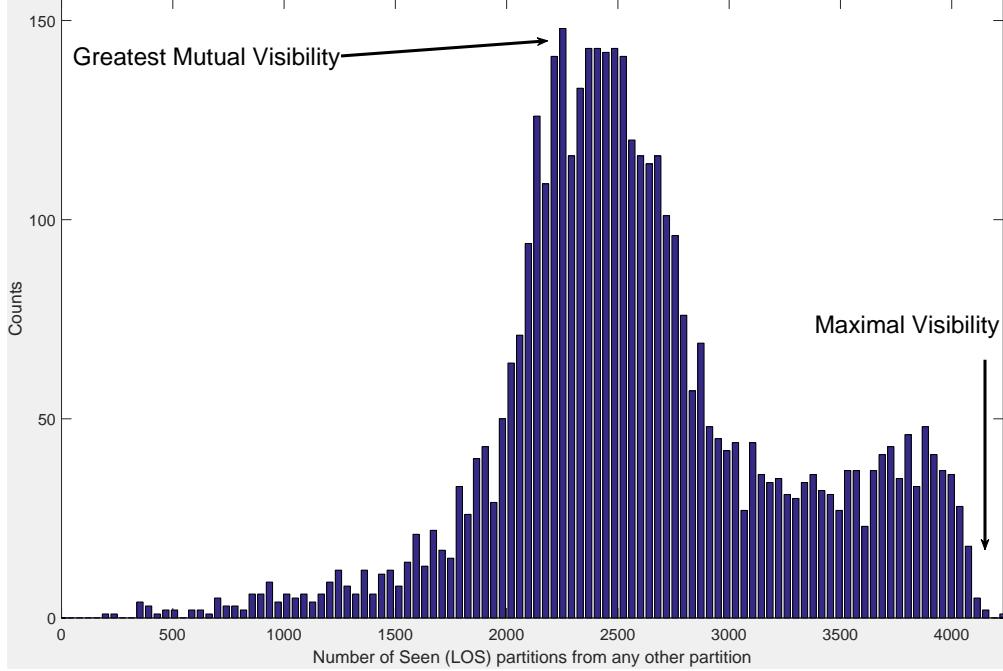


FIG. 5. Histogram of LLL values, the visibility from a partition to all others for the Data 3D-2 case. The horizontal axis are visibility values and the vertical axis is the frequency of partitions for a given visibility. For the Data 3D-2 case, the highest visibility is ≈ 4100 between partitions, where LOS-MAXVIS begins searching at the highest visibility bin and ends at the first minima found in the histogram, for a cluster with visibility from $\approx 3300 - 4100$. A LOS-MUTUAL search seeks a cluster with the greatest mutual visibility by beginning the search at the tallest peak around ≈ 2200 whose set size is ≈ 150 partitions. The cluster is formed between the two minima found on each side of the tallest peak. In both cases, the starting point for the cluster search defines which other partitions are near to the goal, either maximal visibility or greatest mutual visibility. The clusters are formed by grouping LOS partitions around the feature sought in the histogram, where peaks are separated by the basins.

361 This analysis employs all three clustering techniques in the eigenspace as well as explores using two
 362 differing sets of eigenvectors, the lowest pair (1,2) as well as the next lowest pair (2,3) as a bases. Spectral
 363 clustering finds clusters of partitions which are connected subdomains; however, when only a single connected
 364 domain is found (clean case), the eigenvectors reveal a modal structure within the connected domain. When
 365 showing the modal structure for the first case using eigenvectors (1,2), the first eigenmode accentuates a
 366 single large feature within the eigenspace, where the second eigenvector segments the space into a small
 367 number of symmetric regions. When using the next lowest pair of eigenvectors (2,3), surpassing the lowest
 368 eigenmode, the modal structure segregates the partitions differently, clustering the partitions into *evenly*
 369 *distributed groups* of data. Once the eigenspace has been populated with the partitions, k-means, k-medoids
 370 as well as traditional 2D histograms can be used to collect the partitions and assign them to cluster IDs.
 371 K-means and k-medoids have been discussed earlier in Sec. 7.1 as to their strengths and weaknesses. As an
 372 alternative approach to finding the clusters within the eigenspace, simply histogram the 2D eigenspace and
 373 assign to each non-zero bin a different cluster ID (2DHIST). This approach has the advantage of simplicity
 374 and finds exactly the number of clusters that fill bins within the eigenspace, not requiring an initial guess as
 375 the number of possible clusters, as in the case of k-means or k-medoids, however a maximum possible count
 376 of clusters is set by the number of bins of the 2D histogram, typically set at $(\sqrt{k} + 2) \times (\sqrt{k} + 2)$ so that the
 377 k-means and k-medoid searches are comparable to the size of the clusters sought.

378 **7.7. Clustering By Coarse Position (LMH-POS).** The most obvious form of clustering is to
 379 associate a partition solely by its *position* (LMH-POS) using a coarse binning within the partition space. By
 380 setting the number of bins along each dimension to three, the bins are interpreted as being *Low*, *Medium* or
 381 *High* for the values represented along each axis. In this case, the *sequential partition bin index*, k , becomes
 382 the cluster ID, with the maximum number of possible clusters at 3^{N_D} , for the three bins along each axis.

383 This approach is a coarse designation for clustering as it employs no complicated algorithms, and data
384 with similar values are associated irrespective of all other factors. This approach suffers from many problems
385 in that data in one bin will not be clustered with data from a neighboring bin no matter how close in
386 proximity the two are to one another. Clusters from LMH-POS characterize data in the crudest sense with
387 no refinement for the shape of a distribution or even the relative sizes of the distribution. One advantage to
388 this approach is that it is easy to understand, even while spanning multiple dimensions, making it an easy
389 entry point for a discussion of the data. When handling large data sets, this approach allows for a quick
390 look at where the data reside within the larger space.

391 **8. Results.** This sections shows a sampling of results from the application of 26 techniques to 12
392 test cases. The strengths and weaknesses of these techniques are exposed leading to the conclusion that a
393 consensus approach is reasonable. Ideally, all clustering techniques plus the four robust consensus results of
394 each test case would be presented, leading to 360 figures, but due to space limitations, the full set of clustering
395 results are provided in the supplemental material. Throughout this section, the term “noise” refers to data
396 sets where a significant number of isolated small subsets of partitions, including singletons, are present, while
397 the term “clean” refers to data sets without these smaller subsets. Among the supplemental material, for
398 each data set, a high data threshold, $\Theta_{perc} = 2\%$, and a low threshold, $\Theta_{perc} = 0\%$, are applied showing how
399 clustering is achieved in a clean versus noisy environment respectively. Data clustering is also shown at two
400 different bin resolutions to illustrate how too fine of a resolution may not achieve good clustering. Finally,
401 the figures are grouped into fullpage comparisons for a single test case with all 30 techniques shown as well
402 as single page comparison for each technique across all test cases, leading to 2880 figures over 168 pages.

403 **8.1. Discussion of Techniques.** Of the sampled results provided, the first figure in each case shown
404 is for k-medoids clustering using ($k = 16$) unless otherwise stated in order to give a comparison between
405 established clustering and other approaches. The remaining figures are chosen to demonstrate a particular
406 trait of a clustering technique. Figure 12 shows all of the techniques applied to the *Data2d2* test case. When
407 appropriate, a circle is shown within a cluster to indicate the medoid of the cluster.

408 K-means and k-medoids results are well understood for both their strengths and weaknesses. MAXGLOB
409 and MAXPATHL tend to mirror results from k-means and k-medoids with the exception that MAXPATHL
410 is restricted to clustering within a connected set, making it seek clusters following a distributions’ shapes
411 rather than just using proximity between data. The CONN technique clusters data within a connected
412 set regardless of other criteria. LOS-MAXVIS and LOS-MUTUAL cluster according to data within convex
413 hulls, seeking similar visibility features as part of the gathering criteria to form clusters. Spectral techniques
414 form clusters within the eigenspace formed from two eigenvectors. The adjacency matrix used to form the
415 Laplacian determines the nature of the neighbors used, traditionally a first nearest neighbor, however, this
416 study employs both the LOS criteria to define “neighbors” as well as a radial basis. Further, the choice of
417 eigenvectors used to form the eigenspace determines whether a prominent feature is clustered about (using
418 the 1st and 2nd eigenvectors), or a more evenly distributed clustering is achieved using the 2nd and 3rd
419 eigenvectors. Due to the number of variations in spectral clustering, the techniques are identified by an
420 index given in table 3, while in the text, a shorthand will be used: ([1,2],NN1,2DHIST) to represent the
421 use of the 1st and 2nd eigenvectors, utilizing a Laplacian based on an adjacency matrix derived from the
422 first nearest neighbor matrix **NN1** and gathering the partitions into clusters within the eigenspace using
423 a 2D histogram (6x6) bins. Other variations in the notation are: [2,3] for eigenvectors, LOS or RAD for
424 the adjacency matrix, and “kmeans” or “kmedoids” to be used in the gathering of partitions within the
425 eigenspace. Finally, the LMH-POS technique clusters using a coarse resolution (Low-Medium-High valued)
426 for the binning choice, simply separating the domain into three bins per dimension to get a quick look at
427 how the data is distributed.

428 Figures 6-8 show the clustering results for the data sets derived from images, where one datum exists
429 for each pixel turned on. After binning, these test cases generally have flat distributions, so the clustering
430 results reflect geometrical features, useful for showing data grouping. Figures 9-12 show the clustering
431 results for simulated data sets for ellipsoidal distributions, where the data is unevenly distributed, some with
432 overlapping tails, helpful in illustrating data identification.

433 **8.2. Concentric1.** Figure 6 shows clustering for the *Concentric1* case, where a high degree of symmetry
 434 is present as well as a large obstruction in the middle of the data. The data is evenly distributed across
 435 the domain, so the clustering techniques fall into two groups, those that adhere to the symmetry of the
 436 domain and those that break the symmetry. K-medoids (6a) shows the attempt of creating 16 clusters
 437 that almost respect the symmetry of the data set. LOS-MAXVIS (6b) finds clusters based on highest
 438 visibility first, where it assigns clusters to a set of three subsets first, then proceeds to find further clusters
 439 within the remaining set of data, creating an odd symmetry around the ring. LOS-MUTUAL (6c) finds
 440 clusters based on the largest set of partitions with visibility values in common, leading to one large cluster
 441 formed, then the next largest, and so on until all data are clustered, making this a “greedy” algorithm,
 442 taking the largest pieces first. SPECTRAL01 (6d) finds clusters ([1,2],NN1,2DHIST), which finds a single
 443 large subset of the data first based on modes, then proceeds to cluster to smaller subsets. This approach
 444 has the effect of creating clusters that “stripe” the domain starting from the large feature to the smaller
 445 features. SPECTRAL06 (6e) clusters ([2,3],NN1,kmedoids) balance the assignment of clusters to data,
 446 creating more evenly spaced clusters. SPECTRAL07 (6f) clusters ([1,2],LOS,2DHIST) have the effect of
 447 finding large features based on visibility first, then smaller visibility clusters next. SPECTRAL12 (6g)
 448 clusters ([2,3],LOS,kmedoids) attempts to balance the LOS assignments. SPECTRAL15 (6h) finds clusters
 449 ([1,2],RAD,kmedoids) Laplacian based on a gaussian to assign clusters, leading to a set of clusters formed
 450 around a central location - in this case, the center of the ring.

451 **8.3. L and Plus1.** Figure 7 shows clustering for the *L* and *Plus1* cases, where the *L* case is an extension
 452 of the discussion given in the simple *L* example (Sec. 7.5) and the *Plus1* case can be viewed as either an
 453 “intersection” or as a tiling of the *L* case, where a four-fold symmetry exists. K-medoids are shown first
 454 in both cases (7a, 7e) showing typical clustering based on proximity. For *L*, LOS-MUTUAL (7b) shows
 455 the greedy nature of the LOS mutual technique, grabbing all of the partitions along the long axis as the
 456 first cluster, then the remaining short axis partitions that are left behind are clusters. The partitions
 457 in the intersection are not separately identified or shared in the clustering. SPECTRAL01 (7c) clusters
 458 ([1,2],NN1,2DHIST) by striping the “L” with the two largest features at the ends of the L. SPECTRAL07
 459 (7d) finds clusters ([1,2],LOS,2DHIST) where the corner has the highest visibility such that the clusters group
 460 according to visibility - as was shown numerically in the simple *L* example. For the *Plus1* case, SPECTRAL01
 461 (7f) finds a curved arrangement in the clusters ([1,2],NN1,2DHIST), yet still finding two large subsets first,
 462 followed by smaller featured subsets last. SPECTRAL04 (7g) shows the clusters ([2,3],NN1,2DHIST) with
 463 similar curved features, where the overlap of the curved sub-domains break the data into symmetric clusters.
 464 SPECTRAL12 (7h) shows clusters ([2,3],LOS,kmedoids) which extend by symmetry those found in “L” set,
 465 where the large central diamond is the extension of the corner triangle from Fig. 7d.

466 **8.4. Flame1 and Flame3.** Figure 8 shows clustering for the *Flame1* and *Flame3* cases, where sym-
 467 metry is not present yet filamentary features are present with both close and larger gaps as well. K-medoids
 468 are shown first in both cases (8a, 8e) showing typical clustering based on proximity, worth noting is that
 469 when gaps become close, k-medoids will create a single cluster on both sides of the gap due to proximity. For
 470 *Flame1*, LOS-MUTUAL (8b) finds clusters in subsets grouped by visibility, finding the most in common first.
 471 The large group of small clusters can be assigned to nearby larger clusters, however, this study did not focus
 472 on this level of refinement, mainly the viability of the algorithm to seek clusters. SPECTRAL12 (8c) clus-
 473 ters ([2,3],LOS,kmedoids) using visibility as well finding similar groups with small deviations. SPECTRAL15
 474 (8d) clusters ([1,2],NN1,kmedoids) using a radial basis clustering around a central location in the middle
 475 of the flame. For *Flame3*, LOS-MAXVIS (8f) finds clusters according to maximal visibility first exposing
 476 long clusters within the filamentary portions of the flame. SPECTRAL06 (8g) clusters ([1,2],NN1,kmedoids)
 477 larger features first irrespective of gaps in the data, striping to smaller to features.

478 **8.5. Data2D-1.** Figure 9 shows clustering for the *Data2D-1* case, where little symmetry is present in
 479 a noisy environment at a high bin resolution where no distributions overlap. K-medoids with $k = 6$ (9a)
 480 illustrates how data near the tail of an elongated ellipsoid can be associated with a different ellipsoidal
 481 distribution if the center of the second distribution is closer to the datum than the center of the distribution
 482 to which it belongs. The red and black square highlight this situation. MAXGLOB (9b) clusters similar to k-
 483 medoids without the problem discussed, it does not require a connected set to form clusters and is sensitive

484 to weighted partitions. MAXPATHL (9c) is also sensitive to weighted partitions but further requires a
 485 connection between data to form a cluster, but also is sensitive to how far within the connected set a datum is
 486 to the closest maximal density of data. LOS-MAXVIS (9d) find clusters based on visibility which also requires
 487 connectivity within the partitions in order to cluster, leading to many “island” clusters. SPECTRAL01 (9e)
 488 clusters ([1,2],NN1,2DHIST) form around the three ellipsoids, however, it also combines all of the smaller
 489 subsets into a cluster with one the larger subsets, resulting from using 2DHISTO as a gathering mechanism
 490 in the eigenspace, where the location of the large ellipsoid is too close the locations of the smaller subsets to
 491 distinguish them from one another. SPECTRAL02 (9f) clusters ([1,2],NN1,kmeans) shows similar clustering,
 492 where the k-means algorithm identified all of the connected smaller subsets first, leaving the three large
 493 ellipsoids to be given a single cluster ID. This failure is inherent to the spectral techniques in the presence
 494 of noise, where the high multiplicity of the lowest eigenvalue can lead to clustering that is hard to interpret.
 495 A remedy would be to order the eigenvectors with common eigenvalues by the number of non-zero elements,
 496 favoring the larger clusters first. SPECTRAL07 (9g) clusters ([1,2],LOS,2DHIST) identifies the three large
 497 ellipsoids more consistently than SPECTRAL02, among the large set of singleton partitions. SPECTRAL15
 498 (9h) clustering ([1,2],RAD,kmedoids) is less sensitive to singleton and small subsets as the adjacency matrix
 499 correlates partitions from disconnected regions, such the clusters formed are showing the modal structure
 500 rather than the connected structure.

501 **8.6. Data3D-1.** Figure 10 shows clustering for the *Data3D-1* case, where little symmetry is present
 502 in a clean environment at a high bin resolution where no distributions overlap. K-medoids (10a) shows
 503 16 clusters found in three main ellipsoids, with k-means (10b) giving similar yet different results. CONN
 504 (10c) clusters partitions connected to one another, which in a clean environment finds three ellipsoidal
 505 distributions, however, some partitions may have been “cutoff” from the main ellipsoids due to the higher
 506 data threshold placed, leading to singleton clusters. MAXPATHL (10d) shows similar results to CONN,
 507 however, additional clusters are found due to local maxima in the weighted partitions, leading to smaller
 508 clusters near the edge of the subsets along with the three main ellipsoids. LOS-MAXVIS (10e) clusters by
 509 maximal visibility first, finding variation within the three main ellipsoids based on visibility. SPECTRAL01
 510 (10f) ([1,2],NN1,2DHIST) again finds difficulty finding the three ellipsoids in the presence of noise, whereas
 511 SPECTRAL07 (10g) ([1,2],LOS,2DHIST) under the same conditions finds the three clusters. ROBUST2
 512 (10h) clusters according to a consensus of 50% of the algorithms in agreement of cluster IDs.

513 **8.7. Data3D-2.** Figure 11 shows clustering for the *Data3D-2* case, where little symmetry is present
 514 in a clean environment at a high bin resolution where two distributions overlap. K-medoids (11a) shows
 515 16 clusters found in three main ellipsoids, with k-means (11b) giving similar yet different results. CONN
 516 (11c) clusters partitions connected to one another, which in a clean environment finds three ellipsoidal
 517 distributions, however, some partitions may have been “cutoff” from the main ellipsoids due to the higher
 518 data threshold placed, leading to singleton clusters. MAXPATHL (11d) shows similar results to CONN,
 519 however, additional clusters are found due to local maxima in the weighted partitions, leading to one cluster
 520 which follows the contour of the merged ellipsoids. LOS-MAXVIS (11e) clusters by maximal visibility
 521 first, finding the intersection as a cluster first followed by clusters based on lesser visibility. SPECTRAL01
 522 (11f) ([1,2],NN1,2DHIST) finds difficulty again for the case of three ellipsoids in the presence of noise while
 523 SPECTRAL07 (11g) ([1,2],LOS,2DHIST) under the same conditions only finds the two clusters among the
 524 three distributions. ROBUST2 (11h) clustering again performs well over the shortcomings of the individual
 525 techniques by using a consensus of 50% of the algorithms in agreement of cluster IDs.

526 **8.8. Data2D-2.** Figure 12 shows clustering for the *Data2D-2* case, where little symmetry is present
 527 in a noisy environment at a high bin resolution where three distributions overlap. All 26 techniques are on
 528 display allowing for a cross-comparison along with the four robust algorithms in the last row (slightly larger).
 529 Of the 26 algorithms, 14 have been discussed in the previous test cases. Among the spectral techniques in
 530 a noisy environment, SPECTRAL01-SPECRTAL12, when the number of clusters sought is less than the
 531 high multiplicity of the lowest eigenvalue, without additionally specifying the order of the eigenvectors, the
 532 possibility exists that some larger subsets of the data will not be clustered as expected, leading to smaller
 533 subsets assigned to clusters otherwise seen as noise (singletons). The radial basis spectral techniques do
 534 not suffer from this confusion as they do not require a connection between the partitions in order to form

535 clusters, allowing the approach to be sensitive to the larger structure within the domain, creating clusters
536 around centroids within the data.

537 **9. Robust Clustering over Multiple Algorithms.** In this paper, multiple clustering algorithms
538 have been presented and applied to several test cases. Each technique has strengths as well as weaknesses
539 which have been exposed through the cases presented. When using multiple techniques, the possibility exists
540 to leverage the information gathered from all techniques to arrive at a final cluster designation, based on the
541 level of agreement or disagreement found between the algorithms Strehl and Ghosh (2003). This approach
542 is comparable to ensemble modeling used in various fields Hansen (2002); Tebaldi and Knutti (2007). This
543 section proposes four possible robust ways to gather the cluster information and assign new cluster IDs.

544 In each approach taken, the cluster information for the partitions is represented by a matrix of cluster
545 IDs, where each row represents results from a single cluster algorithm and each column is a partition. The
546 values along each row are the cluster IDs assigned to each partition, forming the matrix, $\{\text{CLUS} \in \mathbb{N}^{C \times P}\}$
547 where $C = 26$ and P is the number of partitions. In order to find agreement or disagreement between cluster
548 IDs across many techniques, the rows are sorted so that the cluster IDs are sorted in ascending order along
549 the first column. For any repeated values in the first column, the next column is then sorted in a similar
550 fashion, continuing to sort further columns until all repeated values are addressed. Table 4 illustrates this
551 process for a sample of 40 partitions using six cluster algorithms. The top matrix is the initial partition
552 cluster ID matrix unsorted. The second matrix is the sorted cluster ID matrix described above. Finally,
553 the third matrix from the top shows the differences in cluster IDs *along each row*, where a one represents
554 a change in cluster designation for that rows' technique. The process of assigning cluster IDs to partitions
555 begins with the lowest numbered cluster IDs over all algorithms, and proceeds in increasing cluster ID order.
556 In the table shown, this is equivalent to following the partitions from left to right across the page.

557 As examples of robust clustering, the last four figures from Fig. 12 as well as Tab. 4 are provided to illus-
558 trate the process. These figures show the results from a consensus using all clustering techniques excluding
559 the LMH-POS algorithm for the *Data2D2* test case with no minimal population set for the partitions. The
560 LMH-POS technique was excluded as its partition definitions do not align with the remaining 25 algorithms.
561 In cases where multiple techniques are compared using differing partition sizes, the robust technique is then
562 applied *per datum*, using the same procedures, however, the sorting is performed over all data instead of
563 partitions. The last four figures from Fig. 12 show the following consensus techniques, from left to right,
564 the *Fractured*, *Majority Changed* (75%), *All Changed* (100%) and *No Overlap* cases.

565 The *Fractured* robust designation results by assigning each partition a new cluster ID starting from one
566 and increasing the cluster ID each time *any* technique changes its ID. This results in the largest set of clusters
567 found. This approach is the most sensitive to changes in the cluster designations. The *Majority Changed*
568 robust technique assigns a new cluster ID each time the accumulated number of algorithm cluster ID changes
569 reaches a majority of the total number of algorithms. For each clustering technique, when a change occurs,
570 any further changes from that technique are not registered until a majority is reached, at which point the
571 accumulated sum of changes is reset to zero. This results in a medium sized set of clusters found, where a
572 significant number of algorithms found a change, however, not all algorithms are required to note the change
573 in ID. In the figure, a 75% majority was required, where ideally, the best majority threshold would create the
574 largest number of clusters with the highest average membership. The *All Changed* robust case is equivalent
575 to the *Majority Changed* case with a 100% majority threshold. This results in a small-medium sized set of
576 clusters found, where every algorithm found a change, however, the changes may not have been at the same
577 partition number, merely, that the total set of changes across all algorithms eventually required a change of
578 ID. The *No Overlap* robust case assigns a new cluster ID whenever the total number of algorithms changes
579 designation *simultaneously*. This results in the smallest sized set of clusters found, where every algorithm
580 must find a change for all partitions in a subset. Ideally, this would happen for each disconnected group of
581 partitions, however, several techniques are “global” in scope and do not require a connection to exist to form
582 clusters, leading to a single large cluster.

583 Several of the clustering techniques used in this study require either a guess or fore-knowledge of the
584 number of clusters sought, such as KMEANS and KMEDOIDS. Robust clustering can provide a reasonable
585 guess for the k-value, by first attaining consensus over all techniques that do *not* use a k-value, which are:
586 MAXGLOB, MAXPATHL, CONN, LOS-MAXVIS, LOS-MUTUAL. Using the *Majority Changed* technique

TABLE 4

A sample set of partitions that have had six differing cluster algorithms applied. In each case, the cluster algorithm identified up to nine different clusters. The set contains 40 partitions. The top table represents the data initially unsorted. Each row is a different cluster algorithm and each column is a partition where a cluster ID has been assigned. The second table has sorted each row while maintaining the assignments to each partition. The third table from the top are the differences in cluster ID assignments from one column to the next. The fourth table is the final cluster assignment given to the partitions when any one change occurs (a disagreement) between the cluster algorithms. The fifth table requires a 50% majority of the cluster algorithms to change (cumulatively) before a new cluster assignment is designated. The sixth table only changes the cluster assignment once all cluster algorithms cumulatively have changed. Finally, the last table requires that all algorithms change assignments simultaneously before a new cluster ID is designated (the clusters are disjoint - with no overlap).

with a suitable choice in consensus threshold, the number of clusters found can be used as a k -value, which allows a reasonable guess to re-run the analysis utilizing the full complement of techniques.

10. Conclusions. A study using 26 clustering techniques has been performed over 12 test cases to illustrate both the strengths and weaknesses of clustering algorithms. A robust form of clustering is achieved through consensus over all techniques, helping reduce clustering problems by finding consistent clustering definitions across many approaches. The approach taken by this study utilizes six main ideas to produce a robust clustering analysis:

- Reduce a large data set by binning the space, where the filled bins are the multi-dimensional partitions of the data set, each with a unique serial index, k .
 - Algorithms use the path length between any connected partitions as well as traditional distance metrics (L_1 , L_2 , etc...).
 - A Line-Of-Sight (LOS) algorithm is developed to enhance the probability that two data are associated with one another. LOS also provides a new “super” neighborhood definition to be used in graph-based techniques. Data identification is addressed in two differing ways by LOS.
 - Spectral clustering using the [2,3] eigenvectors addresses data grouping better than other methods.
 - Employ multiple clustering techniques to the set of partitions based on first nearest neighbors, distance weighted factors and geometrical properties of the set.
 - Establish a final cluster ID based on all the consensus of techniques employed.

605 This study shows that high dimensional, big-data analysis can be reduced to a smaller set of partitions where
 606 multiple clustering techniques can be used to sort the data into clusters. While the techniques presented are
 607 all computationally $\mathcal{O}(N_P^2)$, by reducing the data set to partitions, these routines are reasonable to perform.
 608 The introduction of the LOS criteria created new avenues for cluster seeking. The combination of multiple
 609 clustering techniques, various distance metrics and traditional data reduction lead to a robust set of clusters
 610 found, which worked well in addressing issues of data identification, clustering as well as grouping.

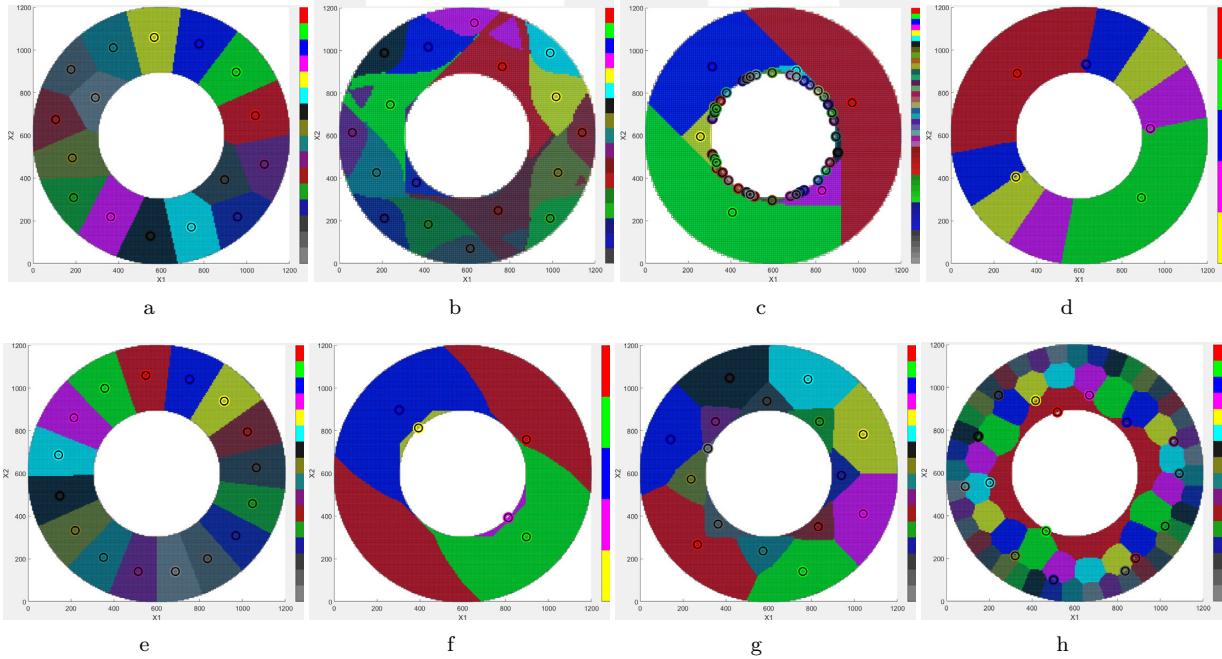


FIG. 6. Test bank case for Concentric1 showing the following techniques: 6a) KMEDOIDS ($k = 16$), 6b) LOS-MAXVIS, 6c) LOS-MUTUAL, 6d) SPECTRAL01, 6e) SPECTRAL06, 6f) SPECTRAL07, 6g) SPECTRAL12, 6h) SPECTRAL15.

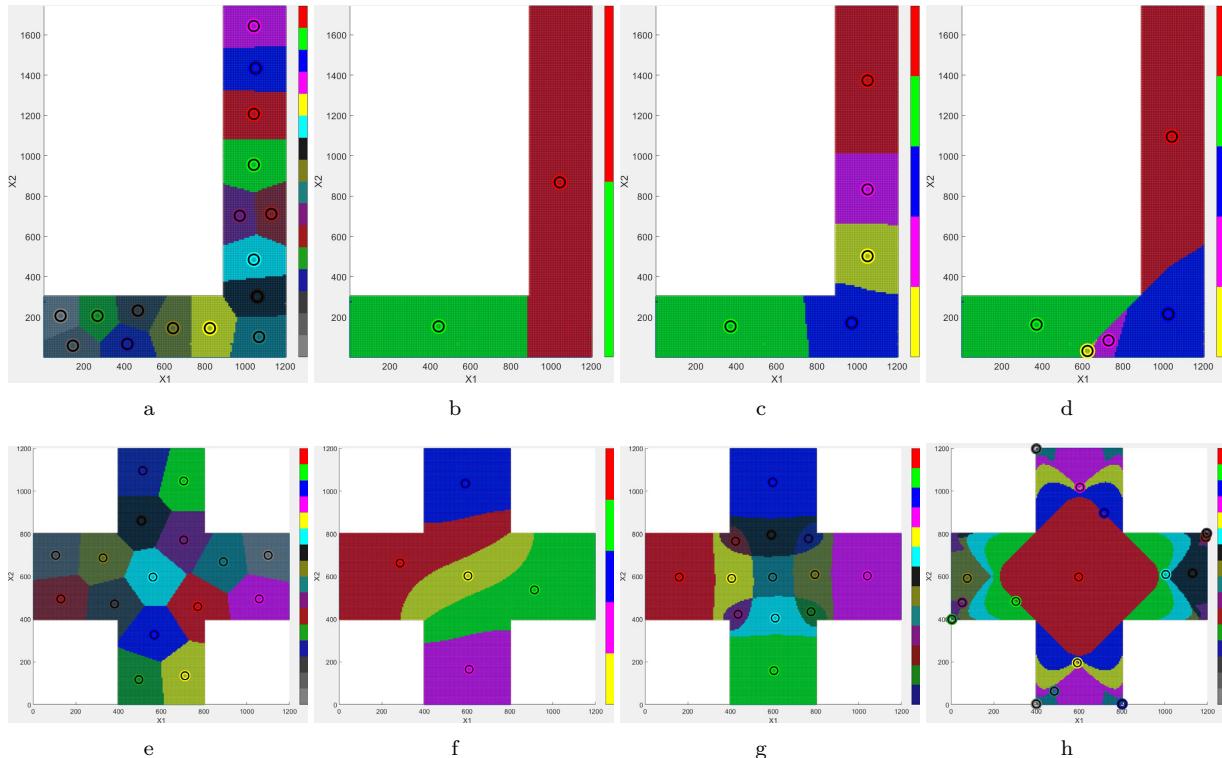


FIG. 7. Test bank case for L showing the following techniques: 7a) KMEDOIDS ($k = 16$), 7b) LOS-MUTUAL, 7c) SPECTRAL01, 7d) SPECTRAL07. Test bank case for Plus1 showing the following techniques: 7e) KMEDOIDS ($k = 16$), 7f) SPECTRAL01, 7g) SPECTRAL04, 7h) SPECTRAL12.

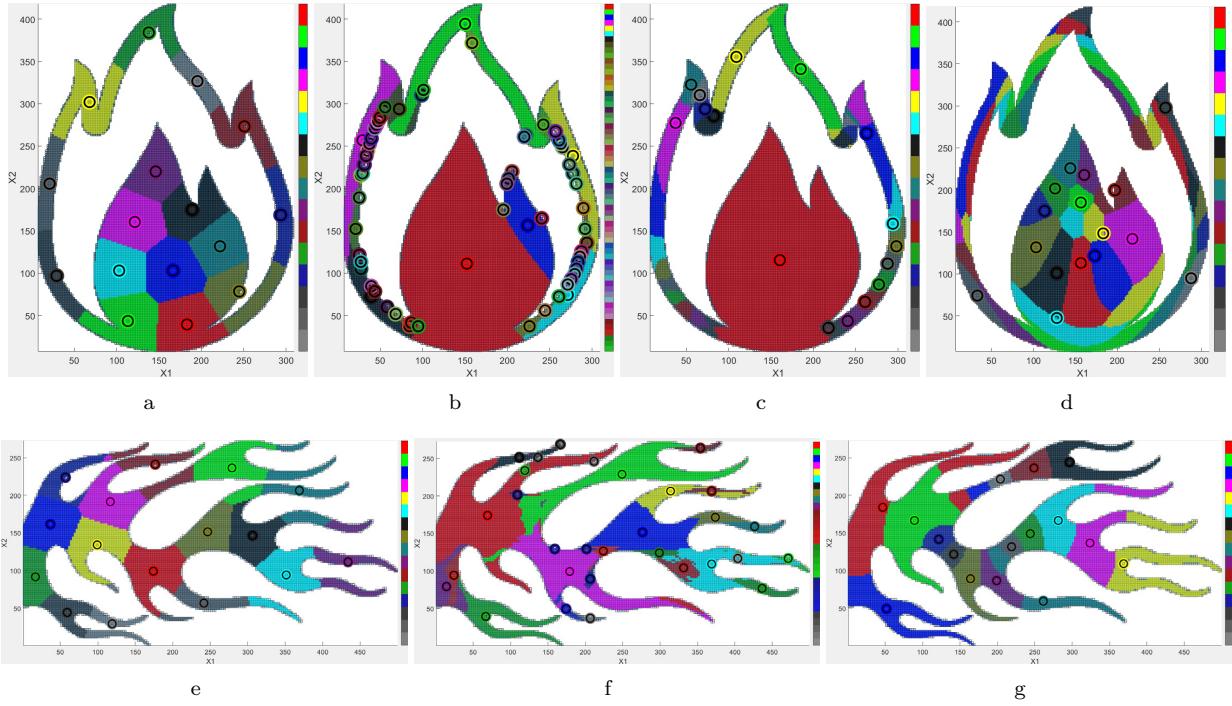


FIG. 8. Test bank case for Flame1 showing the following techniques: 8a) KMEDOIDS ($k = 16$), 8b) LOS-MUTUAL, 8c) SPECTRAL12, 8d) SPECTRAL15. Test bank case for Flame3 showing the following techniques: 8e) KMEDOIDS ($k = 16$), 8f) LOS-MAXVIS, 8g) SPECTRAL06.

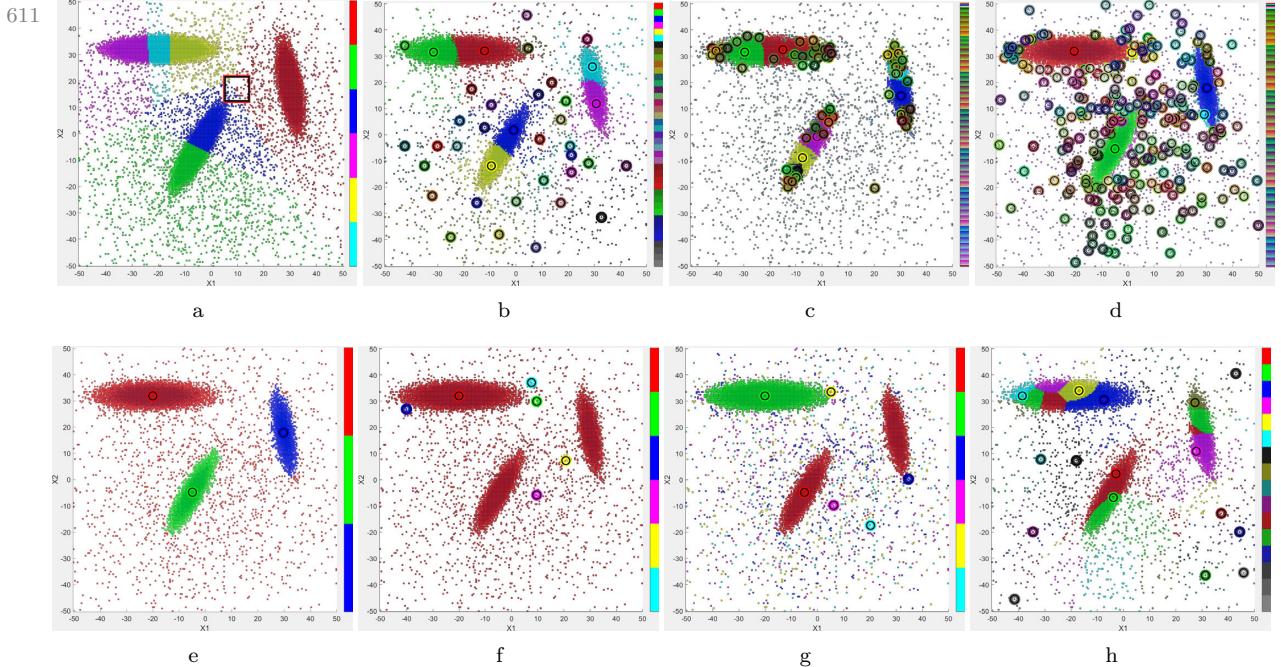


FIG. 9. Test bank case for Data2D1 in a noisy environment, $\Theta_{perc} = 0\%$, using a high number of bins showing the following techniques: 9a) KMEDOIDS ($k = 6$), 9b) MAXGLOB, 9c) MAXPATHL, 9d) LOS-MAXVIS, 9e) SPECTRAL01, 9f) SPECTRAL02, 9g) SPECTRAL07, 9h) SPECTRAL15.

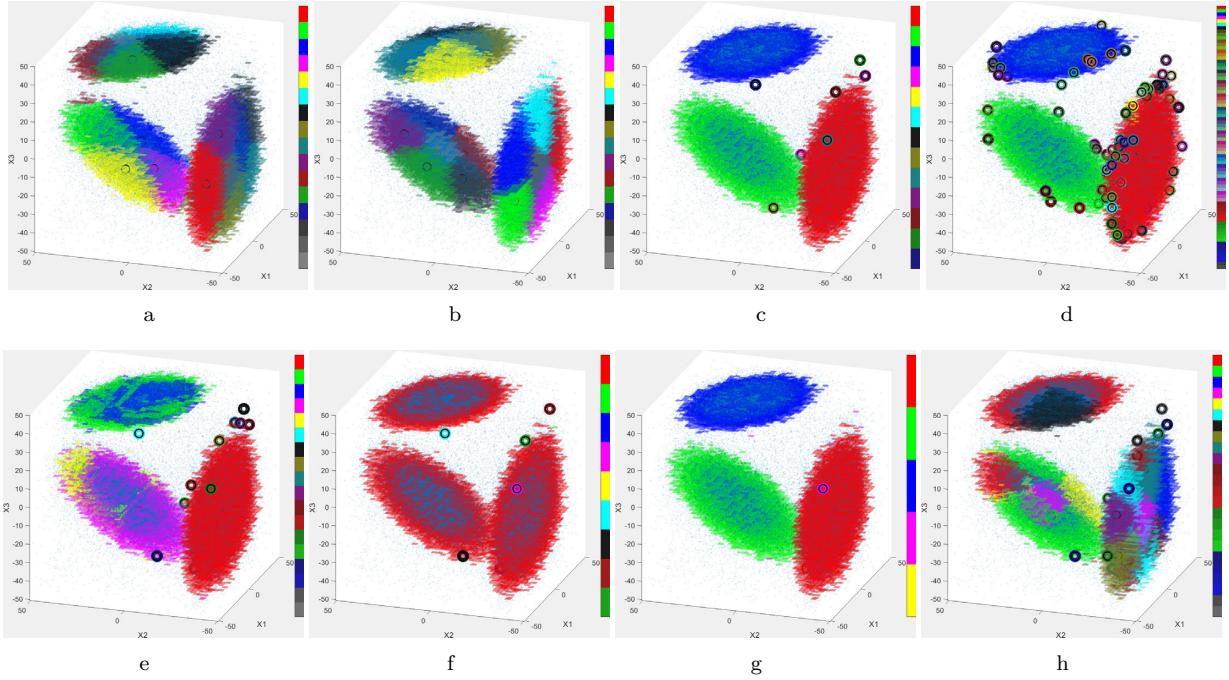


FIG. 10. Test bank case for Data3d1 in a low noise environment, $\Theta_{perc} = 2\%$, showing the following techniques: [10a](#)) KMEDOIDS ($k = 16$), [10b](#)) KMEANS ($k = 16$), [10c](#)) CONN, [10d](#)) MAXPATHL, [10e](#)) LOS-MAXVIS, [10f](#)) SPECTRAL01, [10g](#)) SPECTRAL07, [10h](#)) ROBUST2.

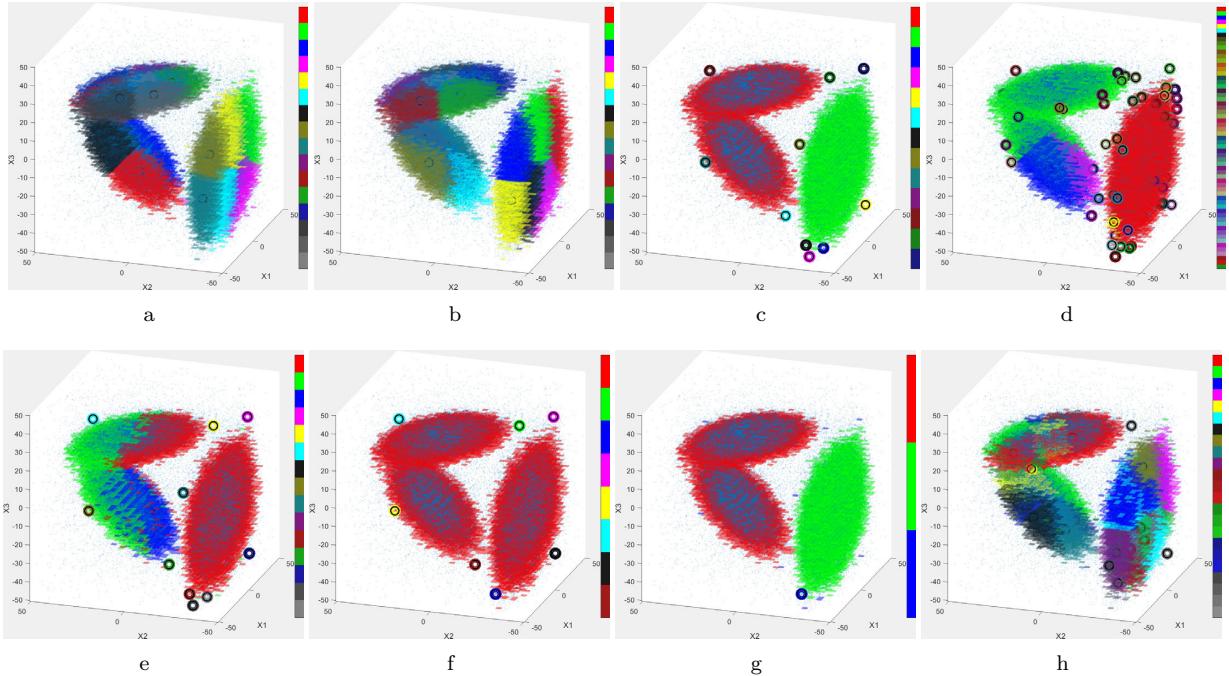


FIG. 11. Test bank case for Data3d2 for a low noise environment, $\Theta_{perc} = 2\%$, showing the following techniques: [11a](#)) KMEDOIDS ($k = 16$), [11b](#)) KMEANS ($k = 16$), [11c](#)) CONN, [11d](#)) MAXPATHL, [11e](#)) LOS-MAXVIS, [11f](#)) SPECTRAL01, [11g](#)) SPECTRAL07, [11h](#)) ROBUST2.

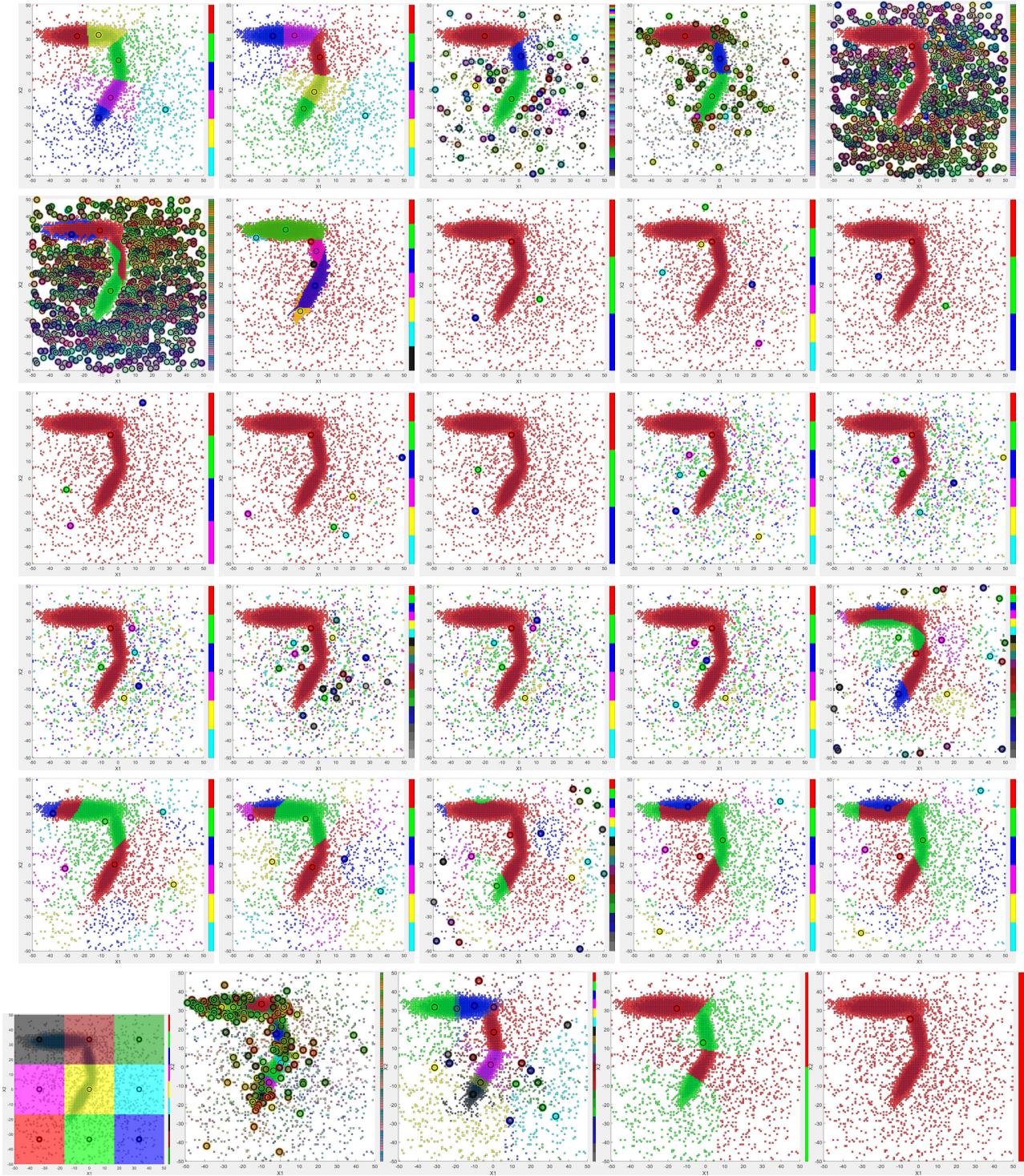


FIG. 12. Clustering techniques applied to the Data2D2 set in a noisy environment, $\Theta_{perc} = 0\%$, at a high number of bins. Clustering techniques are shown in the following order (from left-right, top-bottom): KMEDOIDS, KMEANS, MAXGLOB, MAXPATHL, CONN, LOS-MAXVIS, LOS-MUTUAL, SPECTRAL01-18, LMH-POS. On the last row, robust clustering results are shown slightly larger, from left to right, ROBUST1 - fractured, ROBUST2 - majority, ROBUST3 - all changed and ROBUST4 - no overlap.

- 612 **References.**
- 613 Barbakh, W. A., Wu, Y., and Fyfe, C. (2009). Review of clustering algorithms. In *Non-Standard Parameter*
 614 *Adaptation for Exploratory Data Analysis*, pages 7–28. Springer.
- 615 Binder, D. A. (1978). Bayesian cluster analysis. *Biometrika*, **65**(1), 31–38.
- 616 Chung, F. R. K. (1997). *Spectral Graph Theory*. American Mathematical Society.
- 617 Djikstra, E. (1959). A note of two problems in connexion with graphs. *Numerische Matematik*, **1**, 269–271.
- 618 Forgy, E. (1965). Cluster analysis of multivariate data: efficiency versus interpretability of classifications.
 619 *Biometrics*, **21**, 768–780.
- 620 Hansen, J. A. (2002). Accounting for model error in ensemble-based state estimation and forecasting. *Monthly*
 621 *Weather Review*, **130**(10), 2373–2391.
- 622 Jain, A. K. (2010). Data clustering: 50 years beyond k-means. *Pattern recognition letters*, **31**(8), 651–666.
- 623 Jain, A. K., Murty, M. N., and Flynn, P. J. (1999). Data clustering: a review. *ACM computing surveys*
 624 (*CSUR*), **31**(3), 264–323.
- 625 Kass, R. E. and Raftery, A. E. (1995). Bayes factors. *Journal of the american statistical association*, **90**(430),
 626 773–795.
- 627 Kaufman, L. and Rousseeuw, P. J. (2009). *Finding groups in data: an introduction to cluster analysis*,
 628 volume 344. John Wiley & Sons.
- 629 Lloyd, S. (1982). Least squares quantization in pcm. *IEEE Transactions on Information Theory*, **28**(2),
 630 129–137.
- 631 Ng, A. Y., Jordan, M. I., Weiss, Y., et al. (2002). On spectral clustering: Analysis and an algorithm.
 632 *Advances in neural information processing systems*, **2**, 849–856.
- 633 Park, H.-S. and Jun, C.-H. (2009). A simple and fast algorithm for k-medoids clustering. *36*, 3336–3341.
- 634 Razavi Zadegan, S. M., Mirzaie, M., and Sadoughi, F. (2013). Ranked k-medoids: A fast and accurate
 635 rank-based partitioning algorithm for clustering large datasets. *Know.-Based Syst.*, **39**, 133–143.
- 636 Strehl, A. and Ghosh, J. (2003). Cluster ensembles — a knowledge reuse framework for combining multiple
 637 partitions. *J. Mach. Learn. Res.*, **3**, 583–617.
- 638 Tebaldi, C. and Knutti, R. (2007). The use of the multi-model ensemble in probabilistic climate projections.
 639 *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering*
 640 *Sciences*, **365**(1857), 2053–2075.
- 641 von Luxburg, U. (2007). A tutorial on spectral clustering. *Statistics and Computing*, **17**(4), 395–416.