

Model Architecture and Additional Evaluations for paper "Target Speaker
Extraction through Comparing Noisy Positive and Negative Audio
Enrollments"

April 1, 2025

Table 1: Comparison between our method and the baselines’ SNR (\uparrow), SNRi (\uparrow), SI-SNR (\uparrow), SI-SNRi (\uparrow), PESQ (\uparrow), STOI (\uparrow), DNSMOS (\uparrow), and WER (\downarrow) on samples synthesized from the Librispeech dataset. For the binaural audio extraction task, since target is to predict the reverberant binaural audio of the target speaker, even a perfect prediction may score poorly due to the reverberation effect. We included the WER metric here for completeness, and omit the WER in the rest evaluations when the model is expected to predict reverberant audio.

Method	Audio Type	Metric	2 Speakers	3 Speakers	4 Speakers	5 Speakers	6 Speakers
Ours (Classical-monaural)	Mono	SNR	7.40 \pm 2.17	5.51 \pm 2.23	3.68 \pm 2.47	2.44 \pm 2.43	1.12 \pm 2.35
		SNRi	10.81 \pm 2.43	10.96 \pm 2.74	10.33 \pm 2.75	10.24 \pm 2.73	9.80 \pm 2.74
		SI-SNR	6.48 \pm 3.41	3.87 \pm 4.01	0.71 \pm 5.42	-1.98 \pm 6.22	-5.45 \pm 7.12
		SI-SNRi	9.89 \pm 3.42	9.33 \pm 3.99	7.37 \pm 5.19	5.79 \pm 5.88	3.21 \pm 6.41
		PESQ	1.35 \pm 0.16	1.20 \pm 0.11	1.13 \pm 0.07	1.09 \pm 0.05	1.08 \pm 0.06
		STOI	0.79 \pm 0.10	0.70 \pm 0.13	0.60 \pm 0.15	0.50 \pm 0.17	0.40 \pm 0.18
		DNSMOS	2.98 \pm 0.22	2.87 \pm 0.20	2.81 \pm 0.18	2.75 \pm 0.16	2.71 \pm 0.16
		WER (%)	34.3 \pm 35.7	52.4 \pm 29.7	69.5 \pm 26.2	81.3 \pm 32.4	87.2 \pm 17.1
Ours (Film Fusion)	Mono	SNR	6.59 \pm 2.11	4.53 \pm 2.38	3.27 \pm 2.34	1.78 \pm 2.39	0.82 \pm 2.36
		SNRi	10.01 \pm 2.55	9.98 \pm 2.68	9.93 \pm 2.69	9.55 \pm 2.61	9.38 \pm 2.50
		SI-SNR	5.46 \pm 3.41	2.22 \pm 4.71	0.03 \pm 5.14	-3.05 \pm 5.69	-5.51 \pm 6.39
		SI-SNRi	8.88 \pm 3.52	7.66 \pm 4.50	6.67 \pm 4.92	4.73 \pm 4.94	3.03 \pm 5.61
		PESQ	1.30 \pm 0.14	1.17 \pm 0.09	1.10 \pm 0.06	1.08 \pm 0.05	1.07 \pm 0.04
		STOI	0.75 \pm 0.11	0.65 \pm 0.14	0.52 \pm 0.16	0.43 \pm 0.18	0.38 \pm 0.17
		DNSMOS	2.93 \pm 0.22	2.83 \pm 0.20	2.76 \pm 0.17	2.74 \pm 0.15	2.73 \pm 0.15
		WER (%)	37.8 \pm 27.0	63.1 \pm 38.7	73.9 \pm 37.5	84.7 \pm 18.8	89.7 \pm 29.3
TCE	Mono	SNR	4.89 \pm 2.44	2.68 \pm 2.54	1.22 \pm 2.36	-0.11 \pm 2.30	-0.83 \pm 2.08
		SNRi	8.46 \pm 2.57	8.24 \pm 2.58	8.19 \pm 2.55	7.73 \pm 2.10	7.75 \pm 2.05
		SI-SNR	3.18 \pm 3.82	-0.42 \pm 4.57	-3.28 \pm 5.10	-5.89 \pm 4.84	-8.00 \pm 4.88
		SI-SNRi	6.75 \pm 3.66	5.16 \pm 4.11	3.69 \pm 4.39	1.94 \pm 3.47	0.57 \pm 3.80
		PESQ	1.21 \pm 0.12	1.11 \pm 0.08	1.08 \pm 0.05	1.06 \pm 0.05	1.05 \pm 0.05
		STOI	0.68 \pm 0.12	0.55 \pm 0.15	0.45 \pm 0.14	0.36 \pm 0.15	0.31 \pm 0.13
		DNSMOS	2.68 \pm 0.18	2.65 \pm 0.15	2.65 \pm 0.14	2.65 \pm 0.13	2.63 \pm 0.14
		WER (%)	72.6 \pm 24.2	87.6 \pm 17.4	92.8 \pm 12.2	95.7 \pm 8.9	97.1 \pm 12.7
SpeakerBeam	Mono	SNR	-3.72 \pm 1.34	-3.77 \pm 1.29	-3.74 \pm 1.39	-3.76 \pm 1.32	-3.68 \pm 1.34
		SNRi	-0.42 \pm 2.15	1.51 \pm 2.41	3.10 \pm 2.52	3.85 \pm 2.65	4.91 \pm 2.72
		SI-SNR	-6.32 \pm 6.22	-8.18 \pm 5.33	-9.64 \pm 4.44	-10.37 \pm 4.28	-11.65 \pm 4.15
		SI-SNRi	-3.01 \pm 5.98	-2.90 \pm 5.08	-2.81 \pm 3.96	-2.75 \pm 3.68	-3.03 \pm 3.38
		PESQ	1.15 \pm 0.14	1.07 \pm 0.11	1.05 \pm 0.04	1.05 \pm 0.03	1.04 \pm 0.02
		STOI	0.47 \pm 0.28	0.40 \pm 0.19	0.32 \pm 0.17	0.28 \pm 0.15	0.25 \pm 0.11
		DNSMOS	2.57 \pm 0.18	2.50 \pm 0.18	2.43 \pm 0.14	2.47 \pm 0.13	2.40 \pm 0.12
		WER (%)	92.1 \pm 12.8	95.4 \pm 9.9	95.8 \pm 11.2	98.9 \pm 8.3	97.4 \pm 8.1
Ours (Classical-binaural)	Bi	SNR	5.30 \pm 1.98	3.34 \pm 1.87	2.14 \pm 1.51	1.40 \pm 1.28	0.65 \pm 0.78
		SNRi	8.97 \pm 2.39	8.83 \pm 2.64	9.04 \pm 2.53	9.28 \pm 2.51	10.06 \pm 2.67
		SI-SNR	3.52 \pm 3.27	-0.09 \pm 4.40	-2.73 \pm 4.14	-5.22 \pm 4.41	-7.97 \pm 3.62
		SI-SNRi	7.18 \pm 3.25	5.42 \pm 4.26	4.16 \pm 3.82	2.66 \pm 3.85	1.46 \pm 3.22
		PESQ	1.40 \pm 0.20	1.20 \pm 0.13	1.12 \pm 0.09	1.10 \pm 0.11	1.09 \pm 0.15
		STOI	0.73 \pm 0.09	0.56 \pm 0.15	0.45 \pm 0.15	0.36 \pm 0.14	0.29 \pm 0.13
		DNSMOS	2.78 \pm 0.20	2.68 \pm 0.18	2.61 \pm 0.16	2.56 \pm 0.15	2.52 \pm 0.14
		WER (%)	69.6 \pm 26.4	88.2 \pm 18.0	93.7 \pm 10.7	97.9 \pm 53.7	96.3 \pm 7.3
LookOnceToHear	Bi	SNR	3.58 \pm 2.33	2.30 \pm 2.66	1.51 \pm 2.05	0.91 \pm 1.92	0.12 \pm 1.73
		SNRi	7.10 \pm 3.16	7.72 \pm 3.16	8.34 \pm 2.75	8.62 \pm 2.72	9.28 \pm 2.61
		SI-SNR	0.48 \pm 5.60	-2.55 \pm 6.60	-4.54 \pm 6.14	-6.53 \pm 6.09	-9.29 \pm 5.38
		SI-SNRi	4.05 \pm 5.80	2.93 \pm 6.45	2.32 \pm 6.02	1.34 \pm 5.71	0.09 \pm 5.04
		PESQ	1.30 \pm 0.19	1.18 \pm 0.12	1.14 \pm 0.14	1.11 \pm 0.15	1.11 \pm 0.12
		STOI	0.65 \pm 0.16	0.52 \pm 0.19	0.42 \pm 0.18	0.36 \pm 0.16	0.30 \pm 0.15
		DNSMOS	2.72 \pm 0.23	2.66 \pm 0.20	2.60 \pm 0.18	2.56 \pm 0.17	2.54 \pm 0.15
		WER (%)	70.8 \pm 28.7	87.8 \pm 17.8	95.2 \pm 45.7	95.1 \pm 10.7	96.7 \pm 9.3

Table 2: Comparison of model parameter size, inference time, and memory usage when extracting 1-second audio on an Intel Xeon Silver 4314 @ 2.40GHz CPU. We report the average inference time for extraction in 50 repeated experiments. Since our model replaced the Film fusion method in the LookOnceToHear baseline with a cross-attention based fusion method, our model has significantly smaller parameter size.

Model	Param Size	Inference Time	Inference Memory Usage
Ours	1.88 M	0.37s	2.37 GB
TCE	2.54 M	0.11s	1.94 GB
LookOnceToHear	4.41 M	0.35s	2.38 GB

Table 3: Affect of audio reverberation on model performance. We fine-tune our model on reverberant monaural audio and test its performance to extract reverberant monaural audio. In comparison to our model’s performance on non-reverberant monaural audio (Shown in Table 1), the fine-tuned (Reverb-monoaural) model shows lower performance, demonstrating the significant difficulty in extracting reverberant audio.

Method	Audio Type	Metric	2 Speakers	3 Speakers	4 Speakers	5 Speakers	6 Speakers
Ours (Reverb-monoaural)	Mono	SNR	3.35 ± 1.97	1.56 ± 2.29	0.09 ± 2.28	-0.80 ± 2.43	-1.68 ± 2.31
		SNRi	7.25 ± 2.43	7.49 ± 2.33	7.39 ± 2.13	7.68 ± 1.93	7.69 ± 1.80
		SI-SNR	0.29 ± 4.00	-3.23 ± 5.06	-6.26 ± 5.16	-8.12 ± 5.26	-10.05 ± 5.10
		SI-SNRi	4.18 ± 3.64	2.69 ± 4.14	1.05 ± 4.06	0.34 ± 3.96	-0.69 ± 3.70
		PESQ	1.19 ± 0.12	1.12 ± 0.07	1.09 ± 0.05	1.08 ± 0.06	1.07 ± 0.05
		STOI	0.56 ± 0.11	0.44 ± 0.13	0.34 ± 0.12	0.29 ± 0.12	0.24 ± 0.11
		DNSMOS	2.65 ± 0.18	2.61 ± 0.17	2.60 ± 0.16	2.60 ± 0.16	2.56 ± 0.15

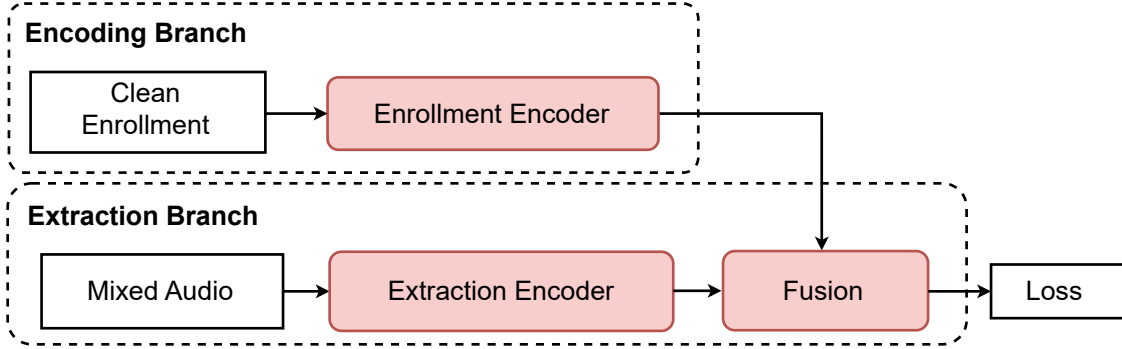
Table 4: Comparison between our monaural model and the baseline model TCE’s SNR (\uparrow), SNRi (\uparrow), SI-SNR (\uparrow), SI-SNRi (\uparrow), PESQ (\uparrow), STOI (\uparrow), DNSMOS (\uparrow), and WER (\downarrow) performance on samples synthesized from the WSJ0 dataset.

Method	Audio Type	Metric	2 Speakers	3 Speakers	4 Speakers	5 Speakers	6 Speakers
Ours (Classical-monaural)	Mono	SNR	5.59 \pm 2.72	3.42 \pm 2.63	1.77 \pm 2.72	0.79 \pm 2.37	0.10 \pm 2.10
		SNRi	9.00 \pm 2.99	8.47 \pm 2.80	8.40 \pm 2.54	8.34 \pm 2.27	8.30 \pm 1.97
		SI-SNR	3.23 \pm 6.49	-0.78 \pm 7.03	-4.20 \pm 7.55	-6.54 \pm 7.25	-7.81 \pm 6.38
		SI-SNRi	6.63 \pm 6.58	4.31 \pm 6.93	2.42 \pm 7.07	1.03 \pm 6.83	0.40 \pm 5.81
		PESQ	1.31 \pm 0.14	1.20 \pm 0.11	1.13 \pm 0.08	1.10 \pm 0.05	1.09 \pm 0.06
		STOI	0.70 \pm 0.16	0.57 \pm 0.19	0.46 \pm 0.18	0.37 \pm 0.19	0.32 \pm 0.16
		DNSMOS	2.75 \pm 0.27	2.66 \pm 0.22	2.59 \pm 0.20	2.56 \pm 0.18	2.51 \pm 0.16
TCE	Mono	SNR	3.45 \pm 2.23	2.08 \pm 2.02	0.75 \pm 1.99	0.08 \pm 1.89	-0.83 \pm 1.77
		SNRi	6.65 \pm 2.60	6.96 \pm 2.10	7.28 \pm 1.97	7.15 \pm 1.49	7.39 \pm 1.24
		SI-SNR	-0.19 \pm 6.10	-3.12 \pm 5.73	-6.45 \pm 5.98	-8.01 \pm 6.00	-9.94 \pm 5.38
		SI-SNRi	3.02 \pm 6.15	1.75 \pm 5.31	0.10 \pm 5.43	-0.90 \pm 4.88	-1.70 \pm 4.26
		PESQ	1.19 \pm 0.11	1.13 \pm 0.08	1.08 \pm 0.06	1.08 \pm 0.05	1.07 \pm 0.08
		STOI	0.56 \pm 0.17	0.46 \pm 0.17	0.36 \pm 0.16	0.31 \pm 0.14	0.26 \pm 0.12
		DNSMOS	2.59 \pm 0.17	2.59 \pm 0.15	2.58 \pm 0.13	2.56 \pm 0.14	2.57 \pm 0.13

Table 5: Affect of enrollment audio quality on model performance. We compare ours and the baseline methods' performance when noisy or clean target speaker enrollment is used.

Method	Audio Type	Metric	2 Speakers (Clean Enroll.)	2 Speakers (Noisy Enroll.)	3 Speakers (Clean Enroll.)	3 Speakers (Noisy Enroll.)
Ours (Classical-monaural)	Mono	SNR	5.99 \pm 2.90	7.40 \pm 2.17	3.72 \pm 2.81	5.59 \pm 2.23
		SNRi	9.27 \pm 2.87	10.81 \pm 2.43	9.01 \pm 2.65	10.90 \pm 2.68
		SI-SNR	4.19 \pm 4.89	6.49 \pm 3.41	0.72 \pm 5.38	3.99 \pm 3.95
		SI-SNRi	7.47 \pm 4.65	9.90 \pm 3.42	6.01 \pm 4.92	9.31 \pm 3.99
		PESQ	1.26 \pm 0.16	1.35 \pm 0.16	1.14 \pm 0.09	1.20 \pm 0.11
		STOI	0.73 \pm 0.14	0.79 \pm 0.10	0.61 \pm 0.16	0.70 \pm 0.13
		DNSMOS	2.94 \pm 0.23	2.98 \pm 0.22	2.85 \pm 0.18	2.87 \pm 0.20
		WER (%)	60.7 \pm 36.2	34.3 \pm 35.7	72.9 \pm 29.9	52.4 \pm 29.7
SpeakerBeam	Mono	SNR	9.29 \pm 3.21	-3.72 \pm 1.34	4.97 \pm 3.85	-3.77 \pm 1.29
		SNRi	12.87 \pm 3.63	-0.42 \pm 2.15	10.36 \pm 3.02	1.51 \pm 2.41
		SI-SNR	8.28 \pm 6.20	-6.32 \pm 6.22	2.97 \pm 6.38	-8.18 \pm 5.33
		SI-SNRi	11.85 \pm 6.36	-3.01 \pm 5.98	8.37 \pm 5.27	-2.90 \pm 5.08
		PESQ	1.52 \pm 0.24	1.15 \pm 0.14	1.23 \pm 0.17	1.07 \pm 0.11
		STOI	0.83 \pm 0.14	0.47 \pm 0.28	0.67 \pm 0.17	0.40 \pm 0.19
		DNSMOS	2.64 \pm 0.18	2.57 \pm 0.18	2.57 \pm 0.18	2.50 \pm 0.18
		WER (%)	30.5 \pm 21.9	101.4 \pm 7.3	39.8 \pm 15.7	101.6 \pm 8.7
Ours (Classical-binaural)	Bi	SNR	4.28 \pm 2.75	5.30 \pm 1.98	2.09 \pm 2.42	3.34 \pm 1.87
		SNRi	7.94 \pm 2.91	8.97 \pm 2.39	7.62 \pm 2.65	8.83 \pm 2.64
		SI-SNR	1.88 \pm 4.73	3.52 \pm 3.27	-1.93 \pm 5.02	-0.09 \pm 4.40
		SI-SNRi	5.55 \pm 4.51	7.18 \pm 3.25	3.61 \pm 4.70	5.42 \pm 4.26
		PESQ	1.31 \pm 0.22	1.40 \pm 0.20	1.16 \pm 0.13	1.20 \pm 0.13
		STOI	0.66 \pm 0.15	0.73 \pm 0.09	0.51 \pm 0.17	0.56 \pm 0.15
		DNSMOS	2.80 \pm 0.18	2.78 \pm 0.20	2.74 \pm 0.16	2.68 \pm 0.18
LookOnceToHear	Bi	SNR	4.01 \pm 2.02	3.58 \pm 2.33	3.00 \pm 1.93	2.30 \pm 2.66
		SNRi	7.69 \pm 3.01	7.10 \pm 3.16	8.38 \pm 2.74	7.72 \pm 3.16
		SI-SNR	1.71 \pm 4.00	0.48 \pm 5.60	-0.53 \pm 4.53	-2.55 \pm 6.60
		SI-SNRi	5.43 \pm 4.44	4.05 \pm 5.80	4.86 \pm 4.60	2.93 \pm 6.45
		PESQ	1.32 \pm 0.19	1.30 \pm 0.19	1.21 \pm 0.14	1.18 \pm 0.12
		STOI	0.68 \pm 0.12	0.65 \pm 0.16	0.57 \pm 0.15	0.52 \pm 0.19
		DNSMOS	2.75 \pm 0.22	2.72 \pm 0.23	2.65 \pm 0.19	2.66 \pm 0.20

a) Prior works' Model Architecture



b) Proposed Model Architecture

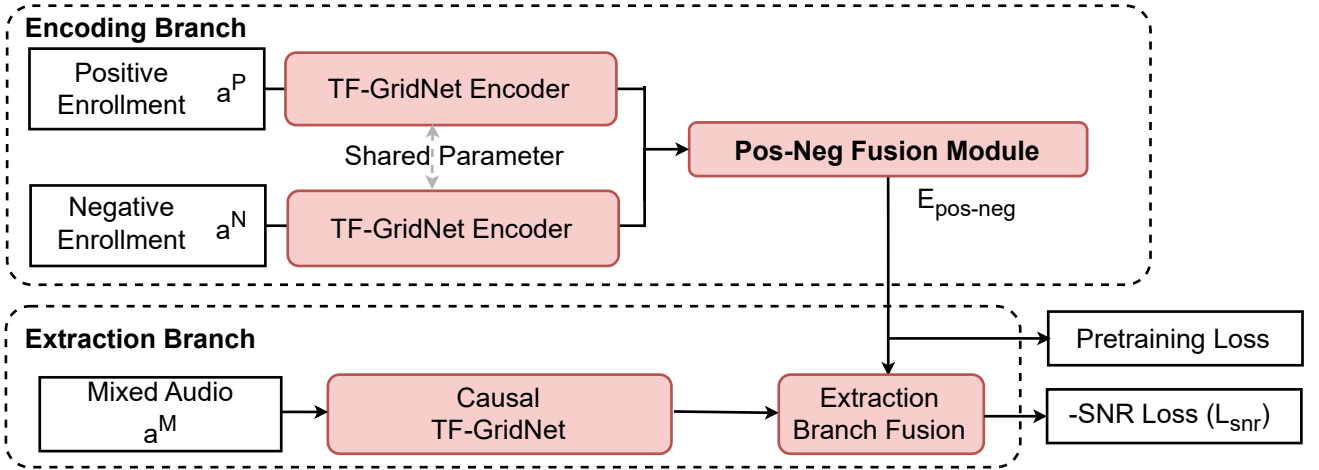


Figure 1: Comparison of the proposed model architecture with the commonly used architecture in prior works (e.g. LookOnceToHear, SpeakerBeam). We encode the Positive and the Negative Enrollments separately using two encoders with shared parameters, and fuse their embeddings using an attention based Pos-Neg Fusion Module.