

TS-SEP: Joint Diarization and Separation Conditioned on Estimated Speaker Embeddings

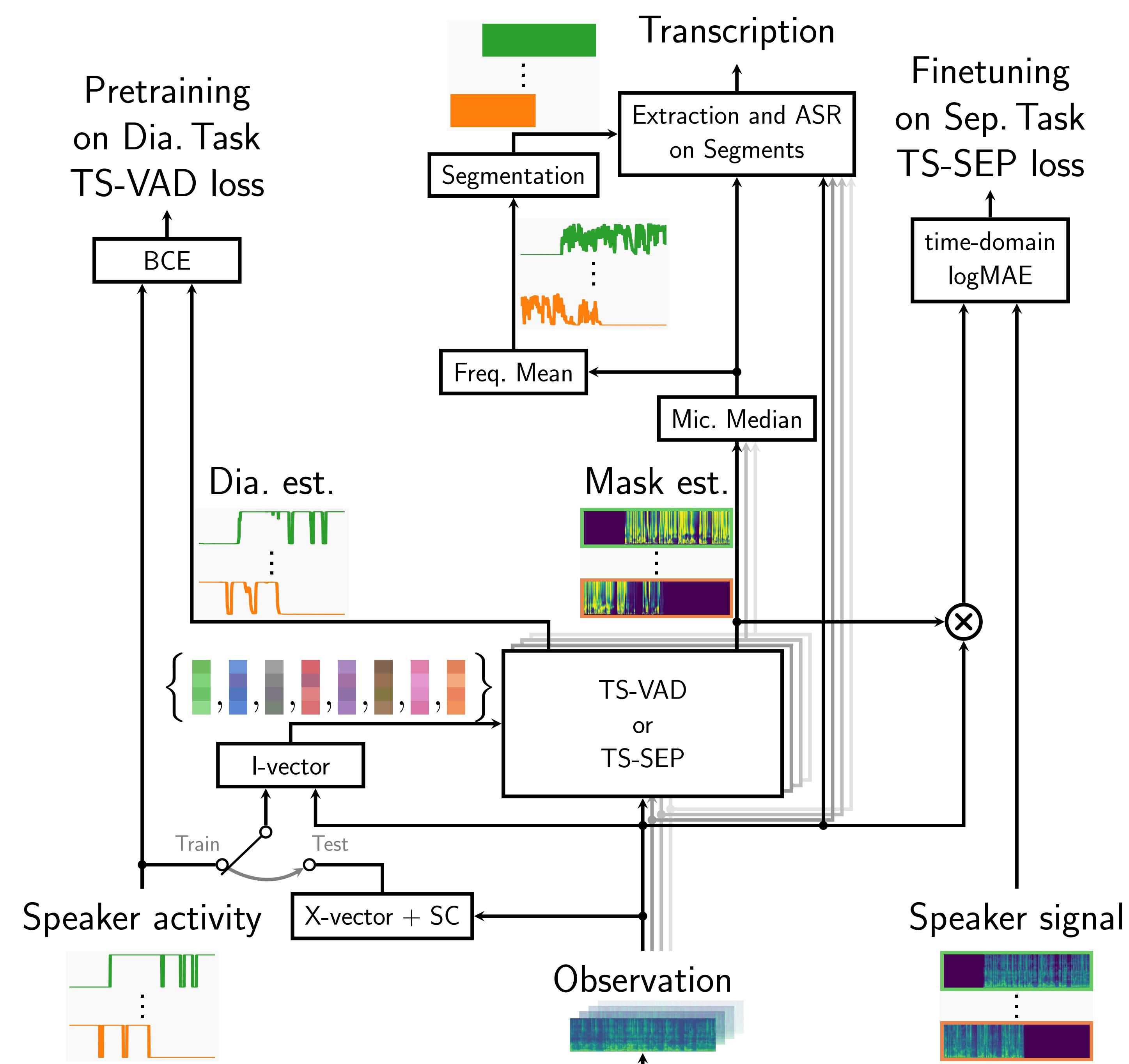
 Christoph Boeddeker^{1,2}, Aswin Shanmugam Subramanian², Gordon Wichern², Reinhold Haeb-Umbach¹, Jonathan Le Roux²
¹Paderborn University, Germany, ²Mitsubishi Electric Research Laboratories (MERL), USA

Highlights

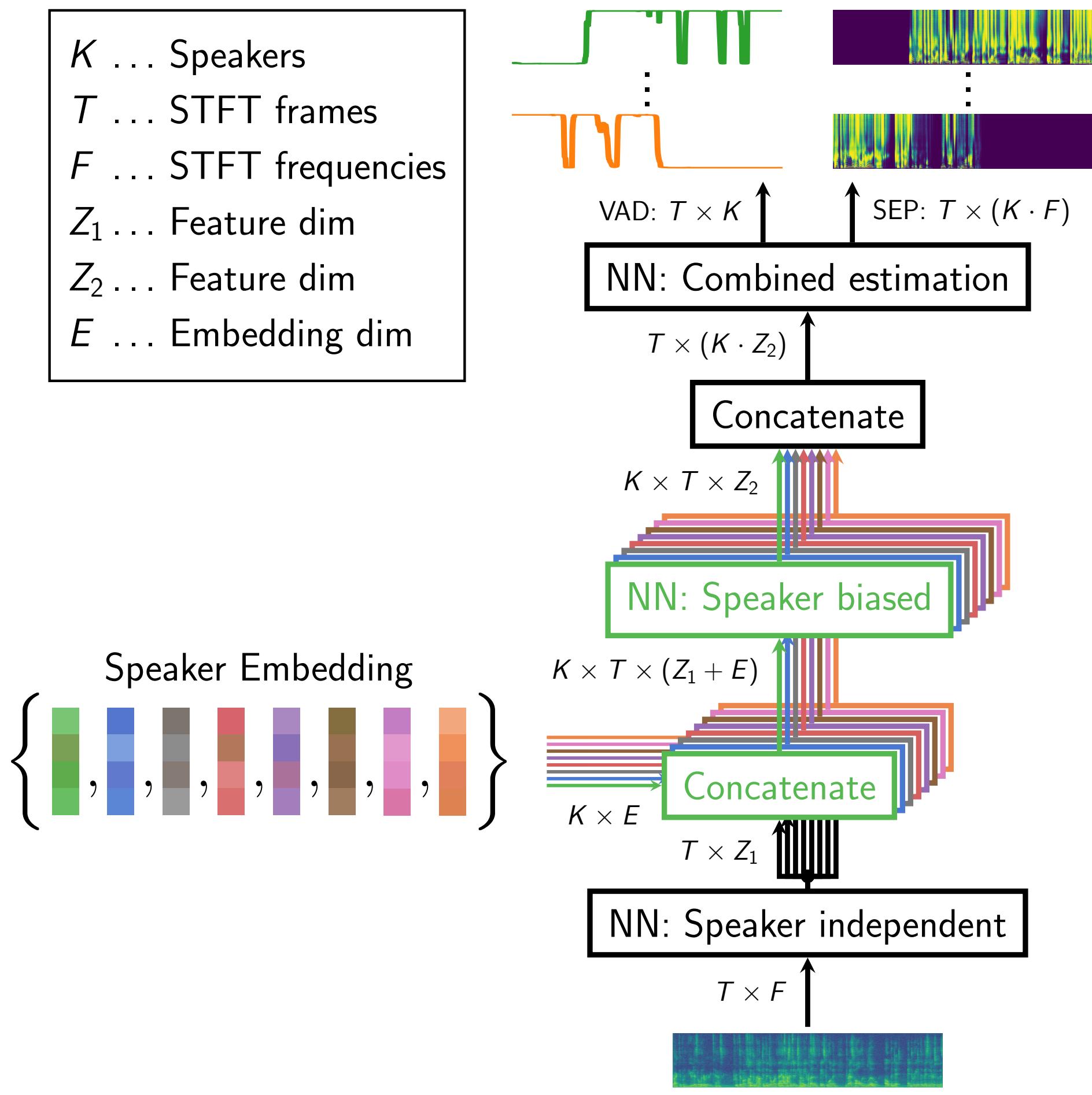
- TS-SEP: Extension of TS-VAD from diarization to joint diarization and separation
- Extensive experimental evaluation
 - ▶ Segmentation analysis
 - ◆ Trade-off between DER and cpWER
 - ◆ WER prefers “overestimation”
 - ▶ Impact of multichannel information on system components
 - ▶ Comparison of extraction techniques: masking, beamforming, GSS
 - ▶ TS-SEP enables faster GSS
- SotA on LibriCSS
 - ▶ Single channel: 11.6% → 7.81% cpWER
 - ▶ Multi channel: 5.9% → 5.36% cpWER
 - ▶ TS-VAD: 11.2% → 5.77% cpWER
- Open Source PyTorch implementation of TS-VAD and TS-SEP



System overview



TS-VAD and TS-SEP structure



Datasets

- Train
 - ▶ jsalt2020 simulate
 - ▶ Simulated meeting data
 - ▶ From LibriCSS authors
- Test
 - ▶ LibriCSS
 - ▶ Meeting scenario
 - ▶ Re-recordings of LibriSpeech sentences
 - ▶ 60 recordings, 10 min each
 - ▶ Overlap ratios: 0 % (short/long pauses), 10 %, 20 %, 30 %, 40 %

Effect of Extraction techniques

Extraction

- - WPE
 - WPE → Masking (0.0)
 - WPE → Masking (0.5)
 - WPE → BF
 - WPE → BF → Masking (0.0)
 - WPE → BF → Masking (0.5)
 - WPE → BF → Masking (0.5)
 - GSS T-init
 - GSS TF-init
- WPE ... Dereverberation technique
 • Masking (ξ) ... Clip values below ξ to ξ before masking
 • BF ... Beamforming (MVDR)
 • GSS ... Guided Source Separation (includes WPE)
 • T-init ... Classical GSS diarization init
 • TF-init ... Init GSS with TF mask

Dereverberate → 24.42

Clip small values to 0.5 → 24.08

• 10.49 → 21.07

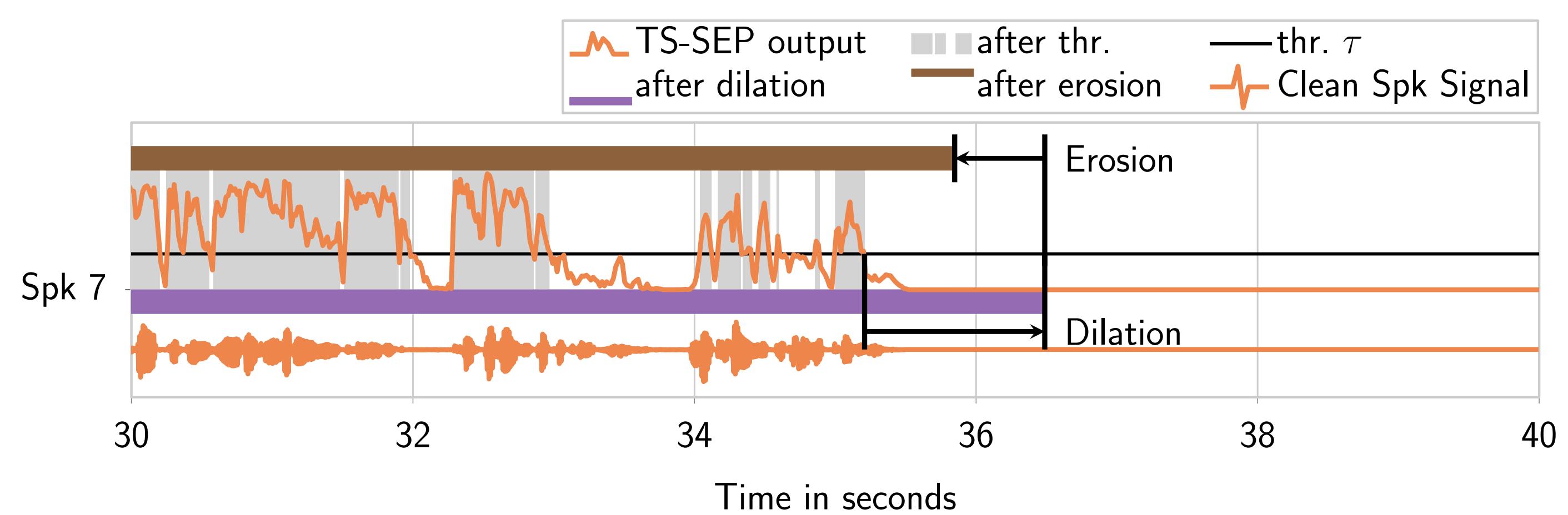
• 8.82 → 9.29 Clipping enables BF+masking

• 8.42 → 7.72 Segm. Thr: 0.6 → 0.3

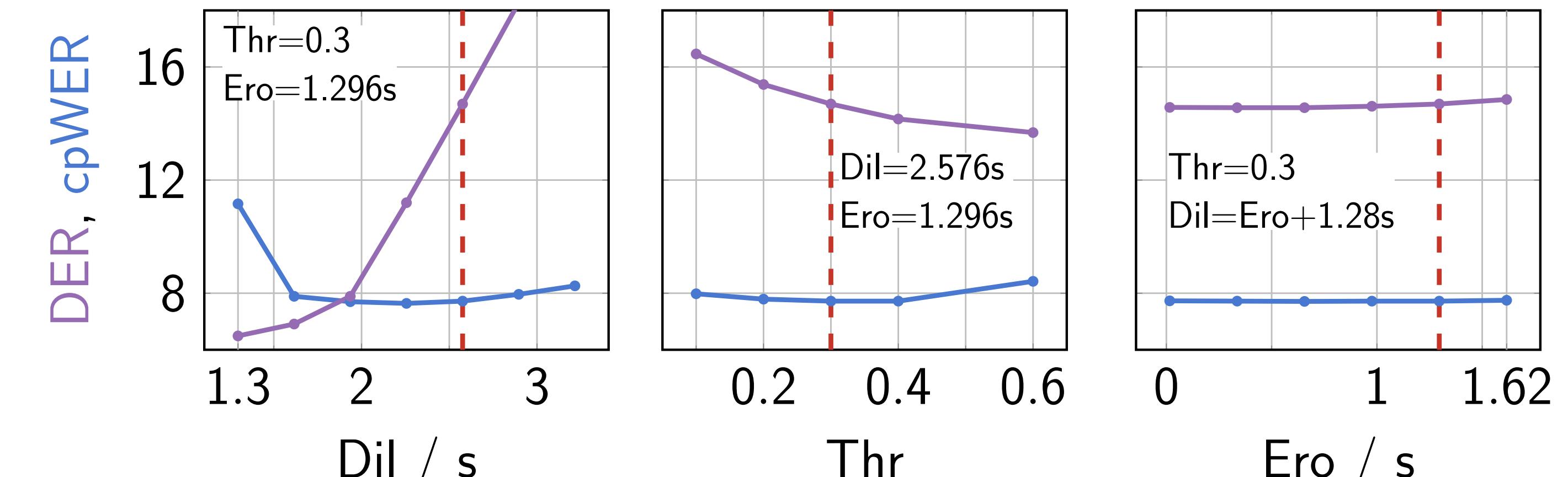
• 6.32 → 6.15

• 6.15 → EVAL cpWER

Segmentation (Threshold, Dilation, Erosion)



Impact of Segmentation



- cpWER ... concatenated minimum-permutation Word Error Rate
- DER ... Diarization Error Rate

Literature comparison (LibriCSS)

System (other)	Style	Single Ch.	cpWER
CSS with DOA Dia [58]	S → D → A	–	12.98
CSS with DOA Dia [58]	S → D → A	–	12.40
SMM [6]	D + S → A	–	5.9
CSS → SC [2]	S → D → A	–	12.7
Transcribe-to-Diarize [61]	E2E	✓	11.6
TS-VAD → Speakerbeam [12]	D → S → A	✓	18.8
TS-VAD → GSS [12]	D → S → A	–	11.2
SC → GSS [62]	D → S → A	–	12.12
System (mixed)			
SC [62] → GSS → WavLM	D → S → A	–	13.85
System (our)			
TS-VAD → GSS → Base	D → S → A	–	6.70
TS-VAD → WavLM	D → A	✓	9.26
TS-VAD → GSS → WavLM	D → S → A	–	5.77
TS-SEP → Mask. → Base	D + S → A	✓	11.61
TS-SEP → GSS → Base	D + S → A	–	6.42
TS-SEP → Mask. → WavLM	D + S → A	✓	7.81
TS-SEP → GSS → WavLM	D + S → A	–	5.36