

# DNN Approach to Speaker Diarisation Using Speaker Channels

Rosanna Milner and Thomas Hain

Machine Intelligence for Natural Interfaces  
Speech and Hearing  
University of Sheffield



The  
University  
Of  
Sheffield.



# Outline

- Introduction
- Background
- DNN approach using speaker channels
  - Fixed or mixed number of channels
  - System overview
- Experiments
  - Test Data
  - Setup
  - Evaluation
  - Results
- Conclusion

# Introduction

Speaker diarisation - 'who speaks when'

- the 3 main tasks are SAD, speaker segmentation and speaker clustering
- **step-by-step**: performs stages separately
- **integrated**: performs some stages together

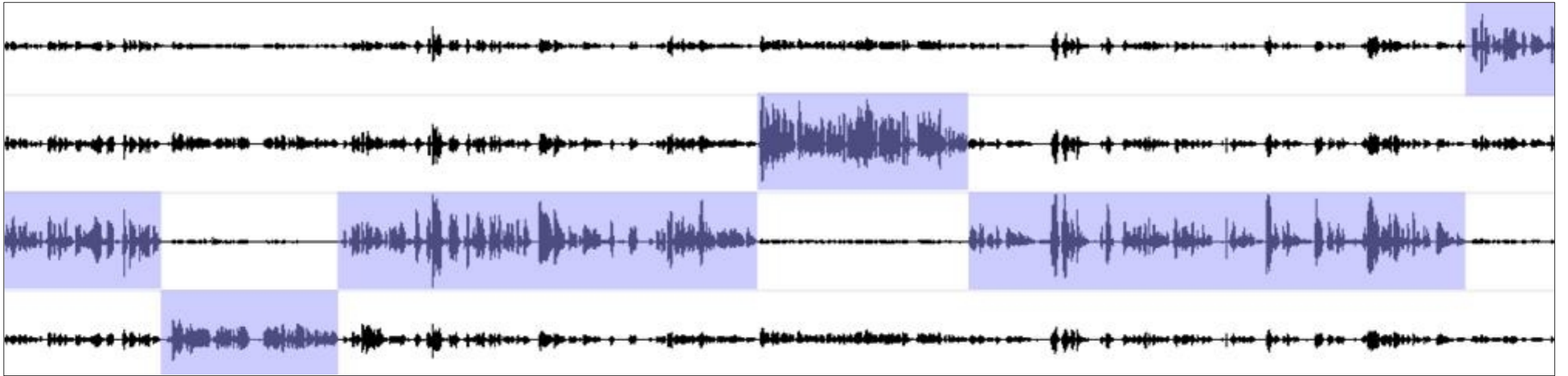
Typically unsupervised

- **unsupervised**: no prior knowledge or information
- **lightly/semi-supervised**: auxiliary information or metadata available
- **supervised**: prior knowledge about test data known

Presenting a semi-supervised integrated method using DNNs trained on concatenated IHM features

- semi-supervised: uses IHM speaker channels (instead of SDM)
- integrated: performs all three tasks together

# Multi-channel Diarisation Approaches



- Single channel (SDM)
  - Segmentation, change detection and clustering
- How close are the channels to the speaker ?
  - Associated speaker channels (IHM++)
  - Distant speaker channels (MDM)

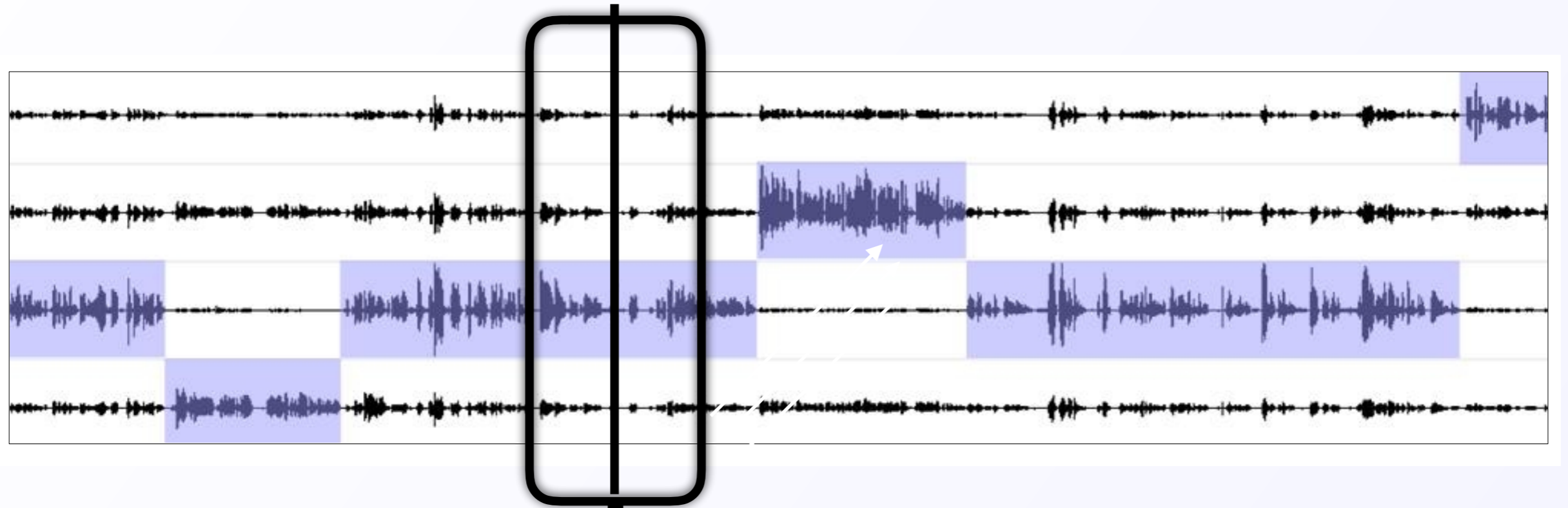
# Scoring Diarisation Output

- Diarisation Error Rate (DER)
  - Frame based metric
  - Collar changes reference
  - No penalty for data fragmentation
- Scoring
  - on individual channels - IHM
    - very low number due to speaker prior !
  - on one 'global channel' - SDM
  - based on true activity - MDM
- Alternative scoring (Milner & Hain, ICASSP'16)

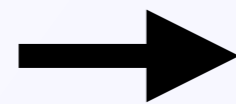
# Previous Work

- **Multichannel diarisation**
  - beamforming focuses on speakers (Anguera et al. 2007)
  - detecting closest speech and disregarding other speech (Dines et al. 2006, Wrigley et al. 2005)
- **DNNs for diarisation**
  - feature transforms using ANNs (Yella et al. 2014, 2015)
  - DNNs trained for SAD (Dines et al., 2006, Milner & Hain, 2015)
  - windowing segmentation method and clustering using AANNs (Jothilakshmi et al. 2009)
  - Clustering by adapting speaker separation DNNs to specific recordings (Milner & Hain, 2016)

# Approach



Channel 2



Rosanna



# Approach: Using speaker assigned channels

## Fixed number of channels

- every recording must contain the same number of speaker channels,  $x$  concatenate  $x$  channels in all permutations:  $x!$  features per recording
- final layer in DNN is  $x+1$ , representing  $x$  speakers and NONSPEECH

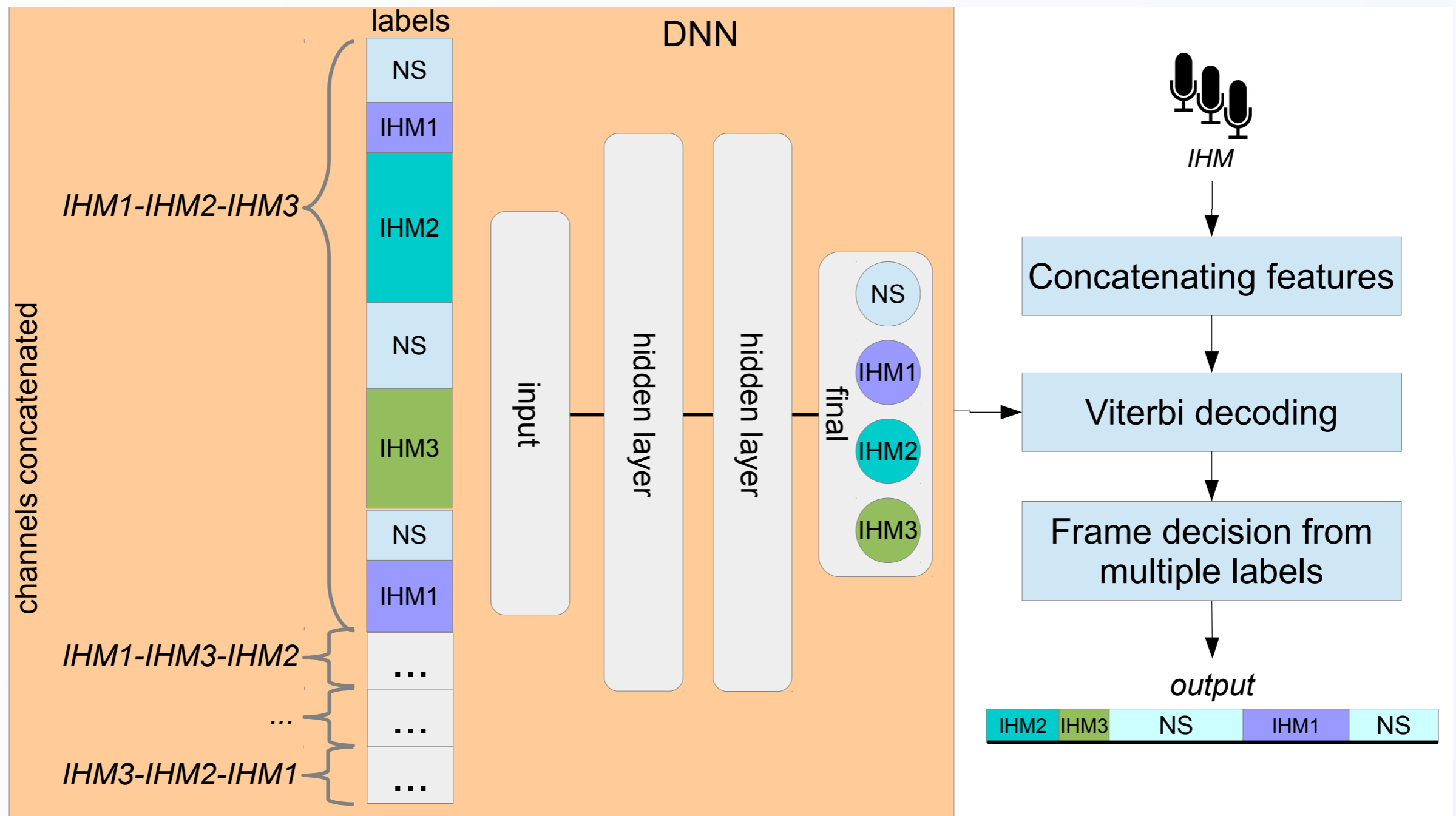
## Mixed number of channels

- recordings contain different number of speaker channels
- concatenate all pairs of channels:  $x(x-1)$  features per recording
- final layer in DNN is 3, representing 2 speakers and NONSPEECH
- a speaker in label file but not in channel pair has NONSPEECH label

*Both methods require every speaker having their own channel*



# Approach



# Frame Decisions

## Combinatorial Voting

FRAME	<i>IHM1-IHM2-IHM3</i>	<i>IHM1-IHM3-IHM2</i>	...	<i>IHM3-IHM2-IHM1</i>	OUTPUT
...	...	...		...	...
204	IHM1	IHM1		IHM1	IHM1
205	IHM1	IHM1		IHM1	IHM1
206	IHM1	IHM1		IHM1	IHM1
207	NS	NS	...	NS	NS
208	NS	NS		NS	NS
209	IHM3	IHM3		IHM2	IHM3
210	IHM2	IHM2		IHM3	IHM2
...	...	...		...	...

- All combinations of feature concatenations used for testing
- results in multiple labels for every frame
- simply count occurrences and choose label which occurs most often
- additionally: apply a prior for NONSPEECH

# Data - Meetings

- NIST RT'07 - meeting data
  - NIST reference and manually transcribed reference (0.1 sec precision)
  - 11144 segments, 35 speakers
  - 8 meetings
  - 6 meetings: 4 speakers, 1 meeting: 5 speakers, 1 meeting: 6 speakers



Improved manual reference on

<http://mini.dcs.shef.ac.uk/resources/dia-improvedrt07reference/>

# Data: Talk Show Radio 4

- The Bottom Line - BBC Radio4
- Topics in Economics
- 3 participants
- 1 interviewer (Evan Davis)



- manually transcribed reference
- 8749 segments, 40 speakers
- 12 train and 10 test programmes



# Features and configuration

- Features
  - Log filterbank (23 coefs, 32 frames, compressed)
  - Cross talk features (Wrigley et al, 2006) - normalised energy, kurtosis, mean/max cross correlation and differentials, 7 per channel
- DNN configurations
  - 2 hidden layers (1000 hidden units)
  - With cross talk features (31 frames)
  - trained with or without overlapping speech (OV) - unique labels - TBL only - overlap 7.5%

# Evaluation

## Diarisation error rate

- $DER = MS + FA + SE$
- does not consider the segmentation quality so all tables show the number of detected segments

## Two scoring settings

- NIST
  - collar 0.25s
  - score specified times only (UEM)
  - NIST provided reference (where possible)
- SHEF
  - collar 0.05s
  - score complete recordings
  - manually transcribed references

# Baseline results

- LIUM SpkrDiarization (Rouvier et al., 2013)
  - tailored for TV and radio broadcasts
  - BIC segmentation with CLR and integer linear programming and i-vector clustering

Channel	#Segs	#Spkrs	NIST DER%	SHEF DER%
Data: TBL				
SDM	2030	82	16.6	27.8
IHM	8478	40	393.9	335.9
Data: RT07				
SDM	2648	72	40.1	66.4
IHM	13070	35	308.1	371.0

Crosstalk on channels which results in high false alarm

Channel	#Segs	#Spkrs	NIST DER%	SHEF DER%
ICSI - SDM	3082	54	21.7	66.2



# Fixed Channel Experiments - TBL

- Only possible on TBL data
  - with or without overlap in training
  - with or without cross talk features

DNN			#Segs	MS%	FA%	SE%	SHEF DER%
Train	OV	CT					
TBL	x		6732	4.3	2.4	1.2	8.0
TBL	x	x	7136	4.3	2.4	1.7	8.4
TBL			7269	4.3	2.5	1.5	8.3
TBL		x	2964	4.6	3.7	1.4	9.7

**DNN TBL+OV gives lowest SHEF DER, crosstalk features do not help**

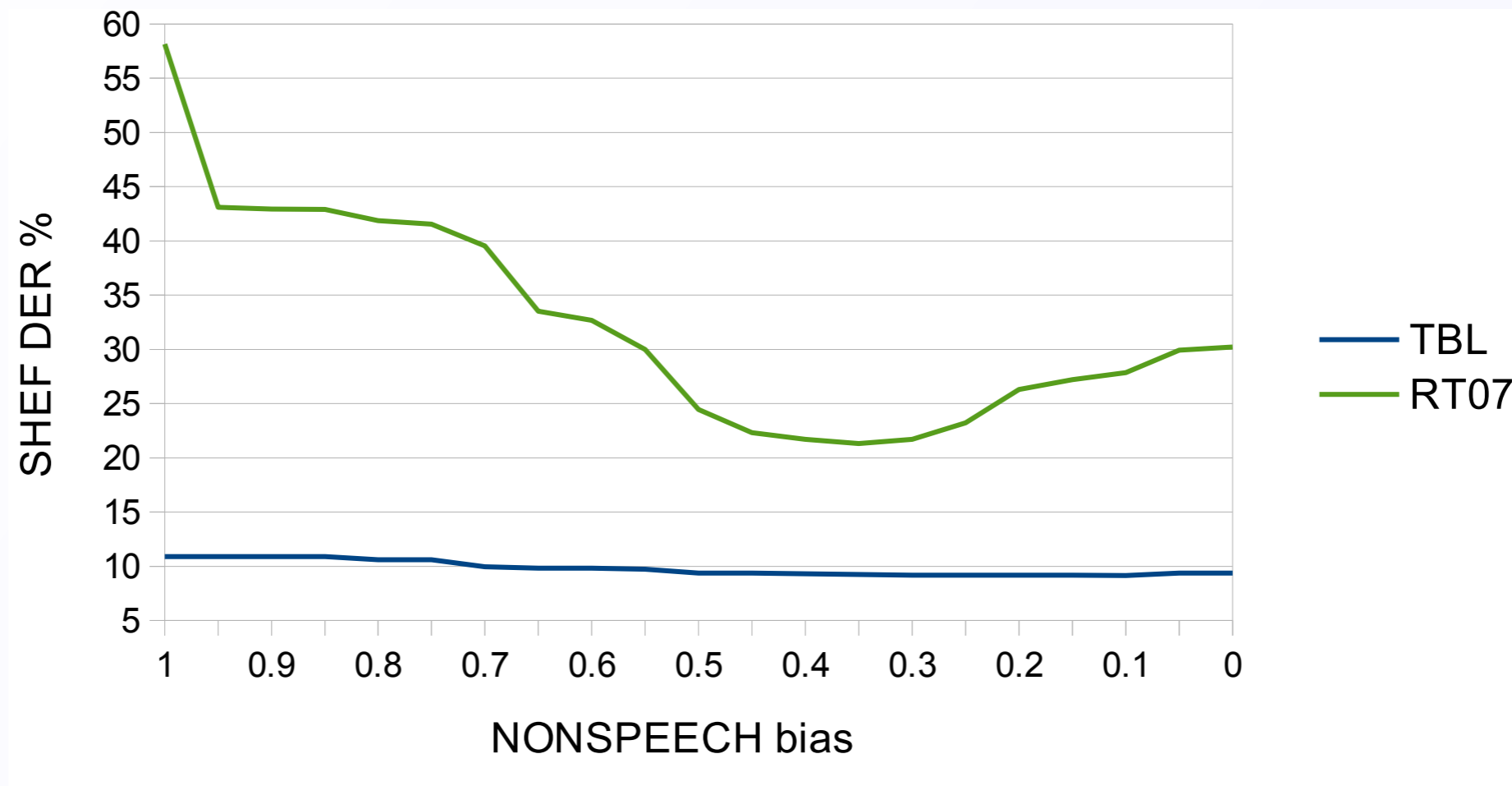
Weight	#Segs	MS%	FA%	SE%	SHEF DER%
0.75	6594	4.3	2.6	1.3	8.2
0.5	6571	4.2	2.7	1.3	8.2
0.25	6569	4.2	2.8	1.4	8.3

# Mixed Channel Numbers - I

Data	DNN			#Segs	MS%	FA%	SE%	SHEF DER%
	Train	OV	CT					
TBL	TBL	x		8295	20.3	1.1	0.9	22.4
	TBL	x	x	10551	34.8	0.7	1.1	36.5
	TBL			8263	17.0	1.4	1.0	19.4
	TBL		x	7932	7.7	0.9	1.2	10.9
	AMI			10354	16.6	1.0	4.9	22.5
	AMI		x	7683	22.9	0.9	5.0	28.8
RT07	TBL	x		7979	60.9	0.8	0.4	62.1
	TBL	x	x	4169	79.6	0.4	0.1	80.1
	TBL			8430	56.5	1.2	0.4	58.2
	TBL		x	5993	59.7	1.3	0.2	61.2
	AMI			8791	58.9	0.5	0.1	59.5
	AMI		x	6873	62.4	0.5	0.1	63.0

- crosstalk features only improve for DNN for TBL
- including overlap in DNN training gives worse performance
- DNNs trained on AMI data do not perform as well as DNNs trained on TBL data without overlap

# Mixed Channel Numbers - II



- applying a weight helps both datasets
- RT07 benefits the most with a large performance increase from 58.2% to 23% SHEF DER

# Best results

Data	SHEF DER%	NIST DER%
TBL	9.2	5.7
RT'07	23.2	15.1

# Conclusions

- presented two approaches for speaker diarisation using only IHM channels
- evaluated on two datasets: RT07 (meeting) and TBL (broadcast media)
- two methods for scoring: NIST and SHEF
- applying a nonspeech bias reduces error in mixed method
- training on OV benefits fixed method but not mixed
- CT only benefit DNN trained on TBL and tested on TBL
- Best result appears to be significantly better than best reported result on SDM

# The End

---

Thank you.