

Multi-class Spectral Clustering with Overlaps for Speaker Diarization

Desh Raj, Zili Huang, Sanjeev Khudanpur

Center for Language and Speech Processing, Johns Hopkins University

Motivation

What is speaker diarization?

Task of “who spoke when”

Input: *recording containing multiple speakers*

Output: *homogeneous speaker segments*

Motivation

What is speaker diarization?

Task of “who spoke when”

Input: *recording containing multiple speakers*

Output: *homogeneous speaker segments*

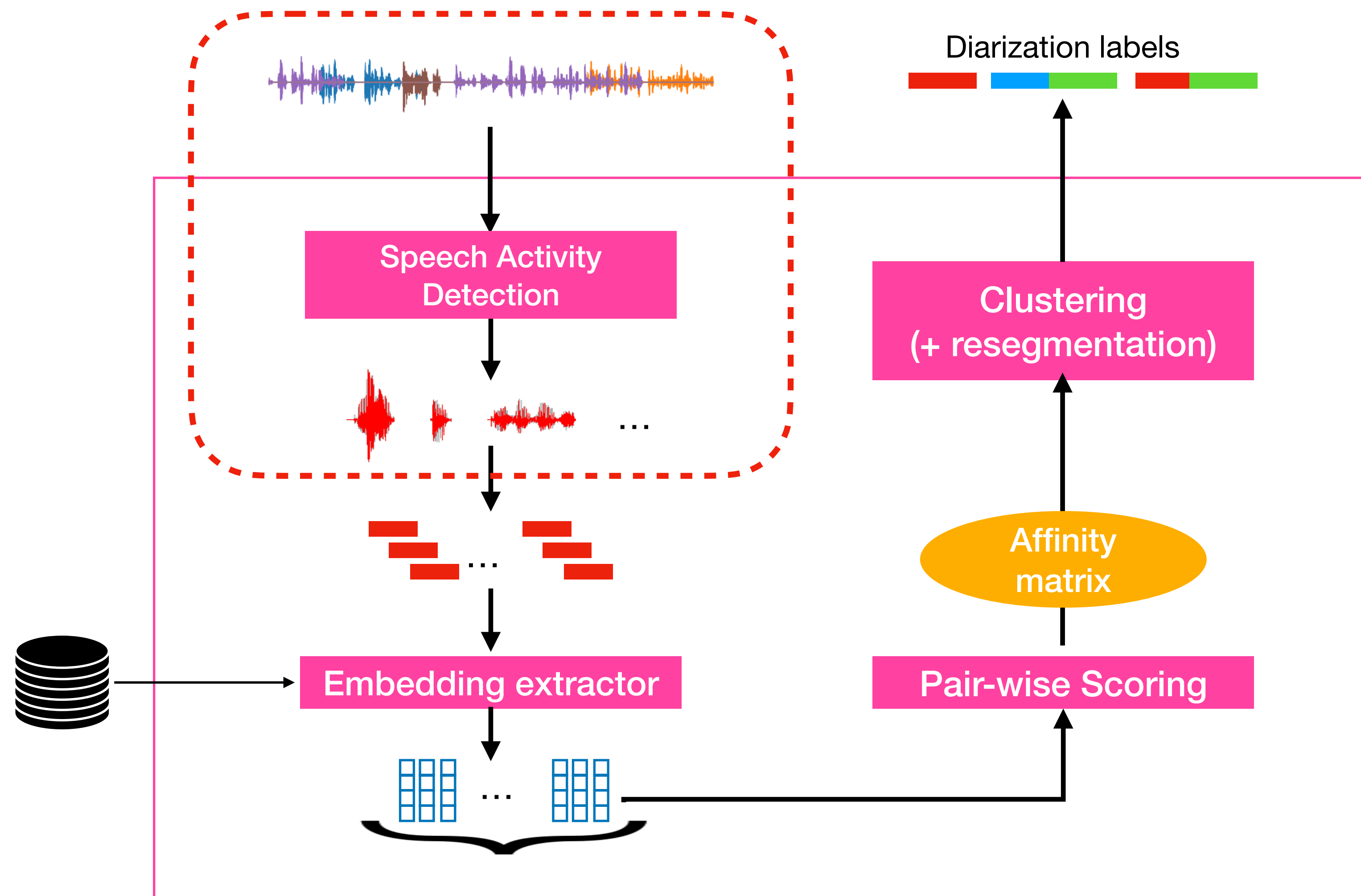
Number of speakers may be unknown

Overlapping speech may be present



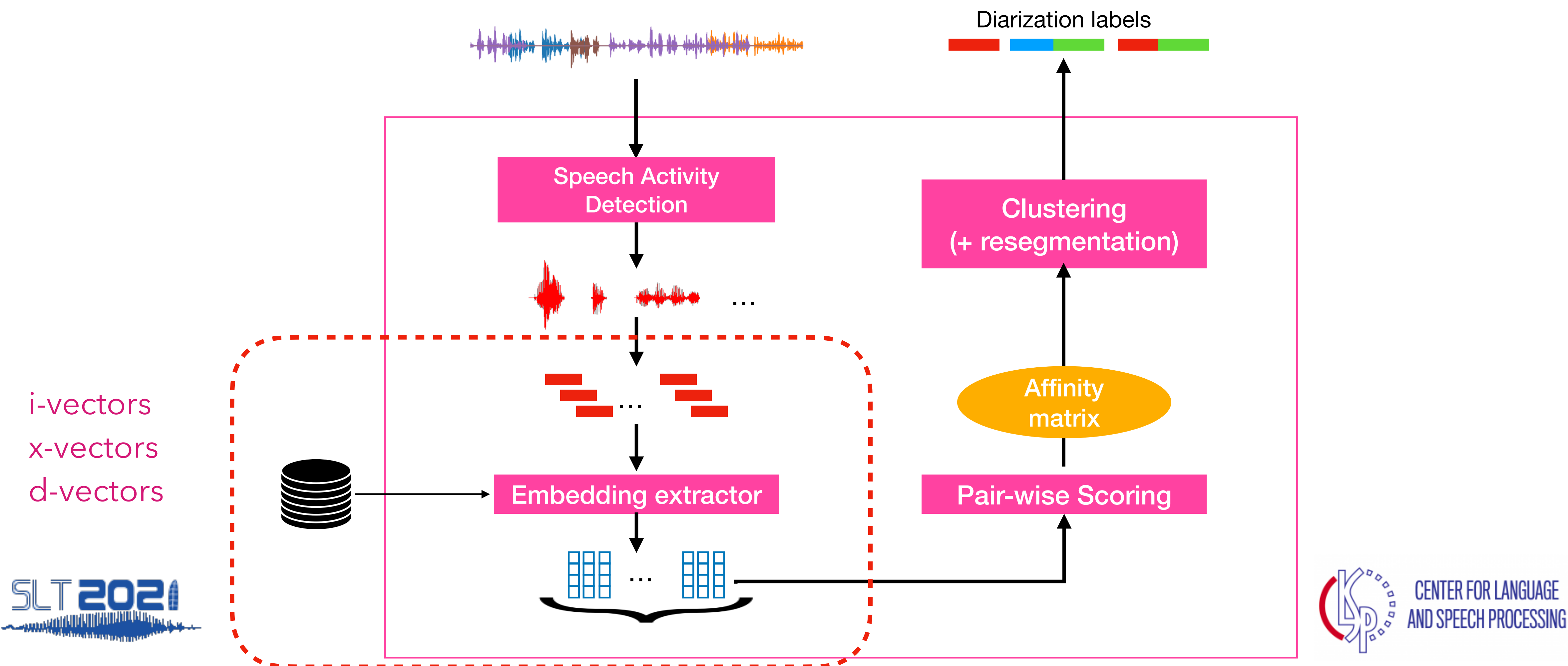
Traditional clustering-based diarization

SAD extracts speech segments from recordings



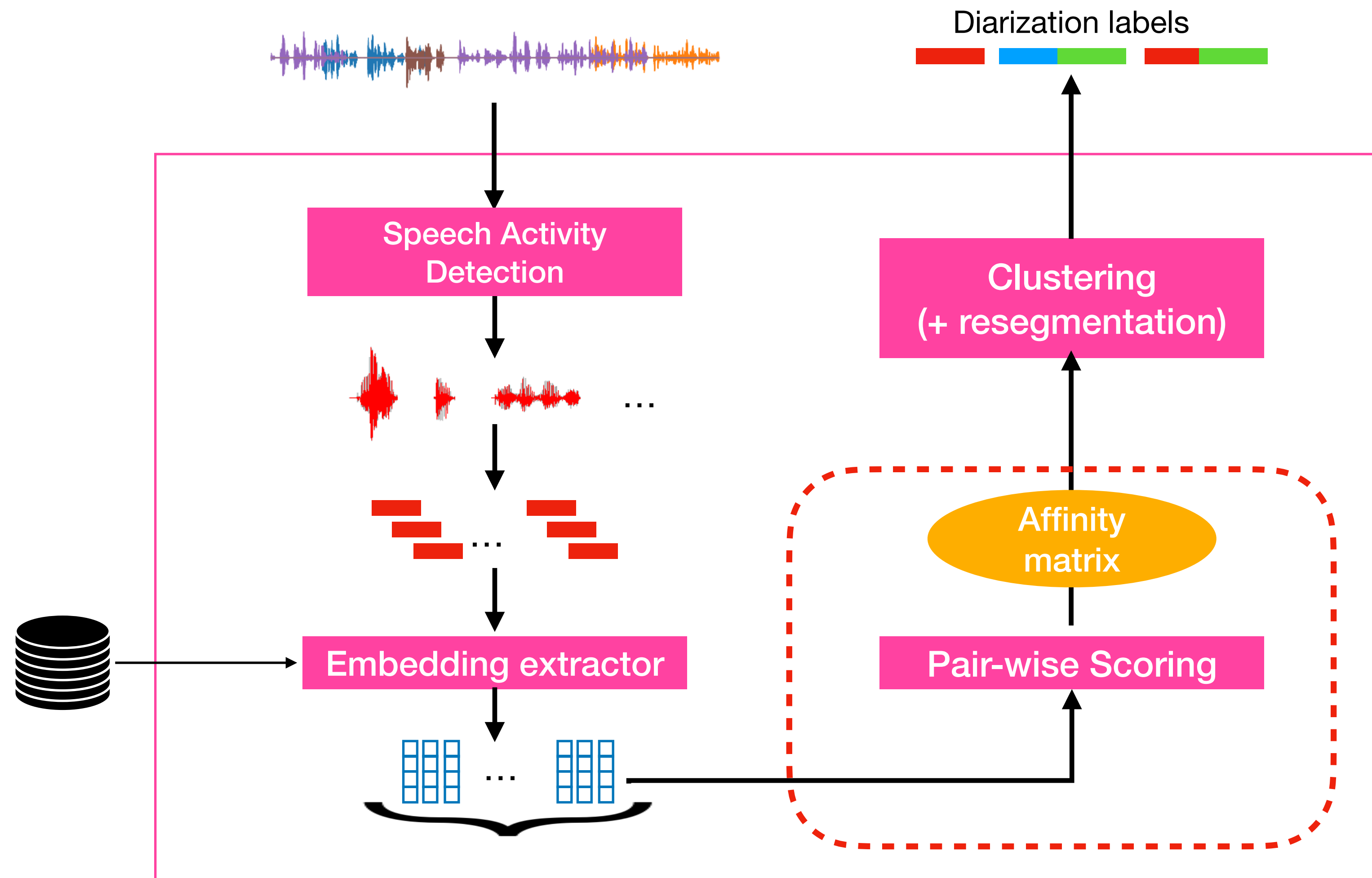
Traditional clustering-based diarization

Embeddings extracted for small subsegments



Traditional clustering-based diarization

Pair-wise scoring of subsegments



PLDA scoring
Cosine scoring

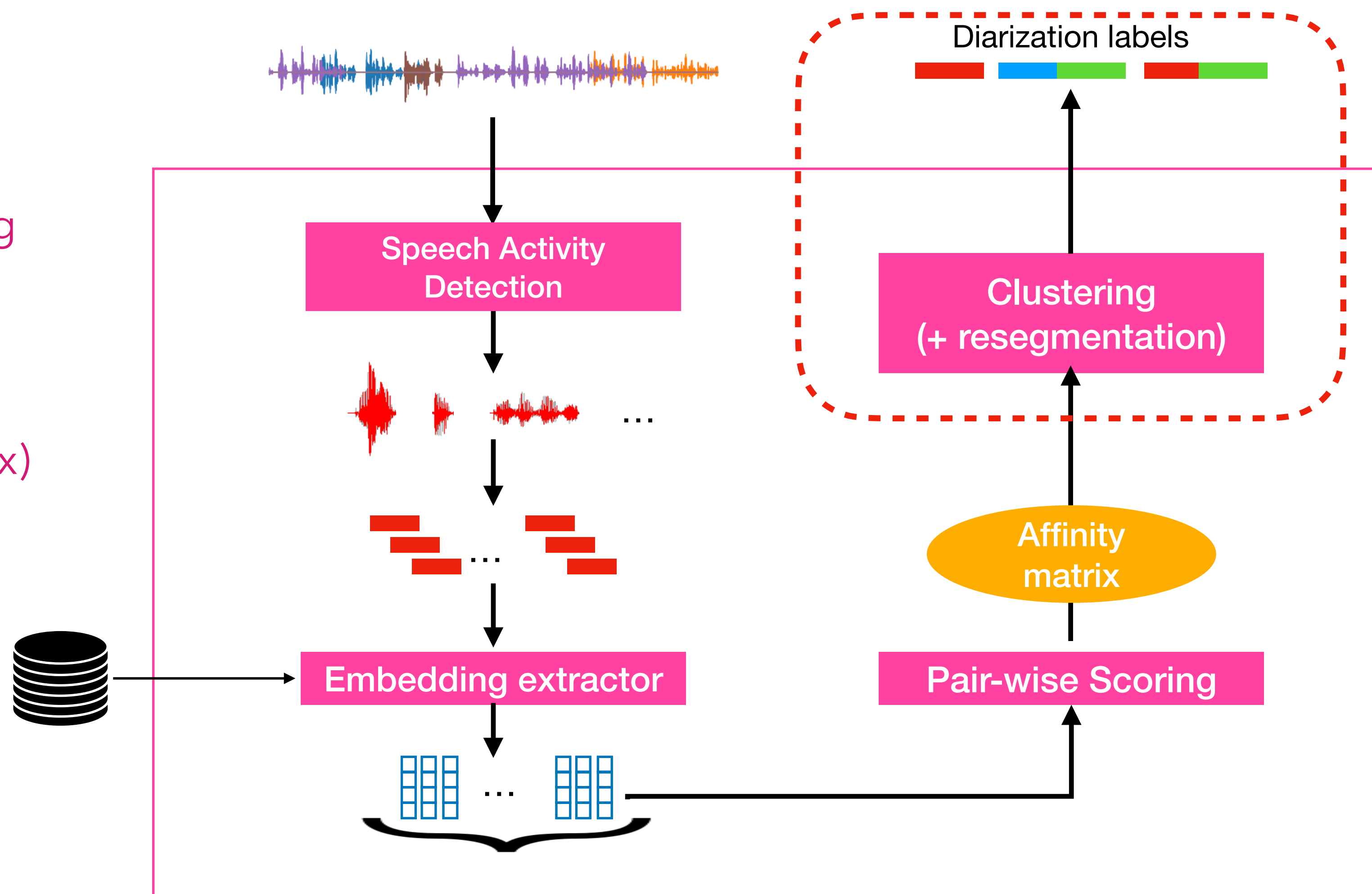
Traditional clustering-based diarization

Clustering based on the affinity matrix, followed by optional resegmentation

Agglomerative
hierarchical clustering

Spectral clustering

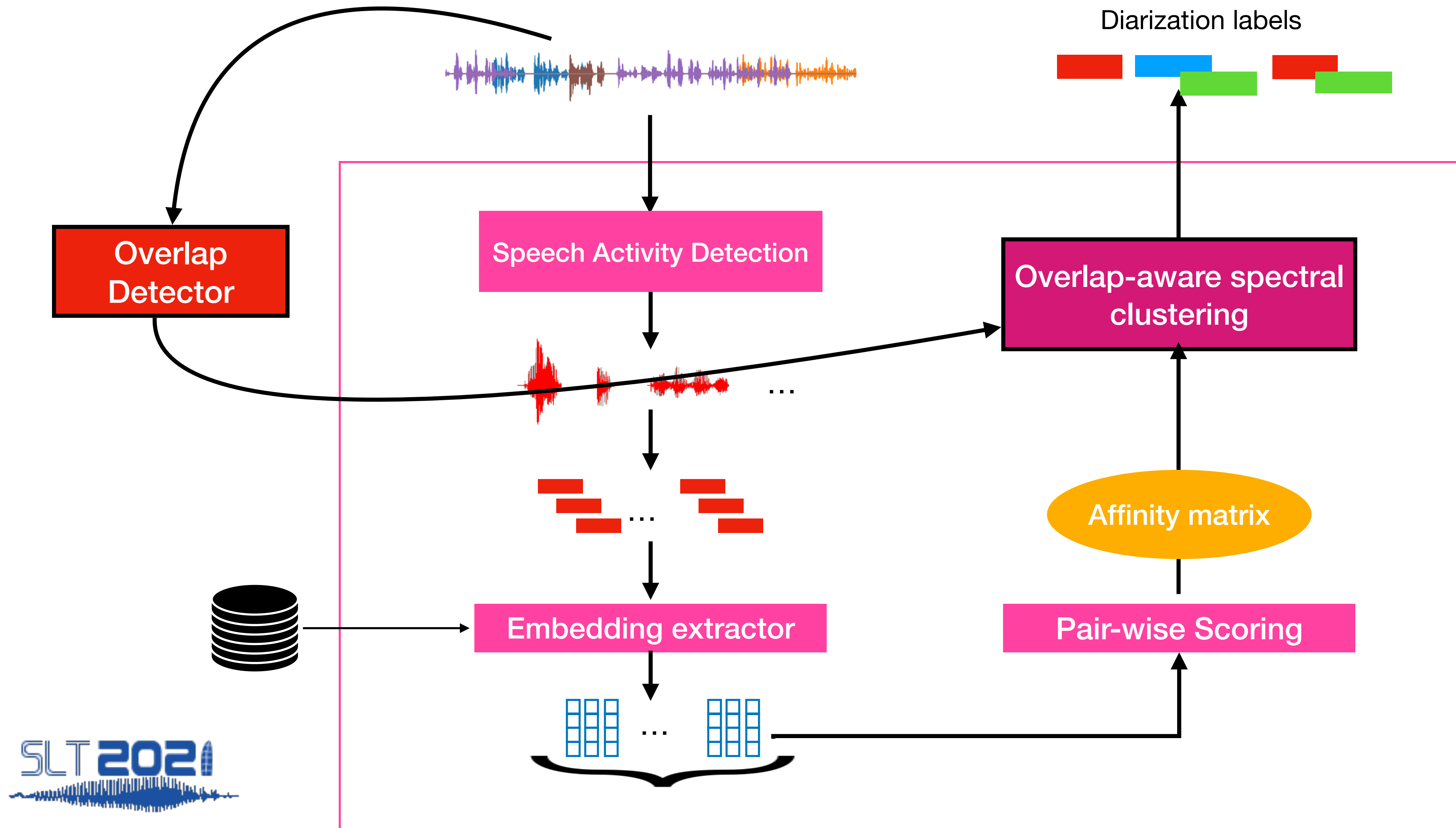
Variational Bayes (VBx)



Clustering paradigm assumes **single-speaker segments**

So overlapping speakers are completely ignored!

Overlap-aware spectral clustering



Results on AMI Mix-Headset eval

12.0% relative improvement over spectral clustering baseline

System	DER
Spectral clustering	26.9
AHC	28.3
VBx	26.2
Overlap-aware SC	24.0

Park et al., “Auto-tuning spectral clustering for speaker diarization using normalized maximum eigengap,” IEEE Signal Processing Letters, 2020.

Garcia-Romero et al., “Speaker diarization using deep neural network embeddings,” ICASSP 2017.

Díez et al., “Speaker diarization based on Bayesian HMM with eigenvoice priors,” Odyssey 2018.

Results on AMI Mix-Headset eval

Comparable with other overlap-aware diarization methods

System	DER
VB-based overlap assignment	23.8
Region proposal networks	25.5
Overlap-aware SC	24.0

Bullock, et al., “Overlap-aware diarization: resegmentation using neural end-to-end overlapped speech detection,” ICASSP 2020.

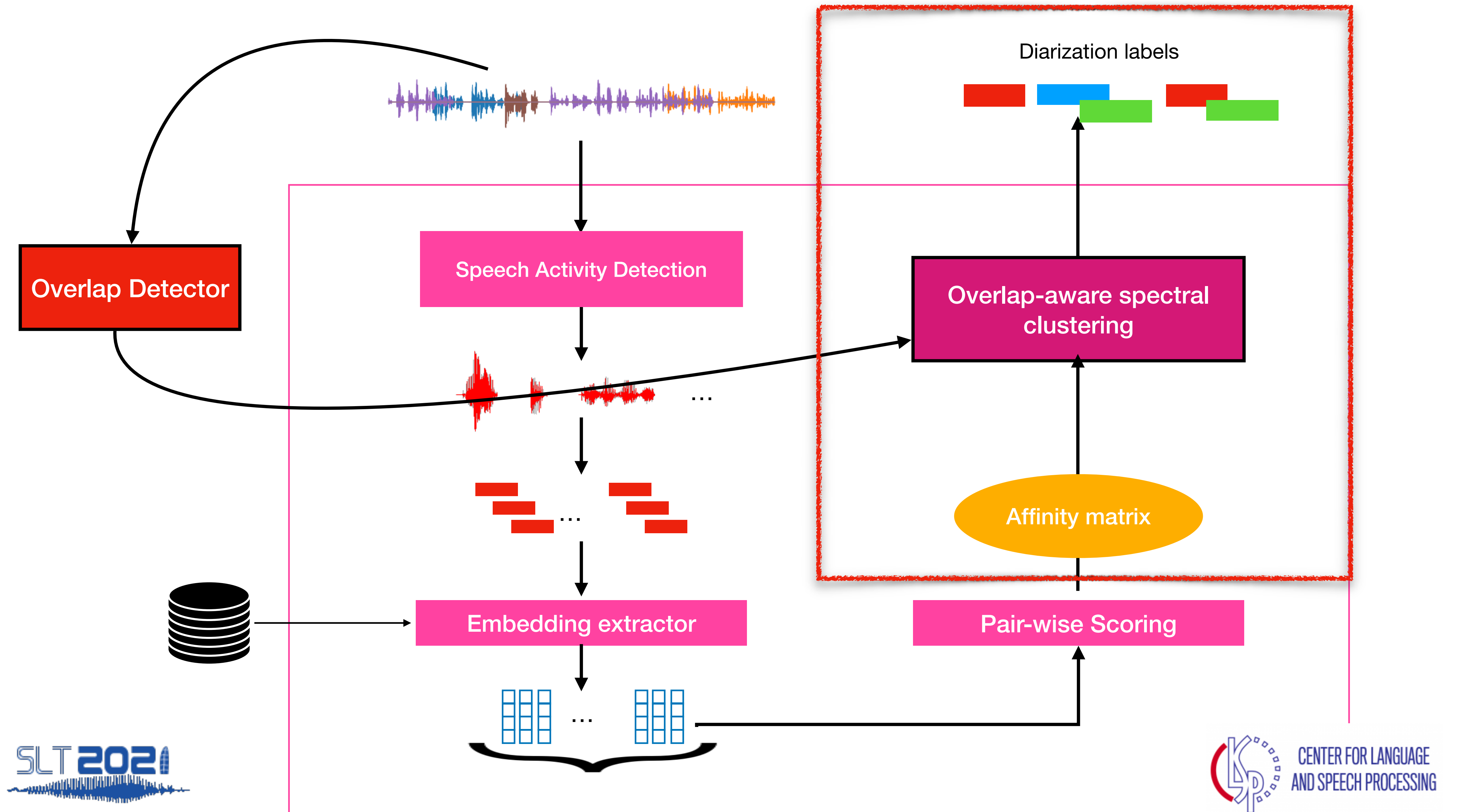
Huang et al., “Speaker diarization with region proposal network,” ICASSP 2020.

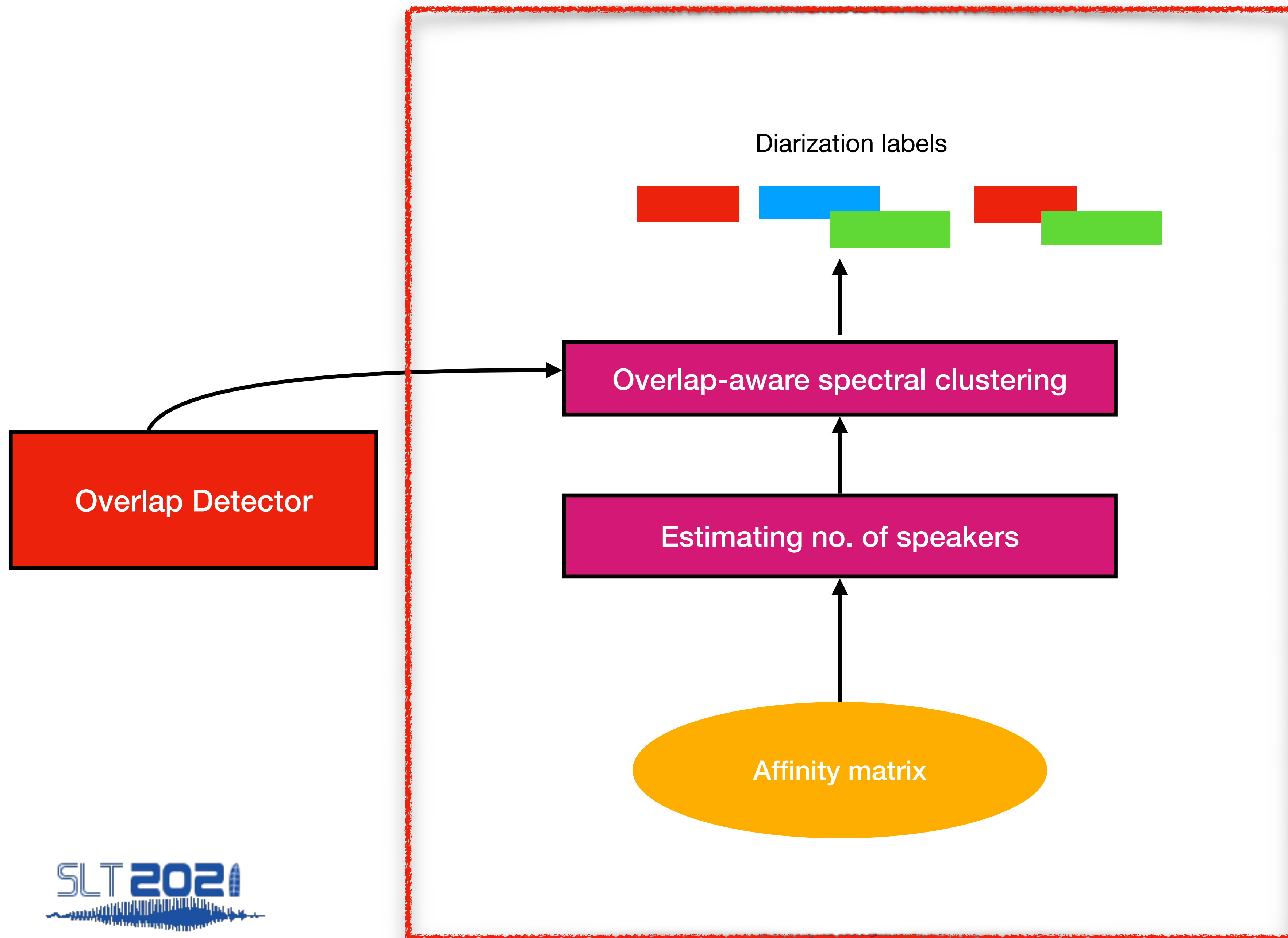
Does not require **matching training data** or **initialization** with other diarization systems.

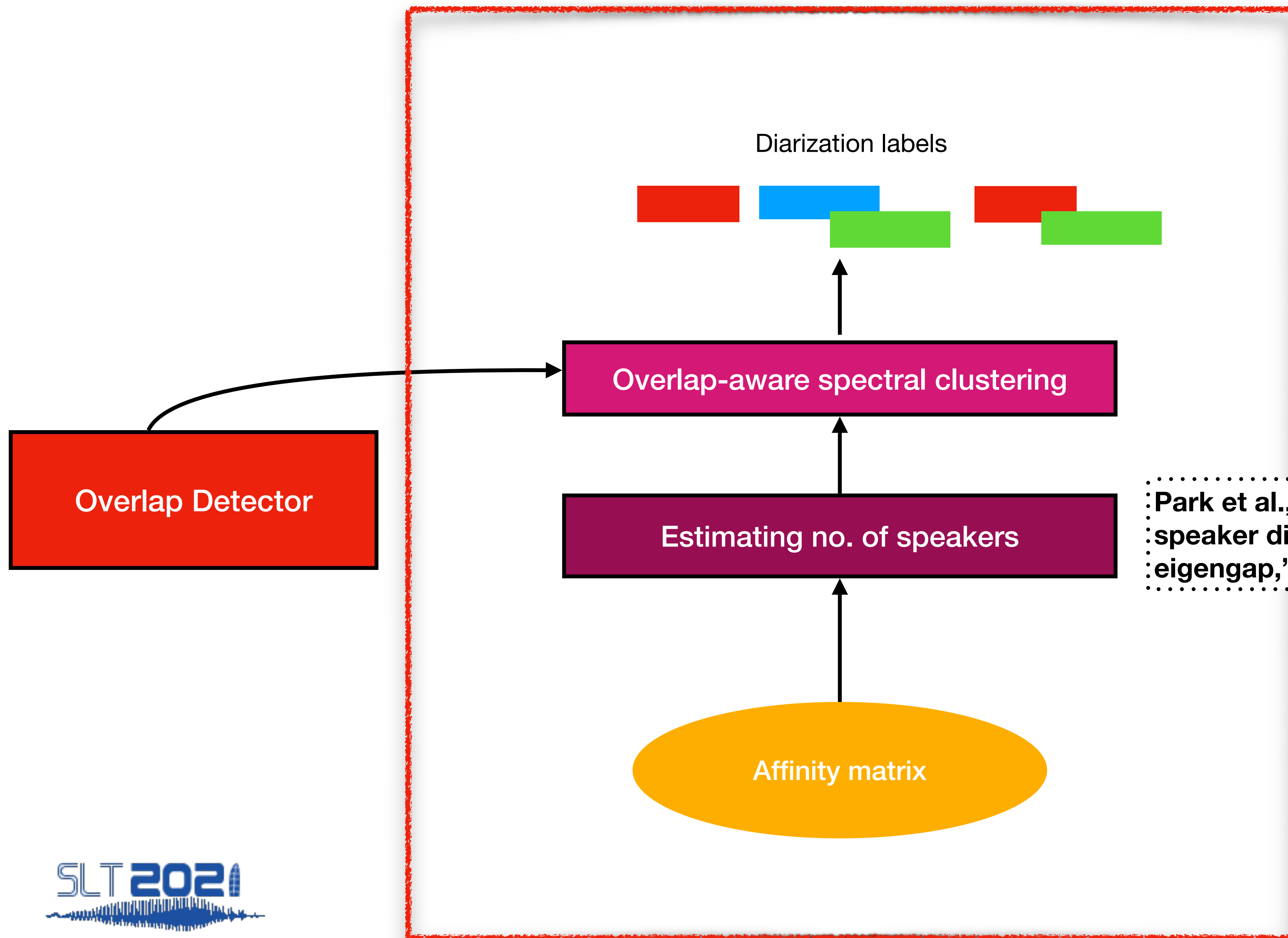
End of Highlight

Overview

- **Overlap-aware spectral clustering**
 - **Estimating number of speakers**
 - **Reformulation of the clustering problem**
 - **Incorporating overlap detector output**
- **HMM-DNN overlap detector: Overview**
- **More Results**
 - **Error analysis on AMI**



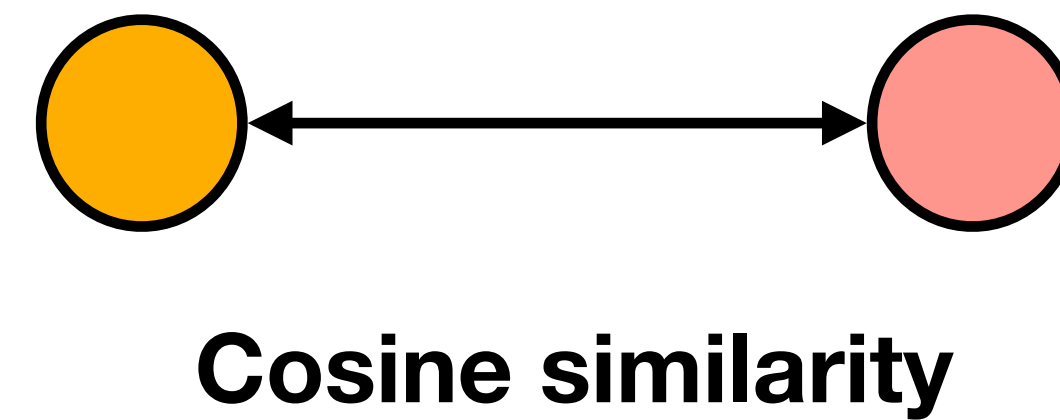
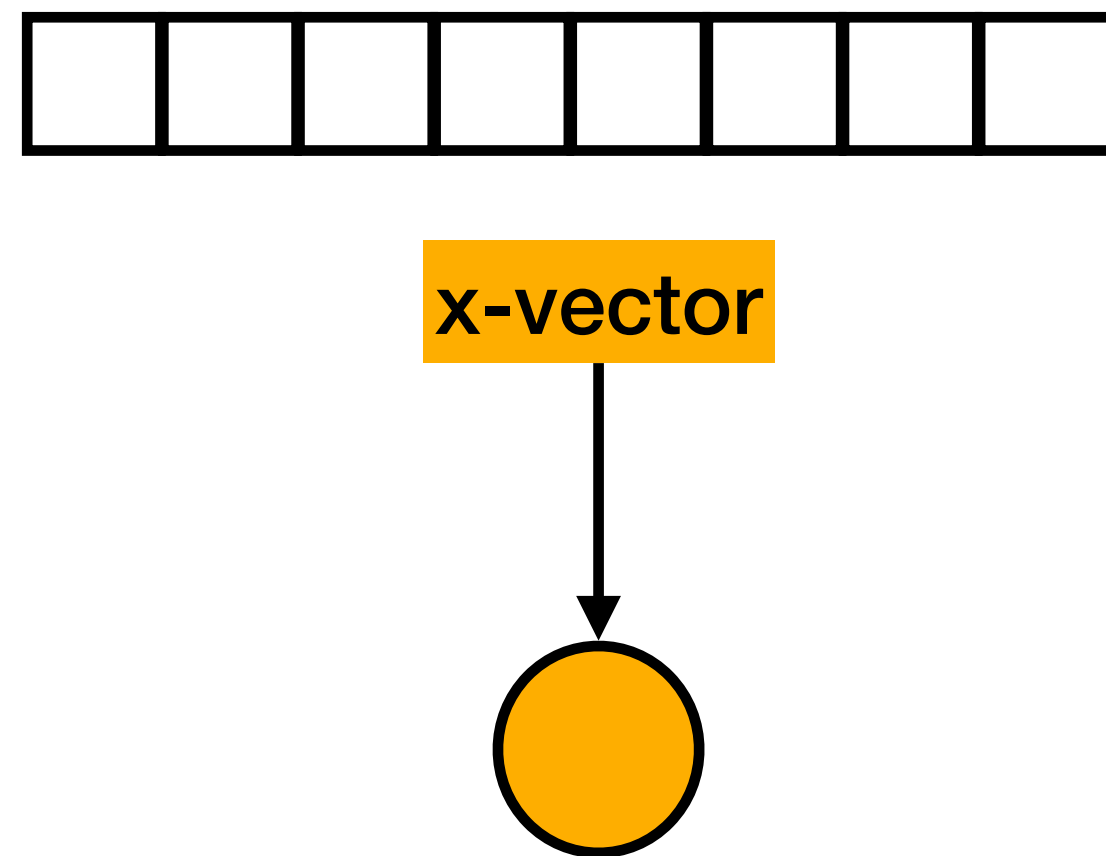




Park et al., “Auto-tuning spectral clustering for speaker diarization using normalized maximum eigengap,” *IEEE Signal Processing Letters*, 2020.

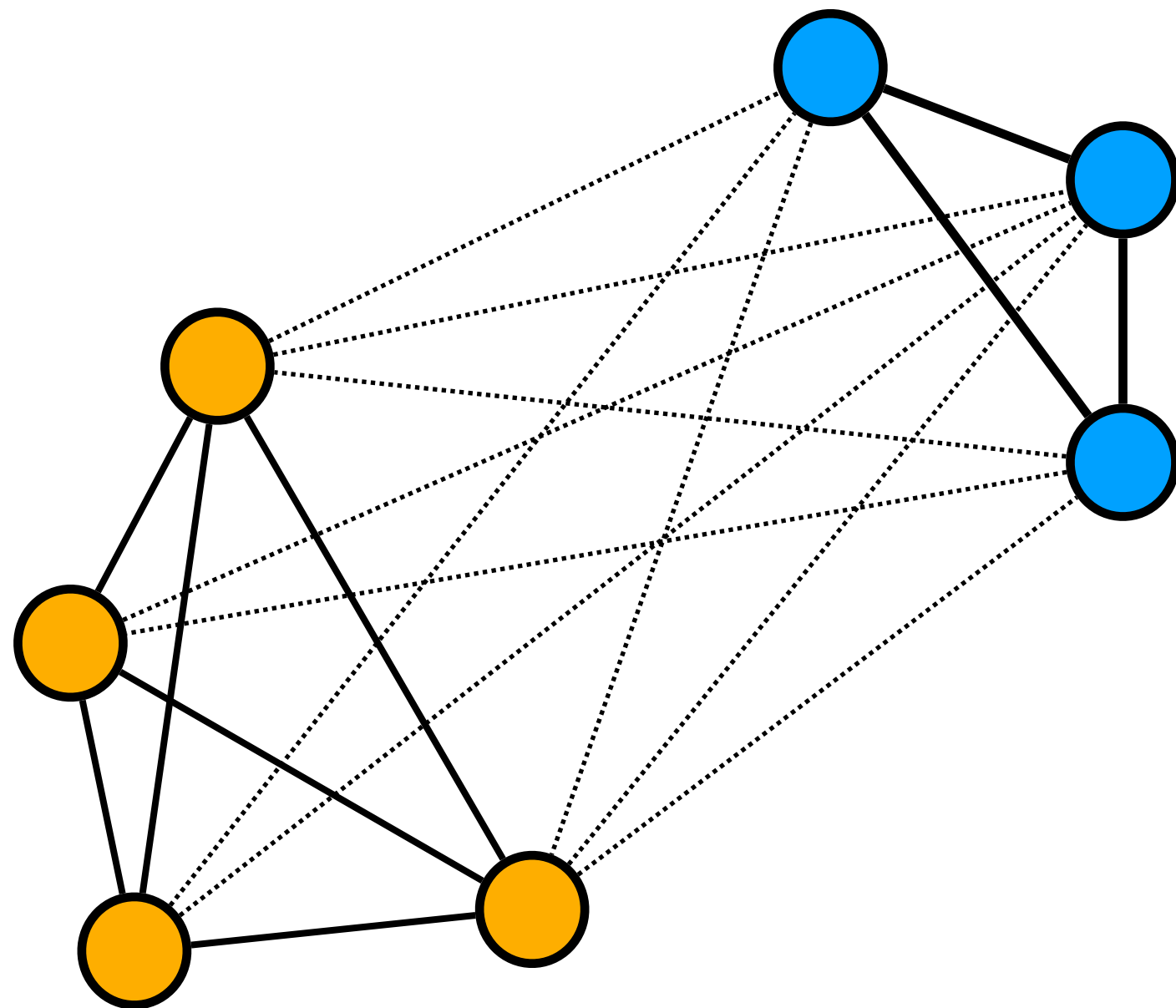
New formulation for spectral clustering

The basic clustering problem: a graph view

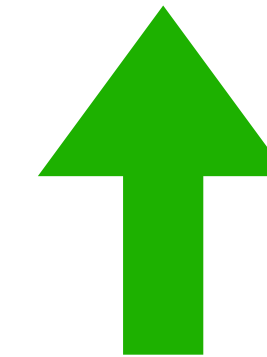


New formulation for spectral clustering

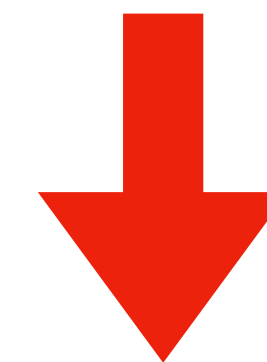
The basic clustering problem: a graph view



Edge weights within a group

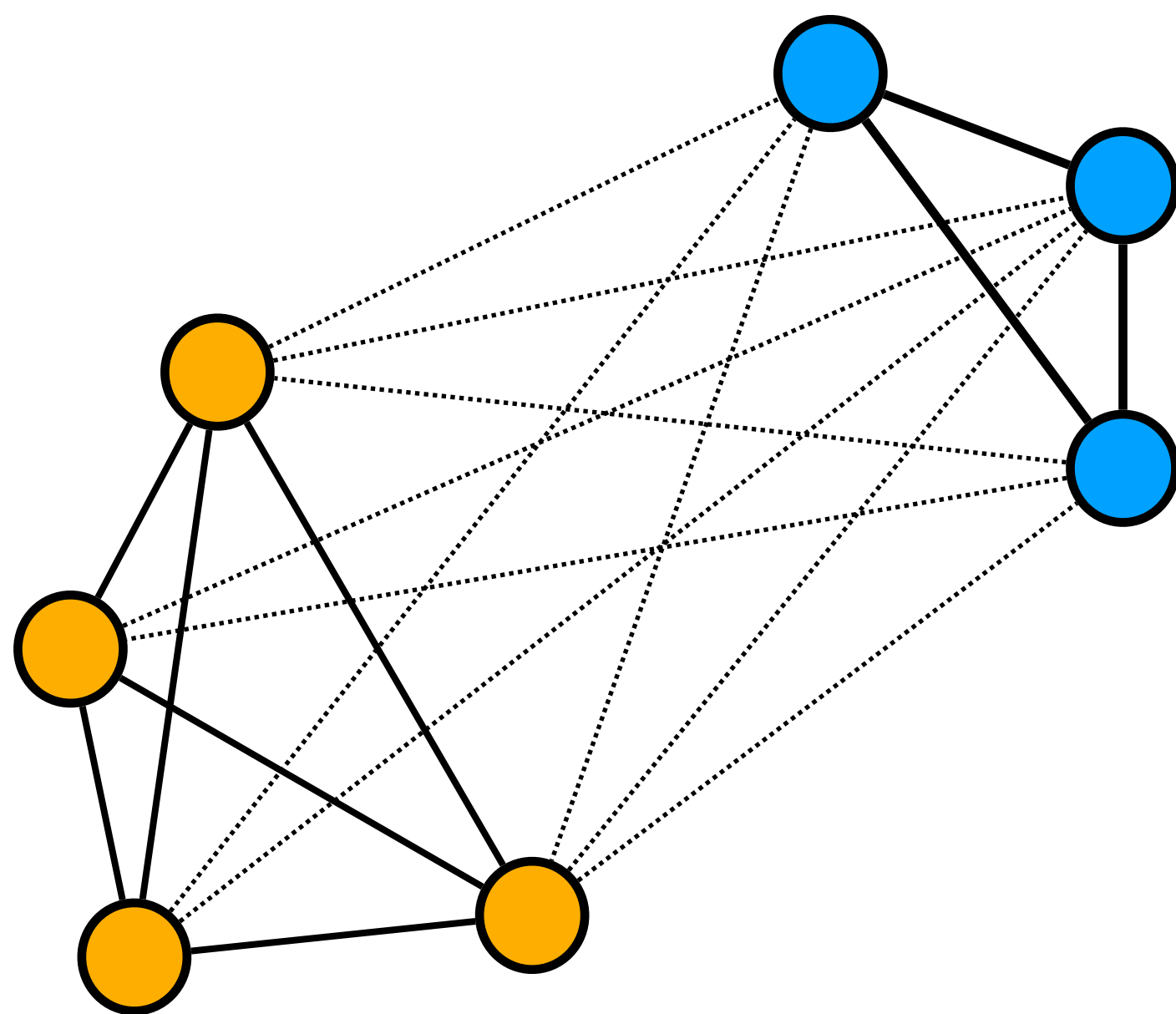


Edge weights across groups



New formulation for spectral clustering

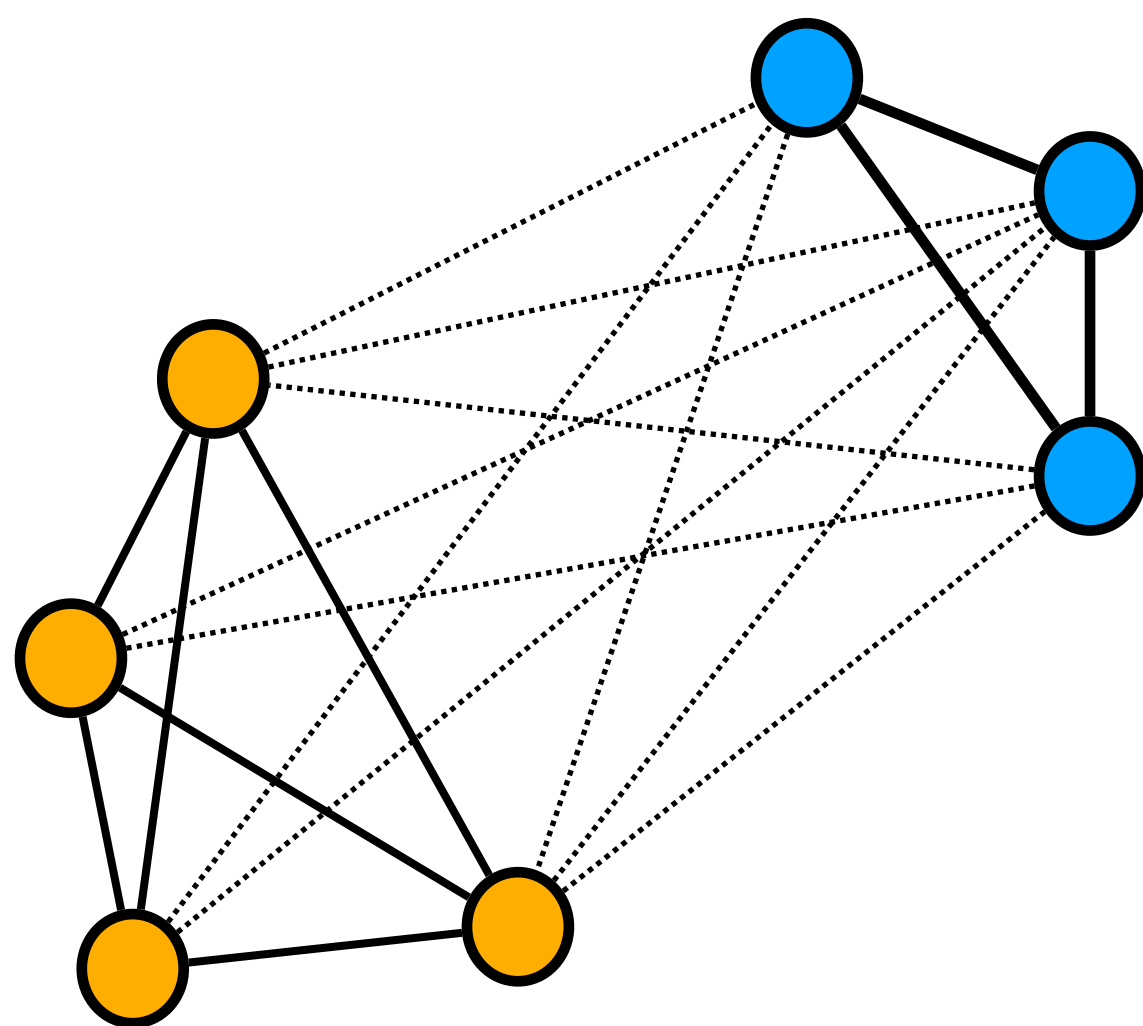
The basic clustering problem: a graph view



$$\textit{maximize} \quad \frac{\text{Edge weights within a group}}{\text{Edge weights across groups}}$$

New formulation for spectral clustering

The basic clustering problem: a graph view



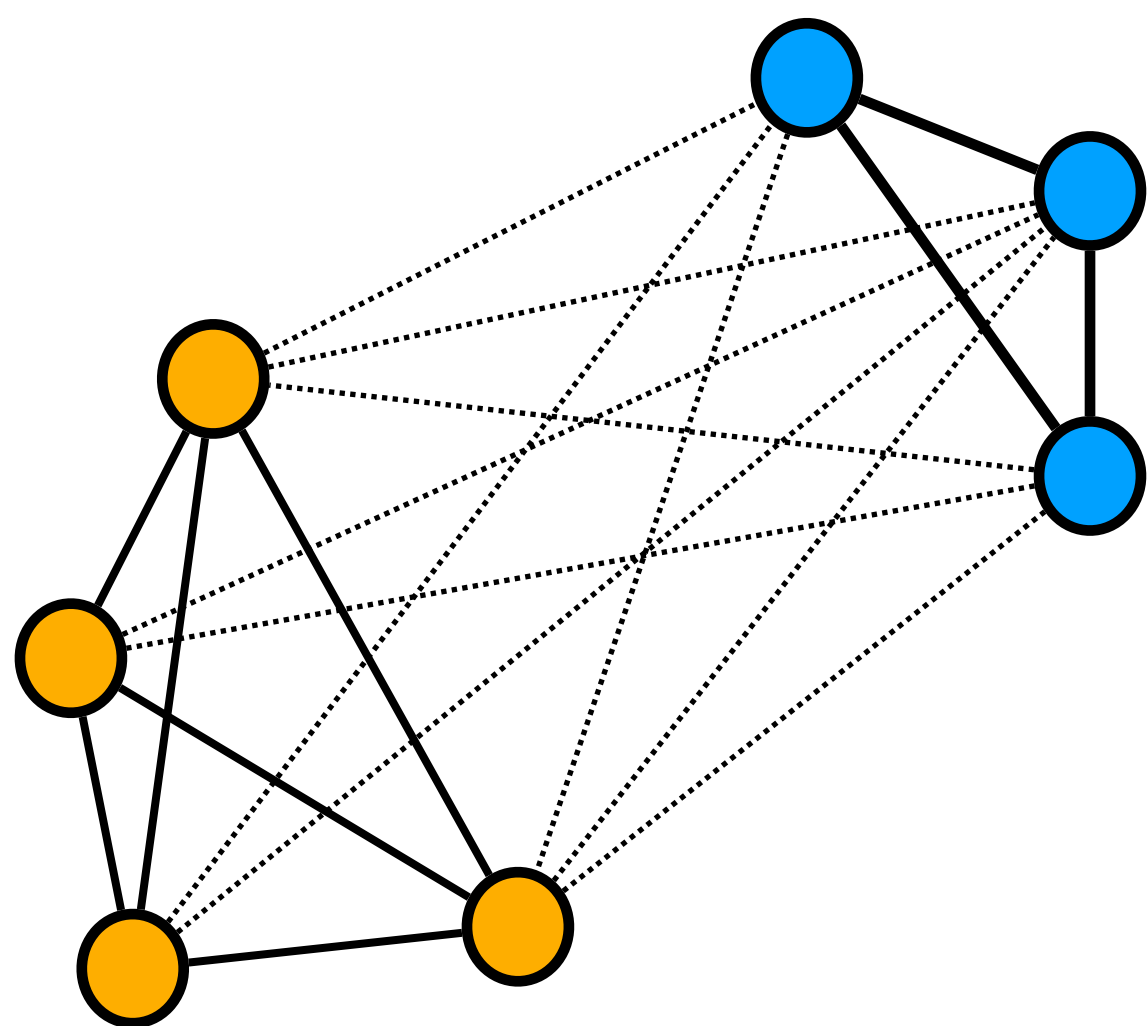
maximize $\frac{\text{Edge weights within a group}}{\text{Edge weights across groups}}$

$$\text{maximize} \quad \epsilon(X) = \frac{1}{K} \sum_{k=1}^K \frac{X_k^T \mathbf{A} X_k}{X_k^T \mathbf{D} X_k}$$

$$\text{subject to} \quad X \in \{0,1\}^{N \times K}, \\ X \mathbf{1}_K = \mathbf{1}_N.$$

New formulation for spectral clustering

The basic clustering problem: a graph view



maximize

$$\epsilon(X) = \frac{1}{K} \sum_{k=1}^K \frac{X_k^T \mathbf{A} X_k}{X_k^T \mathbf{D} X_k}$$

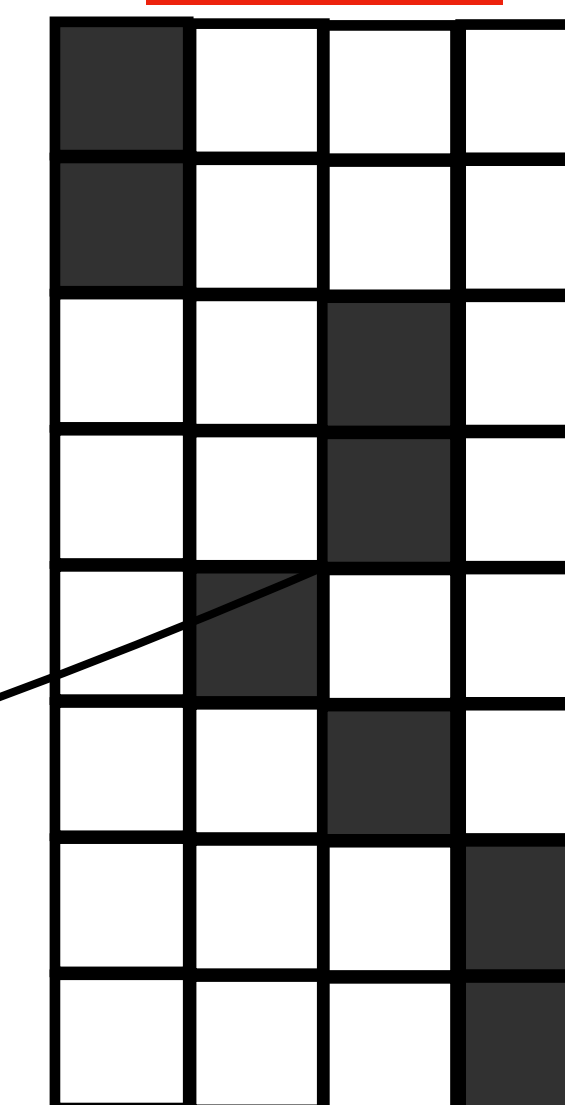
subject to

$$X \in \{0,1\}^{N \times K},$$

$$X \mathbf{1}_K = \mathbf{1}_N.$$

Affinity
matrix

#speakers



#segments

New formulation for spectral clustering

This problem is NP-hard!

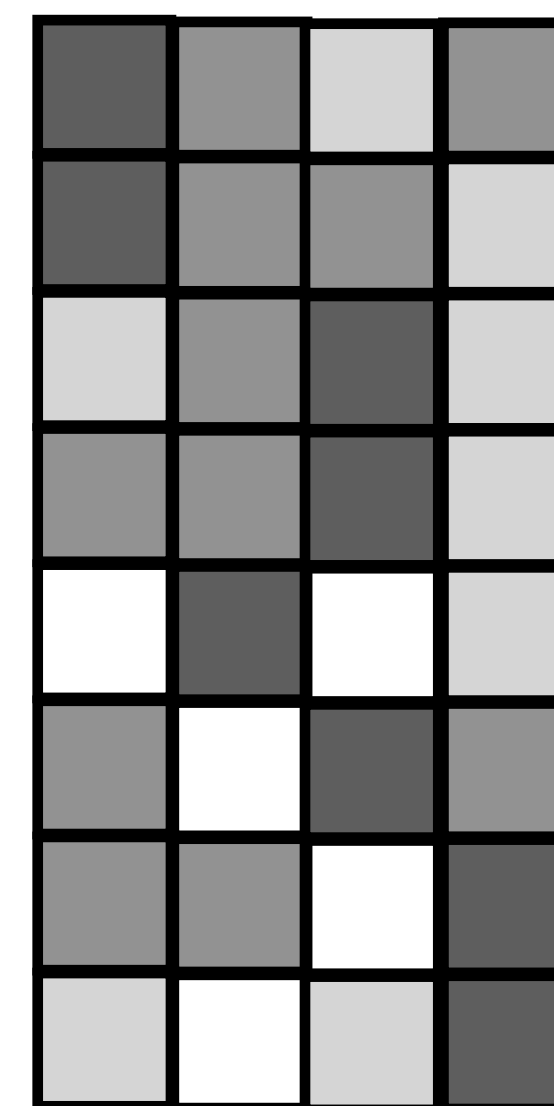
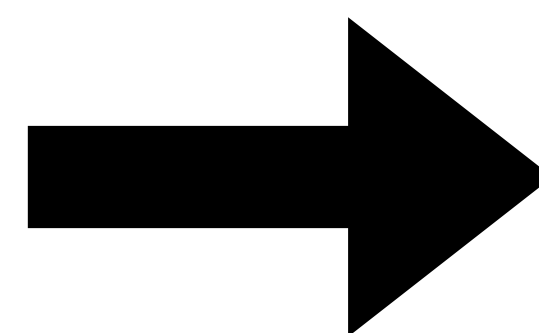
$$\begin{aligned} &\text{maximize} && \epsilon(X) = \frac{1}{K} \sum_{k=1}^K \frac{X_k^T \mathbf{A} X_k}{X_k^T \mathbf{D} X_k} \\ &\text{subject to} && X = \{0, 1\}^{N \times K}, \\ &&& X \mathbf{1}_K = \mathbf{1}_N. \end{aligned}$$

Remove the discrete constraints to make the problem solvable

New formulation for spectral clustering

Relaxed problem has a set of solutions

$$\begin{aligned} &\text{maximize} && \epsilon(X) = \frac{1}{K} \sum_{k=1}^K \frac{X_k^T \mathbf{A} X_k}{X_k^T \mathbf{D} X_k} \\ &\text{subject to} && X = \{x_1, \dots, x_K\}^{N \times K}, \\ &&& X \mathbf{1}_K = \mathbf{1}_N. \end{aligned}$$

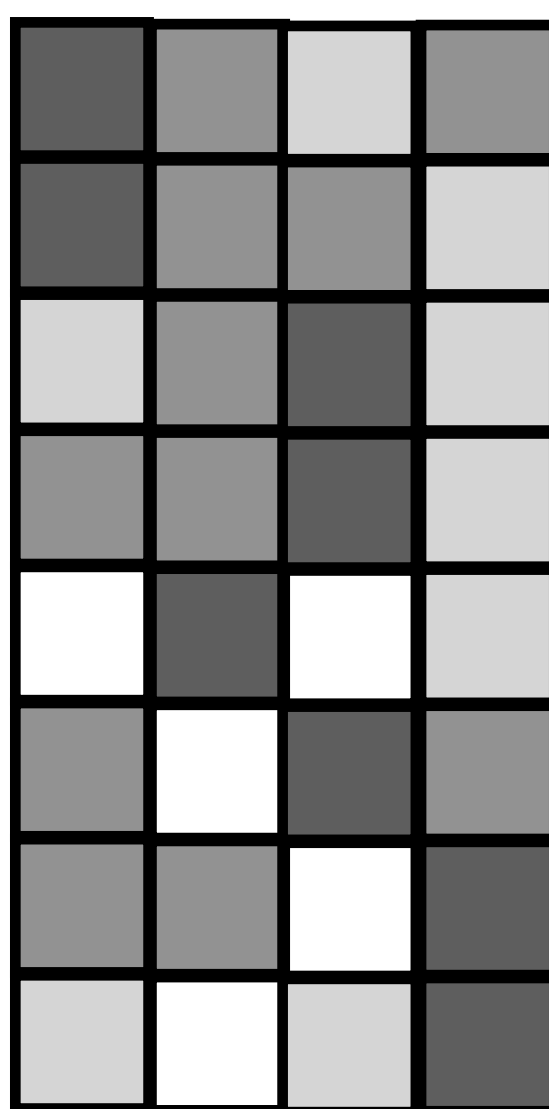


and its orthonormal transforms

Solution to the **relaxed** problem

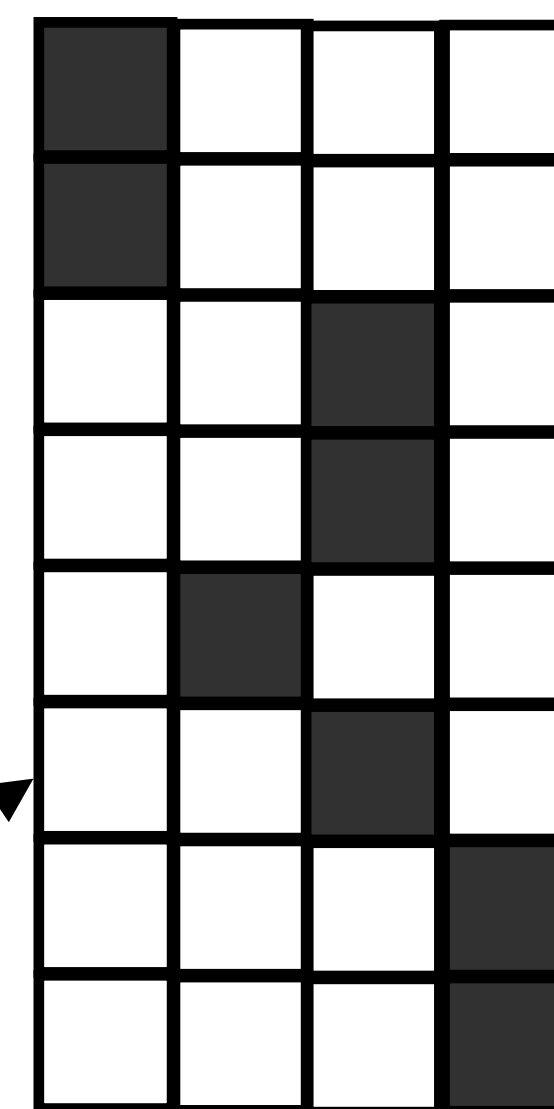
New formulation for spectral clustering

Now we need to **discretize** this solution!



and its orthonormal
transforms

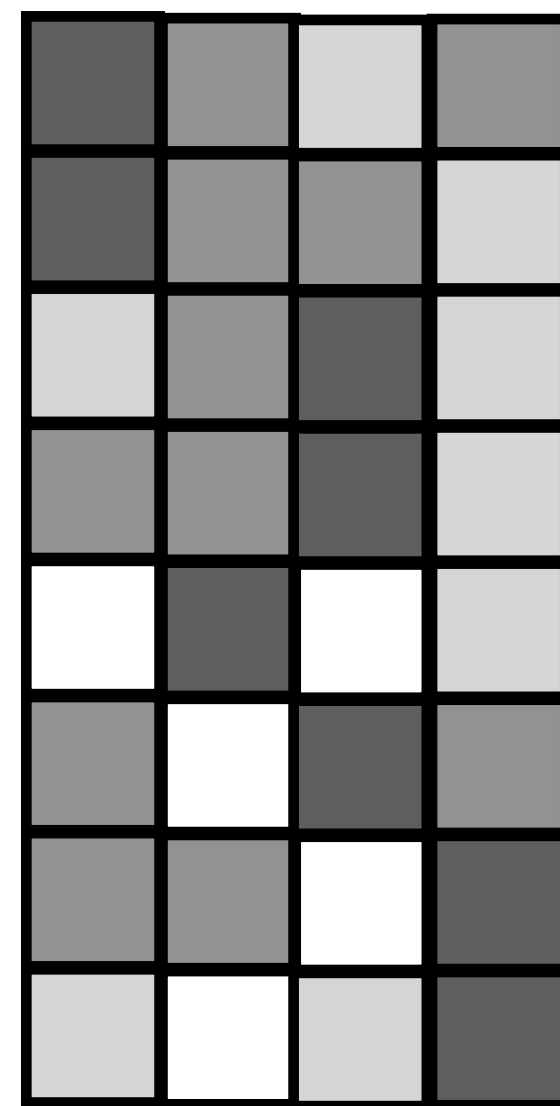
subject to $X \in \{0,1\}^{N \times K},$
 $X \mathbf{1}_K = \mathbf{1}_N.$



Find a matrix which is **discrete** and also close
to any one of the **orthonormal**
transformations of the relaxed solution

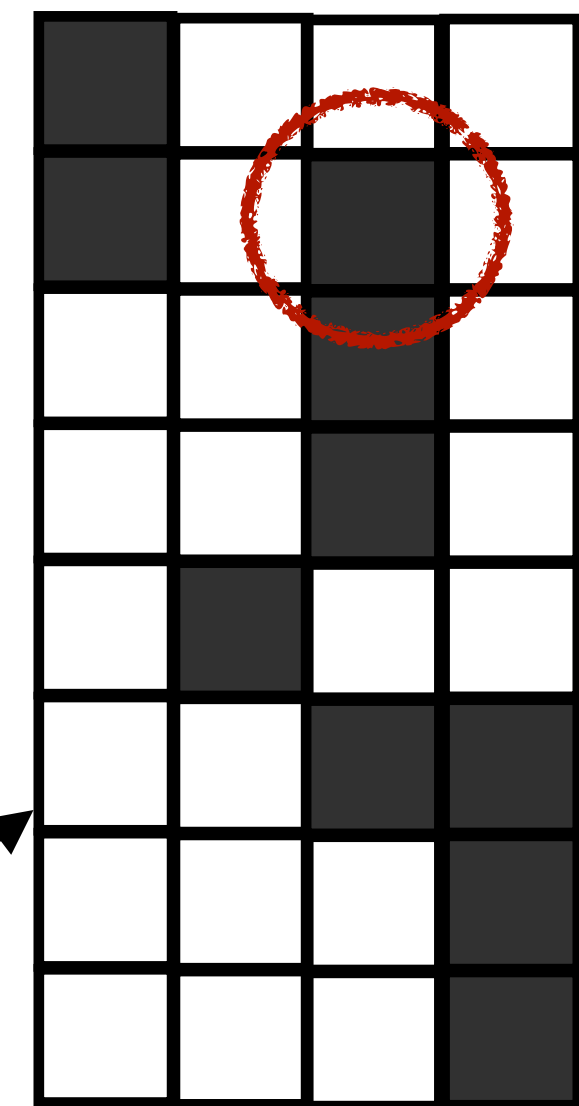
Here we use the overlap detector output

Suppose we have v_{OL}



and its orthonormal
transforms

subject to $X \in \{0,1\}^{N \times K},$
 $X \mathbf{1}_K = \mathbf{1}_N + v_{OL}.$



Non-maximal suppression: for overlapping
segments, select top 2 speakers

How does it compare with ...





... VB-based overlap assignment

Bullock, et al., “Overlap-aware diarization: resegmentation using neural end-to-end overlapped speech detection,” ICASSP 2020.

	Overlap-aware SC	VB based overlap assignment	
	Uses external overlap detector	Uses external overlap detector	
✓	No initialization required	Needs initialization, e.g. from AHC system	✗
✗	Does segment-level assignment (coarse)	Does frame-level assignment (fine-grained)	✓

How does it compare with ...

... EEND, RPN, TS-VAD

	Overlap-aware SC	RPN, EEND, TS-VAD	
	Uses external overlap detector	Includes overlap detection/assignment	
	Matched training data not required	Requires matched (simulated) training data	

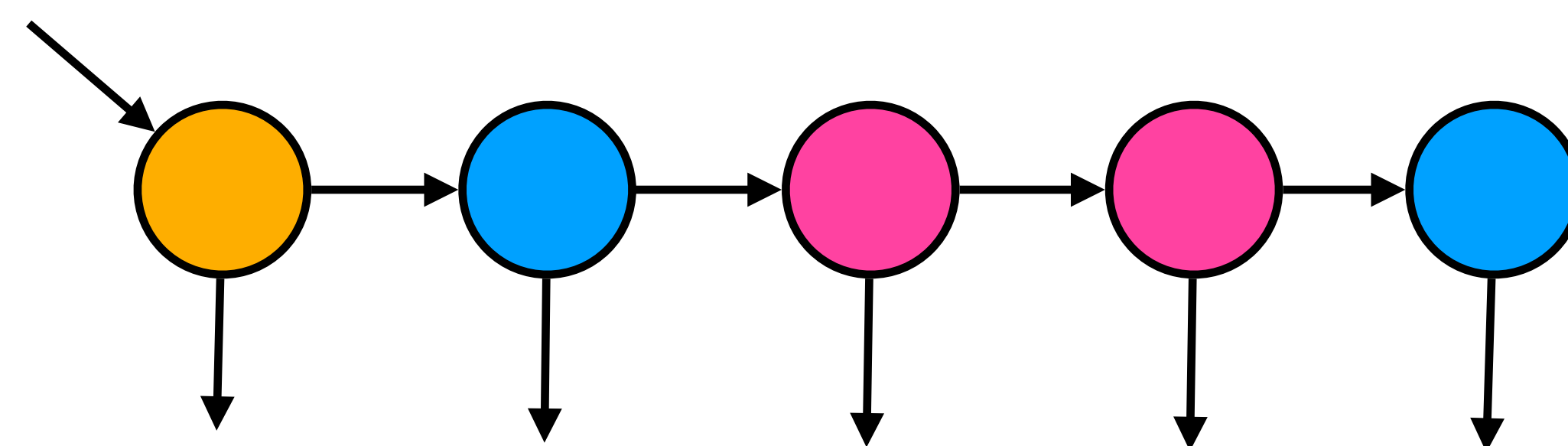
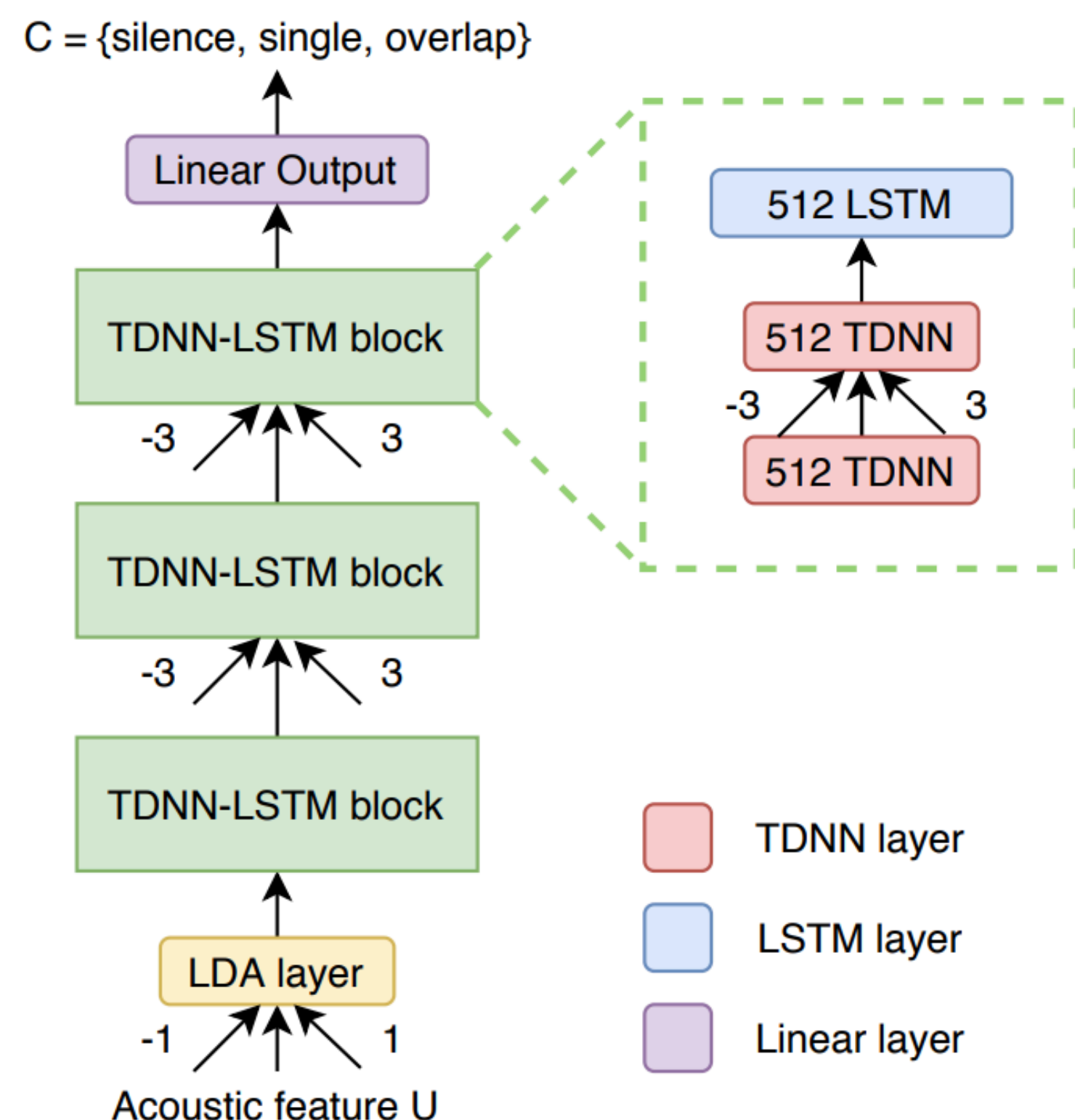
Huang et al., “Speaker diarization with region proposal network,” ICASSP 2020.

Fujita et al., “End-to-end neural diarization: Reformulating speaker diarization as simple multi-label classification,” ArXiv.

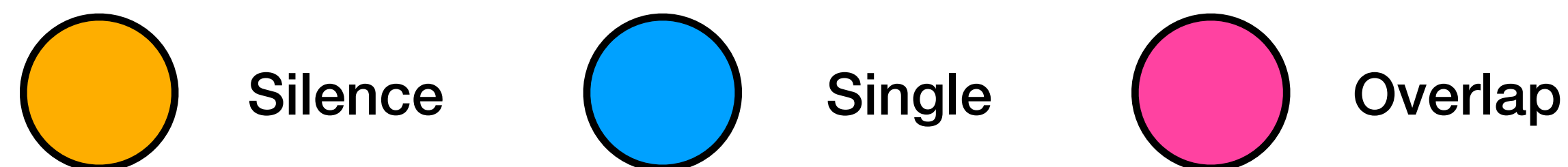
Medennikov, et al., “Target speaker voice activity detection: a novel approach for multispeaker diarization in a dinner party scenario,” Interspeech 2020.

Hybrid HMM-DNN overlap detector

(Now merged as a 🎧 *KALDI* recipe)



Viterbi decoding used for inference



Results: DER breakdown on AMI eval

System	Missed speech	False alarm	Speaker conf.	DER
AHC/PLDA	19.9	0.0	8.4	26.9
Spectral/cosine	19.9	0.0	7.0	28.3
VBx	19.9	0.0	6.3	26.2
VB-based overlap assignment	13.0	3.6	7.2	23.8
RPN	9.5	7.7	8.3	25.5
Overlap-aware SC	11.3	2.2	10.5	24.0

Results: DER breakdown on AMI eval

Missed speech decreases significantly



System	Missed speech	False alarm	Speaker conf.	DER
AHC/PLDA	19.9	0.0	8.4	26.9
Spectral/cosine	19.9	0.0	7.0	28.3
VBx	19.9	0.0	6.3	26.2
VB-based overlap assignment	13.0	3.6	7.2	23.8
RPN	9.5	7.7	8.3	25.5
Overlap-aware SC	11.3	2.2	10.5	24.0

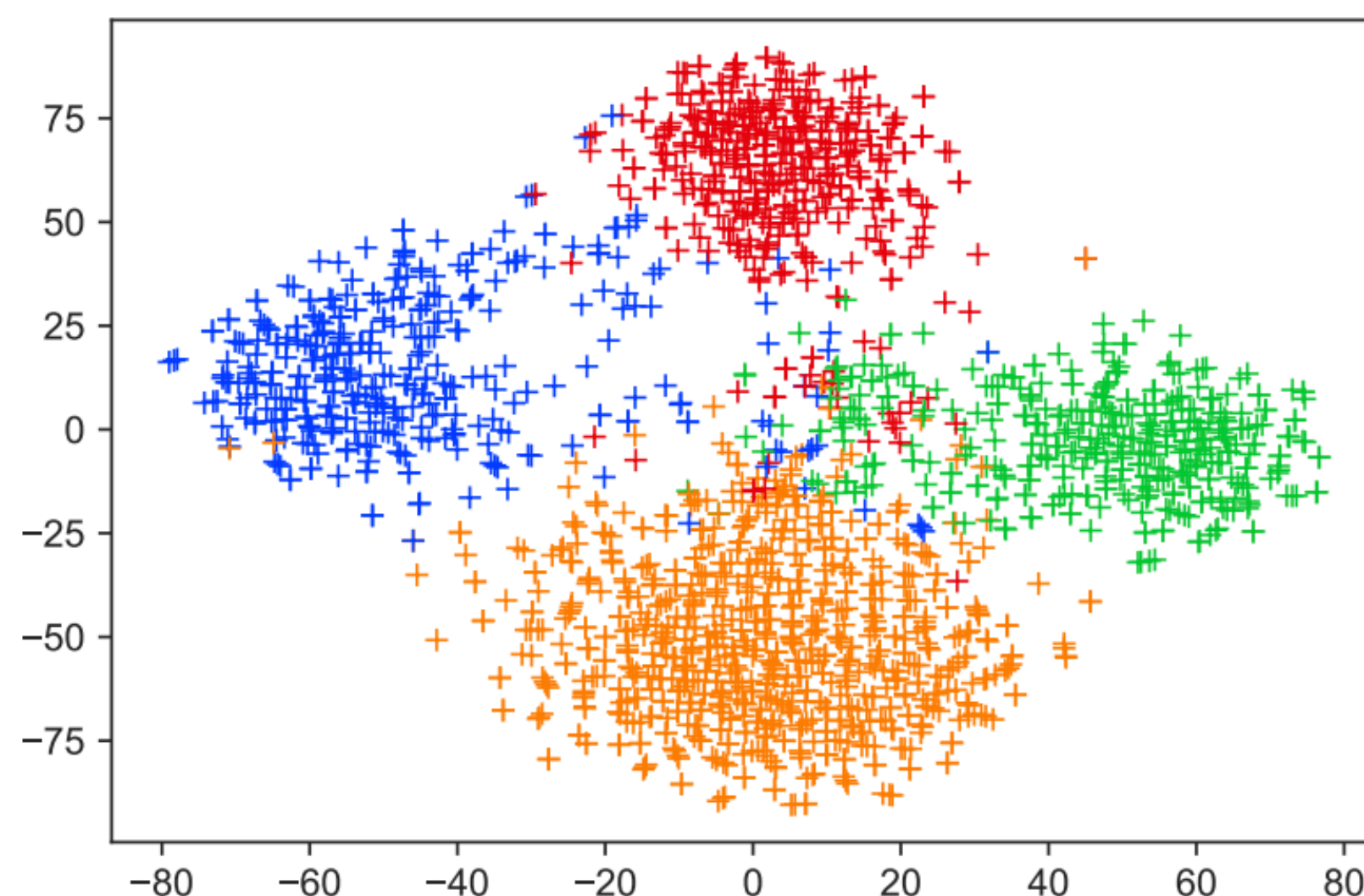
Results: DER breakdown on AMI eval

Speaker confusion increases

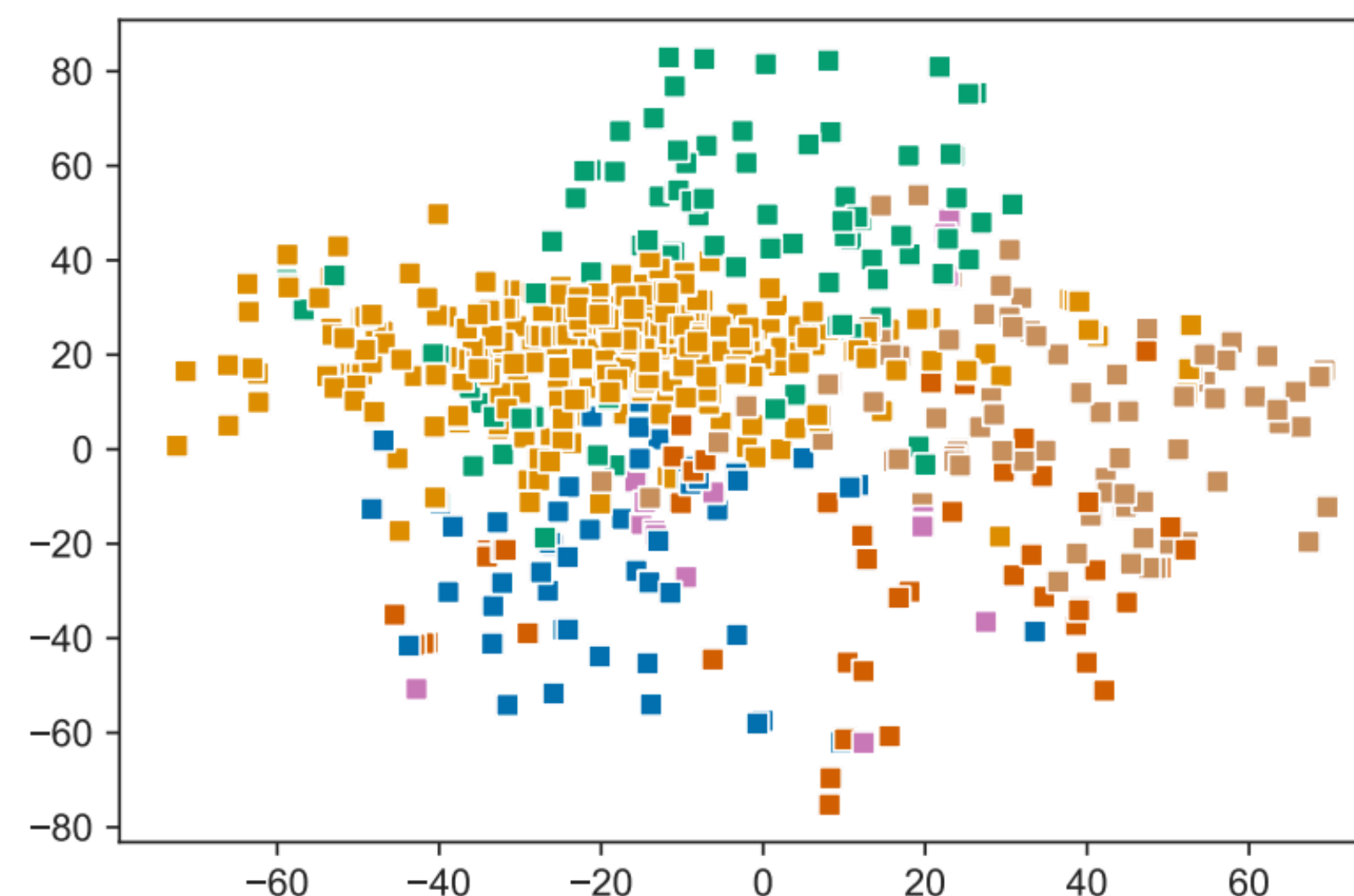


System	Missed speech	False alarm	Speaker conf.	DER
AHC/PLDA	19.9	0.0	8.4	26.9
Spectral/cosine	19.9	0.0	7.0	28.3
VBx	19.9	0.0	6.3	26.2
VB-based overlap assignment	13.0	3.6	7.2	23.8
RPN	9.5	7.7	8.3	25.5
Overlap-aware SC	11.3	2.2	10.5	24.0

Future work: train a more robust x-vector extractor



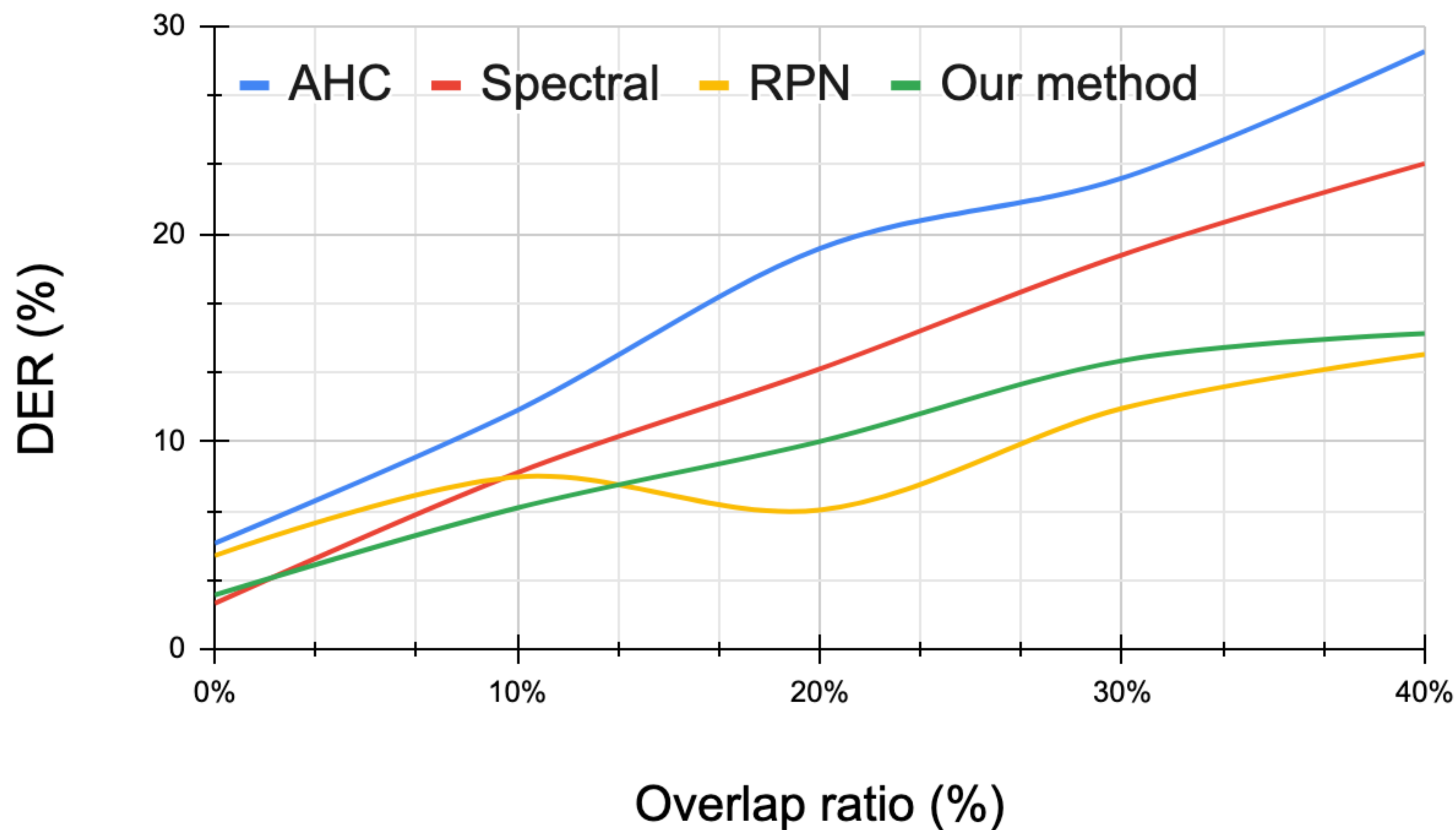
Non-overlapping segments



Overlapping segments

T-SNE plot of x-vector embeddings

More results: DER on LibriCSS



Try out the code!



Acknowledgments:

Paola Garcia, for helpful discussions and insights.

Takuya Yoshioka, for simulation script for generating LibriCSS training data.