

# Diarization in a world full of DNNs (special tutorial)

Leibny Paola Garcia Perera



# Collaborators

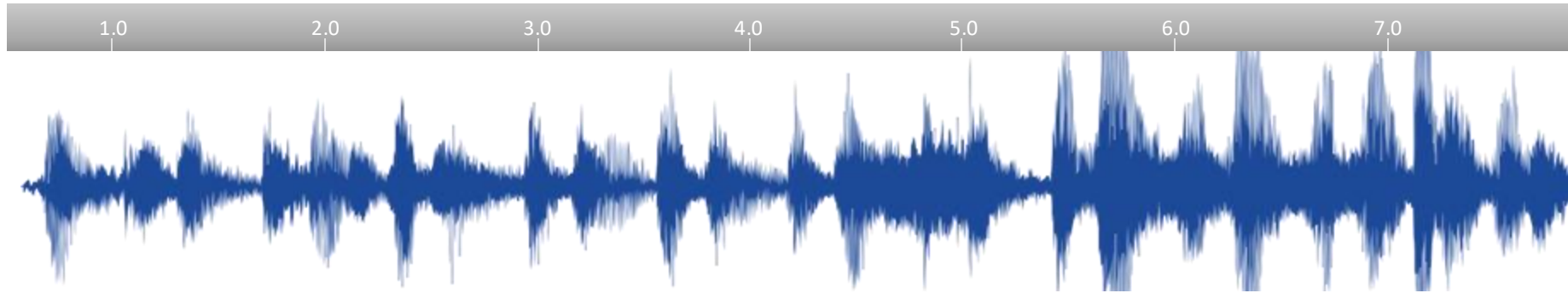
- Thanks to each one of them!
- JHU: Zili Huang, Desh Raj, Matthew Maciejewski, Jesus Villalba, Sanjeev Khudanpur
- NTU: Hexin Xie, Victoria Chua, Suzy Styles, Justin Dauwels
- CMU: Shinji Watanabe
- Some other collaborators: Latane Bullock, Herve Bredin, Marvin Lavenchin
- From Hitachi: Yusuke Fujita (now at Line), Shota Horiguchi, Yawen Xue, Yuki Takashima, Nelson Yalta
- CCWD, CDS project (lots of people behind this amazing project)

# Disclaimer

There a lot of fascinating works out there.

These are just a few ones in which somehow I have been involved.

# What is the goal?



SPK1

SPK1

SPK2

SPK2

SPK3

Spk1: If we want to address the next diarization problems..

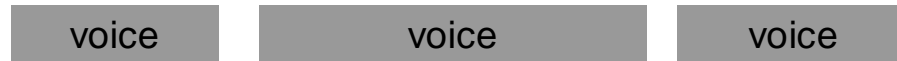
Spk2: You should need to first breakdown those results.

Spk3: Ok, I will put them in a Table or graph. I will also detail the algorithm...

# Where is diarization in the speech world?



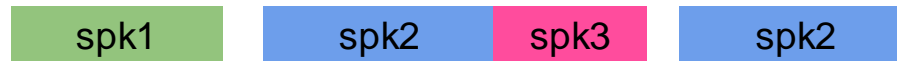
# Who spoke when?



voice activity detection

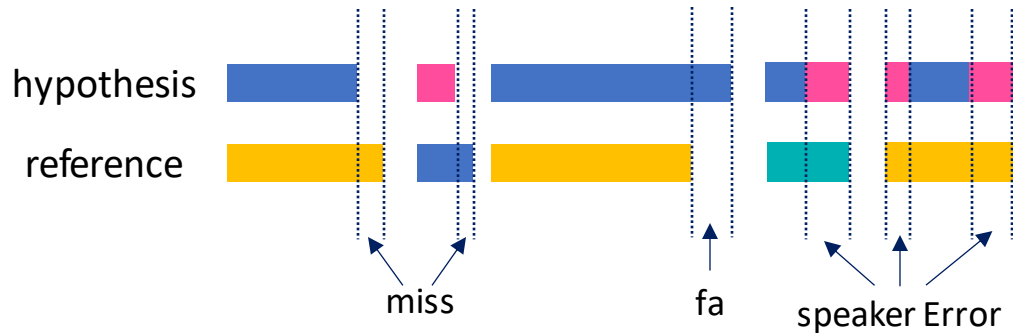


speaker type classification



speaker diarization

# How good the diarization is?



## Diarization Error Rate

$$\text{DER} = \frac{\text{false alarm} + \text{missed detection} + \text{speaker error}}{\text{total}}$$

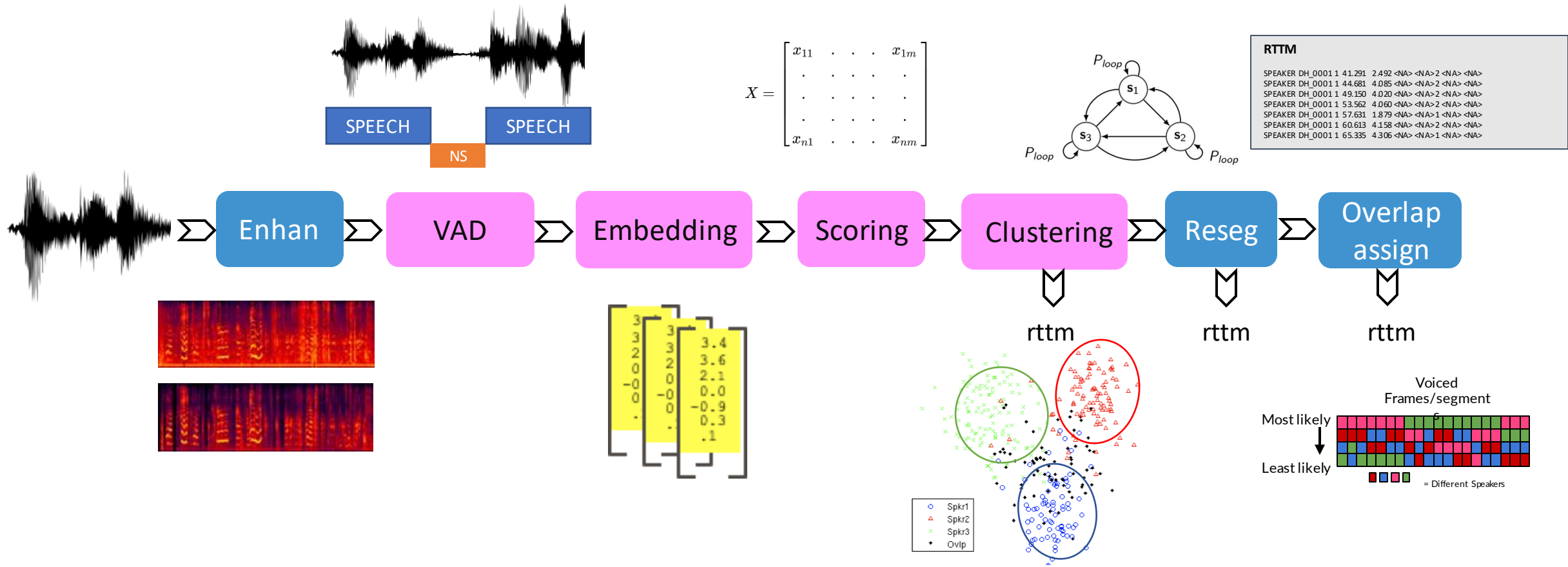
## Jaccard Error Rate

For each reference speaker  $s_i$ , where  $i = 1, \dots, N$

$$\text{JER}_{s_i} = \frac{\text{false alarm}_{s_i} + \text{missed detection}_{s_i}}{\text{total}_{s_i}}$$

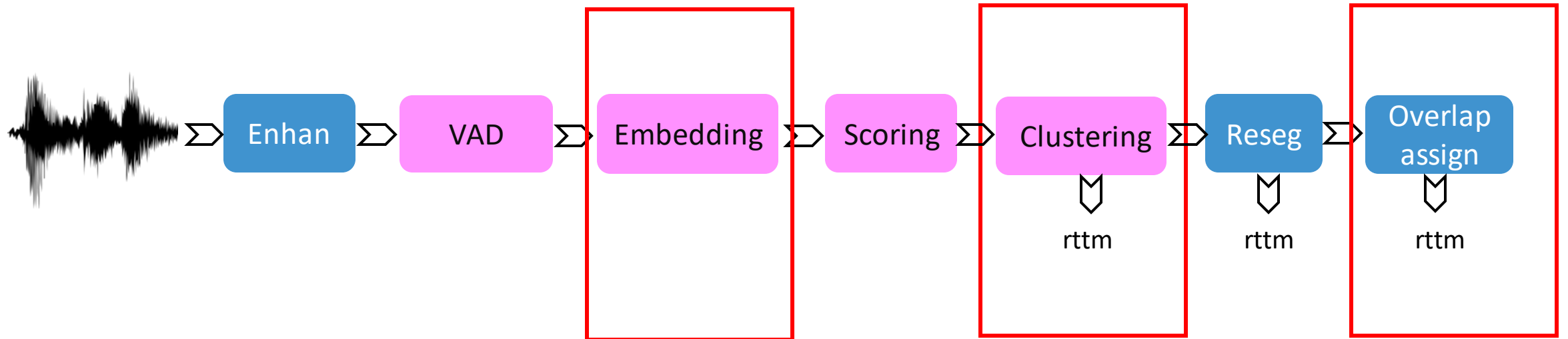
$$\text{JER} = \frac{1}{N} \sum_1^N \text{JER}_{s_i}$$

# Traditional diarization approach

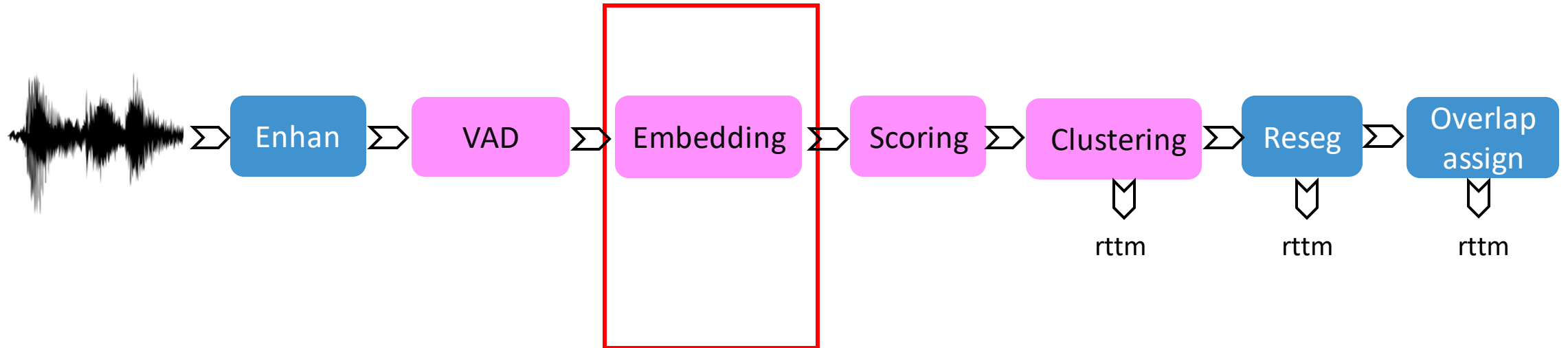




# Key ideas

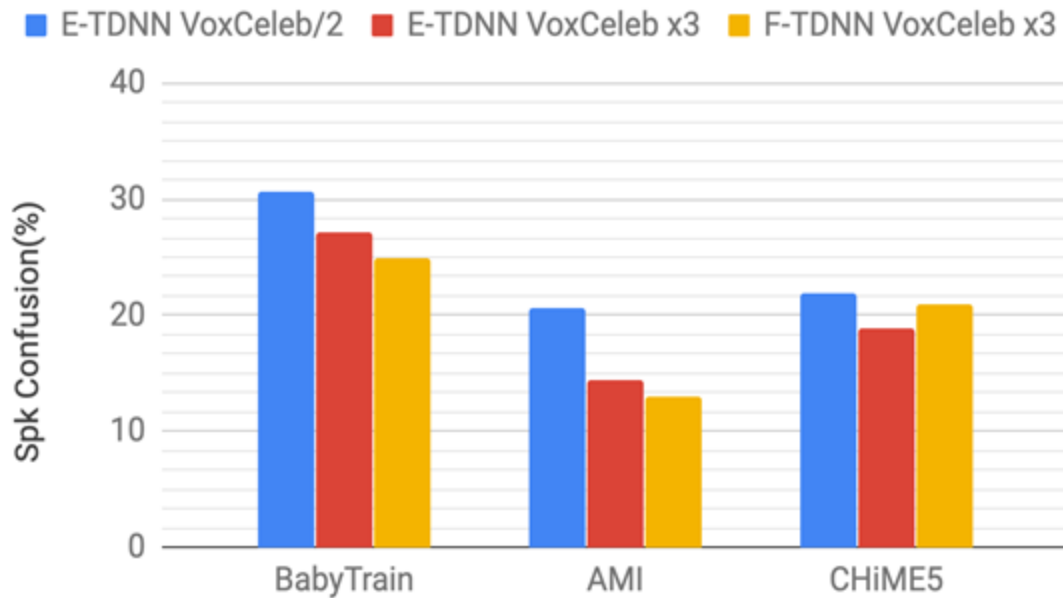


# Key ideas

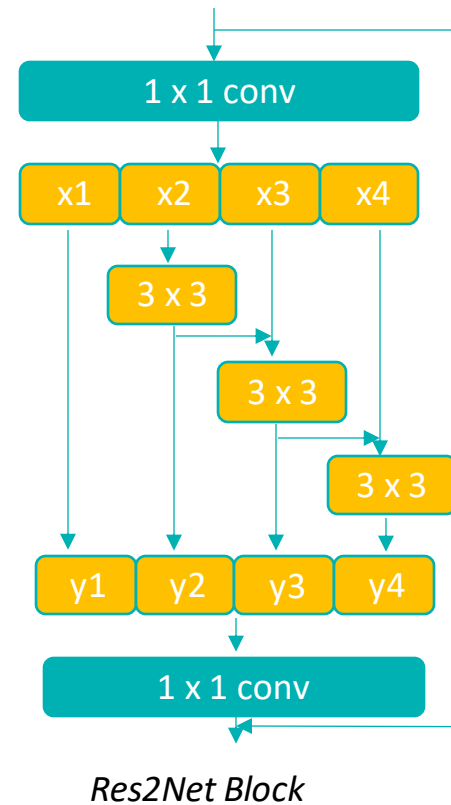


# Embeddings

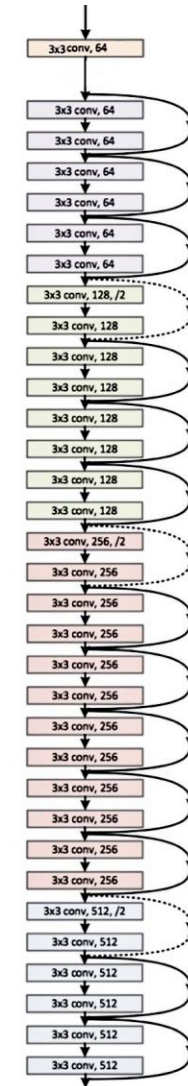
TDNN-xvector



Res2net

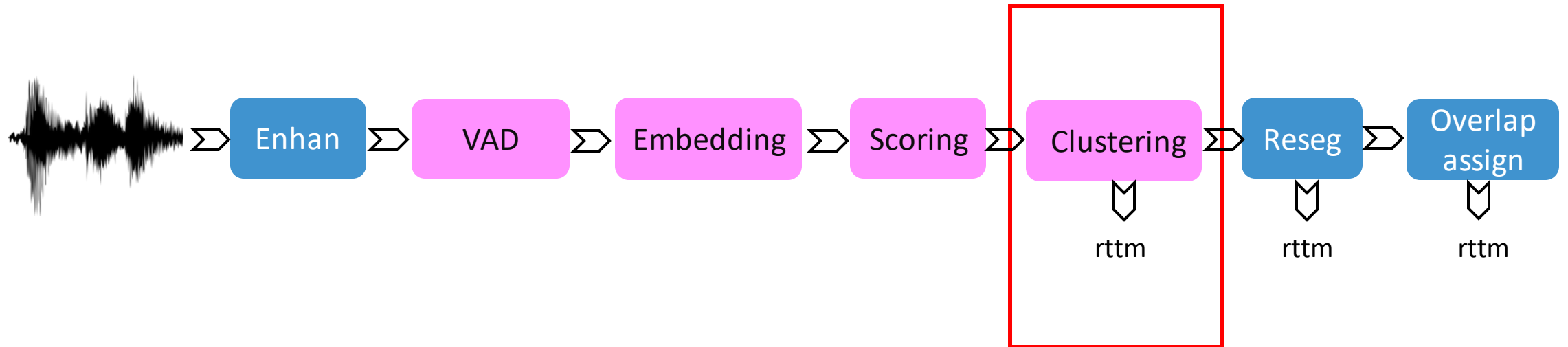


ResNet34



Jesus Villalba, et.al., at JSALT 2019 workshop  
 Shang-Hua Gao, et.al., Res2net: A new multi-scale backbone architecture  
 Xiong, Xiao, et. Al., Microsoft Speaker Diarization system for Voxceleb speaker recognition challenge 2020  
 Wang, Weiking, et.al., The DKU-DukeECE-Lenovo Diarization System for VoxSRC 2021

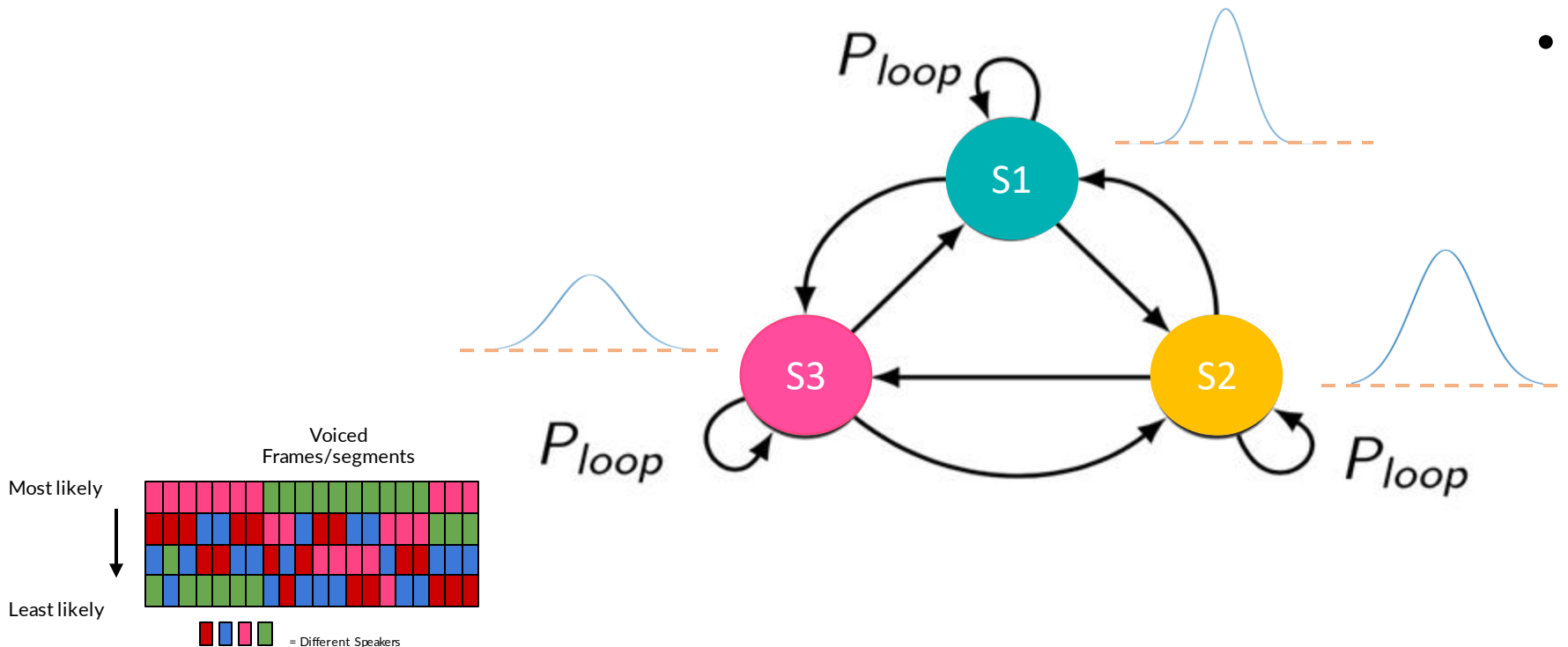
# Key ideas



# Clustering

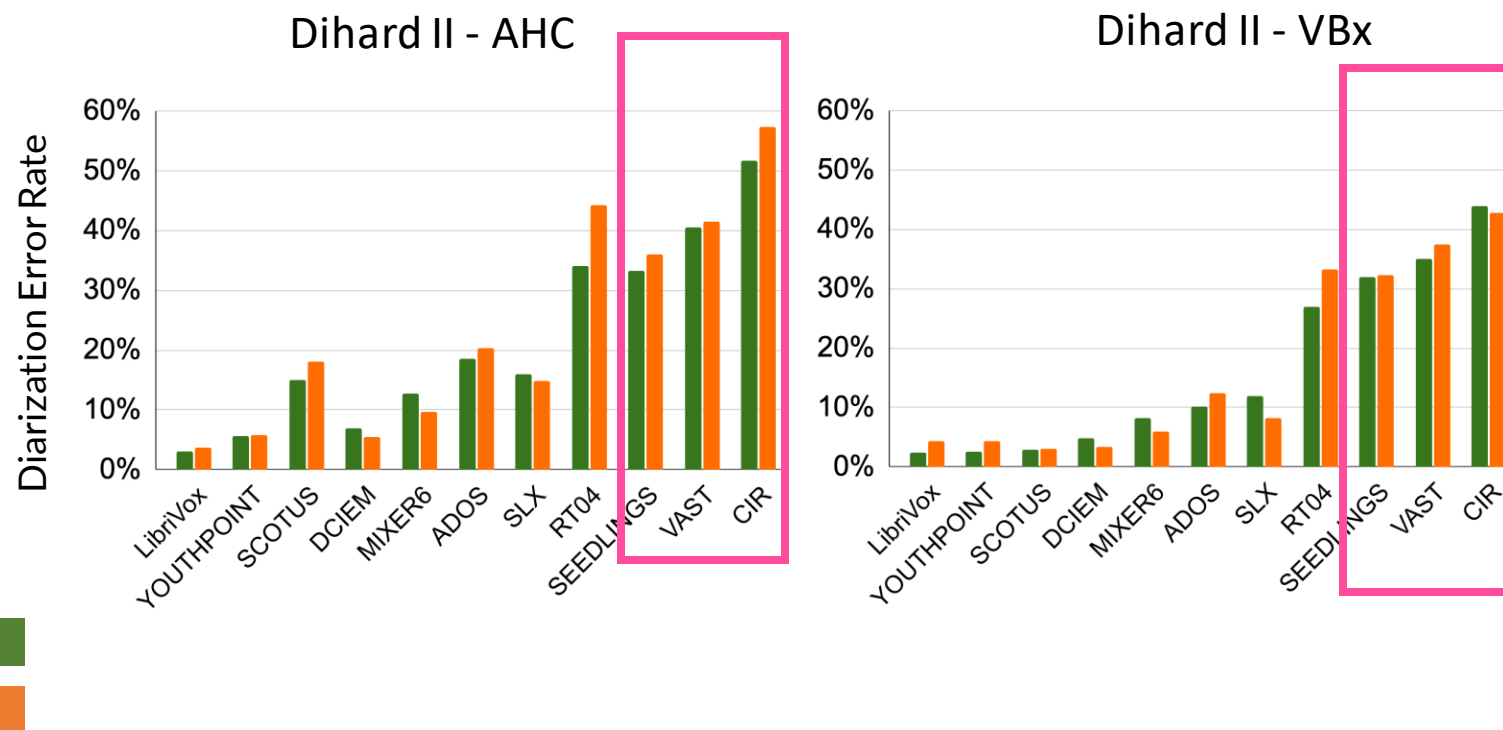
- VB-HMM Clustering - VBx

- Uses x-vectors
- Same model as the Bayesian HMM
- A single Gaussian per state
- Parameters were initialized from pre-trained **PLDA** model.

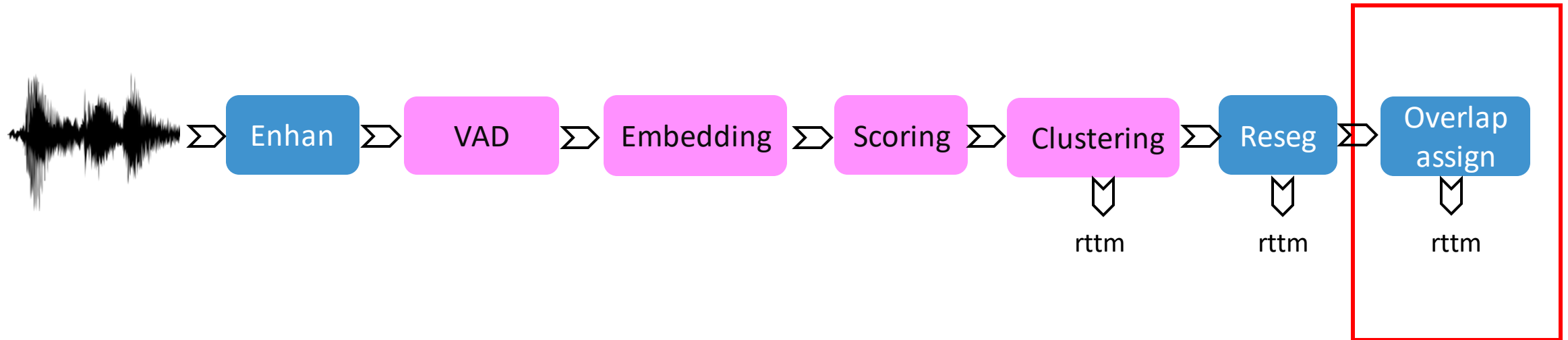


# Clustering

- VB-HMM Clustering

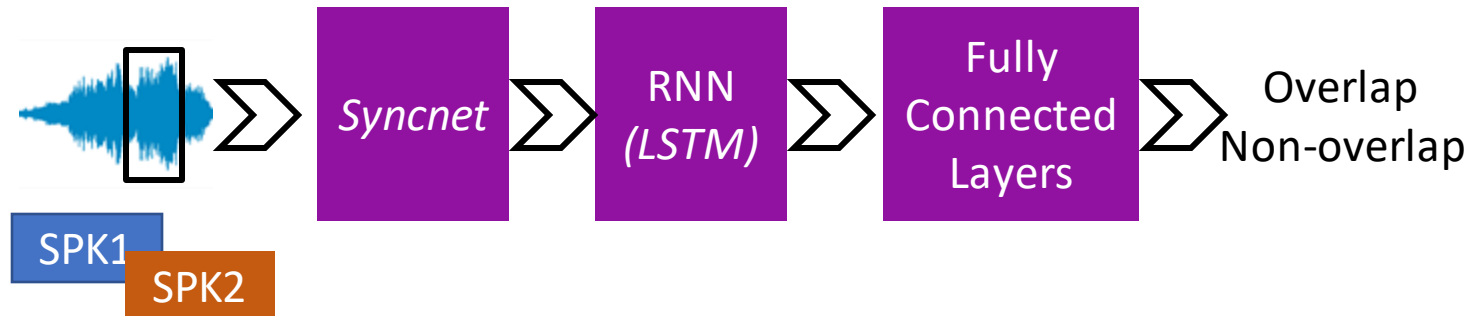


# Key ideas

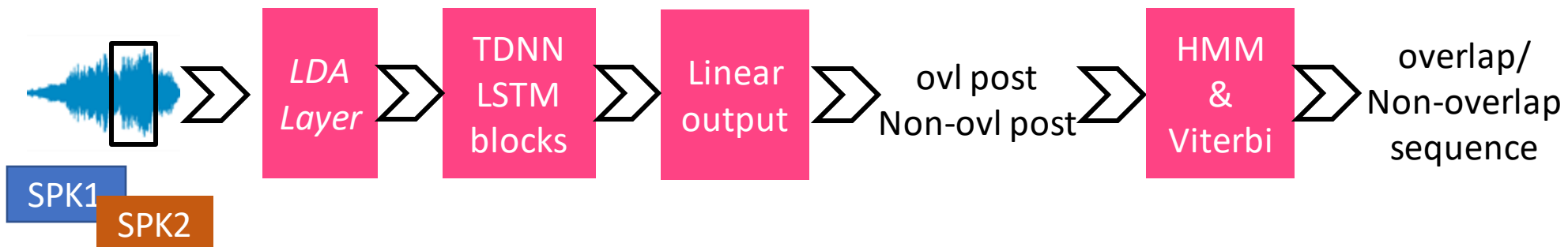


# Overlap Assignment

- Neural network Overlap detector

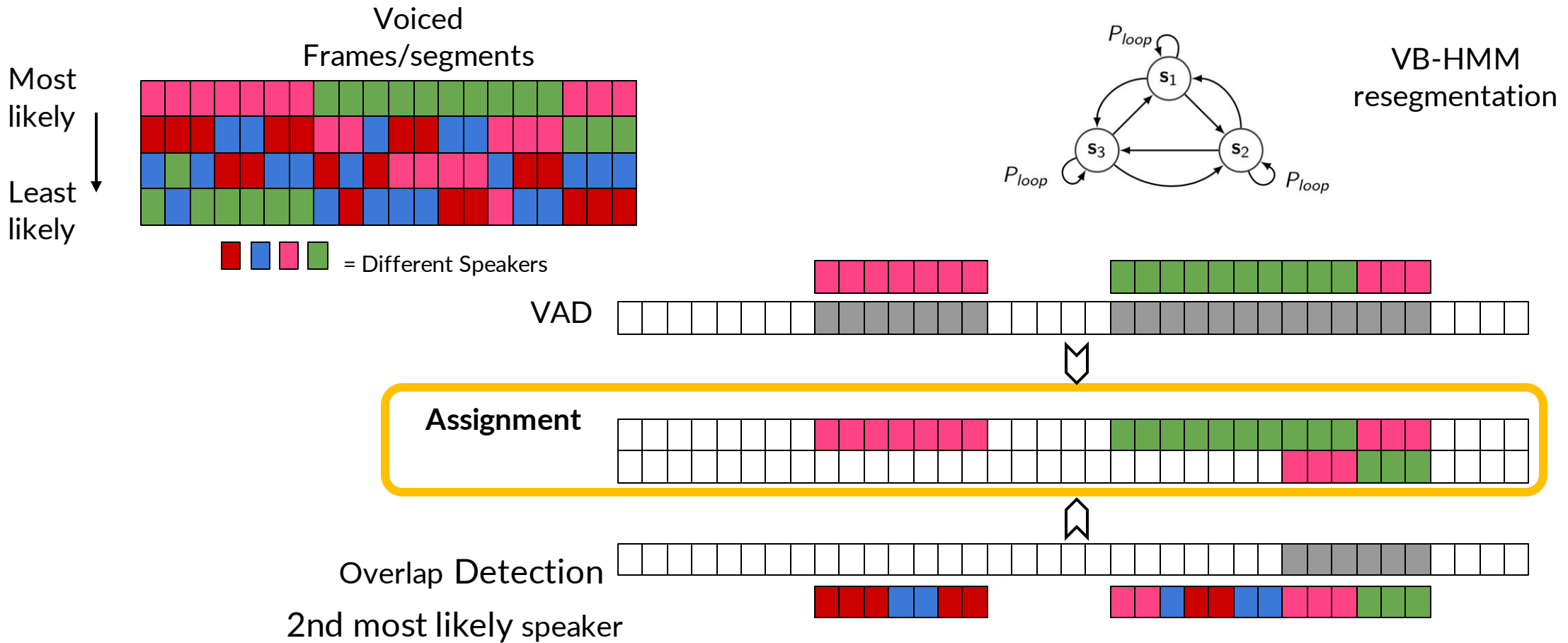


- HMM-DNN based overlap detector

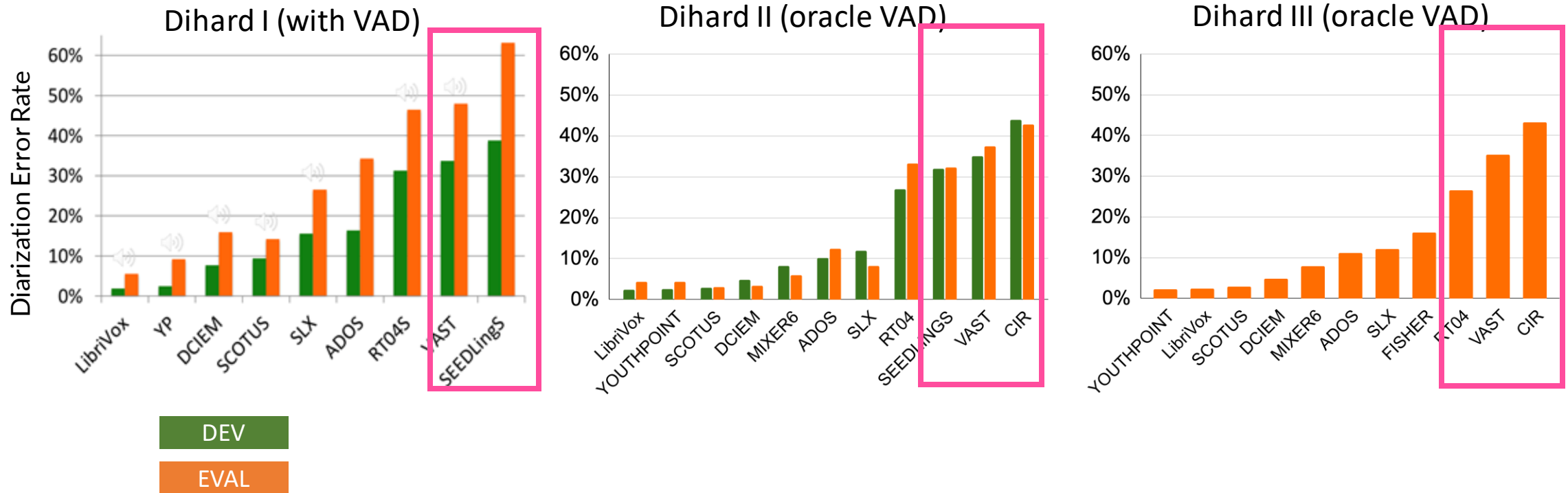




# Overlap Assignment



# Where are we?



Mireia Diez, et.al., Speaker Diarization based on Bayesian HMM with Eigenvoice Priors, 2018.

Greg, Sell, et. Al., Diarization is Hard: Some Experiences and Lessons Learned for the JHU Team in the Inaugural DIHARD Challenge, 2017.

Hitachi-JHU diarization system, DIHARD III, 2020.

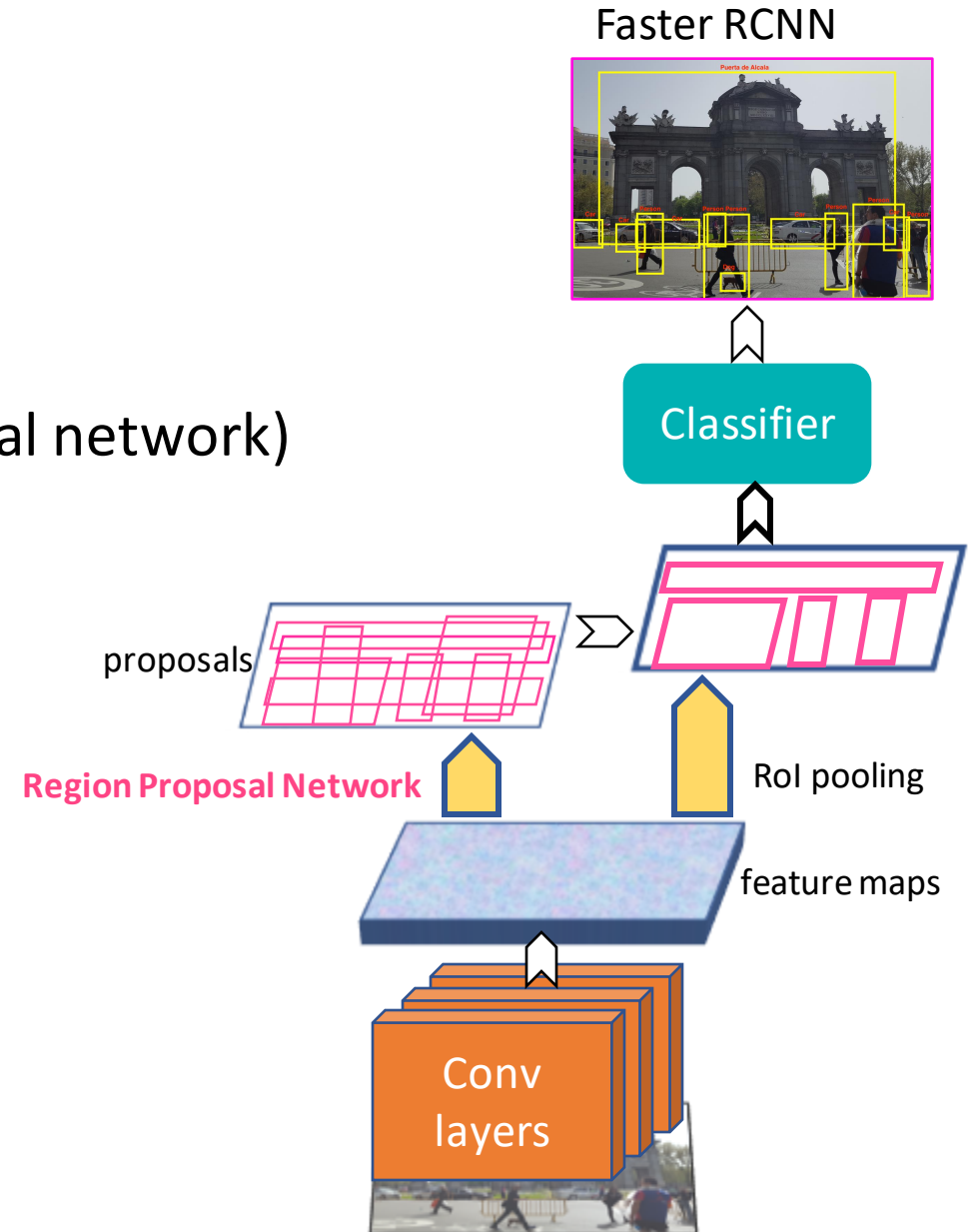
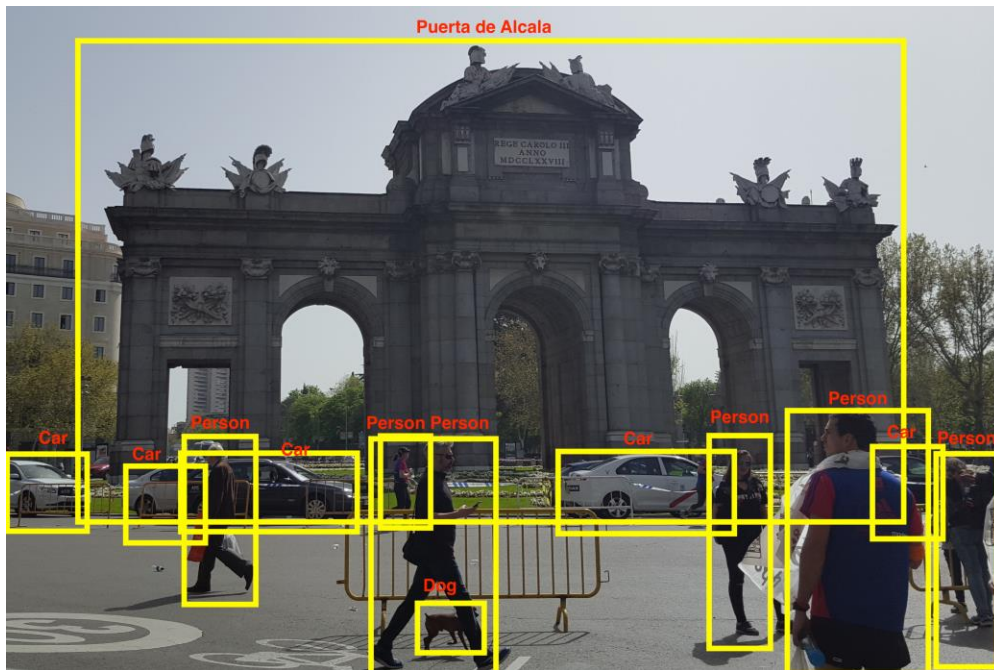
# S.O.S! We have some issues

- We are still trying to handle overlapping speakers
- The system is not designed to minimize the diarization error but optimizes every module separately.

So we look for solutions!

# Region Proposal Network

- One of the first attempts on using NNs
  - Called RPN (inspired by Region proposal network)

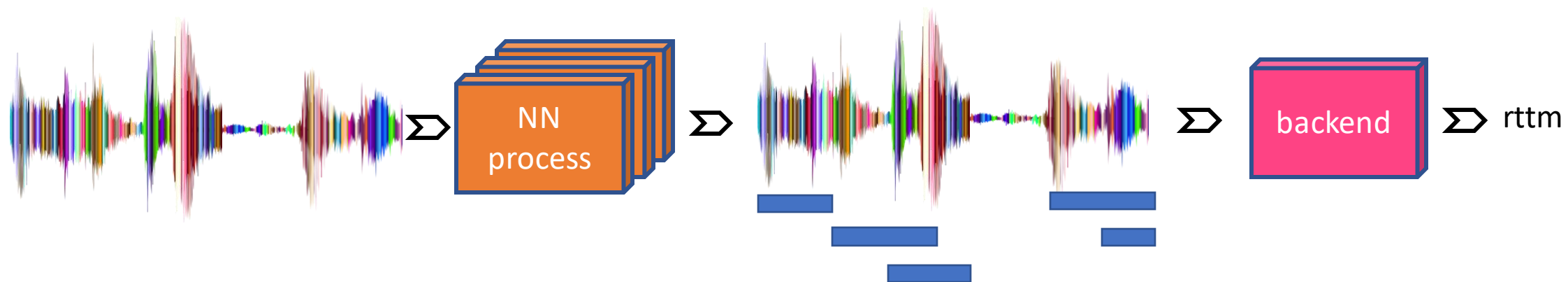


Zili Huang, et.al., Speaker Diarization with Region Proposal Network, 2020.

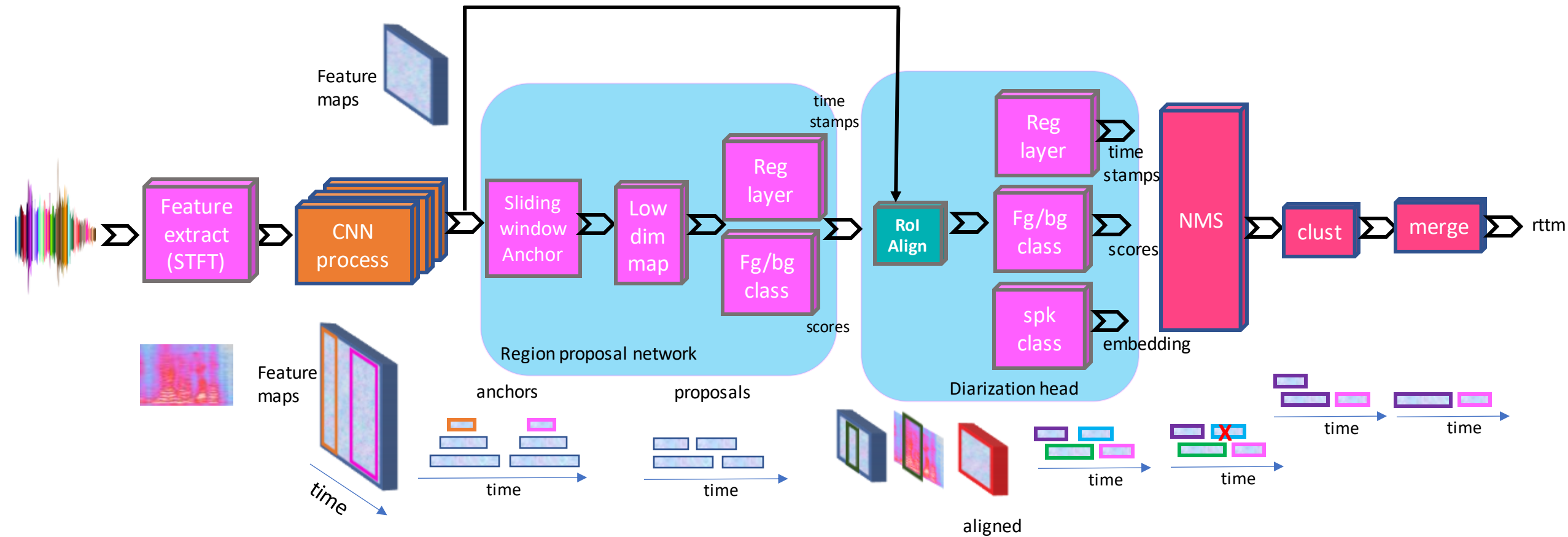
Ren, Shaoqing, et al. "Faster R-cnn: Towards real-time object detection with region proposal networks." *Advances in neural information processing systems*, 2015.

# Region Proposal Network for Speaker Diarization

- Same idea as in image processing but with speech 😊



# Region Proposal Network for Speaker Diarization



# Region Proposal Network for Speaker Diarization

- How to control this pipeline?

$$L = L_{\text{RPN}_{\text{cls}}} + L_{\text{RPN}_{\text{reg}}} + L_{\text{RCNN}_{\text{cls}}} + L_{\text{RCNN}_{\text{reg}}} + \alpha L_{\text{spk}_{\text{cls}}}$$

- $L_{\text{cls}}$  (classification loss): classifies whether a speech segment is foreground or background
- $L_{\text{reg}}$  (regression loss): smooth L1 loss to regress the center point and the length of the speech segments
- $L_{\text{spk}_{\text{cls}}}$  (speaker classification loss): classifies the speaker identity of the speech segments

System	DER (%)	JER (%)
DIHARD baseline	40.86	66.60
DIHARD best VBx	27.11	49.07
RPNSD #oracle num spk	33.12	49.69

We still need to know:

- How to handle overlapping speakers?
- How to design the system in such a way that the diarization error is minimized?
- Using DNNs in a world full of DNNs



# Neural diarization

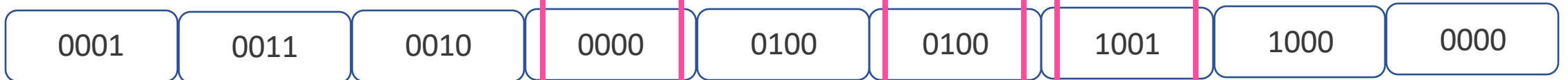


# Diarization as a multi-label classification

- Samples



- Labels



Dorothy

Tinman

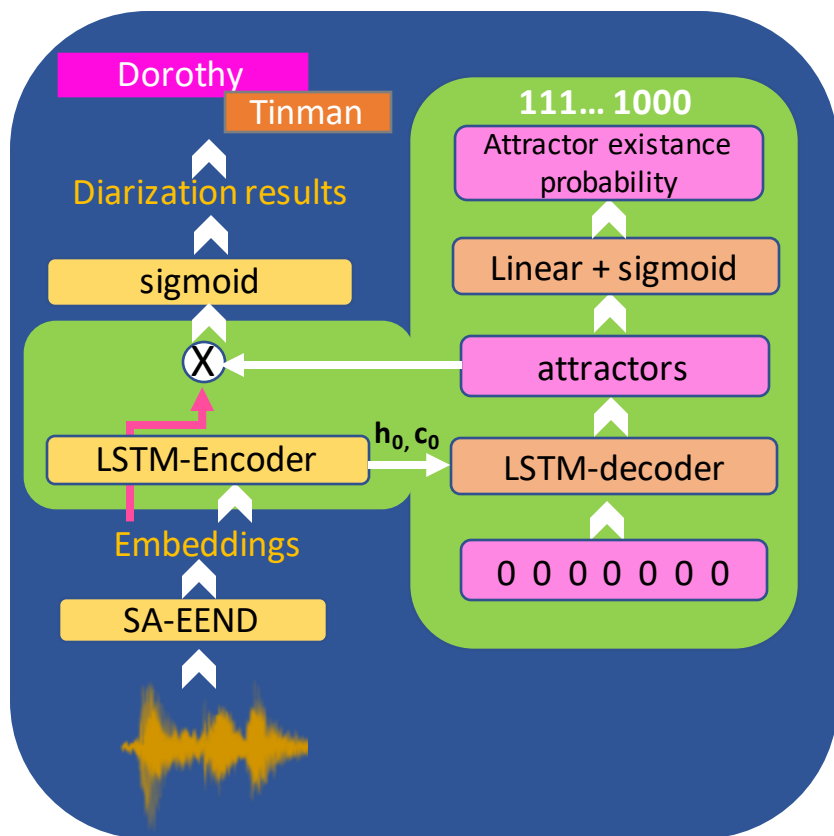
Scarecrow

Dorothy

Lion

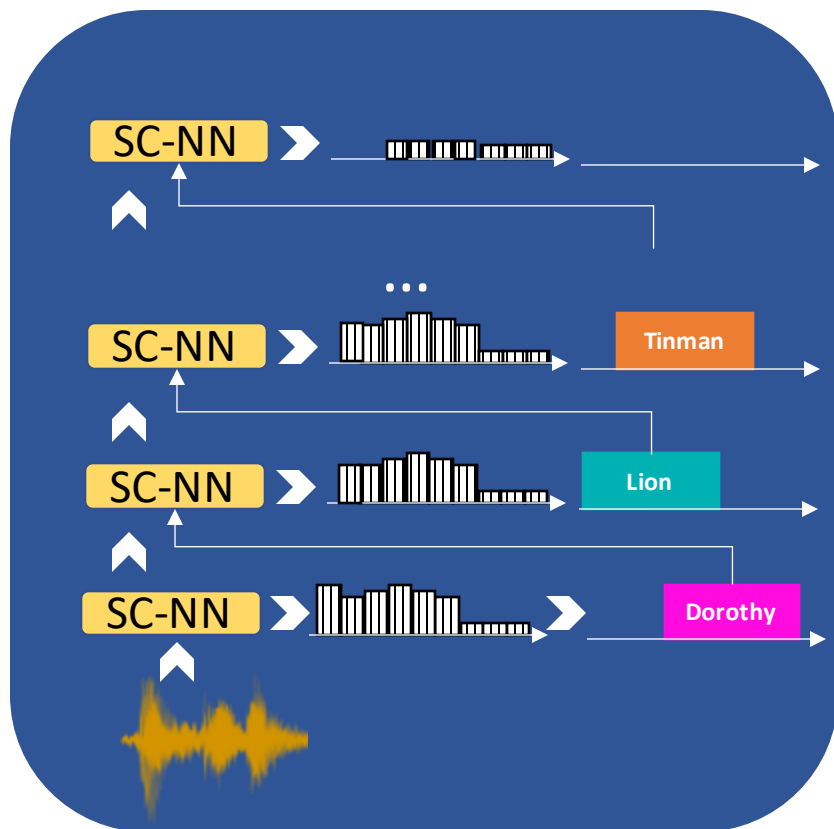


# EEND-EDA



- Encoder decoder attractor for variable number of speakers.
- Starting point is the SA-EEND and the embeddings.
- LSTM-decoder produce attractors
- Dot product between attractors and embeddings produce the diarization results.

# SC-EEND



- Speaker wise conditional EEND
- Deals with variable number of speakers
- Fully conditional model
- Decode speaker-wise sequentially, conditioned on previous speech activities
- Uses teacher–forcing in the training with a modification that takes the appropriate permutation.

# Some results

DIHARD III (track 1- oracle SAD)

System	DEV (DER%)		EVAL (DER%)	
	full	core	full	core
Baseline	19.41	20.25	19.25	20.65
TDNN+VBx+Ovlassign	13.87	14.88	15.65	18.20
EEND-EDA	<b>12.92</b>	13.95	<b>13.95</b>	<b>17.28</b>
SC-EEND	13.13	<b>13.13</b>	15.16	19.14

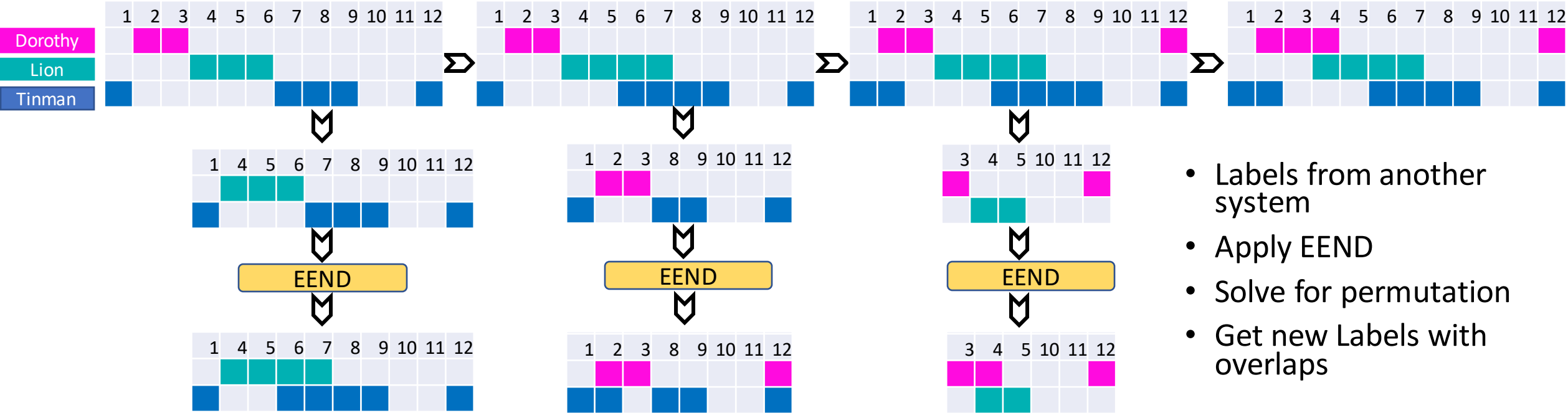
What if we have a system that does not consider overlap, is it possible to fix it?

Nice output



# EENDasP

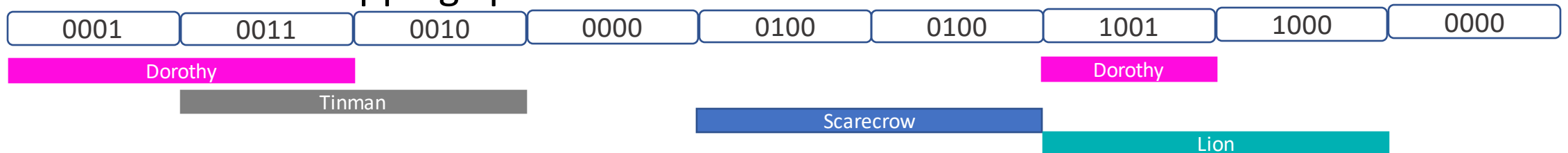
Nice output





- We are considering the two problems now:

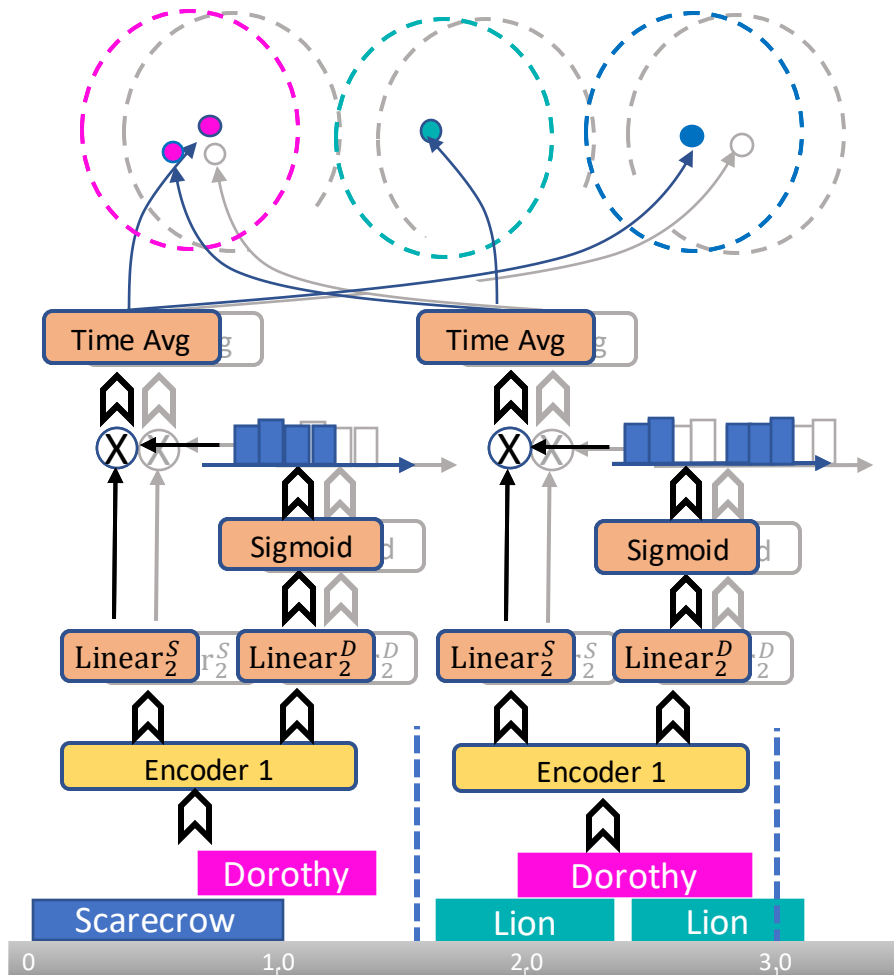
- Handle overlapping speakers



- Designed to minimize the diarization error.

$$\mathcal{L}_{diar} = \frac{1}{ST} \min_{\phi \in \text{perm}(S)} \sum_{t=1}^T \text{BCE}(\hat{y}_t, y_t^{\phi})$$

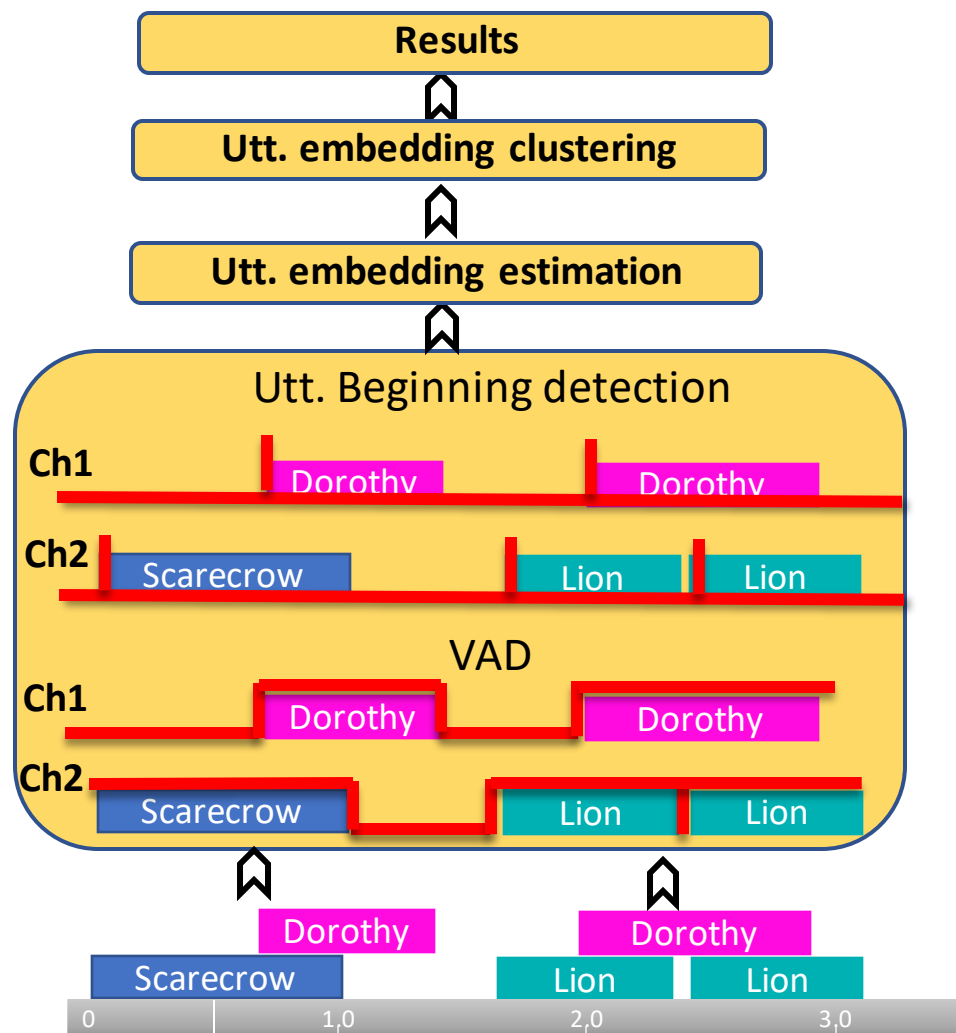
# EEND-vector clustering



- EEND-vector clustering
- Hybrid system for overlapped speech, long recordings and different number of speakers

System	Test duration (min)			
	3	5	10	20
EEND	8.0	8.7	9.3	N/A
Propos+chunking+clust	7.4	6.5	5.9	5.5

# EEND (utterance by utterance)



- Graph-PIT –based VAD
- Segmentation is not longer a limitation
- Utterance by utterance diarization
- Callhome

Method	Number of speakers					
	2	3	4	5	6	Avg
EEND-VC-5s	7.0	14.2	16.7	31.6	29.9	13.7
Graph-PIT-EEND-VC	7.1	12.6	18.3	31.1	30.7	13.5



# More results

DIHARD II (track 1- oracle SAD)

System	DER (%)
Baseline (offline)	26.0
UIS-RNN-SML*	27.3
EEND-EDA w/STB	<b>25.9</b>
SC-EEND w/STB	<b>25.3</b>

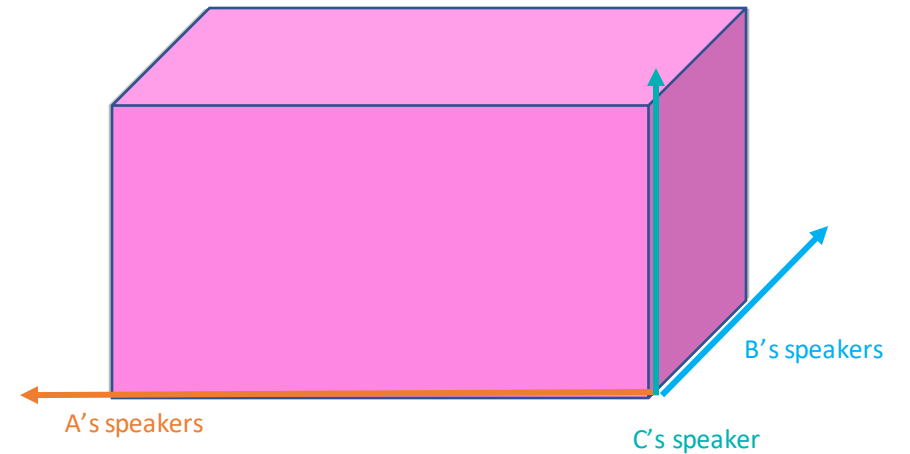
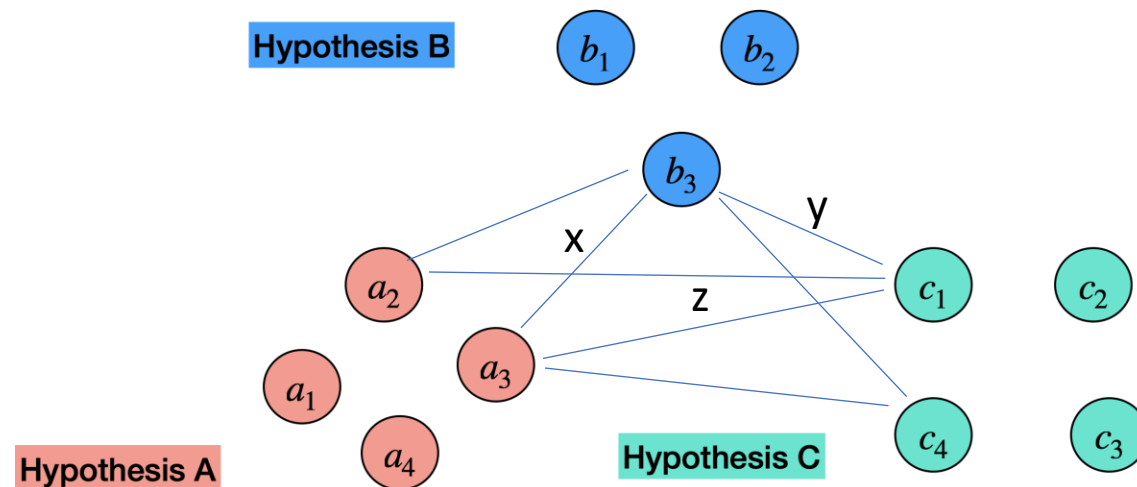
\*Enrico Fini, et.al., Supervised online diarization with sample mean loss for multi-domain data, 2020.  
Yawen Xue, et.al., Online End-to-End Neural Diarization with Speaker Tracing Buffer, 2021.

Let's say that we have four or five different systems with different scores for overlapping regions. Is there a way to deal with them?

DOVER-lap

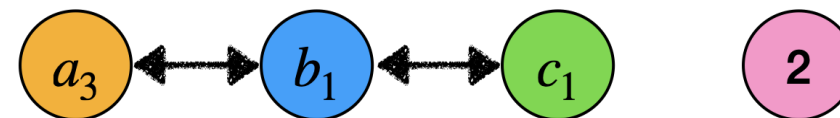
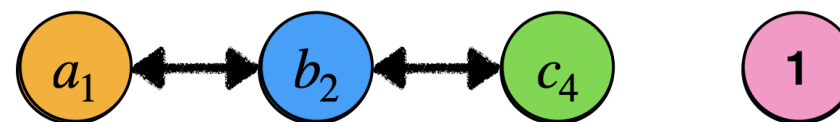
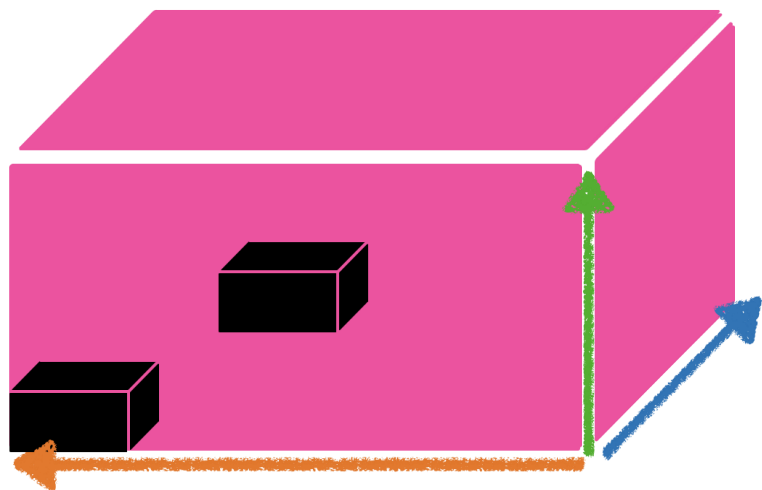
# DOVER-Lap

- Dover with overlap handling
- Two stages:
  - Label mapping

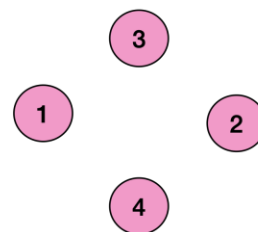


# DOVER-Lap

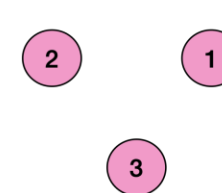
- Tuple with lowest cost and assign the same label



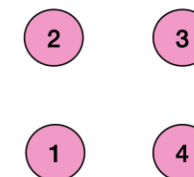
Hypothesis A



Hypothesis B



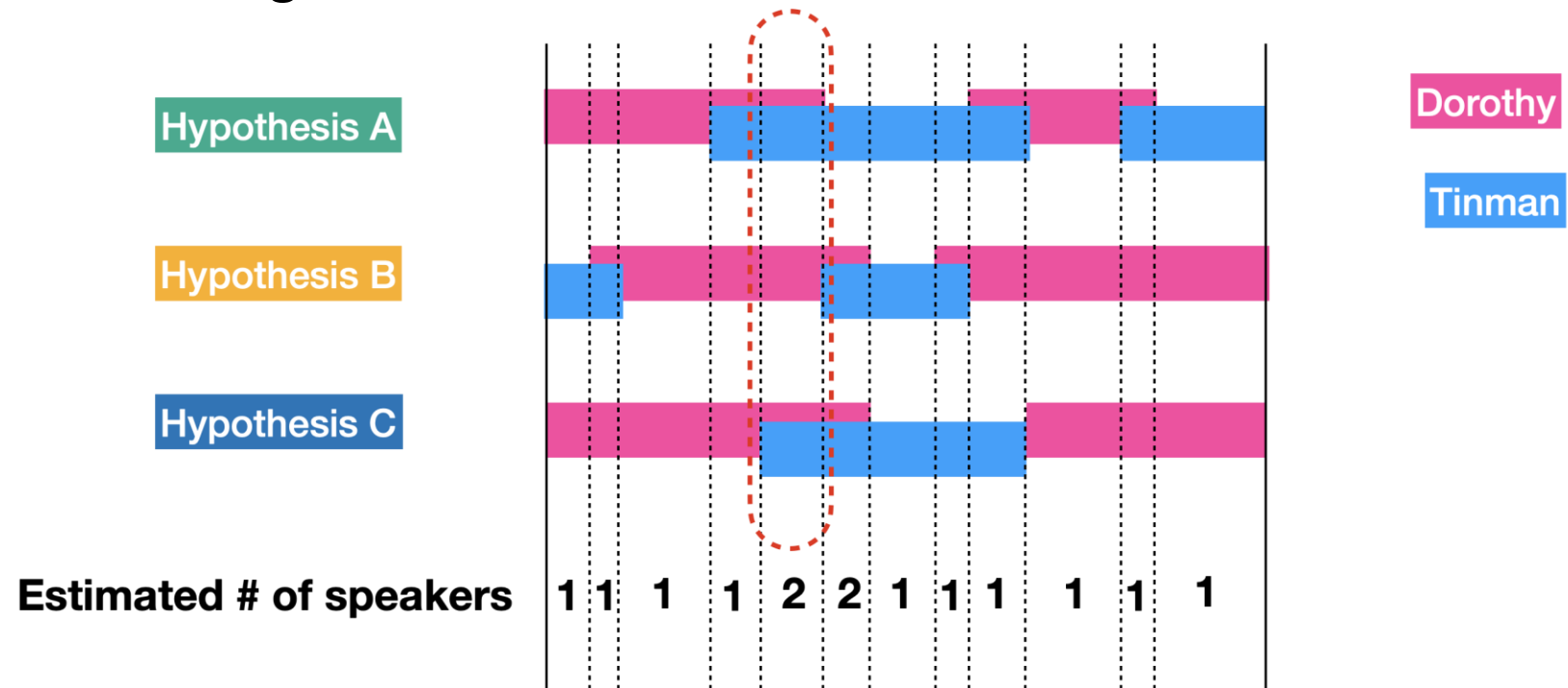
Hypothesis C





# Dover-Lap

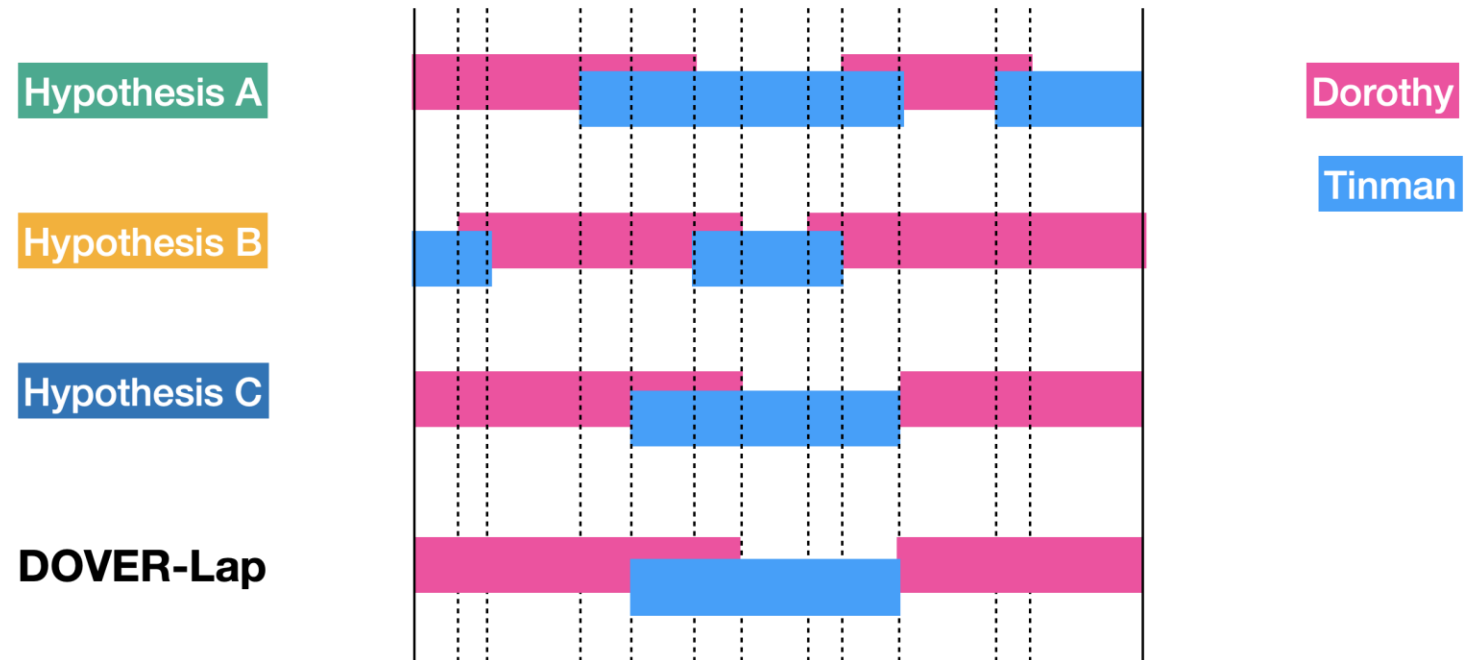
- Label voting



**# speakers** = weighted mean of # speakers in hypotheses


**Weights** -> obtained by ranking hypotheses by **total cost**

# Dover-Lap



# Some results

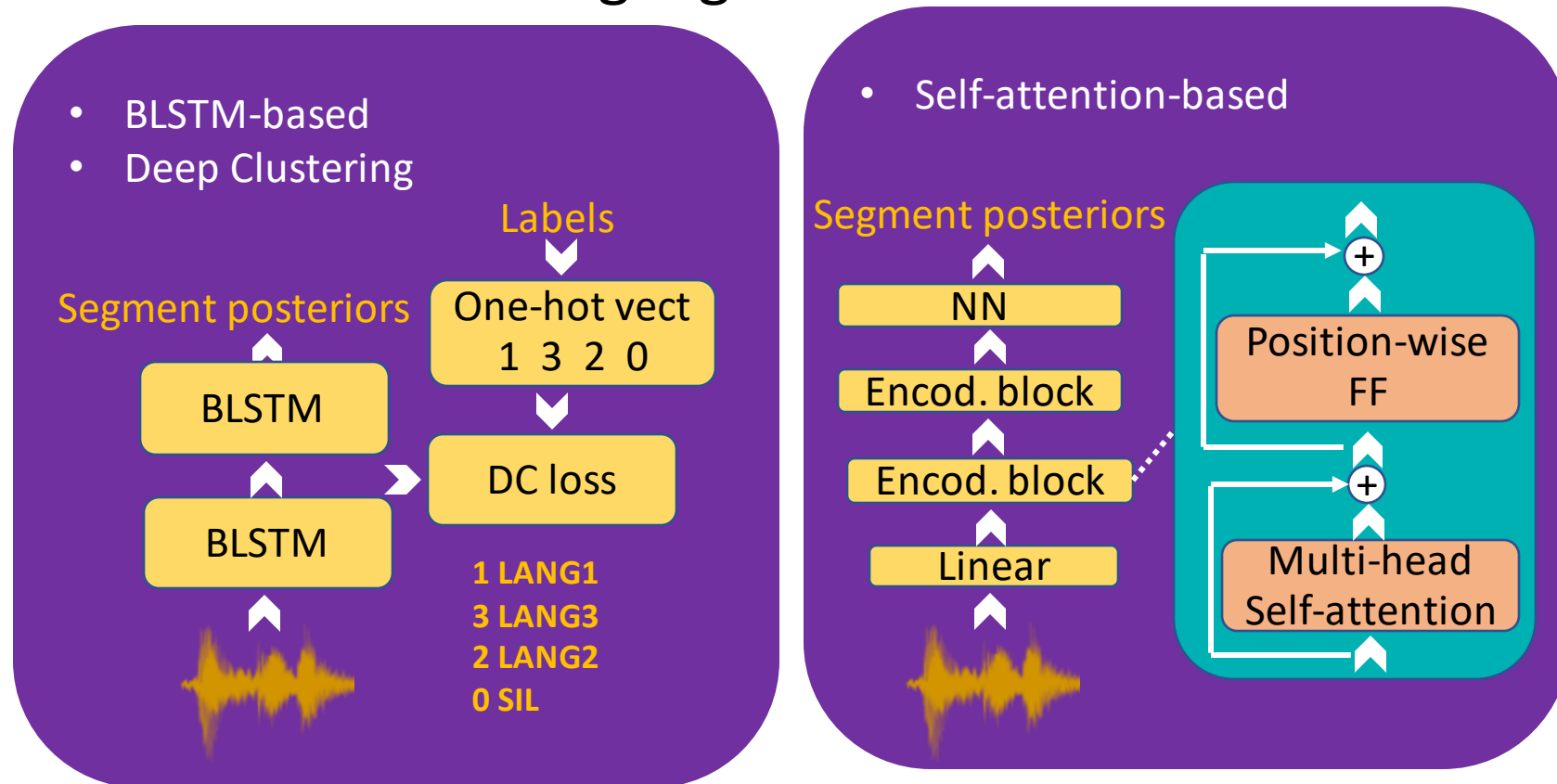
- Dihad III (track 1 – oracle SAD)

System	DEV (DER%)		EVAL (DER%)	
	full	core	full	core
Baseline	19.41	20.25	19.25	20.65
TDNN+VBx+Ovlassign	13.87	14.88	15.65	18.20
EEND-EDA	12.92	13.95	13.95	17.28
SC-EEND	13.13	13.13	15.16	19.14
TDNN+VBx+EENDasP	12.63	14.61	13.30	15.92
DOVER-Lap (  )	<b>10.73</b>	<b>12.56</b>	<b>11.83</b>	<b>14.41</b>

Can this idea be extended to other tasks?

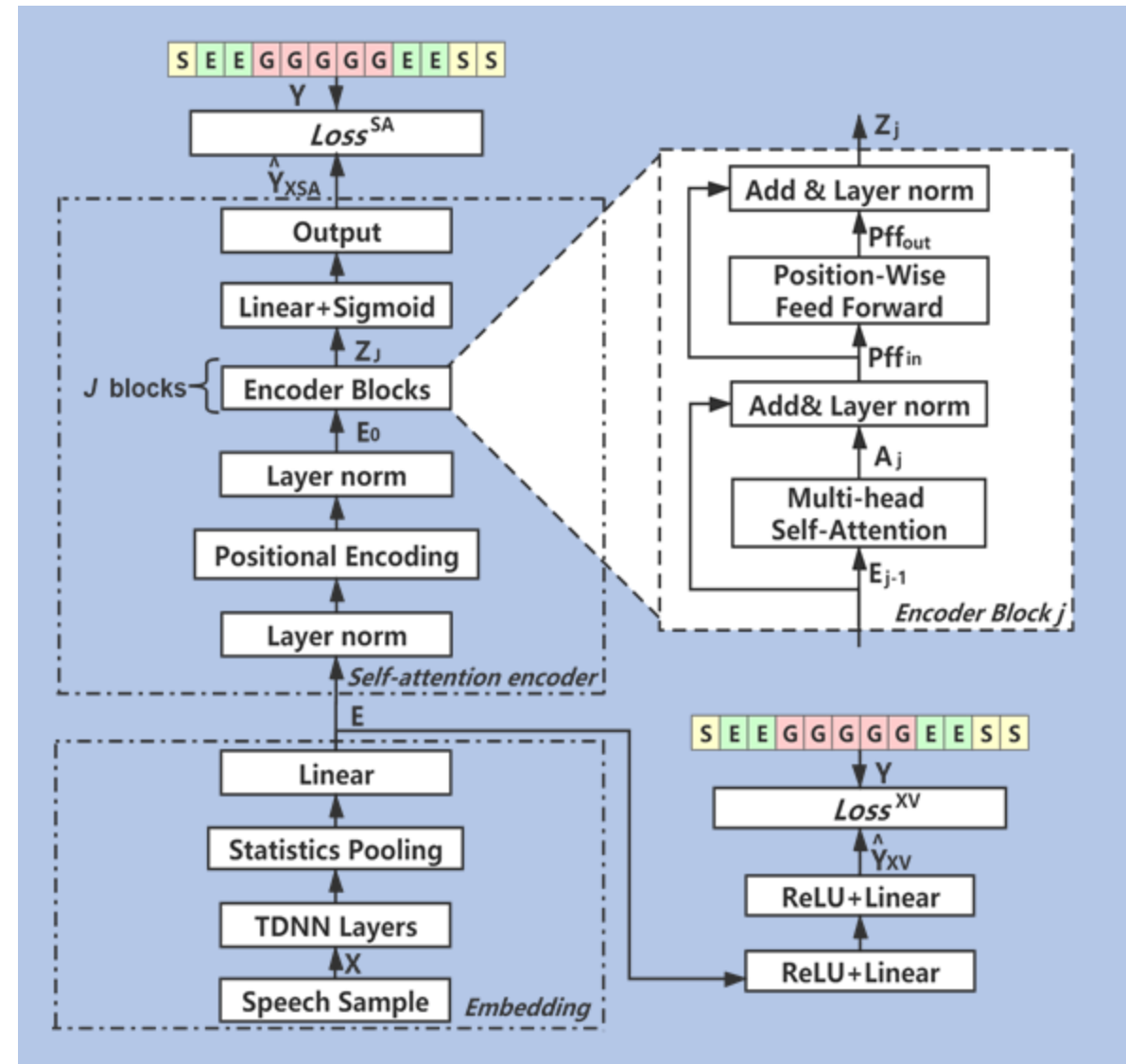
# Language diarization

- End-to-end extension for Language diarization

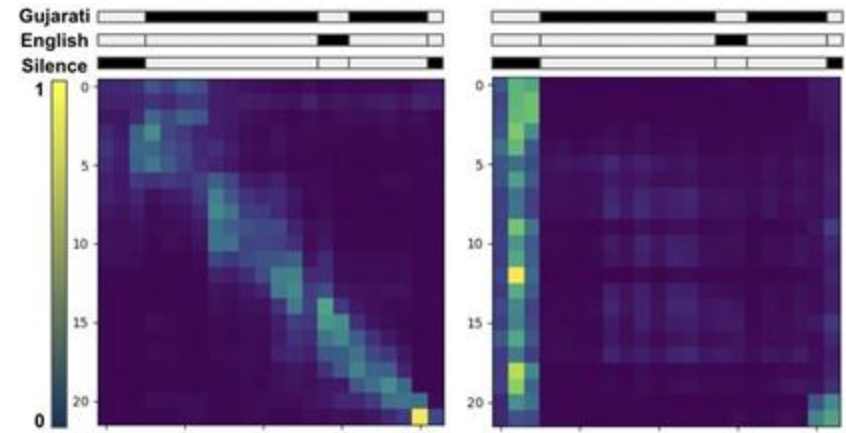


# Language diarization

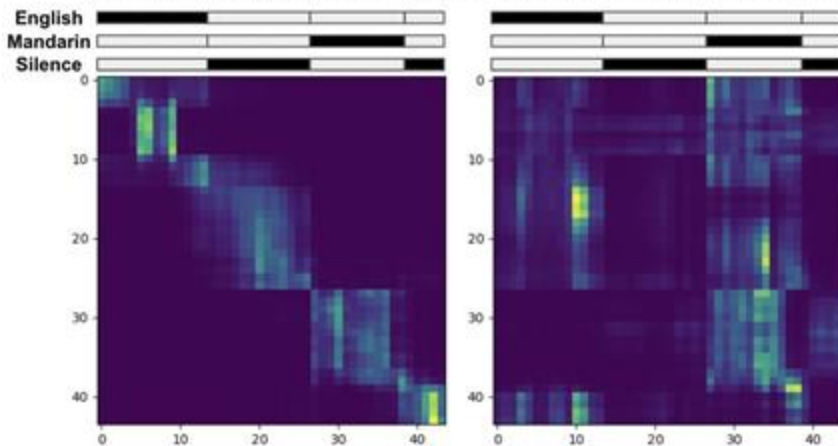
- Similar idea as EEND
  - Bilingual code-switching speech
- End-to-end model with joint training
- The **speech activity detection**, defining silences as a label
- Hierarchical processing:
  - Segment-level (200ms) local language information: x-vector approach
  - Global dependency: self-attention transformer encoder



# Language diarization



(a) Gujarati-English data from the shared task B in WSTCSMC 2020



(b) Simulated data using SEAME dataset

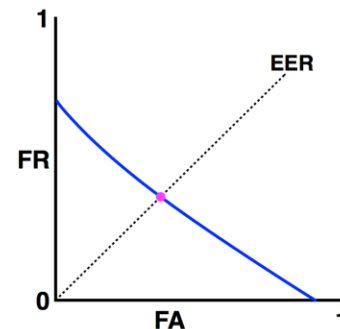
What are the heads doing?

- The left head shows a linear transformation
- The right head highlights the different classes.

# Some results on Language diarization

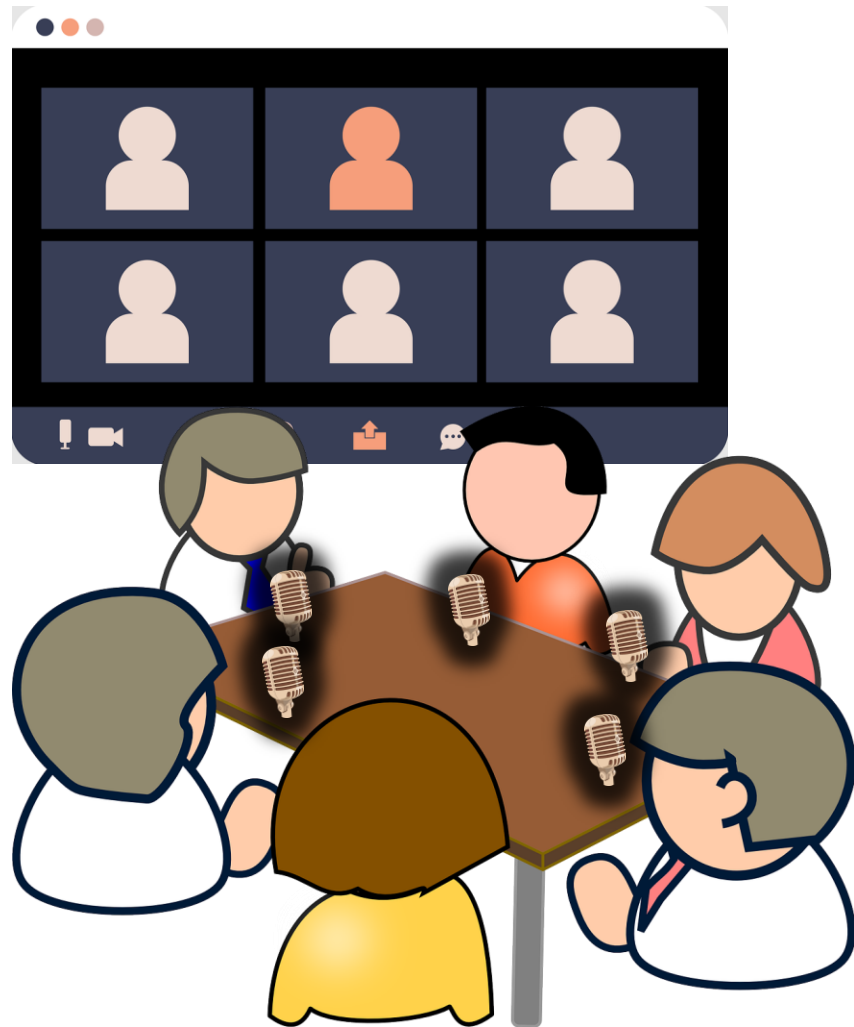
- Comparison of our approaches on 3-language-pair code-switching data in WSTCSMC 2020 (First Workshop on Speech Technologies for Code-switching in Multilingual Communities 2020)

Method	EER (%)					Accuracy(%)
	en	gu	ta	te	silence	
BLSTM-E2E	6.27	3.94	3.55	3.52	2.97	80.15
SA-E2E	6.33	3.59	3.73	3.65	3.49	79.21
XSA-E2E	5.99	2.98	3.21	3.05	3.56	81.20





# Multi-channel end-to-end Neural Diarization



- EEND using multi-channel signals from distributed microphones
- Transformer encoders in EEND replaced to process a multichannel input:
  - Spatio-temporal encoders
  - co-attention encoders
- Model adaptation method using only single-channel recordings.
- The method works on multi-channel inputs, such as in hybrid meetings

# Why is this effort important?

Because we can use it as part of downstream tasks

- ASR
- Multi-microphone ASR
- Virtual assistants
- Broadcast transcriptions
- Emotion recognition
- ...

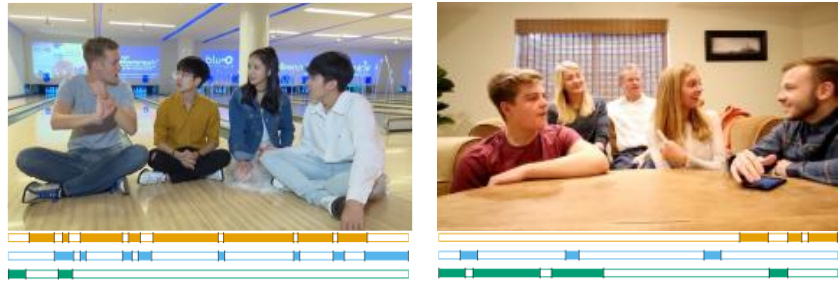
# Multimodal

# Audio-Visual (dataset)

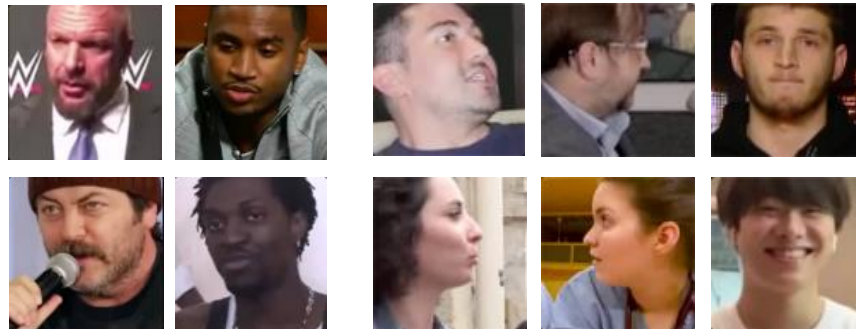


AMI

AVA-AVD



MSDWild (Ours)



Speaking  
VoxCeleb2

Speaking  
MSDWild (Ours)

Not-Speaking

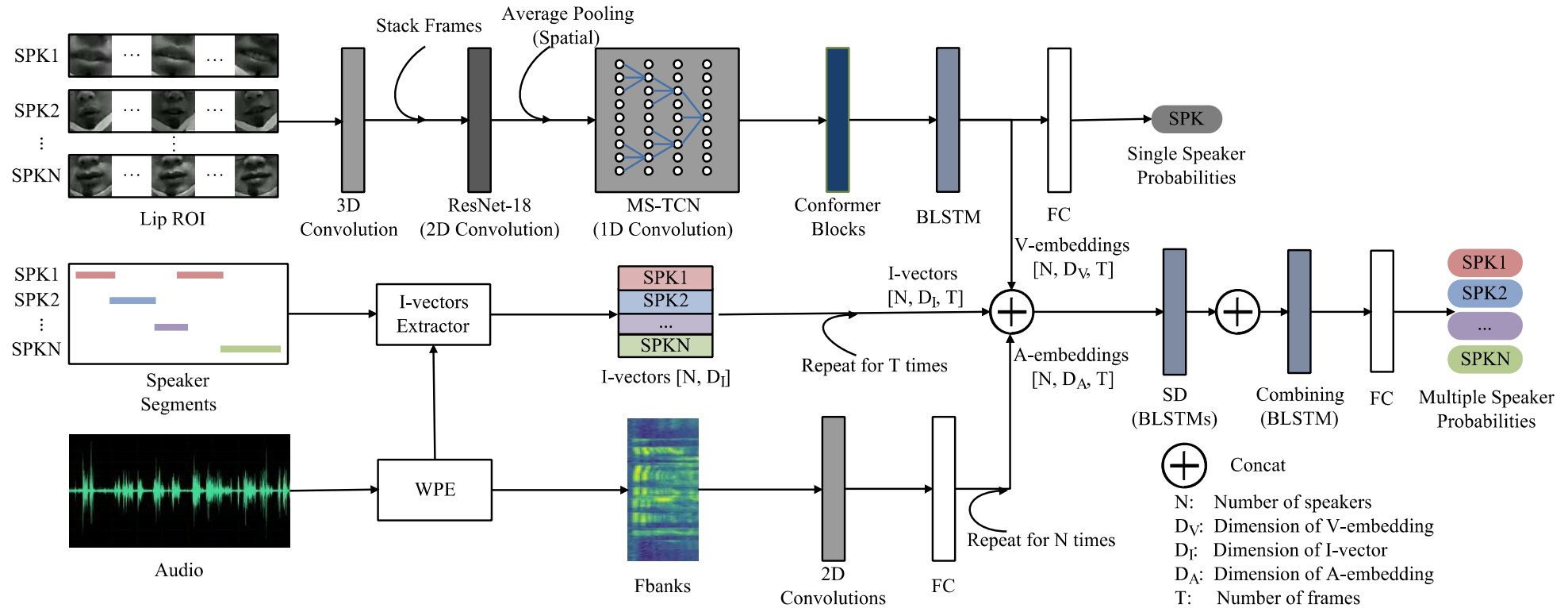
- MSDWILD: MULTI-MODAL SPEAKER DIARIZATION DATASET IN THE WILD
- Benchmark dataset
- Rich real world scenarios
- Different languages
- No overediting
- Overlap speech
- Cocktail party research
- Video and audio released

# Audio-Visual (dataset)

Dataset	source	#videos	duration	Speech %	overlapped	#spk	#SC	noise	continuos	language
Voxconverse	TV-show	448	63h 50min	90.7	3.6	1/5.6/21	3.28	no	yes	en
MSDWild	Daily conversation	3143	80h 3min	91.29	14.01	2/2.7/10	11.8	yes	yes	multi

Method	DER	JER
Audio-only	43.15	84.28
Visual-only	53.71	62.71
Audio-visual (two-stream)	52.6	63.7
Audio visual (fused)	25.86	54.79

# Audio-Visual (end-to-end)



# Audio-Visual (end-to-end)

- End-to-end Audio-visual diarization:
  - audio features, multi-speaker lip (ROI), ivectors.
- Classification output layers produce labels.
- Handle:
  - Overlap, speech vs non-speech
- I-vectors used for alignment
- **MISP** (eval)

Model	DER (w -VAD)	DER (w/o -VAD)
TS-VAD	28.95	-
VSD	13.07	19.64
Audio-visual AVSD w/i-vector	10.05	
Audio-visual AVSD i-vector	10.1	
Audio-visual AVSD i-vector+Joint training	9.49	10.99
Doverlap	8.85	-

# Takeaways

- EEND, VBx, TS-VAD are still good solutions.
- Unsupervised/self-supervised methods were proposed.
- Some methods deal with long recordings (podcasts).
- New methods are out of the traditional (sparse optimization, role labels )
- Online diarization is becoming feasible with good results.
- Some approaches are still on simulated data, some others are on real data





# Research directions

- Multi-modal diarization (text, video, audio)
- Speaker imbalance
- Online diarization
- Real world data
- Long recordings ( like children's speech)
- We should consider speech separation algorithms that combined with diarization can improve the overall performance
- Diarization is only part of more complex scenarios, eg., using diarization for ASR



Why does it really matter?

# Bilingualism in child-centered speech

- Day-long recordings

- In real life we *don't have dev data*, we *only have eval data* 😊



- Diarization error rates for all systems drop dramatically.

# Child centered data

## Speaker Diarization

System (Seedlings)	DER (%)
Baseline AHC	63
VBx+	61.49
EEND-EDA	62.57
Oracle VAD VBx	32.33

System (BLIP)	DER (%)
Baseline AHC	86.26
VBx	65.81
Oracle VAD VBx	41.01

## Language diarization

System (BLIP)	ACC(%)	EER(%)
xvector-LD	80	20
XSA-LD	89	11

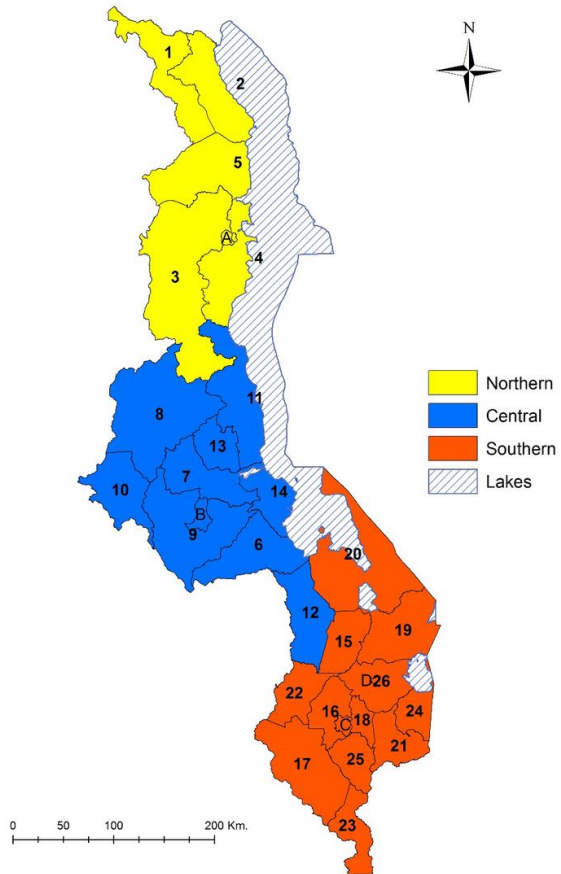
VanDam, Mark, VanDam Public Daylong HomeBank Corpus. doi:10.21415/T5388S,2018.  
<https://media.talkbank.org/homebank/Public/VanDam-5minute/CI40/>  
Victoria Chua, Suzy Styles, et.al., Blip audio private collection provided by NTU, 2020.  
Liu, Hexin, BLIP data results, internal report, 2021.

# Child development study

- The loop: *“child and adolescent under-development, as existing programs struggle to deliver the right intervention at the right time.”* 😞
- Breaking the loop: *“high-frequency data on child and youth development aided by customized technologies to inform timely responses, tailored to the needs of each child and adolescent.”*

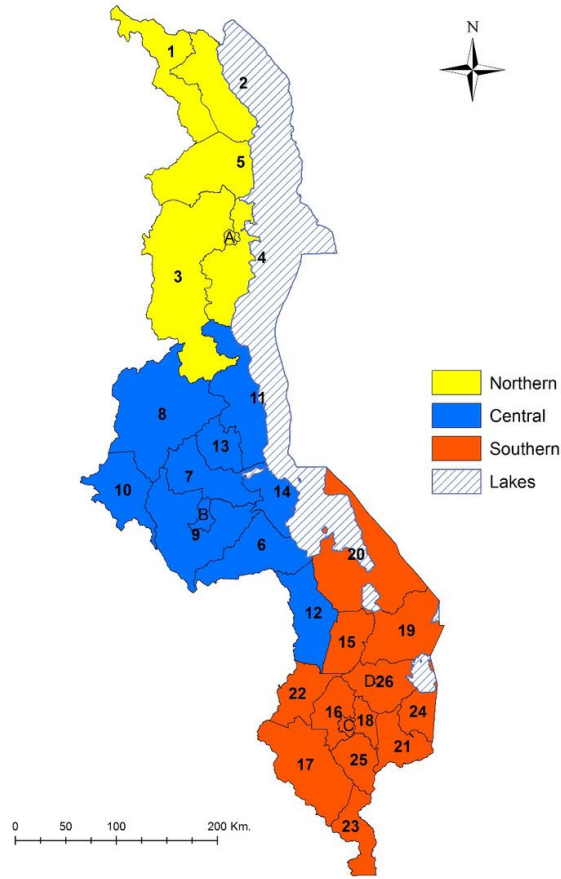


# What are we doing?



- Do you know this country?
- Do you know the language?

# What are we doing?



- Location: Malawi
- Language: Chichewa
- From the recordings, we plan to analyze aspects of adult-child and child-child communication such as:

diarization

- Number and type (adult vs child) of speakers/interactants
- Amount of adult speech
- Amount of child speech
- Number of conversational turns between the target child and other speakers (adults and children)
- Timing of conversational turns

ASR

- Number of different words (requires transcription)
- Complexity of utterances (requires transcription)
- Types of questions (requires transcription)

# Takeaways

- Lots of flavors to choose from 😊
- We have huge improvements, but we are not yet there
- Neural diarization is becoming as good as embedding clustering methods
- Overlap detection is still an ongoing research
- Speaker imbalance is also an ongoing research
- Diarization to help downstream tasks (like ASR)
- Day-long recordings, cocktail party scenarios need diarization solutions
- Next time we can talk about self-supervised learning for diarization as a new direction



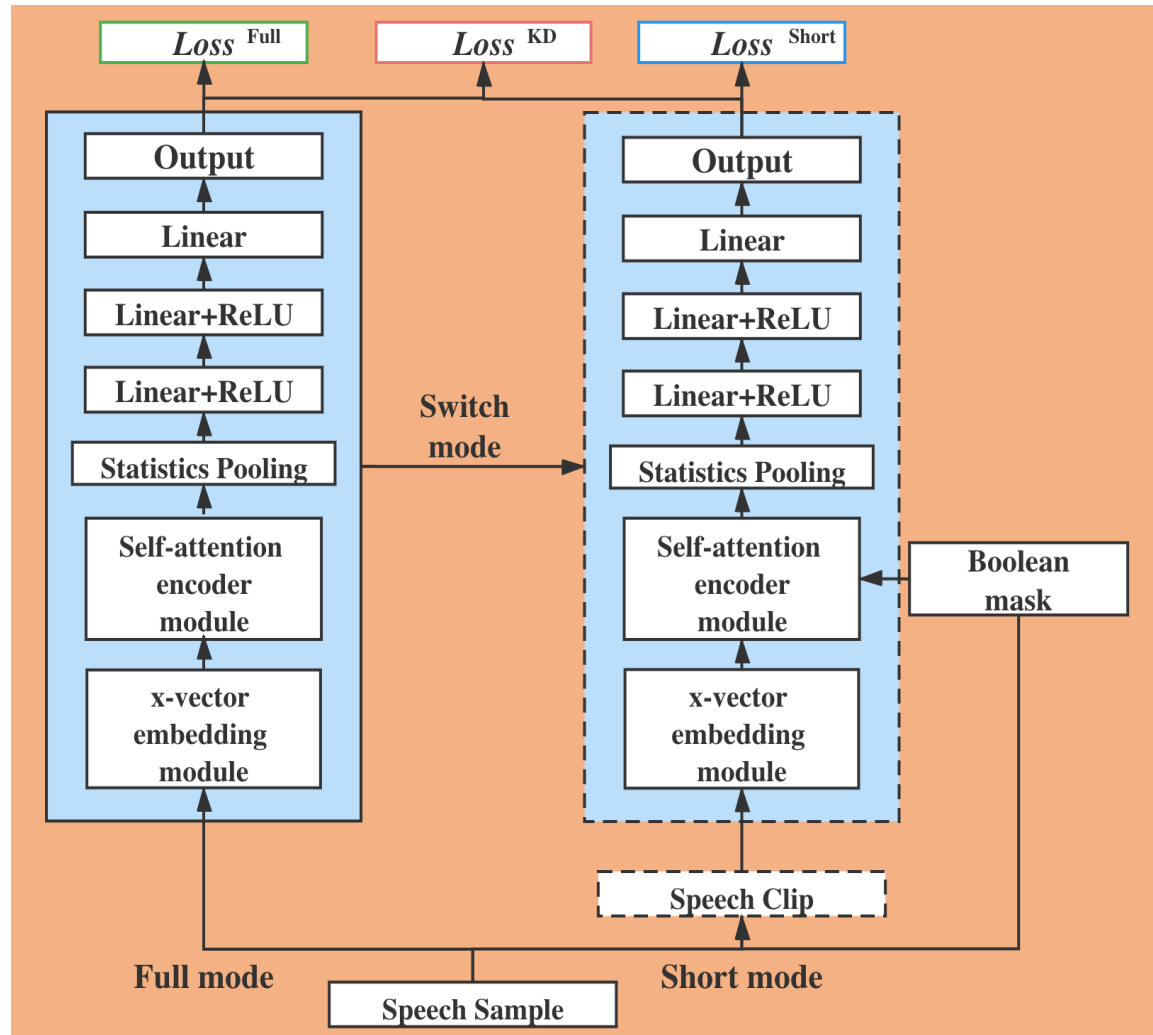
Thank you!



Questions?

- Complementary slides if needed.

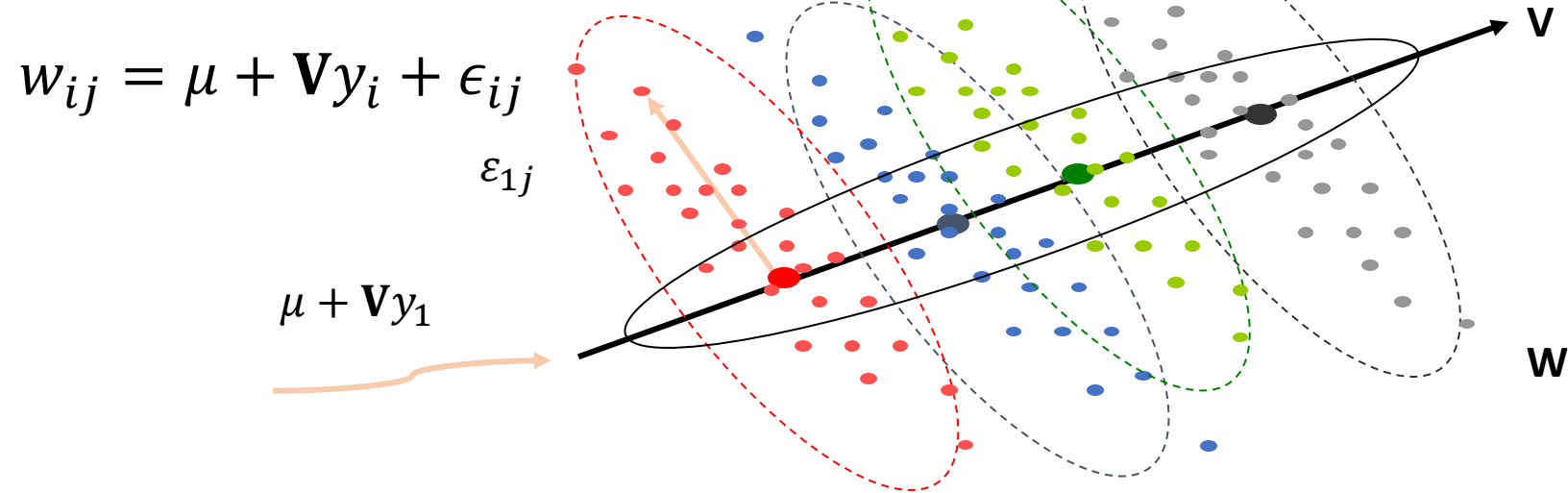
# Dual Mode Language Identification



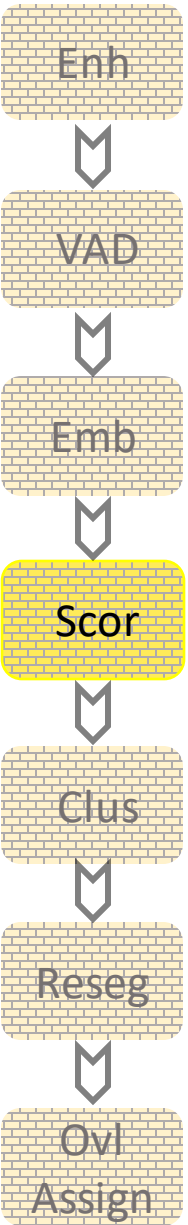
- Addressing short utterances
- XSA-LID model jointly optimizing
  - full-length speech
  - short clip (extracted by a specific Boolean mask)
- We apply knowledge distillation (KD) to boost the performance on short utterances.
- We investigate the impact of clip-wise linguistic variability and lexical integrity for LID.

# Scoring

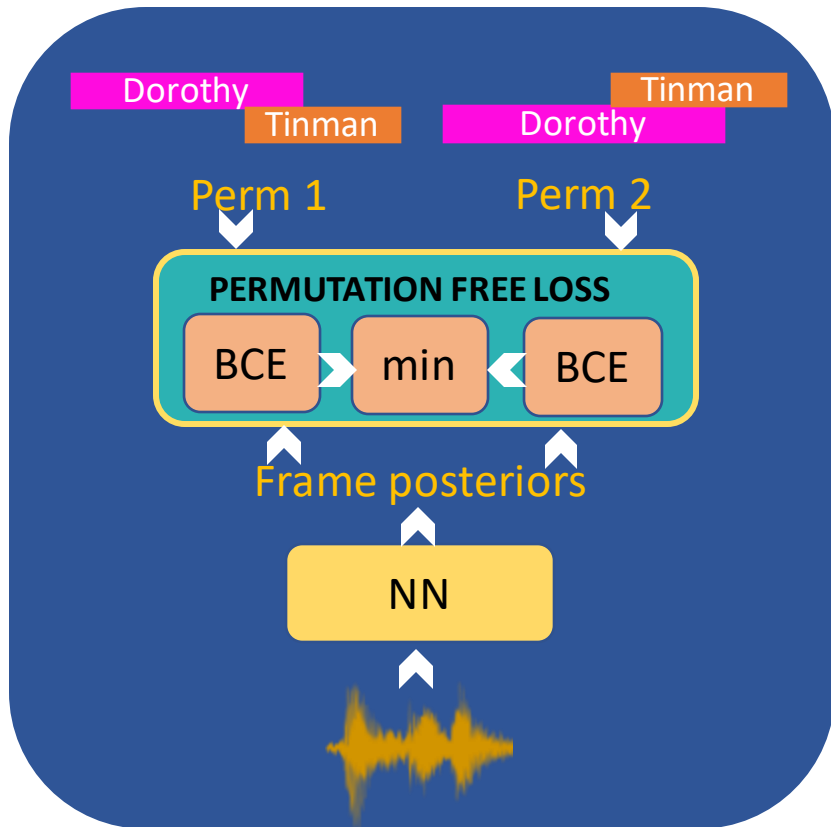
- PLDA



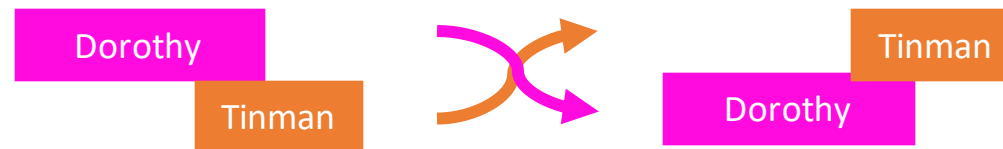
$$\text{LLR} = \log \frac{P(w_1, w_2 | \text{same})}{P(w_1, w_2 | \text{diff})} = w_1^T A w_2 + w_1^T B w_1 + w_2^T B w_2 + C^T w_1 + C^T w_2 + D$$



# EEND

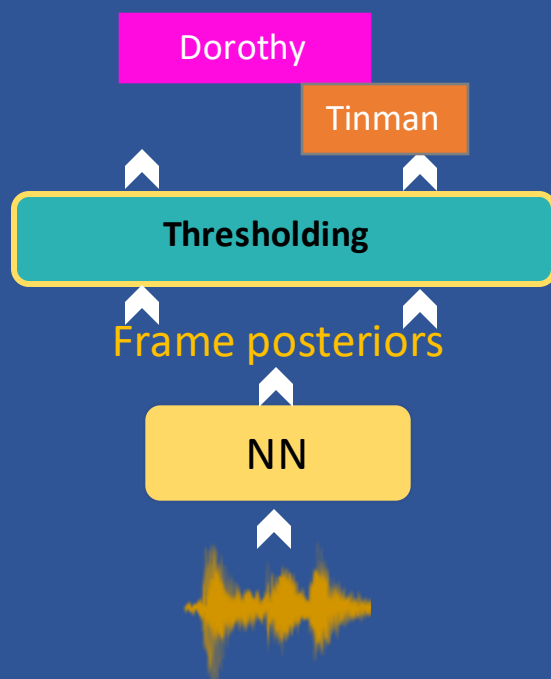


- End-to-end neural diarization
- Single network, supervised
- Two speaker case, proof of concept
- Handles overlapping speech!
- Training uses permutation invariant training (PIT) to prevent the labeling ambiguity.



# EEND

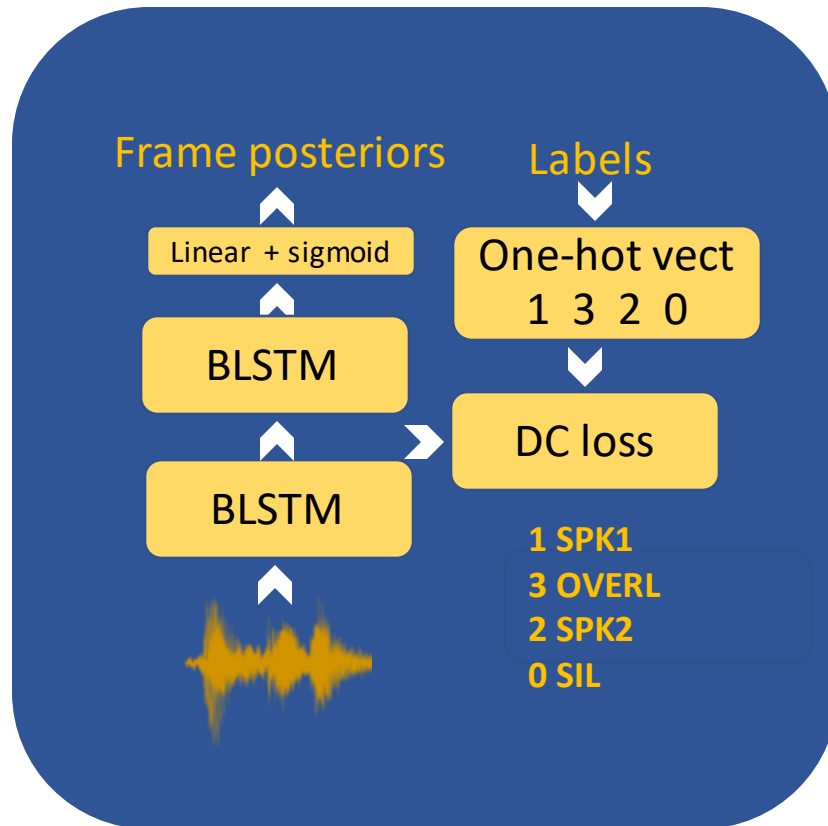
- Two speakers



- During test

- Frame posteriors
- Threshold
- Speaker labels

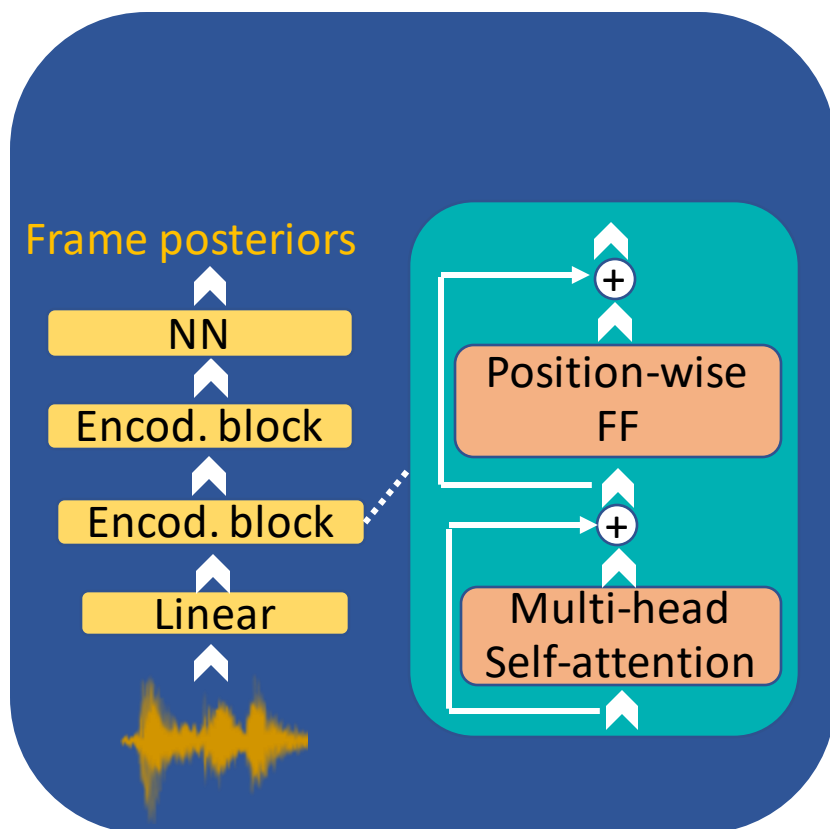
# EEND



- BLSTM
- Embeddings from lower layers
  - Speaker training criterion on middle layer activations
- Deep Clustering to partition the embedding into:
  - Speaker dependent-clusters
  - Overlap
  - silence



# EEND-SA



- BLSTM captures local temporal dynamics
- Self attention captures long context
  - The heads capture different characteristics

