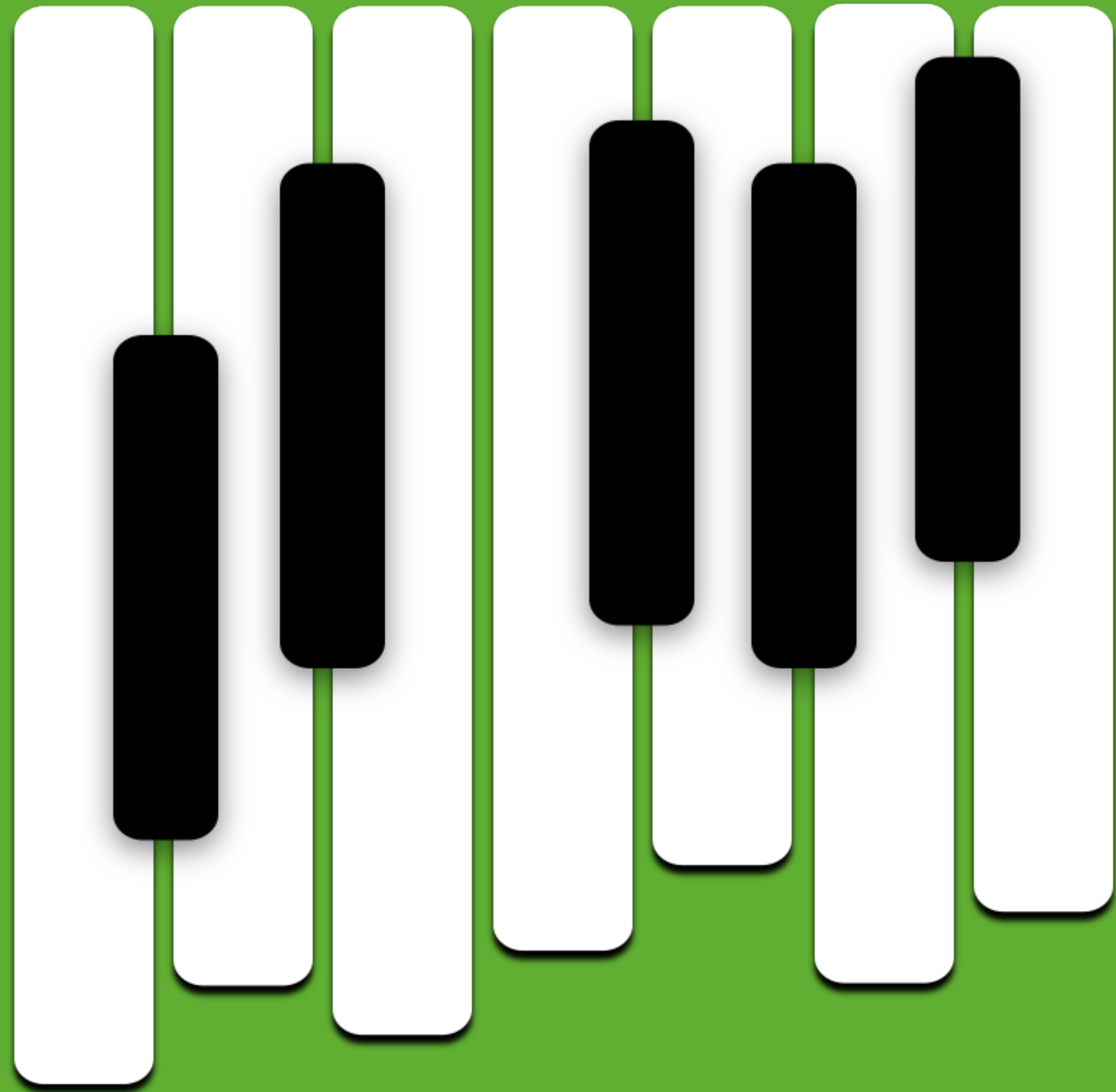


pyannotate



at VoxSRC 2023

Séverin Baroudi
Hervé Bredin
Alexis Plaquet
Thomas Pellegrini



UNIVERSITÉ
TOULOUSE III
PAUL SABATIER

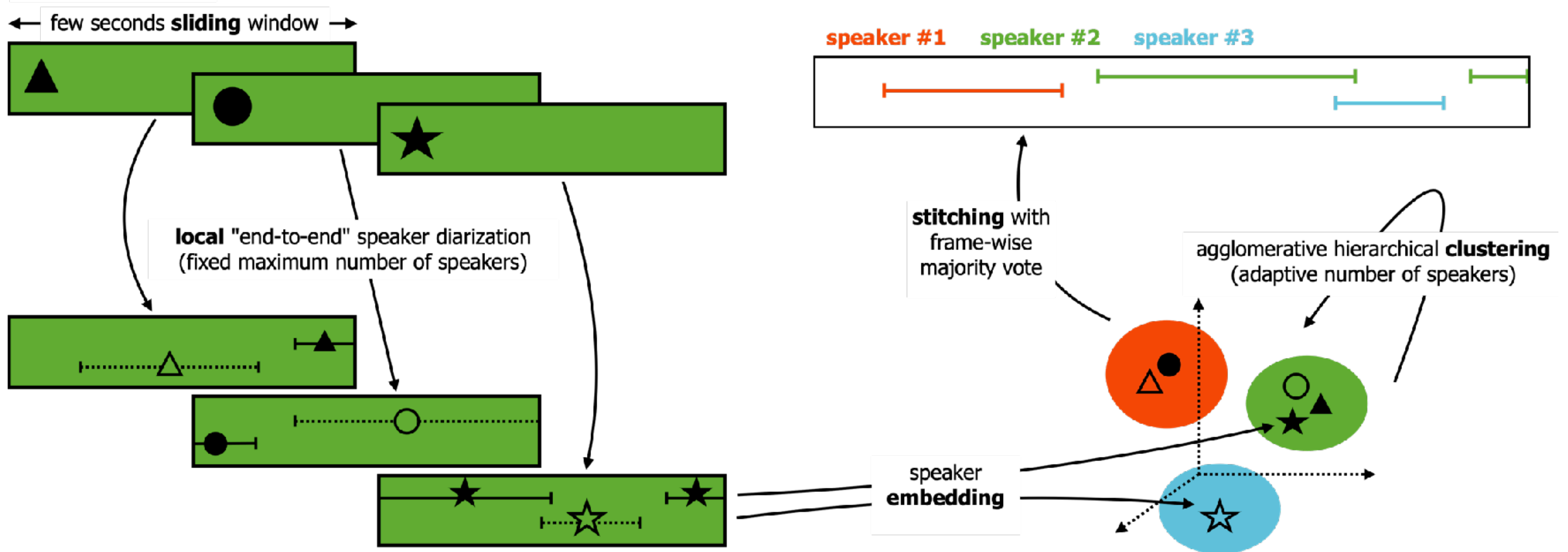


Université
de Toulouse



Main principle

the best of both worlds



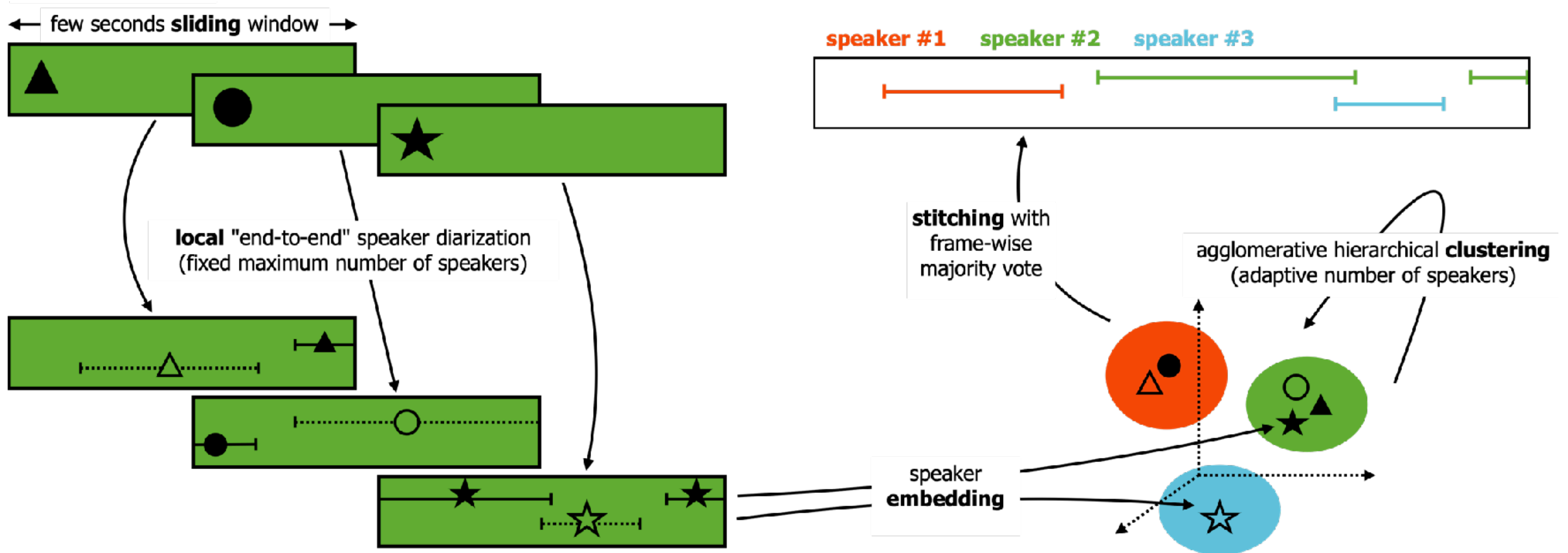


Main principle

the best of both worlds

pyannote.audio 2.1 speaker diarization pipeline
principle, benchmark, and recipe

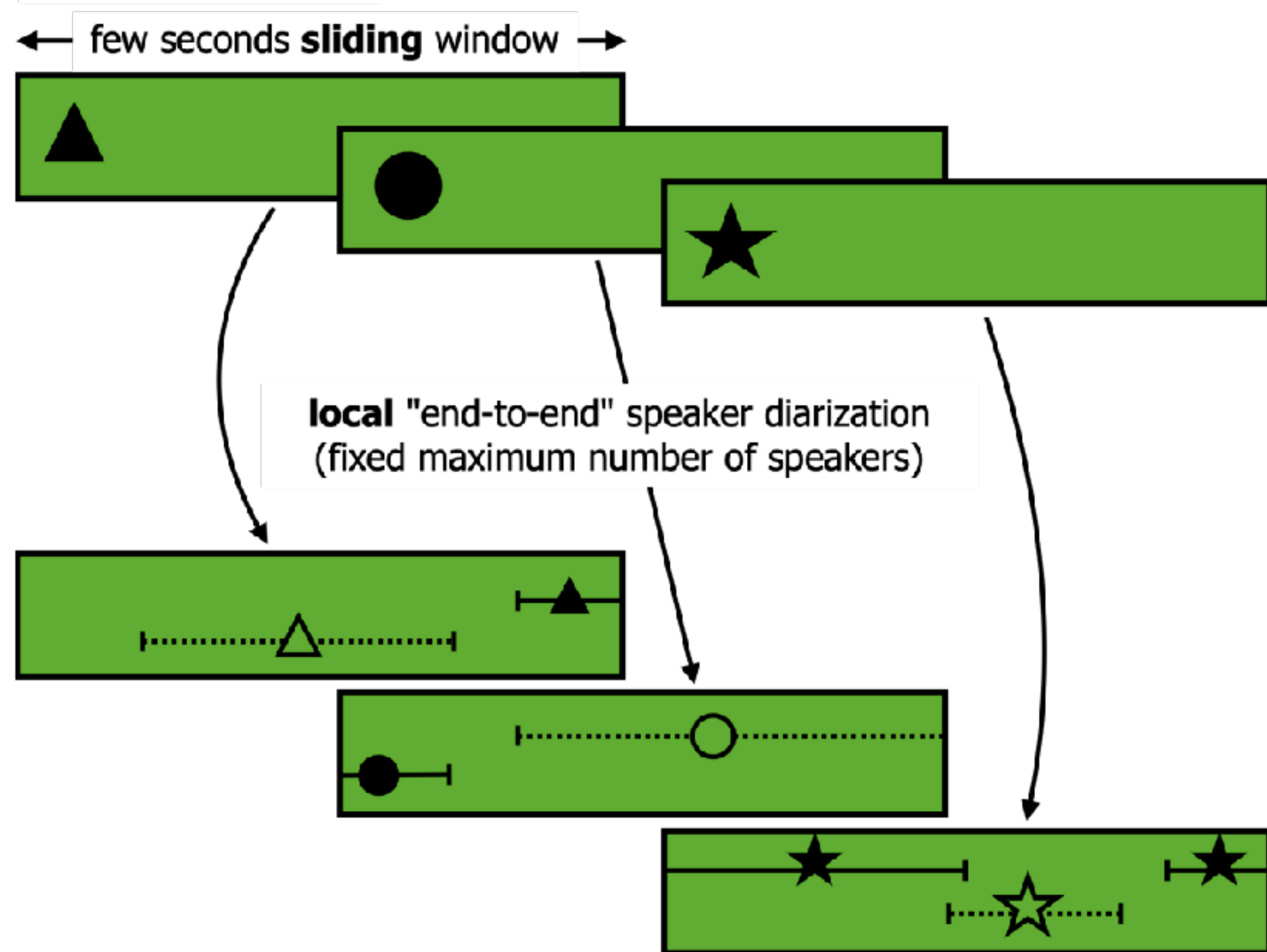
Hervé Bredin





Local "end-to-end" speaker diarization

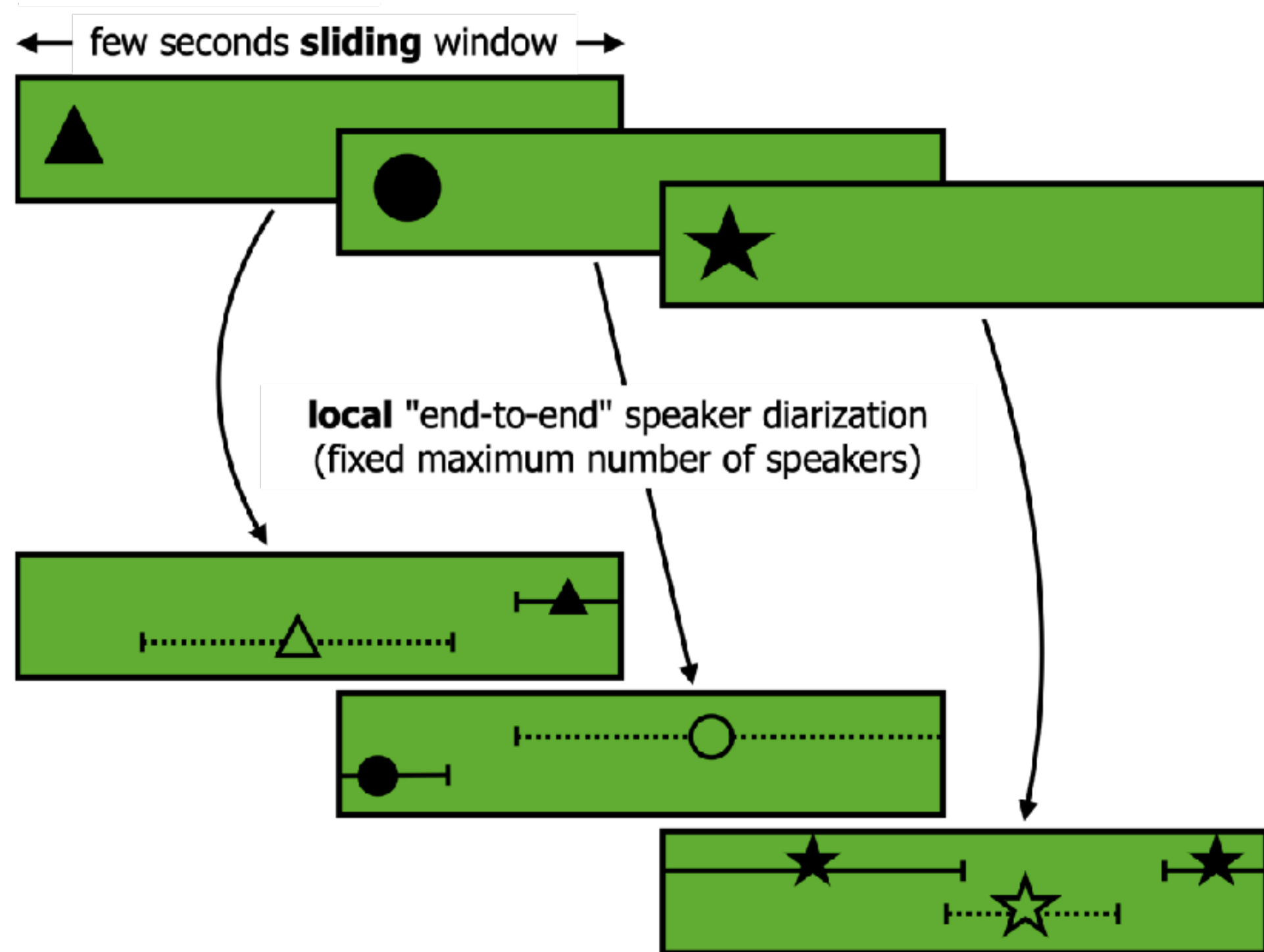
Improvements since last year





Local "end-to-end" speaker diarization

Improvements since last year

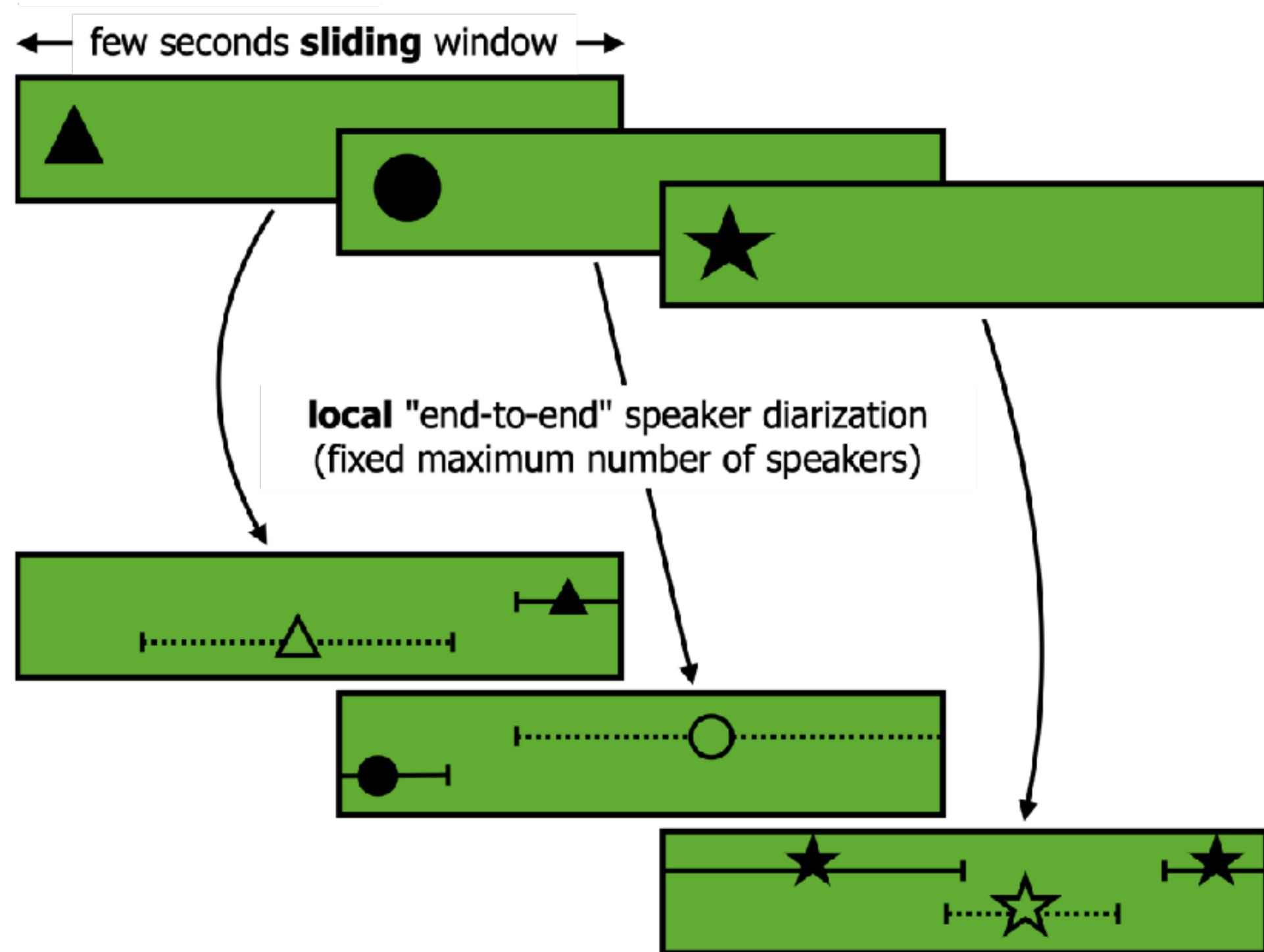


- Larger training set
AMI, DIHARD, VoxConverse
+ *AISHELL4+ AliMeeting*
+ *AVA-AVD + Ego4D*
+ *MSDWild + REPERE*



Local "end-to-end" speaker diarization

Improvements since last year

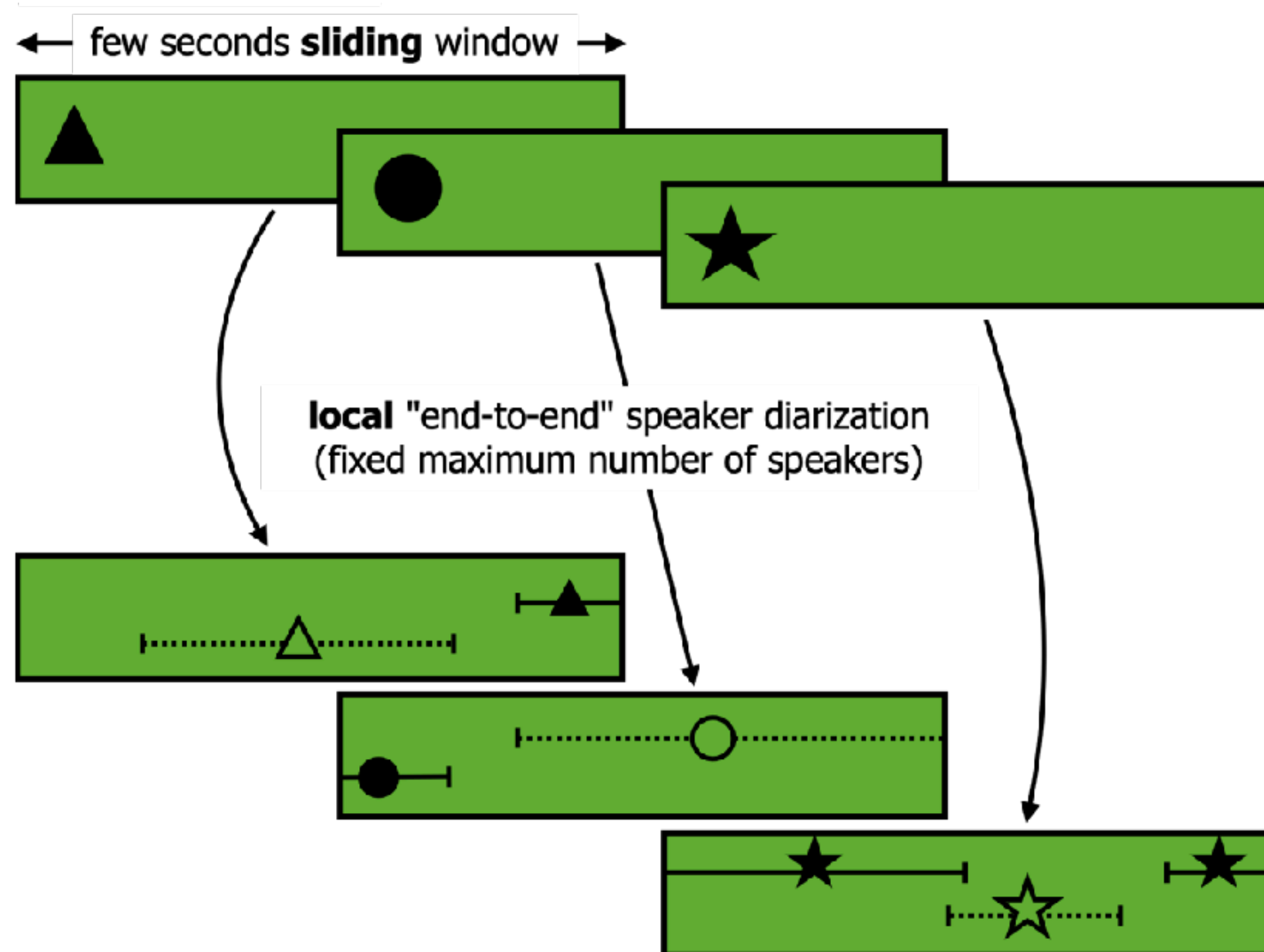


- Larger training set
AMI, DIHARD, VoxConverse
+ AISHELL4+ AliMeeting
+ AVA-AVD + Ego4D
+ MSDWild + REPERE
- Longer windows
from 5s to 10s



Local "end-to-end" speaker diarization

Improvements since last year

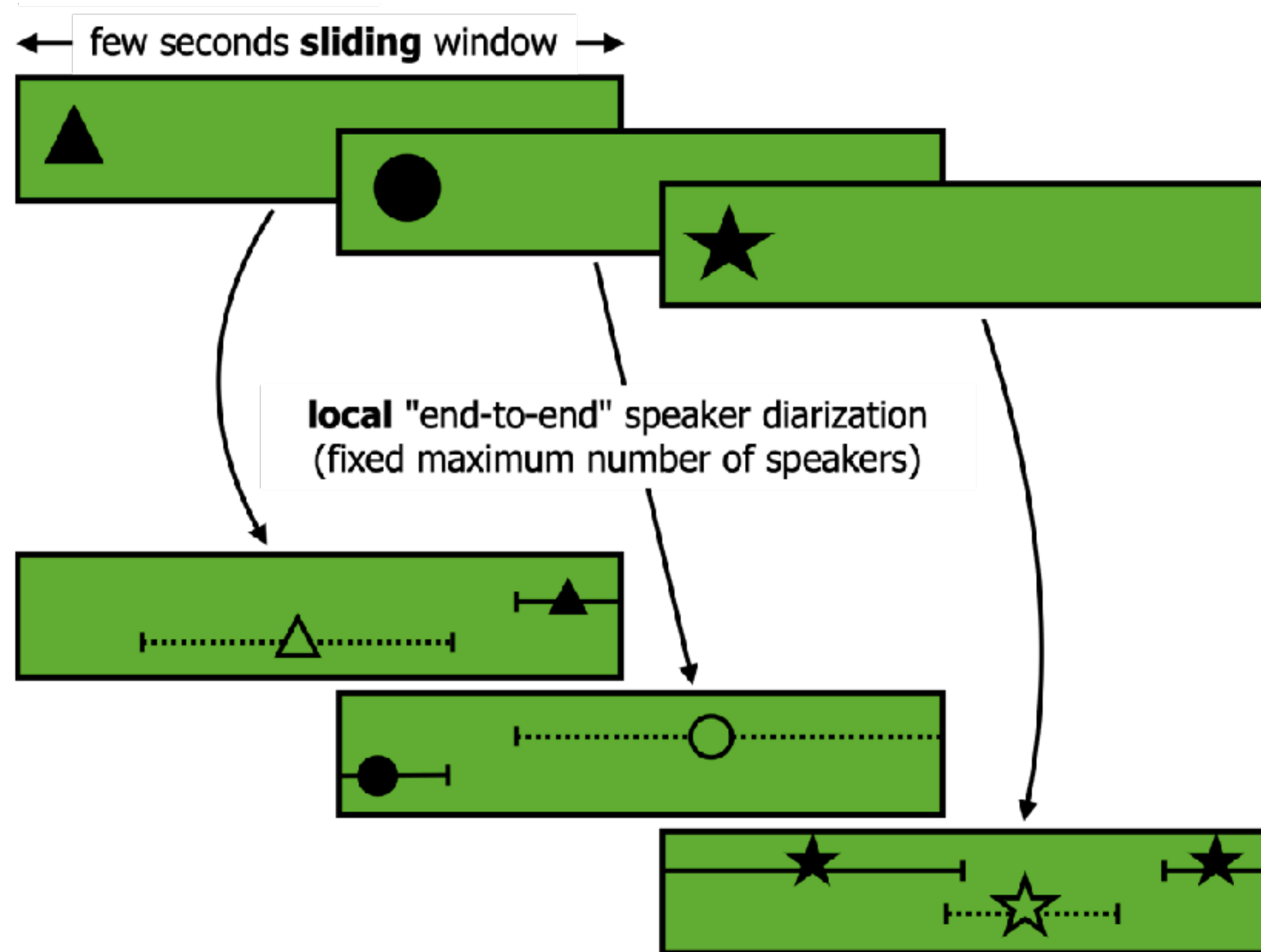


- Larger training set
AMI, DIHARD, VoxConverse
+ AISHELL4+ AliMeeting
+ AVA-AVD + Ego4D
+ MSDWild + REPERE
- Longer windows
from 5s to 10s
- **PowerSet** encoding



Local "end-to-end" speaker diarization

Improvements since last year



- Larger training set
AMI, DIHARD, VoxConverse
+ *AISHELL4+ AliMeeting*
+ *AVA-AVD + Ego4D*
+ *MSDWild + REPERE*
- Longer windows
from 5s to 10s
- **PowerSet** encoding
- **Self-supervised** feature extraction



Powerset encoding

from multi-label (+ threshold) to multi-class (+ argmax)

s_1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0
s_2	0	0	0	1	1	1	1	1	1	0	0	0	1	0	0	
s_3	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1
y_1												y_t	y_T		

|-----| **alice**
|-----| **bob**
|-----| **carol**
|-----| **bob**

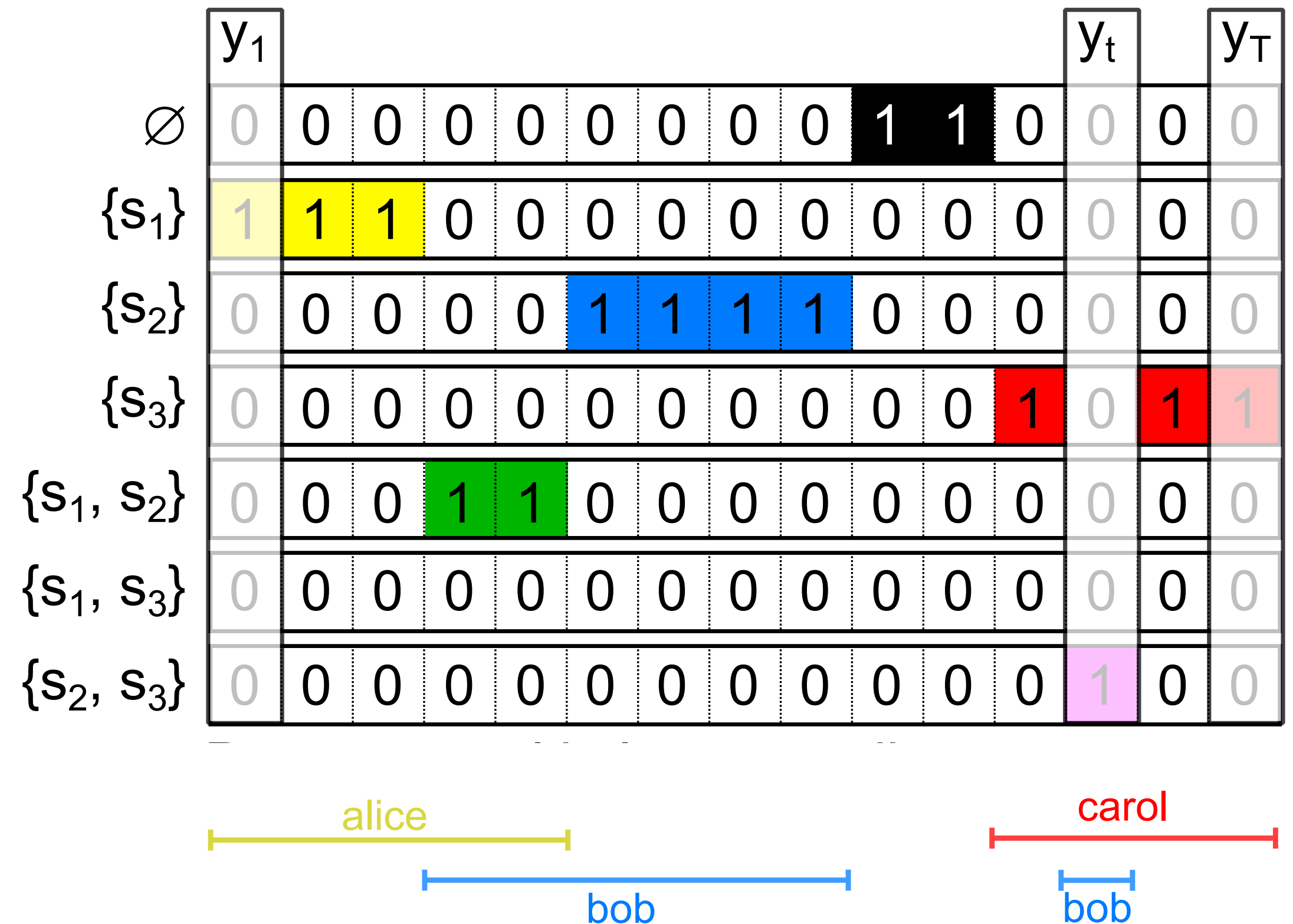
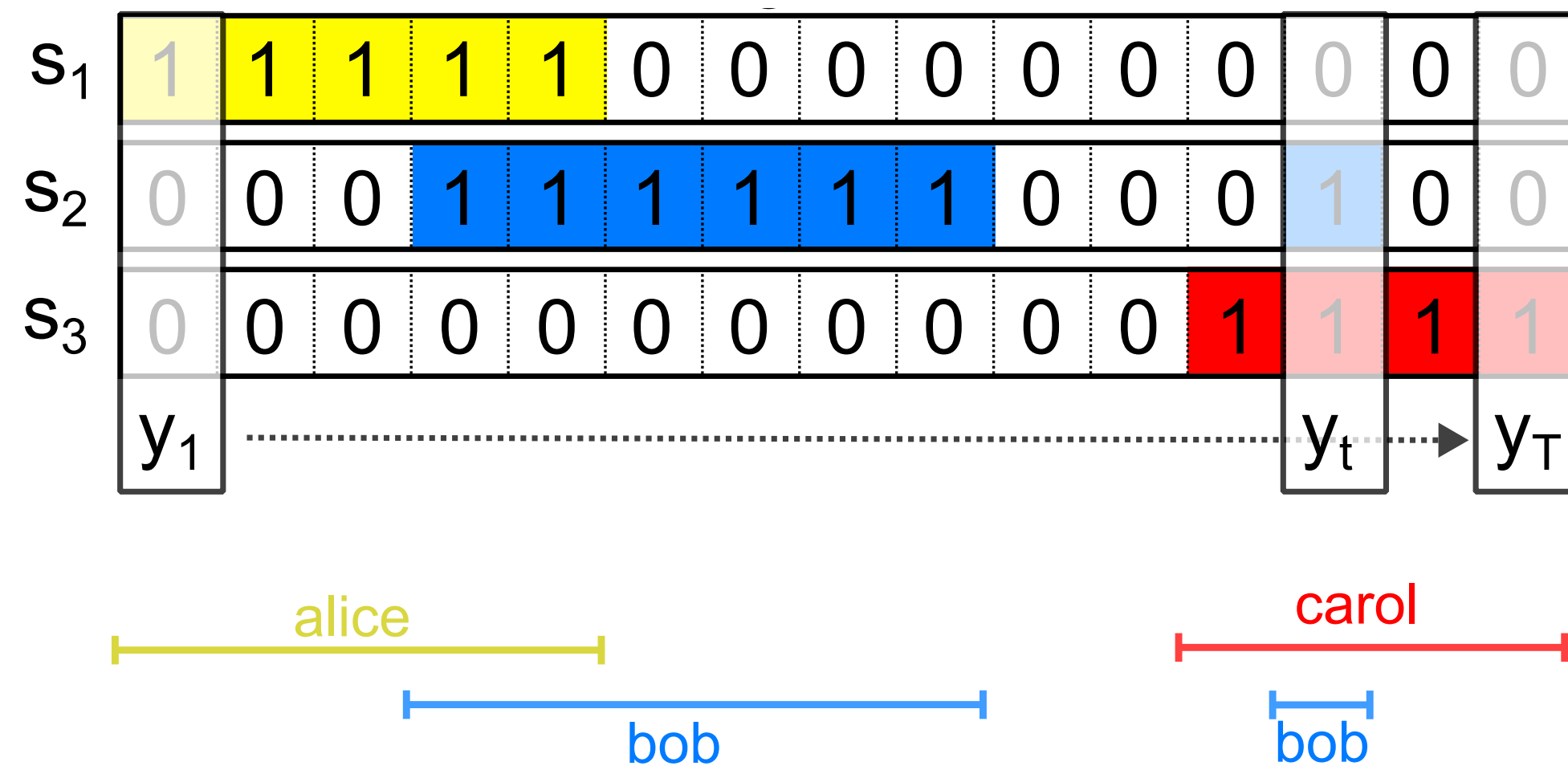
	y_1													y_t	y_T				
\emptyset	0	0	0	0	0	0	0	0	0	0	0	0	1	1	0	0	0	0	
$\{s_1\}$	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
$\{s_2\}$	0	0	0	0	0	1	1	1	1	0	0	0	0	0	0	0	0	0	
$\{s_3\}$	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1	0	1	1
$\{s_1, s_2\}$	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	
$\{s_1, s_3\}$	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
$\{s_2, s_3\}$	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	1	0	0

|-----| **alice**
|-----| **bob**
|-----| **carol**
|-----| **bob**



Powerset encoding

from multi-label (+ threshold) to multi-class (+ argmax)



Powerset multi-class cross entropy loss for neural speaker diarization

Alexis Plaquet & Hervé Bredin

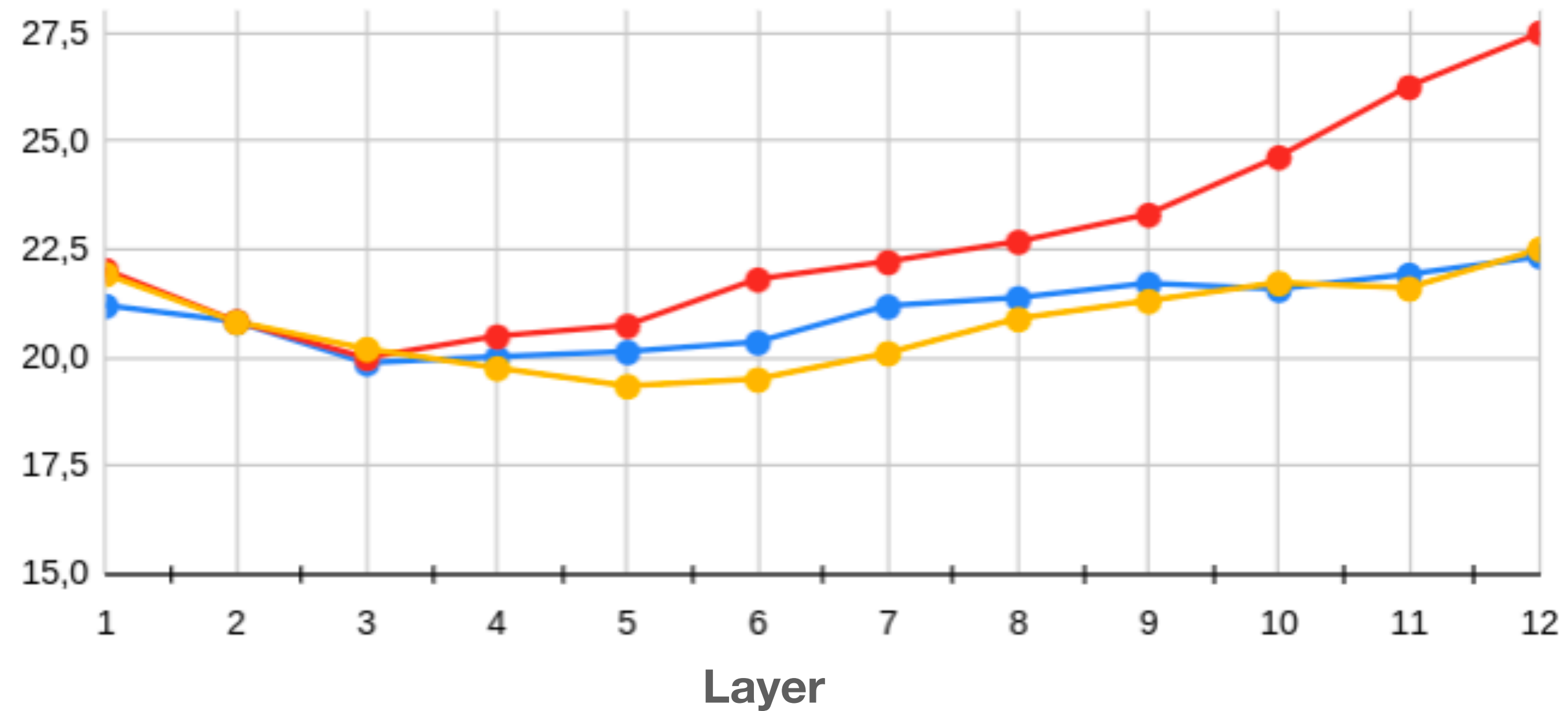




Self-supervised feature extraction

wav2vec 2.0 vs *HuBERT* vs *WavLM*

Diarization error rate on DIHARD



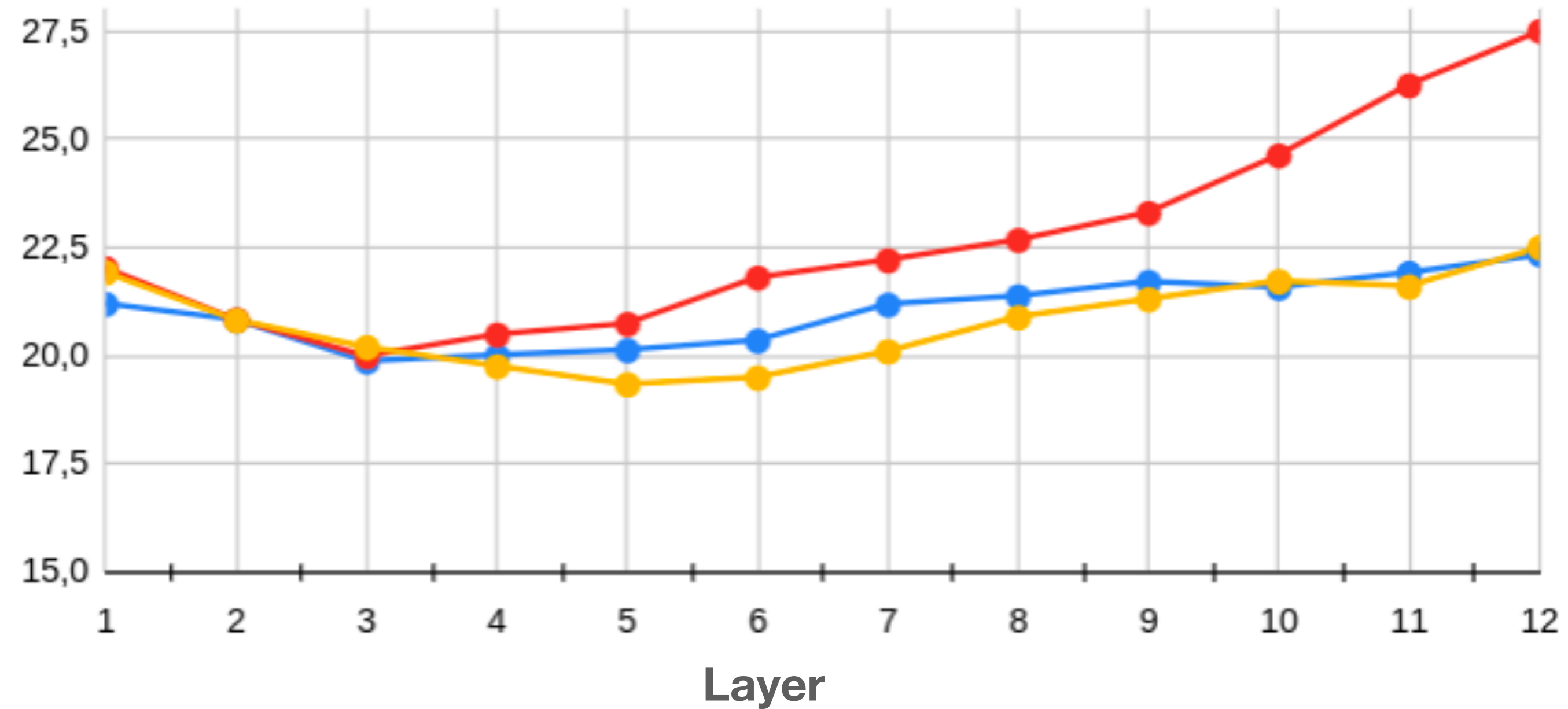
WavLM pretrained
on Librispeech



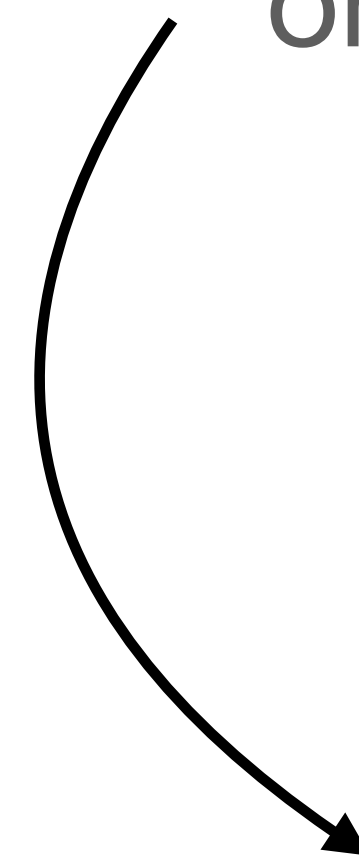
Self-supervised feature extraction

wav2vec 2.0 vs *HuBERT* vs *WavLM*

Diarization error rate on DIHARD



WavLM pretrained on Librispeech



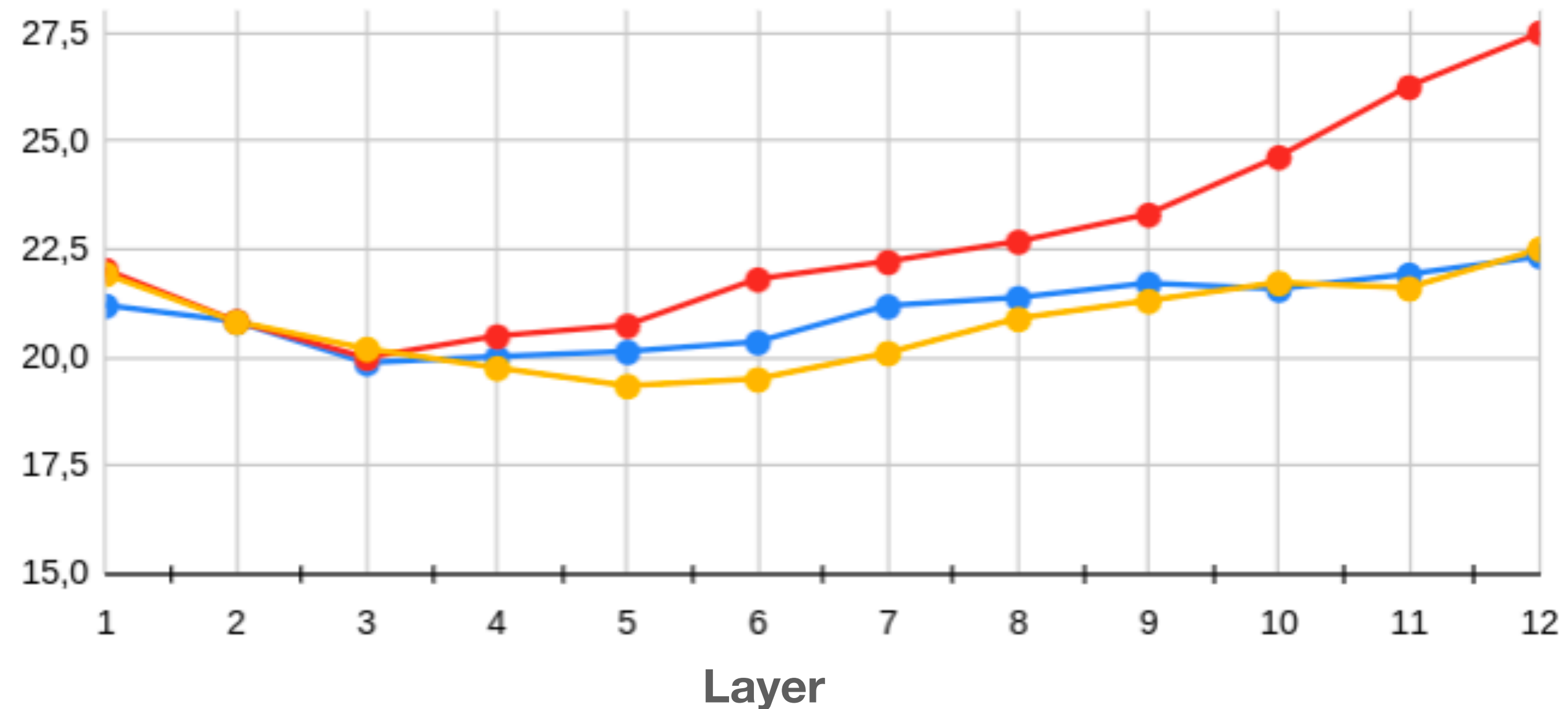
WavLM pretrained on diarization datasets



Self-supervised feature extraction

wav2vec 2.0 vs *HuBERT* vs *WavLM*

Diarization error rate on DIHARD



WavLM pretrained on Librispeech

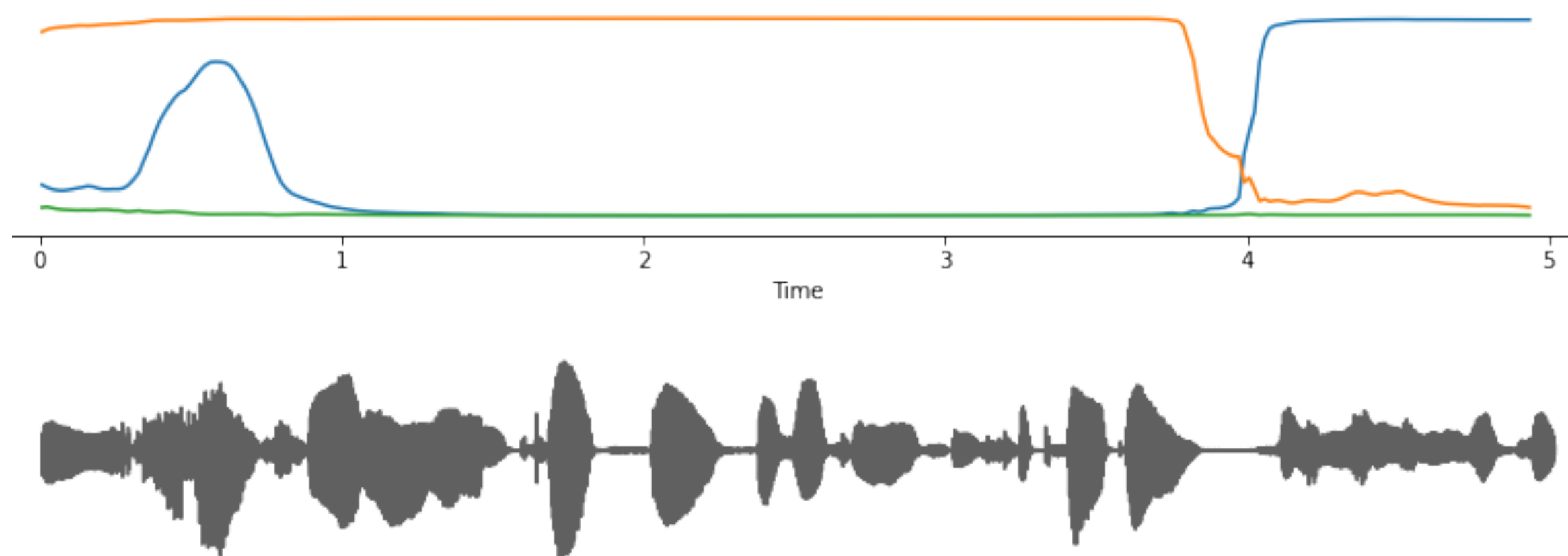
21% relative decrease in diarization error rate

WavLM pretrained on diarization datasets



Speaker embedding

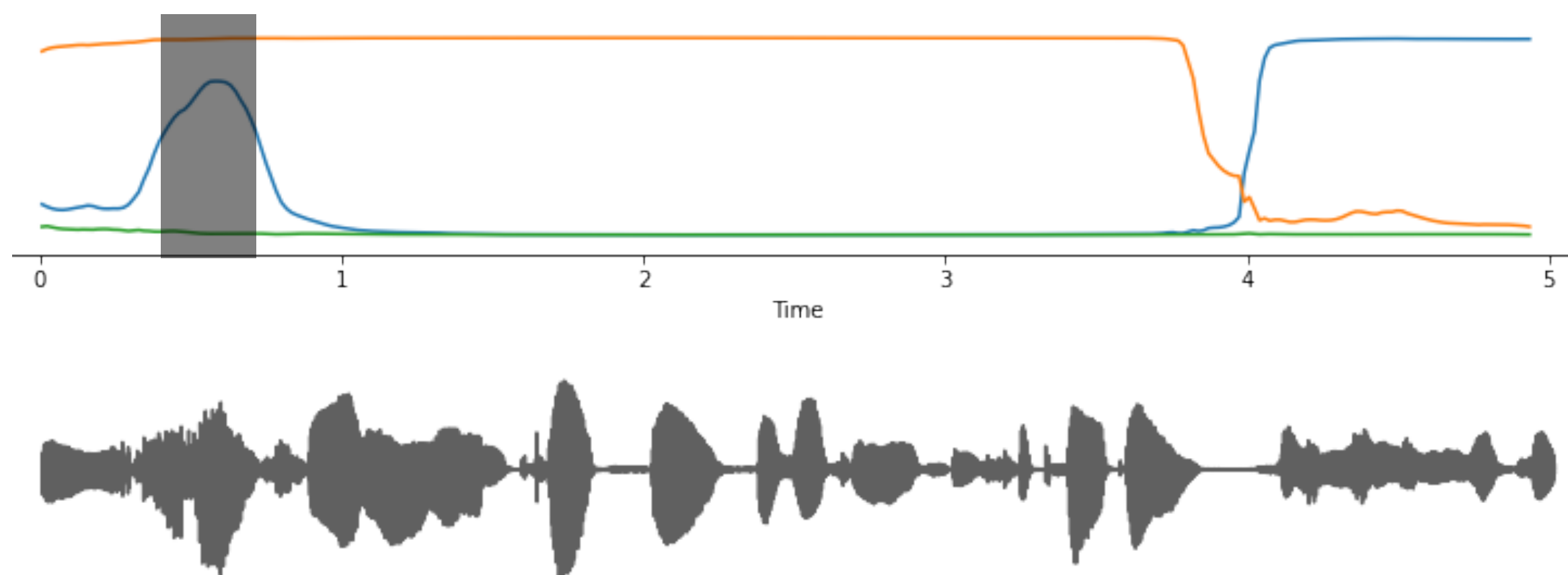
from SpeechBrain ECAPA-TDNN to WeSpeaker ResNet34





Speaker embedding

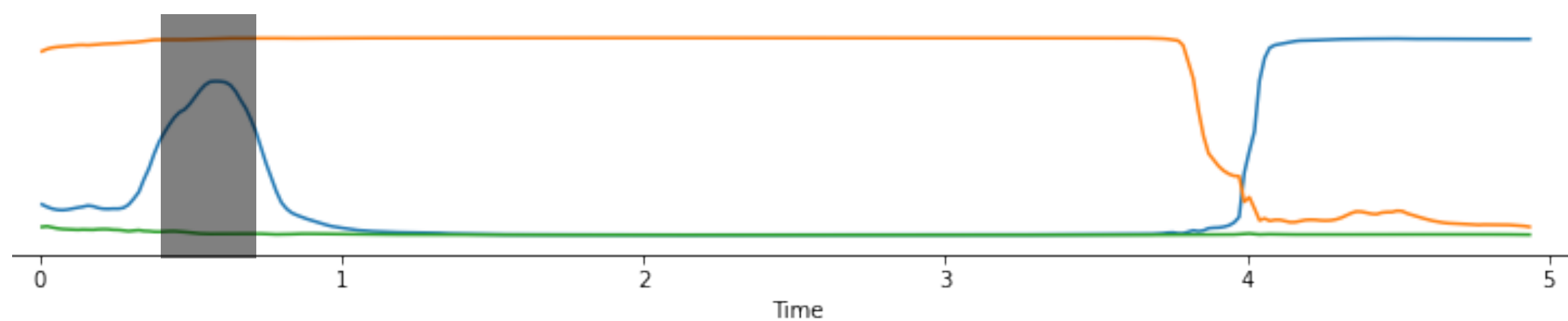
from SpeechBrain ECAPA-TDNN to WeSpeaker ResNet34





Speaker embedding

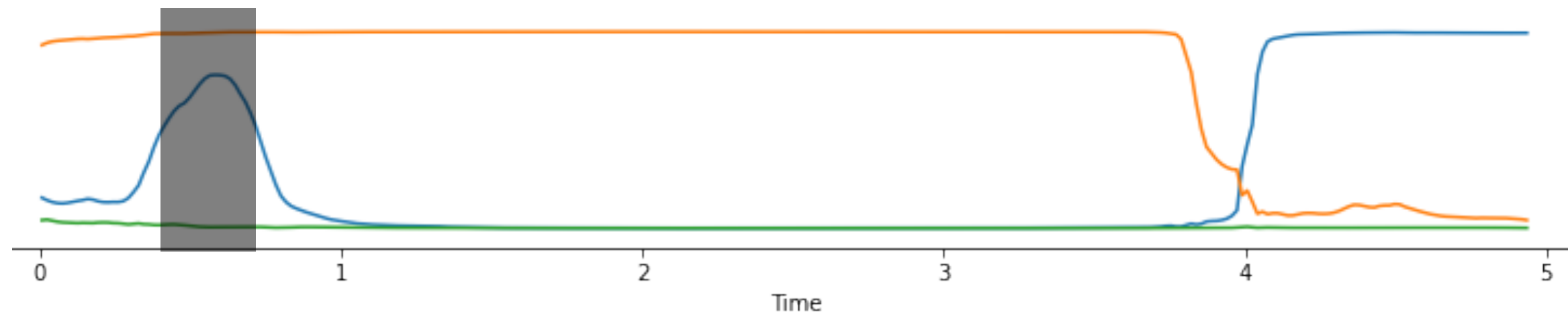
from SpeechBrain ECAPA-TDNN to WeSpeaker ResNet34





Speaker embedding

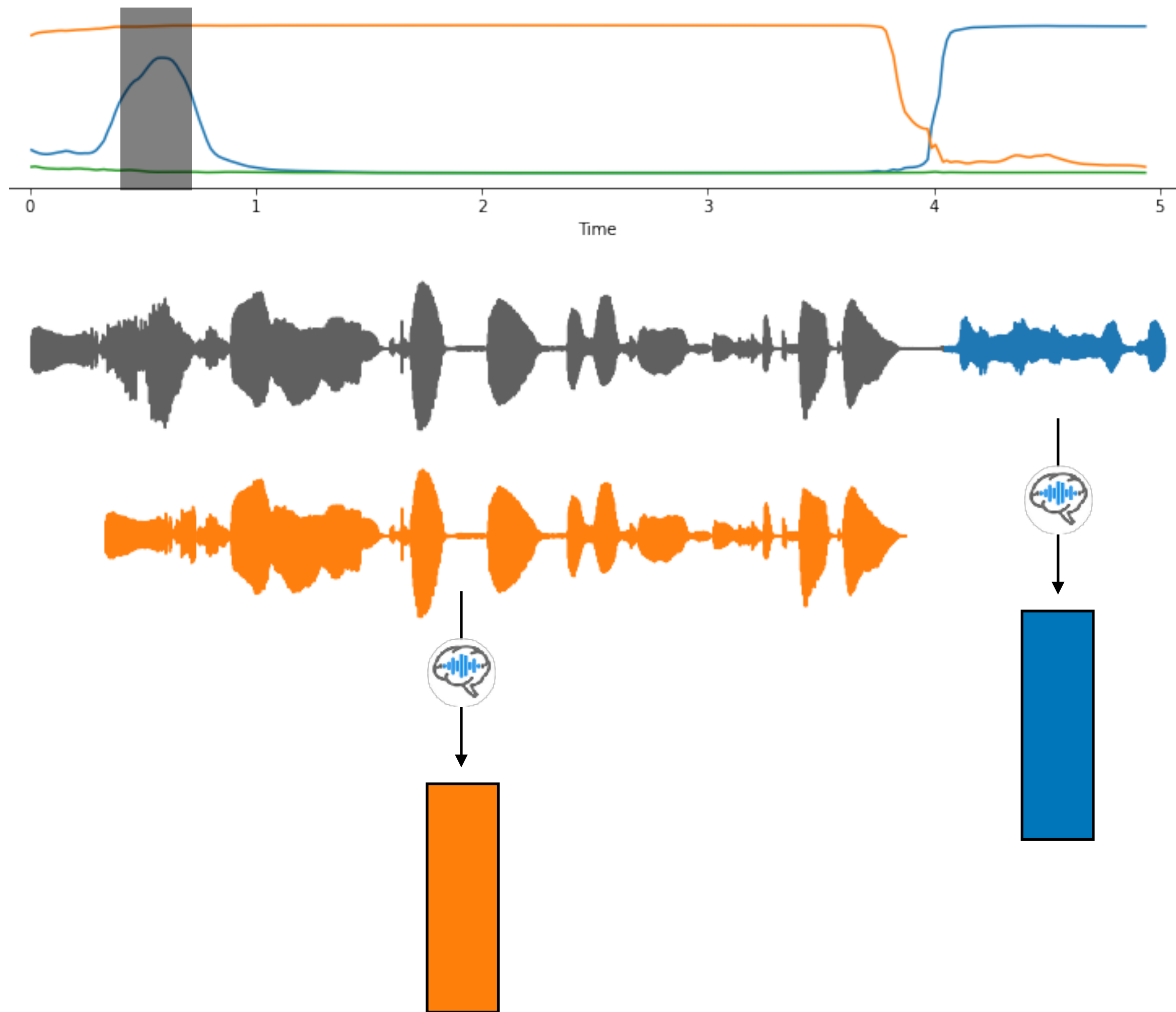
from SpeechBrain ECAPA-TDNN to WeSpeaker ResNet34





Speaker embedding

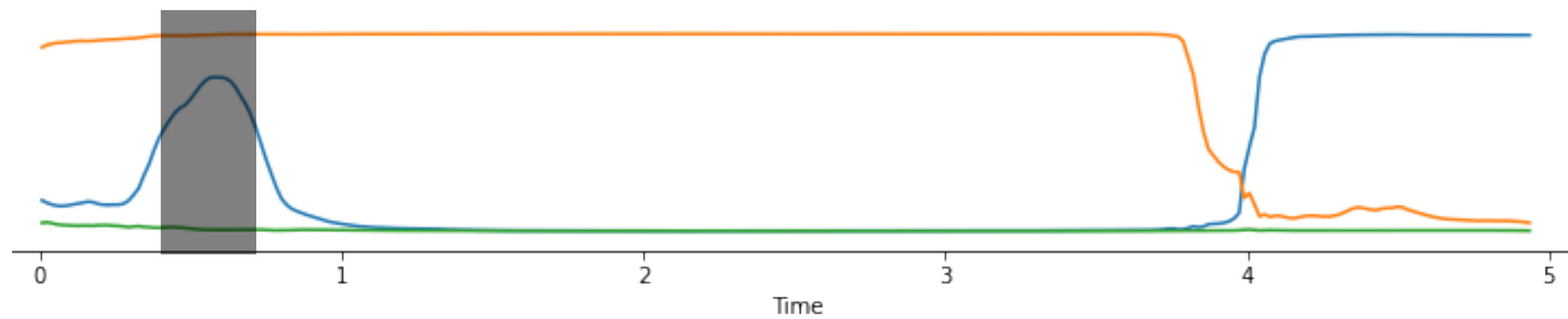
from SpeechBrain ECAPA-TDNN to WeSpeaker ResNet34



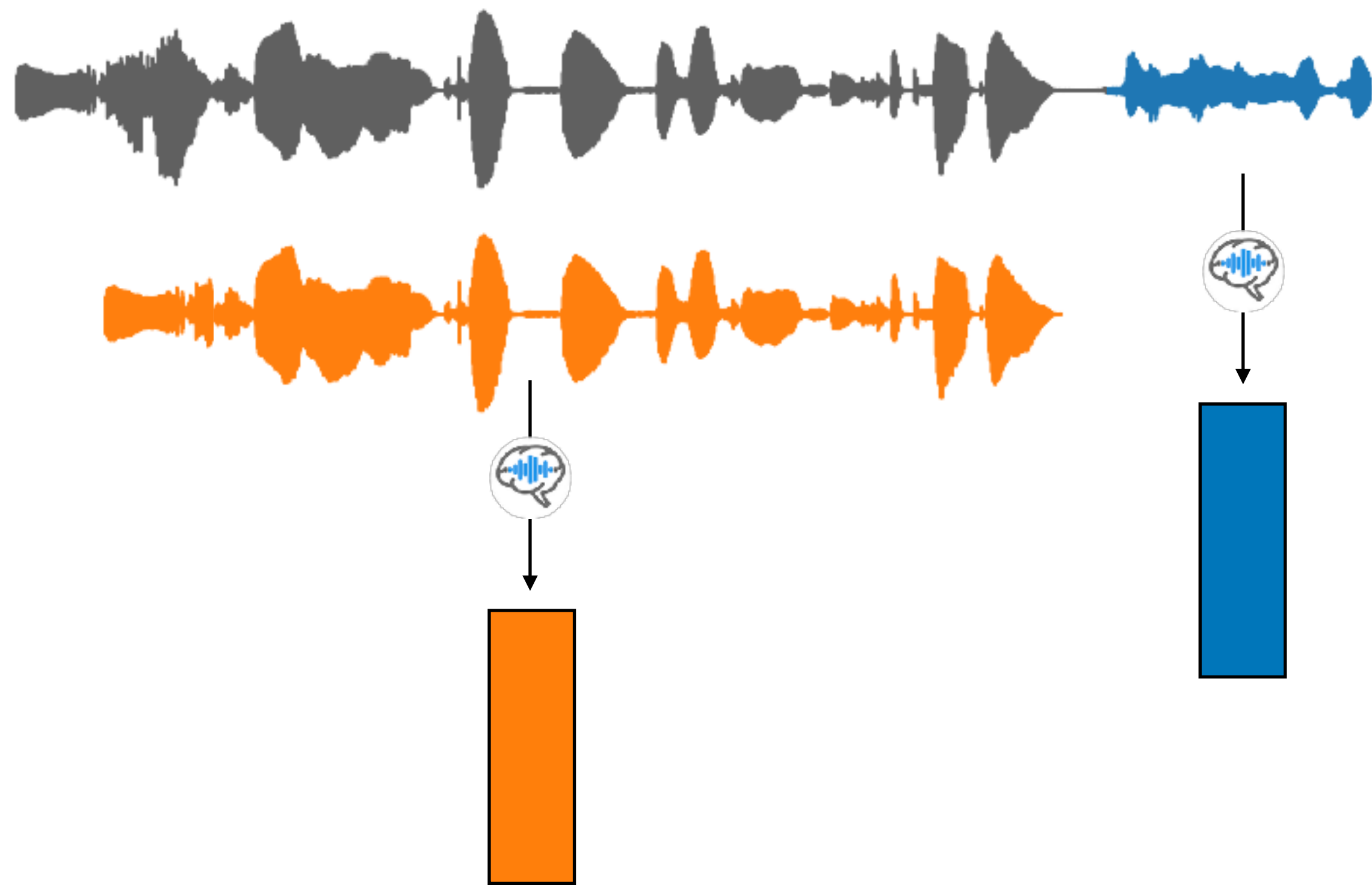


Speaker embedding

from SpeechBrain ECAPA-TDNN to WeSpeaker ResNet34



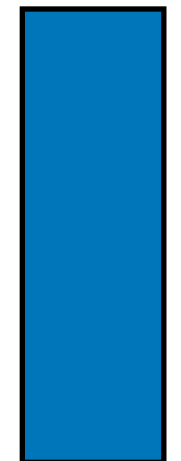
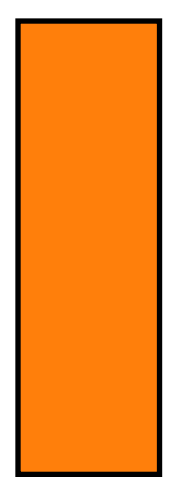
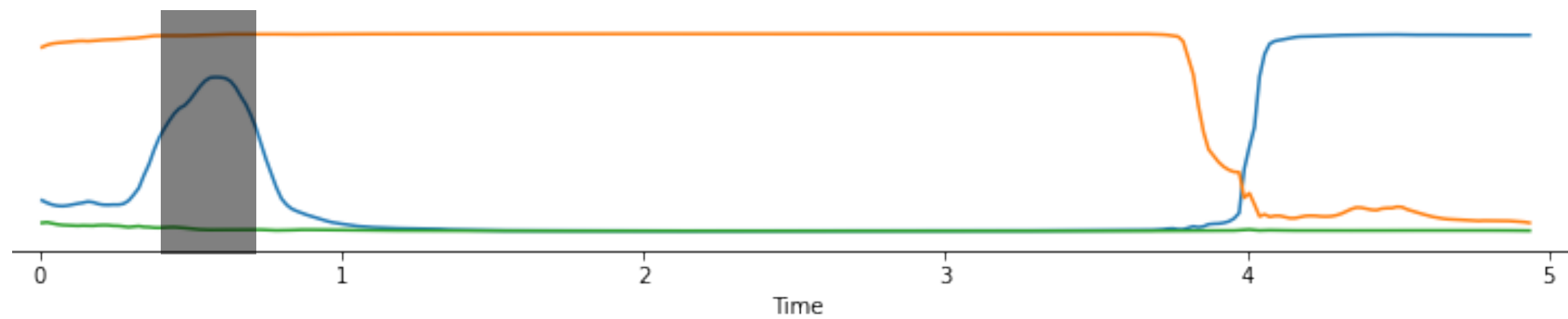
SpeechBrain
ECAPA-TDNN





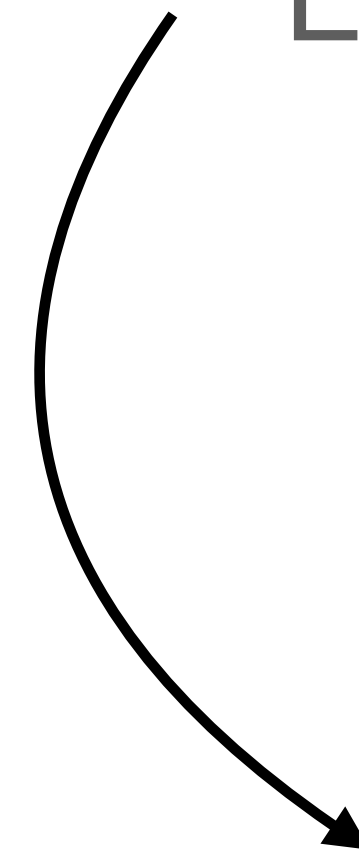
Speaker embedding

from SpeechBrain ECAPA-TDNN to WeSpeaker ResNet34



SpeechBrain
ECAPA-TDNN

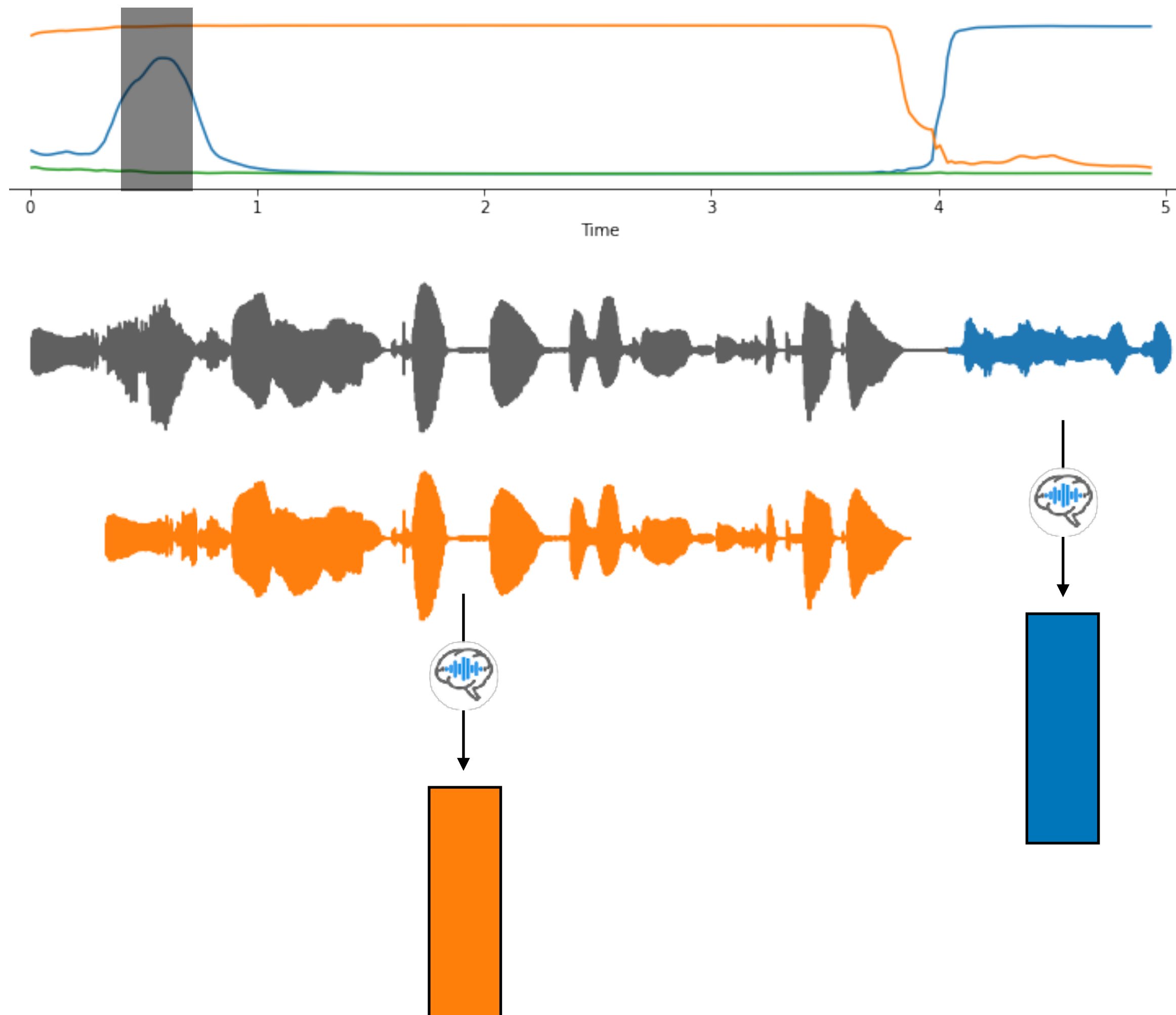
WeSpeaker
ResNet34





Speaker embedding

from SpeechBrain ECAPA-TDNN to WeSpeaker ResNet34



SpeechBrain
ECAPA-TDNN

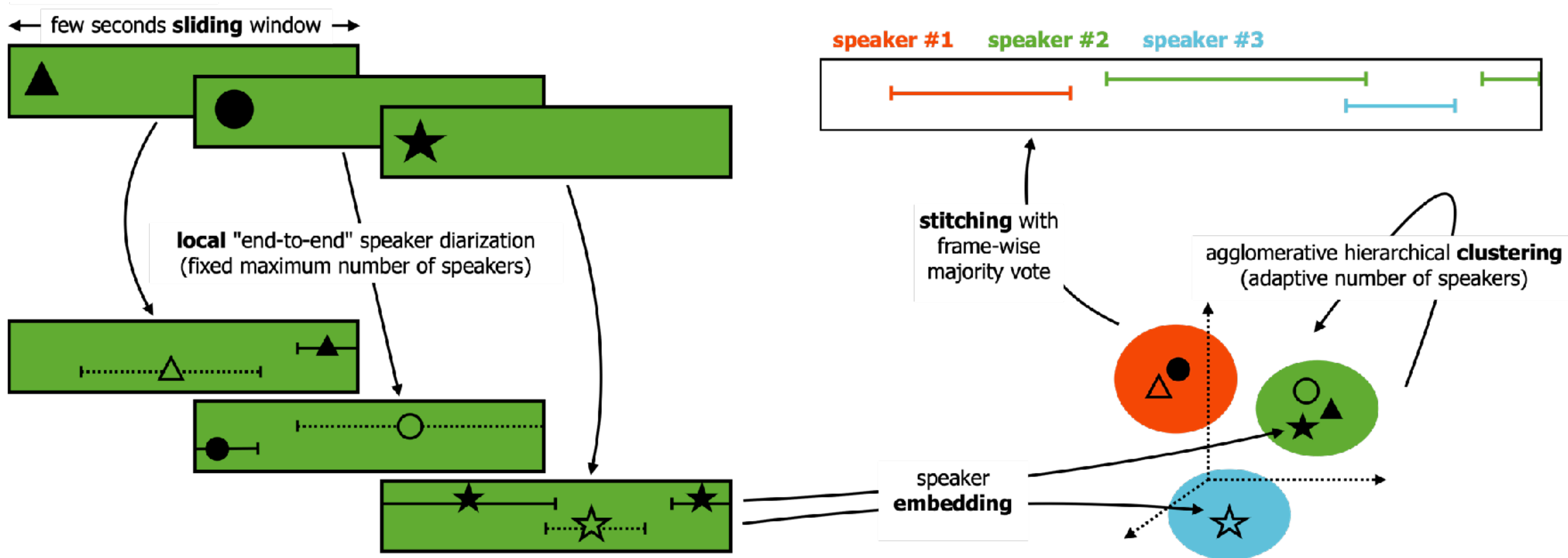
**14% relative decrease
in speaker confusion**

WeSpeaker
ResNet34



Final run

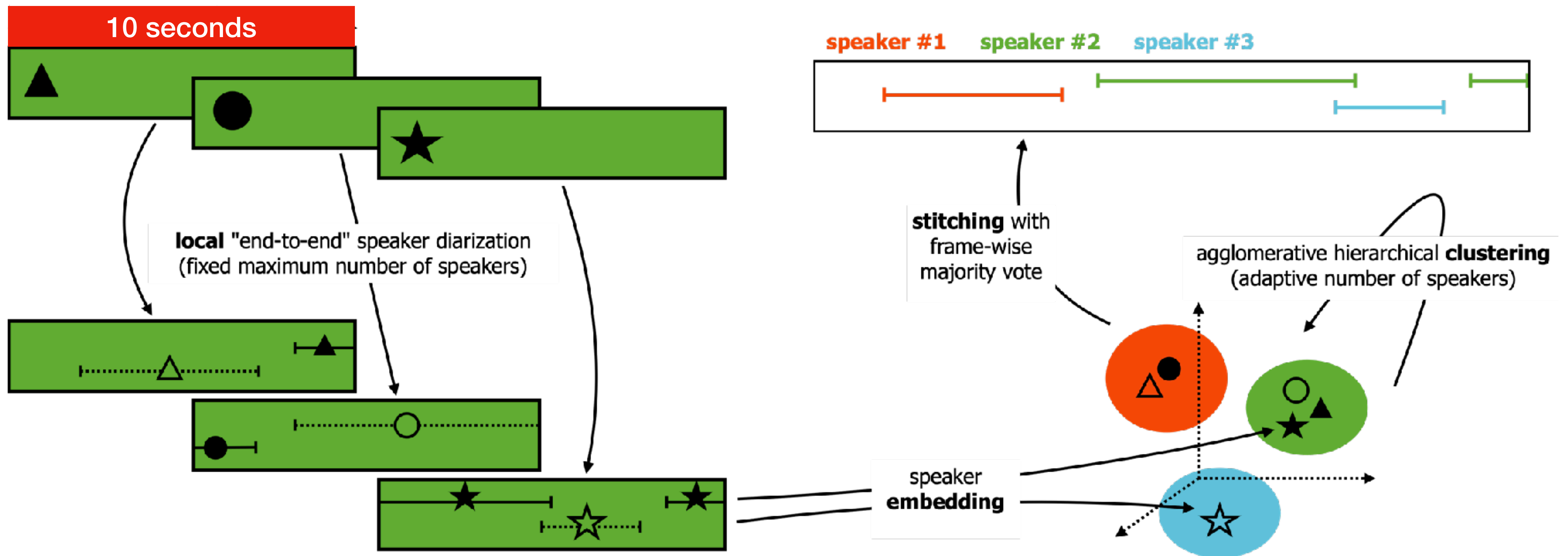
Single system (no ensembling)





Final run

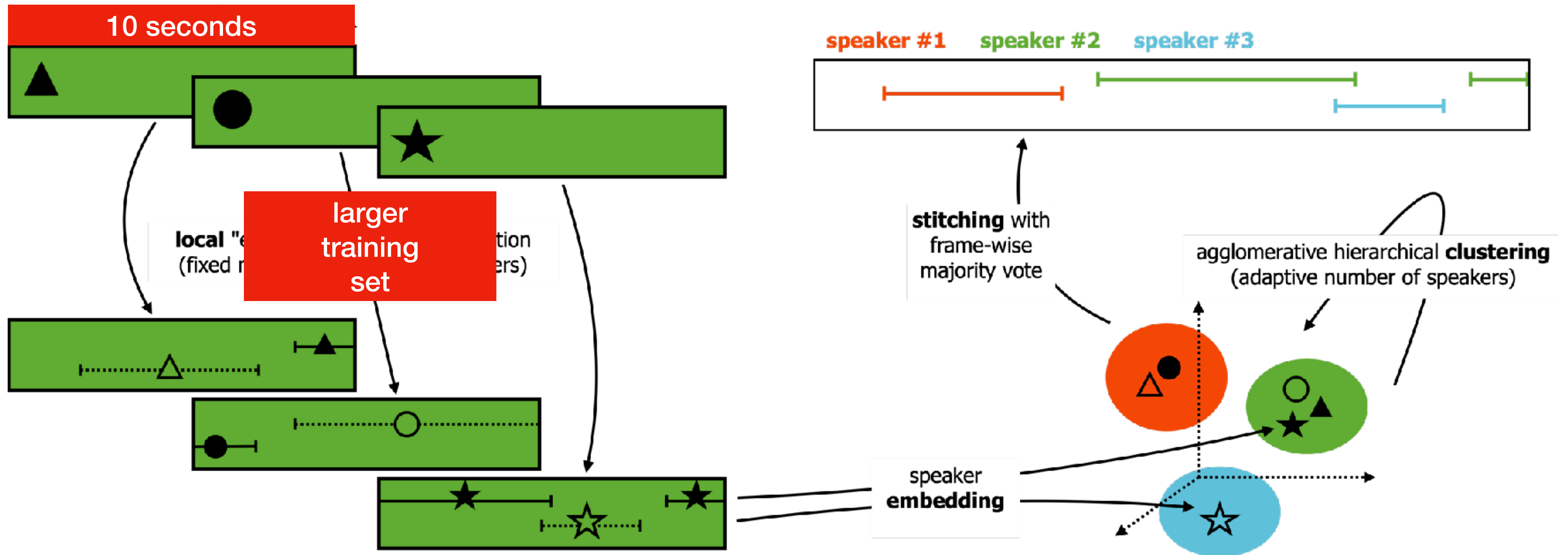
Single system (no ensembling)





Final run

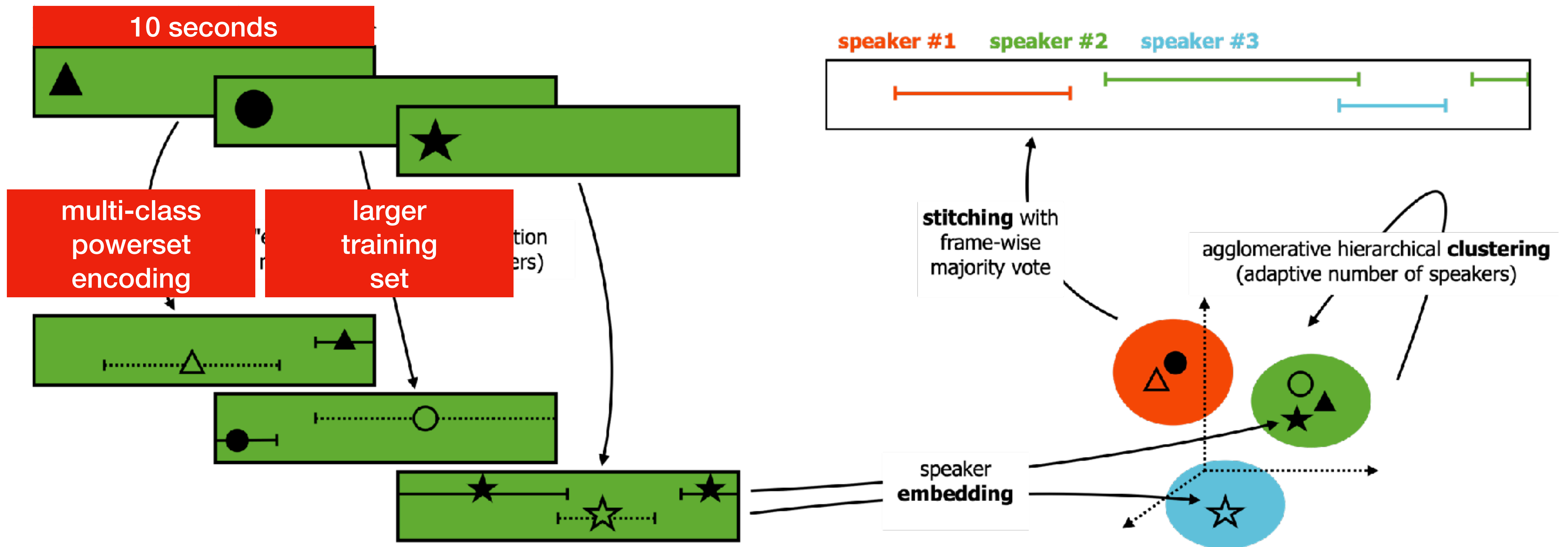
Single system (no ensembling)





Final run

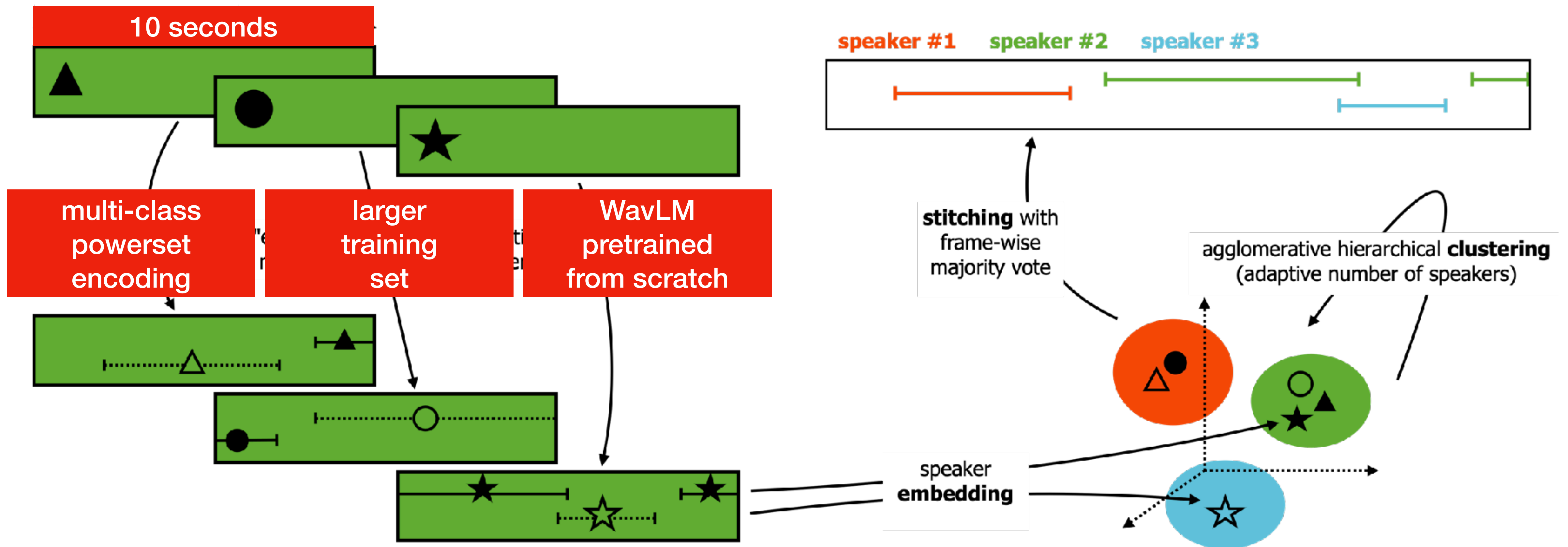
Single system (no ensembling)





Final run

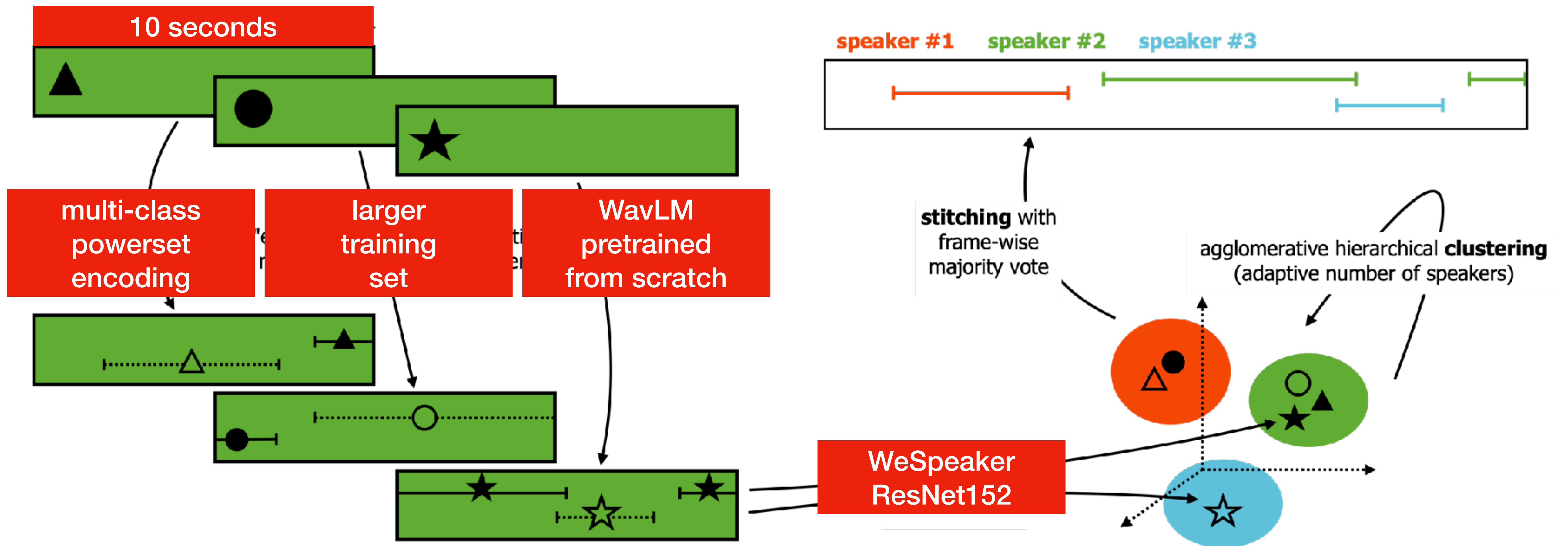
Single system (no ensembling)





Final run

Single system (no ensembling)



Acknowledgments



"Jean Zay" supercomputer



Get your stickers!

Hiring!



**pyannote.audio 2.1 speaker diarization pipeline:
principle, benchmark, and recipe**

Hervé Bredin

**Powerset multi-class cross entropy loss
for neural speaker diarization**

Alexis Plaquet & Hervé Bredin