

Slovenian Language Assistant for Virtual Conversations (SLAVC)

Luka Škodnik, Robert Jutreša, Valter Hudovernik

Abstract

Keywords

Large Language Models, Machine Translation, Conversational AI, Slovene, Question Answering

Advisors: doc. dr. Slavko Žitnik

Introduction

Natural Language Processing (NLP) is an exciting field of Artificial Intelligence that focuses on teaching machines how to understand and respond to human language. In this report, we will discuss our project for the Natural Language Processing course 2022/23, where we aim to develop a Slovenian Conversational Agent called Slovenian Language Assistant for Virtual Conversations or SLAVC for short. We will provide an overview of our preliminary research, current implementation and explore various data sources that we plan to use and discuss possible models to train or finetune our model.

Related Work

Transformers & Large Language Models

Attention [1] is a key component of transformer models, which are widely used in natural language processing tasks such as language translation and text summarization. It is also a part of many modern large language models (LLM) that we will use during this project. In transformer models, attention mechanisms allow the model to focus on the most relevant parts of the input sequence at each step of processing. This is achieved by assigning weights to each element in the input sequence based on its relevance to the current step. By doing so, the transformer can capture long-range dependencies between different parts of the input sequence, which is particularly important for language processing.

SloBERTa [2] is a Slovene large language model. It is a large pre-trained masked language model based on the BERT architecture. The authors trained the model on a large corpus of Slovene text and evaluated it on various downstream

tasks such as text classification, named entity recognition, and part-of-speech tagging. The results show that SloBERTa outperforms existing Slovene language models and achieves competitive performance with state-of-the-art multilingual language models.

We examine the alignment of language model paradigm using reinforcement learning such as InstructGPT [3]. However due to the lack of resources we instead focus our attention to using already compiled high quality datasets for a one-time training of our model. However their work points out the total number of collective prompts used in training their model, which was 77k. This gives us an approximate estimate of how much data to use in the fine-tuning process. Additionally they discuss deduplication, however in future works such as [4] have pointed out that repeating prompts is not a big issue.

Datasets

The P3 dataset is a collection of prompted English datasets covering a diverse set of NLP tasks. The use of prompts allows for the creation of consistent and standardized data examples across different datasets, which can facilitate the development of new models and the comparison of results across different tasks [5].

A corpus of automatically annotated TV series subtitles for dialogue construction was developed in [6]. The authors used a combination of rule-based and machine-learning techniques to identify speaker turns and assign speaker identities. They evaluated their method on a corpus of subtitled TV series episodes and achieved high accuracy in speaker identification and turn-taking recognition. The resulting corpus was used for various downstream tasks such as emotion recognition and dialogue act classification.

GOS [7] is a reference corpus of spoken Slovene language. The methodology used to collect the corpus involved recording conversations of native Slovene speakers in various domains such as business, education, and social interactions. The transcription process involved annotating the recordings with orthographic, phonetic, and prosodic information. The resulting corpus was used for various research tasks such as acoustic modeling, speaker recognition, and speech synthesis.

Alpaca [4] by Stanford University is a set consisting of a data generation procedure, dataset, and training recipe. It is a fine tuned model from 7B LLaMA [8] on 52K instruction-following data generated by the techniques in the Self-Instruct [9] paper, with some modifications. In a preliminary human evaluation, it was found that the Alpaca 7B model behaves similarly to the GPT-3 model on the Self-Instruct instruction-following evaluation suite.

BLOOM [10] a multilingual LLM was trained on the ROOTS corpus [11], amounting to 1.61 terabytes of text that span 46 natural languages and 13 programming languages. Unfortunately Slovenian is not one of the available languages. However, due to vast collections of datasets available such as Hugging Face datasets [12], we focus on developing tools for robust translation in order to translate the large amounts of available data to Slovenian in order to facilitate the development of a conversational LLM.

SuperGLUE [13] is a benchmark for evaluating the performance of natural language understanding models. It consists of eight challenging natural language understanding tasks, including both textual entailment and commonsense reasoning tasks. The benchmark is designed to test the ability of models to handle more complex linguistic phenomena and to generalize to new examples. The authors compare the performance of several state-of-the-art models on SuperGLUE and find that there is still a significant gap between the best models and human performance.

Logar and Šikonja (2022) [14] discuss the challenges of question answering for less-resourced languages and presents an adaptation of the English UnifiedQA approach to the Slovene language. The adaptation uses encoder-decoder transformer models (SloT5 and mT5) to handle different question-answering formats, and existing Slovene adaptations of four datasets, as well as machine translation of the MCTest dataset. The study shows that a general model can perform at least as well as specialized models for answering questions in different formats. However, the performance of the Slovene model still lags behind that of English, and cross-lingual transfer from English is used to improve the results.

The paper "OpenAssistant: Aligning Large Language Models with Human Preferences using Open-Source Conversations" [15] presents a new corpus of human-generated, human-annotated assistant-style conversations called OpenAssistant Conversations. The corpus consists of 161,443 messages distributed across 66,497 conversation trees, in 35 different languages, annotated with 461,292 quality ratings, and was created through a worldwide crowd-sourcing effort

involving over 13,500 volunteers. The authors also release their code and data under fully permissive licenses, making their work easily accessible to the wider research community. We choose to adapt this dataset for training of a Slovene conversational Assistant. The structure of the dataset is depicted in Figure 1, on which we comment further in the following section.

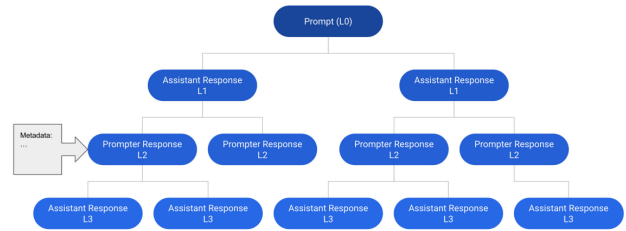


Figure 1. Dataset Structure An example of a conversation tree. All paths from the root node to the leaf prompts represent unique conversations. Additionally any two subsequent prompt and reply levels may also be used as training examples. *Figure taken from the original paper by Köpf and Kilcher et. al [15]*

Methods

Data translation

Recently the OpenAssistant Conversations Dataset (OASST1) has been released. We focus on translating the dataset in the conversation tree form where multiple replies can be nested to form conversations. We select a subset of trees that are “ready for export” as they have deleted spam messages and do not contain low quality messages and trees with only one prompt. Since OASST1 contains messages in various languages we translated the messages only from the 4 most common languages in the dataset: namely English, Spanish, Russian and German.

We construct pipelines for translation using NLLB-200 [16] and GoogleTranslate [17] using the deep translator python library [18]. Some of the messages also contain code which was replaced with a predefined substring at the time of translation and afterwards substituted back into the translated text. In this way we do not translate the code literally (eg. Translating predefined words such as *import* or *let.*). After initial inspection, this approach only works for code in the assistant replies as it has been nicely formatted in markdown quotes `>>>“bash ...“<<<`. However when users post questions about code, they simply paste it into the prompt and thus our method of detecting code snippets was unable to detect it. We decide to keep the code in the prompts translated for this version, as the target sequences are defined correctly. However, for future work developing a simple labeling model and training it on the annotated assistant responses without quotes would be a way to also detect code in the prompts.

We translate 8654 trees which totals to 78474 messages. We observe that the translations obtained using google trans-

late are best after manual inspection of translations.

Data processing

We consider two possible approaches for transforming the obtained translations into features to be fed into models. Since some of the models are limited in the lengths of the sequences they can accept and a lot of the conversations in our conversation trees are fairly long we firstly consider splitting the conversations into pairs. We do this by considering each prompt and reply in a conversation and at the end filter out the pairs where the prompter was the assistant (keep only pairs where assistant is the one replying). For training we define the prefix “*Uporabnik:* ” for all of the input prompts, while the replies do not have a specified prefix. We also tried training for 1 epoch with no prefixes or prefixed both of the sequences and the results were very similar so in the end this should not matter that much (**TODO: keep this?**). Since one prompt can have multiple replies we end up with training examples that have the same input (prompt) sequence but a different output (reply) sequence. We believe that this should not be a problem and could even contribute positively since in reality there can be multiple “correct” replies to the same prompt. In this way the models will hopefully learn multiple ways to respond to a prompt and we have more examples on which the model can learn how to formulate coherent sentences. **TODO: a je bilo smiselno vn fliknt vse prompte kjer ni biu reply role asistent? - mogoče napišemo da bi blo smiselno poskusiti tudi z vsemi ampak ni blo časa**

Another approach we considered is to generate all of the conversations and sub-conversations in the tree. This should also include all of the pairs as described above. We split this conversations into input and output sequences by taking the last response in the conversation as the output sequence and everything else as the input. In the case that the input sequence is too long for the model we only consider the end of it. All of the prompts and replies have the prefixes “*UPORABNIK:* ” or “*ASISTENT:* ” appended to them before the input sequence is constructed. Such data would be better for a larger model that can accept longer lengths and is capable of learning on more complex data.

As we are training sequence to sequence models we only need to tokenize the above described data without the need to perform any additional feature extraction. However we do have some additional information/features in the metadata of the translated dataset but we mostly ignore them except for the rank in the case of RRHF idk what to write here **TODO: keep this? RRHF**

Table 1 shows the size of our dataset and the train, validation and test splits that we created. We split the data with an 80-10-10 split on a shuffled dataset of conversation trees and only after process the data in a way that our models can accept. The ratio of the splits on the processed data is about the same as on the initially split conversation trees as we have a large enough dataset.

Data Type	Train	Val	Test	Total
Conversation Trees	6923	865	866	8,654
Messages	60,400	7,554	7,759	75,713
Unique Prompts	14,342	1,775	1,856	17,973
Prompt - Reply	37,902	4,737	4,875	47,514
Context - Reply	54,836	6,874	7,140	68,850

Table 1. Corpus Description: TODO describe table and reference in text

Models

T5

GPT

RRHF - sloberta We produce multiple replies using our fine-tuned language model and then employ RRHF [?], to select the best reply.



Figure 2. A random visualization. This is an example of a figure that spans only across one of the two columns.

Results

TODO: describe the used metrics - currently rouge and bertscore

Qualitative results

TODO: examples of well generated sequences for prompts, example(s) of bad ones - maybe ones where we can explain why they are bad. Maybe the “Kdo je France Prešeren?” example?

Our plans for evaluation are comparison of responses to ChatGPT in Slovene and OpenAssistant replies translated to Slovene. We plan to do a qualitative comparison of 10 replies for each model, which will be scored by users who will not know which model produced the reply. Additionally we consider inspecting similarity and correlations of BERT embeddings for the model replies.

TODO: chery-pick the examples: THIS IS WHERE WE LOOK GOOD

Discussion

Table 2. TODO.

	Model	Bertscore			ROUGE			
		Precision	Recall	F1	rouge1	rouge2	rougeL	rougeLsum
t5-small	baseline							
	pairs							
	conv							
	pairs-RRHF							
	conv-RRHF							
t5-large	baseline							
	pairs							
	conv							
	pairs-RRHF							
	conv-RRHF							
GPT2	baseline							
	pairs							
	conv							
	pairs-RRHF							
	conv-RRHF							

Table 3. TODO.

Prompt: Kdo je France Prešeren?

Reply: France Prešeren je bil francoski pesnik, ki se je rodil leta 1922 v Parizu.

Reply ChatGPT: TODO

Reply Open Assistant: TODO

Prompt: Kdo je France Prešeren?

Reply: France Prešeren je bil francoski pesnik, ki se je rodil leta 1922 v Parizu. Bil je francoski pisatelj in pisatelj, znan po svojem delu in delu v različnih delih sveta.

Reply ChatGPT: TODO

Reply Open Assistant: TODO

Prompt: Kdo je France Prešeren?

Reply: France Prešeren je bil francoski pesnik, ki se je rodil leta 1922 v Parizu. Bil je francoski pisatelj in pisatelj, znan po svojem delu in delu na področju glasbe, glasbe in glasbe.

Reply ChatGPT: TODO

Reply Open Assistant: TODO

Future work

For future work, combining translations from multiple sources and back-translation could improve the consistency and coherence of the conversational data.

Furthermore, conducting manual inspection to identify specific linguistic nuances and subtle information that may have been lost in the translation process. This involves closely examining the messages for wordplay, understanding error messages, and capturing the underlying semantics accurately. Additionally, we will explore the transformation of message formality by applying lemmatization techniques. By changing the prompts from a more formal style to an informal and conversational tone, we aim to create a dataset that better aligns with natural language interactions and user preferences.

These future research directions will contribute to the refinement and enhancement of translated datasets, facilitating improved performance and user experiences in multilingual conversational applications especially for few resource languages such as Slovene.

TODO: Add more parameters! LLaMA (7B) If we are to use the LLaMA model, we are still considering additional Slovenian pretraining on the following corpora; Gigafida [19] and ccKRES [20], we also consider the Šolar corpus [21], with considerations that the written language of primary and high school students could potentially have an impact on the models performance, and GOS [7] with similar considerations about spoken language.

Conclusion

References

- [1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [2] Matej Ulčar and Marko Robnik-Šikonja. Sloberta: Slovene monolingual large pretrained masked language

- model. *Proceedings of SI-KDD within the Information Society 2021*, pages 17–20, 2021.
- [3] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022.
 - [4] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. <https://github.com/tatsu-lab/stanford-alpaca>, 2023.
 - [5] Stephen H Bach, Victor Sanh, Zheng-Xin Yong, Albert Webson, Colin Raffel, Nihal V Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Fevry, et al. Promptsources: An integrated development environment and repository for natural language prompts. *arXiv preprint arXiv:2202.01279*, 2022.
 - [6] Leilan Zhang and Qiang Zhou. Automatically annotate tv series subtitles for dialogue corpus construction. In *2019 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pages 1029–1035, 2019.
 - [7] Darinka Verdonik, Iztok Kosem, Ana Zwitter Vitez, Simon Krek, and Marko Stabej. Compilation, transcription and usage of a reference speech corpus: the case of the slovene corpus GOS. *Language Resources and Evaluation*, 47(4):1031–1048, January 2013.
 - [8] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023.
 - [9] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions, 2022.
 - [10] Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022.
 - [11] Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Moral, Teven Scao, Leandro Werra, Chenghao Mou, Eduardo Ponferrada, Huu Nguyen, Jörg Froberg, Mario Šaško, Quentin Lhoest, Angelina McMillan-Major, Gerard Dupont, Stella Biderman, Anna Rogers, Loubna allal, Francesco Toni, and Yacine Jernite. The bigscience roots corpus: A 1.6tb composite multilingual dataset. 03 2023.
 - [12] Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, et al. Datasets: A community library for natural language processing. *arXiv preprint arXiv:2109.02846*, 2021.
 - [13] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. Superglue: A stickier benchmark for general-purpose language understanding systems. *Advances in neural information processing systems*, 32, 2019.
 - [14] Katja Logar and Marko Robnik-Šikonja. Unified question answering in slovene, 2022.
 - [15] Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. Openassistant conversations – democratizing large language model alignment, 2023.
 - [16] NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Hefernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Sermarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazıl Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. No language left behind: Scaling human-centered machine translation. 2022.
 - [17] Google LLC. Google Translate. <https://translate.google.com/?hl=sl>, Accessed 2023, March 21.
 - [18] deep-translator 1.10.1.
 - [19] Simon Krek, Tomaž Erjavec, Andraž Repar, Jaka Čibej, Špela Arhar Holdt, Polona Gantar, Iztok Kosem, Marko Robnik-Šikonja, Nikola Ljubešić, Kaja Dobrovoljc, Cyprian Laskowski, Miha Grčar, Peter Holozan, Simon Šuster, Vojko Gorjanc, Marko Stabej, and Nataša Logar. Corpus of written standard slovene gigafida 2.0, 2019. Slovenian language resource repository CLARIN.SI.
 - [20] N Logar Berginc, M Grčar, M Brakus, T Erjavec, Š Arhar Holdt, and S Krek. *Korpusi slovenskega jezika Gigafida, KRES, ccGigafida in ccKRES: gradnja, vsebina, uporaba*. Institut za slovenski jezik Frana Ramovša ZRC SAZU, Ljubljana, 2012.

- [21] Iztok Kosem, Tomaž Rozman, and Mateja Stritar. How do slovenian primary and secondary school students write and what their teachers correct: A corpus of student writing. In *Proceedings of Corpus Linguistics Conference*, pages 20–22, July 2011.