

# LENS-DF: Deepfake Detection and Temporal Localization for Long-Form Noisy Speech

Xuechen Liu, Wanying Ge, Xin Wang, Junichi Yamagishi @ National Institute of Informatics, Japan

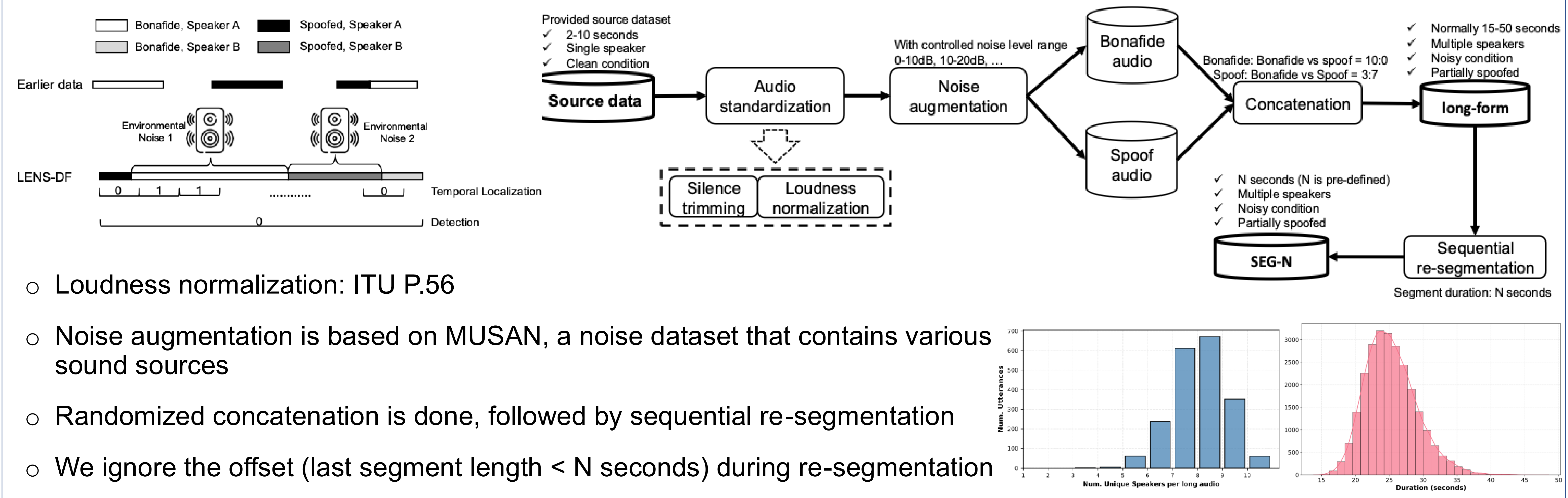
## Main Contributions

- We propose a controllable, comprehensive data complication framework that includes several artefacts that reflects real-world, in-the-wild Deepfake audios
- We benchmark the adaptability and robustness of the state-of-the-art self-supervised learning-based audio Deepfake detectors against those realistic variations
- We perform ablation analysis on those artefacts' impact on model robustness and reliability

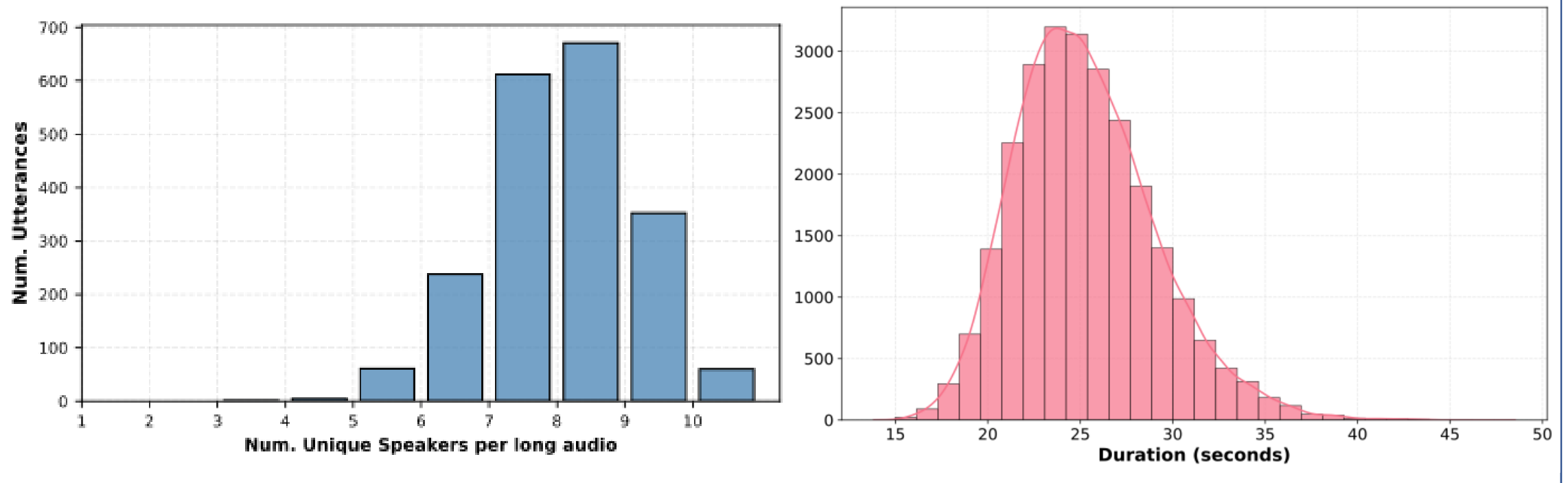
## The Conceptual Novelty

- We are not completely generating “real-world” audio, since the definition is hard and itself unrealistic
- Instead, we include multiple realistic variations/artefacts into the generated audio: longer duration, multiple speakers (mix of Bonafide and spoof), and noises
- The detection is performed per long audio, while the temporal localization is formed here as detection per short segment individually

## LENS-DF Data Complication Framework



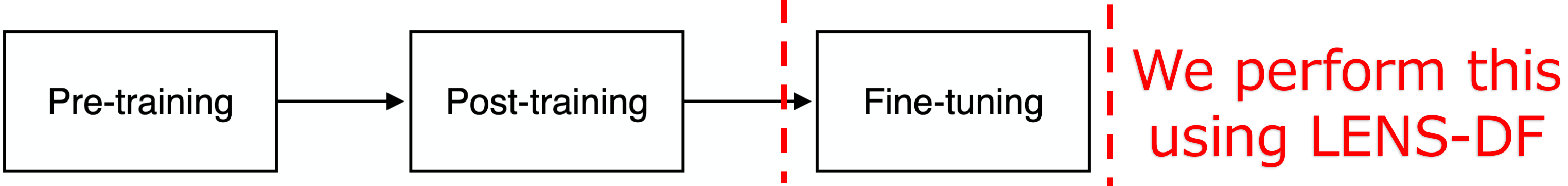
- Loudness normalization: ITU P.56
- Noise augmentation is based on MUSAN, a noise dataset that contains various sound sources
- Randomized concatenation is done, followed by sequential re-segmentation
- We ignore the offset (last segment length < N seconds) during re-segmentation



## Models & Experimental Setup

- Models: MMS-300M and MMS-1B from AntiDeepfake
- Evaluation metric: Equal error rate (EER) & Half-Total error Rate (HTER, threshold tuned by development set)

AntiDeepfake { ~74,000 hours  
>100 languages



## Main Results

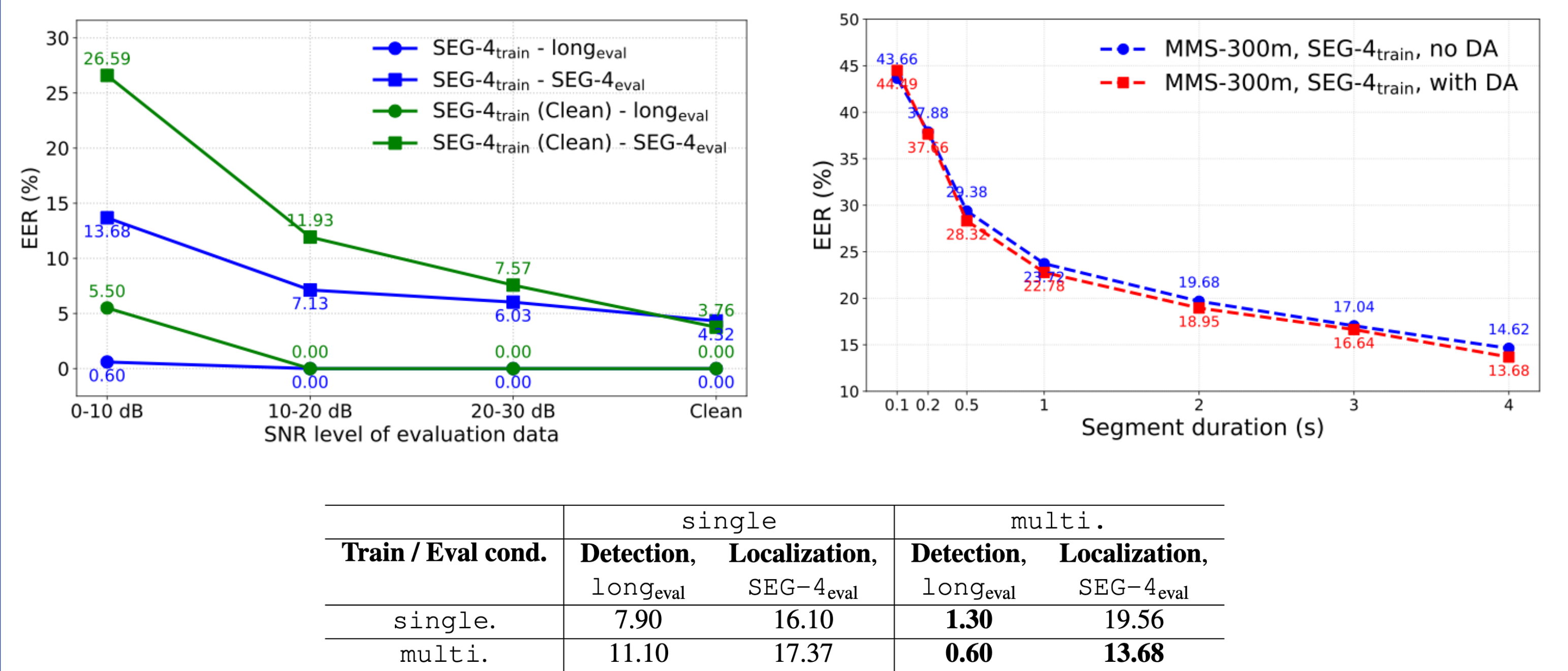
SSL Model	Training Data	Detection, 19LA <sub>eval</sub>		Detection, long <sub>eval</sub>		Localization, SEG-4 <sub>eval</sub>	
		EER (%)	HTER (%)	EER (%)	HTER (%)	EER (%)	HTER (%)
MMS-1B	19LA <sub>train</sub>	0.39	1.97	6.40	9.30	22.10	23.04
MMS-300M	19LA <sub>train</sub>	<b>0.15</b>	<b>0.52</b>	2.90	4.05	21.12	21.09
MMS-1B	long <sub>train</sub>	5.15	4.97	0.90	1.65	19.20	17.86
MMS-300M	long <sub>train</sub>	7.45	5.32	1.30	<b>1.40</b>	17.81	17.26
MMS-1B	SEG-4 <sub>train</sub>	4.62	7.98	<b>0.60</b>	4.20	15.14	14.34
MMS-300M	SEG-4 <sub>train</sub>	4.92	4.66	1.00	8.40	<b>14.62</b>	<b>14.08</b>

SSL Model	Training Data	RawBoost	Detection, 19LA <sub>eval</sub>		Detection, long <sub>eval</sub>		Localization, SEG-4 <sub>eval</sub>	
			EER (%)	HTER (%)	EER (%)	HTER (%)	EER (%)	HTER (%)
MMS-300M	19LA <sub>train</sub>	Yes	0.46	7.80	9.60	14.80	26.32	27.75
	long <sub>train</sub>	Yes	12.06	7.99	0.90	<b>1.90</b>	18.89	18.36
	SEG-4 <sub>train</sub>	Yes	8.31	6.92	<b>0.60</b>	3.80	<b>13.68</b>	<b>13.52</b>
	SEG-4 <sub>train</sub> (Clean)	Yes	6.85	7.68	5.90	12.30	26.66	26.65
XLS-R-300M	SEG-4 <sub>train</sub> (Clean)	No	3.26	8.36	8.90	9.75	29.94	31.41
	19LA <sub>train</sub>	Yes	<b>0.19</b>	<b>0.94</b>	15.70	17.60	30.41	27.30

- Acquiring LENS-DF generated data for training substantially improves adaptability on complex data
- Temporal localization remains as a challenging task even with LENS-DF
- RawBoost addition improves the performance of temporal localization, but not much

## Ablation Study



- Higher noise level in SNR and shorter segments leads to bad performance
- Single speaker cannot generalize on multi-speaker cases
- Those realistic conditions & artefacts may unravel task switching/distraction of the model, which is for future work





