

SUPPLEMENT FOR PAPER ID 3956
UNIFIED STREAMING AND NON-STREAMING MODEL FOR AUDIO-VISUAL SPEECH
RECOGNITION

Table 1. WER (%) of the proposed streaming and non-streaming models in babble noise environments. The chunk size is set to 16 for streaming models which means that the latency of streaming model is 400ms.

Modality	Streaming	SNR					Clean
		-5dB	0dB	5dB	10dB	15dB	
Audio-only	N	65.2	21.8	7.7	4.1	3.1	2.4
	Y	65.9	25.1	10.1	6.0	4.4	3.4
Audio-visual	N	25.9	10.6	4.9	3.3	2.7	2.3
	Y	33.1	14.5	7.4	5.1	4.3	3.6

Table 2. Ablation study of the Ro-conformer encoder. ✓ denotes that the Ro-conformer is used in the corresponding module, while × denotes that the Ro-conformer is not used. The evaluation metric is WER(%).

Audio back-end	Visual back-end	Fusion encoder	Test set
✓	×	✓	2.3
✓	✓	✓	2.5
✓	✓	×	2.6