



kosp2e: Korean Speech to English Translation Corpus

Won Ik Cho¹, Seok Min Kim¹, Hyunchang Cho², Nam Soo Kim¹

¹Dept. of ECE and INMC, Seoul National University, Seoul, Korea

²PAPAGO, NAVER, Seongnam, Korea

{wicho, smkim}@hi.snu.ac.kr, hyunchang.cho@navercorp.com, nkim@snu.ac.kr

Abstract

Most speech-to-text (S2T) translation studies use English speech as a source, which makes it difficult for non-English speakers to take advantage of the S2T technologies. For some languages, this problem was tackled through corpus construction, but the farther linguistically from English or the more under-resourced, this deficiency and underrepresentedness becomes more significant. In this paper, we introduce *kosp2e* (read as ‘kospī’), a corpus that allows Korean speech to be translated into English text in an end-to-end manner. We adopt open license speech recognition corpus, translation corpus, and spoken language corpora to make our dataset freely available to the public, and check the performance through the pipeline and training-based approaches. Using pipeline and various end-to-end schemes, we obtain the highest BLEU of 21.3 and 18.0 for each based on the English hypothesis, validating the feasibility of our data. We plan to supplement annotations for other target languages through community contributions in the future.

Index Terms: Korean, English, speech translation, corpus

1. Introduction

Speech to text (S2T) translation is achieving the text of the target language from the speech of the source language. Speech translation is actively used in international conferences, video subtitling, and real-time translation for tourists. Traditionally, automatic speech recognition (ASR) and machine translation (MT) were adopted as a cascading method. Recently, end-to-end approaches have been widely utilized owing to the success of data-driven methodologies [1, 2].

However, studies on end-to-end speech translation using languages other than English as sources have not been driven actively. For instance, the survey in MuST-C [3] suggests that among the 13 speech translation corpora, 12 handle English as source or target, and about half of the total temporal volume (547 hours over 1080 hours), including the largest ones [4, 5] only deals with English as a source speech.

On one side, the reason English is exploited as a source language is probably that it fulfills large industrial needs [6]. It is well-known as a world-widely used language that facilitates communication. Besides, recent S2T translation corpora [5, 3] usually leverages the ASR databases which were mainly proposed for English speech recognition, such as LibriSpeech [7], which is plausible considering that recording speech in individual languages is extremely costly.

However, speech translation shows a difference in application depending on what the source language is. In international conferences where speakers usually utter English, speech translation modules trained with an English source will obviously be effective. But what about YouTube video subtitling or tourists’ real-time translation? Many videos contain utterances in the casts’ own language. Also, it is challenging for non-English cit-

izens to ask fluent English when they tour foreign countries. In other words, there are more points to consider on where non-English source speech can be utilized.

In this respect, we deemed that speech translation with a non-English speech source also needs attention. Among them, Korean is one of the languages that have been less investigated in S2T translation due to the deficiency of open spoken language corpora [8] and its syntax and writing system being distinct from largely studied Indo-European (IE) languages. We construct a ko-en S2T translation dataset, *kosp2e* (read as ‘kospī’), for the Korean spoken language, and release it publicly. This parallel corpus is the first to be constructed for S2T translation with Korean speech, and is constructed via manual translation of scripts, using or recording Korean speech databases.

Our contribution to the community is as follows:

- We first build the Korean to English speech translation corpus, adopting the open-licensed speech recognition, machine translation, and spoken language corpora.
- We check that the feasibility of using ASR pretrained model in an end-to-end manner is promising in ko-en speech translation as well.

2. Related Work

Automatic speech recognition (ASR) denotes transcribing speech into text with the letters adopted in the language’s writing system. Therefore, speech corpus used for training and evaluation of ASR has audio files and transcripts (sometimes with timestamps) accordingly. In the case of Korean speech recognition, Zeroth¹, KSS [9], and ClovaCall [10] are used as benchmark resources in the literature.

Machine translation (MT) is mapping a text of the source language into the text of the target language. Due to this characteristic, MT corpora are often referred to as a parallel corpus. Storage and inspection of text corpora are easier than those of speech corpora, but bilingual processing is inevitable in making up MT corpora. This often makes construction costly, and there are not sufficient parallel corpora used as open resources in Korean [8]. Although resources such as Korean Parallel Corpora [11] and Open Subtitles² exist, it is status quo that manually translated MT resources lack at this point compared to IE languages [12].

Speech translation is not simply a combination of speech recognition and machine translation database; it refers to a case in which tuples of source speech - source text - target text exist in parallel. The representative one is MuST-C [3], which contains the translation to eight languages with English TED talks as a source. It contains eight IE target language texts, namely

¹<https://github.com/goodatlas/zeroth>

²<https://www.opensubtitles.org/ko>

German, Spanish, French, Italian, Dutch, Portuguese, Romanian, and Russian. There is recently distributed S2T dataset incorporating Korean regarding Duolingo challenge³ [13], presented at the past IWSLT workshop. However, the collecting process was reading the prompts for fluency check rather than speaking natural language sentences, that we deemed difficult to utilize it for Korean train/test as our aim. Also, there is currently no domestic S2T translation resource available, up to our knowledge, whether it allows public access or not.

3. Corpus Construction

We aim at the open-licensed distribution of our data contribution, either in commercial or non-commercial perspective. This allows the potential remix of the dataset, which can be done by not only the authors but also by community contribution. In this regard, we survey the open Korean spoken language corpora currently available with CC-BY license and list up the properties thereof, along with how we set up the principles in treating them in the recording and translation phase.

3.1. Corpora

3.1.1. KSS

KSS [9] is a single speaker speech corpus used for Korean speech recognition and has a license of CC-BY-NC-SA 4.0⁴. It consists of about 13K scripts, which a single voice actress utters and also includes the English translation. The scripts are sentences from the textbook for language education, and contain descriptions of daily life.

We re-recorded this corpus with verified crowd-workers so that diverse voices can be provided. In this process, no specific instruction was given. When it was not clear how to read, the participants were recommended to refer to the given translation or to record after listening to the source speech. A human inspector accompanied every output of the recording.

3.1.2. Zeroth

Zeroth is a Korean speech recognition dataset, and is released under a license of CC-BY 4.0⁵. It consists of a total of about 22K scripts, and 3K unique sentences extracted from the news domain are used as scripts.

Zeroth was adopted for translation. Since it mainly consists of news scripts, it contains information on politics, social events, and public figures. We had the following guideline for the translation:

- Some reporter names are inserted at the very first of the scripts. Translators may take it into account by, for example, placing a comma after the name coming at the front of the sentence.
- Public figures' names are fixed to the format of 'Family name - First name' such as 'Cho Won-ik'.
- Entities such as organizations, locations, or titles, are translated into English if adequate cognate exists, and are romanized if not.

³<https://sharedtask.duolingo.com/>

⁴<https://www.kaggle.com/bryanpark/korean-single-speaker-speech-dataset>

⁵<https://github.com/goodatlas/zeroth>

3.1.3. StyleKQC

StyleKQC [14] is a script of 30K Korean directives that can be used in the AI agent domain, publicly available as CC-BY-SA 4.0⁶. It contains six topics, namely messenger, calendar, weather and news, smart home, shopping, and entertainment, and four speech acts of alternative questions, *wh*-questions, requirements, and prohibition. In addition, every utterance is tagged with either it has a formal or informal style of speaking. Ten sentences with the same intent are provided in a group. The original corpus contains text without punctuation, but in the translation and recording process, the speech act is also provided so that one can disambiguate the semantics and reflect it in translation or recording. Specifically, the following was suggested.

Recording

- Though the original script does not contain punctuation, the speech act types should be taken into account while recording.
- The tone of the recording should be differentiated between formal and informal utterances; formal utterances carefully as when facing elderly or colleagues, and informal utterances as when heading friends.
- Albeit some sentences may seem awkward due to scrambling or fillers, read it as naturally as possible by placing a delay or pronouncing it chunk by chunk.

Translation

- Although the sentences of the same intent may have the same connotation, the difference in style of each utterance should be reflected in the translation (e.g., "How many people are there in the US currently who are diagnosed as Covid 19?" vs. "What is the current number of Covid 19 patients in the US?" as a difference in the sentence structure, or "You know the current number of Covid 19 people in States?" as a difference in the formality)
- Since the original corpus contains free-style spoken language that includes scrambling, filler, or sometimes fragments, the translators are asked to reflect that into the English translation
- Translation on names and entities follow the principles proposed for Zeroth.

3.1.4. Covid-ED

Covid-ED (COVID-19 Emotion Diary with Empathy and Theory-of-Mind Ground Truths Dataset) is a collection of crowdsourced diaries written in pandemic situations. It is labeled with emotions, empathy, personality, and levels of theory-of-mind. The dataset is publicly available as CC-BY-NC-SA 4.0⁷ [15]. Upon writing diaries, workers were told to either exclude or anonymize personally identifiable information (e.g., address, residential number, bank accounts, etc.) Each worker wrote five diaries for five consecutive days. Ground truth emotion labels per document were provided by the writers so that such information can be reflected in translation and recording. We considered the following in the process of recording and translation.

⁶<https://github.com/cynthia/stylekqc>

⁷<https://github.com/humanfactorspsych/covid19-tom-empathy-diary>

Table 1: *kosp2e* subcorpus specification by domain.

Dataset	License	Domain	Characteristics	Volume (Train / Dev / Test)	Tokens (ko / en)	Speakers (Total)
Zeroth	CC-BY 4.0	News / newspaper	DB originally for speech recognition	22,247 utterances (3,004 unique scripts) (21,589 / 197 / 461)	72K / 120K	115
KSS	CC-BY-NC-SA 4.0	Textbook (colloquial descriptions)	Originally recorded by a single speaker (multi-speaker recording augmented)	25,708 utterances = 12,854 * 2 (recording augmented) (24,940 / 256 / 512)	128K / 190K	17
StyleKQC	CC-BY-SA 4.0	AI agent (commands)	Speech act (4) and topic (6) labels are included	30,000 utterances (28,800 / 480 / 720)	237K / 391K	60
Covid-ED	CC-BY-NC-SA 4.0	Diary (monologue)	Sentences are in document level; emotion tags included	32,284 utterances (31,324 / 333 / 627)	358K / 571K	71

Recording

- One participant should record all the diaries of one writer, and in this process, the gender and age of the diary writer and those of the one who records should be aligned as much as possible.
- The recording should be done considering (up to) two emotions that are tagged per diary. However, in this case, emotions do not need to be reflected in every sentence, instead, to relevant sentences.
- The diary is written in monologue format, but is basically web text, that there are various non-Hangul expressions in the script. Therefore, there may be variance in reading English words, numbers, and special symbols. The details of handling those are separately provided⁸.

Translation

- Sentences appear subsequently due to the script being a diary; thus, the subject omission may occur accordingly. In translating these parts, the translator should consider as much as possible the circumstances that the sentences are treated independently.
- English words should be translated as they are. The numbers should be translated according to whether it is a cardinal/ordinal number.
- Leetspeaks such as ㅋㅋ (laughter) and ㅠㅠ (sadness) should be translated to online expressions such as ‘lol’ and ‘T.T’.

3.2. Recording and Translation

3.2.1. Recording

The recording was conducted for KSS, StyleKQC, and Covid-ED. KSS includes speech files, but since it was recorded by a single speaker, additional recordings were performed as mentioned above. Our final corpus includes both original and recorded versions. StyleKQC and Covid-ED were recorded in the same way. The recording was performed by selected workers in the crowd-sourcing group *Green Web*⁹ and *Ever Young*¹⁰, with some considerations regarding device and environment¹¹.

⁸<https://github.com/warnikchow/kosp2e/wiki/English,-Numbers,-and-Symbols-in-Covid-ED>

⁹<https://www.gwebscorp.com/>

¹⁰<http://everyoungkorea.com/>

¹¹<https://github.com/warnikchow/kosp2e/wiki/Considerations-in-Recording>

3.2.2. Translation

The translation was conducted for Zeroth, StyleKQC, and Covid-ED. All the texts were translated by a translation expert and checked by an inspector, managed by Lexcode¹². Except for KSS, the same post-processing rules were applied to scripts that have undergone expert translation (normalizing double spaces, removing non-Hangul and special tokens, etc.). In addition, we refined the translation of KSS referring to this standard.

3.3. Statistics

The specification of our corpus is as follows.

- Utterances: 110,239 / Hours: 198H
- Source language (tokens): Korean (795K)
- Target language (tokens): English (1,272K)
- Speakers: 263 (in total)
- Domains: News, Textbook, AI agent command, Diary

Domain-wise descriptions of the subcorpora are in Table 1. For further usage, we attach the license policy of the original corpora. Two are non-commercial, but all subcorpora are available for remix and redistribution at least for academic purposes.

4. Experiment

We split the proposed corpus into train, dev, and test set. In this process, for each type of subcorpus, speaker or sentence types were distributed with balance. In specific, for Zeroth, where the number of unique sentences is much smaller than that of the utterances, we reformulated the original train-test set so that there are only unseen sentences in the test set. For KSS, the original dataset was not included in the training set since the single speaker voices being dominant can badly affect the acoustic diversity of the corpus. For StyleKQC, where the original dataset has no particular test set for considering topic, act, and style at the same time, we split the whole dataset so that such attributes do not cause a bias. The split on Covid-ED was done in the way that there is no overlap of diary writers between each set.

In total, we have 1,266 and 2,320 utterances each for the dev and test set. The train set contains the rest, namely 106,653 utterances. We performed four ko-en speech translation experiments using the sets, with the evaluation metric of corpus BLEU [16] of Sacrebleu [17] python library. Submodule performances were separately measured, especially using word error rate (WER) for Korean ASR.

¹²<https://lexcode.co.kr/>

Table 2: Pipeline and end-to-end implementation using the constructed corpus. Note that the submodules of ASR-MT are evaluated with our test set, and those of ASR pretraining and warm-up are evaluated with randomly split validation set (not official).

Model	BLEU	Submodules	
		WER (ASR)	BLEU (MT/ST)
ASR-MT (Pororo)	16.6	34.0	18.5 (MT)
ASR-MT (PAPAGO)	21.3	34.0	25.0 (MT)
Transformer (Vanilla)	2.6	-	-
ASR pretraining	5.9	24.0*	-
Transformer + Warm-up	11.6	-	35.7 (ST)*
+ Fine-tuning	18.0	-	-

4.1. ASR-MT Pipeline

In the ASR-MT pipeline, publicly available ASR and MT modules were used. We adopted Speech Recognition [18]¹³ toolkit for ASR, using Korean as option, and MT was performed with Pororo¹⁴ [19] NLP toolkit and PAPAGO translation API¹⁵. The performance of both modules was checked with the utterances of the test set (Table 2).

4.2. End-to-end Manner

For end-to-end implementation, fairseq-S2T [20] based on fairseq [21] was used. Three approaches were implemented; first vanilla transformer [22], second using ASR pretrained model, and the last augmenting pseudo-gold translations for model warm-up. For the vanilla model, we stacked 12 transformers for the encoder and 6 for the decoder, both with 8 attention heads. For the second approach, we adopted a large-scale speech corpus (of 1,000H Korean utterances) publicly available at AI HUB¹⁶ for ASR pretraining. The ASR module was trained based on fairseq script using source language text, and the soft label was concatenated with the decoder transformer trained upon it. The dimension of embeddings and dropouts were fixed to 256 and 0.1 each. For the last, which is inspired by [23], we machine-translated Korean scripts of AI HUB data with PAPAGO API. The transformer is first warmed up 8 epochs by large-scale pseudo-gold speech translation samples, and is further fine-tuned 30 epochs with our data.

4.3. Results

The results of the pipeline and end-to-end baseline models are exhibited in Table 2. In the pipeline, though the translation scheme of the dataset used in the adopted module may differ from our translation scheme, we deemed that the results show the consistency of our corpus with conventional MT benchmarks used for Korean speech recognition and ko-en MT. In an end-to-end manner, we have obtained a mediocre performance for the vanilla transformer but a much-improved result for the ones using ASR pretrained model and pseudo-gold translations. The results suggest that it is not feasible with vanilla end-to-end models to yield transformation between Korean speech and English text at this point (which have different modality and

distinct syntax at the same time), and it necessitates at least soft symbol-level representation or sufficient audio samples to get a satisfying ST performance.

4.4. Discussion

We intended to verify via an experiment that our corpus qualifies as a spoken language translation corpus. That is, we tried to show whether ours can reach the standard performance of ST benchmarks, overcoming the limitation in the gap of syntax and writing system between Korean and English. It seems that the fully end-to-end approach needs improvement, but the ones that leverage the ASR pretrained model or pseudo-gold translations promise the utility of advanced strategies such as large-scale ASR or ST pretrained accompanied with distillation strategies [24].

Our corpus has both strengths and weaknesses. First, we have incorporated the various style of utterances, from formal (news/newspaper) to colloquial (daily descriptions/smart home/diary). Also, we considered the vocabularies of diverse domains, which may contribute to S2T translation in real-world applications. Lastly, we constructed our corpus leveraging open datasets so that the distributed version can be updated and re-distributed with community contributions. This will make the dataset more viable to further annotations, such as tagging emotion or dialog acts.

On the other side, our corpus has some weaknesses in its level of perfection, due to its scripts or recordings being adopted from open Korean language resources, which may have less consistency in between. Also, the recording format is not unified to a single type, which requires slightly more preprocessing (audio reading and resampling) before the training phase. However, the errata or other flaws in the original content were checked and corrected in the inspection process of our construction, and diverse speech files produced in non-fixed environments instead cover the speech inputs of various quality that are probable in real-world applications. Therefore, we deem that the weaknesses do not hinder our contribution as the first fully open resource in from-Korean speech translation, and we expect the issues to be resolved with user reports and remixes.

5. Conclusion

Through this study, we constructed a to-English translation dataset using Korean, which is less studied in speech translation as a source speech. In this process, ASR/MT corpora and spoken language corpora were used, either by means of recording and translation, to make up the whole corpora of about 110K translated utterances. In addition, we evaluated the training results of our corpus via pipeline and end-to-end manner, obtained the best BLEU of 21.3 and 18.0 for each, and discussed where the improvement should be made. We plan to open this corpus publicly to contribute to the speech and translation community¹⁷, and will expand it into a corpus considering more target languages, with further community contributions.

6. Acknowledgements

This work was supported by PAPAGO, NAVER Corp. The authors appreciate Hyoung-Gyu Lee, Eunjeong Lucy Park, Ji-hyung Moon, and Doosun Yoo for discussions and support. Also, the authors thank Taeyoung Jo, Kyubyong Park, and Yoon Kyung Lee for sharing the resources.

¹³https://github.com/Uberi/speech_recognition

¹⁴<https://github.com/kakaobrain/pororo>

¹⁵<https://papago.naver.com/>

¹⁶The dataset is downloadable from <https://aihub.or.kr/aidata/105>, but the detailed usage should refer to <https://github.com/sooftware/KoSpeech> for non-Korean citizens.

¹⁷<https://github.com/warnikchow/kosp2e>

7. References

- [1] A. Bérard, O. Pietquin, C. Servan, and L. Besacier, “Listen and translate: A proof of concept for end-to-end speech-to-text translation,” *arXiv preprint arXiv:1612.01744*, 2016.
- [2] A. Bérard, L. Besacier, A. C. Kocabiyikoglu, and O. Pietquin, “End-to-end automatic speech translation of audiobooks,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 6224–6228.
- [3] M. A. Di Gangi, R. Cattoni, L. Bentivogli, M. Negri, and M. Turchi, “MuST-C: a Multilingual Speech Translation Corpus,” in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 2012–2017. [Online]. Available: <https://www.aclweb.org/anthology/N19-1202>
- [4] J. Niehues, R. Cattoni, S. Stuker, M. Cettolo, M. Turchi, and M. Federico, “The IWSLT 2018 evaluation campaign,” in *IWSLT 2018*.
- [5] A. C. Kocabiyikoglu, L. Besacier, and O. Kraif, “Augmenting LibriSpeech with French translations: A multimodal corpus for direct speech translation evaluation,” *arXiv preprint arXiv:1802.03142*, 2018.
- [6] C. Nickerson, “English as a lingua franca in international business contexts,” *English for Specific Purposes*, vol. 24, pp. 367–380, 2005.
- [7] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “LibriSpeech: an ASR corpus based on public domain audio books,” in *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2015, pp. 5206–5210.
- [8] W. I. Cho, S. Moon, and Y. Song, “Open Korean corpora: A practical report,” in *Proceedings of Second Workshop for NLP Open Source Software (NLP-OSS)*, 2020, pp. 85–93.
- [9] K. Park, “KSS dataset: Korean single speaker speech dataset,” 2018. [Online]. Available: <https://kaggle.com/bryanpark/korean-single-speaker-speech-dataset>,
- [10] J.-W. Ha, K. Nam, J. Kang, S.-W. Lee, S. Yang, H. Jung, H. Kim, E. Kim, S. Kim, H. A. Kim *et al.*, “ClovaCall: Korean goal-oriented dialog speech corpus for automatic speech recognition of contact centers,” 2020, pp. 409–413.
- [11] J. Park, J.-P. Hong, and J.-W. Cha, “Korean language resources for everyone,” in *Proceedings of the 30th Pacific Asia Conference on Language, Information and Computation: Oral Papers*, Seoul, South Korea, Oct. 2016, pp. 49–58. [Online]. Available: <https://www.aclweb.org/anthology/Y16-2002>
- [12] W. I. Cho, J. Kim, J. Yang, and N. S. Kim, “Towards cross-lingual generalization of translation gender bias,” in *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, ser. FAccT ’21. New York, NY, USA: Association for Computing Machinery, 2021, p. 449–457. [Online]. Available: <https://doi.org/10.1145/3442188.3445907>
- [13] S. Mayhew, K. Bicknell, C. Brust, B. McDowell, W. Monroe, and B. Settles, “Simultaneous translation and paraphrase for language education,” in *Proceedings of the ACL Workshop on Neural Generation and Translation (WNGT)*. ACL, 2020.
- [14] W. I. Cho, S. Moon, J. I. Kim, S. M. Kim, and N. S. Kim, “StyleKQC: A style-variant paraphrase corpus for Korean questions and commands,” *arXiv preprint arXiv:2103.13439*, 2021.
- [15] Y. K. Lee, Y. Jung, I. Lee, J. E. Park, and S. Hahn, “Building a psychological ground truth dataset with empathy and theory-of-mind during the COVID-19 pandemic,” Jun 2021. [Online]. Available: psyarxiv.com/mpn3w
- [16] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, “Bleu: a method for automatic evaluation of machine translation,” in *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*. Philadelphia, Pennsylvania, USA: Association for Computational Linguistics, Jul. 2002, pp. 311–318. [Online]. Available: <https://www.aclweb.org/anthology/P02-1040>
- [17] M. Post, “A call for clarity in reporting BLEU scores,” in *Proceedings of the Third Conference on Machine Translation: Research Papers*. Brussels, Belgium: Association for Computational Linguistics, Oct. 2018, pp. 186–191. [Online]. Available: <https://www.aclweb.org/anthology/W18-6319>
- [18] A. Zhang, “Speech Recognition (version 3.8) [software],” 2017. [Online]. Available: https://github.com/Uberi/speech_recognition#readme
- [19] H. Heo, H. Ko, S. Kim, G. Han, J. Park, and K. Park, “PORORO: Platform Of neuRal mOdelS for natuRal language prOcessing,” <https://github.com/kakaobrain/pororo>, 2021.
- [20] C. Wang, Y. Tang, X. Ma, A. Wu, D. Okhonko, and J. Pino, “fairseq S2T: Fast speech-to-text modeling with fairseq,” in *Proceedings of the 2020 Conference of the Asian Chapter of the Association for Computational Linguistics (ACL): System Demonstrations*, 2020.
- [21] M. Ott, S. Edunov, A. Baevski, A. Fan, S. Gross, N. Ng, D. Grangier, and M. Auli, “fairseq: A fast, extensible toolkit for sequence modeling,” in *Proceedings of NAACL-HLT 2019: Demonstrations*, 2019.
- [22] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 6000–6010.
- [23] J. Pino, Q. Xu, X. Ma, M. J. Dousti, and Y. Tang, “Self-training for end-to-end speech translation,” *Proc. Interspeech 2020*, pp. 1476–1480, 2020.
- [24] Y. Liu, H. Xiong, J. Zhang, Z. He, H. Wu, H. Wang, and C. Zong, “End-to-end speech translation with knowledge distillation,” *Proc. Interspeech 2019*, pp. 1128–1132, 2019.