

ViWOZ: A Multi-Domain Task-Oriented Dialogue Systems Dataset For Low-resource Language

Phi Nguyen Van

phin.v@vnu.edu.vn

Tung Cao Hoang

19020055@vnu.edu.vn

Dung Nguyen Manh

18020370@vnu.edu.vn

Quan Nguyen Minh

19020019@vnu.edu.vn

Long Tran Quoc

tqlong@vnu.edu.vn

Abstract

Most of the current task-oriented dialogue systems (ToD), despite having interesting results, are designed for a handful of languages like Chinese and English. Therefore, their performance in low-resource languages is still a significant problem due to the absence of a standard dataset and evaluation policy. To address this problem, we proposed **ViWOZ**, a fully-annotated Vietnamese task-oriented dialogue dataset. **ViWOZ** is the first multi-turn, multi-domain task-oriented dataset in Vietnamese, a low-resource language. The dataset consists of a total of 5,000 dialogues, including 60,946 fully annotated utterances. Furthermore, we provide a comprehensive benchmark of both modular and end-to-end models in low-resource language scenarios. With those characteristics, the **ViWOZ** dataset enables future studies on creating a multilingual task-oriented dialogue system.

1 Introduction

Task-oriented dialogue systems (ToD), a narrowed version of a general dialogue system, allow users to interact with virtual agents to accomplish certain tasks. In recent years, task-oriented has become an active research topic thanks to the development of neural network models. (Ren et al., 2018; Wu et al., 2020; Henderson et al., 2020). Despite the accelerating movement of performances, those works are limited to hands-on languages like English. (Budzianowski et al., 2018) and Chinese (Zhu et al., 2020) due to the lack of large-scale multilingual languages.

In the field of task-oriented dialogue systems, Natural Language Understanding (NLU) is the only module having diverse language datasets. Multi-ATIS++ (Xu et al., 2020) is a dataset translated from ATIS (Price, 1990), which covers 9 languages but only has a single domain of airline travel, and most of the languages are also from the same language family. Along with MultiATIS++, several

researchers have translated the ATIS into individual languages like Chinese (He et al., 2013), Vietnamese (Nguyen and Nguyen, 2021), and many other languages (Susanto and Lu, 2017; Upadhyay et al., 2018; Xu et al., 2020). Extending the dataset into multi-domain, MTOP (Li et al., 2021) and SID (van der Goot et al., 2021) are both multi-domain and multi-lingual dialogue datasets, but those datasets are still limited to NLU annotations only. Despite having multilingual diversity, simply translating and mapping the dataset will make the dialogue become tedious, lacking locale, and trivial, leading to over-optimistic results. (Ponti et al., 2020; Artetxe et al., 2020).

For multi-lingual Dialogue State Tracking (DST), (Mrkšić et al., 2017b) translated the WoZ 2.0 DST dataset (Wen et al., 2017) into German and Italian. DSTC 5 and DSTC 6 (Kim et al., 2016; Hori et al., 2019) show the benchmark of DST models on zero-shot cross-lingual transfer from English to Chinese and transferring dialogue knowledge from English to Japanese. Finally, DSTC 9 (Gunasekara et al., 2020) introduces the first challenge to benchmark cross-lingual DST systems on large scale datasets, focused on transfer between English and Chinese, using MultiWOZ 2.1 (Eric et al., 2020) as the English dataset and CrossWOZ (Zhu et al., 2020) as the Chinese dataset. Those datasets enable various methods in both modular and end-to-end ToD.

As far as we know, there is still no large-scale, multi-domain dataset for task-oriented dialogue systems in Vietnamese, a low-resource language scenario. Therefore, we proposed **ViWOZ**, built upon MultiWOZ 2.1, as the first fully annotated multi-domain Vietnamese task-oriented dialogue dataset. Our dataset has the following advantages compared to previous works:

- **ViWOZ** is the first fully annotated multi-domain task-oriented dialogue dataset in Vietnamese, a low-resource language.

	MultiWOZ	RiSAWOZ	CrossWoz	BiToD	ViWOZ
Language(s)	EN	ZH	ZH	EN, ZH	VI
No. Dialogues	8,438	10,000	5,012	5,787	4014
No. domains	7	12	5	5	7
No. turns	115,434	134,580	84,692	115,638	48,944
Avg turns	13.46	13.5	16.9	19.98	12.19
Slots	25	159	72	68	25
Values	4,510	4,061	7,871	8,206	4,510

Table 1: Comparisons of training set with other ToD datasets.

- In previous ToD datasets, researchers simply translated the original corpora into the target language, which could make the dataset unrepresentative of realistic conversation due to a lack of culture and locale. Furthermore, existing Machine Translation system such as Google Translate has known issues of making translation errors. In our approach, diverse paraphrases from crowd-sourced workers resolve the translations errors and increase the naturalness of conversations. It is the first fully annotated multi-domain task-oriented dialogue dataset in Vietnamese, a low-resource language.
- **ViWOZ** is built upon MultiWOZ 2.3, providing fully annotated dialogue states and dialogue acts for both the system side and user side. We also corrected the errors in the original data when mapping from English to Vietnamese.
- The effectiveness of multilingual/cross-lingual/monolingual language models on low resource languages, as well as the performance of current approaches, remains an open question. (Razumovskaia et al., 2021). In this work, we conduct intensive experiments in different settings as an initial step to uncover the multilingual ToD problem.

2 Related Work

According to (Budzianowski et al., 2018), there are three categories of task-oriented dialogue datasets: human-to-human, human-to-machine and machine-to-machine. Each category of datasets is associated with whether the user and system agent are machine or human. Regardless of dataset category, most of them are single-domain or only available in high-resource languages or both. As far as

our limited knowledge, there is still no large-scale multi-domain dataset for low-resource languages.

Human-to-Human is a type of dataset where both the user and the system are human agents. The crowd-sourced workers are hired to talk to each other and are given some instructions. This setup is widely used in both single-domain (Kim et al., 2016) and large-scale multi-domain scenarios (Eric et al., 2020). The dataset built upon this setup creates a diverse and realistic dialogue, but it also requires intensive human effort in data creation and quality control.

Human-to-Machine is a type of dataset where humans interact with a dialogue system. This setup is already used to create the dataset (Hori et al., 2019) from The Dialogue State Tracking Challenges (DSTC). Despite the automation on the system side, the quality of those datasets is heavily influenced by the quality of the dialogue system.

Machine-to-Machine is a type of dataset where both the user and the system are autonomous agents. This fully automated setup enables the creation of large scale datasets with minimal human effort (Shah et al., 2018; Rastogi et al., 2019). However, these datasets still lack lingual diversity.

Multilingual ToD dataset Despite the efforts of several researchers (Mrkšić et al., 2017b; Zhu et al., 2020; Lin et al., 2021), one of the major barriers to the widespread application of ToD research (Razumovskaia et al., 2021) is a lack of multilingual ToD datasets (Razumovskaia et al., 2021). Most of the multilingual datasets only focus on single-domain or only on a handful of languages like English and Chinese.

3 ViWOZ Dataset

3.1 Language analysis

Linguists have conducted researches about the difference of Vietnamese and English. The majority of linguistics researchers consider Vietnamese a

	Train	Dev	Test	Single-domain	Multi-domain	All
Dialogues	4,014	493	493	2,104	2,896	5000
Turns	48,944	5,936	6,066	17,160	43,786	60,946
Tokens	897,550	109,453	110,897	286,684	831,214	1,117,900
Vocabulary	8,994	3,072	2,908	4,563	8,538	10,136
Avg. turns	12.19	12.04	12.30	8.15	15.12	12.19
Avg. acts	1.24	1.23	1.23	1.83	1.26	1.28
Avg. user-acts	1.52	1.52	1.51	1.39	1.57	1.64
Avg. system-acts	1.64	1.65	1.63	1.58	1.66	1.52

Table 2: ViWOZ dataset statistic. White space tokenization is used to count number of token and vocabulary.

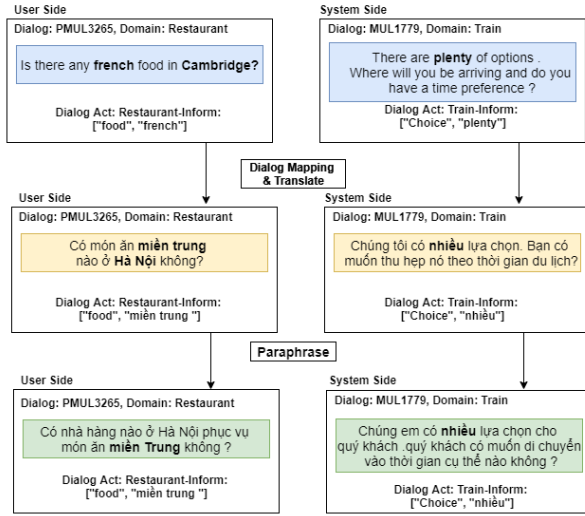


Figure 1: Illustration of data construction for user side (dialogue PMUL3265) and system side (dialogue MUL1779)

Austroasiatic language (Alves, 2006), while (Giang, 2007) has identified a variety of differences between Vietnamese and English, such as lexical, semantic, and grammatical characteristics. For example, when describing an object, a Vietnamese would use the word order $N + Adj$, while an English speaker would use $Adj + N$. Another example is the diversity of pronouns in Vietnamese: a single pronoun I representing the speaker could be translated to *Tôi, tao, mình, tớ, anh, chị, em, cô, dì, chú, bác, etc.*, each has a different level of formality, and some could only be used in certain contexts where the speaker knows exactly who he/she is having a conversation with. These linguistic differences, together with the cultural and language usage variations, could be a significant problem for machine translation systems, as grammar or appropriate language usage may not be fully reflected. Figure 2 demonstrates this problem: In the first example, different mentions referring to



Figure 2: Translation of 2 utterances from dialogue PMUL1043 in MultiWOZ 2.3

" I " were used, which is very uncommon in Vietnamese; in the second example, the translation of the adverb "*please*" remains at the end of the utterance, while it is more common to have it in the beginning, i.e. "*Xin vui lòng cho tôi địa chỉ*"

3.2 ViWOZ

Our dataset is built upon MultiWOZ 2.3 (Han et al., 2020) by translating the dialogues from English to Vietnamese. We also performed manual post-processing to modify the translated data for translation accuracy. A Vietnamese-localized database is also constructed together with the dataset for cultural appropriation. The process of data creation is described in this section.

3.2.1 Ontology construction

Task-oriented dialogue systems require a source of truth, also known as an *ontology*, to provide information to the users. MultiWOZ 2.3's ontology consists of a total of 5,946 entities spanning 6 domains (attraction, bus, hospital, hotel, police, restaurant, train). For exact location domain where each entity corresponds to a location in real life, we gathered the same list of locations in Hanoi by crawling from several local online booking and

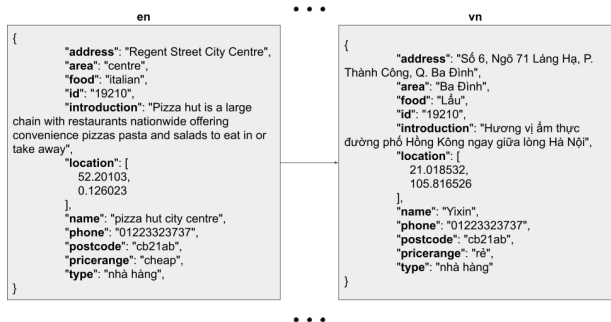


Figure 3: Example of a mapping in the *restaurant* domain of the constructed ontology

travel agent websites (i.e. traveloka, booking.com, etc.). For the remaining train and bus domains where each item contains a travel route information from a city in England to the main city of Cambridge, we manually map each city in the source data to a real province in Vietnam. Cambridge is mapped to Hanoi, the capital city of Vietnam. All the available values for each item in the ontology are translated to Vietnamese properly. For example, *duration: 51 minutes* is converted into *duration: 51 phút*. The result of this process is a localized ontology in Vietnamese where each entity is uniquely mapped from English.

3.2.2 Slot Mapping

Previous datasets like MultiATIS++ simply kept entities intact when translating datasets. This will lead to unnatural dialogues that are unrepresentative of real-life conversations in Vietnamese. In our approach, we use the constructed localized ontology to perform slot mapping based on the utterances and the annotated dialogue acts. Since multiple dialogue act values can represent the same value in the ontology (i.e. *center* and *centre* could both mean the center area of Cambridge), we build a dictionary to map annotated dialogue slots from the original values to new values in Vietnamese.

Manually selecting tens of thousands of values to be mapped in the dictionary can be very time-consuming and inefficient. Therefore, we first build the mapping dictionary by performing a full-text search on the ontology slots. For example, the value in *Restaurant – Name – Inform* slot of an utterance will be searched in every *name* slots in the ontology’s *restaurant* domain, the entity with the lowest Levenshtein distance will be selected, and the *name* slot of the corresponding Vietnamese entity is considered the result of

the correct mapping value. The dictionary is then cross-checked and modified manually by two researchers to mitigate the problems of the automated searching algorithm.

Using the generated dictionary, we replaced all the annotated dialogue slots in the utterance and in the annotation with their new value in Vietnamese. The resultant dialogues are in English, with the annotated spans being replaced by the localized entity from the previously constructed dictionary. We apply the mapping process to 5,000 dialogues, randomly sampled from MultiWOZ 2.3. In summary, we have mapped and checked a total of 75,852 entities, 60,946 dialogue acts, and 30,475 dialogue states.

3.2.3 Dialogue translation and paraphrasing

We utilized the Google Translate API to translate every dialogue, with the entities having been replaced by localized ones, from English to Vietnamese. Twelve native speakers are then selected to paraphrase all the translated utterances. Given a translated utterance, the annotators must rewrite the sentence so that the information about the dialogue acts is kept and the language is as natural as possible. Because the original translated value may not be compatible with the paraphrased utterance, the crowd-sourced workers are also required to re-annotate every slot value to rebuild the dialogue state. This ensures that the dialogue is written naturally, as a simple translation of the dataset creates over-simplistic and unrealistic dialogues. The examples in Figure 1 demonstrate how the paraphrasing process improves the naturalness of the dialogue. The final dataset consists of mapped slots, states, dialogue acts, dialogue goals, and paraphrased utterances for both system and user utterances. Table 2 shows the basic statistics of ViWOZ dataset.

3.2.4 Human Evaluation

Multiwoz is a noisy dataset (Zang et al., 2020; Ye et al., 2021), and after many attempts to rectify the dialogue errors, we still find errors from the original dataset while converting the dialogues from English to Vietnamese. Missing slot value and incorrectly spanned annotation are the two most common errors, which are corrected throughout the annotation of crowd-sourced workers.

Finally, two researchers perform quality assurance by manually reviewing the paraphrased utterances to identify possible syntactic, semantic, and

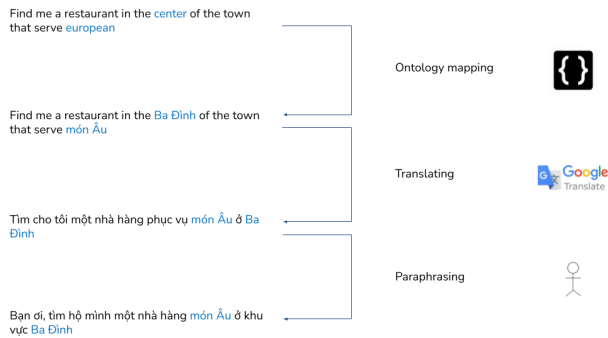


Figure 4: The ontology improves dialogue localization, while the paraphrasing acts as human post-processing to improve dialogue naturalness

slot annotation errors

4 Benchmark and Analysis

In ToD, there are two main approaches: modular and end-to-end. While the modular approach aims to create ToD by breaking down the system into specialized modules like Natural Language Understanding (NLU), Dialogue State Tracking (DST), Dialogue Policy, and Natural Language Generation (NLG), end-to-end models often use large language models to directly map the dialogue history to output utterance. With **ViWOZ**, we can do multiple assessments of those approaches in low resource scenarios. Our main concern is whether current methods, which are optimized for high resource languages like English and Chinese, still work on low-resource language scenarios, especially on the new **ViWOZ** dataset. We provided benchmarks on various models and settings for two sub-tasks of modular ToD systems: NLU and DST. In addition, we also conduct experiments on end-to-end models.

- Performance difference between monolingual, bilingual, and multilingual pre-trained models. This helps us to understand the effectiveness of the monolingual language model in low-resource language ToD.
- Effectiveness of model size to performance: whether bigger language models has any advantage on low-resource language ToD.
- Zero-shot/bilingual training from cross-lingual pre-trained models: whether high-resource language datasets help boost the performance of low resource ToD.

We use the following pre-trained models to conduct the experiments: (1) PhoBERT (Nguyen and Tuan Nguyen, 2020) is a monolingual language model trained only on Vietnamese corpus, (2) EnViBERT (Nguyen et al., 2020) is a bilingual RoBERTa model train on English and Vietnamese, (3) XLM-RoBERTa (Conneau et al., 2020) is a multilingual model trained on 100 languages.

For the dataset, we use three settings:

- **Zero shot.** To see the effectiveness of multilingual LM in a zero-shot setting, we train the models on the original English dataset and evaluate them on a Vietnamese test set.
- **Monolingual - VN/EN.** In this setting, we train and evaluate models only on Vietnamese data. The EN column is added to Table 3 for convenient comparison between the performance of models in English and Vietnamese.
- **Bilingual - VN+EN.** To investigate knowledge transfer from a high-resource to a low-resource language, we combine the English and Vietnamese training sets into a single training set and test the models on a Vietnamese test set.

4.1 Natural Language Understanding

NLU is the first module in modular ToD systems. The main purpose of the module is to identify the intent and extract information from user and system utterances. It often consists of two sub-tasks: intent classification and slot filling. There are two main approaches for this module: joint-model and separate model, where the joint model tries to solve two sub-tasks at once, while the separate model solves each task with individual models. Joint-models often have the potential to create a more compact model, and multi-task training objectives also offer better performance (Razumovskaia et al., 2021). In cross-lingual NLU, the lack of multilingual data makes zero-shot or few-shot learning the default approaches. By utilizing the massively multilingual Transformers models, the zero-shot model shows strong performance on multilingual NLU datasets. (Zhang et al., 2019; Krishnan et al., 2021).

Experiment settings: In this section, we extend BERT NLU (Devlin et al., 2019) by replacing BERT with other alternative pre-trained language models. All model configurations (such as the classification head) are the same throughout the experiments, while the language models act

Task	Model	Pretrained	Metrics	en	zero shot	vn	vn+en
NLU	Joint model	PhoBERT	Dialogue act F1	-	-	84.70	-
		EnViBERT		88.79	44.02	85.04	85.49
		XLM-base		89.50	45.89	84.14	84.46
DST	TRADE		Joint goal accuracy	48.62	-	46.61	44.63
	SUMBT	PhoBERT		-	-	52.95	-
		EnViBERT		61.86*	6.17	54.70	52.09
		XLM-base	61.86*	6.20	52.00	52.98	
End-to-end	MinTL	mT5-small	Inform	80.04	35.5	38.1	43.2
			Success	72.71	21.3	17.8	20.3
			BLEU	19.11	0.0	18.1	16.5
		mT5-base	Inform	82.15	34.7	46.9	51.7
			Success	74.44	18.1	35.7	32.7
			BLEU	18.59	0.0	17.1	17.2

Table 3: Benchmark of a few recent methods on NLU, DST and end-to-end. — cells indicate experiments are not conducted because models do not have the vocabulary of the language. *: result trained on BERT-base from (Ye et al., 2021).

as the utterance encoders in a sequence-tagging and multi-label classification joint-training task. By comparing different language models in different training schemes, we can show the effectiveness of pre-trained language models in Natural Language Understanding. We conduct experiments with PhoBERT (Nguyen and Tuan Nguyen, 2020), EnViBERT (Nguyen et al., 2020) and XLM-RoBERTa (Conneau et al., 2020) for monolingual, bilingual and zero-shot settings, respectively.

Result analysis: The result in the first portion of Table 3 indicates that the bilingual and multilingual models yield a noticeable zero-shot dialogue act F1 score. There is a gap in the monolingual scheme result between ViWOZ and MultiWOZ, since ViWOZ has about half the training samples in comparison to MultiWOZ. In addition, there is little to no difference in the results when training in monolingual data (where the model is trained and validated on Vietnamese only) and in bilingual data (where the model is trained with data from both languages and evaluated on a Vietnamese test set). When taking the model size into account, despite having only about half the number of parameters as XLM-RoBERTa (base), EnViBERT and PhoBERT still achieve competitive results, if not better, suggesting that a well-trained monolingual or bilingual language model could outweigh a larger multi-lingual one.

4.2 Dialogue State Tracking

Dialogue State Tracking (DST) takes responsibility for maintaining a dialogue belief state and contains information about the dialogue throughout the conversations (Mrkšić et al., 2017a). Traditionally, DST often takes implicit input from NLU and previous belief states to produce an updated state of the current utterance. Recently, there has been a trend to utilize long-distance representation of Transformers models to represent and track the whole dialogue.

Experiment setting TRADE (Transferable Dialogue State Generator) (Wu et al., 2019) and SUMBT (Slot-Utterance Matching for Universal and Scalable Belief Tracking) (Lee et al., 2019) are the two models we used to perform experiments in this part. TRADE generates conversation state using an encoder-decoder design, whereas SUMBT tracks dialogue state using contextual representation from a pre-trained language model. We can demonstrate the function of a pre-trained language model in tracking conversation status by comparing two models with different dataset configurations.

Result analysis The second portion of Table 3, shows that SUMBT outperforms TRADE across all datasets, indicating the effectiveness of pre-training LM. Mixed Vietnamese and English data hindered model performance, indicating that a mix-language training method produces noisy data and makes it more difficult for models to monitor crucial information while maintaining multilinguality.

4.3 End-to-End

Modular ToD trains and combines each model individually. This pipeline approach leads to error cascading when the performance of later modules depends largely on the performance of previous modules. End-to-end models try to solve this problem by combining modules into a single LM to generate a system response directly from user utterances.

Experiment setting The state-of-the-art end-to-end task-oriented dialogue models MinTL (Lin et al., 2020) are used. We just change the pre-trained LM from T5 (Raffel et al., 2020) to mT5 (Xue et al., 2021) and leave the rest of the implementation unchanged. From MinTL’s official code, the hyper-parameters are likewise set to default. All models are trained on a single NVIDIA A6000.

Result analysis: The MinTL failed on all dataset experiment settings; the inform and success scores are significantly lower than the English-only setup, although the BLEU score is similar; the model appears to create reasonable responses but fails to convey accurate information. Furthermore, BLEU scores are similar in monolingual and bilingual setups, showing that there is no improvement when high resource language data is included; this can be explained by the language family difference between English and Vietnamese. The zero shot option, on the other hand, produced results that were equivalent to the monolingual data setting in terms of inform and success scores, demonstrating that multilingual pre-training is beneficial when transferring dialogue structure. When we use mT5-base to increase the model size, we observe that performance on providing correct information improves dramatically from 38.1% to 46.9% on monolingual setup and from 43.2% to 51.7% on bilingual setup, indicating that larger LMs tend to perform better on downstream tasks.

5 Conclusion

This study presents ViWOZ, the first multi-domain Vietnamese task-oriented dataset. The dataset contains 5,000 dialogues and 60,946 utterances, all of which are fully annotated on both the system and user sides. We also do extensive experiments on a number of settings, from zero shot through mixed-language training on a variety of architectures and pre-trained language models. Our findings show that doing task-oriented dialogue system in a low-resource language context is

still difficult, especially end-to-end manner. The findings and analysis further show that, on low-resource scenario, while multilingual and bilingual LMs cannot achieve similar performance to that of models trained on high-resource language, multilinguality still has a positive impact on model’s outcomes when compared to monolingual models, especially when pre-trained LMs are not available.

Limitations and future work. The fundamental limitation of our work is that the quantity of dialogues is still lower than in English, but this also corresponds to a real-world problem in which low-resource language data is not easily accessible and it is costly to annotate. In the future, we intend to expand the dataset by including side information such as co-reference resolution data, as well as conduct more experiments using machine translation algorithms for zero-shot and few-shot scenarios. Despite these limitations, we believe the ViWOZ is an essential step in the development of multilingual task-oriented dialogue systems.

References

- Mark Alves. 2006. [Linguistic research on the origins of the vietnamese language: An overview](#). *Journal of Vietnamese Studies*, 1:104–130.
- Mikel Artetxe, Gorra Labaka, and Eneko Agirre. 2020. [Translation artifacts in cross-lingual transfer learning](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7674–7684, Online. Association for Computational Linguistics.
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. [MultiWOZ - a large-scale multi-domain Wizard-of-Oz dataset for task-oriented dialogue modelling](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026, Brussels, Belgium. Association for Computational Linguistics.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). 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)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi, Sanchit Agarwal, Shuyang Gao, Adarsh Kumar, Anuj Goyal, Peter Ku, and Dilek Hakkani-Tur. 2020. **MultiWOZ 2.1: A consolidated multi-domain dialogue dataset with state corrections and state tracking baselines**. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 422–428, Marseille, France. European Language Resources Association.
- Tang Giang. 2007. **Cross-linguistic analysis of vietnamese and english with implications for vietnamese language acquisition and maintenance in the united states**. *Journal of Southeast Asian American Education and Advancement*, 2.
- Chulaka Gunasekara, Seokhwan Kim, Luis Fernando D’Haro, Abhinav Rastogi, Yun-Nung Chen, Mihail Eric, Behnam Hedayatnia, Karthik Gopalakrishnan, Yang Liu, Chao-Wei Huang, et al. 2020. Overview of the ninth dialog system technology challenge: Dstc9. *arXiv preprint arXiv:2011.06486*.
- Ting Han, Ximing Liu, Ryuichi Takanobu, Yixin Lian, Chongxuan Huang, Dazhen Wan, Wei Peng, and Minlie Huang. 2020. **Multiwoz 2.3: A multi-domain task-oriented dialogue dataset enhanced with annotation corrections and co-reference annotation**. *arXiv preprint arXiv:2010.05594*.
- Xiaodong He, Li Deng, Dilek Hakkani-Tur, and Gokhan Tur. 2013. **Multi-style adaptive training for robust cross-lingual spoken language understanding**. In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 8342–8346.
- Matthew Henderson, Iñigo Casanueva, Nikola Mrkšić, Pei-Hao Su, Tsung-Hsien Wen, and Ivan Vulić. 2020. **ConveRT: Efficient and accurate conversational representations from transformers**. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2161–2174, Online. Association for Computational Linguistics.
- Chiori Hori, Julien Perez, Ryuichiro Higashinaka, Takaaki Hori, Y-Lan Boureau, Michimasa Inaba, Yuiko Tsunomori, Tetsuro Takahashi, Koichiro Yoshino, and Seokhwan Kim. 2019. Overview of the sixth dialog system technology challenge: Dstc6. *Computer Speech & Language*, 55:1–25.
- Seokhwan Kim, Luis Fernando D’Haro, Rafael E Banchs, Jason D Williams, Matthew Henderson, and Koichiro Yoshino. 2016. The fifth dialog state tracking challenge. In *2016 IEEE Spoken Language Technology Workshop (SLT)*, pages 511–517. IEEE.
- Jitin Krishnan, Antonios Anastasopoulos, Hemant Purohit, and Huzefa Rangwala. 2021. **Multilingual code-switching for zero-shot cross-lingual intent prediction and slot filling**. In *Proceedings of the 1st Workshop on Multilingual Representation Learning*, pages 211–223, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Hwaran Lee, Jinsik Lee, and Tae-Yoon Kim. 2019. **SUMBT: Slot-utterance matching for universal and scalable belief tracking**. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5478–5483, Florence, Italy. Association for Computational Linguistics.
- Haoran Li, Abhinav Arora, Shuohui Chen, Anchit Gupta, Sonal Gupta, and Yashar Mehdad. 2021. **MTOP: A comprehensive multilingual task-oriented semantic parsing benchmark**. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2950–2962, Online. Association for Computational Linguistics.
- Zhaojiang Lin, Andrea Madotto, Genta Indra Winata, and Pascale Fung. 2020. **MinTL: Minimalist transfer learning for task-oriented dialogue systems**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3391–3405, Online. Association for Computational Linguistics.
- Zhaojiang Lin, Andrea Madotto, Genta Indra Winata, Peng Xu, Feijun Jiang, Yuxiang Hu, Chen Shi, and Pascale Fung. 2021. **Bitod: A bilingual multi-domain dataset for task-oriented dialogue modeling**. *arXiv preprint arXiv:2106.02787*.
- Nikola Mrkšić, Diarmuid Ó Séaghdha, Tsung-Hsien Wen, Blaise Thomson, and Steve Young. 2017a. **Neural belief tracker: Data-driven dialogue state tracking**. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1777–1788, Vancouver, Canada. Association for Computational Linguistics.
- Nikola Mrkšić, Ivan Vulić, Diarmuid Ó Séaghdha, Ira Leviant, Roi Reichart, Milica Gašić, Anna Korhonen, and Steve Young. 2017b. **Semantic specialization of distributional word vector spaces using monolingual and cross-lingual constraints**. *Transactions of the Association for Computational Linguistics*, 5:309–324.
- Dat Quoc Nguyen and Anh Tuan Nguyen. 2020. **PhoBERT: Pre-trained language models for Vietnamese**. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1037–1042, Online. Association for Computational Linguistics.
- Linh The Nguyen and Dat Quoc Nguyen. 2021. **PhoNLP: A joint multi-task learning model for Viet-**

- namese part-of-speech tagging, named entity recognition and dependency parsing. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Demonstrations*, pages 1–7, Online. Association for Computational Linguistics.
- Thai Binh Nguyen, Quang Minh Nguyen, Thi Thu Hien Nguyen, Quoc Truong Do, and Chi Mai Luong. 2020. [Improving Vietnamese Named Entity Recognition from Speech Using Word Capitalization and Punctuation Recovery Models](#). In *Proc. Interspeech 2020*, pages 4263–4267.
- Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Korhonen. 2020. [XCOFA: A multilingual dataset for causal commonsense reasoning](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2362–2376, Online. Association for Computational Linguistics.
- P. J. Price. 1990. [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*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *Journal of Machine Learning Research*, 21(140):1–67.
- Pushpendre Rastogi, Arpit Gupta, Tongfei Chen, and Mathias Lambert. 2019. [Scaling multi-domain dialogue state tracking via query reformulation](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Industry Papers)*, pages 97–105, Minneapolis, Minnesota. Association for Computational Linguistics.
- Evgeniia Razumovskaia, Goran Glavaš, Olga Majewska, Anna Korhonen, and Ivan Vulić. 2021. [Crossing the conversational chasm: A primer on multilingual task-oriented dialogue systems](#). *arXiv preprint arXiv:2104.08570*.
- Liliang Ren, Kaige Xie, Lu Chen, and Kai Yu. 2018. [Towards universal dialogue state tracking](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2780–2786, Brussels, Belgium. Association for Computational Linguistics.
- Pararth Shah, Dilek Hakkani-Tür, Gokhan Tür, Abhinav Rastogi, Ankur Bapna, Neha Nayak, and Larry Heck. 2018. [Building a conversational agent overnight with dialogue self-play](#). *arXiv preprint arXiv:1801.04871*.
- Raymond Hendy Susanto and Wei Lu. 2017. [Neural architectures for multilingual semantic parsing](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 38–44, Vancouver, Canada. Association for Computational Linguistics.
- Shyam Upadhyay, Manaal Faruqui, Gokhan Tür, Hakkani-Tür Dilek, and Larry Heck. 2018. [\(almost\) zero-shot cross-lingual spoken language understanding](#). In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6034–6038.
- Rob van der Goot, Ahmet Üstün, and Barbara Plank. 2021. [On the effectiveness of dataset embeddings in mono-lingual, multi-lingual and zero-shot conditions](#). In *Proceedings of the Second Workshop on Domain Adaptation for NLP*, pages 183–194, Kyiv, Ukraine. Association for Computational Linguistics.
- Tsung-Hsien Wen, David Vandyke, Nikola Mrkšić, Milica Gašić, Lina M. Rojas-Barahona, Pei-Hao Su, Stefan Ultes, and Steve Young. 2017. [A network-based end-to-end trainable task-oriented dialogue system](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 438–449, Valencia, Spain. Association for Computational Linguistics.
- Chien-Sheng Wu, Steven C.H. Hoi, Richard Socher, and Caiming Xiong. 2020. [TOD-BERT: Pre-trained natural language understanding for task-oriented dialogue](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 917–929, Online. Association for Computational Linguistics.
- Chien-Sheng Wu, Andrea Madotto, Ehsan Hosseini-Asl, Caiming Xiong, Richard Socher, and Pascale Fung. 2019. [Transferable multi-domain state generator for task-oriented dialogue systems](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 808–819, Florence, Italy. Association for Computational Linguistics.
- Weijia Xu, Batool Haider, and Saab Mansour. 2020. [End-to-end slot alignment and recognition for cross-lingual NLU](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5052–5063, Online. Association for Computational Linguistics.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. [mT5: A massively multilingual pre-trained text-to-text transformer](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

- Fanghua Ye, Jarana Manotumruksa, and Emine Yilmaz. 2021. Multiwoz 2.4: A multi-domain task-oriented dialogue dataset with essential annotation corrections to improve state tracking evaluation. *arXiv preprint arXiv:2104.00773*.
- Xiaoxue Zang, Abhinav Rastogi, Srinivas Sunkara, Raghav Gupta, Jianguo Zhang, and Jindong Chen. 2020. Multiwoz 2.2: A dialogue dataset with additional annotation corrections and state tracking baselines. *arXiv preprint arXiv:2007.12720*.
- Zhichang Zhang, Zhenwen Zhang, Haoyuan Chen, and Zhiman Zhang. 2019. A joint learning framework with bert for spoken language understanding. *IEEE Access*, 7:168849–168858.
- Qi Zhu, Kaili Huang, Zheng Zhang, Xiaoyan Zhu, and Minlie Huang. 2020. [CrossWOZ: A large-scale Chinese cross-domain task-oriented dialogue dataset](#). *Transactions of the Association for Computational Linguistics*, 8:281–295.