

---

# ASR-GLUE: A New Multi-task Benchmark for ASR-Robust Natural Language Understanding

---

**Lingyun Feng**  
Tsinghua University  
Tencent AI Lab  
Shen Zhen, China  
fly19@mail.tsinghua.edu.cn

**Jianwei Yu**  
Tencent AI Lab  
Shen Zhen, China  
tomasyu@tencent.com

**Deng Cai**  
The Chinese University of Hong Kong  
Hong Kong, China  
thisisjcykcd@gmail.com

**Songxiang Liu**  
Tencent AI Lab  
Shen Zhen, China  
shaunxliu@tencent.com

**Haitao Zheng**  
Tsinghua University  
Shen Zhen, China  
zheng.haitao@sz.tsinghua.edu.cn

**Yan Wang**  
Tencent AI Lab  
Shen Zhen, China  
brandenwang@tencent.com

## Abstract

Language understanding in speech-based systems has attracted attention in recent years with the growing demand for voice interface applications. However, the robustness of natural language understanding (NLU) systems to errors introduced by automatic speech recognition (ASR) is under-examined. In this paper, we propose *ASR-GLUE* benchmark, a new collection of 6 different NLU tasks for evaluating the performance of models under ASR error across 3 different levels of background noise and 6 speakers with various voice characteristics. Based on the proposed benchmark, we systematically investigate the effect of ASR error on NLU tasks in terms of noise intensity, error type and speaker variants. The analysis of this dataset shows that NLU under ASR errors is still very challenging and requires further research.<sup>1</sup>

## 1 Introduction

Language understanding in speech-based systems has attracted much attention in recent years with the growing demand for voice interface applications and devices such as Alexa [1], Siri [2], and Cortana [3]. These speech-based intelligent systems usually comprise an automatic speech recognition (ASR) component which converts audio signals to readable natural language text, and a natural language understanding (NLU) component which takes the output of the ASR component as input and fulfills downstream tasks such as sentiment analysis, natural language inference, and response selection. The upstream ASR error may propagate to the downstream NLU component and degrade the overall performance [4, 5]. In real-world scenarios, ASR error can be ubiquitous due to poor articulation and acoustic variability caused by environment noise and reverberation [6]. The persistence of ASR error indicates a need for ASR-robust natural language understanding.

---

<sup>1</sup>The dataset is available at  
[https://drive.google.com/drive/folders/1slqI6pUiab470vCxQBZemQZN-a\\_ssv1Q](https://drive.google.com/drive/folders/1slqI6pUiab470vCxQBZemQZN-a_ssv1Q)

Previous work in this area is limited to task-oriented language understanding such as hotel reservation and meeting scheduling through human-machine interactions [7, 8, 9, 10]. However, ASR error can affect many other NLU tasks, such as sentiment analysis in voice assistants. A benchmark that enables the comprehensive evaluation of NLU under ASR error on a diverse range of tasks is still missing.

In this paper, to quantitatively investigate how ASR error affects NLU capability, we propose the ASR-robust General Language Understanding Evaluation (*ASR-GLUE*) benchmark: a collection of 6 NLU tasks including sentiment analysis, similarity and paraphrase tasks, and natural language inference (NLI) tasks. We hire 6 native speakers to convert the test data into audio recordings with 3 different levels of environment noise. Each speaker is requested to record all test data to study the variance between individuals.

Finally, we get 18 different types of audio recordings (3 levels of noise \* 6 different speakers) for each of the 6 NLU tasks, varying in noise intensity, error type, and speaker variants. In addition, we also test human performance under different noise levels. We hope it would benefit the research of ASR-robust NLU in the future.

Our contributions are as follows: 1) A new benchmark dataset, *ASR-GLUE*, is proposed to enable a comprehensive evaluation of the robustness of NLU model to ASR error, covering 6 diversified tasks under 6 different speakers and 3 different levels of environment noise. 2) We systematically and quantitatively investigate the sensitivity of state-of-the-art NLU models to ASR error in terms of noise intensity, error type and speaker variants.

## 2 *ASR-GLUE*

### 2.1 Selected NLU Tasks

The proposed *ASR-GLUE* is constructed on the basis of GLUE [11], a popular NLU evaluation benchmark consisting of diverse NLU tasks. We select 5 typical NLU tasks from it, namely: Sentiment classification (SST-2 [12]), Semantic Textual Similarity (STS-B [13]), paraphrase (QQP<sup>2</sup>), Question-answering NLI (QNLI [14]), Recognizing Textual Entailment (RTE [15, 16, 17, 18].) and incorporate with a Science NLI task (SciTail [19]), resulting in 6 tasks in total. They are common and typical tasks of language understanding in speech-based scenarios, making them suitable for *ASR-GLUE*. These tasks are described in detail below.

**SST-2** The Stanford Sentiment Treebank [12] is a single-input understanding task for sentiment classification. The task is to predict the sentiment of a given sentence in the movie reviews domain. Accuracy (ACC) of the binary classification (positive or negative) is used as the metric.

**STS-B** The Semantic Textual Similarity Benchmark [13] consists of sentence pairs drawn from news headlines, video and image captions, and natural language inference data. The task is to predict sentence similarity scores which range from 1 to 5. We evaluate using Pearson and Spearman correlation coefficients.

**QQP** The Quora Question Pairs<sup>3</sup> dataset consists of question pairs in social QA questions domain. The task is to determine whether a pair of questions are semantically equivalent. Accuracy (ACC) is used as the metric.

**QNLI** Question-answering NLI is modified from the Stanford Question Answering dataset [14]. This is a sentence pair classification task that determines whether the context sentence contains the answer to the question. Accuracy (ACC) is used as the metric.

**SciTail** SciTail [19] is a recently released challenging textual entailment dataset collected from the science domain. This is a natural language inference task that determines if a natural language hypothesis can be justifiably inferred from a given premise. Accuracy (ACC) is used as the metric.

**RTE** The Recognizing Textual Entailment (RTE) datasets come from a series of annual textual entailment challenges merged from a collection of [15, 16]. All datasets are combined and converted to two-class classification: entailment and not entailment. Accuracy (ACC) is used as the metric.

---

<sup>2</sup>[data.quora.com/First-Quora-Dataset-Release-Question-Pairs](https://data.quora.com/First-Quora-Dataset-Release-Question-Pairs)

<sup>3</sup>[data.quora.com/First-Quora-Dataset-Release-Question-Pairs](https://data.quora.com/First-Quora-Dataset-Release-Question-Pairs)

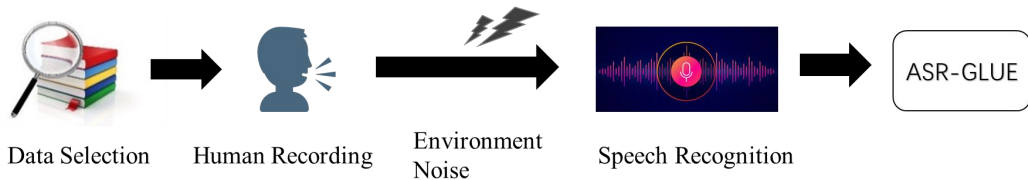


Figure 1: An illustration of the data collection and recording process.

## 2.2 Data Construction

Since the original datasets are presented in clean text form, we manually select instances from their test sets for human audio recording to evaluate the NLU capability in the presence of ASR error. Samples with non-standard words (e.g., abbreviations, currency expressions, strange names, and addresses) or sentences that are too long (more than 100 words) are excluded from these selected test sets. Considering the cost and quality of annotation, we keep the original training set and randomly select a subset of samples from the test set for human audio recording on each task<sup>4</sup>. The statistics of the data is shown in Table 1.

In the human recording process, 6 native speakers are hired to record all test samples. Different levels of environment noise is provided and the audio signals are sent into an ASR system to get the final ASR hypothesis. The overall process is depicted in Figure 1.

Table 1: Statistics of *ASR-GLUE*. The reported hours are the sum of recording time under different levels of noise. We also report Word Error Rates (WER) for test sets under each noise level.

Corpus	#Train	#Dev	#Test	WER (test)			Hours
				Low Noise	Medium Noise	High Noise	Test+Dev
SST-2	67349	2772	2790	18.35%	30.76%	34.86%	8.05
STS-B	5749	3042	3222	12.49%	24.70%	28.53%	10.82
QQP	363846	1476	3996	13.78%	24.34%	27.45%	11.56
QNLI	104743	2718	2718	22.51%	33.61%	37.90%	18.00
RTE	2490	2070	2088	24.54%	39.47%	47.03%	26.19
SciTail	23596	2718	2736	17.55%	30.09 %	34.04%	16.81

**Human Audio Recording** We hire six native speakers to record the test sets for each task. The speakers are 3 males and 3 females with different ages ranging from 18 to 40 from the U.S. In the recording process, each speaker is required to record all text samples independently so that we can study the impact of speaker variation. They are instructed to imagine they are communicating with someone when speaking the sentences of the six tasks. Note that they are allowed to make minor changes to the original text for natural and smooth expression such as change *cannot* to *can't*. To collect high-quality audio, all the original speech signals are recorded in a low-noise environment.

**Environment Noise** In real-world scenarios, human audio often recorded with environment noise and reverberation [20]. The presence of the background interference will lead to substantial performance degradation of current ASR systems [21] and further affect the downstream NLU systems [22]. Therefore, to better evaluate the robustness of NLU models in the noisy acoustic environment, speech data with different levels of noise is further provided in the *ASR-GLUE* corpus.

In *ASR-GLUE*, the widely-adopted simulation approach [23] is used to introduce different levels of noise and reverberation into the low-noise audio signals. Specifically, the background noise caused by such as phone ringing, alarm clocks and incoming vehicles are randomly sampled and added to the original recordings with a signal-to-noise-ratio (SNR) from 10dB to 15dB. Here SNR is defined as:  $SNR := 10 \log \frac{\|s_{target}\|^2}{\|s_{noise}\|^2}$ , where  $s_{target}$  and  $s_{noise}$  denote the acoustic signal of the clean speech and the noise respectively.  $\|s\|^2 = s^T s$  denotes the signal power. In addition, the room reverberation

<sup>4</sup>If there is no public test set, we use their dev set instead.

is also introduced by involving the recorded audio signals with the Room Impulse Responses (RIRs)<sup>5</sup> generated by the image-source method [24]. The simulation process totally covers 843 kinds of different background noise and 417 types of different RIRs.

Finally, for each recorded human audio signal, we get three versions: (1) Low-level noise, same as the original audio. (2) Medium-level noise which introduces reverberation and 15dB SNR level background noise into the original audio. (3) High-level noise, which introduces reverberation and 10dB SNR level background noise into the original audio. Then we build a 6000h trained ASR system based on the widely-used open-source Kaldi toolkit<sup>6</sup> [25, 26] to transcribe these audio files into text. Table 2 shows the WER of this system under different speakers and different noise levels.

Table 2: Detailed kaldi-based ASR WER on ASR-GLUE test set

Corpus	Noisy level	Speaker1	Speaker2	Speaker3	Speaker4	Speaker5	Speaker6
SST-2	Low	11.50	18.53	30.03	19.94	13.96	38.36
STS-B		6.87	11.89	17.76	12.62	7.67	13.62
QQP		10.44	13.05	25.46	13.10	10.77	25.49
QNLI		13.47	17.43	31.05	19.88	14.55	30.42
RTE		23.04	18.98	34.34	22.05	15.60	37.77
SciTail		8.47	13.98	23.49	16.83	9.18	27.53
SST-2		Medium	21.64	29.68	40.12	34.60	23.23
STS-B	20.06		28.16	31.11	21.79	16.49	29.65
QQP	18.85		22.13	39.01	19.90	18.26	40.39
QNLI	22.98		31.75	40.09	32.82	20.98	47.32
RTE	37.21		34.39	47.25	35.63	25.54	59.38
SciTail	18.71		26.72	34.75	33.20	15.89	44.13
SST-2	High		25.98	33.90	42.82	37.77	27.27
STS-B		24.66	32.61	33.22	25.09	18.83	33.71
QQP		21.64	24.45	41.52	22.11	22.23	43.77
QNLI		28.09	36.05	42.96	35.53	24.76	51.73
RTE		44.12	41.20	54.56	50.70	33.33	64.91
SciTail		23.13	29.89	36.60	37.65	19.93	48.55

### 3 Analysis of ASR-GLUE

In this section we give extensive analyses on the proposed *ASR-GLUE* dataset. In Section 3.1, we obtain human performance to measure the ceiling performance on the test set. Then we analyse the performance of recent SOTA NLU models across different levels of environment noise in Section 3.2. Furthermore, in Section 3.3 we categorize ASR errors into four types, and systematically investigate the impact of different error types on NLU performance. Finally, we analyse the effect of voice variation from different speakers in Section 3.4.

#### 3.1 Human Performance

To obtain human performance under environment noise, we hire native annotators to predict labels of each test sample in audio form.

The annotators are hired from a third-party company. To guarantee the annotators fully understand these tasks, we first give the annotators a brief training on each task in *ASR-GLUE*. Then we ask them to take an exam before starting annotation. Only annotators who pass the exam will be hired. Finally, we have 3 annotators to measure the ceiling performance for each task in *ASR-GLUE*. Details about the exam and annotation process are presented below.

In the annotation process, the annotators are required to predict the labels of each test sample according to the corresponding audio signals. Note that for each test sample we have 18 audio signals (3 levels of noise \* 6 speakers), which is a big burden for the annotators. So we randomly select

<sup>5</sup>The noise and RIR files can be found at [http://www.openslr.org/resources/28/rirs\\_noises.zip](http://www.openslr.org/resources/28/rirs_noises.zip)

<sup>6</sup><https://github.com/kaldi-asr/kaldi>

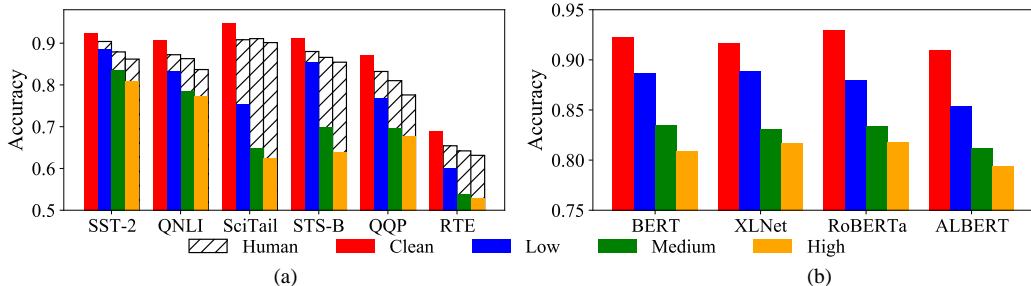


Figure 2: (a) The performance of BERT on different tasks under different levels of noise. The shaded area represents human performance. (b) Accuracy results for different model architecture on SST-2 task. Here “Human” indicates human performance under various noise settings. “Clean” stands for test on clean text data. “Low/Medium/High” stands for test in low-level/medium-level/high-level noise respectively.

Table 3: Examples of speech recognition perturbation in *ASR-GLUE*.

Error Type	Ground Truth	Recognition Result
Similar sounds	The man couldn't <b>lift</b> his son.	The man couldn't <b>lived</b> his son.
	Tommy <b>dropped</b> his ice cream.	Tommy <b>jumped</b> his ice cream.
Liaison	Does Quora <b>stand for</b> question or answer.	Does Quora <b>Stanford</b> question or answer.
	The <b>drain is</b> clogged with hair.	The <b>drains</b> clogged with hair.
Insertion	This afternoon.	This <b>after</b> afternoon.
	A warm funny * engaging film	A warm funny <b>and</b> engaging film.
Deletion	A black and white photo of an old train station.	A black * white * of * train station
	Old style bicycle <b>parked on</b> floor	Old style bicycle * * floor

one audio from the six speakers with each noise level and report the average performance of the annotators on all tasks.

The exam is set to guarantee the annotators fully understand each task. In the exam, we randomly select 50 samples from the original datasets which are in text form for each task. Annotators who achieve at least 90% accuracy on these samples will be hired.

### 3.2 Performance of Existing NLU Models

To demonstrate the significance of the ASR error issue, we leverage typical NLU models such as BERT to test their robustness to different levels of ASR error on different tasks. As shown in Figure 2(a), While BERT yields promising results on error-free text, its performance degrades in the presence of ASR error on six tasks. As the noise increases, the performance of the model drops more severely. In contrast, humans are less affected by the environment noise, which indicates that there still remains a gap between the robustness of models and humans to ASR error.

We also investigate the effect of ASR error with different noise levels on different NLU models. We adopt base version of BERT [27], RoBERTa [28], ALBERT [29], XLNet [30] as the NLU model. As shown in Fig. 2(b) we can observe that all these pretrained language models are sensitive to ASR error and the performance degrades with the increase of the noise level.

### 3.3 Breakdown Analysis by ASR Error Type

We conclude that the ASR error types can be categorized into four-folds, namely similar sounds, liaison, insertion, and deletion. Specifically, (i) Similar sounds happen when the ASR system sometimes wrongly identifies one word as another with similar pronunciation. (ii) Liaison constantly occurs between successive words which have sound fusion across word boundaries. (iii) Insertion happens when the ASR system makes word redundancies. (iv) Deletion happens when there are word omissions in the ASR hypothesis. Examples of each error type are presented in Table 3.

We report the percentage of each aforementioned error type in Figure 3(a). We choose the SST-2 task as an example and observe that similar sounds most commonly happen. As the noise increases, the

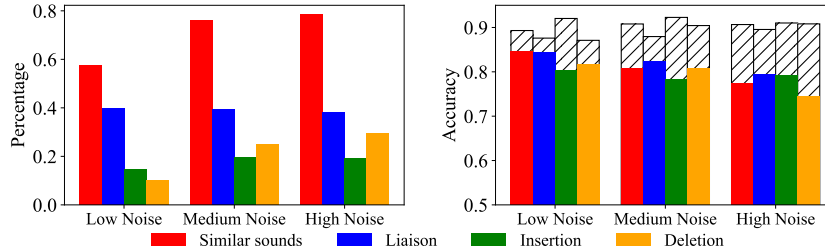


Figure 3: Left: The percentage of each error type under different noise setting in SST-2 dataset. Right: The accuracy of BERT on four subsets under different noise level. Each subset only contains test samples with one specific error type. For example, the red block represents the accuracy of BERT on test samples which contain similar sounds errors. The shaded area represents the performance degradation against clean text caused by a specific error type.

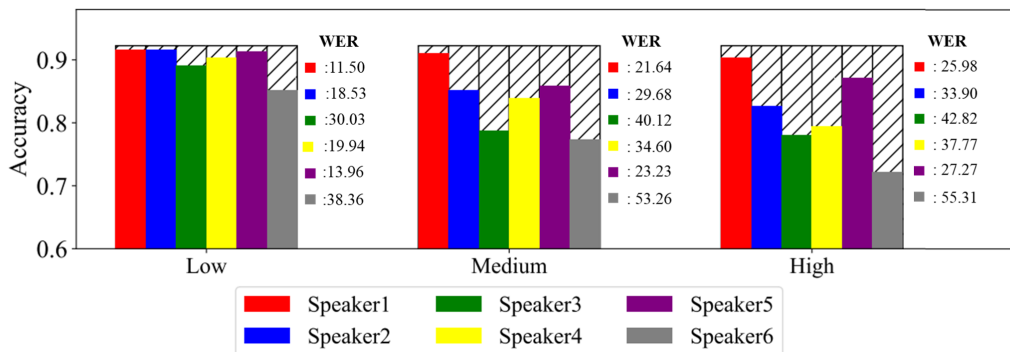


Figure 4: NLU performance under voice variation from different speakers.

percentage of similar sounds and deletion made by the ASR system gradually increases while the percentage of liaison and insertion error types remain relatively stable. Note that the sum percentages of the four error types are over 100%, since different error types may simultaneously exist in one hypothesis.

We also investigate the impact of each error type on NLU model on SST-2 task. We group data according to error types and compare the model performance with the same grouped raw data without errors respectively. As shown in Figure 3 (b), we can observe that BERT can better handle Liaison error for the performance degradation is minimal under varied noise settings. As noise increases, the accuracy of the model decreases more severely for each different type of error.

### 3.4 Effect of Individual Difference

We test the effect of voice variation of the six hired speakers on the test set of SST-2. As shown in Fig. 4 we can observe that the recognition quality and classification accuracy vary greatly across speakers. For example, the accuracy of BERT on S1 (corresponds to Speaker 1) is consistently higher than S6 (corresponds to Speaker 6) with a very large margin (27% ~ 32%). Meanwhile, we can also observe that the drop in classification accuracy is due to the increase of WER. The higher WER means more ASR errors in the test samples, resulting in more misclassification.

## 4 Related Work

Many benchmark datasets are created to facilitate Spoken Language Understanding (SLU) [31, 32, 5, 33, 22, 34], which evaluate the robustness of the downstream NLU model against the error output from the upstream acoustic model [7, 8, 9, 10]. However, they are only designed for a particular domain or a specific task such as intent detection and slot filling. In contrast, the human ability to

understand language is general, flexible, and robust. There is a need to test general-purpose natural language understanding capability on diverse tasks in different domains.

Large-scale pretrained language models have achieved striking performance on NLU in recent years [35, 30]. Recently many works test their robustness by human-crafted adversarial examples [36] or generated examples by adversarial attacks [35, 37, 38, 39]. [40] projects the input data to a latent space by generative adversarial networks (GANs), and then retrieves adversaries close to the original instance in the latent space. [41] propose controlled paraphrase networks to generate syntactically adversarial examples that both fool pre-trained models and improve the robustness of these models to syntactic variation when used to augment their training data. However, the robustness of pre-trained model to speech recognition error in real conditions has not been fully explored.

## 5 Conclusion

We present *ASR-GLUE*, a new benchmark for evaluating general-purpose language understanding under ASR error in speech-based applications. We propose two ways to improve robustness of the NLU system and find that there is still a gap between the NLU capability of the model and humans. *ASR-GLUE* offers a rich and challenging testbed for work developing ASR robust model for general-purpose language understanding. Given the difficulty of *ASR-GLUE*, we expect that further progress in multi-task, multi-model learning techniques will be necessary to approach human-level performance on the benchmark.

## References

- [1] Wang, L., M. Fazel-Zarandi, A. Tiwari, et al. Data augmentation for training dialog models robust to speech recognition errors. [arXiv preprint arXiv:2006.05635](#), 2020.
- [2] Williams, J. D., S. Young. Partially observable markov decision processes for spoken dialog systems. *Computer Speech & Language*, 21(2):393–422, 2007.
- [3] Wang, S., T. Gunter, D. VanDyke. On modelling uncertainty in neural language generation for policy optimisation in voice-triggered dialog assistants. In *2nd Workshop on Conversational AI: Today’s Practice and Tomorrow’s Potential, NeurIPS*. 2018.
- [4] Serdyuk, D., Y. Wang, C. Fuegen, et al. Towards end-to-end spoken language understanding. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5754–5758. IEEE, 2018.
- [5] Wang, P., L. Wei, Y. Cao, et al. Large-scale unsupervised pre-training for end-to-end spoken language understanding. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7999–8003. IEEE, 2020.
- [6] Errattahi, R., A. El Hannani, H. Ouahmane. Automatic speech recognition errors detection and correction: A review. *Procedia Computer Science*, 128:32–37, 2018.
- [7] Schumann, R., P. Angkititrakul. Incorporating asr errors with attention-based, jointly trained rnn for intent detection and slot filling. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6059–6063. IEEE, 2018.
- [8] Weng, Y., S. S. Miryala, C. Khatri, et al. Joint contextual modeling for asr correction and language understanding. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6349–6353. IEEE, 2020.
- [9] Rao, M., A. Raju, P. Dheram, et al. Speech to semantics: Improve asr and nlu jointly via all-neural interfaces. [arXiv preprint arXiv:2008.06173](#), 2020.
- [10] Huang, C.-W., Y.-N. Chen. Learning asr-robust contextualized embeddings for spoken language understanding. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8009–8013. IEEE, 2020.
- [11] Wang, A., A. Singh, J. Michael, et al. Glue: A multi-task benchmark and analysis platform for natural language understanding. [arXiv preprint arXiv:1804.07461](#), 2018.
- [12] Socher, R., A. Perelygin, J. Wu, et al. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642. 2013.

- [13] Cer, D., M. Diab, E. Agirre, et al. Semeval-2017 task 1: Semantic textual similarity-multilingual and cross-lingual focused evaluation. arXiv preprint arXiv:1708.00055, 2017.
- [14] Rajpurkar, P., J. Zhang, K. Lopyrev, et al. Squad: 100,000+ questions for machine comprehension of text. arXiv preprint arXiv:1606.05250, 2016.
- [15] Dagan, I., O. Glickman, B. Magnini. The pascal recognising textual entailment challenge. In Machine Learning Challenges Workshop, pages 177–190. Springer, 2005.
- [16] Haim, R. B., I. Dagan, B. Dolan, et al. The second pascal recognising textual entailment challenge. In Proceedings of the Second PASCAL Challenges Workshop on Recognising Textual Entailment. 2006.
- [17] Giampiccolo, D., B. Magnini, I. Dagan, et al. The third pascal recognizing textual entailment challenge. In Proceedings of the ACL-PASCAL workshop on textual entailment and paraphrasing, pages 1–9. 2007.
- [18] Bentivogli, L., P. Clark, I. Dagan, et al. The fifth pascal recognizing textual entailment challenge. In TAC. 2009.
- [19] Khot, T., A. Sabharwal, P. Clark. Scitail: A textual entailment dataset from science question answering. In Proceedings of the AAAI Conference on Artificial Intelligence, vol. 32. 2018.
- [20] Wölfel, M., J. McDonough. Distant speech recognition. John Wiley & Sons, 2009.
- [21] Barker, J., S. Watanabe, E. Vincent, et al. The fifth ‘chime’ speech separation and recognition challenge: dataset, task and baselines. arXiv preprint arXiv:1803.10609, 2018.
- [22] Henderson, M., B. Thomson, J. D. Williams. The second dialog state tracking challenge. In Proceedings of the 15th annual meeting of the special interest group on discourse and dialogue (SIGDIAL), pages 263–272. 2014.
- [23] Ko, T., V. Peddinti, D. Povey, et al. A study on data augmentation of reverberant speech for robust speech recognition. In 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 5220–5224. IEEE, 2017.
- [24] Habets, E. A. Room impulse response generator. Technische Universiteit Eindhoven, Tech. Rep., 2(2.4):1, 2006.
- [25] Povey, D., A. Ghoshal, G. Boulianne, et al. The kaldi speech recognition toolkit. In IEEE 2011 workshop on automatic speech recognition and understanding, CONF. IEEE Signal Processing Society, 2011.
- [26] Povey, D., V. Peddinti, D. Galvez, et al. Purely sequence-trained neural networks for asr based on lattice-free mmi. In Interspeech, pages 2751–2755. 2016.
- [27] Devlin, J., M.-W. Chang, K. Lee, et al. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018.
- [28] Liu, Y., M. Ott, N. Goyal, et al. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692, 2019.
- [29] Lan, Z., M. Chen, S. Goodman, et al. Albert: A lite bert for self-supervised learning of language representations. arXiv preprint arXiv:1909.11942, 2019.
- [30] Yang, Z., Z. Dai, Y. Yang, et al. Xlnet: Generalized autoregressive pretraining for language understanding. arXiv preprint arXiv:1906.08237, 2019.
- [31] Coucke, A., A. Saade, A. Ball, et al. Snips voice platform: an embedded spoken language understanding system for private-by-design voice interfaces. arXiv preprint arXiv:1805.10190, 2018.
- [32] Lugosch, L., M. Ravanelli, P. Ignoto, et al. Speech model pre-training for end-to-end spoken language understanding. arXiv preprint arXiv:1904.03670, 2019.
- [33] Price, P. Evaluation of spoken language systems: The atis domain. In Speech and Natural Language: Proceedings of a Workshop Held at Hidden Valley, Pennsylvania, June 24-27, 1990. 1990.
- [34] Peng, B., C. Li, Z. Zhang, et al. Raddle: An evaluation benchmark and analysis platform for robust task-oriented dialog systems. arXiv preprint arXiv:2012.14666, 2020.



- [35] Jin, D., Z. Jin, J. T. Zhou, et al. Is bert really robust? a strong baseline for natural language attack on text classification and entailment. In Proceedings of the AAAI conference on artificial intelligence, vol. 34, pages 8018–8025. 2020.
- [36] Nie, Y., A. Williams, E. Dinan, et al. Adversarial nli: A new benchmark for natural language understanding. arXiv preprint arXiv:1910.14599, 2019.
- [37] Madry, A., A. Makelov, L. Schmidt, et al. Towards deep learning models resistant to adversarial attacks. arXiv preprint arXiv:1706.06083, 2017.
- [38] Zhu, C., Y. Cheng, Z. Gan, et al. Freeb: Enhanced adversarial training for natural language understanding. arXiv preprint arXiv:1909.11764, 2019.
- [39] Dong, X., A. T. Luu, R. Ji, et al. Towards robustness against natural language word substitutions. In 9th International Conference on Learning Representations (ICLR). 2021.
- [40] Zhao, Z., D. Dua, S. Singh. Generating natural adversarial examples. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net, 2018.
- [41] Iyyer, M., J. Wieting, K. Gimpel, et al. Adversarial example generation with syntactically controlled paraphrase networks. arXiv preprint arXiv:1804.06059, 2018.

## 6 Datasheet

### 6.1 Dataset Motivation

(1) For what purpose was the dataset created? Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.

This benchmark was created for the purpose of comprehensively investigate how ASR error affect NLU capability. We propose the ASR-robust General Language Understanding Evaluation (ASR-GLUE) benchmark: a new collection of 6 different NLU tasks for evaluating the performance of models under ASR error across 3 different levels of background noise and 6 speakers with various voice characteristics.

(2) Who created this dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?

Tencent AI Lab created and funded this dataset.

### 6.2 Dataset Use Cases

(1) What (other) tasks could the dataset be used for?

The dataset can be used for sentiment analysis, paraphrase identification, natural language inference and many other NLU tasks. More importantly, this dataset can be used to comprehensively investigate how ASR error affect NLU capability.

(2) Are there tasks for which the dataset should not be used? If so, please provide a description

No.

### 6.3 Dataset Maintenance

(1) Who is supporting/hosting/maintaining the dataset? How can the owner/curator/manager of the dataset be contacted?

Tencent AI Lab will support and maintain the dataset. You can contact the owner via following e-mails: Yan Wang, brandenwang@tencent.com Jianwei Yu, tomasyu@tencent.com

(2) Is there an erratum? If so, please provide a link or other access point.

No.

(3) Will the dataset be updated? Will older versions of the dataset continue to be supported/hosted/maintained?

We will update in future if necessary. We will support and maintain all versions.

(4) If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so?

Others may do so. They just need notice the authors of this paper in advance.

### 6.4 License

The authors bear all responsibility in case of violation of rights.