

Spatial consistency for multiple source direction-of-arrival estimation and source counting

Sina Hafezi,^{1, a)} Alastair H. Moore,^{1, b)} and Patrick A. Naylor^{1, c)}

*Department of Electrical and Electronic Engineering, Imperial College London,
SW7 2AZ, UK*

(Ω Dated: 24 October 2019)

1 A conventional approach to wideband Multi-Source (MS) Direction-of-Arrival (DOA)
2 estimation is to perform Single Source (SS) DOA estimation in Time-Frequency (TF)
3 bins for which a SS assumption is valid. The typical SS-validity confidence metrics
4 analyse the validity of the SS assumption over a fixed-size TF region local to the
5 TF bin. The performance of such methods degrades as the number of simultane-
6 ously active sources increases due to the associated decrease in the size of the TF
7 regions where the SS assumption is valid. A SS-validity confidence metric is pro-
8 posed that exploits a dynamic MS assumption over relatively larger TF regions. The
9 proposed metric first clusters the initial DOA estimates (one per TF bin) and then
10 uses the members' spatial consistency as well as its cluster's spread to weight each
11 TF bin. Distance-based and density-based clustering are employed as two alternative
12 approaches for clustering DOAs. A noise-robust density-based clustering is also used
13 in an evolutionary framework to propose a method for source counting and source
14 direction estimation. The evaluation results based on simulations and also with real
15 recordings show that the proposed weighting strategy significantly improves the ac-
16 curacy of source counting and MS DOA estimation compared to the state-of-the-art.

a) s.hafezi14@imperial.ac.uk

b) alastair.h.moore@imperial.ac.uk

c) p.naylor@imperial.ac.uk

17 I. INTRODUCTION

18 Multi-Source (MS) Direction-of-Arrival (DOA) estimation is required for acoustic source
19 separation/localization/tracking, spatial filtering, environment mapping, dereverberation
20 and speech enhancement. It addresses the often-occurring case in real-world scenarios where
21 two or more sources are active simultaneously. As such it can be used in applications such as
22 hearing aids, robot audition, meeting diarization and teleconferencing. The main challenges
23 for MS DOA algorithms include reverberation, sensor or environmental noise as well as the
24 presence of an unknown number of simultaneously active sources¹. In this work we address
25 MS DOA estimation in a reverberant environment using microphone array with arbitrary
26 geometry and configuration for an unknown number of simultaneously active speech sources
27 where the number of these active sources changes over time.

28 Several existing approaches to MS DOA estimation for speech sources use W-Disjoint
29 Orthogonality (WDO)², assuming sparseness of speech in the Time-Frequency (TF) domain,
30 in combination with subspace decomposition to decompose the noisy observation covariance
31 matrix into signal and noise subspaces^{3,40-42}. Such methods often follow three stages: (1)
32 Single Source (SS) TF bin detection in which the TF bins predominantly containing a
33 single source are detected; (2) SS DOA estimation where a DOA estimator based on the SS
34 assumption is applied only at the detected SS bins; (3) Multiple source direction estimation
35 using the set of temporal narrowband DOA estimates.

36 In a SS bin, the observation covariance matrix formed from the microphone array signals
37 is expected to have unit rank. In real-world scenarios, DOA estimation has to be performed

38 in reverberation that is characterized by the combination of direct-path propagation and
39 reflections⁴. In such scenarios, SS dominance with unit rank covariance matrix rarely occurs
40 at a TF bin and therefore some form of SS-validity confidence metric is used to detect the
41 more reliable SS bins for use in DOA estimation. Methods such as the coherence test⁵, SS
42 Zone (SSZ) detection⁶ or Direct Path Dominance (DPD) test⁷ assume the validity of SS
43 assumption over a local TF region in the vicinity of a TF bin of interest and each method
44 defines a specific SS-validity confidence metric. Having obtained the SS validity measure
45 at each bin, two alternative approaches can be used for the selection of reliable bins. The
46 methods in^{6,7} identify the SS bins based on a comparison between the SS validity measures
47 and a fixed user-defined threshold whereas the method in⁸ selects a user-defined percentage
48 of the TF bins with the strongest SS validity measures. In⁵, SS bins are detected using
49 the rank of the correlation matrix at each TF bin. Due to averaging across only local time
50 frames and lack of subspace decomposition in the selection of SS bins, that approach is
51 most effective only for MS DOA estimation in an anechoic environment. In⁶, the average of
52 pairwise correlation coefficients between adjacent sensors is used as a SS validity confidence
53 metric, where the correlation averaging is performed only across the local frequencies of
54 each time frame. It does not use subspace decomposition and is therefore prone to noise and
55 multiple coherent sources. In DPD⁷, Singular Value Decomposition (SVD) is employed and
56 the Singular Value Ratio (SVR), defined as the ratio of the largest to second largest singular
57 values, of the signals' covariance matrix is used as the SS-validity confidence metric. DPD
58 performs the covariance averaging over adjacent frequencies and time-frames. The latter
59 property, along with the use of subspace decomposition makes DPD robust to reverberation

60 as it aims to find TF bins with not just a dominant SS but also a dominant direct path,
61 ignoring bins containing significant reverberation.

62 As the number of simultaneously active sources increases, the performance of the pre-
63 viously mentioned methods degrades³⁸, although presence of the dominant SS still occurs.
64 This is because the WDO assumption is valid in fewer TF bins and in smaller TF regions
65 as the number of simultaneously active sources increases as shown in Section II.

66 Figure 2 shows an overview of the MS DOA estimation system proposed in this paper.
67 The novelties in this work are: (1) the use of density-based clustering in the context of
68 acoustic DOA estimation, as used in DOAs clustering and source counting units in Fig.
69 2 and (2) a novel SS-validity confidence metric for weighting of initial DOA estimates, as
70 used in DOAs weighting unit in Fig. 2. The proposed Multi-Source Estimation Consistency
71 (MSEC) metric is based on a dynamic MS assumption, as opposed to the SS assumption in
72 conventional approaches. MSEC uses a consistently large TF region where the number of
73 simultaneously active sources within the region is autonomously estimated.

74 This paper is structured as follows. Section II demonstrates the problem when the number
75 of simultaneously active sources increases. Section III reviews two alternative distance- and
76 density-based clustering approaches employed to estimate the number of active sources. It
77 then describes a novel SS-validity confidence metric as well as an autonomous source counting
78 method, an earlier version of which is discussed in⁹. Section IV evaluates the performance of
79 the proposed metric against the state-of-the-art under various simulated scenarios. Finally
80 Section V illustrates the performance and accuracy of the evaluated methods using signals
81 recorded in a real room.

82 II. PROBLEM ANALYSIS

83 Consider a reverberant environment containing N_s simultaneously active speech sources
 84 with uniform angular spacing $\Delta\phi$ at 1 m distance from a microphone array. Each source rep-
 85 resents a different speaker speaking different utterances. The received signal at a microphone
 86 in the Short-time Fourier Transform (STFT) domain is

$$X(k, \tau) = \sum_{n=1}^{N_s} (A_n(k, \tau) + \sum_{j=1}^{\infty} R_{n,j}(k, \tau)), \quad (1)$$

87 where A_n denotes the direct-path component from source n and $R_{n,j}$ is the component of
 88 reflection j from source n . The frequency, k , and time frame, τ , indices are subsequently
 89 omitted for notational simplicity.

90 Let Signal-to-Interference Ratio (SIR) at a TF bin be the ratio of the magnitude of the
 91 dominant direct path, $|A_b|$, and the magnitude of the rest of the signals from the mixture
 92 of N_s sources, $|X - A_b|$, which includes all other direct paths and reverberations excluding
 93 the dominant direct path where

$$b = \underset{n}{\operatorname{argmax}}(|A_n|), \quad (2)$$

94 is the index of the dominant direct path. Figure 1 shows in white the TF bins with SIR
 95 ≥ 10 dB in such a scenario for $N_s = \{2, 3, 4, 5\}$ and $\Delta\phi = 50^\circ$. It can be clearly seen that
 96 with the increase of N_s , the number of the bins and the size of the TF regions with valid
 97 WDO assumption decreases. For SS bin detectors based on a fixed-size analysis TF window,
 98 this leads to increasing failure of the SS assumption validity and consequent performance
 99 degradation for the SS bin detectors that rely on this assumption.

100 One solution¹⁰ to this problem is the use of a dynamic MS assumption over a fixed-size TF
101 region where the number of active sources within the processing TF region is autonomously
102 estimated. For such techniques, estimation of the optimum number of sources remains a
103 challenge. In¹⁰, the authors propose the use of the Akaike Information Criterion (AIC)¹¹ to
104 find the optimum number of eigenvectors spanning the signal space for the MS assumption.
105 Although this approach overcomes the problem of N_s estimation, it loses reliability with
106 noisy observations.

108 The use of temporal narrowband DOA estimation based on the SS assumption in a MS
109 scenario is expected to be relatively accurate at the TF bins containing one significantly
110 dominant direct-path component and inaccurate otherwise^{7,38,39}. The direction of the er-
111 ror in erroneous DOA estimates is determined by the relative phase and amplitude of the
112 impinging plane waves as shown in¹². Such variance of directional displacements in DOA
113 estimates at non-SS bins results in spatially inconsistent erroneous DOAs whereas, in prac-
114 tical scenarios, DOA estimates at SS bins are expected to have spatial consistency if the
115 sources are stationary or only slowly moving over time. In¹³ and¹⁴, the authors propose
116 the use of diffuseness of DOA estimates which is based on SS assumption and suffers from
117 the previously-stated problem of SS-based metrics as the number of sources increases. We
118 therefore investigate the use of spatial consistency of SS-based DOA estimates under the
119 MS assumption. We also investigate how to estimate the number of active sources over a
120 TF region using this approach, as well as the validity of the SS assumption at a TF bin.

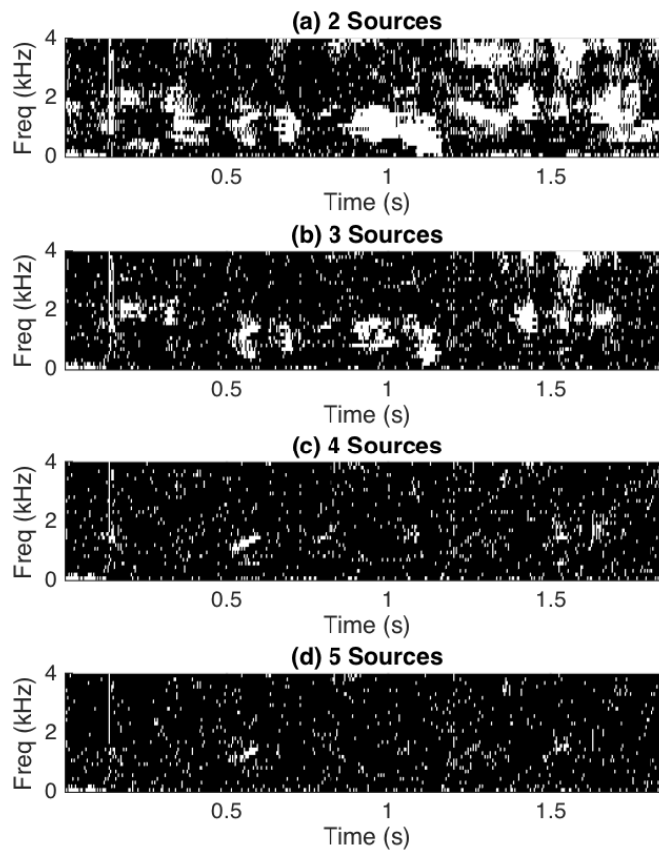


FIG. 1. Illustration showing in white the TF bins with $\text{SIR} \geq 10$ dB considering the signal as the dominant direct path and the interference as the reverberant signals mixture of (a) 2, (b) 3, (c) 4 and (d) 5 sources with $T_{60} = 0.4$ s.

121 III. PROPOSED METHOD

122 Assuming that the initial DOAs (one per TF bin) are provided by any chosen temporal
 123 narrowband SS DOA estimation procedure, a new SS validity confidence metric is proposed
 124 based on spatial consistency of initial DOA estimates and a dynamic MS assumption. Two
 125 alternative distance- and density-based clustering techniques for the dynamic MS assumption

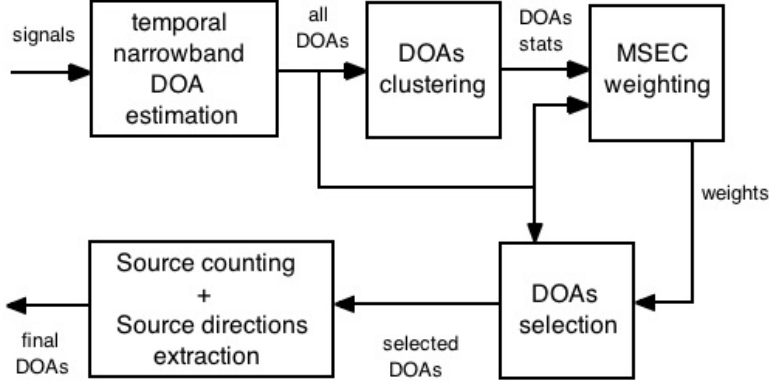


FIG. 2. Block diagram of the proposed system for MS DOA estimation. The MSEC weighting and Source counting blocks are specifically where the proposed methods contribute.

126 are introduced and discussed. The architecture of the proposed system is illustrated in Fig.
 127 2.

128 In order to increase the distinctness of densities between accurate and inaccurate initial
 129 DOAs for the purpose of robust estimation of the number of active sources, we consider all
 130 initial DOA estimates from the previous T frames. Therefore, at each frame τ , we consider
 131 the set of DOA estimates $U(\tau)$ including all initial DOAs from frame τ to $\tau - T$, defined as

$$U(\tau) = \{\hat{\mathbf{u}}(t, k) : \forall k, t \in \{\tau, \tau - 1, \dots, \tau - T\}\}, \tag{3}$$

132 where $\hat{\mathbf{u}}(t, k)$ is the estimated DOA unit vector at time frame t and frequency k and T is a
 133 fixed user-defined temporal window length.

134 To quantify spatial consistency of the multi-modal distribution of DOA estimates from
 135 $U(\tau)$, an adaptive distance-based clustering technique such as K-means and a density-based
 136 noise-robust clustering technique such as Density-based Spatial Clustering of Applications
 137 with Noise (DBSCAN)¹⁵ are used as two alternative approaches. Sections III A and III B re-

138 spectively present adaptive K-means and DBSCAN clusterings applied to DOA distribution.
 139 In the following, τ , t and k are omitted for brevity where unambiguous.

140 **A. Adaptive K-means Clustering**

141 To find the optimum number of clusters, the AIC is calculated as¹¹

$$\text{AIC} = -2Q + 2v, \quad (4)$$

142 in which $-Q$, the negative maximum log-likelihood of the data, represents a measure of
 143 distortion and v , the number of parameters of the model, represents a measure of model
 144 complexity.

145 For K-means with a given K , the first term in (4) is replaced with Residual Sum of
 146 Squared (RSS) of the clustering giving¹⁶

$$\text{AIC}(K) = \text{RSS}(K) + 2JK, \quad (5)$$

147 where $\text{RSS}(\cdot)$ is the sum of squared angular distances of each member to its cluster centroid
 148 and J denotes the number of dimensions of the centroid which leads to JK parameters
 149 for K clusters. Note that with the increase of K , $\text{RSS}(K)$ decreases while $2JK$ increases,
 150 which makes $\text{AIC}(K)$ a penalty factor for a given model where its minimum gives the best
 151 clustering with the minimum number of clusters.

152 Having performed K-means for $K = \{1, \dots, K_{max}\}$ on the set of DOAs U with random
 153 initializations, using (5), the optimum number of clusters, K_c , is chosen as

$$K_c = \arg \min_K [\text{AIC}(K)]. \quad (6)$$

154 **B. DBSCAN Clustering**

155 Unlike distance-based clustering techniques, density-based DBSCAN clustering does not
 156 consider the number of clusters to be known *a priori* but instead is based on a user-defined
 157 minimum density for a cluster. Therefore DBSCAN considers an assumption on the density
 158 of clusters rather than the number of clusters, which makes it robust to noise and suitable
 159 for autonomous cluster counting.

160 The terms used in DBSCAN clustering are defined as follows¹⁵.

161 **1. Neighbourhood DOAs**

162 The set of neighbourhood DOAs for a DOA estimate $\hat{\mathbf{p}}$ is defined as

$$N_\varepsilon(\hat{\mathbf{p}}) = \{\hat{\mathbf{q}} \in U \mid \angle(\hat{\mathbf{p}}, \hat{\mathbf{q}}) \leq \varepsilon\}, \quad (7)$$

163 where $\angle(\hat{\mathbf{p}}, \hat{\mathbf{q}})$ is the angular separation (in degrees) between two DOA estimates $\hat{\mathbf{p}}$ and $\hat{\mathbf{q}}$
 164 and ε is chosen to define the angular extent of the neighbourhood in degrees.

165 **2. Density**

166 The density at a DOA estimate $\hat{\mathbf{p}}$ is defined as the number of DOA estimates (including
 167 $\hat{\mathbf{p}}$ itself) within its neighbourhood $|N_\varepsilon(\hat{\mathbf{p}})|$, where $|\cdot|$ indicates cardinality.

168 **3. Threshold density**

169 The threshold density denoted as MinPts is the minimum density for a potential cluster.

170 **4. *Directly density-reachable***

171 A DOA estimate $\hat{\mathbf{p}}$ is directly density-reachable from another DOA estimate $\hat{\mathbf{q}}$ if

- 172 • $\hat{\mathbf{p}} \in N_\varepsilon(\hat{\mathbf{q}})$ and
- 173 • $|N_\varepsilon(\hat{\mathbf{q}})| \geq \text{MinPts}$ (core point condition).

174 **5. *Density-reachable***

175 A DOA estimate $\hat{\mathbf{p}}$ is density-reachable from another DOA estimate $\hat{\mathbf{q}}$ if there is a chain

176 of DOA estimates $\{\hat{\mathbf{p}}_i\}_{i=1}^L$, where $\hat{\mathbf{p}}_1 = \hat{\mathbf{q}}$ and $\hat{\mathbf{p}}_L = \hat{\mathbf{p}}$, such that $\hat{\mathbf{p}}_{i+1}$ is directly density-

177 reachable from $\hat{\mathbf{p}}_i$.

178 **6. *Density-connected***

179 A DOA estimate $\hat{\mathbf{p}}$ is density-connected to another DOA estimate $\hat{\mathbf{q}}$ if there is a DOA

180 estimate $\hat{\mathbf{m}}$ such that both $\hat{\mathbf{p}}$ and $\hat{\mathbf{q}}$ are density-reachable from it.

181 **7. *Cluster***

182 A cluster S is a subset of U satisfying:

- 183 • $\forall \hat{\mathbf{p}}, \hat{\mathbf{q}} : \text{if } \hat{\mathbf{p}} \in S \text{ and } \hat{\mathbf{q}} \text{ is density-reachable from } \hat{\mathbf{p}}, \text{ then } \hat{\mathbf{q}} \in S \text{ and}$
- 184 • $\forall \hat{\mathbf{p}}, \hat{\mathbf{q}} \in S : \hat{\mathbf{p}} \text{ is density-connected to } \hat{\mathbf{q}}.$

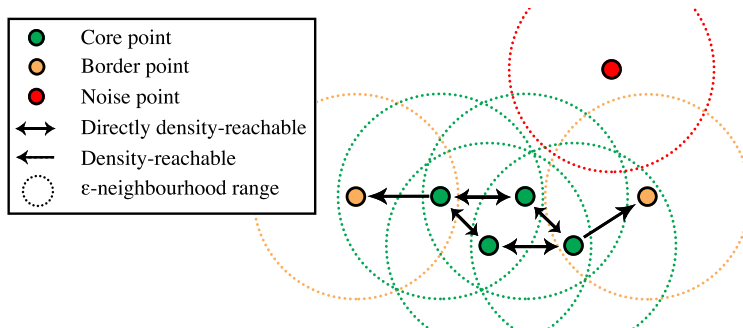


FIG. 3. DOA estimates labelling by DBSCAN with MinPts=3.

185 8. Noise

186 A subset of DOA estimates in U not belonging to any cluster.

187 Figure 3 illustrates the labelling of an example of several DOA estimates by DBSCAN
 188 with MinPts=3. Each core point (green) has at least three DOAs including itself within its
 189 ε -radius neighbourhood while the border (orange) and the noise (red) DOAs do not satisfy
 190 the core point condition.

192 Given the user-defined parameters ε and MinPts, the algorithm first detects all the core
 193 points. A single cluster is identified in two steps: (1) start from an arbitrary core point
 194 and (2) retrieve all points which are density-reachable from it. It then visits the next un-
 195 clustered core point and repeats this process until all core points are clustered. The points
 196 which do not belong to any cluster are labelled as noise.

197 C. MSEC

198 Having performed clustering on data set $U(\tau)$ by either adaptive K-means or DBSCAN,
 199 we obtain the estimated number of clusters $K_c(\tau)$, the clusters $\{S_i(\tau)\}_{i=1}^{K_c(\tau)}$ and the centroids

200 unit vector $\{\hat{\mathbf{c}}_i(\tau)\}_{i=1}^{K_C(\tau)}$ where i is the cluster index. As a representative of the spread of
 201 DOA estimates within each cluster, the average member-to-centroid angular distance $D_i(\tau)$
 202 is calculated for each cluster as

$$D_i(\tau) = \frac{1}{|S_i(\tau)|} \sum_{k \in S_i(\tau)} \angle(\hat{\mathbf{u}}(\tau, k), \hat{\mathbf{c}}_i(\tau)), \quad (8)$$

203 where $\angle(\cdot)$ denotes the angle in degrees between two vectors.

204 The MSEC weight for each DOA estimate is determined from two factors, the cluster
 205 weight and the member weight. For each DOA estimate, the cluster weight, which represents
 206 the normalized measure of concentration in its associated cluster, is

$$\psi(\tau, k) = 1 - \frac{D_i(\tau)}{180}, \quad k \in S_i(\tau), \quad (9)$$

207 and the member weight, which represents the normalized measure of closeness to its associ-
 208 ated centroid, is

$$\lambda(\tau, k) = 1 - \frac{\angle(\hat{\mathbf{u}}(\tau, k), \hat{\mathbf{c}}_i(\tau))}{180}, \quad k \in S_i(\tau). \quad (10)$$

209 The MSEC weight in the TF domain is then formed as

$$w(\tau, k) = \sqrt{\psi(\tau, k)\lambda(\tau, k)}. \quad (11)$$

210 A special case of MSEC with $T = 0$ and $K_{max} = 1$ is proposed in¹⁷, which is based on the
 211 SS assumption within a time-frame.

212 Figure 4 displays DBSCAN and adaptive K-means clusterings of an example distribution
 213 of initial DOA estimates for 5 consecutive frames ($T = 4$). This illustrates that DBSCAN
 214 identifies and ignores the noise DOA estimates due to the use of a static definition of cluster
 215 density while adaptive K-means assigns every DOA estimate to a cluster. Although adaptive

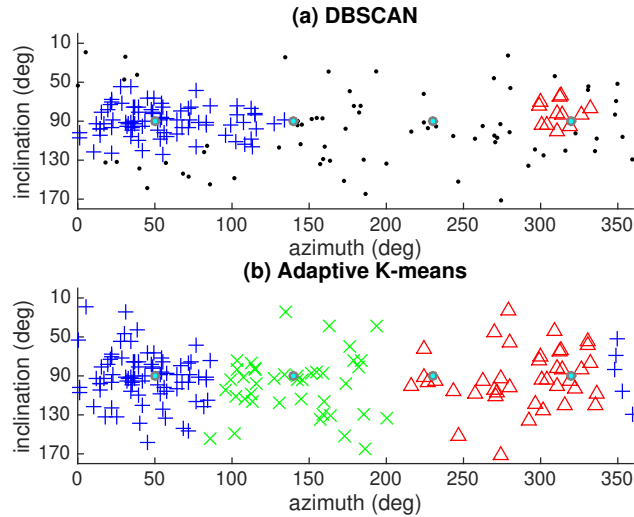


FIG. 4. An example of DOA estimates from 5 consecutive time-frames clustered by (a) DBSCAN with $(\varepsilon, \text{MinPts}) = (20^\circ, 10)$ and (b) adaptive K-means with $K_{max} = 4$. The colours and markers indicate the clusters while the black dots in (a) are the noise DOAs. The true source directions are marked as cyan filled circles.

216 K-means has resulted in detecting more sources, it also includes more erroneous detections
 217 of DOAs and gives less accurate weighting due to the reduced positional accuracy of the
 218 cluster centroid caused by the presence of outliers (erroneous DOAs), as shown next.

219 Figure 5 shows a scatter plot of the normalized MSEC weights versus the normalized
 220 accuracy of the initial DOAs used in the example of Fig. 4. It can be seen that the noise-
 221 robust DBSCAN-MSEC has only weighted strongly the DOAs that have > 0.8 normalized
 222 accuracy, and zero-weighted the inaccurate DOAs with < 0.8 normalized accuracy.

223 Having weighted all the DOA estimates in the TF domain, only the estimates with the $P\%$
 224 strongest weights are selected. One conventional technique to estimate the source directions
 225 from the set of selected DOA estimates is to directly^{18,19} or iteratively^{20,21} find the position

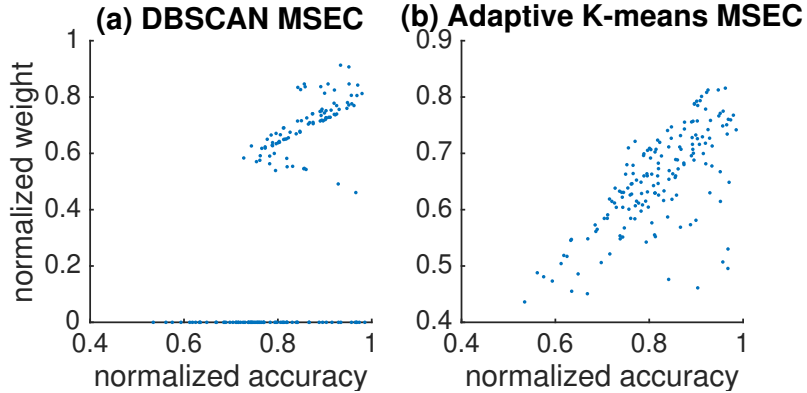


FIG. 5. Normalized weight vs normalized accuracy for MSEC using (a) DBSCAN and (b) adaptive K-means for the example of Fig. 4.

226 of the peaks in the 2D (azimuth \times inclination) smoothed histogram of the selected DOA
 227 estimates. Such techniques assume that the number of sources is known *a priori* and are
 228 sensitive to the smoothing setting (e.g. standard deviation of the smoothing kernel) on the
 229 sources' angular separation, noise level and irregularity in the peaks, which are all assumed
 230 to be unknown in our case. Other clustering-based techniques such as K-means²² or mixture
 231 models using Gaussian²³, Laplacian²⁴ or Von Mises²⁵ distributions are also used in source
 232 direction estimation from the set of initial DOA estimates. These approaches, as well as
 233 peak detection, typically require *a priori* knowledge of the number of sources and are prone
 234 to errors due to outlier DOA estimates.

235 D. Autonomous Robust Source Counting

236 Density-based clustering has received much less attention than distance-based or model-
 237 based clustering techniques in the context of acoustic DOA estimation. We now propose a

238 density-based clustering scheme employing a variant of DBSCAN in an evolutionary frame-
 239 work for source counting and direction estimation from a set of selected DOA estimates.

240 Consider D as the set of selected DOA estimates using MSEC weights (or potentially any
 241 other weighting metrics), which can still include noise DOA estimates. The selection step
 242 ensures that D is significantly more sparse and less noisy than the initial set of all DOA
 243 estimates. In DBSCAN, the threshold density MinPts needs to be chosen and this depends
 244 on the relative density level of the noise DOA estimates in the dataset.

245 As shown in⁹, the original DBSCAN loses reliability in cases with distributions of vary-
 246 ing densities as there may not be a value for MinPts , given ε , for which all densities are
 247 individually clustered. For an example of points distribution in⁹ it is shown that any choice
 248 of MinPts leads either to the erroneous merging of adjacent densities or the missing of the
 249 least dense distribution. Mixtures of distributions with widely varying density often occur in
 250 DOA estimation especially in multi-source acoustic scenarios where one source is less active
 251 or relatively more distant with respect to the microphone array compared to other sources.
 252 Variations of DBSCAN²⁶⁻³⁰ are proposed but all require user intervention for setting pa-
 253 rameters. The DBSCAN employed in MSEC uses an empirically-chosen static MinPts in
 254 order to avoid extremely high computational cost per TF bin. But for a one-run process-
 255 ing of the set of estimates D , the use of dynamic MinPts can improve the performance
 256 of clustering. Unlike DBSCAN, evolutive DBSCAN⁹, for a fixed ε , uses varying MinPts
 257 $\in [\min(|N_\varepsilon(D)|) + 1, \max(|N_\varepsilon(D)|) - 1]$. The step-size for searching MinPts can be either
 258 defined by the user or calculated based on the maximum number of iterations (NumIt)
 259 specified by the user.

At each iteration, after a comparison between the current clustering and the previous clustering, the current clusters are labelled as either ‘dead’ or ‘alive’ each defined as follows

1) Dead: A cluster is dead if one the following two conditions is met. It has a shared member with more than one alive cluster in the last iteration (merge condition) or it has a shared member with any previously dead cluster (re-occurrence of a previously merged cluster). **2) Alive:** A cluster is alive if it is not dead.

At each iteration, the weight and the centroid of the alive clusters are stored where cluster weight is defined as the mean density of the clustered DOAs. The pseudocode for the main part of this algorithm is provided in Algorithm 1.

Having obtained M centroids $\{c_i\}_{i=1}^M$ of alive clusters and their associated weights $\{w_i\}_{i=1}^M$, one final autonomous DBSCAN is applied on the set of centroids which finally estimates the number of detected clusters, L , and their final centroids $\{d_i\}_{i=1}^L$ indicating the estimated number of sources and the final DOA estimates respectively, as shown in Fig. 6(c). Assuming densities of DOAs in D have non-radical skewness, we expect very low spatial variance for the centroids belonging to a repetitively alive cluster at consecutive iterations. Therefore a small value of $\varepsilon_f = 5^\circ$ is defined in the final DBSCAN while MinPts is autonomously determined as follows. A sorted weighted-density graph is built for the centroids $\{c_i\}_{i=1}^M$ using their weights $\{w_i\}_{i=1}^M$ and densities $\{|N_{\varepsilon_f}(c_i)|\}_{i=1}^M$. The use of the weights exaggerates the dynamic range and so the angle of the ‘knee’ in the graph. This is because the outlier centroids are expected to be from low density clusters of outlier DOA estimates that might have been clustered due to low MinPts at the end of evolutionary process. The estimated MinPts for the final DBSCAN is the density at the position of the ‘prominent’ knee with

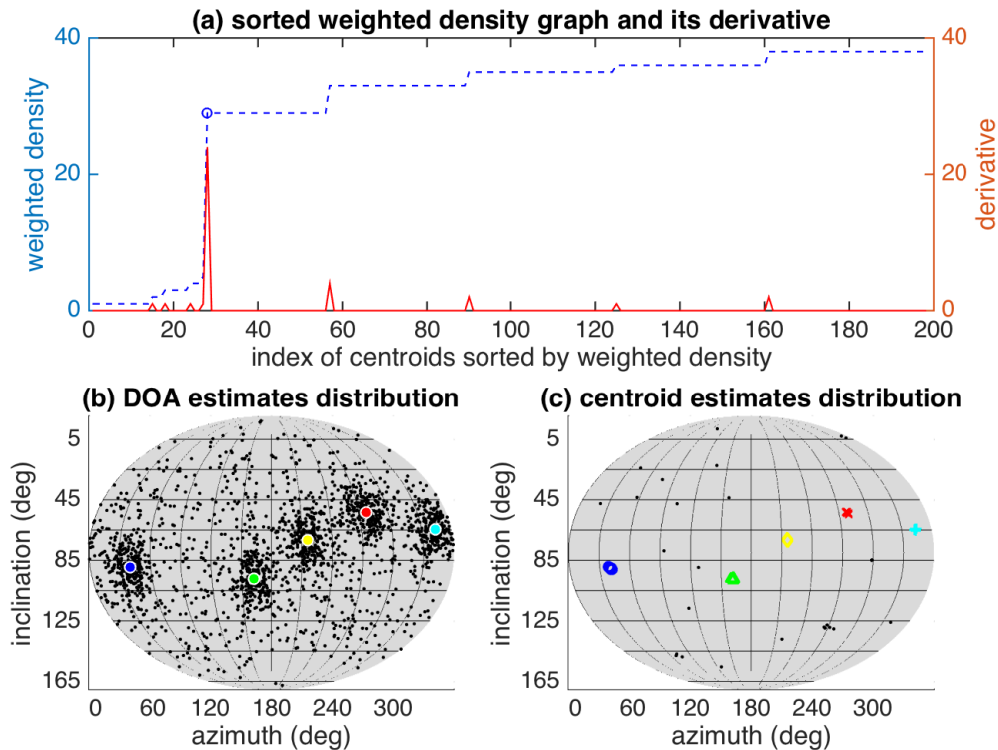


FIG. 6. Evolutive DBSCAN on an example set of selected DOA estimates D . (a) sorted weighted density (blue dashed) graph and its derivative (solid red). Position of the knee marked as blue circle. Distribution of (b) DOA estimates (c) centroid estimates for 5 sources. Final centroids are marked by coloured circles.

282 the lowest weighted density. This is derived as the position of the first peak (excluding the
 283 peaks with less than 10% of the highest peak) in its derivative function as shown in Fig.
 284 6(a). Note that if $\min(\{|N_{\epsilon_f}(c_i)|\}_{i=1}^M) > 1$, no knee detection is needed and MinPts is the
 285 minimum density.

Algorithm 1 Evolutive DBSCAN

```

function EVOLUTIVE_DBSCAN(points)

    centroids=[]; %holds alive centroids and weights

    MinPts=max(|Nε(.)|) - 1; cntr=1;

    while (MinPts ≥ min(|Nε(.)|) + 1) OR (cntr ≤ NumIt)

        C=DBSCAN(points,ε,MinPts);

        if isEmpty(centroids) then

            centroids += C(all).centroid; %initialization

            C.dead_members=[]; %all dead members

        else

            C=LABEL_CLUSTERS(C,C_last);

            if anyClusterAlive(C) then

                centroids += C(alive_ones).centroid;

                C=RemoveDead(C,dead_ones);

            end if

        end if

        C_last=C; MinPts -= step; cntr += 1;

    end

    return centroids

end function

```

287 IV. EVALUATIONS

288 The performance of the proposed method is first evaluated using recorded anechoic speech
 289 convolved with simulated room impulse response for Spherical Microphone Array (SMA)^{43–45}
 290 in the presence of reverberation and sensor noise. Performance using real speakers in a
 291 reverberant room is considered in Section V. The evaluation is performed for a varying
 292 number of sources and angular separation. The DPD method is used as a baseline for
 293 comparison. Without loss of generality, the inclination of sources is fixed at 90° , for simulated
 294 data, so as to place them in the same horizontal plane as the microphone array for clarity
 295 of systematic evaluation of the effect of source separation. However, inclination is varied in
 296 the experimental verification using real data in Section V.

297 The room impulse responses of a 32-element rigid SMA with radius of 4.2 cm (corre-
 298 sponding to the em32 Eigenmike[®]) in a $5 \times 6 \times 4$ m shoebox room with $T_{60} = 0.4$ s³¹ were
 299 simulated using the image method³² implemented by³³. N_s sources were randomly placed
 300 with azimuth interval of $\Delta\phi$ degrees at a distance of 1 m from the centre of the SMA on the
 301 same horizontal plane as SMA. For each N_s and $\Delta\phi$, 100 random trials were used in each of
 302 which the first azimuth was randomly selected from a uniform circular distribution around
 303 the SMA. The source signals consisted of different anechoic speech signals randomly selected
 304 for each trial from the APLAWD database³⁴. The active level of each speech source accord-
 305 ing to ITU-T P.56³⁵, as measured for the omnidirectional eigenbeam, is set to be equal across
 306 all trials. Spatio-temporally white Gaussian noise was added to the microphone signals to
 307 produce an SNR of 25 dB for each source. A sampling frequency of 8 kHz was used with

308 50% overlapping time-frames of 8 ms duration. Any narrowband method can be used for
 309 the DOA estimator but for fast computation, the efficient Pseudointensity Vectors (PIVs)³⁶
 310 method was used in these test as an example SS DOA estimator to obtain the initial DOA
 311 estimates. PIVs use eigenbeams up to the first-order spherical harmonic³⁷.

In DPD⁷ using SMAs, the covariance matrix is approximated as the average covariance matrix over a local TF region^{7,12}

$$\mathbf{R}(\tau, k) = \frac{1}{J_\tau J_k} \sum_{j_\tau=0}^{J_\tau-1} \sum_{j_k=0}^{J_k-1} \mathbf{a}(\tau - j_\tau, k + j_k) \times \mathbf{a}^H(\tau - j_\tau, k + j_k), \tag{12}$$

312 where $J_\tau = 6$ and $J_k = 4$ are the widths (number of bins) of the averaging windows over time
 313 and frequency respectively. This gives 32 ms and 500 Hz window-size in the TF domain based
 314 on our time and frequency resolution. The column vector \mathbf{a} contains spherical harmonic
 315 eigenbeams up to the third order and $(.)^H$ denotes the Hermitian transpose.

316 MSEC has a temporal window size of $T = 4$ frames in (3), which is chosen to be small
 317 enough to decompose the problem of N sources into $L < N$ sources over the interval and
 318 wide enough to form distinguishable densities for consistent DOAs. For clusterings used in
 319 variations of MSEC, adaptive K-means has $K_{max} = 4$ with random initialization per K and
 320 DBSCAN has $\varepsilon = 10^\circ$ and MinPts= 10 which is approximately 5% of the number of the
 321 estimates in dataset $U(\tau)$. These values for the setting parameters of the evaluated methods
 322 are empirically chosen. Both MSEC alternatives mainly rely on two user-tuned parameters,
 323 T for MSEC weighting in addition to K_{max} or MinPts for adaptive K-means and DBSCAN
 324 respectively. Note that the pair of ε and MinPts in DBSCAN are not independent since

325 the user can use a fixed $\varepsilon = 10^\circ$ and adjust MinPts only for optimum results. Therefore
 326 the number of user-tuned parameters in both MSEC approaches is the same as in the DPD
 327 approach, which is also based on two J_τ and J_k user-defined parameters.

328 A uniform weighting strategy, in which all DOA estimates are selected, is also included
 329 in the evaluation as a reference. For the purpose of evaluating the performance of the
 330 weighting metrics only, a fixed selection percentage of $P = 25\%$ is empirically suggested⁸
 331 and used for DPD and both variations of MSEC. Therefore DPD and MSEC both select
 332 an equal number of DOA estimates, which is the top 25% DOAs with the highest weights
 333 while uniform weighting selects all DOA estimates. The error (in degrees) for each selected
 334 DOA estimate is calculated as the angular distance between the estimate and the nearest
 335 true DOA.

336 A. Accuracy of the selected DOAs

337 In this section the accuracy of the DOA estimates selected by the weights is evaluated.
 338 Figure 7 shows the mean error of the DOA estimates selected by each method for $\Delta\phi =$
 339 $\{45^\circ, 90^\circ\}$ and incremental $N_s = \{2, 3, 4\}$. It can be seen that MSEC variations select
 340 significantly more accurate DOA estimates compared to DPD and uniform weighting, which
 341 validates the advantage of MS over SS assumption in MS scenario. DBSCAN-based MSEC
 342 has 92% to 129% mean accuracy improvement in these tests compared to DPD due to the
 343 dynamic MS assumption and noise-robustness. It can also be seen that as N_s increases the
 344 mean accuracy of the uniform and DPD weights improves. This is due to the decrease in
 345 the least possible error as the number of sources increases.

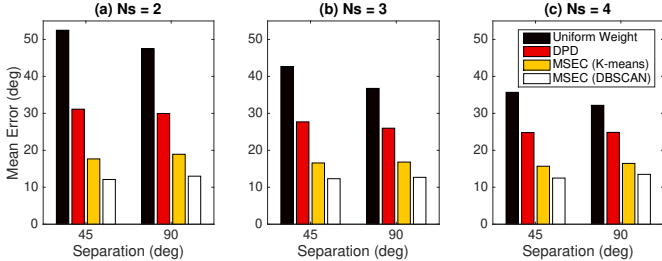


FIG. 7. The overall mean error of the DOAs for varying separation and incremental number of sources.

346 Figure 10 shows the top view and side view of the normalized smoothed histogram of
 347 DOA estimates selected by each method for an example experimental trial. The perfor-
 348 mance benefits of MSEC are shown and can be explained by observing the distinctness and
 349 sharpness of the peaks. It can be seen that MSEC variations have well defined peaks around
 350 each of the source positions especially for the fourth source (from the left) where DPD fails
 351 due to the oversize processing TF region at the TF bins with a significantly dominant fourth
 352 source resulting in selection of inaccurate DOA estimates. The reason for such failure is
 353 visualised and further discussed in the TF domain in Section IV C.

354 **B. Correlation between weights and DOA estimate accuracy**

355 Figure 8 shows the mean correlation between the normalized weights and the normalized
 356 accuracy of their DOA estimate. The normalized accuracy is

$$1 - error/180, \tag{13}$$

357 where *error* (in degrees) is the spherical angle between the DOA estimate and the nearest
 358 true DOA. DPD weights show low correlation with accuracy. On the other hand, MSECs

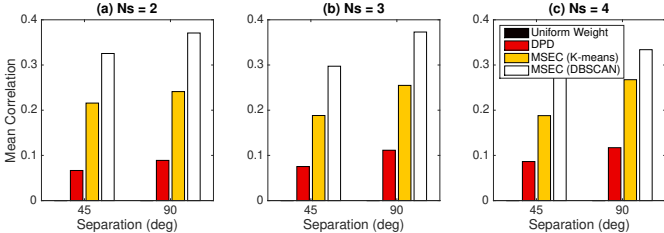


FIG. 8. The overall mean correlation between the normalized weight and accuracy for varying separation and incremental number of sources.

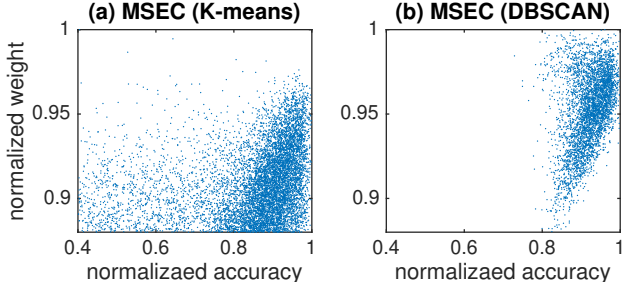


FIG. 9. Distribution of the normalized weights and their DOA estimate accuracy for an example trial with $(N_s, \Delta\phi) = (2, 90^\circ)$.

359 are, at least by a factor of 4, more linearly correlated with DOA estimate accuracy. This
 360 is due to two reasons. (1) MSEC is calculated using the DOA estimates and is therefore
 361 expected to be directly impacted by DOA accuracy unlike DPD which uses eigenbeams. (2)
 362 The MSEC metric is calculated in the spatial domain using angular distances which has the
 363 same unit and nature as the DOA estimate accuracy whereas the DPD metric uses the SVR
 364 of the eigenbeams.

365 Figure 9 illustrates a scatter plot of the selected normalized weights versus normalized
 366 accuracy of their DOA estimates for K-means and DBSCAN-based MSEC for an example

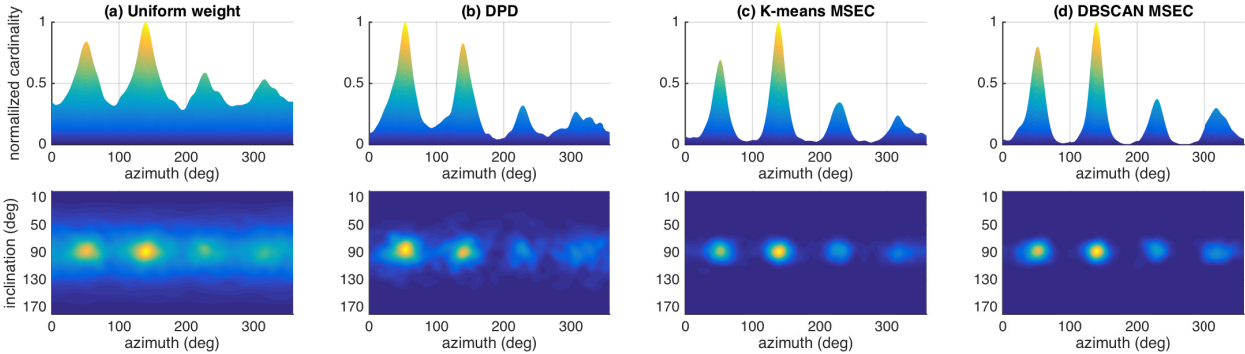


FIG. 10. The side view (top row) and the top view (bottom row) of the normalized smoothed histogram of the selected DOA estimates using (a) Uniform weights including all DOAs, (b) DPD, (c) adaptive K-means MSEC and (d) DBSCAN MSEC for an example trial with $(N_s, \Delta\phi) = (4, 90^\circ)$.

367 trial. It can be seen that DBSCAN-based weighting has significantly fewer inaccurate DOA
 368 estimates which are falsely weighted high compared to K-means. This is due to two reasons.
 369 (1) DBSCAN is a noise-robust clustering technique and is more capable of ignoring the in-
 370 accurate DOA estimates. (2) The outcome clustering of K-means is stochastic for each run
 371 because of random initialization and dependency of the outcome on the initialization, while
 372 DBSCAN does not require initialization and its outcome is therefore deterministic. During
 373 an experimental analysis it was observed that different trials of K-means on the same dataset
 374 with the same choice of K sometimes led to inconsistent clusterings and therefore incon-
 375 sistent estimation of $K_c(\tau)$. Such inconsistent behaviour can sometimes lead to erroneous
 376 clustering and so erroneous weighting.

C. Effect of weightings on counting and direction estimation of sources

In this section the performance of each SS-validity confidence metric is evaluated in the context of source direction estimation and source counting using evolutive DBSCAN⁹ presented in Section III D. In⁹ it is shown that the evolutive DBSCAN outperforms the conventional histogram peak picking as well as adaptive K-means and original DBSCAN techniques and is therefore chosen as our source counting and source direction extraction technique in this paper. The choice of NumIt=50 was empirically found to be a good trade-off between reliability and computational efficiency for our proposed evolutive DBSCAN. MSEC based on K-means is excluded from the evaluation in this section since DBSCAN-MSEC has a better performance as shown in the previous sections.

The two performance metrics Successful Localization Rate (SLR) and Mean Error respectively represent the source counting and DOA estimation accuracy. SLR is the percentage of trials for which the correct number of sources was detected and all the best case data associated pairs of estimate-true DOA are less than 20° , which is half of the minimum source separation used in the evaluation. The mean error is calculated for the successfully localized cases where all sources are detected.

Figure 11 shows the mean error and SLR of DPD and DBSCAN-MSEC, abbreviated to MSEC in this section, for varying $\Delta\phi$ and N_s . It can be seen that MSEC outperforms DPD in all cases. In terms of DOA estimation accuracy, although MSEC and DPD perform very closely, MSEC slightly leads by 1° at 45° separation with 4 sources. In terms of source counting accuracy, MSEC significantly leads especially for $\Delta\phi = 45^\circ$ as N_s increases. MSEC

398 also shows strong robustness to separation and number of sources as its SLR drops only to
 399 75% while DPD’s SLR is reduced to 20% with the decrease in $\Delta\phi$ and increase in N_s . Such
 400 results match with the observation in Fig. 10. It is seen that the peaks of the multi-modal
 401 distributions, which affect the accuracy of DOA estimation, remain approximately at the
 402 same position for DPD and MSEC while the sharpness and distinctness of the peaks, which
 403 affect the source counting, are significantly different.

404 Figure 13 shows the TF bins with the top $P = 25\%$ strongest MSEC and DPD weights
 405 as well as the bins with PIV DOA estimates, which have $\leq 10^\circ$ error and are considered
 406 as accurate DOAs, for an example trial. As shown in Fig. 13(c), accurate DOAs occur at
 407 varying-size TF regions and even at isolated TF bins. It can be clearly seen that MSEC
 408 has been more successful in detecting varying-size TF regions and isolated TF bins due to
 409 dynamic MS assumption over relatively large analysis window-size compared to DPD, which
 410 is based on SS assumption over small analysis window-size.

411 V. EXPERIMENTAL VERIFICATION USING REAL-WORLD DATA

412 In this section the performance of each method is evaluated using real recordings in a
 413 reverberant room. Recordings of 4 s speech utterances in a room with approximate dimen-
 414 sions of $10 \times 9 \times 2.5$ m and $T_{60} = 0.4$ s were obtained using an Eigenmike 32-channel rigid
 415 SMA with radius of 4.2 cm placed close to the centre of the room. Four talkers were simul-
 416 taneously active and were located 1.5 m away from the centre of the array at approximately
 417 60° intervals while their inclinations alternated to be above or below the horizontal plane
 418 of the array. Figure 12 shows the normalized smoothed histograms for uniform weighting

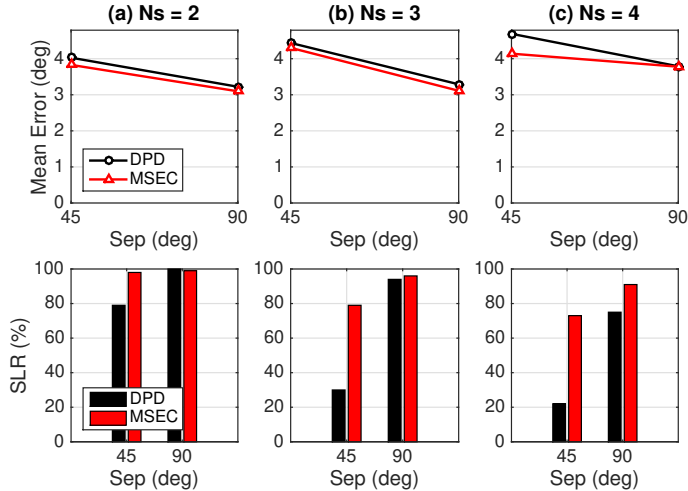


FIG. 11. Mean error (top row) and SLR (bottom row) for (a) 2, (b) 3 and (c) 4 sources with varying source separation.

419 using all DOA estimates, DPD and DBSCAN-MSEC using $P = 25\%$ of the DOA estimates
 420 with the strongest weights, where DOA estimates are obtained using PIVs³⁶. Due to only
 421 approximate knowledge of the ground-truth position of sources and array in the physical
 422 room, accurate numerical estimation error cannot be obtained. The approximate mean es-
 423 timation error for all methods is 4° . All methods successfully estimate peaks corresponding
 424 to all four sources due to wide separation of sources. In order to provide a numerical eval-
 425 uation, for each peak a measure of ‘peak strength’ as suggested in³⁸ is used which is the
 426 ratio of the peak height over the peak smoothness where the peak smoothness is defined as
 427 the average height in the normalized peak distribution within its range of $r_p = 30^\circ$ (half of
 428 source separation) neighbourhood. Table I presents the peak strength of each peak for all
 430 methods. The smoothed histograms in Fig. 12 and the peak strengths in Table I show
 434 that MSEC significantly outperforms the baseline and the state-of-the-art methods using

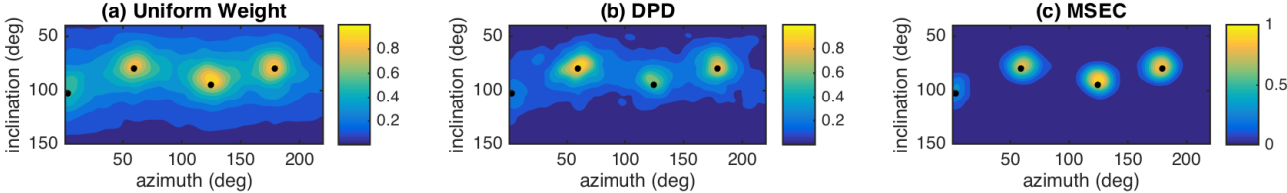


FIG. 12. Normalized smoothed histogram for uniform weighting (all DOA estimates), DPD and DBSAN-MSEC (both based on $P = 25\%$) using real recording. The black dot represents the approximate true DOA.

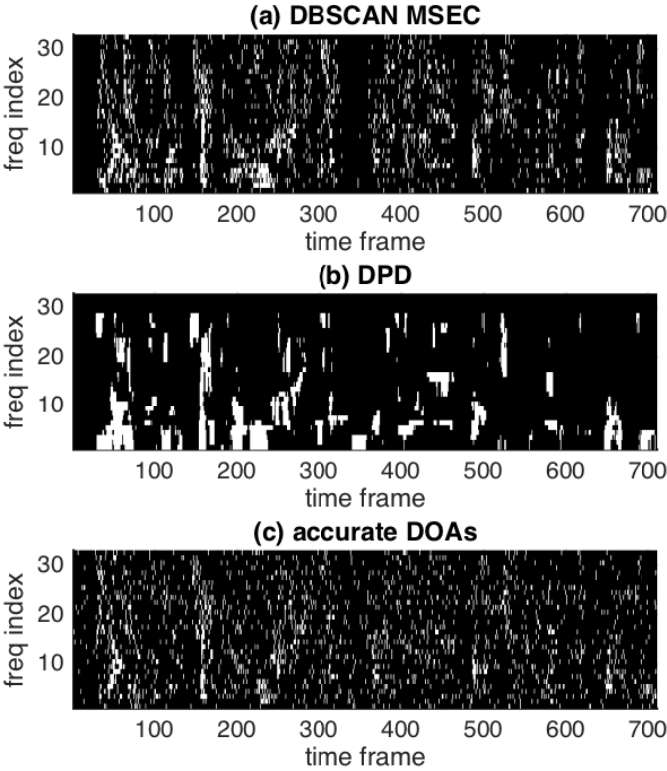


FIG. 13. TF bins with top $P = 25\%$ strongest (a) DBSCAN-MSEC weights, (b) DPD weights and (c) $\leq 10^\circ$ DOA error for $(N_s, \Delta\phi) = (3, 90^\circ)$.

Peak	Uniform Weight	DPD	MSEC
1	2.08	2.94	6.04
2	1.96	2.59	6.03
3	1.82	1.75	4.97
4	0.99	0.67	4.31
Mean	1.71	1.99	5.33

TABLE I. Peak Strength of each peak for all methods

435 real recordings and serves towards validation of the evaluation results based on simulation
 436 in the previous section.

437 VI. CONCLUSION

438 A confidence metric for validity of SS assumption in a TF bin has been proposed using
 439 spatial consistency of initial DOA estimates. It employs adaptive K-means based on AIC
 440 or noise-robust DBSCAN clusterings to group spatially consistent initial DOA estimates,
 441 which are derived by a SS-based DOA estimator. Each DOA estimate is weighted using its
 442 distance-to-centroid and cluster's spread, and finally, the DOA estimates with the strongest
 443 weights are selected to be used in source counting and source direction estimation. The
 444 proposed metric is based on MS assumption over a relatively large TF region compared to
 445 conventional metrics, which are based on SS assumption over a small-size TF region. A
 446 novel use of density-based DBSCAN clustering in the context of source localization has also

447 been used to propose an autonomous evolutionary method for source counting and final
448 source direction estimation. The evolutive DBSCAN uses DBSCAN iteratively for varying
449 density threshold. Such variation makes it robust to a mixture of distributions with varying
450 density. The evaluations using simulation and real recordings show that our proposed metric
451 significantly improves the performance of source counting, compared to the baseline and the
452 state-of-the-art metrics, and provides at least the same accuracy as the state-of-the-art for
453 source direction estimation.

454 REFERENCES

455 ¹H. W. Löllmann, C. Evers, A. Schmidt, H. Mellmann, H. Barfuss, P. A. Naylor, and
456 W. Kellermann, “The LOCATA Challenge Data Corpus for Acoustic Source Localiza-
457 tion and Tracking,” in *IEEE Sensor Array and Multichannel Signal Processing Workshop*
458 (*SAM*), Sheffield, UK (2018).

459 ²S. Rickard and Z. Yilmaz, “On the approximate W-disjoint orthogonality of speech,”
460 in *Proc. IEEE Intl. Conf. on Acoustics, Speech and Signal Processing (ICASSP)* (2002),
461 Vol. 1, pp. 529–532.

462 ³R. O. Schmidt, “Multiple emitter location and signal parameter estimation,” *IEEE Trans.*
463 *Antennas Propag.* **34**(3), 276–280 (1986).

464 ⁴P. A. Naylor and N. D. Gaubitch, eds., *Speech Dereverberation* (Springer, 2010).

465 ⁵S. Mohan, M. E. Lockwood, M. L. Kramer, and D. L. Jones, “Localization of multiple
466 acoustic sources with small arrays using a coherence test,” *J. Acoust. Soc. Am.* **123**,
467 2136–2147 (2008).

468 ⁶D. Pavlidi, M. Puigt, A. Griffin, and A. Mouchtaris, “Real-time multiple sound source
469 localizaation using a circular microphone array based on single-source confidence mea-
470 sures,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Processing (ICASSP)*, Kyoto,
471 Japan (2012), pp. 2625–2628.

472 ⁷O. Nadiri and B. Rafaely, “Localization of multiple speakers under high reverberation
473 using a spherical microphone array and the direct-path dominance test,” *IEEE Trans.*
474 *Audio, Speech, Lang. Process.* **22**(10), 1494–1505 (2014).

- 475 ⁸S. Hafezi, A. H. Moore, and P. A. Naylor, “Multiple DOA estimation based on estima-
476 tion consistency and spherical harmonic multiple signal classification,” in *Proc. European*
477 *Signal Processing Conf. (EUSIPCO)*, Kos, Greece (2017), pp. 1280–1284.
- 478 ⁹S. Hafezi, A. H. Moore, and P. A. Naylor, “Robust source counting for DOA estimation
479 using density-based clustering,” in *IEEE Sensor Array and Multichannel Signal Processing*
480 *Workshop (SAM)*, Sheffield, UK (2018).
- 481 ¹⁰H. Wang and M. Kaveh, “Coherent signal-subspace processing for the detection and esti-
482 mation of angles of arrival of multiple wide-band sources,” *IEEE Trans. Acoust., Speech,*
483 *Signal Process.* **33**(4), 823–831 (1985).
- 484 ¹¹H. Akaike, “A new look at the statistical model identification,” *IEEE Trans. Autom.*
485 *Control* **AC-19**(6), 716–723 (1974).
- 486 ¹²A. H. Moore, C. Evers, and P. A. Naylor, “2D direction of arrival estimation of multiple
487 moving sources using a spherical microphone array,” in *Proc. European Signal Processing*
488 *Conf. (EUSIPCO)* (2016), (submitted).
- 489 ¹³J. Ahonen and V. Pulkki, “Diffuseness estimation using temporal variation of intensity
490 vectors,” in *Proc. IEEE Workshop on Applications of Signal Processing to Audio and*
491 *Acoustics (WASPAA)*, New Paltz, NY, USA (2009), pp. 285–288.
- 492 ¹⁴D. P. Jarrett, O. Thiergart, E. A. P. Habets, and P. A. Naylor, “Coherence-based dif-
493 fuseness estimation in the spherical harmonic domain,” in *Proc. IEEE Convention of*
494 *Electrical & Electronics Engineers in Israel (IEEEI)*, Eilat, Israel (2012).

- 495 ¹⁵M. Ester, H. P. Krigel, J. Sander, and X. Xu, “A density-based algorithm for discovering
496 clusters in large spatial database with noise,” in *Proc. of the 2nd International Conference*
497 *on Knowledge Discovery and Data Mining*, Portland, WA (1996), pp. 226–231.
- 498 ¹⁶C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval* (Cam-
499 bridge University Press, Cambridge, UK, 2008).
- 500 ¹⁷S. Hafezi, A. H. Moore, and P. A. Naylor, “Multiple source localization using estimation
501 consistency in the time-frequency domain,” in *Proc. IEEE Int. Conf. Acoust., Speech,*
502 *Signal Processing (ICASSP)*, New Orleans, LA, USA (2017).
- 503 ¹⁸S. Hafezi, A. H. Moore, and P. A. Naylor, “3D acoustic source localization in the spherical
504 harmonic domain based on optimized grid search,” in *Proc. IEEE Int. Conf. Acoust.,*
505 *Speech, Signal Processing (ICASSP)*, Shanghai, China (2016).
- 506 ¹⁹S. Hafezi, A. H. Moore, and P. A. Naylor, “Multiple source localization in the spheri-
507 cal harmonic domain using augmented intensity vectors based on grid search,” in *Proc.*
508 *European Signal Processing Conf. (EUSIPCO)*, Budapest, Hungary (2016).
- 509 ²⁰A. Griffin, D. Pavlidi, M. Puigt, and A. Mouchtaris, “Real-time multiple speaker doa
510 estimation in a circular microphone array based on matching pursuit,” in *Proc. European*
511 *Signal Processing Conf. (EUSIPCO)*, Bucharest, Romania (2012).
- 512 ²¹D. Pavlidi, S. Delikaris-Manias, V. Pulkki, and A. Mouchtaris, “3D localization of multiple
513 sound sources with intensity vector estimates in single source zones,” in *Proc. European*
514 *Signal Processing Conf. (EUSIPCO)* (2015).

- 515 ²²B. Loesch and B. Yang, “Source number estimation and clustering for underdetermined
516 blind source separation,” in *IWAENC*, Seattle, WA, USA (2008).
- 517 ²³W. Zhang and B. Rao, “A two microphone-based approach for source localization of multi-
518 ple speech sources,” *IEEE Trans. Audio, Speech, Lang. Process.* **18**(8), 1913–1928 (2010).
- 519 ²⁴M. Cobos, J. J. Lopez, and D. Martinez, “Two-microphone multi-speaker localization
520 based on a laplacian mixture model,” *Digital Signal Processing* **18**(1), 66–76 (2011).
- 521 ²⁵C. Kim, C. Khawand, and R. M. Stern, “Two-microphone source separation algorithm
522 based on statistical modeling of angle distributions,” in *Proc. IEEE Int. Conf. Acoust.,
523 Speech, Signal Processing (ICASSP)*, Kyoto, Japan (2012), pp. 4629–4632.
- 524 ²⁶A. Ram, A. Sharma, A. S. Jalall, R. Singh, and A. Agrawal, “An enhanced density based
525 spatial clustering of applications with noise,” in *IEEE International Advance Computing
526 Conferece (IACC)*, Patiala, India (2009).
- 527 ²⁷P. Liu, D. Zhou, and N. Wu, “Vdbscan: Varied density based spatial clustering of appli-
528 cations with noise,” in *Proc. of IEEE International Conference on Service Systems and
529 Service Management*, Chengdu, China (2007), pp. 1–4.
- 530 ²⁸C. Xiaoyun, M. Yufang, Z. Yan, and W. Ping, “Gmdbscan: Multi-density dbscan cluster
531 based on grid,” in *IEEE International Conference on e-Business Enginerring (ICEBE)*
532 (2008).
- 533 ²⁹Z. Xiong, R. Chen, Y. Zhang, and X. Zhang, “Multi-density dbscan algorithm based on
534 density levels partitioning,” *Journal of Information and Computational Science* **9**, 2739–
535 2749 (2012).

- 536 ³⁰O. Uncu, W. A. Gruver, D. B. Kotak, D. Sabaz, Z. Alibhai, and C. Ng, “Gridbscan:
537 Grid density-based spatial clustering of applications with noise,” in *IEEE Intl. Conf. on*
538 *Systems, Man and Cybernetics*, Taipei, Taiwan (2006).
- 539 ³¹P. A. Naylor and N. D. Gaubitch, “Speech dereverberation,” in *Proc. Intl. Workshop*
540 *Acoust. Echo and Noise Control (IWAENC)*, Eindhoven, The Netherlands (2005).
- 541 ³²J. B. Allen and D. A. Berkley, “Image method for efficiently simulating small-room acous-
542 tics,” *J. Acoust. Soc. Am.* **65**(4), 943–950 (1979).
- 543 ³³D. P. Jarrett, E. A. P. Habets, M. R. P. Thomas, and P. A. Naylor, “Simulating room im-
544 pulse responses for spherical microphone arrays,” in *Proc. IEEE Intl. Conf. on Acoustics,*
545 *Speech and Signal Processing (ICASSP)*, Prague, Czech Republic (2011), pp. 129–132.
- 546 ³⁴G. Lindsey, A. Breen, and S. Nevard, “SPAR’s archivable actual-word databases,” Tech-
547 nical Report, University College London (1987).
- 548 ³⁵ITU-T, “Objective measurement of active speech level,” (2011), [http://www.itu.int/](http://www.itu.int/rec/T-REC-P.56-201112-I/en)
549 [rec/T-REC-P.56-201112-I/en](http://www.itu.int/rec/T-REC-P.56-201112-I/en).
- 550 ³⁶D. P. Jarrett, E. A. P. Habets, and P. A. Naylor, “3D source localization in the spherical
551 harmonic domain using a pseudointensity vector,” in *Proc. European Signal Processing*
552 *Conf. (EUSIPCO)*, Aalborg, Denmark (2010), pp. 442–446.
- 553 ³⁷D. P. Jarrett, E. A. Habets, and P. A. Naylor, *Theory and Applications of Spherical Micro-*
554 *phone Array Processing*, Springer Topics in Signal Processing (Springer, Berlin Heidelberg,
555 2016).

- 556 ³⁸S. Hafezi, A. H. Moore, and P. A. Naylor, “Augmented intensity vectors for direction of
557 arrival estimation in the spherical harmonic domain,” *IEEE/ACM Transactions on Audio,
558 Speech, and Language Processing* **25**(10), 1956–1968 (2017).
- 559 ³⁹E. Sato, and Y. Tatakura, “Fast multiple moving sound sources localization utilizing
560 sparseness of speech signals,” *The Journal of the Acoustical Society of America* **140**(4),
561 3061 (2016).
- 562 ⁴⁰D. Khaykin, and B. Rafaely, “Acoustic analysis by spherical microphone array processing
563 of room impulse responses,” *The Journal of the Acoustical Society of America* **132**(1),
564 261–270 (2012).
- 565 ⁴¹J. L. Yuxiang Hu, and X. Qiu, “A maximum likelihood direction of arrival estimation
566 method for open-sphere microphone arrays in the spherical harmonic domain,” *The Jour-
567 nal of the Acoustical Society of America* **138**(2), 791–794 (2015).
- 568 ⁴²D. Levin, E. A. P. Habets, and S. Gannot, “On the angular error of intensity vector based
569 direction of arrival estimation in reverberant sound fields,” *The Journal of the Acoustical
570 Society of America* **128**(4), 1800–1811 (2010).
- 571 ⁴³H. Sun, E. Mabande, K. Kowalczyk, and W. Kellermann, “Localization of distinct reflec-
572 tions in rooms using spherical microphone array eigenbeam processing,” *The Journal of
573 the Acoustical Society of America* **131**(4), 2828 (2012).
- 574 ⁴⁴K. Haddad, and J. Hald, “3D localization of acoustic sources with a spherical array,” *The
575 Journal of the Acoustical Society of America* **123**(5), 3311 (2008).

576 ⁴⁵M. R. Bai, and Y. H. Yao, “Source localization and signal extraction using spherical mi-
577 crophone arrays,” , The Journal of the Acoustical Society of America **137**(4), 2232 (2015).