# OVERLAP-AWARE DIARIZATION:
## RESEGMENTATION USING NEURAL END-TO-END OVERLAPPED SPEECH DETECTION

### Latané Bullock

Undergraduate
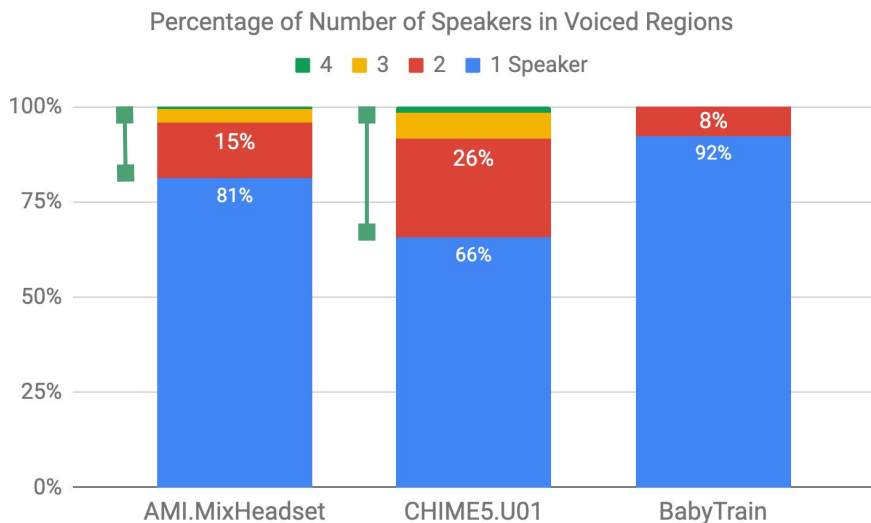Rice University

### Hervé Bredin

Researcher
LIMSI, CNRS France

### Leibny Paola Garcia-Perera

Assistant Research Scientist
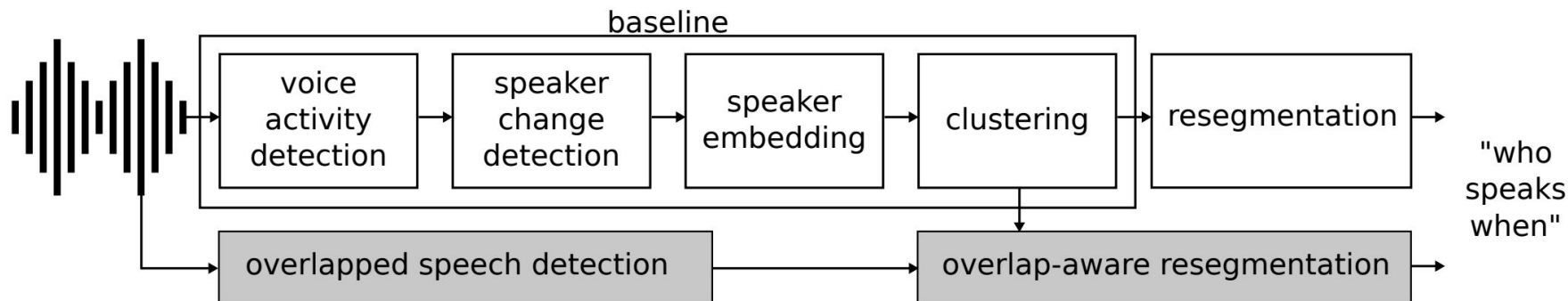CLSP, Johns Hopkins University

ICASSP2020
Barcelona

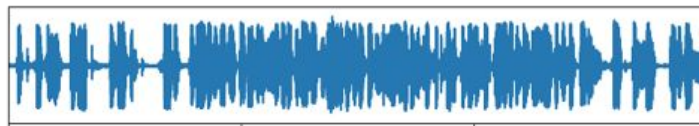# Diarization - "Who spoke when?"



In adverse audio recordings

- Large proportion of overlapped speech
  across all datasets

- Leads to high
  missed detection rate

- May lead to high
  speaker confusion rate

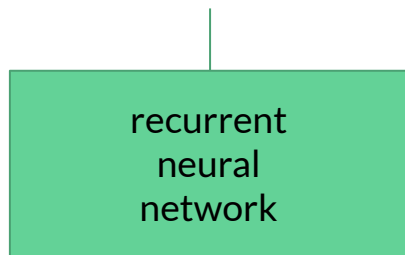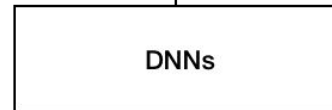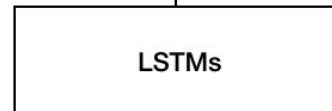# Contributions: Overlap Detection and Overlap-Aware Resegmentation
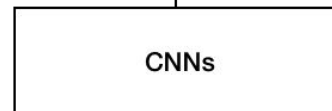
# Overlap Detection with

py**annote**

2 seconds

recurrent
neural
network

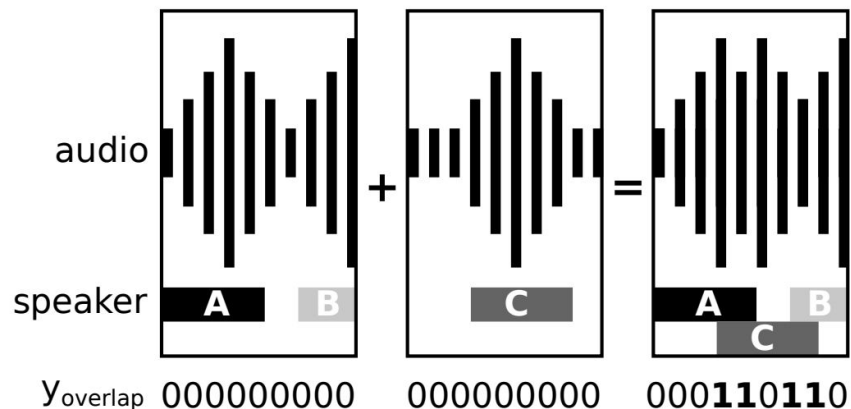1 = overlap
0 = no overlap

SincNet

CNNs

LSTMs

DNNs

Framewise scores

# Data augmentation

- Training a network directly does not work
  - Class imbalance
  - Lack of variability

- Two types of training samples
  - **50% regular**
    2-second chunks extracted from the training set randomly
  - **50% made-up**
    weighted sum of two chunks



audio

speaker  A  B  +  C  =  A  B  C

$y_{overlap}$  000000000   000000000   000**110110**0

To increase the number of positive training samples for overlapped speech detection, artificial audio chunks are created by summing two random audio chunks

# At test time...



**threshold**

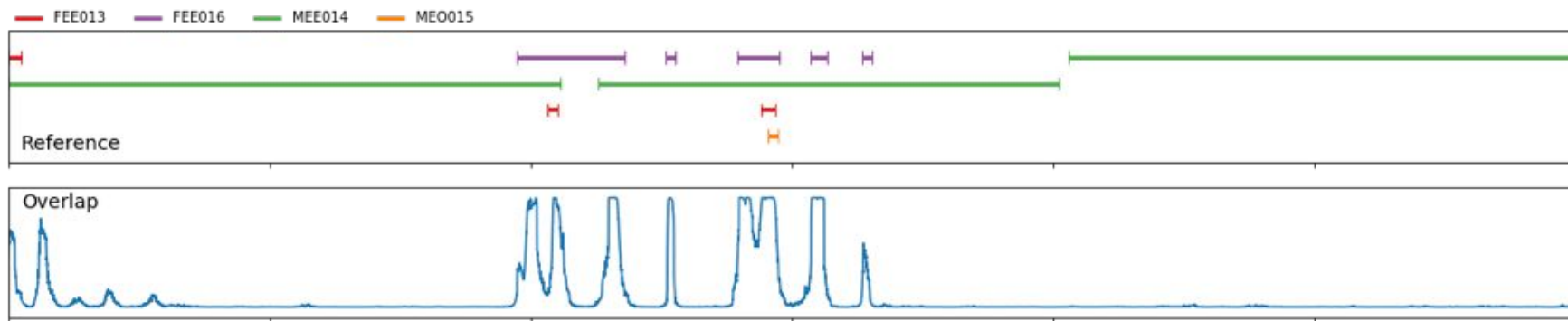# Overlap Detection: Results



| | AMI | | DIHARD | | ETAPE | |
|---|---|---|---|---|---|---|
| | Precision | Recall | Precision | Recall | Precision | Recall |
| Baseline | 75.8 80.5 [8] | 44.6 50.2 [8] | | | 60.3 [20] | 52.7 [20] |
| Proposed (MFCC) | 91.9 90.0 | 48.4 52.5 | 58.0 73.8 | 17.6 14.0 | 67.1 55.0 | 57.3 55.3 |
| Proposed (waveform) | 86.8 90.0 | 65.8 63.8 | 64.5 75.3 | 26.7 24.4 | 69.6 60.0 | 61.7 63.6 |

# Overlap-Aware Resegmentation

# Variational Bayes HMM-GMM Resegmentation

Diez et al, 2018
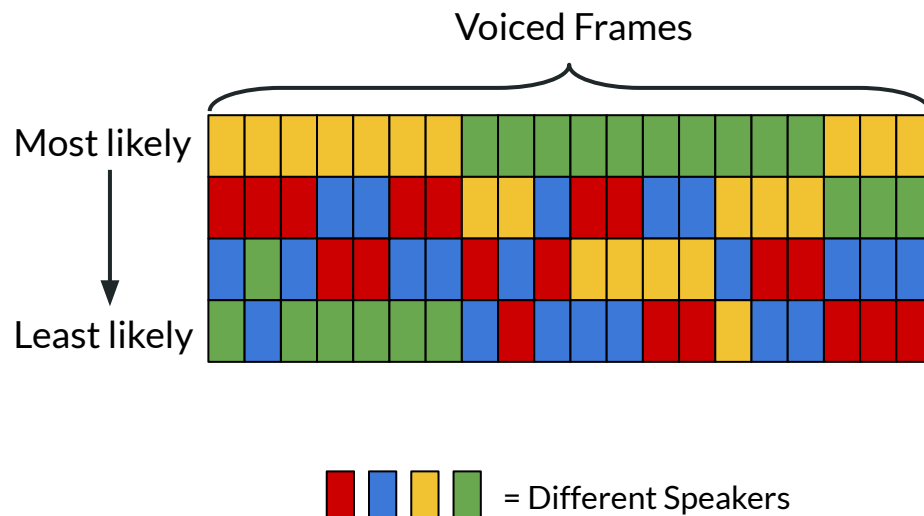


Hidden Markov model where:
- state represents a speaker
- state distributions are GMMs constrained by eigenvoice priors
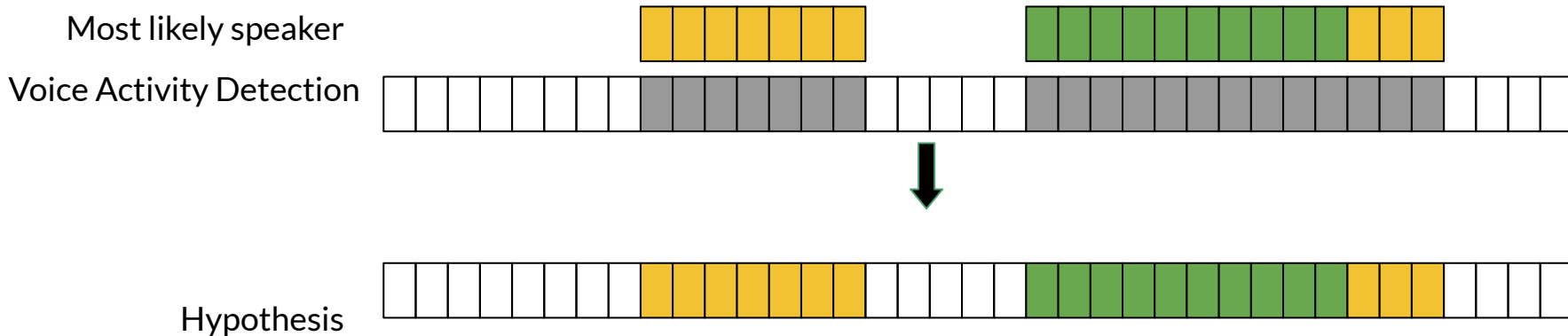
Single model (theoretically) infers:
- Speaker distributions
- **Number of speakers**
- **Speaker sequence**

… with the Variational Bayes

Remnant of VB-HMM resegmentation, the **speaker attribution matrix**:

Voiced Frames

Most likely

Least likely

= Different Speakers

# Two-speaker Assignment in Overlap Regions

Most likely speaker

Voice Activity Detection

Hypothesis

# Two-speaker Assignment in Overlap Regions

Most likely speaker

Voice Activity Detection

Hypothesis

Overlap Detection

# Two-speaker Assignment in Overlap Regions

Most likely speaker

Voice Activity Detection

Hypothesis

Overlap Detection

2nd most likely speaker

# Two-speaker Assignment in Overlap Regions



Most likely speaker

Voice Activity Detection

**Final** Hypothesis

Overlap Detection

2nd most likely speaker

# Results: Overlap Assignment on AMI Headset Mix



Diarization Error Rates with Resegmentation and Overlap Assignment

■ confusion  ■ false alarm  ■ missed detection

Baseline: 29.6% (confusion 5.8%, false alarm 3.0%, missed detection 20.8%)

+Overlap Assignment: 23.8% (confusion 7.2%, false alarm 3.6%, missed detection 13.0%)

+Oracle Detection: 22.3% (confusion 13.2%, false alarm 3.1%, missed detection 6.0%)

+Oracle Assignment: 11.8% (confusion 0.0%, false alarm 0.6%, missed detection 11.2%)

# Conclusions and Takeaways

Overlap detector

- State-of-the-art performance on AMI and ETAPE, sets standard for future comparison on DIHARD II
- Primary gains from decreased missed detection

Overlap-aware Resegmentation

- Results in large decreases in DER
- BUT at the cost of increases in confusion error

# Acknowledgements