

Machine Translation

Romanian-to-Romani Translation

1st Semester of 2024-2025

Zibileanu Sabin

sabin.zibileanu@s.unibuc.ro

Grigore Dragoş - Gabriel

dragos.grigore@s.unibuc.ro

Abstract

We present the steps and the tools necessary to build a machine translation deep learning model for Romanian to Romani. Even though it is a low - resource language, the Romani language is a very popular one around the world and this project represents a starting point for future research and improvements on building machine translation models with this language.

How to use the project

1. the project is open-source and it can be found in this git repository: [MTRepo](#).
2. in the jupyter notebooks available in the repository, all the important parts of the project can be found: data acquisition, tokenization, training, loading the data, etc.
3. knowing the fact that we are working with a low - resource language (Romani), we decided that it was best to upload the whole project for future research and why not, future usage.
4. the dataset is structured in two txt files: romani.txt and romanian.txt containing the acquired samples and also 2 tokenizer directories generated from creating a custom Romanian-to-Romani tokenizer with a base MarianMT Model and another directory containing an extended pre-trained tokenizer.
5. there are also 2 directories in the repository: "data" and "tokenizer", where you can find the code for obtaining the dataset and the tokenizer
6. if you want to reproduce the results, you can run the notebook in the main branch (MTPProject.ipynb).

1. Data Acquisition

working with a low-resource language, we took into account the fact that we had to build a corpora from scratch almost.

however, we found a corpus online with around 7900 samples from Romanian to Romani [from HelsinkiNLP](#) (Christodoulopoulos and Steedman, 2014), which to our knowledge is in the Carpathian Romani dialect.

the available corpora was a great starting point, as it provided a considerable amount, but we wanted to add more data.

to add more data, we scraped samples from [Glosbe](#) as it provided Carpathian Romani samples translated into Romanian.

we chose the dialect Carpathian Romani as it is one that provides the most examples.

because Glosbe is a dictionary website, we had to search for keywords first in order to scrape the resulted data. Those keywords can be: conjunctions, verbs, adverbs, prepositions etc.

it is important to mention that some samples were added manually, because of the fact that certain keywords had multiple available phrases split into different pages and we could not load more pages with our scraper.

next on, we manually gathered a set of samples from different dictionaries, conversational guides, and a Romani course support. We added those samples manually, due to the fact that we mostly chose phrases and also the dialectal mix between the already existing samples in Carpathian Romani and these samples in Kalderash Romani.

a high mix of dialects between the samples could negatively affect the performance of the model.

2. Tokenization

073	the romanian-to-romani tokenizer was	initialization of the model. Precisely, we used	117
074	created using SentencePiece, with a vocab-	a MarianMT model with our custom tokenizer	118
075	ulary size of 4200.	and a pre-trained MarianMT model with the	119
076	after it was created, the tokenizer	already existing tokenizer that we combined	120
077	was loaded using MarianTokenizer (Junczys-	with the custom one. We extended the already	121
078	Dowmunt et al., 2018)	existing pre-trained tokenizer with our custom	122
079	the second tokenizer was an extended	made Romanian to Romani tokenizer.	123
080	version of a pre-trained tokenizer from		
081	HelsinkiNLP.		
082	3. Evaluation	• we chose this approach because the Mari-	124
083	for a low resource language like Romani,	anMT model supports multiple romance lan-	125
084	human evaluation would be the best.	guages and it is also a very well structured	126
085	in our project, however, we present the	Sequence2Sequence model for machine trans-	127
086	results obtained with the BLEU metric (Pap-	lation.	128
087	ineni et al., 2002).	• to our knowledge there are no papers covering	129
088		this task.	130
089	1 Introduction	• Zibileanu Sabin learned:	131
090	The Romani language, although a low - resource	– more about augmenting techniques for	132
091	one is, according to some estimates in the top 3%	datasets	133
092	of the world’s most spoken languages, containing	– how a tokenizer can be built	134
093	multiple dialects: Carpathian, Kalderash, Balkan	– how a Sequence2Sequence model works	135
094	etc. Considering the fact that there are no available	– more about the metrics used to evaluate	136
095	Romanian - Romani machine translation models	a machine translation system	137
096	(to our knowledge), we:	– how to build a corpora almost from	138
097		scratch	139
098	• tried building a machine translation model for	• Grigore Dragoş - Gabriel learned:	140
099	Romanian to Romani	– how a tokenizer can be built	141
100		– how a Sequence2Sequence model works	142
101	• The contributions of Zibileanu Sabin were:	– more about the metrics used to evaluate	143
102	– building the dataset and then cleaning it	a machine translation system	144
103	for duplicates		
104	– augmenting the dataset using 3 tech-	2 Approach	145
105	niques: back translation, random swap	This section covers our approach on this project	146
106	and sentence insertion	divided into detailed steps	147
107	– building the tokenizer		
108	– training the MarianMT models (Junczys-	1. Firstly we gathered the data: merged 2	148
109	Dowmunt et al., 2018)	datasets : one built from scratch using	149
110		Glosbe, dictionaries, conversational guides	150
111	• The contributions of Grigore Dragoş - Gabriel	and a romani course and the other one from	151
112	were	HelsinkiNLP . Afterwards, the dataset was	152
113	– training the MarianMT models (Junczys-	cleaned for duplicate samples.	153
114	Dowmunt et al., 2018)		
115	– preprocessing the text	2. Moving on, we felt that it was necessary to	154
116	– building the tokenizer vocabulary	try and use some data augmentation tech-	155
	– sanity checking the dataset after cleaning	niques, especially working with a low re-	156
		source language. The techniques that we used	157
		were: back translation , where the Roma-	158
		nian sentence was translated into English and	159
		then back to Romanian ; sentence insertion	160
		where we inserted a Romanian token and a	161

Romani token at the begining of a sentence, "Po del chavo - " for Romani and "Începutul propoziției -" for Romanian. Finally, the last technique was: **random char insertion** where a random character was introduced in a word from the sentence, simulating a typo. The augmentation was made on 15% of the training set.

3. The training proces took around 1 to 2 hours depending on the model and was done in Google Colab with a T4 GPU.

4. The deep learning tool used, was: MarianMT Transformer Model. For augmentation (random char) we used nlpaug, and Google - Translator for back translation

5. The results can be seen in Table 1

4 Conclusions and Future Work

This work presented the steps and the tools necessary to build a machine translation model for a low resource language. We liked this project because it allowed us to build something from the ground up and learn more about machine translation in the process. We believe that this project will surely leave room for improvement for the future and is a starting point for this task.

Results			
Data	Model	BLEU	Epochs
Augmented - data	Base MarianMTModel	1.969900	10
Augmented - data	Pre-trained MarianMTModel (en - RO-MANCE)	3.726000	3
Non-augmented - data	Base MarianMTModel	0.823531	10
Non-augmented - data	Pre-trained MarianMTModel (en - RO-MANCE)	2.415000	3

Table 1: BLEU Scores for Augmented and Non-Augmented Data Using Different Models

3 Limitations

The most notable limitation that can be seen is the quality of the results. This is caused by the low amount of data that was gathered. In order to build a high - standard machine translation model for Romanian to Romani, a large corpus is needed first (millions of sentences) that can be verified by a Romani speaker. It can be built from books, dictionaries, conversational guides and even songs why not. While working on this project, we found that there are not that many Romani books, that are translated in Romanian (excluding dictionaries).

Finally, a computational limitation to our end was the fact that we could not fine-tune the pre - trained model for more than 3 epochs because of time restrictions with the free T4 GPU in Google Colab.

4.1 References

References

- Christos Christodoulopoulos and Mark Steedman. 2014. [A massively parallel corpus: the bible in 100 languages](#). *Language Resources and Evaluation*, 49:1–21.
- Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Aji, Nikolay Bogoychev, André Martins, and Alexandra Birch. 2018. [Marian: Fast neural machine translation in c++](#). pages 116–121.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ACL '02, page 311–318, USA. Association for Computational Linguistics.