# Supporting Information for "Multiscale Spatial Patterns in Giant Dike Swarms Identified through Objective Feature Extraction"

A. Kubo Hutchison[1], L. Karlstrom[1], T. Mittal[2]

[1]Department of Earth Sciences, 1272 University of Oregon, Eugene, Oregon 97403, USA

[2]Department of Geosciences, 503 Deike Building, University Park, PA 16802

**Contents of this file**

  1. Text S1 to S3

  2. Figures S1 to S4

**Additional Supporting Information (Files uploaded separately)**

  1. Captions for Datasets S1 to S4

**Introduction**

**Text S1.**

## 1.  Parallel Line Segments

For two parallel line segments with slope, $m$, and intercepts $b_1$ and $b_2$ the distance, $d_c$ between the two lines compared to the same difference in Hough space $dHT = |\rho_2 - \rho_1|$ regardless of origin. In the following, we show this mathematically.

In cartesian space the distance between the two parallel lines is

$$d = (b_1 - b_2)sin(\theta). \tag{1}$$

December 23, 2022, 4:28am

where

$$\theta = arctan(\frac{1}{m}) = arctan(\frac{x_a - x_b}{y_a - y_b}) \qquad (2)$$

where the subscripts a,b reflect the end points of the each respective segment.

In Hough space the distance between the two points is $|\rho_2 - \rho_1|$ where

$$\rho_1 = (x_1 - x_0)cos\theta + (y_1 - y_0)sin\theta. \qquad (3)$$

$$\rho_2 = (x_2 - x_0)cos\theta + (y_2 - y_0)sin\theta. \qquad (4)$$

with an origin of $(x_0, y_0)$. The intercept of the segments, which is the only difference between them can be written as:

$$b = y_a - m * x_a \qquad (5)$$

for $b_1$ and $b_2$ and each segment's respective endpoints. However, for $b'$ modified for changing the origin the equation is

$$b' = (y_a - y_0) - m * (x_a - x_0) \qquad (6)$$

Rewriting Equation 5 using rearranging Equations 2-4 we find:

$$y_a = \frac{\rho_1}{sin\theta} + \frac{x_a - x_0}{tan\theta} + y_0 \qquad (7)$$

Substituting Equation 6 and 2 into Equation 5 yields

$$b_1 = \frac{\rho_1}{sin\theta} + \frac{x_0}{tan\theta} + y_0 \qquad (8)$$

$$b_2 = \frac{\rho_2}{sin\theta} + \frac{x_0}{tan\theta} + y_0 \qquad (9)$$

December 23, 2022, 4:28am

but

$$b' = \frac{\rho}{sin\theta} \tag{10}$$

Finally substituting Equations 7 and 8 into 1 we prove that

$$d_C = \rho_2 - \rho_1 = d_{HT} \tag{11}$$

the distance between two parallel lines in Cartesian space is equal to distance between the two points in the Hough space. For truly parallel line segments, there is no $\rho$ distortion between Hough and Cartesian space and it is not dependent on angle or origin.

## 2. Subparallel Line Segments

For two subparallel line segments, in the same quadrant, with some small difference in angle $\theta_e$, what is the $\rho$ distortion and is it dependent on angle or origin choice? We analyze this case by taking two line segments (red and blue respectively in Fig S1) which are some small $\theta_e$ apart in angle with $\rho_2 > \rho_1$ .

With two different angles we find the perpendiculars for each respective line.

$$\rho_2 = (x_1 - x_0)cos\theta + (y_1 - y_0)sin\theta. \tag{12}$$

$$\rho_1 = (x_2 - x_0)cos(\theta + \theta_e) + (y_2 - y_0)sin(\theta + \theta_e) \tag{13}$$

In subparallel lines, there are two distances of interest as defined by the point formed by the intersection of the perpendicular and line formed by the line segment. We define two distances $d_1$ and $d_2$ which are analogous to $d$ above and the ratio of interest $\frac{d_1}{d_2}$.

Using the right triangles shown in the inset of Figure S1, we can find that

$$d_1 = \frac{\rho_2}{cos(\theta_e)} - \rho_1 \tag{14}$$

and

$$d_2 = \rho_2 - \frac{\rho_1}{cos(\theta_e)} \tag{15}$$

which yields the ratio

$$\frac{d_1}{d_2} = \frac{\rho_2 - \rho_1 cos\theta_e}{\rho_2 cos\theta_e - \rho_1} \tag{16}$$

which as $\theta_e \to 0$ , $\frac{d_1}{d_2} \to 1$ and the distance is equal to $\rho_2 - \rho_1$.

Plugging in 12 and 13 into 14 and 15 respectively and expanding the trigonometry identity using $(\theta + \theta_e)$, we find that the origin choice does affect the $d_1$ and $d_2$ values but disappears as $\theta_e \to 0$.

$$d_1 = (x_2 - x_0)cos\theta + (y_2 - y_0)sin\theta - cos\theta_e \Big( (x_1 - x_0)[cos\theta cos\theta_e - sin\theta sin\theta_e]+ \\ (y_1 - y_0)[sin\theta cos\theta_e + cos\theta sin\theta_e] \Big) \tag{17}$$

$$d_2 = (x_2 - x_1)cos\theta cos\theta_e + (y_2 - y_1)sin\theta cos\theta_e+ \\ cos\theta_e((x_1 - x_0)sin\theta + (y_1 - y_0)cos\theta) \tag{18}$$

which both simplify to $(x_2 - x_1)cos\theta + (y_2 - y_1)sin\theta$ when $\theta_e \to 0$ and the origin terms cancel out.

The distance in Hough space $d_{HT}$ is singular and equal to :

$$d_{HT} = \sqrt{(\theta_e)^2 + (\rho_2 - \rho_1)^2} \tag{19}$$

where

$$\rho_2 - \rho_1 = (x_2 - x_0)cos\theta + (y_2 - y_0)sin\theta - ((x_1 - x_0)[cos\theta cos\theta_e - sin\theta sin\theta_e] \\ + (y_1 - y_0)[sin\theta cos\theta_e + cos\theta sin\theta_e]) \tag{20}$$

This shows that for almost subparallel line segments, there is a distortion in $\rho$ space of the Hough Transform which is dependent on both angle and origin location. As the

73 angle difference betwene the two lines decreases, the dependence on angle and origin

74 decreases.The ratio of $\frac{d_1}{d_2}$ can measure the amount of distortion in between two segments.

75 When $\frac{d_1}{d_2}$ is near one the distortion is small and the difference in angle is small.

76     For the limit when $\theta_e$ is small, we can simplify Equation 20 to

$$d_{HT} = \sqrt{d_p + [(x_1 - x_0)sin\theta - (y_1 - y_0)cos\theta] * \theta_e} \tag{21}$$

## 3. Crossing Negative to Positive Angles

78     In the HT space, lines with $-90°$ and $90°$ have the same horizontal orientation (E-

79 W from a map point of view). Thus we would wish to cluster some dikes which have

80 this orientation. However, there is additional distortion which occurs when we cluster

81 from negative to positive angles. We can illustrate this with an example of two points

82 in HT space $p_1 = [-90°, 100]$ and $p_1 = [90°, 100]$. However, in Cartesian space, when

83 we calculate the intercept between these two lines using Equation 10 the intercepts are

84 far apart $b_1 = -100$ and $b_2 = 100$ leading to an overly wide cluster violating our desired

85 constraints for high aspect clusters. In this case, a third point $p_3 = [90°, -100]$ would

86 have the same orientation as $p_1$ since $b_3 = -100$.

87     To solve this issue in HT space clustering we simply rotate the dataset so that the

88 median angle is centered on 20 degrees which aims to minimize the amount of clusters

89 which would need to cross $-90°$ and $90°$ and $0°$.

90 **Text S2.**

## 4. Clustering Sensitivity Analysis

91     To investigate dike cluster lengths and act as data reduction we apply a clustering al-

92 gorithm to the datasets transformed into Hough Space. We chose to use Agglomerative

Hierarchical clustering (Everitt, 1980) which has two parameters which effect the algorithm: first, the linkage criterion of which we chose complete linkage; and secondly the distance parameter over which clusters will not be merged. In the two dimensional space of the Hough Transform there are two different scales over which this distance parameter calculates Euclidean distance, the angle which varies from $-90$ to $90$ and $\rho$ value which varies over several orders of magnitude from positive to negative values and is measured in kilometers. To account for this, we scale our datasets by the $\theta$ and $\rho$ cutoff values. For all datasets, we use the same $\theta$ threshold of $2°$ while the $\rho$ threshold varies for each dataset. For the $\rho$ threshold, we have chosen the smallest scale length available in the data which is the mean dike segment length. Our choices for these parameters are explained in the main text.

To examine the effect of choosing different parameters and the robustness of our clustering analysis we performed a sensitivity analysis on the Chief Joseph Dike Swarm dataset testing the effects of changing three clustering parameters: $\rho$ threshold, $\theta$ threshold, and linkage criterion (Figure S2). These tests were performed with the Scipy Hierarchical clustering algorithms (Virtanen et al., 2020). When changing the $\rho$ threshold (Figure S2 a,d,g), we applied a $\theta$ threshold of $2°$ using complete linkage. When changing the $\theta$ threshold (Figure S2 b,e,h), we applied a $\rho$ threshold of 400 m using complete linkage. When changing the linkage (Figure S2 c,f,i), we applied a $\rho$ threshold of 400 m and a $\theta$ threshold of $2°$. For reference, the $\rho$ threshold of 400 m is what was used in the main text. Each result was then put through the same post-processing steps and compared. We applied the same filtering step described in the Methods section of the main text.

We see that increasing it significantly to 1000 m has little affect on median cluster sizes although it slightly increases the outlier clusters. Increasing it again to 5000 m, a 10X, does introduce significantly larger clusters and raises median cluster significantly for the filtered datasets. As would be expected it does raise the cluster width however the ranges ranges seen in the 400 m and 1000 m cut off are similar.

We tested three linkage types; first complete linkage in which the distance is determined by the furthest most data points in a cluster; second, single linkage in distance is determined by the closest objects in a cluster only; finally, average or unweighted pair-group method aproach (UPGMA) in which the distance between two clusters is the average distance between all objects in those clusters. Single linkage is the most effected by single data points while complete linkage creates more compact clusters. There are a few extremely large clusters using single linkage ($< 600$ segments) this is likely due to the "chaining" phenomenon in which single segments cluster and this creates more "outlier" clusters however the ranges for dike cluster length and width remain consistent. This behavior can be desirable for some research applications but not for the issue of linking segments with like orientation. After filtering for cluster size, the single linkage also reveals the largest deviations from complete and average linkage likely due to the presence of chaining clusters which are also more likely to be $> 3$ segments. Average linkage behaved similarly to complete linkage but showed slightly larger, longer, and wider clusters.

Changing the $\theta$ threshold has little effect on median cluster size until it reaches $10°$ which does show significantly larger clusters in the filtered dataset and higher extremes in size. Changing angle from $1°$ to $2°$ effects the filtered clusters but has little effect on the overall dataset. Increasing it again from $2°$ to $4°$ increases the median size from 2 to

138  3. When it comes to cluster length, changing the $\theta$ threshold does not effect the range of

139  cluster lengths seen in the cluster data base although it does effect the range seen in the

140  filtered database likely due to the changing median size. This behavior would be expected

141  since the Hough transform has no information about how segments are arranged relative

142  to one another and changes are likely due to increasing cluster sizes. It does however

143  increase cluster width linearly. This is described above in Text S1 since increasing the

144  $\theta$ threshold increases the values of $\theta_e$ and increases possible widths between subparallel

145  segments. Overall, changing $\theta$ threshold has relatively larger effects on the cluster length

146  then changing $\rho$ threshold.

147  **Text S3.**

## 5.   Identification of Radial Swarms

148  Due to the distortion described above and the nature of the Hough Transform as a set

149  of dikes gets further from the choosen origin of the Hough Transform it will appear more

150  and more like a radial swarm despite being a random array of orientations. This is a

151  function of the size of the radial swarm, it's range and standard deviation of $\rho$ rather than

152  it's Cartesian length. To test this we created several completely random dike swarms with

153  no radial pattern of different scales (1, 10, 100 km) and placed them at different distances

154  from the origin then fit the radial swarm equation (Eq. 9 in the main text) to the datasets.

155  To evaluate how well it is fit, erroneously, by a radial swarm we look at the goodness of fit

156  ($R_{sq}$) value (Figure S3). We find that when the distance from the HT origin is scaled by

157  the standard deviation in $\rho$ ($\mu$), radial swarm start to appear erroneously at approximately

158  $2 - 2.5\mu$. Assuming the dike swarms are normally distributed in $\rho$ this indicates that the

159  majority of dikes ($> 95\%$) which will fall at distances less than $\pm 2\mu$ can be correctly

identified as radial if that is there pattern and few swarms will be incorrectly identified as radial. Additionally, using this knowledge one can move the HT center accordingly to account for features that are in the far ranges of $\rho$ if necessary and run the radial identification as many times as necessary then compile it into one Cartesian dataset.

**Data Set S1.** Linked dike clusters for the Columbia River Flood Basalt group including the four identified subswarms: Chief Joseph, Monument, Ice Harbor, and Steens as compiled in Morriss, Karlstrom, Nasholds, and Wolff (2020). This dataset was produced using the Agglomerative Clustering algorithms using the parameters set in Table 1. This dataset is in the format of a CSV file but can be read into GIS programs using Well Known Text (WKT) linestring. This dataset uses the a UTM Zone 11N projection (EPSG:26911). The file includes the start and end points of the average line in the cluster and it's mid points, cluster length and width (Xstart, Xend, Xmid, Ymid, in meters and UTM coordinates, Dike Cluster Width or R_Width, Dike Cluster Length or R_Length all in meters); calculated average $\rho$ and $\theta$ for the Hough Transform $\rho$ units measured in meters, $\theta$ units measured in degrees, unless otherwise stated); the origin used for the Hough Transform which is different for each subswarm ($xc$,$yc$, meters in UTM coordinates); average slope and intercept (AvgSlope, AvgIntercept meters); range and standard deviation for $\rho$ and $\theta$ for all objects in the cluster ($\rho$ units measured in meters, $\theta$ units measured in degrees); cluster size (Size); sum of segment lengths in a cluster (SegmentLSum, meters); whether the cluster crosses between negative and positive values (ClusterCrossesZero, boolean); overlap as calculated in the main text where the length of overlap is normalized by the sum of segment lengths in a cluster; maximum number of overlapping segments (nOverlapingSegments); twist angle which is the difference in angle betweeen the average cluster

line and the average line formed by cluster midpoints (EnEchelonAngleDiff, degrees); the

p-value for the midpoint line fit of the segments where $p < 0.05$ is considered to be a sig-

nificant fit (EEPValue); the maximum, median, and minimum segment nearest neighbors

distances in the cluster which is calculated using the cartesian midpoints of each segment

and normalized by the Cluster Length (MaxSegNNDist, MedianSegNNDist, MinSegN-

NDist); characterization of each cluster as filtered or not, filtered clusters are of size greater

than 3 and have a MaxSegNNDist of less than 0.5 (TrustFilter, boolean); the date edited

(Date_Changed), and the clustering parameters used for each cluster (Rho_Threshold in

meters, Theta_Threshold in degrees) and a unique identification calculated based on the

start and endpoints (ClusterHash).

**Data Set S2.** Linked dike clusters for the Deccan Traps including the four identified

subswarms: Saurashtra, Narmada-Tapi, Central and Coastal. Due to their overlap Cen-

tral and Coastal Swarms have been combined in this dataset into the Central Swarm.

This dataset was produced using the Agglomerative Clustering algorithms using the pa-

rameters set in Table 1. This dataset is in the format of a CSV file but can be read into

GIS programs using Well Known Text (WKT) linestring. This dataset uses the a WGS

84 projection (EPSG:3857). The file includes the start and end points of the average

line in the cluster and it's mid points, cluster length and width (Xstart, Xend, Xmid,

Ymid, in meters and UTM coordinates, Dike Cluster Width or R_Width, Dike Cluster

Length or R_Length all in meters); calculated average $\rho$ and $\theta$ for the Hough Transform

$\rho$ units measured in meters, $\theta$ units measured in degrees, unless otherwise stated); the

origin used for the Hough Transform which is different for each subswarm ($xc,yc$, me-

ters in UTM coordinates); average slope and intercept (AvgSlope, AvgIntercept meters);

range and standard deviation for $\rho$ and $\theta$ for all objects in the cluster ($\rho$ units measured in meters, $\theta$ units measured in degrees); cluster size (Size); sum of segment lengths in a cluster (SegmentLSum, meters); whether the cluster crosses between negative and positive values (ClusterCrossesZero, boolean); overlap as calculated in the main text where the length of overlap is normalized by the sum of segment lengths in a cluster; maximum number of overlapping segments (nOverlapingSegments); twist angle which is the difference in angle betweeen the average cluster line and the average line formed by cluster midpoints (EnEchelonAngleDiff, degrees); the p-value for the midpoint line fit of the segments where $p < 0.05$ is considered to be a significant fit (EEPValue); the maximum, median, and minimum segment nearest neighbors distances in the cluster which is calculated using the cartesian midpoints of each segment and normalized by the Cluster Length (MaxSegNNDist, MedianSegNNDist, MinSegNNDist); characterization of each cluster as filtered or not, filtered clusters are of size greater than 3 and have a MaxSegNNDist of less than 0.5 (TrustFilter, boolean); the date edited (Date_Changed), and the clustering parameters used for each cluster (Rho_Threshold in meters, Theta_Threshold in degrees) and a unique identification calculated based on the start and endpoints (ClusterHash).

**Data Set S3.** Dike segment data for Spanish Peaks and Dike Mountain located in the Rio Grande Rift of Colorado. This dataset was digitized using QGIS based on the map by (Johnson, 1961). This dataset is in the format of a CSV file but can be read into GIS programs using Well Known Text (WKT) linestring. This dataset uses the a UTM Zone13N projection (EPSG:32613). The file includes the start, end points, and midpoints of the dikes; segment length; calculated $\rho$ and $\theta$ for the Hough Transform; the origin used for the Hough Transform which is different for each subswarm ($xc,yc$); dike rock type

229 if known; and a unique identification calculated based on the start and endpoints. This

230 dataset has been preprocessed to remove curving dikes and is the data set used to produce

231 later products (Data set S6).

232 **Data Set S4.** Linked dike clusters for the Spanish Peaks and Dike Mountain. This

233 dataset was produced using the Agglomerative Clustering algorithms using the parame-

234 ters set in Table 1. This dataset is in the format of a CSV file but can be read into GIS

235 programs using Well Known Text (WKT) linestring. This dataset uses the a UTM Zone

236 13N projection (EPSG:32613). The file includes the start and end points of the average

237 line in the cluster and it's mid points, cluster length and width (Xstart, Xend, Xmid,

238 Ymid, in meters and UTM coordinates, Dike Cluster Width or R_Width, Dike Cluster

239 Length or R_Length all in meters); calculated average $\rho$ and $\theta$ for the Hough Transform

240 $\rho$ units measured in meters, $\theta$ units measured in degrees, unless otherwise stated); the

241 origin used for the Hough Transform which is different for each subswarm ($xc$,$yc$, me-

242 ters in UTM coordinates); average slope and intercept (AvgSlope, AvgIntercept meters);

243 range and standard deviation for $\rho$ and $\theta$ for all objects in the cluster ($\rho$ units measured

244 in meters, $\theta$ units measured in degrees); cluster size (Size); sum of segment lengths in a

245 cluster (SegmentLSum, meters); whether the cluster crosses between negative and posi-

246 tive values (ClusterCrossesZero, boolean); overlap as calculated in the main text where

247 the length of overlap is normalized by the sum of segment lengths in a cluster; maximum

248 number of overlapping segments (nOverlapingSegments); twist angle which is the differ-

249 ence in angle betweeen the average cluster line and the average line formed by cluster

250 midpoints (EnEchelonAngleDiff, degrees); the p-value for the midpoint line fit of the seg-

251 ments where $p < 0.05$ is considered to be a significant fit (EEPValue); the maximum,

median, and minimum segment nearest neighbors distances in the cluster which is calculated using the cartesian midpoints of each segment and normalized by the Cluster Length (MaxSegNNDist, MedianSegNNDist, MinSegNNDist); characterization of each cluster as filtered or not, filtered clusters are of size greater than 3 and have a MaxSegNNDist of less than 0.5 (TrustFilter, boolean); the date edited (Date_Changed), and the clustering parameters used for each cluster (Rho_Threshold in meters, Theta_Threshold in degrees) and a unique identification calculated based on the start and endpoints (ClusterHash).

**Figure S1.**    Figure S1 : **A.** shows two line in Cartesian space that are subparallel and the distances ($d_1$ and $d_2$) between them relative to the Hough transform $\rho$ distances. **B.** shows the same view with the origin at the top and intersections labeled $A, B, C, D$.
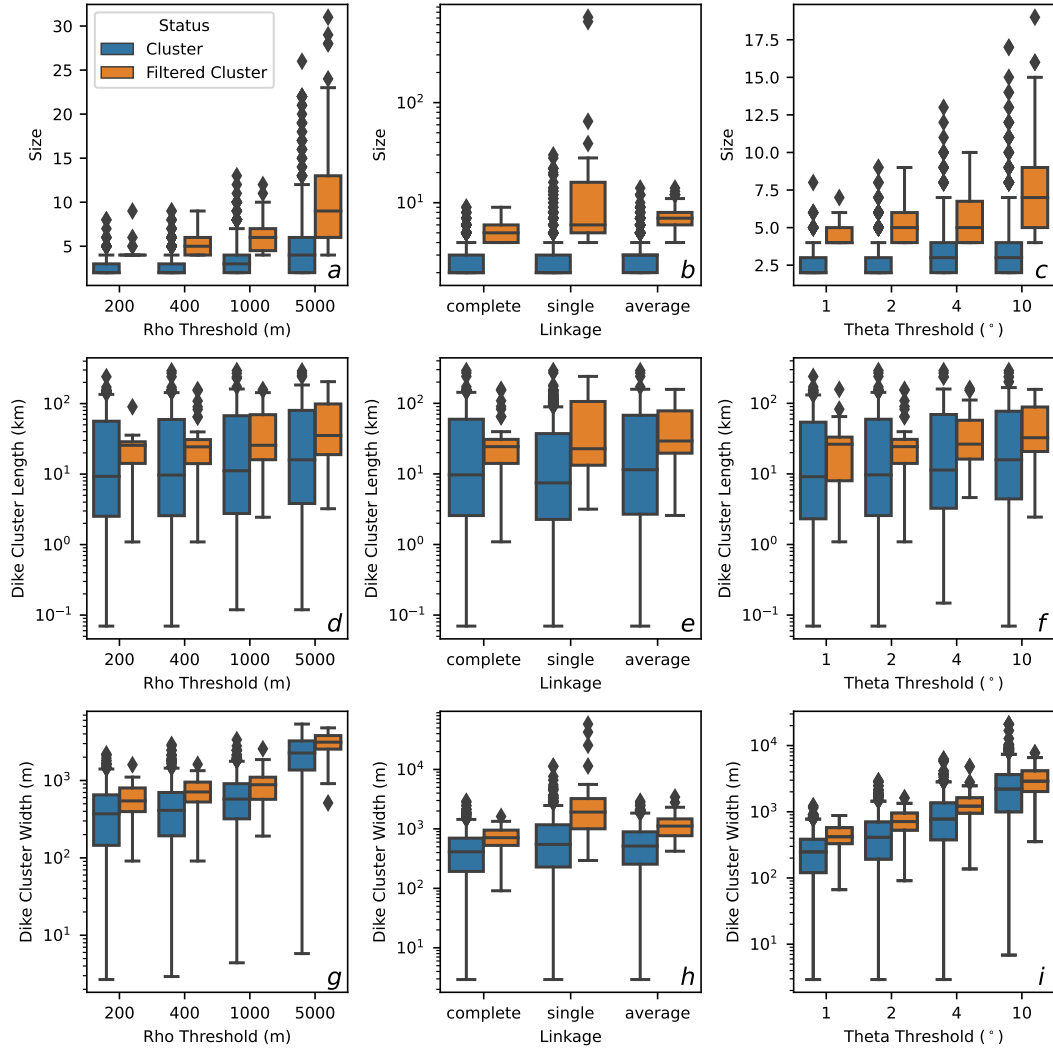
**Figure S2.** Figure S2 : **A.-C.** show a box and whisker plot of cluster size for the whole linked dataset and filtered clusters when changing $\rho$ threshold, linkage, and $\theta$ threshold respectively. **D.-F.** show a box and whisker plots of cluster length in log scale. **G.-I.** show a box and whisker plots of cluster width in log scale.
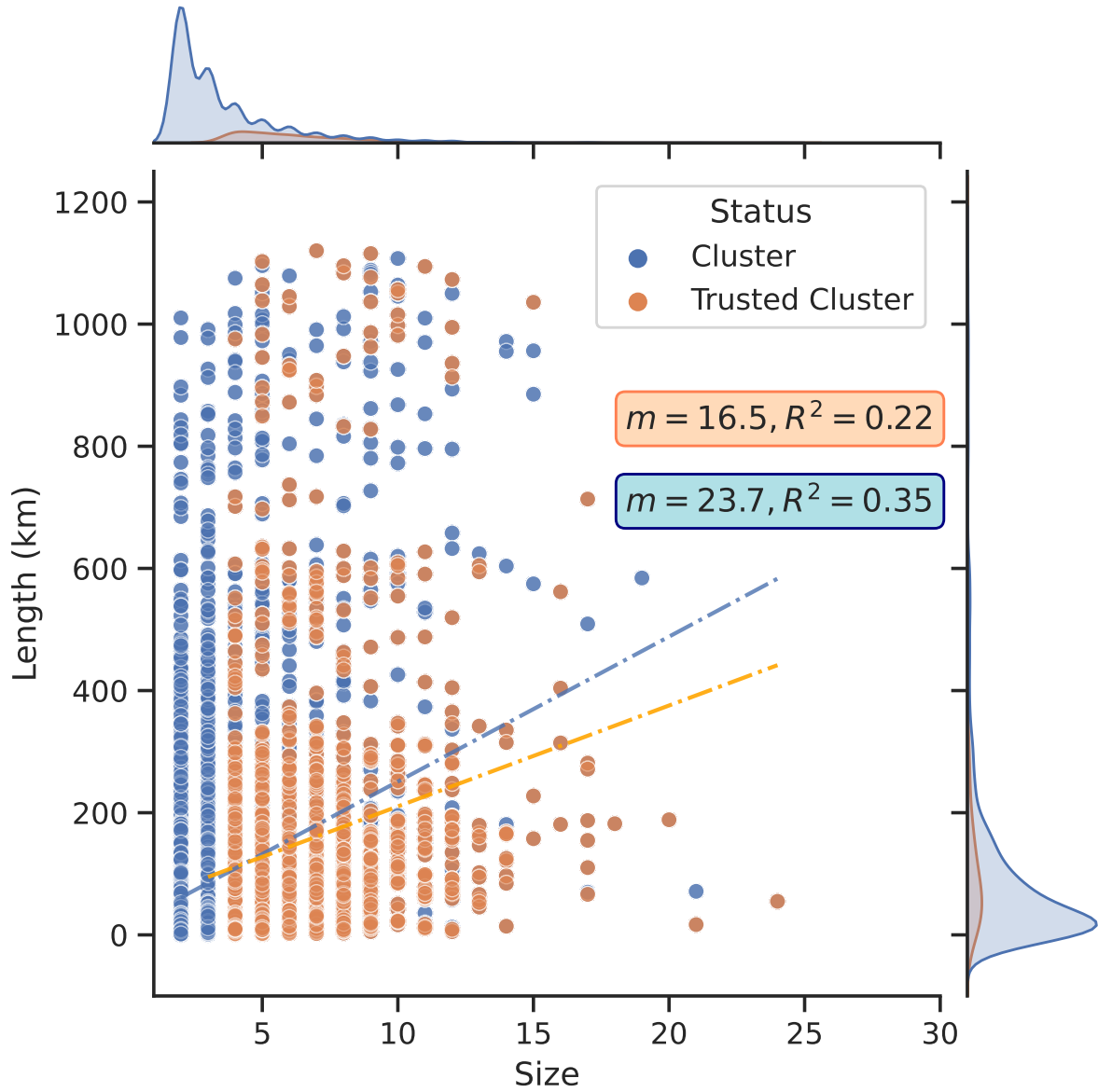
**Figure S3.**    Cluster length as a function of cluster size for all dike segment datasets with all clusters in blue and filtered clusters in orange. Overlaid are two regression curves (dotted dashed lines) and their respective fits. Overall neither the full cluster database nor the filtered database show strong relationship between cluster size and cluster length. This is to be expected since clustering is done in the Hough transform space irrespective of Cartesian location. Very long dike clusters ($> 200km$) are seen with cluster sizes of 2-18 although clusters of over 5 are relatively rare.
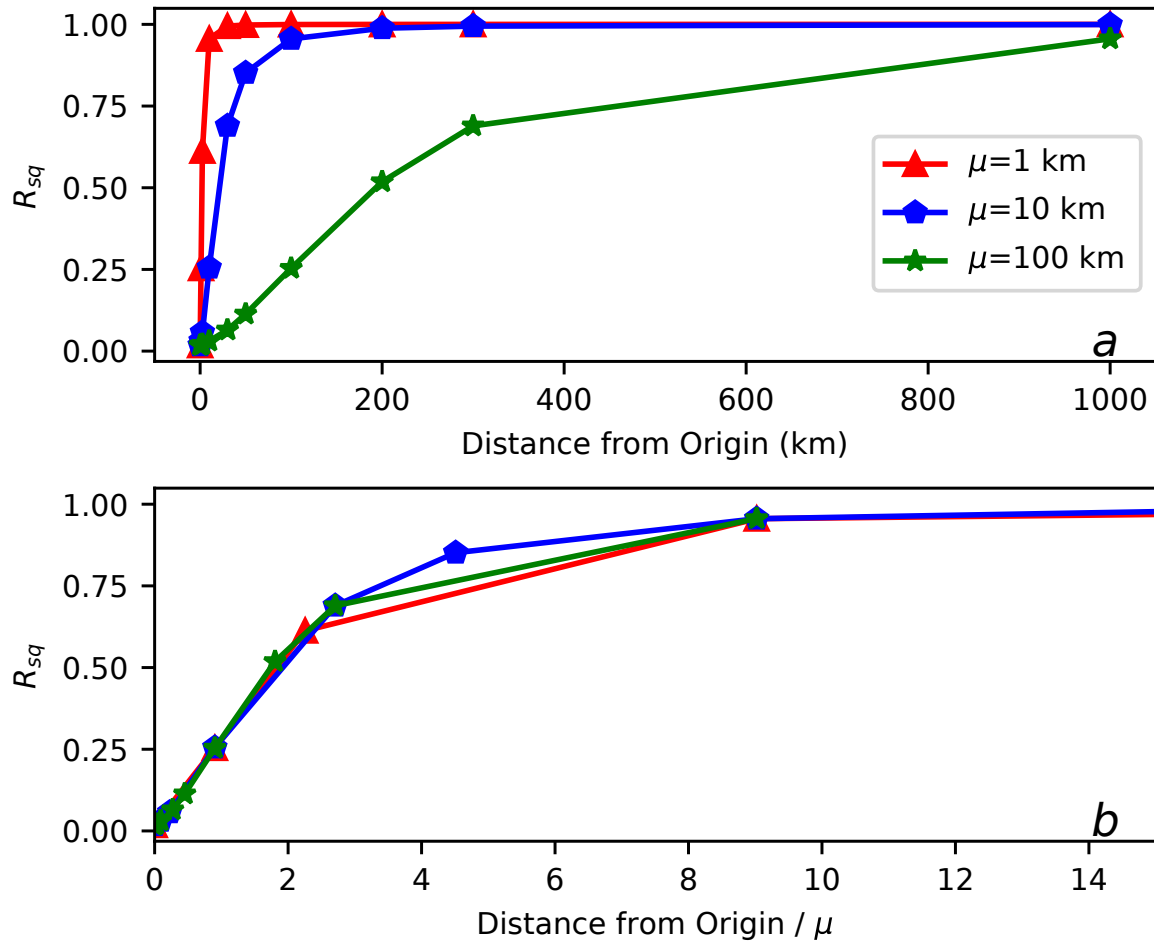
**Figure S4.** **A** For three synthetic randomly oriented swarms of different scales, the distance from the Hough Transform Origin is plotted against the goodness of fit for a radial swarm. **B** shows the same but with distance normalized by the scale of the swarm ($\mu$), the standard deviation of $\rho$.

# References

Everitt, B. (1980). *Cluster analysis* (Vol. 14) (No. 1). doi: 10.1007/BF00154794

Johnson, R. (1961). Patterns and Origin of Radial Dike Swarms Associated with West Spanish Peak and Dike Mountain, South-Central Colorado. *Geological Society of America Bulletin*, *72*(April), 579–590.

Morriss, M. C., Karlstrom, L., Nasholds, M. W., & Wolff, J. A. (2020). The chief Joseph dike swarm of the Columbia river flood basalts, and the legacy data set of William H. Taubeneck. *Geosphere*, *16*(4), 1793–1817. doi: 10.1130/GES02173.1

Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., ... Contributors, S. . (2020). SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, *17*, 261–272.