

Leveraging BERT for Natural Language Understanding of Domain-Specific Knowledge

Vasile Ionut Iga

Business Informatics Research Center

Babes-Bolyai University

Cluj-Napoca, Romania

vasile.iga@stud.ubbcluj.ro

Gheorghe Cosmin Silaghi

Business Informatics Research Center

Babes-Bolyai University

Cluj-Napoca, Romania

Gheorghe.Silaghi@ubbcluj.ro

Abstract—Natural Language Understanding (NLU) is a core task when building conversational agents, fulfilling the objectives of understanding the user’s goal and detecting any valuable information regarding it. NLU implies Intent Detection and Slot Filling, to semantically parse the user’s utterance. One caveat when training a Deep Learning model for domain-specific NLU is the lack of specific datasets, which leads to poorly performing models. To overcome this, we experiment with fine-tuning BERT to jointly detect the user’s intent and the related slots, using a custom-generated dataset built around an organization-specific knowledge base. Our results show that well-constructed datasets lead to high detection performances and the resulting model has the potential to enhance a future task-oriented dialogue system.

Index Terms—natural language understanding, intent detection, slot filling, BERT, task oriented dialogue system

I. INTRODUCTION

Artificial Intelligence has become a hot topic domain in recent years, and the ability to understand natural language is a true requirement for such systems [1]. Conversational agents try to implement this ability, in the form of chatbots or task-oriented dialogue systems. Nonetheless, each agent solves the task of understanding the user’s utterance by enhancing the Natural Language Understanding module. NLU covers Intent Detection (ID) and Slot Filling (SF); the former deals with recognizing the user’s goal from a sentence, while the latter identifies key information accompanying it. Treated as two separate goals, recent literature presents neural networks that jointly solve them, as they are highly correlated and tied [2]. Recent advances [3]–[7] make use of transfer learning to construct ID and SF classifiers starting from pre-trained models such as BERT [8], T5 [9], or GPT-3 [10].

Task-oriented dialogue (TOD) systems help users to solve particular tasks related to a specific context [11], [12]. That

means tailoring all constructed models for the domain knowledge describing that context. Our long-term goal is to construct a TOD system for automating tasks within a business organization characterized by a knowledge base (KB), persisted in a machine-readable format. Therefore, we can leverage this knowledge base to produce the requested dataset for training the NLU component of a future TOD system. Initially, our TOD system [13] provides support with the essential Create-Retrieve-Update-Delete (CRUD) procedures for the management of the organizational KB. Therefore, the NLU module needs to identify the intents related to those procedures and the associated slots, strongly tied with the specific instances stored within the organizational KB.

Our work focuses on fine-tuning a BERT instance to jointly detect intent and associated slots, trained using a custom-generated dataset for recognizing domain-specific actions. With this approach, we can develop a TOD system that is able to understand the needs of a company, without giving away private data to third-party dialogue systems (such as ChatGPT), and comply with GDPR regulations. We assess the quality of the simple architecture built on top of BERT by first separately fine-tuning it on two popular datasets, ATIS [14] and SNIPS [15], and compare results with other state-of-the-art BERT models for NLU. The custom-generated dataset was built using a dialogue simulator [13], that exploits a given organizational knowledge base and generates dialogues between two machines, simulating human-like discussions regarding the management of the KB. Results show that minimal adjustments have to be done to pre-trained models like BERT, if the given dataset is properly fitting the desired domain, leading to high model performance.

The paper is structured as follows. Section II presents related work that influenced our experimental research. Section III shows our methodology, including the dataset, model’s architecture, performance evaluation procedure and metrics. Section IV discusses the obtained results and section V concludes the paper.

II. RELATED WORK

ID can be thought of as a sentence classification problem, while SF is considered a sequence labeling task. Let’s analyze the user’s utterance “Insert a project with code as 123 and

The present work received financial support through the project: *Integrated system for automating business processes using artificial intelligence*, POC/163/1/3/121075 - a project co-financed by the European Regional Development Fund (ERDF) through the Competitiveness Operational Programme 2014-2020.

V.I. Iga thanks Babes-Bolyai University for granting him the special scholarship for scientific activity during the academic year 2022-2023.

Computational resources were provided by the high-performance computational facility of the Babes-Bolyai University (projects 48801/1862 and POC/398/1/1/124155) co-financed by the European Regional Development Fund of the European Union.

class as Python”. A NLU model would consider “insert” as the intention, while “code = 123” and “class = Python” are the related slots. Usually, all the detectible intents and slots are defined before constructing the classification models, and they are fit for the target domain of the NLP problem under investigation.

Initially, ID and SF were not solved together. Ravuri and Stolcke [16] used RNN and LSTM models to detect user’s intention, while Mesnil et al. [17] implemented and compared a variety of RNN architectures, including Elman and Jordan type, to detect slots. Later, the two tasks were combined to exploit their connections. Hakkani-Tur et al. [18] leveraged the power of combining bidirectional RNN with LSTM networks, to solve both tasks at the same time. Another interesting approach was introduced by Zhang et al. [19], who introduces the Graph LSTM to tackle the shortcomings of sequential models and exploit the correlation between slots and intents.

Recently, pre-trained DL models helped researchers to achieve new state-of-the-art results, by using the concept of transfer learning. This means one does not need to start from a vanilla model but can inherit the already learned knowledge and fine-tune it for its desired task. Examples of popular pre-trained models are BERT [8], T5 [9], or GPT-3 [10]. Other studied paradigms are cross-lingual [20], or low-resource NLU [21], for transferring knowledge to low-resource languages.

BERT [8] is a popular pre-trained model among researchers, therefore a large number of NLU versions were fine-tuned. Chen et al. [3] only add a softmax layer on top of BERT to predict intent and slots, but prove that it is enough to achieve great results on datasets such as ATIS or SNIPS. Castellucci et al. [4] attach two layers on top of BERT, one for intent detection and the other for slot filling, but design a single loss function for both tasks. Also, they extend their research by generating an Italian dataset and prove that training a model over multiple languages, even with noticeable syntax differences, can increase the overall performance of a model.

Zhang et al. [5] propose an encoder-decoder framework, leveraging BERT as an encoder and implementing a two-stage decoding process for ID and SF. An interesting mechanism is the insertion of the intent encoding into the encodings of each associated slot, to increase the detection of intent-specific slots and leave out the unimportant ones.

Krone et al. [6] propose a few-shot joint learning task, to improve the performance on labels not seen at training time. They prove that their algorithm, together with a pre-trained network such as BERT, is complementary and yields further gains. Qin et al. [7] modify the standard attention mechanism from vanilla transformers with a co-interactive module that considers the bidirectional connection between intents and slots, rather than the classical approach which focuses only on the intents-slots direction. Additionally, they connect their approach with BERT, proving that pre-trained models can increase the final performances of a system.

III. METHODOLOGY

In this section, we introduce the architecture of our model, the datasets, and the training and testing procedures. The code¹ was written in Python 3.10 and computational resources were provided by high-performance computational facility of the Babeş-Bolyai University [22]. The tokenizer and the base BERT transformer are loaded using the Hugging Face libraries, while the joint NLU component is described using the Tensorflow library. We followed the suggestions of Chen et al. [3] and Shawon Ashraf², who implemented their models using Tensorflow.

A. Datasets

The performance of a deep learning model strongly depends on the quality of the datasets used for training. Therefore, to increase the model’s capacity of recognizing entities of interest, one needs to use specific data, well covering the search space of the target problem. In general, public datasets are tailored for solving general problems, but for specific tasks, such as the one introduced by this paper, those datasets are of little help.

To accomplish our needs, i.e. creating a dataset with conversations related to our specific concepts and instances within our organization, we designed a Machine-to-Machine (M2M) conversation generator with the overall architecture presented in fig. 1. The system consists of a rule-based TOD system fully described in [13], an user simulator and a prompt.

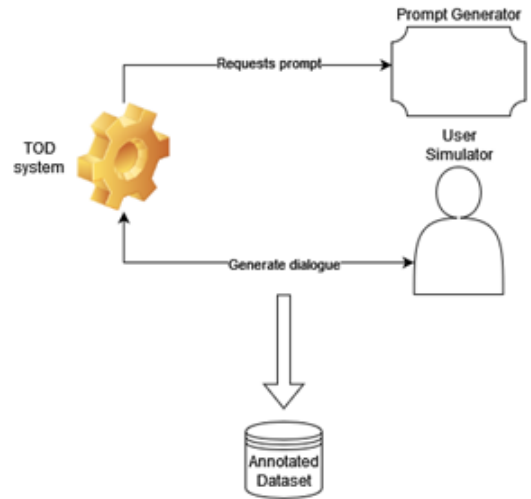


Fig. 1. M2M conversation generator system

The system generates conversations about specific concepts from a provided ontology, growing an organizational knowledge base (KB) related to the concepts supplied in the ontology. The final goal of each conversation is to perform CRUD (Create-Retrieve-Update-Delete) operations on the organizational KB, thus collecting important information

¹Code and data are available at <https://github.com/IonutIga/Domain-Specific-NLU-BERT>

²<https://github.com/ShawonAshraf/nlu-jointbert-dl2021>

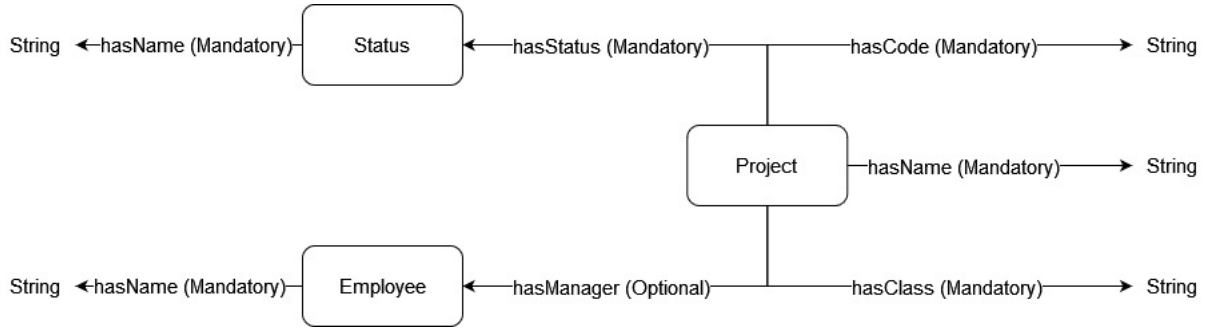


Fig. 2. The input ontology

regarding the target organization, meantime generating the properly annotated dataset for the NLP tasks, in our case, intent detection and slot filling.

The starting ontology is presented in Fig. 2. It describes three concepts (Project, Status, and Employee) and the relationships between them or other literal values such as strings or integers. Each relationship represents a parameter assigned to a concept. In each conversation the user simulator asks the TOD system to perform one or more CRUD operations over the organizational KB and if the requested operation is confirmed by the system, it is persisted as valid in the organizational KB. Each user utterance could request one of the 13 intents, with a total of 27 slots. User intents recognized by the TOD system helps fulfilling specific CRUD procedures on the organizational KB (like *insert*, *select*, *delete* or *update*), discard the value of a parameter (*remove*) confirm or reject a given response (*agree* and *disagree*) or maintain the conversation flow (intents *hello*, *goodbye* and *thank*). The TOD system responds with 27 actions, such that to maintain a natural conversation flow. All user intents and TOD system actions are described in detail in [13]).

Generated data is presented in the JSON format depicted in Fig.3. Each phrase has a unique ID and the key elements stored are the user utterance, the intent of the phrase, important slots, and their positions.

```

{
  'utterance ID': {
    'text': 'user utterance',
    'slots': {'slot name': 'slot value' etc.}
    'positions': {'slot name': [start index, end index], etc.}
    'intent': 'detected intent'
  }...
}

```

Fig. 3. Generated dataset format

We can generate datasets with a variable number of conversations, thus, we generated multiple datasets with a number of conversations ranging from 625 to 5000, equivalent to 2500 to 17500 user utterances (phrases) per dataset, letting us select

the best dataset alternative considering the tradeoff between training performance and training time.

Table I presents the statistics of the generated datasets (GD) used for model training and performance assessment. In each case, the validation procedure used for model fine-tuning was the hold-out procedure with 20% of the training data kept for validation. For the testing dataset we selected about 1000 conversations or a number of conversations no more than 25% of the training data.

TABLE I
STATISTICS REGARDING THE GENERATED DATASETS

Statistics / Dataset	GD1	GD2	GD3	GD4
No of conversations	625	1250	3250	5000
No of user utterances (phrases)	2458	4677	9750	17515
Training dataset (no. phrases)	2000	3800	8750	16500
Test dataset (no. phrases)	458	877	1000	1015

B. The model's architecture

The architecture of our model adapts the one supplied by Castellucci et al. [4] to our specific ID and SF problems with 13 intents and 27 slots. It starts from a vanilla version of BERT, with two added layers on top of it. One dense layer is used for deciding about the intent, while another dense layer detects important slots associated with it. Fig. 4 depicts the model's architecture.

Text processing starts by tokenizing the input with a BERT-specific tokenizer, that uses the sub-word technique. Next, the tokens are fed into a pre-trained instance of BERT. A dropout layer with a rate of 0.1 is added in order to prevent overfitting and increase the overall performance of the system. The last step is predicting the intent and slots of the input utterance. Both classifiers contain a number of perceptrons equal to the number of intents and slots, respectively. Finally, the model outputs logit values, in the form of tensors. For intent detection, the output tensor contains 13 logit values. For slot filling, the model outputs a tensor with 27 logit values for each token in the sentence. To predict the final label, we choose the highest logit from each tensor, get its ID and convert it into natural language using a dictionary of intents/slots.

The optimizer selected for the perceptrons is Adam, with a learning rate of $3e-5$, and $\epsilon = 1e-8$, to avoid weights division

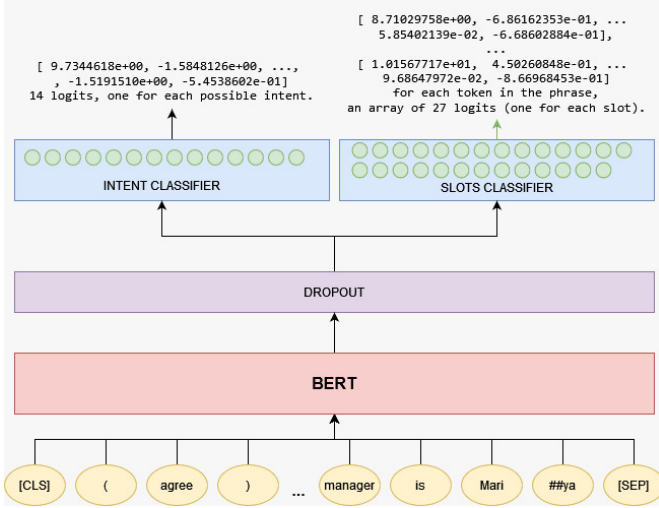


Fig. 4. Model’s architecture

by zero. Categorical cross-entropy is selected as the loss function for each task. The metric observed during training is sparse categorical accuracy. The batch size is 32 and the number of epochs is set to 4, as decided from the experiments performed in subsection III-D. All parameter values are set to respect the guidelines from Devlin et al. [8] when dealing with fine-tuning operations of BERT for NLP tasks.

C. Model’s performance evaluation

The main metrics that assess the quality of a model are intent accuracy, slot F1 score, and overall accuracy. The intent accuracy is a simple metric that tells us the percentage of times a model predicted the right intent, while the slot F1 score assures that the model keeps a high balance between precision and recall, to prevent wrong label predictions from propagating to further tasks. The overall accuracy measures whether both intent and slots were correctly predicted in one utterance, which is the desired behavior of a model.

D. Model’s baseline performance

The first step of the training process is to assess the model’s baseline capacity and determine the suitability of the selected architecture for our specific ID and SF tasks. Therefore, we have separately fine-tuned our BERT-based architecture presented in subsection III-B on ATIS [14] and SNIPS [15] datasets, comparing the results with other state-of-the-art BERT models.

ATIS consists of transcripts of audio recordings between humans and some automated airline travel inquiry systems. The training dataset contains 4478 utterances, while the test set has 893. It has 120 slot labels and 21 intent types in the training data. SNIPS includes phrases generated by the interaction between humans and a smart assistant that can execute a variety of desired tasks. The train set has 13084 phrases, while the test set consists of 700. There are 72 slot labels and 7 intent types.

The training procedure for ATIS and SNIPS is not straightforward as our data format is JSON, while ATIS and SNIPS data are texts. We have converted the ATIS and SNIPS datasets in our JSON format, preserving all important data and labeling. We run the training procedure with a number of epochs varying from 2 to 32.

Table II presents the performance of our fine-tuned BERT-based architecture for intent detection and slot filling on ATIS and SNIPS datasets for various number of epochs, together with the necessary time to train up to a given epoch. On the last column of table II we report a training efficiency measure computed according with eq. 1: the additional time needed to complete the training reported to the performance gain, measured with the help of the overall accuracy.

$$RT = \frac{T_x - T_2}{OA_x - OA_2} \quad (1)$$

where $x \in [2, 4, 8, 16, 32]$ is the number of epochs, T_x stands for the training time of the model fine-tuned for x epochs and OA_x represents the overall accuracy of that model.

On the test sets, compared to the 2 epochs variant, results have shown that using 4 epochs has the best time-to-train / overall accuracy increase ratio among all variants. On ATIS, the overall accuracy increased by 16.23%, needing only 12.8 seconds per gained percentage. Similarly, on SNIPS, the overall accuracy raised by 8%, leading to 30.4 seconds per gained percentage, 6 seconds faster than the 8 epochs variant. Although a higher number of epochs generated higher values for the overall accuracy, in order to prevent overfitting and waste of resources, we have concluded that it is best to keep the number of epochs equal to 4, which also follows the guidelines from BERT [8].

Comparing our results with the ones reported in the literature, we notice that our model yields good results on the individual tasks, but lower overall accuracy compared to the other ones. Nonetheless, as ATIS and SNIPS are more complex datasets than our custom-generated ones, our proposed architecture seems promising to solve our specific intent detection and slot-filling tasks.

IV. RESULTS

In this section we present the results of the model fine-tuned on the generated datasets presented in subsection III-A, with the above-described methodology.

First, we fine-tuned our architecture for the generated datasets described in table I. Results are presented in table III.

Training happens surprisingly quickly, as for a dataset of 7000 instances (GD3) it took 11 minutes and 48 seconds to finish the 4 training epochs. Increasing the number of instances in the dataset leads to a more than proportional increase in the training time. Models trained on a large enough number of conversations are good candidates for selection, as the intent and slot detection accuracy surpasses 99%. On the test set, all models reached 100% intent detection, while the slot F1 score ranged between 88%-99.6%. Although all models seem powerful enough to complete the task, the overall accuracy is

TABLE II
JOINT MODEL PERFORMANCE ON INTENT DETECTION AND SLOT FILLING
FOR BERT-BASED MODELS

Dataset	Epochs	Training time (s)	Slot F1	Intent Acc	Overall Acc	RT (s./%)
ATIS	2	175	87.42	95.97	54.54	-
	4	384	92.63	97.10	70.77	12.88
	8	515	94.59	97.08	78.61	14.13
	16	979	95.13	96.97	82.41	28.85
	32	1909	95.76	97.42	85.56	55.90
	BERT-SLU [5]		99.76	98.75	93.89	
	Qin et al. [7]		98.00	96.10	88.80	
	Joint BERT [3]		97.50	96.10	88.20	
	BERT-Joint [4]		97.80	95.70	88.20	
SNIPS	2	302	89.00	97.28	61.42	-
	4	545	91.26	97.71	69.42	30.38
	8	1032	93.34	98.00	74.71	36.64
	16	1982	93.74	97.85	77.71	58.32
	32	3900	93.20	98.29	77.50	119.28
	BERT-SLU [5]		98.96	98.78	96.76	
	Qin et al. [7]		98.80	97.10	93.10	
	Joint BERT [3]		98.60	97.00	92.80	
	BERT-Joint [4]		99.00	96.20	91.60	

TABLE III
PERFORMANCE METRICS FOR MODELS FINE-TUNED ON GENERATED
DATASETS WITH VARIOUS NUMBER OF CONVERSATIONS

Metrics / Training Dataset	GD1	GD2	GD3	GD4
Training intent accuracy	98.59	99.47	99.83	99.98
Validation intent accuracy	98.64	99.58	99.83	99.99
Training slot accuracy	100	99.97	100	100
Validation slot accuracy	99.75	100	100	100
Time to train (min.sec)	2.37	3.54	11.48	20.30
Test intent accuracy	100	100	100	100
Test overall slot F1 score	88.78	96.00	98.72	99.58
Test overall accuracy	80.40	90.99	96.80	99.40

the decisive metric for choosing the best model out of them. The GD1 model only has around 80% overall accuracy, which is the lowest among all. Models trained on GD3 and GD4 took significantly more time to train than the first one, but their overall accuracy is high enough to compensate for it. Finally, between the last two models, a comparison of the metrics for each slot was used to determine the best model. GD4 model has two important advantages over GD3 model: it finds slots that the latter one was not able to (ex. *remove_param*) and does not confuse new values and old values slots between themselves. Therefore, the GD4 model was chosen as the best-performing one. This proves that a powerful pre-trained model, such as BERT, can be fine-tuned with a well-designed dataset to obtain excellent results, leading to the possibility of integrating it into a conversational agent's architecture.

Table IV presents an extended analysis of the performance achieved by the GD4 model for each slot. Each line of the table presents the performance for a given slot of the user utterance. *new values* refer to novel values proposed by the user for a given slot and *old values* refer to slot values corresponding either to instances already existing in the organizational KB or values that the user tries to search (or retrieve) from the existing KB.

Within our general TOD architecture, correctly detecting the

TABLE IV
TEST DATASET PERFORMANCE METRICS FOR EACH SLOT

Slot / Metric	Precision	Recall	F1 score
entity	100	100	100
procedure	100	100	100
parameter to remove	100	100	100
hasClass	100	99.81	99.91
hasName	100	100	100
hasCode	99.84	100	99.92
hasStatus	100	100	100
hasManager	100	100	100
hasRole	100	100	100
<i>new values</i> hasClass	100	100	100
<i>new values</i> hasName	100	99.11	99.55
<i>new values</i> hasCode	100	100	100
<i>new values</i> hasStatus	98.94	100	99.47
<i>new values</i> hasManager	99.77	99.77	99.77
<i>new values</i> hasRole	98.51	100	99.25
<i>old values</i> hasClass	100	100	100
<i>old values</i> hasName	97.33	100	98.65
<i>old values</i> hasCode	100	100	100
<i>old values</i> hasStatus	100	100	100
<i>old values</i> hasManager	100	98.08	99.03
<i>old values</i> hasRole	100	96.00	97.96

entity type and the *procedure* is crucial. These let the system properly locate the corresponding relevant information in the KB that accompanies our TOD system. We indeed achieve this, as performance metrics for *entity* and *procedure* are equal to 100%.

The model is able to successfully detect the parameter slots (*hasClass*, *hasCode*, *hasName*, *hasManager*, *hasStatus*, and *hasRole* - see Fig. 2), crucial for the correct insertion / retrieval of instances in the organizational knowledge base.

The *new values* and *old values* slots were more difficult to detect. These slots may appear only in update requests, where the user wants to update existing instances that fit certain filters (hence, *old values*) with new information for some slots (hence, *new values*). The property of slot value being *new* or *old* is differentiated by the presence of other words in the utterance, thus the model has a harder job in identifying those values, and this ends in a lower performance. Although this problem, the GD4 model was the only one to successfully identify both types of slots, without any confusion between them, leading to high performance metrics. This highlights once again that a simple, yet powerful architecture using a pre-trained neural network such as BERT, together with enough, domain-related data can solve the task of Natural Language Understanding.

V. CONCLUSION

In this paper, we have successfully proven that a popular pre-trained model such as BERT, with minimal architectural additions, can solve the task of NLU on a custom-generated dataset. We have tested the proposed architecture on well-known datasets for intent detection and slot-filling tasks, such as ATIS and SNIPS, to establish a baseline verification of the architectural effectiveness. The in-depth analysis of the performance metrics for each slot in the test split has shown

that the model can yield high performance, when sufficient number of instances were seen during training.

Future work will focus on better-generated datasets, that may contain several variations of phrases to include more instances for each slot. Finally, the model will be included in a task-oriented dialogue system, to solve the NLU task and we will test its performance in real-life scenarios.

REFERENCES

- [1] D. Adiwardana, M. Luong, D. R. So, J. Hall, N. Fiedel, R. Thoppilan, Z. Yang, A. Kulshreshtha, G. Nemade, Y. Lu, and Q. V. Le, "Towards a human-like open-domain chatbot," *CoRR*, vol. abs/2001.09977, 2020. [Online]. Available: <https://arxiv.org/abs/2001.09977>
- [2] L. Qin, T. Xie, W. Che, and T. Liu, "A survey on spoken language understanding: Recent advances and new frontiers," in *30th International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021*, Z. Zhou, Ed. ijcai.org, 2021, pp. 4577–4584. [Online]. Available: <https://doi.org/10.24963/ijcai.2021/622>
- [3] Q. Chen, Z. Zhuo, and W. Wang, "BERT for joint intent classification and slot filling," *CoRR*, vol. abs/1902.10909, 2019. [Online]. Available: <https://arxiv.org/abs/1902.10909>
- [4] G. Castellucci, V. Bellomaria, A. Favalli, and R. Romagnoli, "Multi-lingual intent detection and slot filling in a joint BERT-based model," *CoRR*, vol. abs/1907.02884, 2019. [Online]. Available: <https://arxiv.org/abs/1907.02884>
- [5] Z. Zhang, Z. Zhang, H. Chen, and Z. Zhang, "A joint learning framework with BERT for spoken language understanding," *IEEE Access*, vol. 7, pp. 168 849–168 858, 2019. [Online]. Available: <https://doi.org/10.1109/ACCESS.2019.2954766>
- [6] J. Krone, Y. Zhang, and M. Diab, "Learning to classify intents and slot labels given a handful of examples," in *2nd Workshop on Natural Language Processing for Conversational AI*. Association for Computational Linguistics, 2020, pp. 96–108. [Online]. Available: <https://aclanthology.org/2020.nlp4convai-1.12>
- [7] L. Qin, T. Liu, W. Che, B. Kang, S. Zhao, and T. Liu, "A co-interactive transformer for joint slot filling and intent detection," in *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2021, Toronto, ON, Canada, June 6-11, 2021*. IEEE, 2021, pp. 8193–8197. [Online]. Available: <https://doi.org/10.1109/ICASSP39728.2021.9414110>
- [8] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: pre-training of deep bidirectional transformers for language understanding," in *2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, J. Burstein, C. Doran, and T. Solorio, Eds. Association for Computational Linguistics, 2019, pp. 4171–4186. [Online]. Available: <https://doi.org/10.18653/v1/n19-1423>
- [9] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," *J. Mach. Learn. Res.*, vol. 21, pp. 140:1–140:67, 2020. [Online]. Available: <http://jmlr.org/papers/v21/20-074.html>
- [10] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., 2020. [Online]. Available: <https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html>
- [11] H. Chen, X. Liu, D. Yin, and J. Tang, "A survey on dialogue systems: Recent advances and new frontiers," *ACM SIGKDD Explorations Newsletter*, vol. 19, no. 2, pp. 25–35, 2017. [Online]. Available: <https://doi.org/10.1145/3166054.3166058>
- [12] Z. Li, H. Wang, A. Albalak, Y. Yang, J. Qian, S. Li, and X. Yan, "Making something out of nothing: Building robust task-oriented dialogue systems from scratch," in *Alexa Prize TaskBot Challenge Proceedings*, 2022. [Online]. Available: <https://www.amazon.science/alexa-prize/proceedings/making-something-out-of-nothing-building-robust-task-oriented-dialogue-systems-from-scratch>
- [13] V. I. Iga and G. C. Silaghi, "Ontology-based dialogue system for domain-specific knowledge acquisition," in *31st International Conference on Information Systems Development, ISD 2023, Lisbon, Portugal, August 30 - September 1, 2023*. AIS Library, 2023, to appear.
- [14] C. T. Hemphill, J. J. Godfrey, and G. R. Doddington, "The ATIS spoken language systems pilot corpus," in *Speech and Natural Language: Proceedings of a Workshop Held at Hidden Valley, Pennsylvania, USA, June 24-27, 1990*. Morgan Kaufmann, 1990. [Online]. Available: <https://aclanthology.org/H90-1021/>
- [15] A. Coucke, A. Saade, A. Ball, T. Bluche, A. Caulier, D. Leroy, C. Doumouro, T. Gisselbrecht, F. Caltagirone, T. Lavril, M. Primet, and J. Dureau, "Snips voice platform: an embedded spoken language understanding system for private-by-design voice interfaces," *CoRR*, vol. abs/1805.10190, 2018. [Online]. Available: <http://arxiv.org/abs/1805.10190>
- [16] S. V. Ravuri and A. Stolcke, "Recurrent neural network and LSTM models for lexical utterance classification," in *INTERSPEECH 2015, 16th Annual Conference of the International Speech Communication Association, Dresden, Germany, September 6-10, 2015*. ISCA, 2015, pp. 135–139. [Online]. Available: <https://doi.org/10.21437/Interspeech.2015-42>
- [17] G. Mesnil, X. He, L. Deng, and Y. Bengio, "Investigation of recurrent-neural-network architectures and learning methods for spoken language understanding," in *INTERSPEECH 2013, 14th Annual Conference of the International Speech Communication Association, Lyon, France, August 25-29, 2013*, F. Bimbot, C. Cerisara, C. Fougeron, G. Gravier, L. Lamel, F. Pellegrino, and P. Perrier, Eds. ISCA, 2013, pp. 3771–3775. [Online]. Available: <https://doi.org/10.21437/Interspeech.2013-596>
- [18] D. Hakkani-Tür, G. Tür, A. Celikyilmaz, Y. Chen, J. Gao, L. Deng, and Y. Wang, "Multi-domain joint semantic frame parsing using bi-directional RNN-LSTM," in *INTERSPEECH 2016, 17th Annual Conference of the International Speech Communication Association, San Francisco, CA, USA, September 8-12, 2016*, N. Morgan, Ed. ISCA, 2016, pp. 715–719. [Online]. Available: <https://doi.org/10.21437/Interspeech.2016-402>
- [19] L. Zhang, D. Ma, X. Zhang, X. Yan, and H. Wang, "Graph LSTM with context-gated mechanism for spoken language understanding," in *34th AAAI Conference on Artificial Intelligence, AAAI 2020, New York, NY, USA, February 7-12, 2020*. AAAI Press, 2020, pp. 9539–9546. [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/6499>
- [20] W. Xu, B. Haider, and S. Mansour, "End-to-end slot alignment and recognition for cross-lingual NLU," in *2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, B. Webber, T. Cohn, Y. He, and Y. Liu, Eds. Association for Computational Linguistics, 2020, pp. 5052–5063. [Online]. Available: <https://doi.org/10.18653/v1/2020.emnlp-main.410>
- [21] Y. Hou, W. Che, Y. Lai, Z. Zhou, Y. Liu, H. Liu, and T. Liu, "Few-shot slot tagging with collapsed dependency transfer and label-enhanced task-adaptive projection network," in *58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, D. Jurafsky, J. Chai, N. Schluter, and J. R. Tetreault, Eds. Association for Computational Linguistics, 2020, pp. 1381–1393. [Online]. Available: <https://doi.org/10.18653/v1/2020.acl-main.128>
- [22] D. Bufeana, V. Niculescu, G. Silaghi, and A. Sterca, "Babes-Bolyai University's high performance computing center," *Studia Universitatis Babes-Bolyai, Informatica Series*, vol. 61, no. 2, pp. 54–69, 2016.