

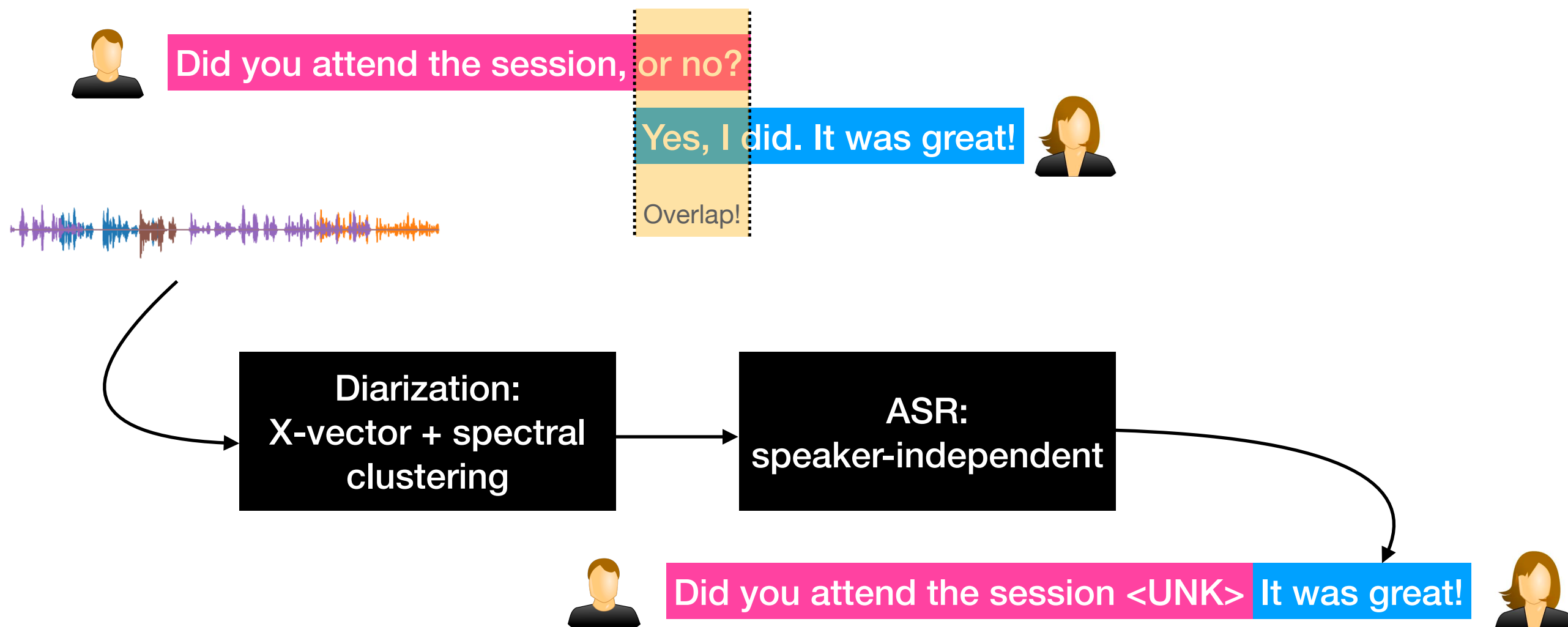
Informed Target Speaker ASR

JSALT 2020

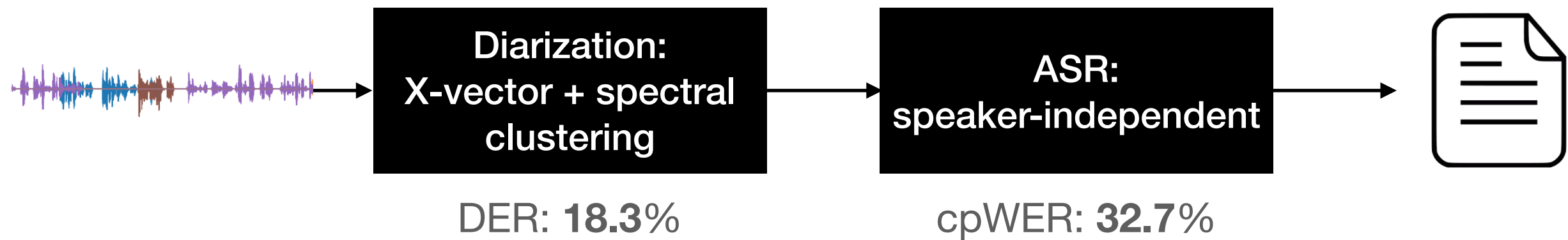
DESH RAJ

Collaborators: Marc Delcroix, Shinji Watanabe, Kateřina Žmolíková, Pavel Denisov, Zili Huang, Sanjeev Khudanpur

The single-stream baseline system...

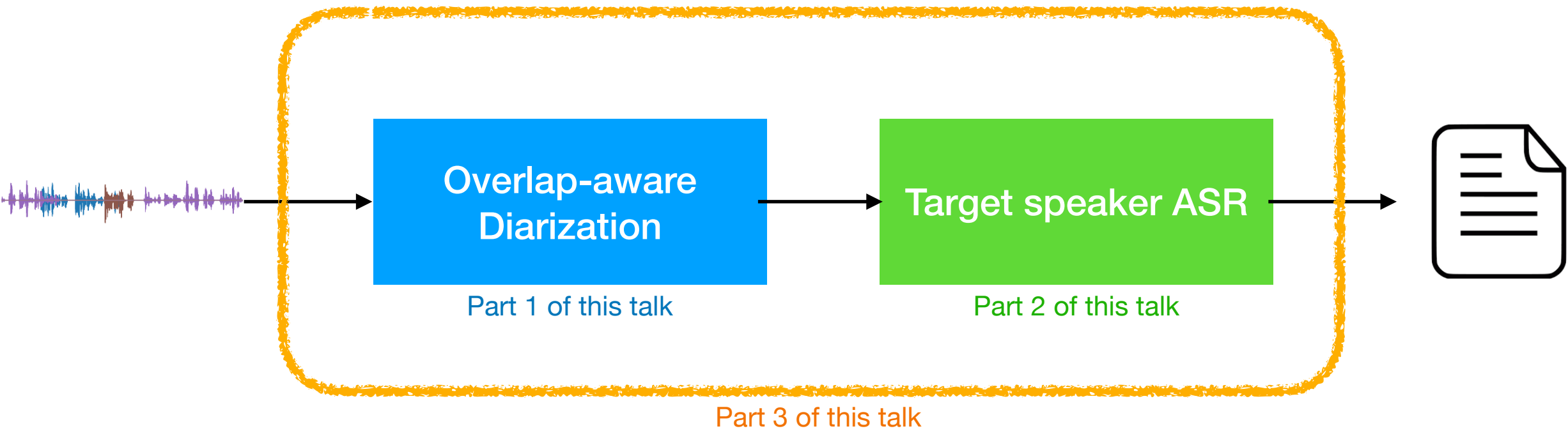


The single-stream baseline system...

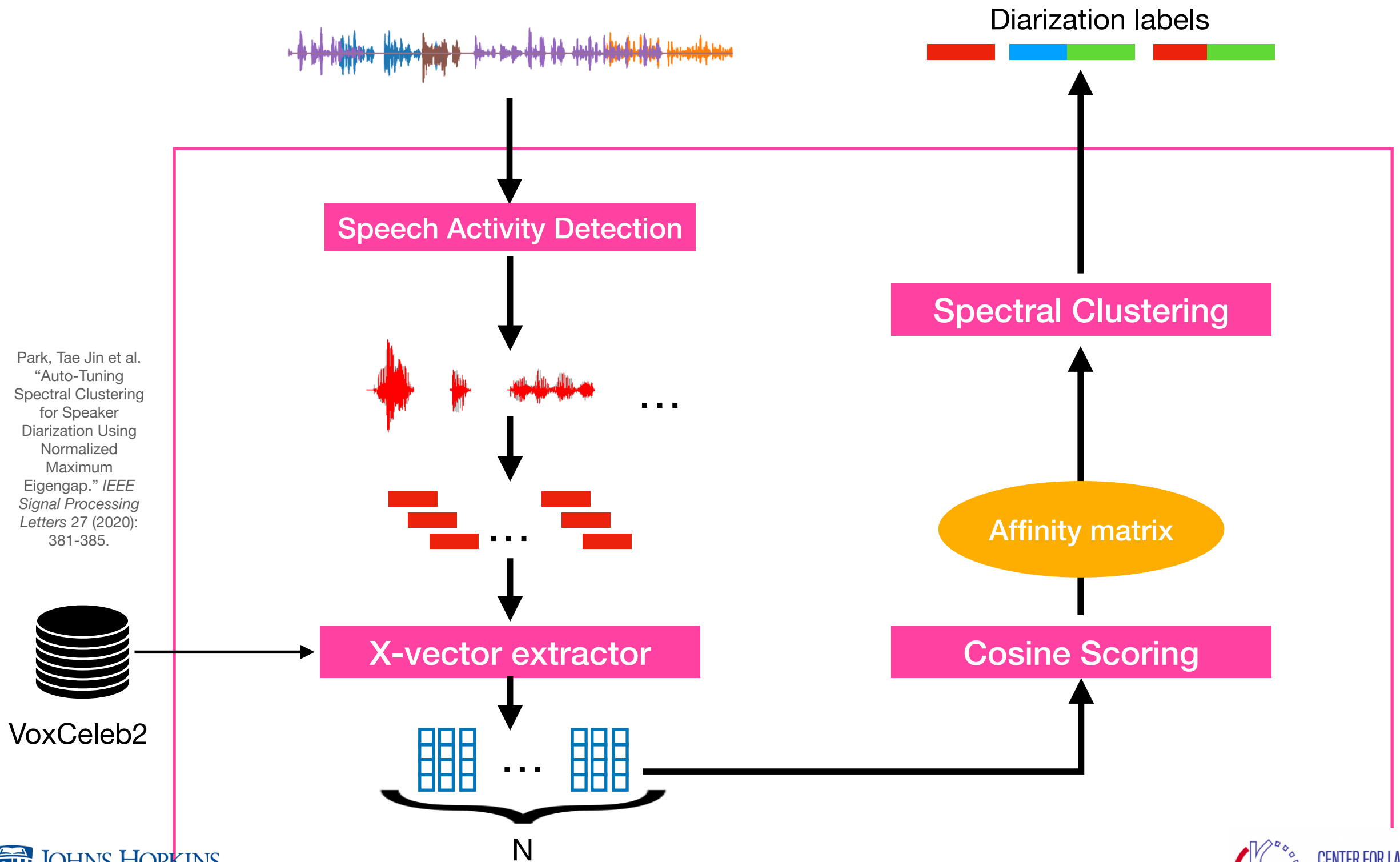


...cannot handle **overlapping** speech

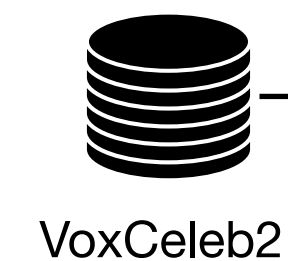
To solve this problem:



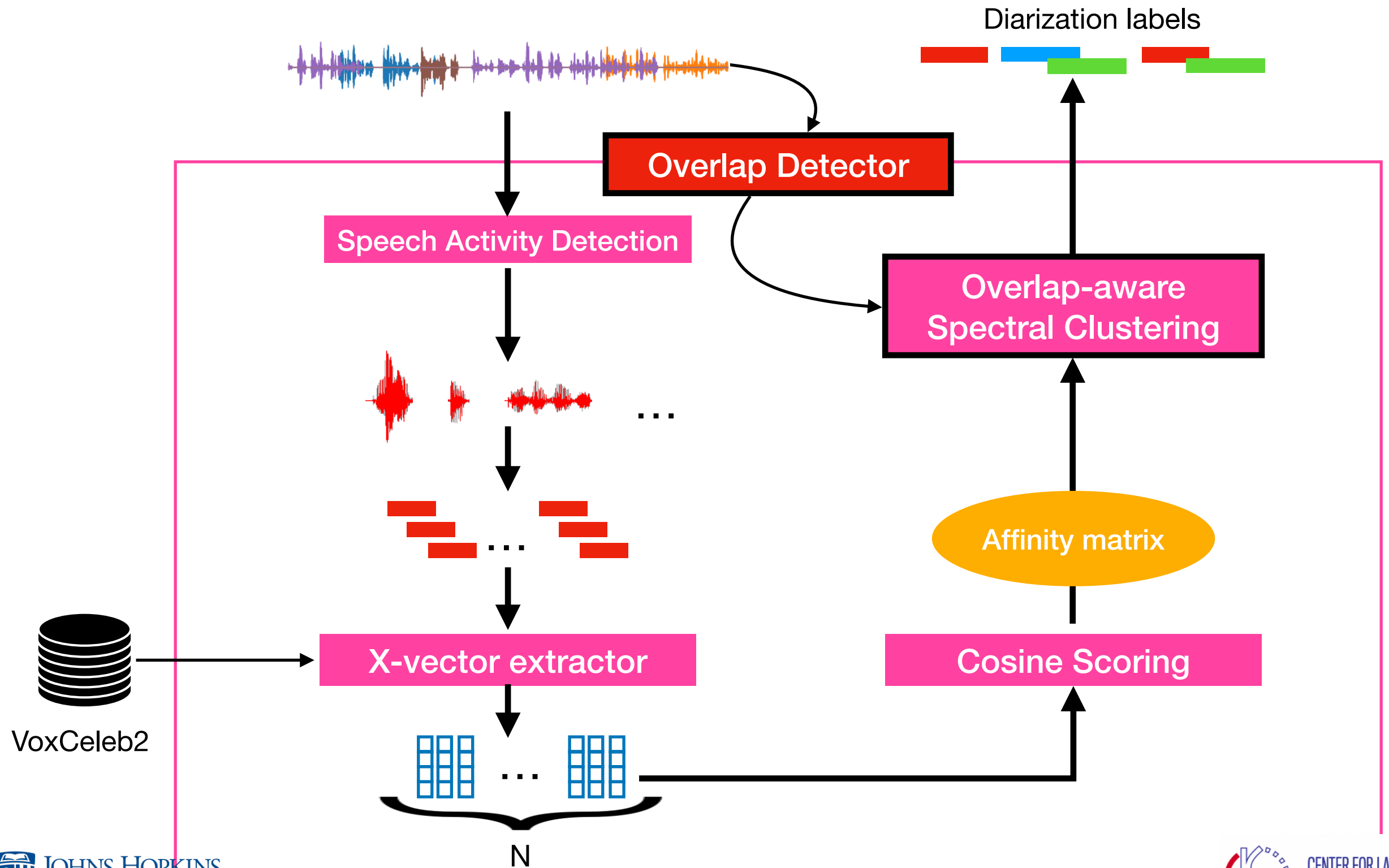
Baseline spectral clustering



Park, Tae Jin et al.
 "Auto-Tuning Spectral Clustering for Speaker Diarization Using Normalized Maximum Eigengap." *IEEE Signal Processing Letters* 27 (2020): 381-385.



Overlap-aware spectral clustering...



...using new problem formulation

$$\begin{aligned} \max \epsilon(X) &= \frac{1}{K} \sum_{k=1}^K \frac{X_k^T \mathbf{A} X_k}{X_k^T \mathbf{D} X_k} \\ \text{s.t. } X &\in \{0, 1\}^{N \times K} \\ X \mathbf{1}_K &= \mathbf{1}_N + \mathbf{v}_{OL} \end{aligned}$$



Next, we seek to obtain a discrete approximation for Z . First, we note from equation (7) that

$$X = f^{-1}(Z) = \text{Diag}(\text{diag}^{-\frac{1}{2}}(ZZ^T)) Z, \quad (11)$$

Using this transformation, we can characterize the solution obtained in equation (10) as

$$\{\tilde{X}^* R : R^T R = I_K, \tilde{X}^* = f^{-1}(Z^*)\}. \quad (12)$$

Now, our discretization problem is to find an X which approximates $\tilde{X}^* R$ for some orthonormal R , such that the discrete constraints from problem (6). This problem is formulated as

$$\begin{aligned} \min \phi(X, R) &= \|X - \tilde{X}^* R\|^2 \\ \text{s.t. } X &\in \{0, 1\}^{N \times K} \\ X \mathbf{1}_K &= \mathbf{1}_N \end{aligned} \quad (13)$$

It is difficult to solve this problem, so we optimize over X and R . Suppose we are given some R , then problem (13) reduces to

$$\begin{aligned} \min \phi(X) &= \|X - \tilde{X}^* R\|^2 \\ \text{s.t. } X &\in \{0, 1\}^{N \times K}, \\ X \mathbf{1}_K &= \mathbf{1}_N. \end{aligned} \quad (14)$$

The optimal solution to this problem is given by non-maximal suppression, i.e.,

$$X^*(i, l) = \left\langle l = \arg \max_{k \in [K]} \tilde{X}(i, k) \right\rangle, \quad i \in \{1, \dots, N\}. \quad (15)$$

Intuitively, we set the largest entry in each row as 1 and zero out all the others. This ensures that each sample belongs to exactly 1 cluster. In the next section, we will see how to reformulate problem (14) for the case when some samples can belong to more than one clusters.

Next, we fix X^* and solve the following problem for R^* :

$$\begin{aligned} \min \phi(R) &= \|X^* - \tilde{X}^* R\|^2 \\ \text{s.t. } R^T R &= I_K. \end{aligned} \quad (16)$$

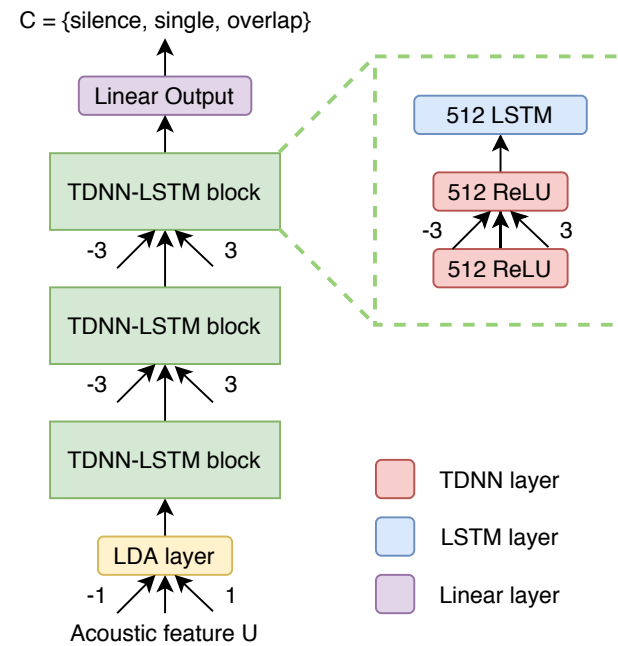
The solution to this problem is given by

$$R^* = \tilde{U} U^T, \quad (17)$$


$$\begin{aligned} X^*(i, l) &= \left\langle l = \arg \max_{k \in [K]} \tilde{X}(i, k) \right\rangle \\ &+ \mathbf{v}_{OL}^{(i)} \times \left\langle l = \arg \max_{k \in [K]} \tilde{X}(i, k) \right\rangle \end{aligned}$$

*Manuscript in preparation for SLT 2021

Diarization results on LibriCSS eval set



Overlap detector:

Precision = 96.3%

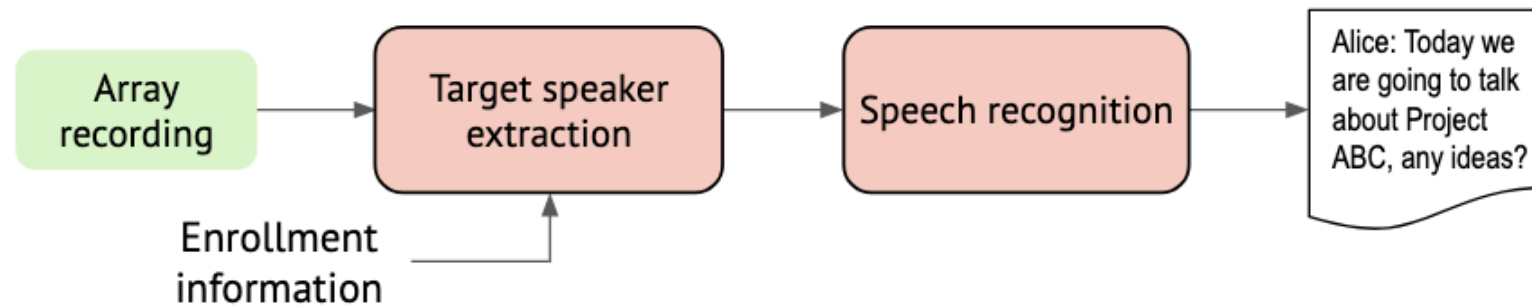
Recall = 83.8%

| Method | Avg. DER | DER on 40% |
|----------------------------------|----------|------------|
| AHC/PLDA | 16.3 | 28.8 |
| Spectral/cosine | 12.6 | 23.4 |
| Overlap-aware SC | 9.3 | 15.2 |
| + oracle overlap detector | 8.8 | 14.4 |

*All DER reported using oracle VAD

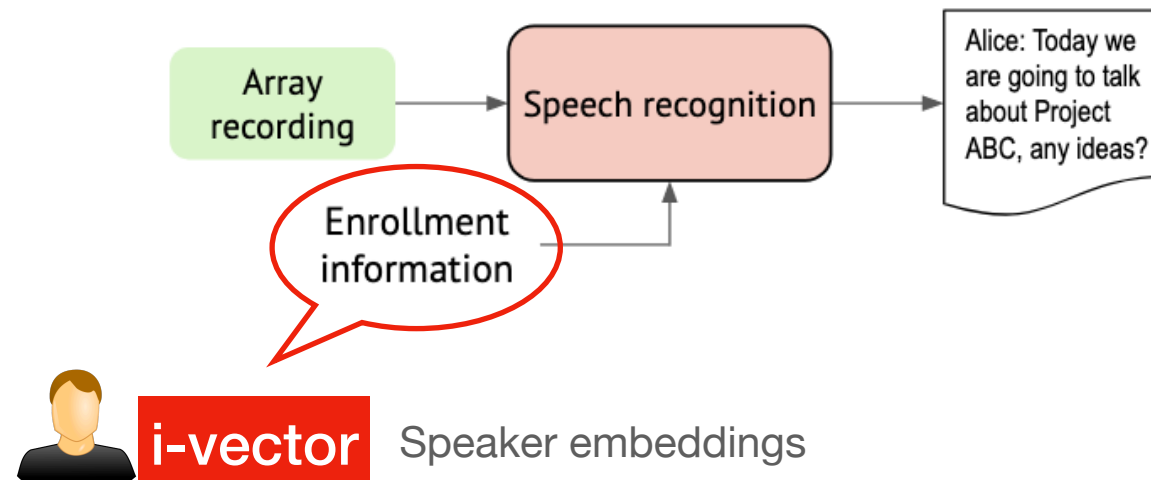
Target speaker ASR

Target-speaker extraction (Katka's talk):



Žmolíková, Kateřina et al. "SpeakerBeam: Speaker Aware Neural Network for Target Speaker Extraction in Speech Mixtures." *IEEE Journal of Selected Topics in Signal Processing* 13 (2019): 800-814.

Target-speaker ASR (this talk):



How to train your Target speaker ASR

| Method | Details | Avg. WER | WER on 40% overlap region |
|------------------------|-----------------------|----------|---------------------------|
| Speaker-independent AM | Librispeech with RIRs | 27.88 | 47.7 |
| | | | |
| | | | |
| | | | |
| | | | |



**17-layer
TDNNF**

LF-MMI
objective

Hybrid HMM-DNN with RNNLM rescoring

Povey, Daniel et al. "Purely Sequence-Trained Neural Networks for ASR Based on Lattice-Free MMI." *INTERSPEECH* (2016).

Xu, Hainan et al. "A Pruned RNNLM Lattice-Rescoring Algorithm for Automatic Speech Recognition." *2018 IEEE ICASSP* (2018).

i. Train with simulated overlapping data

| Method | Details | Avg. WER | WER on 40% overlap region |
|------------------------|------------------------------|----------|---------------------------|
| Speaker-independent AM | Librispeech with RIRs | 27.88 | 47.7 |
| Speaker-independent AM | + simulated overlapping data | 18.00 | 25.44 |
| | | | |
| | | | |
| | | | |

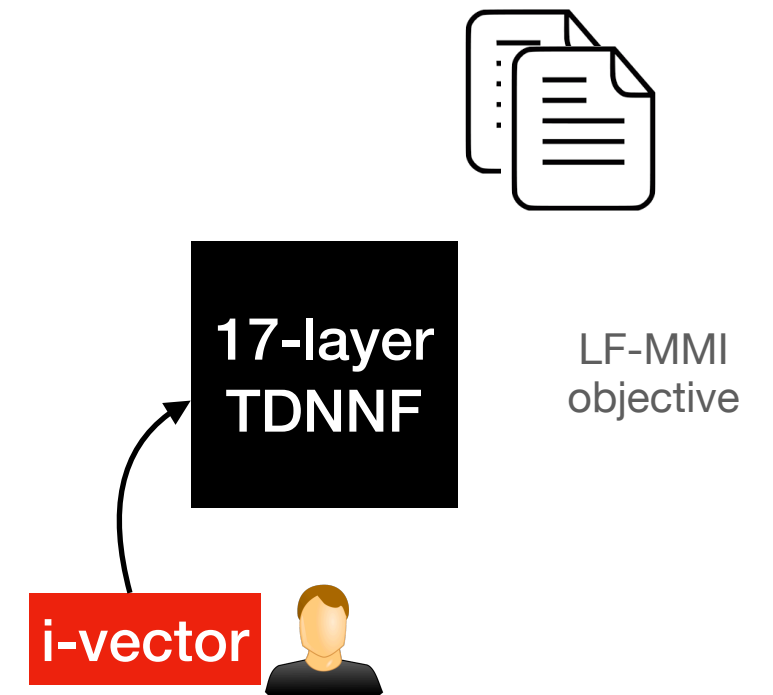


**17-layer
TDNNF**

LF-MMI
objective

ii. Add speaker i-vectors to guide ASR

| Method | Details | Avg. WER | WER on 40% overlap region |
|------------------------|------------------------------|----------|---------------------------|
| Speaker-independent AM | Librispeech with RIRs | 27.88 | 47.7 |
| Speaker-independent AM | + simulated overlapping data | 18.00 | 25.44 |
| + speaker i-vectors | - | 16.99 | 21.74 |
| | | | |
| | | | |

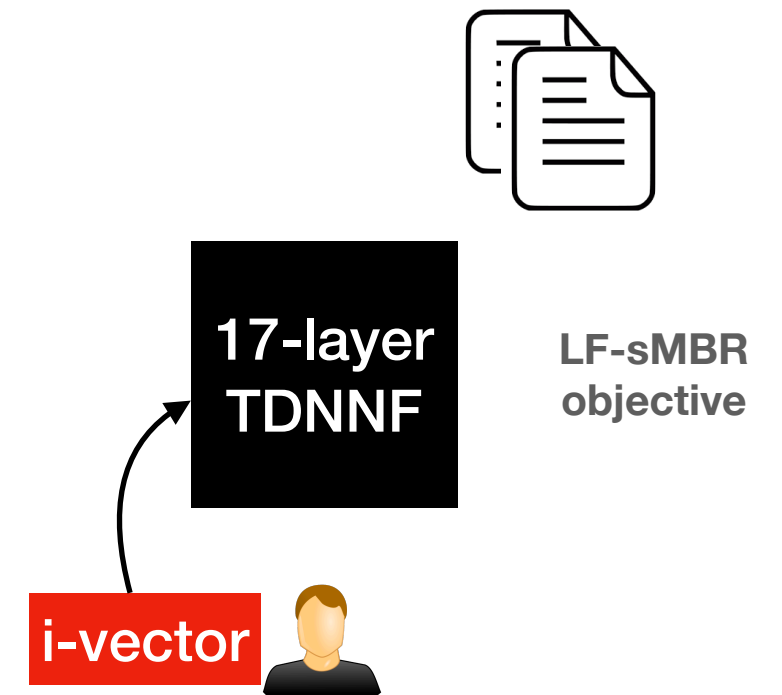


*i-vectors estimated from clean enrollment utterances

iii. Discriminative training

| Method | Details | Avg. WER | WER on 40% overlap region |
|------------------------|------------------------------|----------|---------------------------|
| Speaker-independent AM | Librispeech with RIRs | 27.88 | 47.7 |
| Speaker-independent AM | + simulated overlapping data | 18.00 | 25.44 |
| + speaker i-vectors | - | 16.99 | 21.74 |
| + speaker i-vectors | Discriminative training | 16.96 | 22.85 |
| | | | |

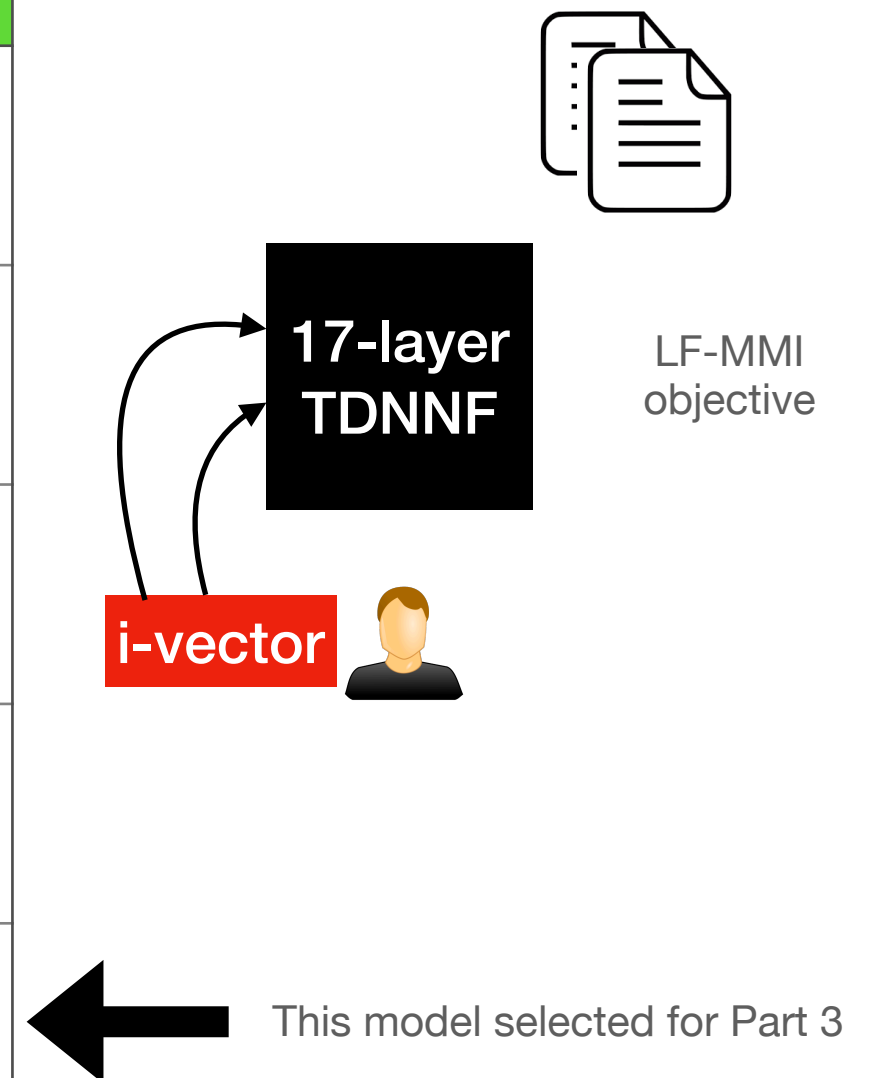
*i-vectors estimated from clean enrollment utterances



Kanda, Naoyuki et al. "Lattice-free State-level Minimum Bayes Risk Training of Acoustic Models." *INTERSPEECH* (2018).

iv. Add i-vectors to several layers

| Method | Details | Avg. WER | WER on 40% overlap region |
|------------------------|------------------------------|--------------|---------------------------|
| Speaker-independent AM | Librispeech with RIRs | 27.88 | 47.7 |
| Speaker-independent AM | + simulated overlapping data | 18.00 | 25.44 |
| + speaker i-vectors | - | 16.99 | 21.74 |
| + speaker i-vectors | Discriminative training | 16.96 | 22.85 |
| + deeper integration | - | 16.56 | 21.12 |

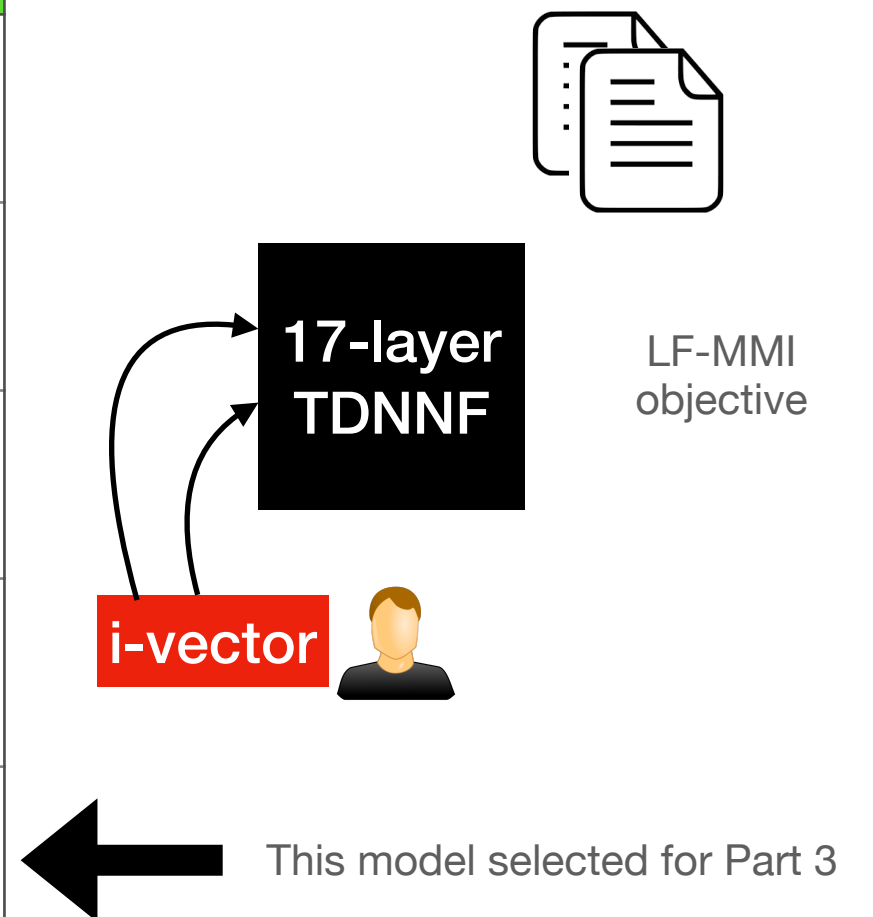


*i-vectors estimated from clean enrollment utterances

v. Next steps?

| Method | Details | Avg. WER | WER on 40% overlap region |
|------------------------|------------------------------|--------------|---------------------------|
| Speaker-independent AM | Librispeech with RIRs | 27.88 | 47.7 |
| Speaker-independent AM | + simulated overlapping data | 18.00 | 25.44 |
| + speaker i-vectors | - | 16.99 | 21.74 |
| + speaker i-vectors | Discriminative training | 16.96 | 22.85 |
| + deeper integration | - | 16.56 | 21.12 |
| + deeper integration | Discriminative training | ? | ? |

*i-vectors estimated from clean enrollment utterances



Combining Parts 1 and 2

Results for TS-ASR are on oracle segments with i-vectors estimated from enrollment utterances.

- i. Estimate i-vectors from **recording** instead of **enrollment** utterances
- ii. Use segments obtained from **Diarization** instead of **oracle** segments

Combining Parts 1 and 2

High overlap



| ASR | i-vectors obtained from | Segments | DER | cpWER | cpWER (OV40) |
|-----------------------------------|-------------------------|-----------------------------------|-------|-------|--------------|
| Baseline (speaker independent AM) | no i-vectors used | Baseline Diarization (no overlap) | 18.28 | 32.72 | 48.05 |
| Baseline (speaker independent AM) | no i-vectors used | oracle | 0.00 | 27.88 | 47.70 |
| TS-ASR | enrollment | oracle | 0.00 | 16.56 | 21.12 |
| | | | | | |
| | | | | | |
| | | | | | |

*All results are on evaluation set

a) Using recording to estimate i-vectors

| ASR | i-vectors obtained from | Segments | DER | cpWER | cpWER (OV40) |
|-----------------------------------|-------------------------|-----------------------------------|-------|-------|--------------|
| Baseline (speaker independent AM) | no i-vectors used | Baseline Diarization (no overlap) | 18.28 | 32.72 | 48.05 |
| Baseline (speaker independent AM) | no i-vectors used | oracle | 0.00 | 27.88 | 47.70 |
| TS-ASR | enrollment | oracle | 0.00 | 16.56 | 21.12 |
| TS-ASR | recording | oracle | 0.00 | 19.29 | 24.55 |
| | | | | | |
| | | | | | |

*All results are on evaluation set

b) Using baseline diarization to get segments

| ASR | i-vectors obtained from | Segments | DER | cpWER | cpWER (OV40) |
|-----------------------------------|-------------------------|-----------------------------------|-------|-------|--------------|
| Baseline (speaker independent AM) | no i-vectors used | Baseline Diarization (no overlap) | 18.28 | 32.72 | 48.05 |
| Baseline (speaker independent AM) | no i-vectors used | oracle | 0.00 | 27.88 | 47.70 |
| TS-ASR | enrollment | oracle | 0.00 | 16.56 | 21.12 |
| TS-ASR | recording | oracle | 0.00 | 19.29 | 24.55 |
| TS-ASR | recording | Baseline Diarization (no overlap) | 18.28 | 32.07 | 42.99 |
| | | | | | |

*All results are on evaluation set

c) Using overlap-aware diarization

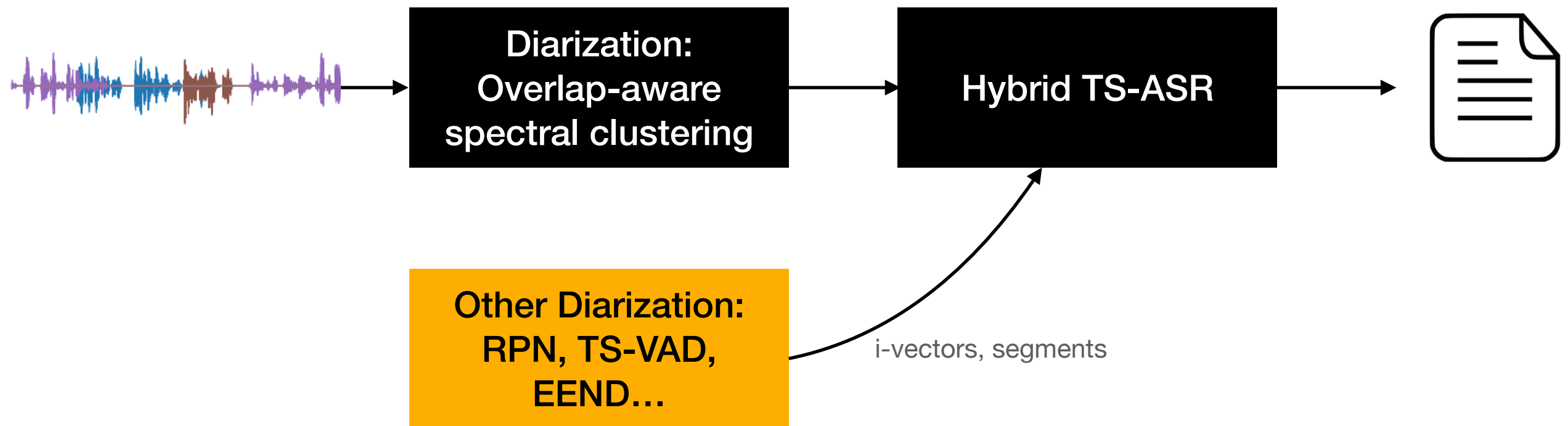
| ASR | i-vectors obtained from | Segments | DER | cpWER | cpWER (OV40) |
|-----------------------------------|-------------------------|-----------------------------------|-------|-------|--------------|
| Baseline (speaker independent AM) | no i-vectors used | Baseline Diarization (no overlap) | 18.28 | 32.72 | 48.05 |
| Baseline (speaker independent AM) | no i-vectors used | oracle | 0.00 | 27.88 | 47.70 |
| TS-ASR | enrollment | oracle | 0.00 | 16.56 | 21.12 |
| TS-ASR | recording | oracle | 0.00 | 19.29 | 24.55 |
| TS-ASR | recording | Baseline Diarization (no overlap) | 18.28 | 32.07 | 42.99 |
| TS-ASR | recording | overlap-aware diarization | 15.15 | 30.36 | 39.91 |
| | | + oracle VAD | 9.34 | 25.36 | 35.68 |

*All results are on evaluation set

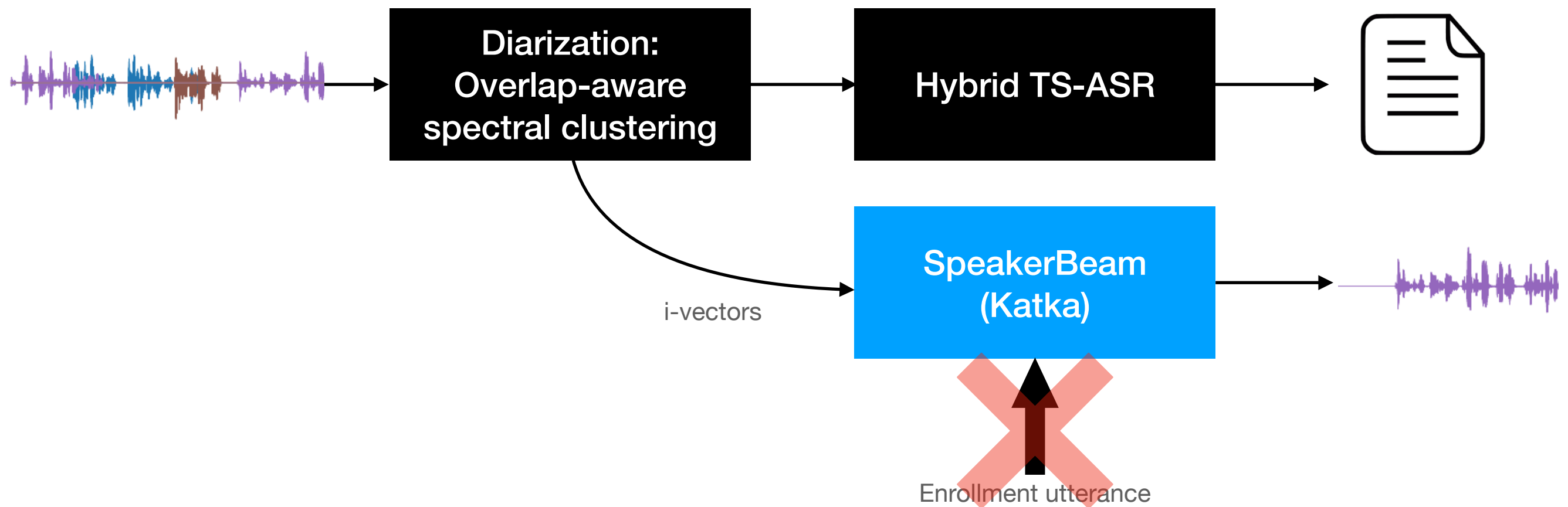
How this fits in the bigger picture



How this fits in the bigger picture



How this fits in the bigger picture



How this fits in the bigger picture

