

1. Overview & Proposed system

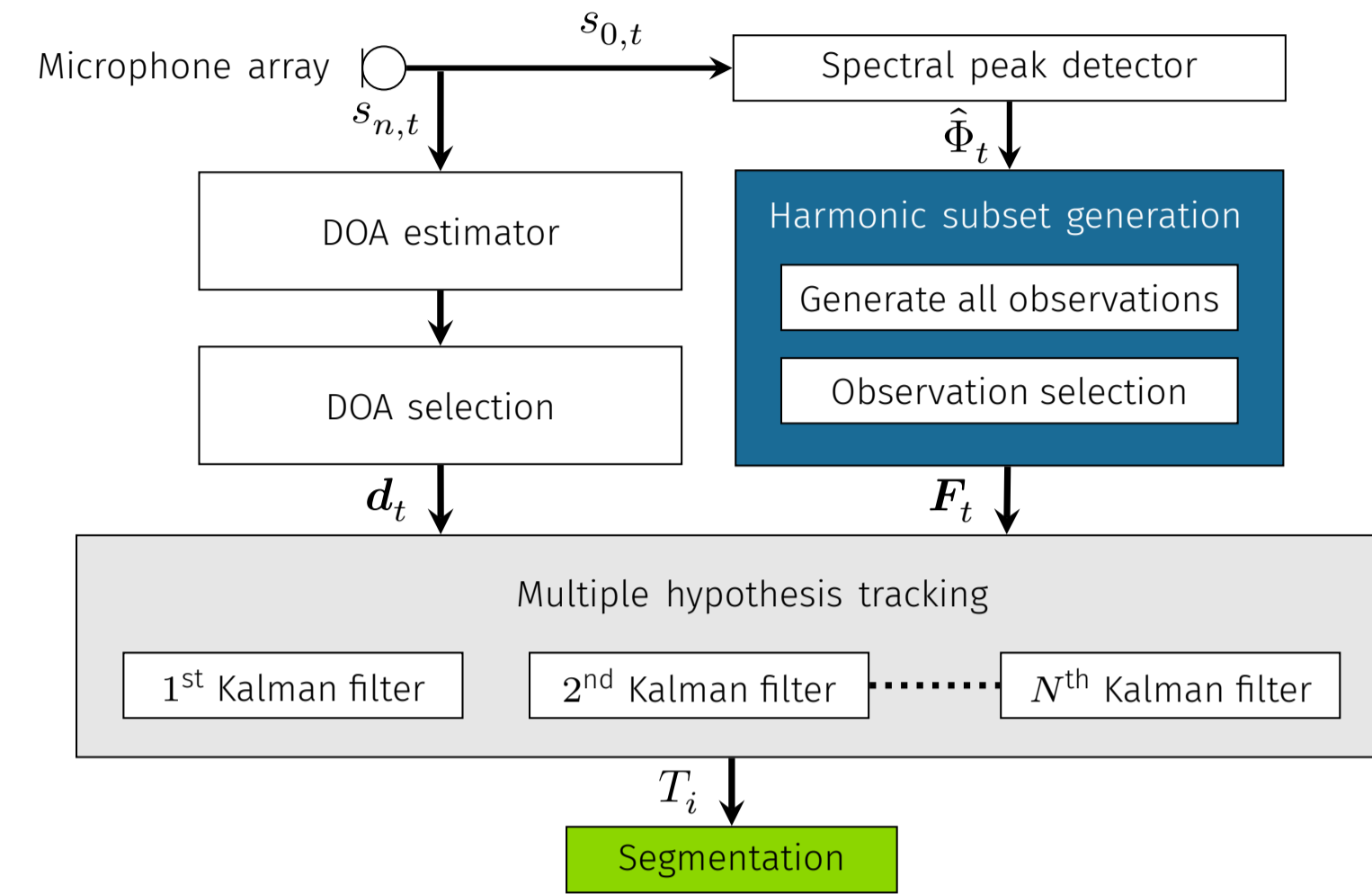
What does this paper propose?

A new multimodal approach for overlapping speaker segmentation that tracks simultaneously both the fundamental frequencies (F_0 s) and the direction of arrivals (DOAs) of multiple, simultaneously active speakers.

How well does this proposed method perform?

The proposed multimodal method shows an improvement in segmentation performance compared to tracking features separately.

Segmentation performance comparable to a deep learning approach is achieved but without the need for labelled training data and handles overlapping talkers.



2. Motivation

What is speaker diarization?

Answers the question "who spoke when?" and is required for applications such as, speaker indexing and automatic speech recognition (ASR).

Is performing segmentation before clustering beneficial?

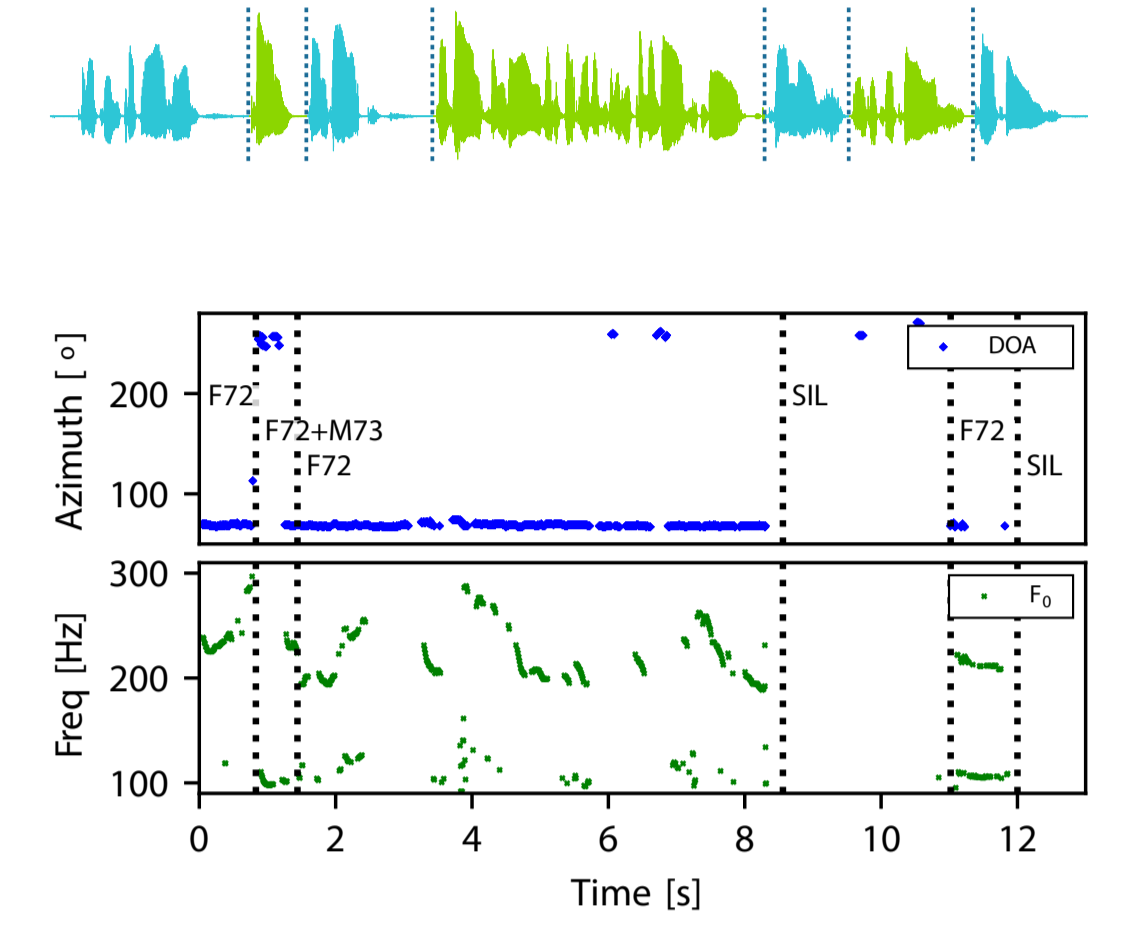
If segmentation is performed before clustering then each segment will contain the maximum amount of information possible on a speaker's identity.

Why use direction of arrival for segmentation?

It has been shown in the past that spatial features can help improve speaker segmentation even in the context of overlapping speech.

Why use pitch for segmentation?

A previous study of meetings in the AMI corpus has shown that abrupt variations in voice pitch estimates are indicative of speaker changes.



3. Kalman Filter For DOA & F_0 Tracking

All the harmonics of voiced speech are tracked along with the DOA estimates so that overlapping speech can be processed.

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:

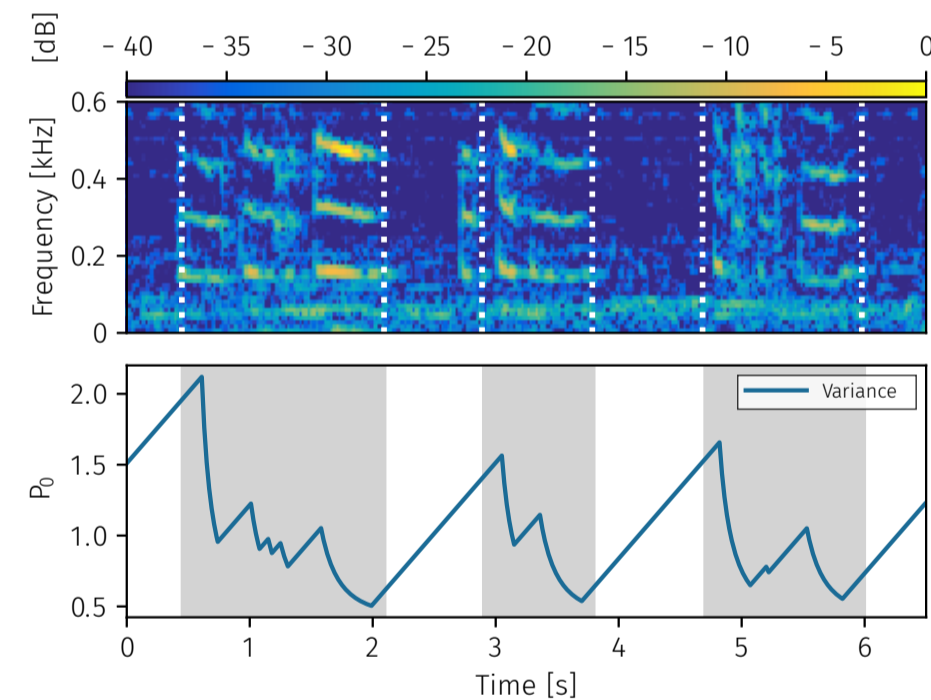
$$\mathbf{x}_{i,t} = \mathbf{x}_{i,t-1} + \mathbf{w}_t, \quad \mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_t),$$

where

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

with observation:

$$\mathbf{z}_{t,n} = \mathbf{H}_{t,n} \mathbf{x}_{i,t} + \mathbf{v}_t, \quad \mathbf{v}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_t), \quad \text{where } \mathbf{H}_{t,n} = \begin{bmatrix} h_{t,n}(0) & h_{t,n}(1) & \dots & h_{t,n}(K) & 0 \\ 0 & 0 & \dots & 0 & 1 \end{bmatrix}^T,$$



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

\mathbf{R}_t is the covariance of the observation noise.

Prediction step:

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

\mathbf{Q}_t is the covariance of the process noise.

Update step:

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

Optimal Kalman gain:

$$\mathbf{k}_{i,t} = \mathbf{P}_{i,t|t-1} \mathbf{H}_{t,n}^T \mathbf{S}_{i,t}^{-1}, \quad \text{where } \mathbf{S}_{i,t} = \mathbf{H}_{t,n} \mathbf{P}_{i,t|t-1} \mathbf{H}_{t,n}^T + \mathbf{R}_t.$$

Estimation error:

$$\mathbf{e}_{i,t|t} = \mathbf{z}_{t,n} - \mathbf{H}_{t,n} \hat{\mathbf{x}}_{i,t|t}.$$

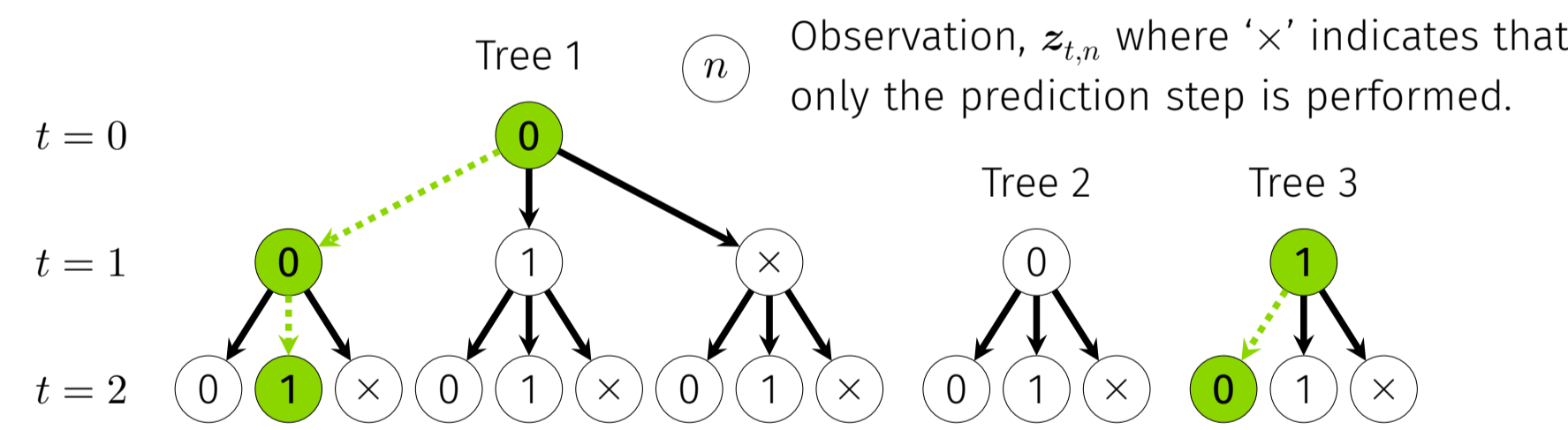
4. Multiple Hypothesis Tracking

A MHT framework is exploited to simultaneously track both the F_0 and DOA features.

Tracking observations

- Observations in time
- Possible track
- Best track

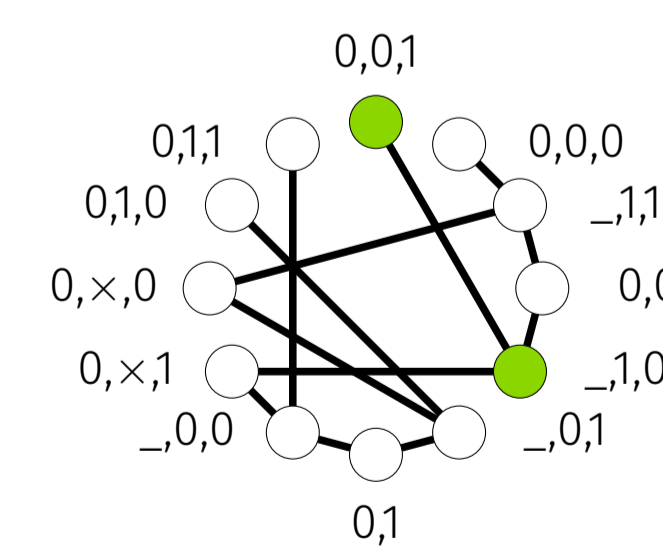
Multiple hypothesis tracking



Maximum weighted clique (MWC)

Used to find the most likely set of tracks which do not conflict.

Each node is a track hypothesis and each edge connects 2 tracks which do not conflict. A score is assigned which is calculated by taking the average value of all previous estimation errors.



Track generation

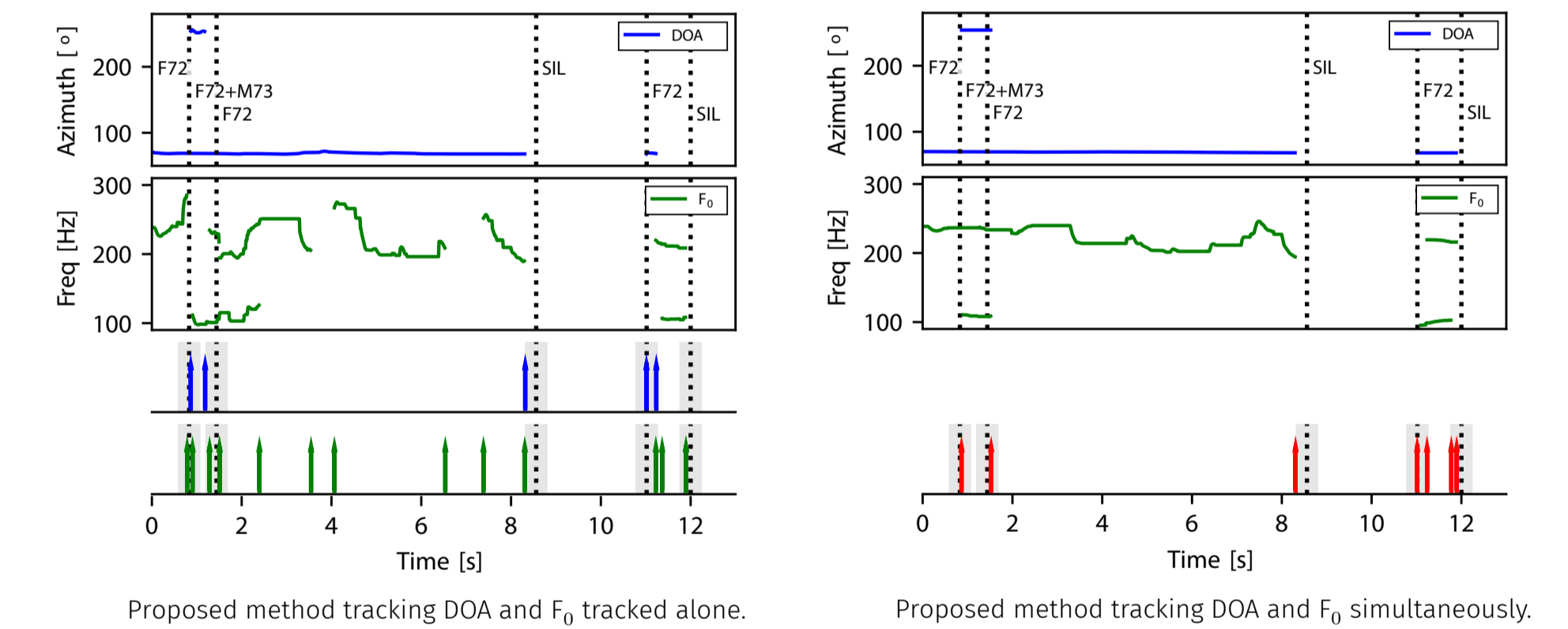
At each time-frame three possible tracks can be generated: a track only containing the F_0 observation; a track only containing the DOA observation and a track that fuses information from the F_0 and DOA observations.

The prediction step is executed for every time-frame. However, the update step is only performed when new observations emerge.

5. Results & Conclusion

Illustrative (AMI) example

Part of a meeting from the AMI corpus that contains overlapping speech.



Evaluation on AMI corpus

Performance comparison of the proposed method on the AMI corpus compared against the performance achieved by only using DOA or F_0 features alone as well as a bidirectional long short term memory networks (BLSTM) approach [2].

Method	Hit	Miss	Multi-Hit	FA
Proposed	81.2%	18.8%	36.0%	65.3%
F_0 Only	82.2%	17.8%	52.3%	72.0%
DOA Only	76.5%	23.6%	50.4%	75.6%
BLSTM	67.1%	32.9%	49.4%	43.7%

Mean values across 12 meetings from the AMI corpus.

[2] H. Bredin et al, "pyannote.audio: neural building blocks for speaker diarization," in *IEEE Int. Conf. on Acoust., Speech, and Signal Process. (ICASSP)*, 2020.

Conclusion

- A novel method has been proposed that uses a MHT framework to track the F_0 and DOA of multiple speakers simultaneously.
- MHT of both the DOA and F_0 can lead to an improved speaker segmentation performance on the AMI corpus over tracking just one of these features alone.