

Exponential Smoothing of Noisy Embeddings for Phoneme-Aware Speech Enhancement

Anonymous ACL submission

Abstract

Speech enhancement (SE) improves the robustness of downstream speech technologies under noisy conditions. Self-supervised models such as Wav2Vec 2.0 produce robust frame-level representations that capture phonetic information, but directly conditioning on noisy embeddings can propagate errors. In this work, we propose a temporal abstraction strategy that applies exponential smoothing to Wav2Vec 2.0 embeddings before conditioning a diffusion-based SE network via FiLM modulation. This approach reduces sensitivity to noise. We evaluate the method on a noisy dataset and demonstrate improvements in PESQ and STOI across multiple SNRs and model configurations. Ablations show that exponential smoothing outperforms both naive averaging and unconditioned diffusion models by a PESQ score of 0.35.

1 Introduction

Robust speech enhancement (SE) is a cornerstone front-end for many downstream speech technologies: improving signal quality before recognition or downstream processing often yields large gains in automatic speech recognition (ASR), keyword spotting, and other speech understanding tasks, especially in adverse acoustic conditions. Prior work (Nayem and Williamson, 2021) has shown that incorporating textual or linguistic information into SE can substantially improve intelligibility and downstream recognition performance by providing high-level constraints that guide denoising and spectral reconstruction.

Recent advances in large-scale weakly-supervised ASR, most notably OpenAI’s Whisper family (Radford et al., 2023) demonstrate that models trained on massive, diverse, and partly noisy web-scale audio–transcript pairs learn representations and transcription robustness that generalize well across noise types, languages, and recording conditions. This robustness arises from

scale and exposure to naturally noisy real-world data, which makes Whisper-style transcriptions an attractive source of supervision when clean transcripts are unavailable.

Motivated by these facts, several lines (Yang et al., 2023) of work have explored conditioning SE models on linguistic cues extracted from the speech signal, by fusing text with noisy waveform through an attention mechanism. Early text-informed SE methods (Wang et al., 2022) showed that supplying word-level transcripts can guide enhancement toward linguistically plausible spectra and improve intelligibility when the text is accurate. More recent work (Lu et al., 2023) shifts from lexical information to intermediate phonetic representations, such as phoneme posteriors and broad phonetic classes, which are easier to synchronize with acoustic features and remain reliable under higher noise levels. These studies collectively demonstrate both the promise and the limitations of linguistic conditioning: linguistic structure significantly helps enhancement, but real-world transcripts are often noisy or imperfectly aligned, which can propagate errors if used directly.

In this paper we leverage the robustness of a large-scale feature extractor like wav2vec 2.0 (Baeovski et al., 2020) but intentionally avoid relying directly on noisy word-level transcripts at inference time. Instead, we extract noisy frame level phonetic information and use these noisy phonemes as conditioning for a neural speech enhancement model. Conditioning on phonemes (rather than raw text) reduces sensitivity to word-level errors, exploits phoneme-specific spectral priors as documented in prior work (Lu et al., 2020), and keeps the conditioning compact and aligned with typical frame-level SE architectures. They suggest that phonetic conditioning is effective even when phonetic inputs are imperfect, motivating our design choice.

We systematically investigate the use of large-

scale self-supervised speech representations for phoneme-aware speech enhancement. In particular, we propose a temporal abstraction strategy that applies exponential smoothing to Wav2Vec 2.0 derived embeddings before conditioning a diffusion-based SE network via FiLM (Perez et al., 2018). Unlike naive averaging, this smoothing preserves phonetic information across time, reducing sensitivity to noisy embeddings while minimizing additional computation. Through extensive experiments, we show that this simple yet effective design consistently improves perceptual quality (PESQ) and intelligibility (STOI) across diverse SNR conditions, highlighting the importance of temporal aggregation in embedding-based SE. This study frames temporal smoothing as a compute-efficient, phonetic-aware enhancement strategy, providing insights that generalize to other embedding-conditioned models.

2 Methodology

2.1 The Wav2Vec 2.0 Feature Extractor

Wav2Vec 2.0 (Baevski et al., 2020) is a self-supervised speech representation model that operates directly on raw audio sampled at 16 kHz and contains a convolutional feature encoder followed by a 12-layer (base, 95 M parameters) or 24-layer (large, 317 M parameters) transformer context network. The feature encoder reduces the temporal resolution by a factor of 320, producing latent representations at approximately 50 Hz. During pretraining, around 50% of the latent time steps are randomly masked, and the model is optimized with a contrastive loss that requires identifying the correct quantized target from a set of typically 100 negative samples, along with a codebook diversity loss. When fine-tuned with a CTC objective, Wav2Vec 2.0 achieves word error rates below 5% on LibriSpeech test-clean using only 10 minutes of labeled data and below 2% with full supervision, while its intermediate transformer layers yield strong phoneme-level discrimination for tasks such as forced alignment and pronunciation quality assessment.

2.2 Feature-wise Linear Modulation (FiLM).

Feature-wise Linear Modulation (FiLM) (Perez et al., 2018) conditions intermediate neural network representations on an external signal. Given an intermediate activation tensor $\mathbf{h} \in \mathbb{R}^{C \times T}$ and a conditioning variable \mathbf{z} , FiLM applies a channel-

wise affine transformation

$$\text{FiLM}(\mathbf{h}_c) = \gamma_c(\mathbf{z}) \odot \mathbf{h}_c + \beta_c(\mathbf{z}), \quad (1)$$

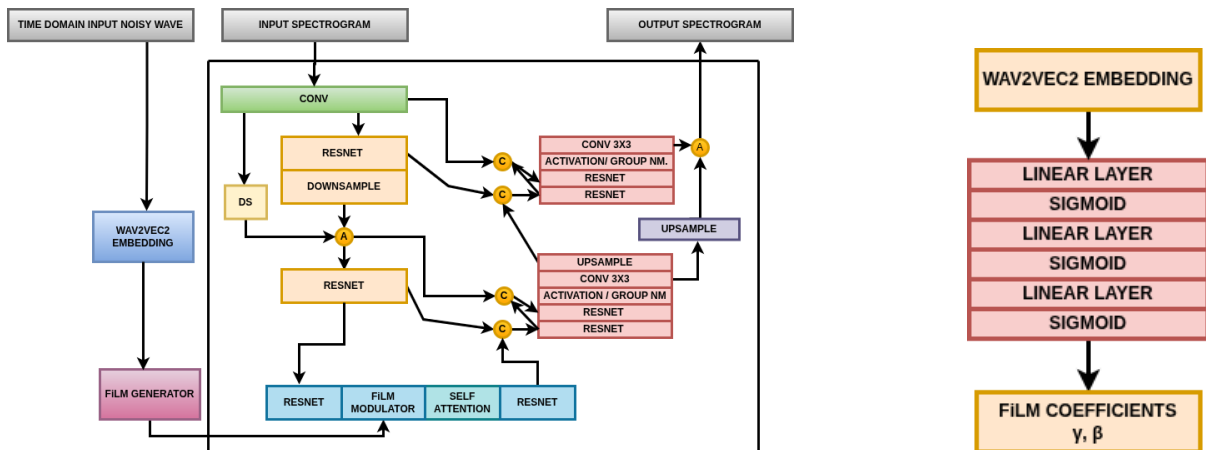
where $c \in \{1, \dots, C\}$ indexes feature channels, \odot denotes element-wise multiplication, and $\gamma_c(\cdot)$ and $\beta_c(\cdot)$ are learnable functions that generate scaling and bias parameters from the conditioning signal. This multiplicative-additive modulation allows the conditioning signal to selectively amplify, suppress, or shift feature activations at different layers, providing a parameter-efficient mechanism for injecting side information without altering the base network architecture.

2.3 The Proposed Model

The backbone of the speech enhancement model that we propose in this work is inspired by the U Net style predictor-generator architecture employed by StoRM (Lemercier et al., 2023). Similar architectures have also been used by image segmentation models to separate the foreground from background. Our design is shown in Figure (a). The model can be broadly divided into three stages, the encoder, the bottleneck and the decoder. On the encoder side, the model cascades a Resnet layer followed by a downsampling stage to form an encoder block. At the end of this block, the block output is added to a downsampled version of the original noisy input spectrogram. The Resnet block consists of a linear layer, group norm/ activation, a dropout of 0.1 and a final 1 x 1 convolution layer to adjust the number of channels.

The bottleneck layer consists of two Resnet layers which are structured as described above, enclosing a standard attention layer on both sides. The bottleneck layer leaves the size of the input tensor unchanged, not contributing to the receptive field size but only using its constituent attention layer to derive context from nearby acoustic features. The attention layer employed by this model is a single head attention which implements the projection operation as a linear layer.

On the decoder side, each block consists of five layers two Resnets, the swish activation function followed by Group Norm, a 3 x 3 convolution layer to reduce the number of channels in the final output to two and finally an upsampler to compensate for the downsampling on the encoder side. As with the encoder end, the final block does not have an upsampling stage. In the version of the model used in this work, we have four blocks on both the



(a) Shared speech enhancement backbone used in this work by both the predictive and generative stages.

(b) FiLM generator module

180 encoder and decoder ends. After the fourth upsampling
 181 block, the output of the Conv 3 x 3 layer in
 182 each block is added together, after appropriate up-
 183 sampling to produce the final estimate of the clean
 184 waveform, to be used in the MSE loss function.
 185 While training, the model truncates a random por-
 186 tion of the training clean and noisy waveforms and
 187 takes a STFT of this portion to create spectrograms
 188 with a pre-specified number of time frames and
 189 frequency bins.

190 2.4 Augmentation with FiLM modulated 191 features

192 In addition to the description in Section 2.3 our
 193 model also consists of an augmentation with a pre-
 194 trained wav2vec 2.0 checkpoint and a FiLM gen-
 195 erator module (Fig 1b). We proceed in two stages.
 196 In the first stage, we use the wav2vec 2.0 on a truncated
 197 chunk of noisy waveform to derive phonetic infor-
 198 mation of the relevant parts. In the second stage,
 199 the wav2vec 2.0 output feature is input to a FiLM
 200 generator module (Fig 1a) to derive the channel wise
 201 coefficients γ_c^t and β_c^t for each 20 ms
 202 time slice t , as described in Section 2.2. Since we
 203 need only one $\tilde{\gamma}_c$ and $\tilde{\beta}_c$ per channel we perform
 204 exponential averaging along time as follows:

$$205 \tilde{\gamma}_c^t = \alpha \gamma_c^t + (1 - \alpha) \tilde{\gamma}_c^{t-1}, \quad (2)$$

$$206 \tilde{\beta}_c^t = \alpha \beta_c^t + (1 - \alpha) \tilde{\beta}_c^{t-1}, \quad 0 < \alpha \leq 1 \quad (3)$$

207 with initialization $\tilde{\gamma}_c^0 = \gamma_c^0$ and $\tilde{\beta}_c^0 = \beta_c^0$. The
 208 final values of $\tilde{\gamma}_c$ and $\tilde{\beta}_c$ are given by $\tilde{\gamma}_c^T$ and $\tilde{\beta}_c^T$,
 209 respectively, with T being the total number of time
 210 slices in the wav2vec 2.0 output. Subsequently,
 211 $\tilde{\gamma}_c$ and $\tilde{\beta}_c$ are used to modulate the input h to the

212 attention layer in the U-Net bottleneck as $\tilde{\gamma}_c \odot h +$
 213 $\tilde{\beta}_c$ (Fig 1a).

214 The key idea behind modulating with wav2vec
 215 2.0 features along the channel dimension is to treat
 216 phonetic information as an additional view of the
 217 acoustic signal, parallel to the spectrogram. In a
 218 standard spectrogram-based U-Net, the network re-
 219 ceives a 4-dimensional tensor (B, C, F, T) consist-
 220 ing of a batch dimension, a channel dimension, and
 221 two spatial dimensions corresponding to frequency
 222 bins and time frames. By modulating the bottle-
 223 neck feature channels with phonetic information,
 224 we give the model an augmented input representa-
 225 tion where each time frame is annotated not only
 226 with its acoustic energy distribution (from STFT)
 227 but also with a phonetic likelihood distribution.

228 An important property of this approach is that
 229 it remains strictly local in time, just like stan-
 230 dard masking-based enhancement. The modula-
 231 tion scheme does not change the architecture in a
 232 disruptive way. They simply provide a richer repre-
 233 sentation of each time frame, grounded in linguistic
 234 structure, which allows the network to operate with
 235 a higher-level prior about which frequency patterns
 236 are speech-like versus noise-like. This is especially
 237 useful in challenging noise conditions or speech
 238 distortions, where the raw spectrogram alone does
 239 not reliably indicate what should be preserved.

240 The method also leverages the fact that phonetic
 241 states change more slowly than raw spectral coeffi-
 242 cients. Wav2vec 2.0 representation carry informa-
 243 tion about phoneme identity, which persists across
 244 tens of milliseconds, while short-term spectral en-
 245 ergy fluctuates rapidly.

246 A further practical benefit is that this fusion does
 247 not require architectural changes to the U-Net's en-

SNR	Ours (128 ch.)		StoRM (128 ch.)		Ours (32 ch.)		StoRM (32 ch.)	
	STOI	PESQ	STOI	PESQ	STOI	PESQ	STOI	PESQ
-3 dB	0.8757	1.7874	0.7258	1.2611	0.7821	1.3743	0.8090	1.3210
0 dB	0.9108	1.9972	0.7870	1.3178	0.8390	1.5557	0.8626	1.4696
3 dB	0.9360	2.2059	0.8464	1.4218	0.8820	1.7714	0.9017	1.6457
6 dB	0.9553	2.3958	0.8885	1.4984	0.9146	1.9854	0.9306	1.8348
12 dB	0.9778	2.7690	0.9506	1.8613	0.9571	2.4080	0.9660	2.2392

Table 1: PESQ and STOI scores for different configurations of the proposed and baseline models.

248 coder–decoder or skip connections. As long as the
249 phoneme information is temporally aligned with
250 the STFT frames (and broadcast across frequency),
251 they plug directly into the model. This makes
252 the approach lightweight, architecture-neutral, and
253 easy to scale. Different linguistic or articulatory
254 side features (e.g., phoneme embeddings, bottle-
255 neck features, articulatory features) can be swapped
256 in without redesigning the backbone.

257 3 Dataset

258 To train the speech enhancement models, 19,992
259 audio waveforms were taken from the train split
260 of Librispeech *train-clean-100* (Panayotov et al.,
261 2015). These files were divided equally among 27
262 noises, which ranged from everyday noises such
263 as baby cries to ambient noises from military set-
264 tings such as helicopters. For each noise type, an
265 equal number of waveforms were randomly cho-
266 sen to be mixed with it at 5 different SNRs: -3dB,
267 0dB, 3dB, 6dB and 12dB. These SNRs were se-
268 lected to model the wide range of scenarios that a
269 general-purpose speech enhancement system might
270 encounter. Owing to the relatively large inference
271 time for diffusion models such as StoRM, we kept
272 our validation set to a minimum. The validation set
273 consisted of 10 waveforms from the Librispeech
274 validation set. The set was divided into 5 groups of
275 2 each for one SNR. The noises used for each set
276 were from the office and crowd domains, as back-
277 ground chatter is considered to be one of the most
278 challenging in the domain of speech enhancement
279 due to the irregular structure of these noises. For
280 the test set, we used 8 noises from the training set
281 and 2 unseen noises. Each noise was combined
282 with 20 utterances from the Librispeech test set
283 at the same five SNRs. Thus, the total number of
284 utterances in the test set was 1000, totaling around
285 4 hours in duration. Model checkpoints, noisy and
286 enhanced test set can be found [here](#).

287 4 Experiments

288 In this section, we evaluate the effectiveness of
289 phoneme-based conditioning by comparing the pro-
290 posed model against a baseline enhancement sys-
291 tem that does not incorporate phonetic cues. In
292 particular, we analyze both (i) the perceptual qual-
293 ity of the reconstructed speech (PESQ scores) and
294 (ii) the intelligibility by measuring how well the en-
295 hanced signal preserves short-time temporal envel-
296 ope (STOI). For the smoothing scheme in Section
297 2.4, we experimented with various values of the
298 hyperparameter α and found that $\alpha = 0.1$ performs
299 best. We experiment with different widths of the
300 proposed model, defining width as the number of
301 output channels of the green convolution layer in
302 Figure 1a. The PESQ and STOI scores obtained for
303 different channel counts and SNR are summarized
304 in Table 1.

305 Further, we experiment with a different aver-
306 aging scheme (naive averaging) and present the
307 enhancement results in Table 2 in Appendix A. We
308 also train and test on the VoiceBank+DEMAND
309 dataset (Valentini-Botinhao et al., 2016). The de-
310 tails and results of this experiment can be found in
311 Appendix B.

312 5 Conclusions

313 In this work, we proposed a diffusion-based speech
314 enhancement framework augmented with phonetic
315 representations derived from wav2vec 2.0. We in-
316 vestigated the impact of different channel configu-
317 rations as well as multiple strategies for fusing pho-
318 netic information within the enhancement pipeline.
319 Experiments conducted on both a standard and a
320 non-standard dataset demonstrate the effectiveness
321 of the proposed approach across diverse acoustic
322 conditions. We evaluate performance using objec-
323 tive metrics, including PESQ for perceptual quality
324 and STOI for intelligibility, and observe consistent
325 improvements over baseline configurations.

326 Limitations

327 While our proposed speech enhancement system
328 demonstrates improved robustness and perceptual
329 quality, it has several limitations. First, the use of
330 a diffusion-based architecture results in high com-
331 putational cost and longer inference times com-
332 pared to conventional feedforward models, which
333 may limit real-time deployment. Second, condition-
334 ing on Wav2Vec 2.0 features increases the overall
335 parameter count, even though these features are
336 frozen during training. Reducing the number of
337 channels can mitigate the parameter and memory
338 footprint; however, this also leads to a measurable
339 degradation in the quality of the enhanced wave-
340 form. Third, our experiments are conducted on a
341 non-standard dataset, which may limit compar-
342 ability with other methods. We plan to extend this
343 study to publicly available speech enhancement
344 datasets such as DNS Challenge, and LibriSpeech-
345 derived noisy corpora to better evaluate general-
346 ization and reproducibility. Finally, as with most
347 supervised speech enhancement methods, perfor-
348 mance may vary across unseen noise types, speak-
349 ers, and recording conditions.

350 Ethics Statement

351 This work studies speech enhancement methods
352 aimed at improving speech quality and intelligibil-
353 ity in noisy conditions. The proposed approach
354 does not involve the collection of new human
355 subject data. Speech enhancement systems may
356 present ethical risks if performance varies across
357 speakers, accents, languages, or recording condi-
358 tions. To mitigate potential bias, we train our model
359 on diverse speakers and noise conditions where pos-
360 sible, and we report aggregate performance rather
361 than claims about individual speakers. Nonethe-
362 less, residual biases may remain, particularly for
363 underrepresented accents or speaking styles, and
364 should be considered when deploying such systems
365 in real-world applications. The proposed method
366 is not intended for surveillance, speaker identifica-
367 tion, or forensic use. While enhanced speech could
368 potentially be misused to improve the intelligibil-
369 ity of recordings without a speaker’s consent, this
370 risk is inherent to speech enhancement technolo-
371 gies broadly and not specific to our approach. We
372 encourage responsible use consistent with privacy
373 laws and ethical guidelines.

References 374

- Alex Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. [wav2vec 2.0: A framework for self-supervised learning of speech representations](#). In *Advances in Neural Information Processing Systems (NeurIPS)* 33. 375-379
- Jean-Marie Lemerrier, Julius Richter, Simon Welker, and Timo Gerkmann. 2023. [Storm: A diffusion-based stochastic regeneration model for speech enhancement and dereverberation](#). *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31:2724–2737. 380-385
- Y.-J. Lu, C.-F. Liao, X. Lu, J.-w. Hung, and Y. Tsao. 2020. [Incorporating broad phonetic information for speech enhancement](#). In *Proc. Interspeech*, pages 2417–2421. 386-389
- Yen-Ju Lu, Chia-Yu Chang, Cheng Yu, Ching-Feng Liu, Jieih-weih Hung, Shinji Watanabe, and Yu Tsao. 2023. [Improving speech enhancement performance by leveraging contextual broad phonetic class information](#). *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31:2738–2750. 390-395
- Khandokar Md. Nayem and Donald S. Williamson. 2021. [Towards an asr approach using acoustic and language models for speech enhancement](#). In *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7123–7127. 396-401
- Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. [Librispeech: An asr corpus based on public domain audio books](#). In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5206–5210. 402-406
- Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron C. Courville. 2018. [Film: Visual reasoning with a general conditioning layer](#). In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*. 407-411
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine Mcleavey, and Ilya Sutskever. 2023. [Robust speech recognition via large-scale weak supervision](#). In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 28492–28518. PMLR. 412-418
- C. Valentini-Botinhao, X. Wang, S. Takaki, and J. Yamagishi. 2016. [Speech enhancement for a noise-robust text-to-speech synthesis system using deep recurrent neural networks](#). In *Proc. Interspeech*. 419-422
- Wei Wang, Wangyou Zhang, Shaoxiong Lin, and Yanmin Qian. 2022. [Text-informed knowledge distillation for robust speech enhancement and recognition](#). In *2022 13th International Symposium on Chinese Spoken Language Processing (ISCSLP)*, pages 334–338. 423-428

SNR	Ours (128 ch.)		StoRM (128 ch.)		Ours (32 ch.)		StoRM (32 ch.)	
	Naive	Exp.	Naive	Exp.	Naive	Exp.	Naive	Exp.
-3 dB	1.4351	1.7874	1.2611	-	1.2062	1.3743	1.3210	-
0 dB	1.6087	1.9972	1.3178	-	1.3229	1.5557	1.4696	-
3 dB	1.8268	2.2059	1.4218	-	1.4938	1.7714	1.6457	-
6 dB	2.0365	2.3958	1.4984	-	1.6824	1.9854	1.8348	-
12 dB	2.4065	2.7690	1.8613	-	2.0755	2.4080	2.2392	-

Table 2: PESQ scores for exponential smoothing and naive averaging for different configurations of the proposed and baseline models. Since baseline StoRM has no averaging we report it’s score under the "Naive" column.

	Ours (128 ch.)	StoRM (128 ch.)	Ours (32 ch.)	StoRM (32 ch.)
PESQ	2.06	2.41	2.29	2.49
STOI	0.84	0.84	0.82	0.82

Table 3: PESQ and STOI scores for VoiceBank+DEMAND noisy test set.

Yufeng Yang et al. 2023. [Time-domain speech enhancement for robust automatic speech recognition](#). In *Proceedings of Interspeech*, pages 4913–4917.

A Exponential Smoothing v/s Naive Averaging

In this section, we replace the running average scheme described in Section 2.4 with a naive averaging scheme. More specifically, we have

$$\begin{aligned}\tilde{\gamma}_c &= \frac{1}{T} \sum_{t=1}^T \gamma_c^t \\ \tilde{\beta}_c &= \frac{1}{T} \sum_{t=1}^T \beta_c^t\end{aligned}\quad (4)$$

A distinct disadvantage of this scheme is that it cannot be pipelined unlike the exponential averaging discussed earlier. Hence it precludes real time usage of the enhancement model since we require the knowledge of γ_c^t, β_c^t for all t to calculate $\tilde{\gamma}_c$ and $\tilde{\beta}_c$ for modulation. Nevertheless, we enhance using a trained checkpoint from naive averaging and compare PESQ scores with exponential smoothing in Table 2. It is evident that exponential smoothing provides superior enhancement performance across all SNRs and channel configurations.

B Enhancement Performance on VoiceBank+DEMAND dataset

In the interest of reproducibility and comparison with other well known models we train the proposed model on VoiceBank+DEMAND 28 speaker training dataset (Valentini-Botinhao et al., 2016)

and enhance the noisy test set. The PESQ and STOI scores are summarized in Table 3. We down sample the training and test splits of the dataset to 16kHz from 48kHz as the pretrained wav2vec 2.0 checkpoint we are using in our model cannot process 48kHz waveforms. We trained and inferred using both 128 and 32 channel variants of our model and compare it’s enhancement capabilities with the baseline model.

Although the baseline model beats us on PESQ, we are equally good on STOI. This is in contrast with the performance trends observed with our custom dataset which we attribute to the different noise types and SNR levels used in preparing this dataset. Nevertheless, we observe an interesting phenomenon from this experiment. It is evident from Table 3 that a larger model does not guarantee superior enhancement performance. The 32 channel variant of the proposed and baseline models beat the 128 channel variant, respectively, on PESQ scores.