

Incorporating Error Level Noise Embedding for Improving LLM-Assisted Robustness in Persian Speech Recognition

Anonymous ACL submission

Abstract

Automatic Speech Recognition (ASR) performance degrades severely under noise, a challenge that is particularly pronounced for low-resource languages such as Persian. Even state-of-the-art systems like Whisper exhibit substantial accuracy loss at low signal-to-noise ratios. We propose a noise-aware ASR error correction framework that combines multiple transcription hypotheses with explicit modeling of linguistic noise. From noisy Persian speech, we generate 5-best hypotheses using Whisper and introduce Error Level Noise (ELN), a representation that captures sentence- and token-level disagreement across hypotheses as a proxy for noise-induced uncertainty. ELN vectors are used to condition a fine-tuned LLaMA-2-7B model during post-hoc correction. Experiments show that ELN conditioning significantly reduces Word Error Rate (WER). On the Mixed Noise test set, our model lowers WER from 31.10% (Raw Whisper) to 24.84%, outperforming a text-only fine-tuned baseline (30.79%), while a zero-shot LLaMA-2 model fails to correct Persian ASR outputs. These results demonstrate the effectiveness of noise-aware multi-hypothesis correction for robust Persian ASR.

1 Introduction

ASR has become an integral part in the interaction of humans with computers in the modern world, and it finds applications in smart assistants, speech translation, and automatic captioning. While recent advances in end-to-end models and large language models (LLMs) have substantially improved ASR performance for high-resource languages, robustness under noisy conditions and for low-resource languages remains a major challenge. In particular, ASR systems for Persian still suffer significant degradation when exposed to environmental noise.

Previous studies on the robustness of automatic speech recognition (ASR) systems have primarily been centered around noise-aware training, multi-

condition training, and data augmentation using synthetic noise (Ko et al., 2015; Watanabe et al., 2018; Li et al., 2016). While such methods have proven to be quite effective, they mainly work at the acoustic modeling level and usually involve the retraining or modification of the ASR system. Moreover, most existing methods are developed for high-resource languages and rely on single best hypotheses, ignoring uncertainty induced by noise. Post-processing approaches based on textual correction also generally lack explicit noise modeling and rarely leverage multi-hypothesis reasoning with large generative models.

In this work, we introduce a noise-robust, post-hoc generative error correction framework for Persian ASR that avoids retraining the underlying recognizer. We generate 5-best hypotheses from a Whisper-based Persian ASR system and propose *Error Level Noise* (ELN) embeddings, which capture sentence- and token-level disagreement across hypotheses as a proxy for noise-induced uncertainty. These embeddings act as inputs to a LLaMA-2-7B model (Touvron et al., 2023) that has been fine-tuned, allowing the noise-aware correction to be informed not only by the linguistic context but also by the variability of the hypothesis. Noisy Persian speech experiments derived from Common Voice (Ardila et al., 2019) and MUSAAN (Snyder et al., 2015) show that the WER of the proposed approach is considerably decreased in comparison to both the ASR baseline and text-only correction models, thus it is a noise-robust ASR solution that is effective and scalable in low-resource scenarios.

2 Related Work

Automatic speech recognition (ASR) error correction has been approached from various angles, including rule-based methods, statistical language models, and more recently, neural and large lan-

082 guage model (LLM)-based methods, with continu- 132
083 ous difficulties in low-resource and noisy environ- 133
084 ments.

085 Retrieval-augmented generation (RAG) has been 134
086 used in ASR correction by relying on external ex- 135
087 amples. GEC-RAG (Robotian et al., 2025) helps 136
088 Persian ASR by fetching closely matching error pat- 137
089 terns for LLM-based correction. On the other hand, 138
090 LA-RAG (Li et al., 2024a) and DARAG (Ghosh 139
091 and et al., 2025) mix retrieval with generative 140
092 modeling to increase resilience and generalization 141
093 across different domains and acoustic scenarios. 142

094 Some of the works are dedicated entirely to 143
095 the Persian ASR and spelling correction. PER- 144
096 CORE (Dashti et al., 2024) brings phonetic features 145
097 to the table for tackling real-word and non-word 146
098 spelling errors, whereas PSRB (Sedghiyeh et al., 147
099 2025) presents a benchmark dataset for assessing 148
100 Persian ASR in varying linguistic and acoustic con- 149
101 ditions. 150

102 Generative Error Correction (GER) (Ma et al., 151
103 2023) sees ASR post-processing as a sequence-to- 152
104 sequence rewriting task and has spawned devel- 153
105 opments like Denoising GER (Liu et al., 2025) 154
106 and RobustGER (Hu et al., 2024b), which in- 155
107 crease noise robustness through language-space 156
108 modeling and knowledge distillation respectively. 157
109 ClozeGER (Hu et al., 2024a) also takes acoustic 158
110 information into account by changing the correc- 159
111 tion problem into a cloze task with speech logits. 160
112 HyParadise (Chen et al., 2023) showcases the capa- 161
113 bility of LLMs for generative speech recognition, 162
114 and multilingual LLMs have been adjusted as well 163
115 for direct 1-best ASR correction without N-best 164
116 hypotheses, exhibiting cross-lingual transfer advan- 165
117 tages for low-resource languages (Li et al., 2024b). 166

118 Differently from previous research, our method 167
119 utilizes disagreement among several ASR hypothe- 168
120 ses to represent noise-induced uncertainty explic- 169
121 itly and uses an LLM to generate structured text- 170
122 level noise representations without the need for 171
123 acoustic inputs or changes to the ASR system. 172

124 3 Methodology 173

125 We propose **GER + Text Denoising (Ours)**, a 174
126 noise-aware ASR error correction framework that 175
127 conditions a large language model (LLM) on struc- 176
128 tured text-level noise representations. The method 177
129 targets low-resource Persian ASR and operates as 178
130 a post-hoc correction module without modifying or 179
131 retraining the underlying recognizer. As shown in 180

Figure 1 (Panel d), the framework uses the top-5 132
ASR hypotheses (*5-best list*) as input. 133

Our approach is inspired by RobustGER (Hu 134
et al., 2024b) but is designed to be computationally 135
lightweight. Unlike RobustGER, which relies 136
on knowledge distillation and optional audio noise 137
modeling, we focus exclusively on *text noise*. We 138
introduce *Error Level Noise* (ELN) vectors that cap- 139
ture sentence- and token-level disagreement across 140
hypotheses and use them as conditioning signals 141
for the LLM. This enables noise-aware correction 142
directly in the language space, without audio fea- 143
tures or mutual information estimation. 144

145 3.1 System Overview 146

Figure 1 compares four architectures used for 146
ASR noise Robustness. Panels (a) through (c) are 147
adapted from (Hu et al., 2024b): (a) GER (Chen 148
et al., 2023; Yang et al., 2023), a text-only genera- 149
tive correction model; (b) GER with audio denois- 150
ing (Liu et al., 2025); (c) RobustGER (Hu et al., 151
2024b), which models language noise via knowl- 152
edge distillation; and (d) our proposed method, 153
which injects text-level noise embeddings into the 154
LLM to improve robustness across noise condi- 155
tions. 156

157 3.2 Hypothesis Generation 158

Given noisy speech, we generate an N-best list us- 158
ing the Whisper large ASR model (Radford and 159
et al., 2023), employing a fine-tuned Persian vari- 160
ant¹. Beam search decoding produces up to five 161
unique hypotheses per utterance; samples with 162
fewer hypotheses are padded to $n = 5$. 163

All hypotheses and references undergo Persian 164
text normalization, including Unicode standardiza- 165
tion, digit conversion, spacing correction, punctua- 166
tion removal, and whitespace cleanup. This prepro- 167
cessing ensures consistent evaluation using Word 168
Error Rate (WER) and supports reliable computa- 169
tion of ELN vectors. 170

171 3.3 ELN Vector Extraction 172

For an utterance, let $\mathcal{H} = \{H_1, \dots, H_n\}$ denote 172
the $n = 5$ ASR hypotheses. ELN consists of 173
sentence-level and token-level components, which 174
together quantify linguistic disagreement across 175
hypotheses and serve as text-level noise representa- 176
tions. 177

¹<https://huggingface.co/vhdm/whisper-large-fa-v1>

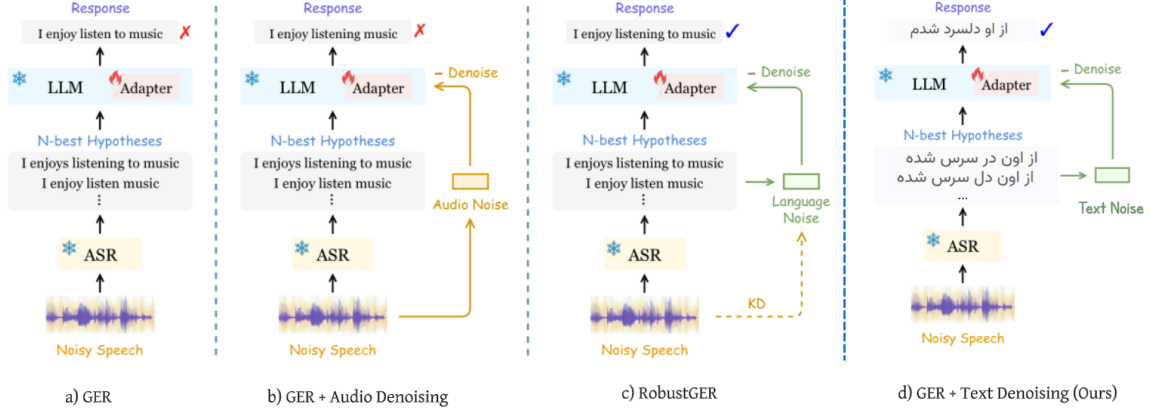


Figure 1: Overview of a) GER (Chen et al., 2023) b) GER + Audio Denoising (Liu et al., 2025) c) RobustGER (Hu et al., 2024b) and d) GER + Text Denoising (Ours).

Sentence-level ELN Each hypothesis is embedded using a sentence encoder (e.g., Sentence-BERT (Reimers and Gurevych, 2019)):

$$\mathbf{e}_i = \text{Embed}_{\text{sent}}(H_i) \in R^d.$$

The sentence-level ELN vector is computed as:

$$\mathbf{v}_{\text{sent}} = \frac{2}{n(n-1)} \sum_{i=1}^n \sum_{j=i+1}^n (\mathbf{e}_i - \mathbf{e}_j)^2$$

capturing semantic variance among hypotheses.

Token-level ELN Each hypothesis is tokenized: $H_i = [t_{i,1}, t_{i,2}, \dots, t_{i,L_i}]$, and padded to the maximum length $L_{\text{max}} = \max_i L_i$. Tokens are embedded into vectors of dimension d' :

$$\mathbf{t}_{i,k} = \text{Embed}_{\text{tok}}(t_{i,k}) \in R^{d'}.$$

Token-level ELN is defined as:

$$\mathbf{v}_{\text{tok}} = \frac{1}{L_{\text{max}}} \sum_{k=1}^{L_{\text{max}}} \frac{2}{n(n-1)} \sum_{i=1}^n \sum_{j=i+1}^n (\mathbf{t}_{i,k} - \mathbf{t}_{j,k})^2$$

Final ELN Representation The final ELN vector concatenates both components:

$$\mathbf{v}_{\text{ELN}} = \mathbf{v}_{\text{sent}} \parallel \mathbf{v}_{\text{tok}},$$

and is injected into the LLM via an adapter-based conditioning mechanism (Zhang et al., 2023).

ELN Magnitude We define ELN magnitude as the L_2 norm:

$$\|\mathbf{v}_{\text{ELN}}\|_2,$$

which provides a scalar measure of noise-induced disagreement, with larger values indicating greater transcription uncertainty.

4 Experimental Setup

We evaluate the proposed method using a unified prompt that instructs the large language model (LLM) to perform ASR error correction based on the top-5 hypotheses generated by the recognizer and considering the noise represented by ELN vector, while preserving meaning and avoiding synonym substitution.

4.1 Datasets

Common Voice 16.1 (Persian) (Ardila et al., 2019) was the source of clean speech in the creation of noisy Persian speech, whereas the noise corpus was from MUSAN (Snyder et al., 2015). The Persian part of Common Voice has nearly 90 hours of transcribed speech (40k utterances) recorded at 16 kHz with various speakers and accents. When it comes to noise, we have taken the ambient noise part of MUSAN which is a collection of around 6 hours of environmental sounds as well as artificial noises.

Noise Augmentation Noisy utterances were generated by adding MUSAN noise and low-amplitude Gaussian noise to clean speech: MUSAN noise was added at SNRs uniformly sampled between 0–15 dB, and Gaussian noise amplitudes were drawn from $[0.001, 0.015]$, yielding a wide range of realistic acoustic conditions.

4.2 Evaluation Metrics and Noise Conditions

Performance is measured using Word Error Rate (WER).

We report results under four conditions: **Clean** (5000 samples), **Mixed Noise** (random MUSAN noise at 0–15 dB, 5000 samples), **SNR = 5 dB**

Method	Clean	Mixed Noise	SNR = 5 dB	SNR = 10 dB
Raw Whisper	24.80	31.10	42.70	38.30
Base Model (Zero-shot)	62.43	64.58	70.63	67.75
Fine-tuned (No ELN)	24.06	30.79	39.76	31.59
Fine-tuned + ELN (Ours)	24.39	24.84	32.34	28.02

Table 1: WER (%) comparison of Raw Whisper, Base LLaMA2 (zero-shot), LLaMA2 fine-tuned without ELN, and the ELN-conditioned model across four acoustic conditions.

(random MUSAN noise at 5 dB, 1000 samples), and **SNR = 10 dB** (random MUSAN noise at 10 dB, 1000 samples). All noisy settings combine MUSAN ambient noise with low-level Gaussian noise.

4.3 Compared Models

We compare three configurations:

1. Zero-shot: LLaMA-2-7B without fine-tuning.
2. Fine-tuned (Text-only): LLaMA-2-7B fine-tuned with LoRA (Hu et al., 2022) on 5-best hypotheses.
3. Fine-tuned + ELN (Ours): LLaMA-2-7B fine-tuned with LoRA incorporating ELN vectors. ELN vectors were mapped through a small MLP to match the LLM embedding dimension and prepended as prefix embeddings.

4.4 Training Details

All models were fine-tuned using 4-bit quantization, a learning rate of 2×10^{-4} for 3 epochs, gradient accumulation and checkpointing, and a cosine learning-rate scheduler with weight decay.

5 Results and Discussion

5.1 Quantitative Results

Table 1 reports the Word Error Rate (WER, %) of all evaluated models under four noise conditions: Clean, Mixed Noise (covering 0–15 dB SNR plus Gaussian noise), SNR = 5 dB, and SNR = 10 dB. For reference, the Raw Whisper baseline follows Radford et al. (Radford and et al., 2023), while the Base Model (Zero-shot) corresponds to LLaMA-2 (Touvron et al., 2023). All WER values are expressed in percentage.

5.2 Comparison with RobustGER on VB-DEMAND

We further evaluate noise robustness and cross-lingual generalization on the VoiceBank-DEMAND (VB-DEMAND) benchmark (Veaux

Method	Baseline WER	Improved WER
RobustGER	13.00	10.70
Ours (Fine-tuned + ELN)	7.93	3.96

Table 2: Comparison of WER (%) between the proposed ELN-conditioned model and RobustGER (Hu et al., 2024b) on the VB-DEMAND dataset.

et al., 2017; Thiemann et al., 2013) using the 16 kHz version. Whisper small.en is used to generate N-best hypotheses, and the ELN-conditioned LLaMA2-7B model is trained on the official training split and evaluated on the test set.

As shown in Table 2, our method reduces WER from 7.93% to 3.96%, substantially outperforming RobustGER (Hu et al., 2024b) (10.70%). This demonstrates that ELN-based text denoising generalizes effectively beyond Persian and provides strong robustness under real-world noise.

For clarity, in Table 2, Baseline WER refers to the Word Error Rate of the Whisper ASR system prior to any correction, while Improved WER indicates the WER after applying the respective error correction model.

5.3 Key Findings

Our results show that: (i) task-specific fine-tuning is essential for Persian ASR correction; (ii) ELN conditioning significantly reduces WER, particularly under moderate and severe noise; and (iii) ELN magnitude strongly correlates with recognition difficulty, validating its role as a noise-aware signal for generative ASR error correction.

6 Conclusion

We explored using large language models (LLMs) for post-hoc correction of ASR outputs, focusing on low-resource Persian. We proposed a computationally efficient framework that conditions a fine-tuned LLaMA-2-7B on an Error Level Noise (ELN) vector derived from n-best ASR hypotheses.

By explicitly modeling linguistic noise via ELN, our approach improves transcription robustness without modifying the original ASR. Conditioning LLMs with structured noise thus offers a practical, scalable way to enhance ASR reliability in low-resource, noisy settings.

Limitations

Although the results were encouraging, a number of the study’s limitations should be considered. First, the experiments used the Common Voice (Persian)

313	dataset which was most probably augmented with		
314	synthetic MUSAN noise. This might not be a per-		
315	fect representative of the diversity and the com-		
316	plexity of the real Persian speech, especially when		
317	it comes to conversational, dialectal, and domain-		
318	specific contexts. Second, the evaluation mainly		
319	depended on the Word Error Rate (WER) metric.		
320	Of course, WER is a valuable quantitative tool, but		
321	the next research should involve human-based as-		
322	sessments or semantic similarity metrics to more		
323	deeply evaluate changes in fluency, coherence, and		
324	meaning preservation. Lastly, the current method		
325	only models noise through textual proxies with the		
326	help of Error Level Noise (ELN) vectors, without		
327	any direct acoustic evidence. It will also be very		
328	useful to add explicit acoustic features to the model		
329	that will enable a complete understanding of the		
330	influence of noise and thus better robustness of the		
331	language space noise model could be achieved.		
332	References		
333	Rosana Ardila, Megan Branson, Kelly Davis, Michael		
334	Henretty, Michael Kohler, Josh Meyer, Reuben		
335	Morais, Lindsay Saunders, Francis M. Tyers,		
336	and Gregor Weber. 2019. Common voice: A		
337	massively-multilingual speech corpus . <i>CoRR</i> ,		
338	abs/1912.06670.		
339	Chen Chen, Yuchen Hu, Chao-Han Huck Yang,		
340	Sabato Marco Siniscalchi, Pin-Yu Chen, and		
341	Eng Siong Chng. 2023. Hyporadise: An open base-		
342	line for generative speech recognition with large lan-		
343	guage models. In <i>Proceedings of the Thirty-seventh</i>		
344	<i>Conference on Neural Information Processing Sys-</i>		
345	<i>tems (Datasets and Benchmarks Track)</i> .		
346	Seyed Mohammad Sadegh Dashti, Amid Khatibi Bard-		
347	siri, and Mehdi Jafari Shahbazzadeh. 2024. Percore:		
348	A deep learning-based framework for persian spelling		
349	correction with phonetic analysis. <i>International Jour-</i>		
350	<i>nal of Computational Linguistics and Applications</i> ,		
351	17(1):1–20.		
352	Suman Ghosh and et al. 2025. Darag: Improving gener-		
353	ative error correction for asr with synthetic data and		
354	retrieval-augmented generation. <i>Findings of the As-</i>		
355	<i>sociation for Computational Linguistics</i> , 3:125–134.		
356	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan		
357	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and		
358	Weizhu Chen. 2022. Lora: Low-rank adaptation of		
359	large language models. In <i>Proceedings of the Inter-</i>		
360	<i>national Conference on Machine Learning (ICML)</i> .		
361	Yuchen Hu, Chen Chen, Chengwei Qin, Qiushi Zhu,		
362	Eng Siong Chng, and Ruizhe Li. 2024a. Listen again		
363	and choose the right answer: A new paradigm for		
364	automatic speech recognition with large language		
	models . In <i>Findings of the Association for Computa-</i>		
	<i>tional Linguistics: ACL 2024</i> . 365		366
	Yuchen Hu, Chen Chen, Chao-Han Huck Yang, Ruizhe		
	Li, Chao Zhang, Pin-Yu Chen, and Eng Siong Chng.		
	2024b. Large language models are efficient learn-		
	ers of noise-robust speech recognition. In <i>Inter-</i>		
	<i>national Conference on Learning Representations</i>		
	<i>(ICLR) 2024</i> . 370		371
	Tom Ko, Vijayaditya Peddinti, Daniel Povey, and San-		
	jeev Khudanpur. 2015. Audio augmentation for		
	speech recognition . In <i>Proceedings of Interspeech</i> ,		
	pages 3586–3589. 374		375
	Qiang Li, Yan He, and Li Deng. 2016. <i>Robust Auto-</i>		
	<i>matic Speech Recognition: A Bridge to Practical</i>		
	<i>Applications</i> . Academic Press. 377		378
	Shaojun Li, Hengchao Shang, Daimeng Wei, Jiaxin		
	Guo, Zongyao Li, Xianghui He, Min Zhang, and		
	Hao Yang. 2024a. La-rag: Enhancing llm-based asr		
	accuracy with retrieval-augmented generation. <i>arXiv</i>		
	<i>preprint arXiv:2409.08597</i> . 380		381
	Sheng Li, Chen Chen, Chin Yuen Kwok, Chenhui Chu,		
	Eng Siong Chng, and Hisashi Kawai. 2024b. In-		
	vestigating asr error correction with large language		
	model and multilingual 1-best hypotheses . In <i>Proc.</i>		
	<i>Interspeech</i> , Kos, Greece. 385		386
	Yanyan Liu, Minqiang Xu, Yihao Chen, Liang He, Lei		
	Fang, Sian Fang, and Lin Liu. 2025. Denoising ger:		
	A noise-robust generative error correction with llm		
	for speech recognition . <i>Preprint</i> , arXiv:2509.04392. 390		391
	Rao Ma, Mengjie Qian, Potsawee Manakul, Mark J. F.		
	Gales, and Kate M. Knill. 2023. Can generative		
	large language models perform asr error correction?		
	<i>CoRR</i> , abs/2307.04172. 394		395
	Alec Radford and et al. 2023. Robust speech		
	recognition via whisper. In <i>arXiv preprint</i>		
	<i>arXiv:2212.04356</i> . 398		399
	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert:		
	Sentence embeddings using siamese bert-networks.		
	In <i>Proceedings of the 2019 Conference on Empirical</i>		
	<i>Methods in Natural Language Processing</i> . Associa-		
	tion for Computational Linguistics. 401		402
	Amin Robatian, Mohammad Hajipour, Moham-		
	mad Reza Peyghan, Fatemeh Rajabi, Sajjad Amini,		
	Shahrokh Ghaemmaghami, and Iman Gholampour.		
	2025. Gec-rag: Improving generative error cor-		
	rection via retrieval-augmented generation for au-		
	tomatic speech recognition systems. <i>arXiv preprint</i>		
	<i>arXiv:2501.10734</i> . 406		407
	Nima Sedghiyeh, Sara Sadeghi, Reza Khodadadi, Farzin		
	Kashani, Omid Aghdaei, Somayeh Rahimi, and Mo-		
	hammad Sadegh Safari. 2025. Psrb: A comprehen-		
	sive benchmark for evaluating persian asr systems.		
	<i>arXiv preprint arXiv:2505.21230</i> . 413		414

- 418 David Snyder, Guoguo Chen, and Daniel Povey. 2015.
419 [Musan: A music, speech, and noise corpus](#). *CoRR*,
420 abs/1510.08484.
- 421 Joachim Thiemann, Nobutaka Ito, and Emmanuel Vincent. 2013. The DEMAND dataset of realistic noise recordings for speech processing. In *Proceedings of the 21st International Conference on Acoustics (IWAENC)*, pages 1–4. IEEE.
- 426 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 others. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- 434 Christophe Veaux, Junichi Yamagishi, and Kirsten Mac-Donald. 2017. Cstr vctk corpus: English multi-speaker corpus for cstr voice cloning toolkit. Version 0.92.
- 438 Shinji Watanabe, Takaaki Hori, Shigeki Karita, Tomoki Hayashi, Jiro Nishitoba, Yuya Unno, Nelson Enrique Yalta Soplín, Jahn Heymann, Matthew Wiesner, Nanxin Chen, Adithya Renduchintala, and Tsubasa Ochiai. 2018. [ESPnet: End-to-end speech processing toolkit](#). In *Proceedings of Interspeech*, pages 2207–2211.
- 445 Chao-Han Huck Yang, Yile Gu, Yi-Chieh Liu, Shalini Ghosh, Ivan Bulyko, and Andreas Stolcke. 2023. [Generative speech recognition error correction with large language models and task-activating prompting](#). In *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 1–8. IEEE.
- 451 Rui Zhang, Songlin Yang, Chenguang Zhu, and Yue Zhang. 2023. [LLM-Adapters: Parameter-efficient fine-tuning of large language models with adapter layers](#). *arXiv preprint arXiv:2304.01933*.
- 452
453
454