



**MULTICHANNEL OVERLAPPING SPEAKER SEGMENTATION USING MULTIPLE
HYPOTHESIS TRACKING OF ACOUSTIC AND SPATIAL FEATURES**

June, 2021

Aidan Hogg, Christine Evers and Patrick Naylor

Electrical and Electronic Engineering, Imperial College London, UK

Electronics and Computer Science, University of Southampton, UK

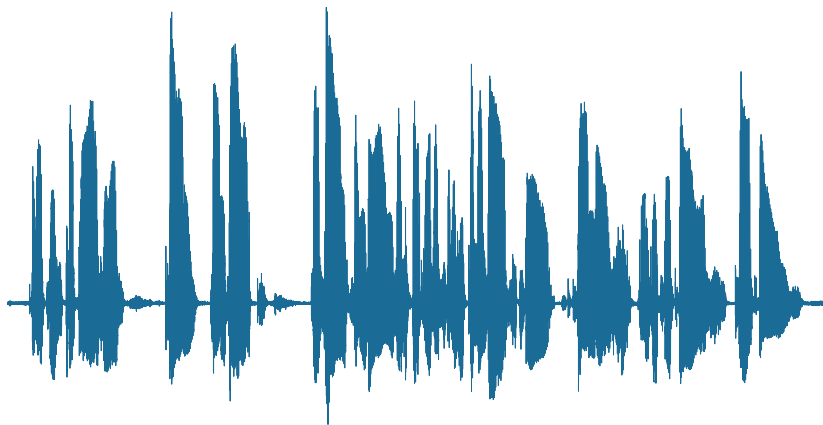
What is speaker diarization?

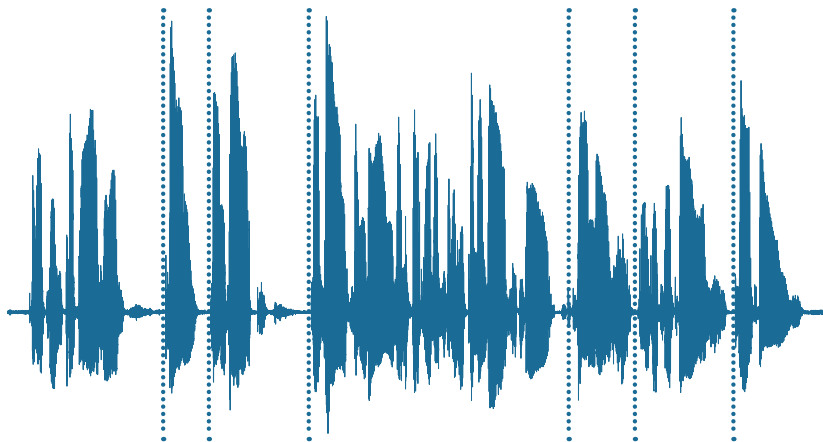
Answers the question “who spoke when?” in an audio recording.

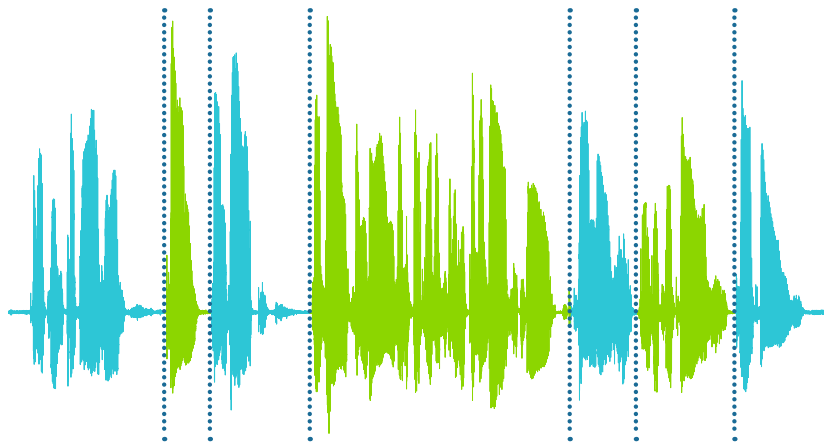
Is diarization really that useful?

- Speaker indexing and rich transcription
- Speaker segmentation and clustering helping Automatic Speech Recognition (ASR) systems
- Preprocessing modules for single speaker-based algorithms

DIARIZATION METHOD

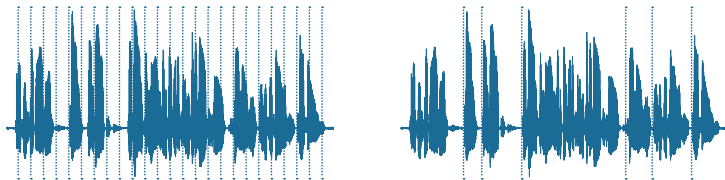






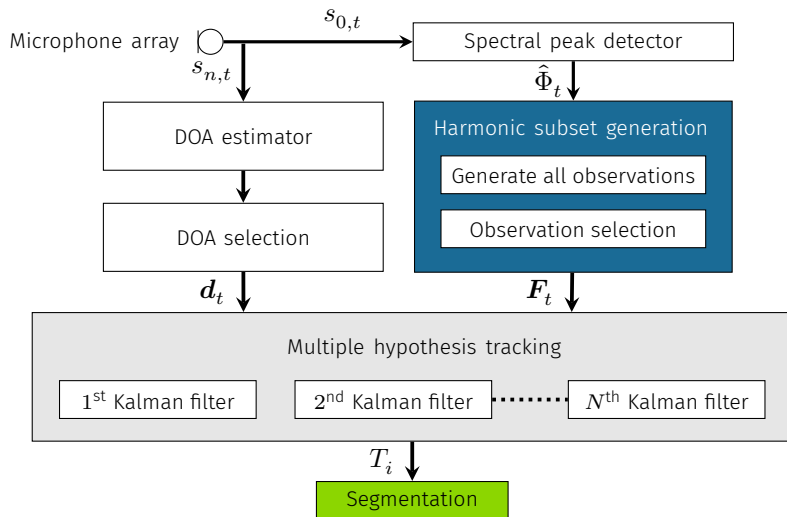
Is good segmentation really that useful?

Why not just segment the audio stream into small uniform segments and cluster with realignment?



If the speech segments are small then each segment only contains a small amount of information that can be used for clustering.

PROPOSED SYSTEM ARCHITECTURE



In this work all the harmonics of voiced speech are tracked along with the DOA estimates so that overlapping speech can be considered.

F_0 harmonics and DOA observation:

$$\mathbf{z}_{t,n} = [\mathbf{f}_{t,m}, d_{t,v}]^T,$$

The state equation for the i^{th} speaker:

$$\mathbf{x}_{i,t} = \mathbf{x}_{i,t-1} + \mathbf{w}_t, \quad \mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_t),$$

where

$$\mathbf{x}_i = [x_f, x_d]^T, \quad \mathbf{Q}_t = \text{diag}(\sigma_w^2, \sigma_w^2).$$

Observation:

$$\mathbf{z}_{t,n} = \mathbf{H}_{t,n} \mathbf{x}_{i,t} + \mathbf{v}_t, \quad \mathbf{v}_t \sim \mathcal{N}(\mathbf{o}, \mathbf{R}_t),$$

where

$$\mathbf{H}_{t,n} = \begin{bmatrix} h_{t,n}(0) & h_{t,n}(1) & \dots & h_{t,n}(K) & 0 \\ 0 & 0 & \dots & 0 & 1 \end{bmatrix}^T,$$

and

$$\mathbf{R}_t = \text{diag}(\sigma_v^2, \dots, \sigma_v^2).$$

Prediction step:

$$\hat{\mathbf{x}}_{i,t|t-1} = \hat{\mathbf{x}}_{i,t-1|t-1},$$

$$\mathbf{P}_{i,t|t-1} = \mathbf{P}_{i,t-1|t-1} + \mathbf{Q}_t.$$

Update step: (performed if an observation exists)

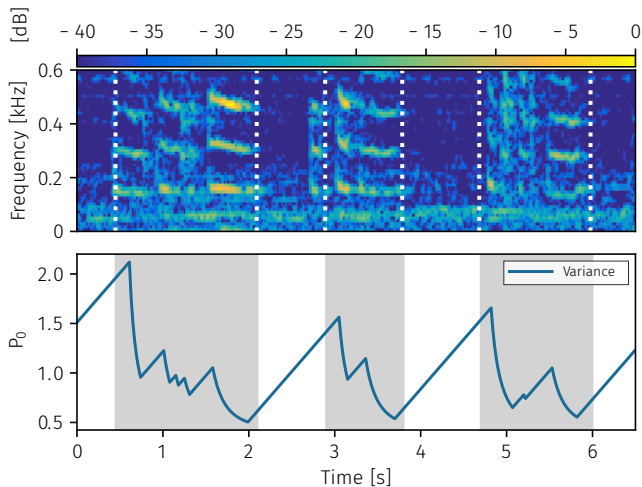
$$\hat{\mathbf{x}}_{i,t|t} = \hat{\mathbf{x}}_{i,t|t-1} + \mathbf{k}_{i,t}(\mathbf{z}_{t,n} - \mathbf{H}_{t,n}\hat{\mathbf{x}}_{i,t|t-1}),$$
$$\mathbf{P}_{i,t|t} = (\mathbf{I} - \mathbf{k}_{i,t}\mathbf{H}_{t,n})^2 \mathbf{P}_{i,t|t-1} + \mathbf{k}_{i,t}\mathbf{R}_t\mathbf{k}_{i,t}^T.$$

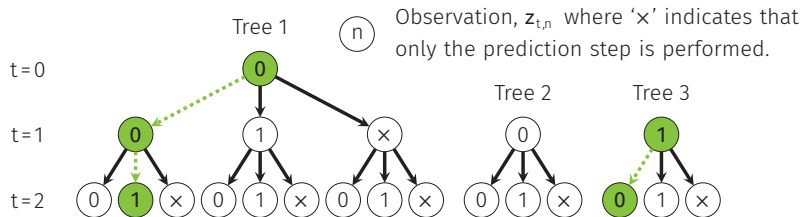
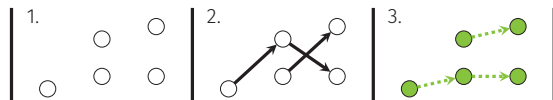
If $\mathbf{k}_{i,t} = [1, \dots, 1]$ $\implies \hat{\mathbf{x}}_{i,t|t} = \mathbf{z}_{t,n}$ (just the measurement)

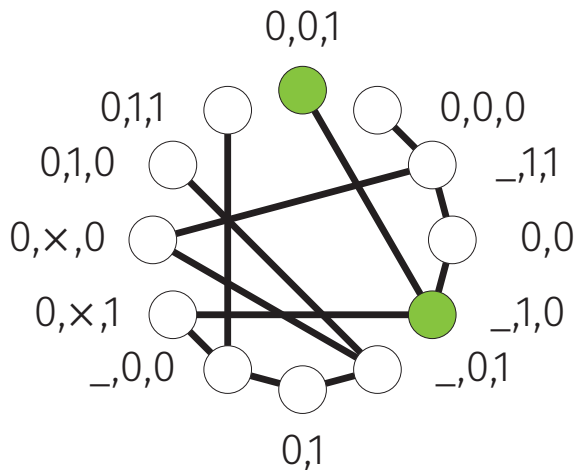
If $\mathbf{k}_{i,t} = [0, \dots, 0]$ $\implies \hat{\mathbf{x}}_{i,t|t} = \hat{\mathbf{x}}_{i,t|t-1}$ (just the prediction)

Optimal Kalman gain:

$$\mathbf{k}_{i,t} = \frac{\mathbf{P}_{i,t|t-1}\mathbf{H}_{t,n}^T}{\mathbf{S}_{i,t}}$$



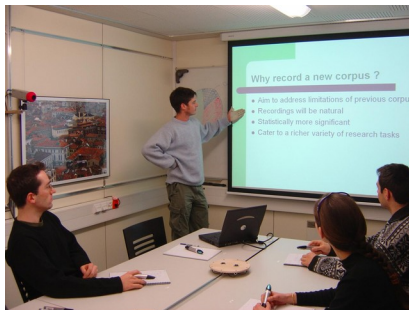


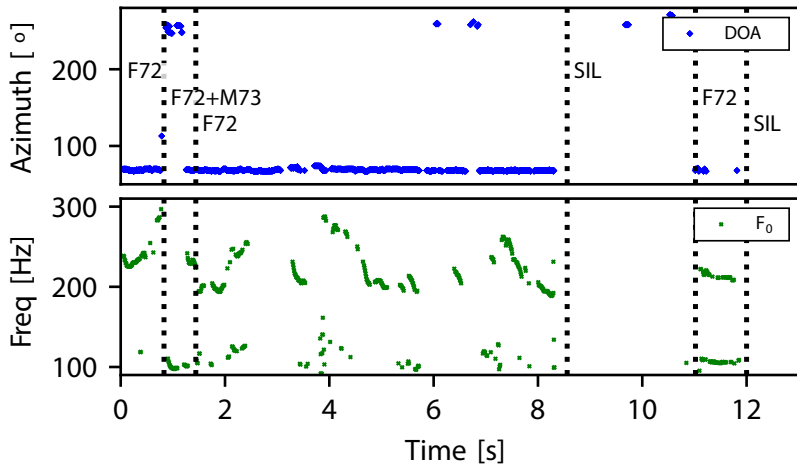


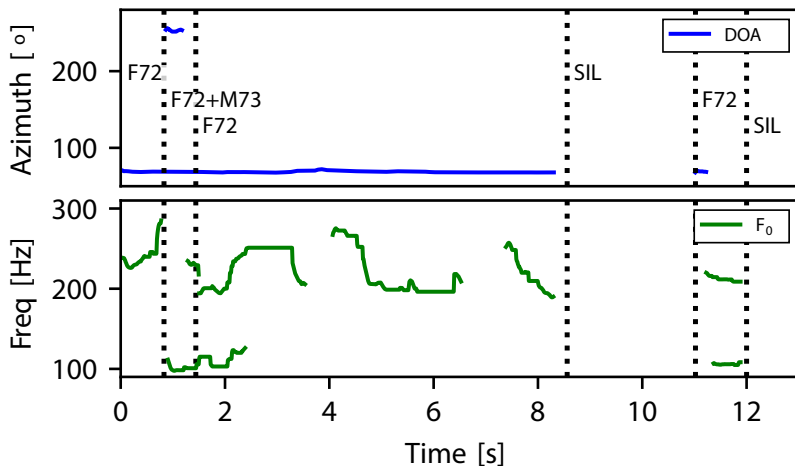
ILLUSTRATIVE EXAMPLE

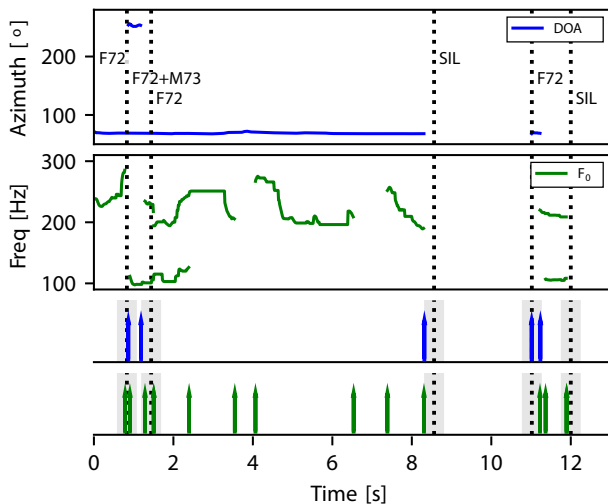
Multi-modal data set consisting of 100 hours of meeting recordings.

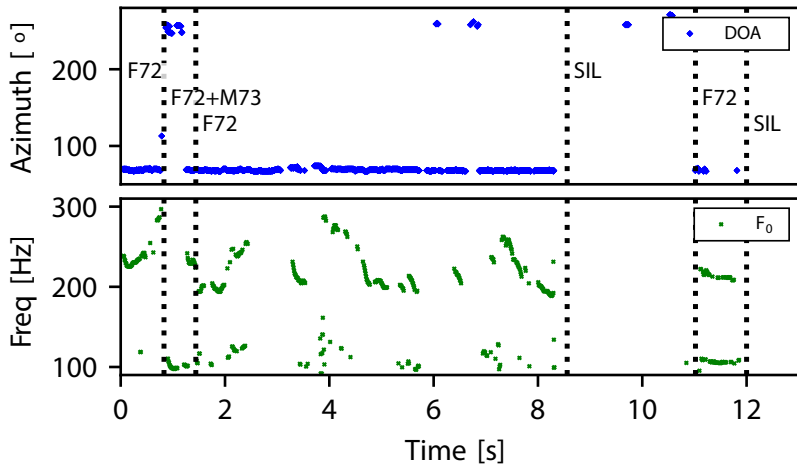
Recorded in English using three different rooms with different acoustic properties and includes mostly non-native speakers.

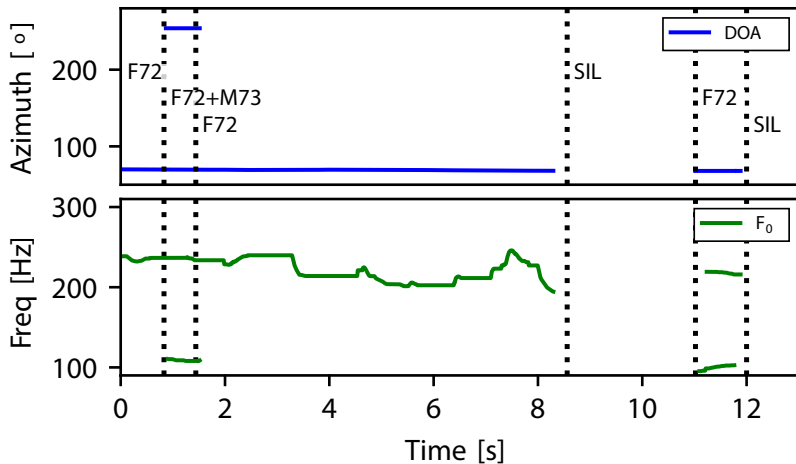


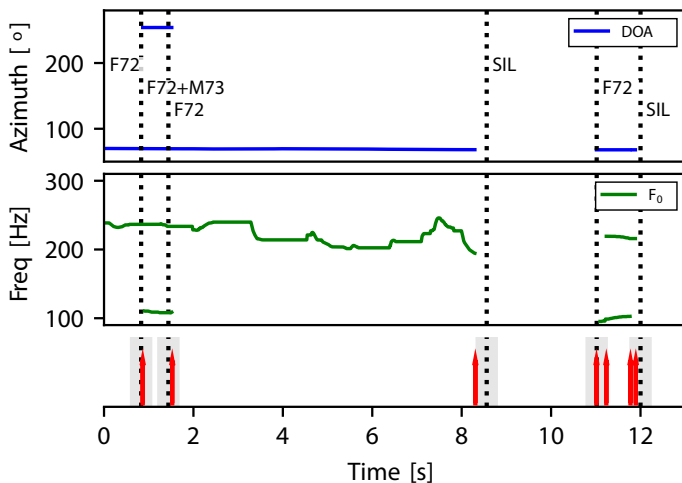




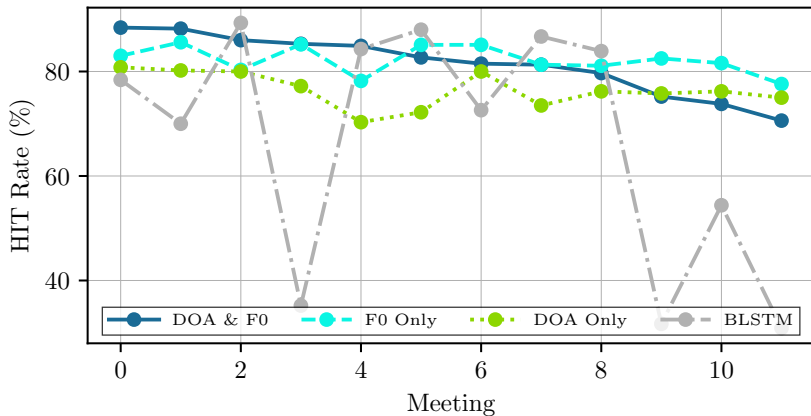


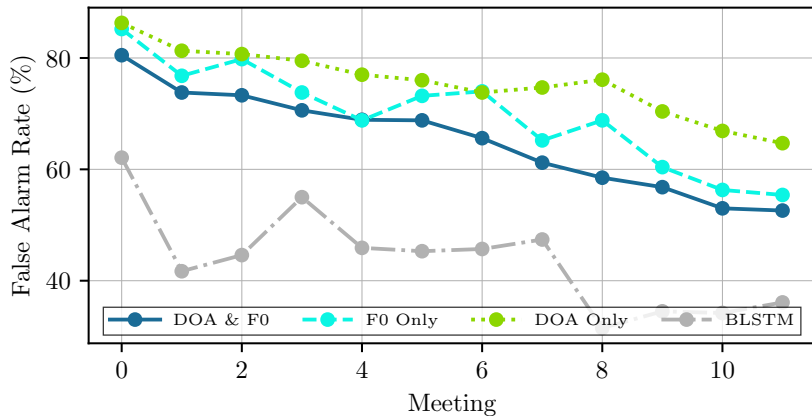


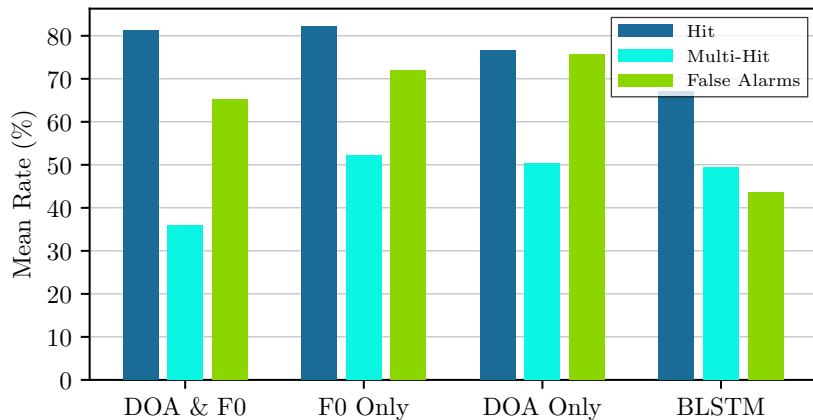




RESULTS







CONCLUSION

In this paper we have...

...proposed a novel method that uses a MHT framework to track the F0 and DOA of multiple speakers simultaneously.

...shown that MHT of both the DOA and F0 can lead to an improved speaker segmentation performance over tracking just one of these features alone.

QUESTIONS?

PLEASE EMAIL: AIDAN.HOGG13@IMPERIAL.AC.UK

Imperial College
London

UNIVERSITY OF
Southampton

EPSRC
Engineering and Physical Sciences
Research Council